

# AfroBench: How Good are Large Language Models on African Languages?

Anonymous ACL submission

## Abstract

Large-scale multilingual evaluations, such as MEGA, often include only a handful of African languages due to the scarcity of high-quality evaluation data and the limited discoverability of existing African datasets. This lack of representation hinders comprehensive LLM evaluation across a diverse range of languages and tasks. To address these challenges, we introduce AFROBENCH—a multi-task benchmark for evaluating the performance of LLMs across 64 African languages, 15 tasks and 22 datasets. AFROBENCH consists of nine natural language understanding datasets, five text generation datasets, five knowledge and question answering tasks, and one mathematical reasoning task. We present results comparing the performance of prompting LLMs to fine-tuned baselines based on BERT and T5-style models. Our results suggest large gaps in performance between high-resource languages, such as English, and African languages across most tasks; but performance also varies based on the availability of monolingual data resources. Our findings confirm that performance on African languages continues to remain a hurdle for current LLMs, underscoring the need for additional efforts to close this gap.

## 1 Introduction

Large language models (LLMs) have risen to the fore of natural language processing (NLP) and also become increasingly commercially viable. These models have empirically demonstrated strong performance across a variety of NLP tasks and languages (Brown et al., 2020; Lin et al., 2021; Chowdhery et al., 2022; Chung et al., 2022). However, their performance on low-resource languages (LRLs), such as African languages, is largely understudied. This is problematic because there is a great disparity in the coverage of languages by NLP technologies. Joshi et al. (2020) note that over 90% of the world’s 7000+ languages are under-studied

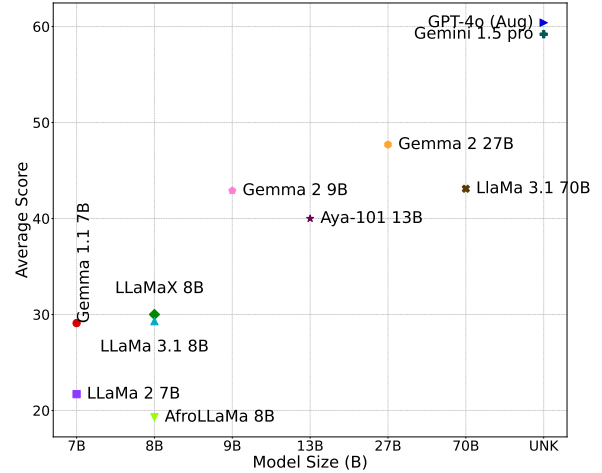


Figure 1: AFROBENCH average score on various LLMs

by the NLP community. Ideally, approaches to enhance language understanding should be applicable to all languages.

While there have been some recent evaluation of the performance of LLMs on several languages (Ahuja et al., 2023a; Lai et al., 2023; Robinson et al., 2023), the evaluation is focused on *closed models* like GPT-3.5 (Ouyang et al., 2022) and GPT-4 (OpenAI, 2023). Megaverse (Ahuja et al., 2023b) extended the evaluation to more models such as PaLM 2 (Anil et al., 2023) and LLaMa 2 (Touvron et al., 2023), Mistral (Jiang et al., 2023), Gemma (Mesnard et al., 2024) and Gemini Pro (Team et al., 2023). However, previous evaluation faces two main issues: (1) they cover only few tasks for African languages, for example, Megaverse only evaluated on part-of-speech, named entity recognition, and cross-lingual question answering for African languages, primarily due to *poor discoverability* of African languages benchmarks, *limited available evaluation data*, and *bias in the selection* of languages covered in the evaluation.<sup>1</sup> (2) Evaluation of LLMs needs to be continuous

<sup>1</sup>Belebele (Bandarkar et al., 2024) covers over 28 African languages, but Megaverse did not include any in their evaluation.

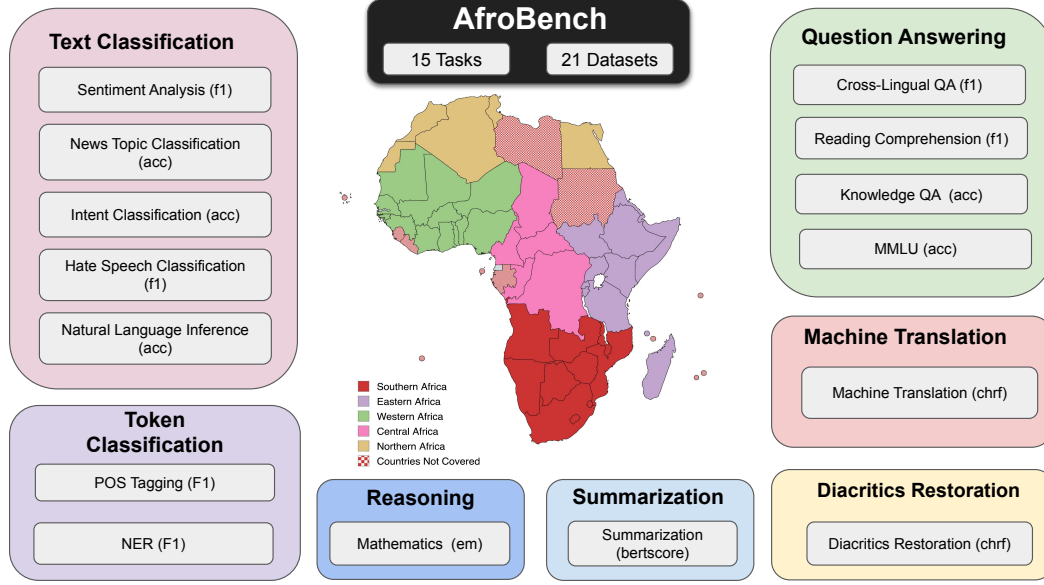


Figure 2: AFROBENCH: A comprehensive benchmark for evaluating performance of Large Language models on African Language tasks. The benchmark features 15 distinct tasks across 22 datasets and 64 indigeneous African languages. The benchmark covers diverse tasks with geographical coverage spanning different regions in Africa.

since many new LLMs have been released with improved multilingual abilities, but a comprehensive evaluation is not available for African languages.

In this paper, we address the challenges of previous large-scale LLM evaluation by introducing a new carefully curated benchmark known as **AFROBENCH** that comprises of 15 tasks, 21 evaluation data, and 64 indigeneous African languages. AFROBENCH consists of nine natural language understanding tasks, five text generation tasks, five knowledge and question answering tasks, and one mathematical reasoning task. Finally, we created a **new evaluation data**, AFRIADR for diacritic restoration of tonal marks and accents on African language texts. Leveraging AFROBENCH, we conduct an extensive analysis of the performance of LLMs for African languages from different language families and geographical locations.

For our evaluation, we compute the average performance score over the 15 tasks covered in AFROBENCH. Additionally, we introduce AFROBENCH-LITE that only cover a subset of seven tasks and 14 diverse languages in AFROBENCH which reduces the evaluation cost for a newly introduced LLM on our leaderboard. Figure 1 shows our evaluation on AFROBENCH, we find that proprietary models such as GPT-4o and Gemini-1.5 pro achieve +13 score improvement over Gemma 2 27B, our best-performing open model. We also compared the performance of English language to 14 African languages, finding that

GPT-4o and Gemma 2 27B achieve better performance than African languages by more than +25 and +40 score improvements respectively. This shows that the gap in the multilingual abilities of open models is wider than that of proprietary models. Finally, we compare the performance of LLMs to fine-tuned models based on AfroXLMR (Alabi et al., 2022), AfriTeVa V2 T5 model (Oladipo et al., 2023) and NLLB (NLLB Team et al., 2022) whenever training data is present. Results show that prompting LLMs often yields lower average performance than the fine-tuned baselines. Our findings show that more effort is needed to close the gap between the performance of LLMs for high-resource languages and African languages.

## 2 Related Work

**Large Language Model Evaluation:** Accurate and reproducible evaluation of language models is important as more and more models are being released. As these models are integrated into various applications, developing robust evaluation frameworks becomes paramount for understanding their true capabilities and limitations. As a result, the community has worked on developing evaluation frameworks (Gao et al., 2024; Fourrier et al., 2023; Liang et al., 2023), leaderboards (Chiang et al., 2024; bench authors, 2023; Fourrier et al., 2024) and benchmarks (Adelani et al., 2024b; Zhou et al., 2023; Hendrycks et al., 2021). While each of these evaluation tools focuses on assessing specific as-

Benchmark	# Tasks	# Datasets	# African Lang.	# LLMs	Closed LLMs evaluated	Dominant task(s)
ChatGPT-MT (Robinson et al., 2023)	1	1	57	1	GPT-3.5	MT
Mega (Ahuja et al., 2023a)	10	16	11	4	GPT-3, GPT-3.5-Turbo, GPT-4	POS, NER
Megaverse (Ahuja et al., 2024)	16	22	16	8	PaLM, GPT-3.5, GPT-4, Gemini Pro	POS, NER, XQA
SIB-200 (Adelani et al., 2024a)	1	1	57	2	GPT-3.5, GPT-4	Topic classification
Belebele (Bandarkar et al., 2024)	1	1	28	6	GPT-3.5-Turbo	QA
Uhura (Bayes et al., 2024)	1	2	6	6	Claude-3.5-Sonnet, GPT-4, 4o, o1-preview	QA
IrokoBench (Adelani et al., 2024b)	3	3	16	16	GPT-3.5, 4.0, Gemini-1.5-Pro, Claude OPUS	NLI, MMLU, Math.
AFROBENCH(Ours)	15	21	60	12	Gemini-1.5-Pro, GPT-4o	several

Table 1: **Overview of Related works that evaluated on African languages.** We included the number of tasks, datasets, African languages, LLMs evaluated, and the dominant tasks covering at least three African languages.

pects of language model capabilities - from basic linguistic understanding to complex reasoning tasks - the development of truly comprehensive benchmarks remains a significant challenge (Ruder, 2021; Biderman et al., 2024). These challenges stem from complex nature of language understanding and the stochastic nature of language models

**Multilingual LLM Benchmarks:** Benchmarks serve as a standard for measuring how systems have improved over time on across specific tasks and metrics. In the context of LLMs, multilingual benchmarks are crucial to assessing both the quality and practical utility of these models across diverse languages and tasks. Our primary focus lies in understanding LLM performance specifically for African languages, with several notable benchmarks (Table 1) having emerged in recent years to address this need. ChatGPT-MT (Robinson et al., 2023) evaluated the translation capability of GPT-4 and they find that it’s demonstrates strong performances on high-resource languages, the performance on low-resource languages is subpar. Belebele (Bandarkar et al., 2024) is a question answering task in 122 languages including 28 African languages for assessing reading comprehension abilities of LLMS. Mega (Ahuja et al., 2023a) and Megaverse (Ahuja et al., 2024) are multi-task multilingual and multimodal benchmarks in 83 languages including 16 African languages.

While these existing benchmarks have provided valuable insights, they collectively highlight a pressing need for more comprehensive evaluation that encompass a broader range of African languages and diverse tasks. Our research, through the development of AFROBENCH, addresses this gap by building upon and complementing existing work. We create a robust evaluation framework that assesses LLM performance across 64 African languages, evaluating capabilities across 15 distinct tasks. This expanded scope allows for a more nuanced and thorough understanding of LLM capabilities in African language contexts.

### 3 AfroBench

AFROBENCH is a comprehensive LLM evaluation benchmark designed to assess both proprietary and open LLMs across diverse Natural Language Processing (NLP) tasks in African languages. As shown in Figure 2, the benchmark encompasses 15 distinct tasks, spanning Natural Language Generation (NLG) and Natural Language Understanding (NLU), incorporating 21 curated datasets in 64 African languages. These evaluation tasks extend beyond traditional NLP benchmarks, such as text classification and named entity recognition, to include more challenging benchmarks such as mathematical reasoning and knowledge QA.

Each task within AFROBENCH has been carefully selected to assess different aspects of language model capabilities, from basic linguistic competency to more complex reasoning abilities. AFROBENCH also provides valuable insights into model behavior across different African language families and their unique linguistic features. All tasks and sub-tasks within AFROBENCH are evaluated using both zero-shot and few-shot prompting to guide model responses. To ensure consistent and reliable evaluation, we implement task-specific response constraints to facilitate systematic extraction and analysis of model outputs. For completion, we compare against existing SoTA encoder-only and encoder-decoder architectures that have previously demonstrated superior performance on individual tasks within the benchmark. This enables us to directly compare the performance of specialized models to general-purpose LLMs.

#### 3.1 Languages

We cover 64 African languages from seven language families (Afro-Asiatic, Atlantic-Congo, Austronesian, Indo-European, Mande, Nilotic, and English-Creole). 40 languages are from the Atlantic-Congo family, 12 from the Afro-Asiatic family, seven from Nilotic family, 2 Indo-European, 2 Creole languages, and 1 Austronesian language.

Figure 2 shows the geographical distribution of the languages covered in Afrobench and the full list of languages can be found in Appendix C.

### 3.2 Evaluation tasks and datasets

Our evaluation spans multiple datasets across 15 NLP tasks. While some of the selected datasets are multilingual and contain lots of languages, we focus specifically on the African language subsets, along with select high-resource languages (English, French, Portuguese and Arabic), due to their widespread use across different African regions.

We model all tasks as text-generation problems, where we combine inputs with prompts to guide language models in generating outputs under specific constraints. To ensure robust evaluation, we employ multiple prompts for each task with few- and zero-shot examples, which helps maintain consistency and minimize potential biases across different models. Our evaluation framework integrates the Eleuther LM Evaluation Harness (Gao et al., 2024) with custom evaluation scripts. The evaluation methodology varies by task type: text classification and multiple-choice tasks are assessed using log-likelihood evaluation, which measures the probability of a prompt-generated continuation containing the expected response, while all other tasks utilize free-form generation approaches.

Next, we present a breakdown of the tasks, sub-tasks and specific datasets contained in Afrobench.

#### 3.2.1 Text Classification

**Sentiment Classification:** We evaluate NOLLYSENTI (Shode et al., 2023) and AFRISENTI (Muhammad et al., 2023). AFRISENTI evaluates sentiment analysis of tweets across 14 African languages, while NOLLYSENTI focuses on movie review sentiment in four African languages.

**Topic Classification:** We evaluate SIB-200 and MASAKHANEWS (Adelani et al., 2023) that covers 57 and 16 African languages, respectively. The topic categories could be general topic such as *business*, *entertainment*, *health*, *politics* etc.

**Intent Classification:** INJONGOINTENT<sup>2</sup> is an intent classification task in 16 African languages. The goal is to classify an utterance into one of 40 intent types from different domains.

**Hate Speech detection:** AFRIHATE (Muhammad et al., 2025) is a multilingual hate speech and

abusive language datasets in 15 African languages for tweets. Each tweet can be categorized into one of *abusive*, *hate* or *neural* label.

**Natural Language Inference:** AFRIXNLI (Adelani et al., 2024b) is a dataset collection in 16 African languages where each datapoint is a pair of sentences (a premise and a hypothesis) and the task is to classify each pair as an *entailment*, *contradictor* or *neural* pair.

#### 3.2.2 Token Classification

**Named Entity Recognition (NER):** We evaluate entity recognition for 20 African languages on MASAKHANER-X (Ruder et al., 2023)—an extension of MASAKHANER dataset (Adelani et al., 2021, 2022b) that converts NER tags from CoNLL format into a text generation task of predicting entities with a delimiter, “\$” between them.

**POS Tagging:** MASAKHAPOS (Dione et al., 2023) is a part-of-speech tagging dataset in 20 African languages created from news articles. Each token is categorized into one of the 17 POS tags.

#### 3.2.3 Reasoning:

**Mathematical reasoning** We evaluate on AFRIMGSM (Adelani et al., 2024b), an extension of the MGSM dataset to 16 African languages. The question is a grade school level question, and a single digit answer.

#### 3.2.4 Question Answering

**Cross-Lingual Question Answering (XQA):** AFRIQA (Ogundepo et al., 2023) is a cross-lingual QA task with questions in 10 African languages and context passages in English or French. The goal is to extract the span with the right answer from the text.

**Reading Comprehension:** We evaluate on NAIJARC (Aremu et al., 2024), a multi-choice reading comprehension dataset in three African languages and BELEBELE (Bandarkar et al., 2024), a multi-choice reading comprehension task for 122 languages including 28 African languages.

**Knowledge QA:** We focus on two human-translated MMLU datasets: OPENAI-MMLU<sup>3</sup> and AFRIMMLU (Adelani et al., 2024b) that covers 3 and 17 African languages respectively.

<sup>2</sup><https://github.com/masakhane-io/masakhane-nlu>

<sup>3</sup><https://huggingface.co/datasets/openai/MMMLU>



Both tasks span multiple subjects and follow a four-option multiple-choice format. Although, the subjects covered by AFRIMMLU are only five. We also extend our evaluation to the human translation of *scientific Arc-Easy* benchmark in six African languages UHURA (Bayes et al., 2024).

### 3.2.5 Text generation

**Machine translation (MT):** Our MT benchmark includes the following datasets: MAFAND (Ade-lani et al., 2022a), FLORES (Goyal et al., 2022), NTREX-128 (Federmann et al., 2022) and SALT (Akeru et al., 2022) covering 57, 25 and 21 translation direction to African languages. All translations are from English except for the MAFAND benchmark with a few languages whose source is French.

**Summarization:** Given a news article, our goal is to generate its summary based on the popular XL-SUM dataset (Hasan et al., 2021) covering 10 African languages.

**Automatic Diacritics Restoration (ADR):** This is a **new benchmark** we introduce called **AFRI-ADR**. Given a sentence in a language, say “*Sugbon sibesibe, Mama o gbagbo*” (in Yorùbá), the model’s goal is to add the missing tonal marks and accents, say “*Ṣùgbọ̀n síbẹ̀síbẹ̀, Màmá ò gbàgbọ̀*”. We cover five African languages for this task.

## 4 Experimental setup

### 4.1 Fine-tuned baselines

For the tasks with available training data, we use available task-specific trained models, such as NLLB-200 3.3B for MT, and fine-tuned multilingual encoders or encoder-decoder T5 models on applicable datasets. We fine-tune AfroXLMR (Alabi et al., 2022)—one of the SoTA BERT-style encoders for African languages on each of the NLU tasks. For summarization and ADR, we fine-tune AfriTeVa V2 Large (Oladipo et al., 2023) on the available training data of each task. While AfriTeVa V2 outperformed mT5 (Xue et al., 2021) overall, its tokenization failed for Fon language, so we fine-tune mT5-large, which as a more diverse tokenizer, for the language.

### 4.2 LLMs Evaluated

We evaluate two broad categories of Large Language Models (LLMs): **Open Models** and **Closed Models**. We evaluate 10 open models: Llama 2 7B (Touvron et al., 2023), Gemma 1.1 7B (Mesnard

et al., 2024), LLAMA 3 series (3 8B, 3.1 8B and 3.1 70B) (Dubey et al., 2024), LlamaX 8B (Lu et al., 2024) (an adapted LLaMa 3 8B to 100 languages), AfroLlama 8B <sup>4</sup> (an adapted LLaMa 3 8B to Swahili, Xhosa, Zulu, Yoruba, Hausa and English languages), GEMMA 2 (9B & 27B) (Riviere et al., 2024), and Aya-101 (an instruction-tuned mT5 encoder-decoder model on massively multilingual prompted dataset). Finally, we evaluate on two popular proprietary models: GPT-4o and Gemini-1.5 pro (Reid et al., 2024). We provide full description of the LLMs in Appendix B.

**Prompts used for evaluation** We make use of *five* different prompts in the evaluation of each task except the text generation tasks, and we report the best prompt in the paper. For the text generation tasks, we reduce the number of prompts to *three* since the generation is often time consuming and expensive especially for summarization tasks. Moreover, we find that performance is less sensitive to prompt templates, unlike the NLU tasks. The prompt templates are provided in Appendix D.

**Few shot evaluation** We restrict the few shot evaluation to the best closed and open models. We fixed the number of examples to *five*, except for AfriMGSM whose number of examples is *eight*.

## 5 Results

### 5.1 AfroBench Evaluation

Table 2 shows the overall results across all the 15 tasks and 22 datasets. Our **first** observation is that closed models such as GPT-4o and Gemini-1.5 pro achieve better performance than the best open model, Gemma 2 27B with differences of +12 or more points on average performance. This shows that the gap in performance is wider for low-resource African languages than for high-resource languages, such as English, when using open models. **Secondly**, we find that performance gap varies across different tasks. Knowledge intensive and reasoning tasks such as ARC-EASY, MMLU, MATH have the largest gaps of +29.4, +19.9, +22.6 respectively, when we compare the performance of GPT-4o to Gemma 2 27B. In general, performance gets better with newer versions of LLMs (e.g. Gemma 1.1 7B vs. Gemma 2 9B and Llama 2 7B vs. Llama 3.1 8B ) and model sizes (Gemma 2 9B and Gemma 2 27B). This suggests that newer iterations of models are getting better

<sup>4</sup><https://huggingface.co/Jacaranda/AfroLlamaV1>

Tasks Metrics	natural language understanding							QA		knowledge		reasoning	text generation					
	POS acc	NER F1	SA F1	TC acc	Intent acc	Hate F1	NLI acc	XQA F1	RC F1	Arc-E acc	MMLU acc	Math EM	MT ChrF	Summ BertScore	ADR ChrF	ALL AVG	FT. AVG	
Fine-tuned baselines																		
AfroXLMR	89.4	84.6	72.1	74.4	93.7	77.2	61.4						en/fr-xx	xx-en/fr				
mT5/AfriTeVa V2 1B								52.5	N/A	N/A	N/A	N/A			72.3	79.4	70.4	
NLLB 3.3B													40.4	47.8				
Prompt-based baselines																		
open models																		
Gemma 1.1 7B	38.6	27.9	43.3	45.3	9.4	24.3	34.0	17.4	38.1	32.2	28.6	4.6	11.7	9.7	49.1	50.8	29.1 29.7	
LLaMa 2 7B	27.9	15.6	42.3	19.4	1.5	21.9	33.8	13.7	24.3	23.3	25.6	2.0	10.5	20.3	46.9	30.4	22.5 22.2	
LLaMa 3 8B	48.5	22.7	43.6	37.0	2.1	27.8	35.4	12.6	27.6	32.0	27.4	5.1	15.9	27.7	66.2	26.1	28.6 28.6	
LLaMaX 8B	41.6	0.0	51.9	49.8	5.6	28.6	40.8	2.2	29.7	39.9	28.3	4.0	22.7	35.0	50.7	49.4	30.0 29.0	
LLaMa 3.1 8B	47.1	11.5	50.5	46.7	6.0	23.6	36.6	21.8	39.5	32.8	31.4	6.8	16.4	28.5	43.7	25.9	29.3 28.1	
AfroLLaMa 8B	0.0	3.5	43.4	19.8	0.8	18.4	35.9	21.8	24.1	37.2	25.8	3.7	8.4	9.5	50.8	5.2	19.3 17.6	
Gemma 2 9B	51.9	40.3	60.0	56.0	29.2	29.9	40.3	45.9	51.6	53.4	37.1	18.7	24.8	29.1	66.1	51.6	42.9 42.9	
Aya-101 13B	0.0	0.0	63.4	70.3	42.4	31.0	51.5	62.5	60.7	59.6	30.9	4.4	23.4	37.9	52.4	50.4	40.1 37.7	
Gemma 2 27B	55.1	50.8	63.4	62.4	33.0	45.5	42.8	50.5	53.9	56.3	40.5	27.0	27.9	32.9	66.4	55.1	47.7 48.3	
LlaMa 3.1 70B	54.1	14.4	52.2	57.7	34.0	49.0	38.0	44	49.7	54.9	39.9	23.2	25.1	37.9	67.6	51.7	43.3 42.6	
proprietary models																		
Gemini 1.5 pro	60.8	41.8	68.3	76.7	74.3	62.1	62.0	40.5	52.7	84.8	57.6	52.3	37.6	41.7	66.7	55.6	58.5 58.9	
GPT-4o (Aug)	62.8	40.7	68.0	74.8	74.0	63.5	64.3	43.4	69.2	85.7	60.4	49.6	35.1	40.7	66.5	54.9	59.6 58.1	

Table 2: **AfroBench Evaluation results on fine-tuned models and LLMs.** We cover 15 tasks, 21 datasets, and 60 African languages in the evaluation. The best closed and open LLMs are highlighted in Cyan . We **bolden** the best result per task in each column. We provide average on **ALL** tasks and on those with fine-tuned baselines (**FT.**)

on low-resource languages, although with limited improvements on knowledge intensive tasks. **Finally**, while LLMs have made significant progress, they still fall behind their *fine-tuned baselines* (**FT. AVG**) when training data is available for a task. The gap in performance is around +11.5 on average, showing that curating annotated datasets for low-resource languages is still beneficial since the capabilities of LLMs lags behind. We provide task and per-language results in Appendix A and E.

## 5.2 AfroBench-LITE evaluation

In the AFROBENCH-LITE evaluation, we restrict the evaluation to seven LLMs, and seven tasks, and compare performance gap to English.

**Large gap in performance when compared to English** One striking observation is that open models such as LLaMa 3.1 70B and Gemma 2 27B have competitive performance to closed models on English language with  $-5$  to  $-2$  performance gap. However, when compared to African languages, GPT-4o and Gemini-1.5 pro achieves an average score better than Gemma 2 27B by more than 20 points on AFROBENCH-LITE. These results suggest that current LLMs especially the open models, are more biased towards *English* and a few high-resource languages. Adapting LLMs for a region of African languages could help bridge the gap. For instance, we see that continually pre-training Llama 3 8B, that resulted in LlamaX 8B shows slight overall performance of +1.4 or more over vanilla Llama 3 8B in Table 2. However, to further boost performance, better adaptation techniques are

Model	Lang	Intent	TC	NLI	RC	MMLU	Math	MT	
								en/fr-xx	AVG
Gemma 1.1 7B	eng	72.1	86.3	59.2	87.9	44.6	20.8	26.1	56.7
	africa	10.2	42.0	34.6	34.1	27.3	5.1	10.9	23.5
Gemma 2 9B	eng	36.3	82.5	70.7	<b>93.7</b>	69.8	68.8	67.9	70.0
	africa	27.8	64.0	40.9	49.3	36.1	21.7	37.2	39.6
Aya-101 13B	eng	78.0	82.8	67.0	86.1	42.8	11.6	64.2	61.8
	africa	40.2	76.0	52.4	59.7	30.3	4.9	31.8	42.2
Gemma 2 27B	eng	84.0	<b>89.3</b>	67.8	93.4	75.6	85.6	68.5	80.6
	africa	31.4	66.6	43.7	52.1	40.8	30.6	39.1	43.5
LLaMa 3.1 70B	eng	84.5	88.3	59.5	93.2	76.4	86.8	<b>71.6</b>	80.0
	africa	36.9	61.9	38.4	45.3	40.6	26.5	29.6	39.9
Gemini 1.5 pro	eng	<b>86.8</b>	88.7	88.5	69.6	<b>88.8</b>	86.8	69.1	82.6
	africa	75.6	<b>81.3</b>	63.6	54.4	<b>62.6</b>	<b>57.7</b>	<b>44.2</b>	62.8
GPT-4o (Aug)	eng	86.2	89.2	<b>89.2</b>	84.3	88.0	<b>88.8</b>	70.2	<b>85.1</b>
	africa	78.4	83.0	66.3	70.3	63.1	57.3	43.6	<b>66.0</b>

Table 3: **AfroBench-LITE Evaluation:** LLM baselines on 7 datasets spanning 14 African languages. Tasks were selected for broad NLP coverage, prioritizing language consistency. The best score per task is in **bold**.

required.

**Performance varies across languages** Figure 3 shows the results for per-language performance scores of 14 languages in AFROBENCH-LITE. Our result shows that performance correlates with the available monolingual texts on the web (Kudugunta et al., 2023). We find that Swahili (swa) with over 2.4GB of monolingual texts has the highest performance among the African languages, while Wolof with the smallest monolingual data (5MB) has the lowest performance. While this data size estimates are approximate, it shows that there is a need to invest more on developing language texts for many African languages for them to benefit in the LLM age. For most languages, GPT-4o gives the best overall results except for Amharic (amh) where

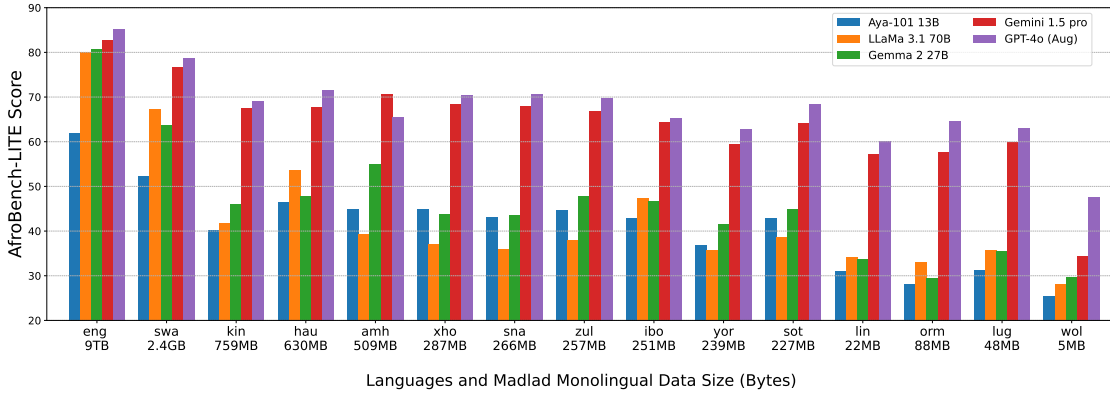


Figure 3: AfroBench-LITE performance of various models across African languages, plotted against the availability of monolingual data (MADLAD byte size).

Tasks	# shots	POS	NER	SA	TC	Intent	Hate	NLI	XQA	RC	MMLU	Math	MT en/fr-xx	MT xx-en/fr	SUMM	ADR	AVG
Gemma 2	0-shot	<u>55.1</u>	<u>50.8</u>	58.6	57.3	35.2	45.5	42.8	50.5	53.6	39.9	27.0	32.4	32.4	66.4	55.1	46.8
	5-shot	43.9	14.5	<u>59.7</u>	<u>62.5</u>	<u>56.7</u>	<u>57.3</u>	<u>56.0</u>	<u>52.4</u>	<u>58.3</u>	<u>44.8</u>	14.4	22.7	<u>34.9</u>	55.5	31.2	43.3
Gemini 1.5	0-shot	<u>60.8</u>	<u>41.8</u>	62.6	74.5	<u>74.3</u>	62.1	62.0	40.5	<u>53.0</u>	<u>60.2</u>	52.3	35.4	41.7	66.7	55.6	56.2
	5-shot	33.2	37.4	<b>64.5</b>	<b>77.3</b>	73.4	<u>64.1</u>	<u>35.9</u>	28.7	24.4	46.0	<b>61.4</b>	<u>37.4</u>	<u>43.1</u>	<b>70.4</b>	<b>63.4</b>	50.7
GPT-4o	0-shot	<u>62.8</u>	40.7	<u>62.6</u>	72.5	<u>74.0</u>	63.5	<b>64.3</b>	43.4	69.1	<u>60.0</u>	49.8	31.5	41.0	66.5	54.9	57.1
	5-shot	62.4	<u>45.0</u>	62.3	<u>72.9</u>	71.6	<b>69.3</b>	64.2	40.0	<b>71.9</b>	59.7	<u>56.1</u>	<b>33.9</b>	<b>43.3</b>	<u>67.9</u>	<u>62.7</u>	<b>58.9</b>

Table 4: **Few-shot Evaluation.** The better score between each model’s 0-shot and 5-shot is in underlined.

Gemini-1.5 pro was better. For the open models, Gemma 2 27B achieves better performance on eight out of the 14 languages, even better than LLaMa 3.1 70B that is more than twice its number of parameters. Although Aya-101 covers 100 languages in its pre-training and often achieves better performance on NLU tasks in AFROBENCH-LITE, it often struggles with math reasoning and MMLU, leading to worse overall results.

### 5.3 Few-shot results

Table 4 shows the result of zero-shot and few-shot evaluation on three LLMs: Gemma 2 27B, Gemini-1.5 pro and GPT-4o. The benefit of few-shot varies for different LLMs and tasks. For GPT-4o, we find that across all tasks, there is an average improvement of +1.8 while the other LLMs dropped in performance on average. The tasks that benefits the most from the few-shot examples are hate speech detection, math reasoning and ADR with +5.8, +6.3 and +7.8 respective points improvement. The result shows that few-shot examples are important for teaching LLM a new task it is unfamiliar of such as ADR since the rules of adding diacritics are not provided during the 0-shot, therefore, 5-examples, provides some demonstration to the LLMs on how to perform the task especially for low-resource languages such as Ghomálá’ and Fon with small monolingual data on the web. These two

languages improved by +16.4 and 7.2 respectively, while the other languages such as Igbo, Wolof and Yorùbá achieved more than +5.0 boost in ChrF scores. Similarly, for **Gemini-1.5 pro**, we observed consistent performance boost for ADR with 5 demonstration examples.

For both GPT-4o and Gemini-1.5 pro, math reasoning improved significantly by more than +6.0 points showing additional benefit of few-shot examples on reasoning tasks as shown in several evaluations such as Dubey et al. (2024). For Hate speech, we provided detailed explanation on the distinction between “abusive” content and “hate” in the prompt, but this is often confusing even for native speakers of the language, who often need examples of such sentences to improve annotation. We found that LLMs also require such additional examples to be able to better predict if a tweet is offensive. In general, Gemma 2 27B improved for several NLU but did not benefit from additional examples for the token classification, math reasoning, summarization and ADR tasks.

## 6 Discussion

### 6.1 Prompt variability

In our evaluation, we present results for the Best prompt rather than the Average results over several prompts to ensure no LLM is at a disadvantage due to their sensitivity to prompt templates. Here,

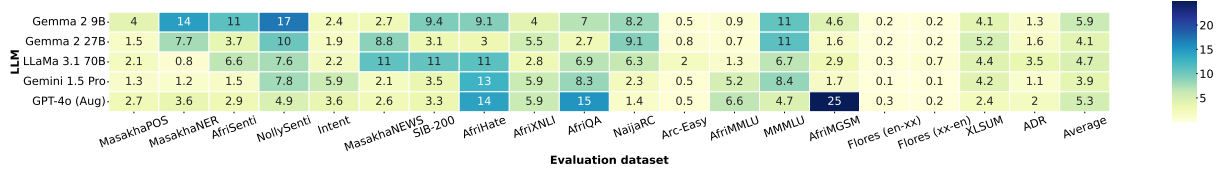


Figure 4: **Prompt Variability:** Heatmap of the difference between the Best and Average prompt results.

we analyze the difference in the performance scores between the Best prompt and the average over five prompts (or three prompts for the NLG tasks).

Figure 4 shows the result of our analysis across 18 tasks. Our **first observation** is that LLMs are not sensitive to different prompts when evaluating text generation tasks, all LLMs have lower than 6 point difference, and the task that is the least sensitive is machine translation (FLORES ). The **second observation** is that Gemini-1.5 pro is the least sensitive LLM to different prompt templates on average. The gap in performance across different prompts is often small for several NLU tasks. Interestingly, we find that GPT-4o is very sensitive to prompts for a few tasks such as hate speech, cross-lingual QA and math reasoning—which explains the large difference in performance scores. This analysis shows the benefit of using several prompts in evaluation, although, the benefit for text generation tasks are limited. Finally, we find that the largest variability is by a small sized Gemma 2 9B , which shows that, smaller LLMs requires more prompt template search than bigger models as shown that Gemma 2 27B is less sensitive.

## 6.2 Qualitative Analysis

Table 5 shows the benefit of few-shot examples on ADR, hate speech and math reasoning—the three tasks that improved the most with few-shot examples. For the ADR evaluation on Ghomáá’, we saw more than 60.0 ChrF point improvement, and noticed that only few characters have the wrong diacritics unlike the zero-shot setting. Similarly, for hate speech, without the few-shot example, the LLM focused on the abusive word “oloriburuku” (i.e. brainless), however, when we consider the *target* to which the tweet is addressed to, it is obvious that it was referring to an entire tribe in Nigeria, which is “hate”. In the definition of “hate” provided in the prompt, and some examples provided, this is clearer to the model than without any demonstration examples. Finally, for the math reasoning, without few-shot, the LLM reason in Yorùbá, which often gives poor result, however, when pro-

Lang.	Shot	Input / Output
ADR: Gemini-1.5 pro		
bbj	<b>Input:</b>	A jwó guŋ tsó awé a ló nǎŋ kwitá
	<b>Target:</b>	À jwó guŋ tsó awé a ló nǎŋ kwitá
	<b>0-shot</b>	à jwé gúŋ tsó áwé à ló nǎŋ kwitá (ChrF: 21.4)
	<b>5-shot</b>	à jwó guŋ tsó awé a ló nǎŋ kwitá (ChrF: 81.6)
Hate speech: Gemini-1.5 pro		
yor	<b>Input:</b>	oloriburuku ni pupo ninu awon igbo yii
	<b>English:</b>	these igbo people are brainless
	<b>Label:</b>	hate
	<b>0-shot</b>	abuse
	<b>5-shot</b>	hate
Math reasoning: GPT-4o		
yor	<b>Input:</b>	Ryan gbin òdòdó 2 ní ojúmọ́ sí inú oḡbaà rè. Léyìn ojọ́ 15, òdòdó mèlòb ní ó ní tí 5 ò bá wù?
	<b>English:</b>	Ryan plants 2 flowers a day in his garden. After 15 days, how many flowers does he have if 5 did not grow?
	<b>Answer:</b>	25
	<b>0-shot</b>	ryan ní òdòdó 30 tí ó bá ní gbin 2 ní ojúmọ́
	<b>8-shot</b>	the number of flowers remaining is 30 - 5 = 25.

Table 5: **Qualitative Analysis** comparison of the 0-shot and 5-shot samples on ADR, Hate speech and Math.

vided few-shots in Yorùbá, we observe that the LLM suddenly switched to English to answer the question, because the Chain-of-Thought answers in the few-shots are in English and it also reasons better in English, further boosting performance.

## 7 Conclusion

In this paper, we introduce a new benchmark, AFROBENCH, that aggregates existing evaluation datasets for African languages, and added a *new* dataset focused on diacritics restoration. AFROBENCH comprises of 15 NLP tasks, 22 datasets, and 64 African languages under-represented in NLP. We evaluate the performance of several closed and open LLMs on these tasks, showing that they all fall behind of fine-tuning baselines. We also show large performance gap compared to English, although we notice the gap is smaller for closed models such as GPT-4o and Gemini-1.5 pro. Through this benchmark, we have created a leaderboard focusing on LLM evaluation for African languages, which will be maintained going forward with additional tasks, LLMs and languages. We will be releasing our prompts and tasks configurations to Eleuther *lm-eval*. We hope this encourages the development of more African-centric LLMs for African languages.



## 8 Limitation

In today’s NLP landscape, large language models are generalist models that are capable of performing multiple NLP tasks without the need for special training on these tasks. These models are often multilingual and are able to perform tasks in multiple languages. Our research examines how these models perform specifically with African languages, revealing performance disparities when compared to more resourced languages. In this section, we discuss some of the limitations of our research methodology and findings.

**1. Training Data Transparency and Contamination:** One of the challenges in evaluating large language models lies in the limited visibility into their training data composition. While organizations frequently publish training documentation, many reports lack comprehensive details about data mixtures and language distributions across different training stages. There are multiple ways that this lack of transparency impacts the findings of our research. Without knowledge of the data mixture, we cannot determine whether or by how much our evaluation sets overlap with the training dataset. Thus, we cannot conclude that superior performance on certain tasks is a true demonstration of generalization or merely the models exposure to similar content during training. In the context of African languages, knowledge about the training data helps us access other factors such as cross-lingual transfer that might help us understand and better analyze evaluation results. A clear understanding of training data composition serves as a crucial foundation for meaningful model evaluation. It helps establish the validity of performance metrics and provides essential context for interpreting results across different languages and tasks.

**2. Limited Selection of LLMs and Evaluation Costs:** We are only able to evaluate a limited set of LLMs due to the computational and financial costs associated with model access and inference. Language models are accessed using two primary methods; loading the pretrained checkpoints directly or via an API service. While providers like Together AI offer access to open-source models and companies like OpenAI provide proprietary model access, both approaches incur considerable costs that directly impact the scope of evaluation studies. In our evaluation, the costs were substantial, requiring approximately \$2,500 each for Gemini-1.5 pro and GPT-4o model access, with an additional

\$1,200 for utilizing the Together.AI platform. The total evaluation costs manifests in two key dimensions; First when running the models locally, the GPU requirements for larger models is substantial and secondly while utilizing API services, the cost scales directly with the size of the evaluation dataset and number of models. These cost implications impose a limitation on the breadth and depth of our evaluation studies. We had to make strategic decisions about which models to include in our benchmark and how extensively to test them. This financial constraint introduces a selection bias on which models and tasks to prioritize which limits the scope of our evaluation

**3. Long-tail Distribution of Languages Across Tasks & Datasets:** Another limitation of AFROBENCH is the uneven distribution of languages across tasks and datasets. While our evaluation covers 64 languages in total, the coverage across tasks and datasets exhibits a long-tail distribution. As shown in Table C, 60% of the languages appear in fewer than 5 of the 21 datasets. This poses two challenges; first, it limits our ability to properly assess the performance of LLMs across these underrepresented languages. Secondly, it highlights the gap in the availability of evaluation datasets even among low-resource languages. Without extensive dataset coverage for these languages, conclusions about LLM capabilities across these languages remains tentative.

**4. Constraints in Machine Translation Metrics:** Machine translation is often evaluated using BLEU and ROUGE, which rely on word-level recall and precision, and chrF, which operates at the character level. Research has shown these metrics sometimes demonstrate poor correlation with human judgments of translation quality. Other evaluation metrics that utilize embedding similarity, such as BERTScore (Zhang\* et al., 2020) and COMET (Rei et al., 2020) / AfriCOMET (Wang et al., 2024), which leverage pretrained encoder models to generate scores by comparing translations against reference texts, are promising alternatives. However, these neural evaluation models have limited language coverage, making them unsuitable for many of the languages in our study. As a result, we rely on chrF++, which combines unigram and character n-gram overlap measurements. While this metric provides broader language coverage, it is a compromise between evaluation quality and practical applicability.

## References

- David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajudeen Gwadabe, et al. 2022a. [A few thousand translations go a long way! leveraging pre-trained models for African news translation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3053–3070, Seattle, United States. Association for Computational Linguistics.
- David Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba Alabi, Yanke Mao, Haonan Gao, and En-Shiun Lee. 2024a. [SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian’s, Malta. Association for Computational Linguistics.
- David Adelani, Graham Neubig, Sebastian Ruder, Shruti Rijhwani, Michael Beukman, Chester Palen-Michel, Constantine Lignos, Jesujoba Alabi, Shamsuddeen Muhammad, Peter Nabende, Cheikh M. Bamba Dione, et al. 2022b. [MasakhaNER 2.0: Africa-centric transfer learning for named entity recognition](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4488–4508, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D’souza, Julia Kreutzer, Constantine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, Stephen Mayhew, Israel Abebe Azime, Shamsuddeen H. Muhammad, Chris Chinenye Emezue, Joyce Nakatumba-Nabende, Perez Ogayo, Aremu Anuoluwapo, et al. 2021. [MasakhaNER: Named entity recognition for African languages](#). *Transactions of the Association for Computational Linguistics*, 9:1116–1131.
- David Ifeoluwa Adelani, Marek Masiak, Israel Abebe Azime, Jesujoba Alabi, Atnafu Lambebo Tonja, Christine Mwase, Odunayo Ogundepo, Bonaventure F. P. Dossou, Akintunde Oladipo, et al. 2023. [Masakhanews: News topic classification for african languages](#). *Preprint*, arXiv:2304.09972.
- David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe Azime, Zhuang Yun Jian, Jesujoba Oluwadara Alabi, Xuanli He, Millicent Ochieng, Sara Hooker, Andiswa Bukula, En-Shiun Annie Lee, Chiamaka Chukwuneke, Happy Buzaaba, et al. 2024b. [Irokobench: A new benchmark for african languages in the age of large language models](#). *ArXiv*, abs/2406.03368.
- Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. 2023a. [MEGA: Multilingual evaluation of generative AI](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4232–4267, Singapore. Association for Computational Linguistics.
- Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma, Ishaan Watts, Ashutosh Sathe, Millicent Ochieng, Rishav Hada, Prachi Jain, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. 2024. [MEGAVERSE: Benchmarking large language models across languages, modalities, models and tasks](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2598–2637, Mexico City, Mexico. Association for Computational Linguistics.
- Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma, Ishaan Watts, Ashutosh Sathe, Millicent Ochieng, Rishav Hada, Prachi Jain, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. 2023b. [Mega-verse: Benchmarking large language models across languages, modalities, models and tasks](#). *Preprint*, arXiv:2311.07463.
- Meta AI. 2024. [Meta ai announces llama 3.1](#). Accessed: Feb 1, 2025.
- Benjamin Akera, Jonathan Mukiibi, Lydia Sanyu Nagayi, Claire Babirye, Isaac Owomugisha, Solomon Nsumba, Joyce Nakatumba-Nabende, Engineer Bainomugisha, Ernest Mwebaze, and John Quinn. 2022. [Machine translation for african languages: Community creation of datasets and models in uganda](#). In *3rd Workshop on African Natural Language Processing*.
- Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. 2022. [Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4336–4349, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua

786	Howland, Andrea Hu, Jeffrey Hui, Jeremy Hur-	846
787	witz, Michael Isard, Abe Ittycheriah, Matthew Jagiel-	847
788	ski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun,	848
789	Sneha Kudugunta, Chang Lan, Katherine Lee, Ben-	
790	jamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li,	849
791	Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu,	850
792	Frederick Liu, Marcello Maggioni, Aroma Mahendru,	851
793	Joshua Maynez, Vedant Misra, Maysam Moussalem,	852
794	Zachary Nado, John Nham, Eric Ni, Andrew Nys-	853
795	strom, Alicia Parrish, Marie Pellat, Martin Polacek,	854
796	Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif,	855
797	Bryan Richter, Parker Riley, Alex Castro Ros, Au-	856
798	rko Roy, Brennan Saeta, Rajkumar Samuel, Renee	857
799	Shelby, Ambrose Slone, Daniel Smilkov, David R.	858
800	So, Daniel Sohn, Simon Tokumine, Dasha Valter,	859
801	Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang,	860
802	Pidong Wang, Zirui Wang, Tao Wang, John Wiet-	861
803	ing, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting	862
804	Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven	
805	Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav	863
806	Petrov, and Yonghui Wu. 2023. <a href="#">Palm 2 technical</a>	864
807	<a href="#">report</a> . <i>Preprint</i> , arXiv:2305.10403.	865
808	Anuoluwapo Aremu, Jesujoba Oluwadara Alabi, Daud	866
809	Abolade, Nkechinyere Faith Aguobi, Shamsud-	867
810	deen Hassan Muhammad, and David Ifeoluwa Ade-	868
811	lani. 2024. <a href="#">NaijaRC: A multi-choice reading com-</a>	869
812	<a href="#">prehension dataset for nigerian languages</a> . In <i>5th</i>	870
813	<i>Workshop on African Natural Language Processing</i> .	
814	Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel	871
815	Artetxe, Satya Narayan Shukla, Donald Husa, Naman	872
816	Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and	873
817	Madian Khabisa. 2024. <a href="#">The belebele benchmark: a</a>	874
818	<a href="#">parallel reading comprehension dataset in 122 lan-</a>	875
819	<a href="#">guage variants</a> . In <i>Proceedings of the 62nd Annual</i>	876
820	<i>Meeting of the Association for Computational Lin-</i>	
821	<i>guistics (Volume 1: Long Papers)</i> , pages 749–775,	877
822	Bangkok, Thailand. Association for Computational	878
823	Linguistics.	879
824	Edward Bayes, Israel Abebe Azime, Jesu-	880
825	joba Oluwadara Alabi, Jonas Kgomo, Tyna	881
826	Eloundou, Elizabeth Proehl, Kai Chen, Imaan	882
827	Khadir, Naome A. Etori, Shamsuddeen Hassan	883
828	Muhammad, Choice Mpanza, Igneciah Pocia Thete,	884
829	Dietrich Klakow, and David Ifeoluwa Adelani.	885
830	2024. <a href="#">Uhura: A benchmark for evaluating scientific</a>	886
831	<a href="#">question answering and truthfulness in low-resource</a>	887
832	<a href="#">african languages</a> . <i>ArXiv</i> , abs/2412.00948.	888
833	BIG bench authors. 2023. <a href="#">Beyond the imitation game:</a>	889
834	<a href="#">Quantifying and extrapolating the capabilities of lan-</a>	890
835	<a href="#">guage models</a> . <i>Transactions on Machine Learning</i>	891
836	<i>Research</i> .	892
837	Stella Biderman, Hailey Schoelkopf, Lintang Sutawika,	893
838	Leo Gao, Jonathan Tow, Baber Abbasi, Alham Fikri	894
839	Aji, Pawan Sasanka Ammanamanchi, Sidney Black,	895
840	Jordan Clive, Anthony DiPofi, Julien Etzaniz, Ben-	896
841	jamin Fattori, Jessica Zosa Forde, Charles Foster,	897
842	Jeffrey Hsu, Mimansa Jaiswal, Wilson Y. Lee, Hao-	898
843	nan Li, Charles Lovering, Niklas Muennighoff, Ellie	899
844	Pavlick, Jason Phang, Aviya Skowron, Samson Tan,	900
845	Xiangru Tang, Kevin A. Wang, Genta Indra Winata,	901
	François Yvon, and Andy Zou. 2024. <a href="#">Lessons from</a>	902
	<a href="#">the trenches on reproducible evaluation of language</a>	903
	<a href="#">models</a> .	904
	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	905
	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	
	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	
	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	
	Gretchen Krueger, Tom Henighan, Rewon Child,	
	Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens	
	Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-	
	teusz Litwin, Scott Gray, Benjamin Chess, Jack	
	Clark, Christopher Berner, Sam McCandlish, Alec	
	Radford, Ilya Sutskever, and Dario Amodei. 2020.	
	<a href="#">Language models are few-shot learners</a> . In <i>Ad-</i>	
	<i>vances in Neural Information Processing Systems</i> ,	
	volume 33, pages 1877–1901. Curran Associates,	
	Inc.	
	Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anas-	
	tasios N. Angelopoulos, Tianle Li, Dacheng Li,	
	Banghua Zhu, Hao Zhang, Michael I. Jordan,	
	Joseph E. Gonzalez, and Ion Stoica. 2024. Chat-	
	bot arena: an open platform for evaluating llms by	
	human preference. In <i>Proceedings of the 41st Inter-</i>	
	<i>national Conference on Machine Learning, ICML’24</i> .	
	JMLR.org.	
	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,	
	Maarten Bosma, Gaurav Mishra, Adam Roberts,	
	Paul Barham, Hyung Won Chung, Charles Sutton,	
	Sebastian Gehrmann, et al. 2022. Palm: Scaling	
	language modeling with pathways. <i>arXiv preprint</i>	
	<i>arXiv:2204.02311</i> .	
	Hyung Won Chung, Le Hou, Shayne Longpre, Barret	
	Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi	
	Wang, Mostafa Dehghani, Siddhartha Brahma, Al-	
	bert Webson, Shixiang Shane Gu, Zhuyun Dai,	
	Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-	
	ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,	
	Dasha Valter, Sharan Narang, Gaurav Mishra, Adams	
	Yu, Vincent Zhao, Yanping Huang, Andrew Dai,	
	Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-	
	cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,	
	and Jason Wei. 2022. <a href="#">Scaling instruction-finetuned</a>	
	<a href="#">language models</a> . <i>arXiv preprint</i> .	
	Cheikh M. Bamba Dione, David Ifeoluwa Adelani,	
	Peter Nabende, Jesujoba Alabi, Thapelo Sindane,	
	Happy Buzaaba, Shamsuddeen Hassan Muhammad,	
	Chris Chinenye Emezue, Perez Ogayo, Anuoluwapo	
	Aremu, Catherine Gitau, Derguene Mbaye, Jonathan	
	Mukiibi, Blessing Sibanda, Bonaventure F. P. Dos-	
	sou, Andiswa Bukula, Rooweither Mabuya, Allah-	
	sera Auguste Tapo, Edwin Munkoh-Buabeng, Vic-	
	toire Memdjokam Koagne, Fatoumata Ouoba Ka-	
	bore, Amelia Taylor, Godson Kalipe, Tebogo	
	Macucwa, Vukosi Marivate, Tajuddeen Gwadabe,	
	Mboning Tchiaze Elvis, Ikechukwu Onyenwe, Gra-	
	tien Atindogbe, Tolulope Adelani, Idris Akinade,	
	Olanrewaju Samuel, Marien Nahimana, Théogène	
	Musabeyezu, Emile Niyomutabazi, Ester Chimbenga,	
	Kudzai Gotosa, Patrick Mizha, Apelete Agbolo, Sey-	
	dou Traore, Chinedu Uchechukwu, Aliyu Yusuf,	



906	Muhammad Abdullahi, and Dietrich Klakow. 2023.	Jay Patel Sathy Rajasekharan Lyvia Lusiji Francesco	963
907	<a href="#">MasakhaPOS: Part-of-speech tagging for typolog-</a>	Piccino Mfoniso Ukwak Ellen Sebastian	964
908	<a href="#">ically diverse African languages</a> . In <i>Proceedings</i>	Jacaranda Health, Stanslaus Mwongela. 2024.	965
909	<i>of the 61st Annual Meeting of the Association for</i>	<a href="#">Afrollama v1: An instruction-tuned llama model for</a>	966
910	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	<a href="#">african languages</a> . Accessed: Feb 12, 2025.	967
911	pages 10883–10900, Toronto, Canada. Association		
912	for Computational Linguistics.		
913	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur	968
914	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	Mensch, Chris Bamford, Devendra Singh Chap-	969
915	Akhil Mathur, Alan Schelten, Amy Yang, Angela	lot, Diego de Las Casas, Florian Bressand, Gi-	970
916	Fan, Anirudh Goyal, Anthony S. Hartshorn, Aobo	anna Lengyel, Guillaume Lample, Lucile Saulnier,	971
917	Yang, Archi Mitra, Archie Sravankumar, Artem Ko-	L’elio Renard Lavaud, Marie-Anne Lachaux, Pierre	972
918	renev, and et al. 2024. <a href="#">The llama 3 herd of models</a> .	Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang,	973
919	<i>ArXiv</i> , abs/2407.21783.	Timothée Lacroix, and William El Sayed. 2023. <a href="#">Mis-</a>	974
		<a href="#">tral 7b</a> . <i>ArXiv</i> , abs/2310.06825.	975
920	Christian Federmann, Tom Kocmi, and Ying Xin. 2022.	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika	976
921	<a href="#">NTREX-128 – news test references for MT evalua-</a>	Bali, and Monojit Choudhury. 2020. <a href="#">The state and</a>	977
922	<a href="#">tion of 128 languages</a> . In <i>Proceedings of the First</i>	<a href="#">fate of linguistic diversity and inclusion in the NLP</a>	978
923	<i>Workshop on Scaling Up Multilingual Evaluation</i> ,	<a href="#">world</a> . In <i>Proceedings of the 58th Annual Meeting of</i>	979
924	pages 21–24, Online. Association for Computational	<i>the Association for Computational Linguistics</i> , pages	980
925	Linguistics.	6282–6293, Online. Association for Computational	981
		Linguistics.	982
926	Clémentine Fourrier, Nathan Habib, Hynek Kydlíček,	Sneha Kudugunta, Isaac Rayburn Caswell, Biao	983
927	Thomas Wolf, and Lewis Tunstall. 2023. <a href="#">Lighteval:</a>	Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati,	984
928	<a href="#">A lightweight framework for llm evaluation</a> .	Romi Stella, Ankur Bapna, and Orhan Firat. 2023.	985
		<a href="#">MADLAD-400: A multilingual and document-level</a>	986
929	Clémentine Fourrier, Nathan Habib, Alina Lozovskaya,	<a href="#">large audited dataset</a> . In <i>Thirty-seventh Conference</i>	987
930	Konrad Szafer, and Thomas Wolf. 2024. Open	<i>on Neural Information Processing Systems Datasets</i>	988
931	llm leaderboard v2. <a href="https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard">https://huggingface.</a>	<i>and Benchmarks Track</i> .	989
932	<a href="https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard">co/spaces/open-llm-leaderboard/open_llm_</a>		
933	<a href="https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard">leaderboard</a> .		
934	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman,	Viet Lai, Nghia Ngo, Amir Pouran Ben Veyseh, Hieu	990
935	Sid Black, Anthony DiPofi, Charles Foster, Laurence	Man, Franck Dernoncourt, Trung Bui, and Thien	991
936	Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li,	Nguyen. 2023. <a href="#">ChatGPT beyond English: Towards</a>	992
937	Kyle McDonell, Niklas Muennighoff, Chris Ociepa,	<a href="#">a comprehensive evaluation of large language mod-</a>	993
938	Jason Phang, Laria Reynolds, Hailey Schoelkopf,	<a href="#">els in multilingual learning</a> . In <i>Findings of the As-</i>	994
939	Aviya Skowron, Lintang Sutawika, Eric Tang, An-	<i>sociation for Computational Linguistics: EMNLP</i>	995
940	ish Thite, Ben Wang, Kevin Wang, and Andy Zou.	2023, pages 13171–13189, Singapore. Association	996
941	2024. <a href="#">A framework for few-shot language model</a>	for Computational Linguistics.	997
942	<a href="#">evaluation</a> .		
943	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	998
944	Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Kr-	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	999
945	ishnan, Marc’Aurelio Ranzato, Francisco Guzmán,	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku-	1000
946	and Angela Fan. 2022. <a href="#">The Flores-101 evaluation</a>	mar, Benjamin Newman, Binhang Yuan, Bobby Yan,	1001
947	<a href="#">benchmark for low-resource and multilingual ma-</a>	Ce Zhang, Christian Alexander Cosgrove, Christo-	1002
948	<a href="#">chine translation</a> . <i>Transactions of the Association for</i>	pher D Manning, Christopher Re, Diana Acosta-	1003
949	<i>Computational Linguistics</i> , 10:522–538.	Navas, Drew Arad Hudson, Eric Zelikman, Esin	1004
		Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren,	1005
950	Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Is-	Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel	1006
951	lam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang,	Orr, Lucia Zheng, Mert Yuksekogonul, Mirac Suzgun,	1007
952	M. Sohel Rahman, and Rifat Shahriyar. 2021. <a href="#">XL-</a>	Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar	1008
953	<a href="#">sum: Large-scale multilingual abstractive summariza-</a>	Khattab, Peter Henderson, Qian Huang, Ryan An-	1009
954	<a href="#">tion for 44 languages</a> . In <i>Findings of the Association</i>	drew Chi, Sang Michael Xie, Shibani Santurkar,	1010
955	<i>for Computational Linguistics: ACL-IJCNLP 2021</i> ,	Surya Ganguli, Tatsunori Hashimoto, Thomas Icard,	1011
956	pages 4693–4703, Online. Association for Computa-	Tianyi Zhang, Vishrav Chaudhary, William Wang,	1012
957	tional Linguistics.	Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Ko-	1013
		reeda. 2023. <a href="#">Holistic evaluation of language models</a> .	1014
958	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	<i>Transactions on Machine Learning Research</i> . Fea-	1015
959	Arora, Steven Basart, Eric Tang, Dawn Song, and	tured Certification, Expert Certification.	1016
960	Jacob Steinhardt. 2021. <a href="#">Measuring mathematical</a>		
961	<a href="#">problem solving with the math dataset</a> . <i>Preprint</i> ,	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	1017
962	<a href="#">arXiv:2103.03874</a> .	Wang, Shuohui Chen, Daniel Simig, Myle Ott, Na-	1018
		man Goyal, Shruti Bhosale, Jingfei Du, Ramakanth	1019
		Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav	1020



1021	Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona T. Diab, Veselin Stoyanov, and Xian Li. 2021. <a href="#">Few-shot learning with multilingual language models</a> . <i>CoRR</i> , abs/2112.10668.	1080
1022		1081
1023		
1024		
1025		
1026	Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. <a href="#">LLaMAX: Scaling linguistic horizons of LLM by enhancing translation capabilities beyond 100 languages</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10748–10772, Miami, Florida, USA. Association for Computational Linguistics.	
1027		
1028		
1029		
1030		
1031		
1032		
1033	Gemma Team Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, L. Sifre, Morgane Rivière, Mihir Kale, J Christopher Love, Pouya Dehghani Tafti, and et al. 2024. <a href="#">Gemma: Open models based on gemini research and technology</a> . <i>ArXiv</i> , abs/2403.08295.	
1034		
1035		
1036		
1037		
1038		
1039	Shamsuddeen Hassan Muhammad, Idris Abdulmumin, Abinew Ali Ayele, David Ifeoluwa Adelani, Ibrahim Said Ahmad, Saminu Mohammad Aliyu, Nelson Odhiambo Onyango, Lilian D. A. Wanzare, Samuel Rutunda, Lukman Jibril Aliyu, Esubalew Alemneh, Oumaima Hourrane, Hagos Tesfahun Gebremichael, Elyas Abdi Ismail, Meriem Beloucif, Ebrahim Chekol Jibril, Andiswa Bukula, Rooweither Mabuya, Salomey Osei, Abigail Opong, Tadesse Destaw Belay, Tadesse Kebede Guge, Tesfa Tegegne Asfaw, Chiamaka Ijeoma Chukwunke, Paul Röttger, Seid Muhie Yimam, and Nedjma Ousidhoum. 2025. <a href="#">Afrihate: A multilingual collection of hate speech and abusive language datasets for african languages</a> . <i>Preprint</i> , arXiv:2501.08284.	
1040		
1041		
1042		
1043		
1044		
1045		
1046		
1047		
1048		
1049		
1050		
1051		
1052		
1053		
1054	Shamsuddeen Hassan Muhammad, Idris Abdulmumin, Abinew Ali Ayele, Nedjma Ousidhoum, David Ifeoluwa Adelani, Seid Muhie Yimam, Ibrahim Sa’id Ahmad, Meriem Beloucif, Saif Mohammad, Sebastian Ruder, Oumaima Hourrane, Pavel Brazdil, Felermimo Dário Mário António Ali, Davis Davis, Salomey Osei, Bello Shehu Bello, Falalu Ibrahim, Tajuddeen Gwadabe, Samuel Rutunda, Tadesse Belay, Wendimu Baye Messelle, Hailu Beshada Balcha, Sisay A dugna Chala, Hagos Tesfahun Gebremichael, Bernard Opoku, and Steven Arthur. 2023. <a href="#">AfriSenti: A Twitter Sentiment Analysis Benchmark for African Languages</a> .	
1055		
1056		
1057		
1058		
1059		
1060		
1061		
1062		
1063		
1064		
1065		
1066		
1067	NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang.	
1068		
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		
	2022. No language left behind: Scaling human-centered machine translation.	1080
		1081
	Odunayo Ogundepo, Tajuddeen R. Gwadabe, Clara E. Rivera, Jonathan H. Clark, Sebastian Ruder, David Ifeoluwa Adelani, Bonaventure F. P. Dossou, Abdou Aziz DIOP, Claytone Sikasote, Gilles Hacheme, Happy Buzaaba, Ignatius Ezeani, et al. 2023. <a href="#">Afriqa: Cross-lingual open-retrieval question answering for african languages</a> . <i>Preprint</i> , arXiv:2305.06897.	1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
	Akintunde Oladipo, Mofetoluwa Adeyemi, Orevaoghene Ahia, Abraham Toluwalase Owodunni, Odunayo Ogundepo, David Ifeoluwa Adelani, and Jimmy Lin. 2023. <a href="#">Better quality pre-training data and t5 models for African languages</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 158–168, Singapore. Association for Computational Linguistics.	1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> . <i>ArXiv</i> , abs/2303.08774.	1098
		1099
	OpenAI. 2024. <a href="#">Gpt-4o system card</a> . Accessed: February 13, 2025.	1100
		1101
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . In <i>Advances in Neural Information Processing Systems</i> .	1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. <a href="#">COMET: A neural framework for MT evaluation</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2685–2702, Online. Association for Computational Linguistics.	1111
		1112
		1113
		1114
		1115
		1116
	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, and et al. 2024. <a href="#">Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context</a> . <i>ArXiv</i> , abs/2403.05530.	1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
	Gemma Team Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, L’eonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ram’e, Johan Ferret, Peter Liu, Pouya Dehghani Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, and et al. 2024. <a href="#">Gemma 2: Improving open language models at a practical size</a> . <i>ArXiv</i> , abs/2408.00118.	1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134

1135	Nathaniel Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. <a href="#">ChatGPT MT: Competitive for high- (but not low-) resource languages</a> . In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 392–418, Singapore. Association for Computational Linguistics.	1195
1136		1196
1137		
1138		
1139		
1140		
1141	Sebastian Ruder. 2021. Challenges and Opportunities in NLP Benchmarking. <a href="http://ruder.io/nlp-benchmarking">http://ruder.io/nlp-benchmarking</a> .	
1142		
1143		
1144	Sebastian Ruder, J. Clark, Alexander Gutkin, Mihir Kale, Min Ma, Massimo Nicosia, Shruti Rijhwani, Parker Riley, Jean Michel A. Sarr, Xinyi Wang, John Wieting, Nitish Gupta, Anna Katanova, Christo Kirov, Dana L. Dickinson, Brian Roark, Bidisha Samanta, Connie Tao, David Ifeoluwa Adelani, Vera Axelrod, Isaac Caswell, Colin Cherry, Dan Garrette, R. Reeve Ingle, Melvin Johnson, Dmitry Panteleev, and Partha Pratim Talukdar. 2023. <a href="#">Xtreme-up: A user-centric scarce-data benchmark for under-represented languages</a> . <i>ArXiv</i> , abs/2305.11938.	
1145		
1146		
1147		
1148		
1149		
1150		
1151		
1152		
1153		
1154		
1155	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2022. <a href="#">Language models are multilingual chain-of-thought reasoners</a> . <i>arXiv preprint</i> .	
1156		
1157		
1158		
1159		
1160	Iyanuoluwa Shode, David Ifeoluwa Adelani, Jing Peng, and Anna Feldman. 2023. <a href="#">NollySenti: Leveraging transfer learning and machine translation for Nigerian movie sentiment classification</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 986–998, Toronto, Canada. Association for Computational Linguistics.	
1161		
1162		
1163		
1164		
1165		
1166		
1167		
1168	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	
1169		
1170		
1171		
1172		
1173		
1174	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>ArXiv</i> , abs/2307.09288.	1195
1175		1196
1176		
1177		
1178		
1179		
1180		
1181		
1182		
1183		
1184		
1185		
1186		
1187		
1188		
1189		
1190		
1191		
1192		
1193		
1194		
	Jiayi Wang, David Ifeoluwa Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Tosin Adewumi, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgoh, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Hassan Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwuneke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Abdi Mohamed, Hassan Ayinde, Oluwabusayo Olufunke Awoyomi, Lama Alkhaled, Sana Al-azzawi, Naome A. Etori, Millicent Ochieng, Clemencia Siro, Njoroge Kiragu, Eric Muchiri, Wangari Kimotho, Lyse Naomi Wamba Momo, Daud Abolade, Simbiat Ajao, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Nasir Iro, Saheed S. Abdullahi, Stephen E. Moore, Bernard Opoku, Zainab Akinjobi, Abee Afolabi, Nnaemeka Obiefuna, Onyekachi Raphael Ogbu, Sam Ochieng', Verrah Akinyi Otiende, Chinedu Emmanuel Mbonu, Sakayo Toadoun Sari, Yao Lu, and Pontus Stenertorp. 2024. <a href="#">AfriMTE and AfriCOMET: Enhancing COMET to embrace under-resourced African languages</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 5997–6023, Mexico City, Mexico. Association for Computational Linguistics.	1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. <a href="#">mT5: A massively multilingual pre-trained text-to-text transformer</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.	1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
	Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. <a href="#">Bertscore: Evaluating text generation with bert</a> . In <i>International Conference on Learning Representations</i> .	1235
		1236
		1237
		1238
	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. <a href="#">Instruction-following evaluation for large language models</a> . <i>Preprint</i> , arXiv:2311.07911.	1239
		1240
		1241
		1242
	Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. <a href="#">Aya model: An instruction finetuned open-access multilingual language model</a> . <i>Preprint</i> , arXiv:2402.07827.	1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
	<b>A Task Based Results</b>	1251
	We group tasks using similar evaluation metrics to analyze model performance systematically.	1252
		1253

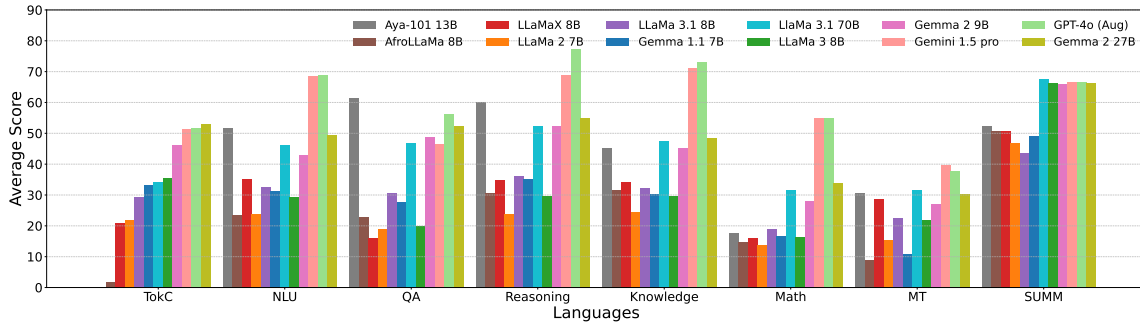


Figure 5: Performance of models across various NLP tasks, grouped by metric-based evaluation categories. Tasks include Token Classification (TokC), Natural Language Understanding (NLU), Question Answering (QA), Reasoning, Knowledge, Mathematics, Machine Translation (MT), and Summarization (SUMM).

## B LLMs evaluated

Models are selected to cover a range of open and closed-source LLMs with diverse parameter sizes, multilingual capabilities, and recent advancements. We prioritize models with strong multilingual support, accessibility for research, and relevance to African languages.

### B.0.1 Open Models

These are LLMs whose architectures, weights, and often training datasets are publicly available, allowing researchers and practitioners to fine-tune or adapt them to specific use cases. These models promote transparency, replicability, and accessibility, particularly for low-resource language tasks.

**Aya-101.** Aya-101 (Üstün et al., 2024) is a T5-style encoder-decoder model specifically fine-tuned for low-resource multilingual applications, including African languages. It was fine-tuned on a curated dataset, consisting of public multilingual corpora, and machine & human translated datasets from more than 100 languages. The model adopts a text-to-text paradigm and emphasizes cross-lingual transfer learning, allowing for robust generalization across various multilingual text-based tasks

**Llama 2 7B Chat.** Llama 2 (Touvron et al., 2023) is a collection of open-source pretrained and fine-tuned generative text models developed by Meta, ranging from 7 billion to 70 billion parameters. The 7B Chat variant allows for dialogue use cases. It employs an auto-regressive transformer architecture and has been fine-tuned using supervised fine-tuning (SFT) and reinforcement learning with human feedback (RLHF). They are pretrained on multiple languages, but has limited coverage of African languages.

**Llama 3 8B Instruct** Llama 3 (Dubey et al., 2024) is an updated variant of Llama 2 (Touvron

et al., 2023) series. They are instruction-fine-tuned to handle a wide range of text-based tasks. Similar to LLaMa 2, it also supports multiple languages but coverage of African languages remains limited. The number of parameters ranges from 8B to 70B; we make use of the 8B for this evaluation.

**Llama 3.1 Instruct (Bb, 70B)** Llama 3.1 (AI, 2024) is an updated variant of the Llama 3 series. Compared to Llama 3 (Dubey et al., 2024), Llama 3.1 (AI, 2024) introduces improvements in multilingual capabilities and general instruction-following. We use the instruction-tuned variants, fine-tuned for a broad range of NLP tasks. While it supports multiple languages, coverage of African languages remains limited. The model is available in parameter sizes ranging from 8B to 405B; due to computational cost, we evaluate only the 8B and 70B variants.

**Gemma 1.1 7B IT.** (Mesnard et al., 2024) is a lightweight open model from Google, built from the same research and technology used to create the Gemini models. They are text-to-text, decoder-only large language models, available in English, with open weights, pre-trained variants, and instruction-tuned variants. However, it does not have strong multilingual support. We evaluate the 7B instruction-finetuned variant of this model.

**Gemma 2 IT (9B, 27B).** Gemma 2 is an improved iteration of the Gemma model series optimized for efficiency. Compared to Gemma 1, Gemma 2 incorporates enhanced instruction-following capabilities and more robust parameter scaling. We evaluate the instruction-tuned variants of Gemma 2 at 9B and 27B parameter scales.

**AfroLlama-V1.** (Jacaranda Health, 2024) is a decoder-only transformer model, optimized for African language applications. It leverages pro-

proprietary datasets, including text from social media, newspapers, and government publications in African languages. Its architecture is based on Llama 3 8B (Dubey et al., 2024), but it incorporates additional pretraining on African-centric text.

### B.0.2 Proprietary Models

These are proprietary systems developed and maintained by organizations. Their training data and architectures are typically undisclosed.

**GPT-4o (Aug)** GPT-4o (OpenAI, 2024) is an optimized version of OpenAI’s GPT-4 model (OpenAI, 2023). It is an autoregressive omni model, trained end-to-end across text, vision, and audio on both public and proprietary data. While specific details about its architecture and datasets are not publicly disclosed, the GPT series is designed to adapt effectively to various language tasks, making it suitable for applications involving African languages. We evaluated the August 2024 version of this model

**Gemini 1.5 Pro 002.** Gemini (Reid et al., 2024) is a cutting-edge proprietary model with strong multilingual capacity. It is a compute-efficient multimodal model with training data tailored for diverse linguistic contexts, including low-resource languages. While specific details about its architecture and datasets are not publicly disclosed, Gemini is designed to adapt effectively to various language tasks, making it suitable for applications involving African languages.

## C Languages covered in the evaluation

Table 6 shows the languages and tasks we evaluated on.



	Language	Branch	Region (of Africa)	Script	# speakers
Afro-Asiatic	Algerian Arabic (arq)	Semitic	North	Arabic	36M
	Amharic (amh)	Ethio-Semitic	East	Ge'ez	57M
	Egyptian Arabic (arz)	Semitic	North	Arabic	41M
	Hausa (hau)	Chadic	West	Latin	77M
	Kabyle (kab)	Berber	North	Arabic	3M
	Oromo (orm)	Cushitic	East	Latin	37M
	Moroccan Arabic (ary)	Semitic	North	Arabic	29M
	Somali (som)	Cushitic	East	Latin	22M
	Tamasheq (taq)	Berber	East	Latin	1M
	Tamazight (tzm)	Berber	East	Latin	-
	Tigrinya (tig)	Ethio-Semitic	East	Ge'ez	9M
	Tunisian Arabic (aeb)	Semitic	North	Arabic	12M
Niger-Congo	Akan (aka)	Tano	West	Latin	10M
	Bambara (bam)	Mande	West	Latin	14M
	Bemba (bem)	Bantu	South, East & Central	Latin	4M
	Chichewa (nya)	Bantu	South-East	Latin	14M
	chiShona (sna)	Bantu	Southern	Latin	11M
	Chokwe (cjk)	Bantu	South & Central	Latin	1M
	Dyula (dyu)	Mande	West	Latin	3M
	Éwé (ewe)	Kwa	West	Latin	7M
	Fon (fon)	Volta-Niger	West	Latin	14M
	Ghomálá' (bbj)	Grassfields	Central	Latin	1M
	Igbo (ibo)	Volta-Niger	West	Latin	31M
	isiXhosa (xho)	Bantu	Southern	Latin	19M
	isiZulu (zul)	Bantu	Southern	Latin	27M
	Kabiyè (kbp)	Gur	West	Latin	1M
	Kamba (kam)	Bantu	East	Latin	5M
	Kikongo (kon)	Bantu	South & Central	Latin	5M
	Kikuyu (kik)	Bantu	East	Latin	8M
	Kimbundu (kmb)	Bantu	Southern	Latin	2M
	Kinyarwanda (kin)	Bantu	East	Latin	10M
	Kiswahili (swa)	Bantu	East & Central	Latin	71M-106M
	Lingala (lin)	Bantu	Central	Latin	40M
	Luba-Kasai (lua)	Bantu	Central	Latin	6M
	Luganda (lug)	Bantu	Central	Latin	11M
	Lugbara (lgg)				
	Mossi (mos)	Gur	West	Latin	8M
	Nigerian Fulfulde (fuv)	Senegambia	West	Latin	15M
	N'Ko (nqo)	Mande	West	Latin	-
	Northern Sotho (nso)	Bantu	Southern	Latin	4M
	Rundi (run)	Bantu	East	Latin	11M
	Runyankole (nyn)				
	Sango (sag)	Ubangian	Central	Latin	5M
	Setswana (tsn)	Bantu	Southern	Latin	14M
	Southern Sotho (sot)	Bantu	Southern	Latin	7M
	Swati (ssw)	Bantu	Southern	Latin	1M
	Twi (twi)	Kwa	West	Latin	9M
	Tumbuka (tum)	Bantu	South & East	Latin	2M

*Continued on next page*

	Language	Branch	Region (of Africa)	Script	# speakers
	Umbundu (umb)	Bantu	Southern	Latin	7M
	Xitsonga (tso)	Bantu	Southern	Latin	7M
	Wolof (wol)	Senegambia	West	Latin	5M
	Yoruba (yor)	Volta-Niger	West	Latin	46M
Nilo-Saharan	Acholi (ach)	Nilotic	East	Latin	1.5M
	Ateso (teo)	Nilotic	East	Latin	2.8M
	Dinka (dik)	Nilotic	Central	Latin	4M
	Kanuri (knc)	Saharan	West/Central	Latin	10M
	Kanuri (knc)	Saharan	West/Central	Arabic	10M
	Luo (luo)	Nilotic	East	Latin	4M
	Neur (nus)	Nilotic	Central	Latin	2M
Austronesian	Malagasy (plt)	Malayo-Polynesian	Southern	Latin	25M
Indo-European	Afrikaans (afr)	Germanic	Southern	Latin	7M
	Mozambican Portuguese (pt-MZ)	Italic	South East	Latin	13M
Creoles	Nigerian Pidgin (pcm)	English-based	West	Latin	121M
	Kabuverdianu (kea)	Portuguese-based	West	Latin	1M

Table 6: **Languages covered in each of our evaluation tasks:** language family, region, script, number of L1 & L2 speakers

Lang.	Classification								Reasoning	Question Answering						Generation						# Tasks
	AFriHATE	AFriSENTI	AFriXNLI	INJONGINTENT	NOLLYSENTI	MASAKHANEWS	MASAKHANER	MASAKHAPOS		AFriMMLU	AFriQA	BELEBELE	NAIJARC	OPENAI-MMLU	UHURA	AFriADR	FLORES	MAFAND	NTREX-128	SALT	XL-SUM	
aeb								✓								✓						2
ach																				✓		1
afr												✓				✓			✓			3
aka								✓								✓						2
amh	✓	✓	✓	✓		✓	✓	✓	✓	✓	✓				✓	✓	✓	✓				14
ara														✓							✓	2
arq	✓	✓																				2
ary	✓	✓										✓					✓					5
arz												✓					✓					3
bam							✓	✓	✓			✓					✓	✓				6
bbj							✓	✓	✓							✓		✓				4
bem									✓		✓						✓		✓			4
cjk									✓								✓					2
dik									✓								✓					2
dyu									✓								✓					2
ewe			✓	✓		✓	✓	✓	✓	✓							✓	✓	✓			10
fon								✓	✓		✓					✓	✓	✓				6
fuv									✓													1
gaz									✓								✓					2
hau	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓		✓	19
ibo	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			✓	✓	✓	✓	✓	✓	19
kab									✓								✓					2
kam									✓								✓					2
kbp									✓								✓					2
kea									✓								✓					2
kik									✓								✓					2
kin	✓	✓	✓	✓		✓	✓	✓	✓	✓	✓					✓	✓	✓	✓			13
kmb									✓								✓					2
knc									✓								✓					2
kon									✓								✓					2
lgg																				✓		1
lin			✓	✓		✓		✓	✓	✓		✓				✓						8
lua									✓								✓					2
lug			✓	✓		✓	✓	✓	✓	✓						✓	✓			✓		11
luo							✓	✓	✓							✓	✓	✓				5
mos							✓	✓	✓							✓		✓				5
nde																			✓			1
nso									✓							✓		✓				3
nus									✓							✓						2
nya						✓	✓	✓								✓	✓	✓				6

Continued on next page

Lang.	Classification								Reasoning	Question Answering						Generation					# Tasks		
	AFriHATE	AFriSENTI	AFriXNLI	INJONGINTENT	NOLLYSENTI	MASAKHANNEWS	MASAKHANER	MASAKHAPOS	SIB-200	AFRIMGSM	AFRIMMLU	AFRIQA	BELEBELE	NAIJARC	OPENAI-MMLU	UHURA	AFRIADR	FLORES	MAFAND	NTREX-128		SALT	XL-SUM
nyn																							1
orm	✓	✓	✓	✓		✓				✓	✓											✓	9
pcm	✓	✓				✓	✓	✓											✓			✓	7
plt									✓									✓		✓			3
run						✓			✓									✓					3
sag									✓									✓					2
sna			✓	✓		✓	✓	✓	✓	✓	✓	✓						✓	✓	✓			12
som	✓					✓			✓									✓		✓		✓	6
sot			✓	✓						✓	✓							✓					5
ssw									✓									✓		✓			3
swa	✓	✓	✓	✓		✓	✓	✓	✓	✓	✓	✓			✓	✓	✓	✓	✓	✓	✓	✓	18
taq																		✓					1
teo																					✓		1
tir	✓	✓				✓			✓				✓					✓		✓		✓	8
tsn							✓	✓										✓	✓	✓			5
tso		✓							✓				✓										3
tum									✓									✓					2
twi		✓	✓	✓		✓	✓	✓	✓	✓	✓	✓						✓	✓				11
tzm									✓									✓					2
umb									✓									✓					2
ven																				✓			1
wol			✓	✓			✓	✓	✓		✓		✓				✓	✓	✓	✓	✓		12
xho	✓		✓	✓		✓	✓	✓	✓	✓	✓	✓						✓	✓	✓			13
yor	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		✓	21
zul	✓		✓	✓		✓	✓	✓	✓	✓	✓	✓	✓		✓		✓	✓	✓	✓			14

Table 7: **Languages covered in each of our evaluation tasks:** check marks (✓) indicate that a language is covered by the task in that column. While 13 languages are covered by  $\geq 10$  tasks, 44 languages are covered by  $\leq 5$  tasks. SIB-200 and FLORES have the broadest coverage of African languages. In general, classification and generation tasks have better coverage of African languages than reasoning and question answering tasks.



## D Prompt Bank

In this section, we list all prompts used in our experiments. We use zero-shot cross-lingual prompts, where the context and query are in English, while the input text is in the target African language. This approach leverages LLMs’ stronger instruction-following in English (Lin et al., 2021; Shi et al., 2022). We display the prompts grouped by the task category shown in Figure 2.

### D.1 Natural Language Understanding

#### POS prompts:

##### Listing 1: MasakhaPOS Prompt 1

Please provide the POS tags for each word in the input sentence. The input will be a list of words in the sentence. The output format should be a list of tuples, where each tuple consists of a word from the input text and its corresponding POS tag label from the tag label set: ['ADJ', 'ADP', 'ADV', 'AUX', 'CCONJ', 'DET', 'INTJ', 'NOUN', 'NUM', 'PART', 'PRON', 'PROPN', 'PUNCT', 'SCONJ', 'SYM', 'VERB', 'X']. Your response should include only a list of tuples, in the order that the words appear in the input sentence, including punctuations, with each tuple containing the corresponding POS tag label for a word.

Sentence: {{text}}  
Output:

##### Listing 2: MasakhaPOS Prompt 2

You are an expert in tagging words and sentences in {{language}} with the right POS tag.

Please provide the POS tags for each word in the {{language}} sentence. The input is a list of words in the sentence. POS tag label set: ['ADJ', 'ADP', 'ADV', 'AUX', 'CCONJ', 'DET', 'INTJ', 'NOUN', 'NUM', 'PART', 'PRON', 'PROPN', 'PUNCT', 'SCONJ', 'SYM', 'VERB', 'X']. The output format should be a list of tuples, where each tuple consists of a word from the input text and its corresponding POS tag label from the POS tag label set provided.

Your response should include only a list of tuples, in the order that the words appear in the input sentence, including punctuations, with each tuple containing the corresponding POS tag label for a word.

Sentence: {{text}}  
Output:

##### Listing 3: MasakhaPOS Prompt 3

Acting as a {{language}} linguist and without making any corrections or changes to the text, perform a part of speech (POS) analysis of the sentences using the following POS tag label annotation ['ADJ', 'ADP', 'ADV', 'AUX', 'CCONJ', 'DET', 'INTJ', 'NOUN', 'NUM', 'PART', 'PRON', 'PROPN', 'PUNCT', 'SCONJ', 'SYM', 'VERB', 'X']. The input will be a list of words in the sentence. The output format should be a list of tuples, where each tuple consists of a word from the input text and its corresponding POS tag label from the POS tag label set provided.

Your response should include only a list of tuples, in the order that the words appear in the input sentence, including punctuations, with each tuple containing the corresponding POS tag label for a word.

Sentence: {{text}}  
Output:

##### Listing 4: MasakhaPOS Prompt 4

Annotate each word in the provided sentence with the appropriate POS tag. The annotation list is given as: ['ADJ', 'ADP', 'ADV', 'AUX', 'CCONJ', 'DET', 'INTJ', 'NOUN', 'NUM', 'PART', 'PRON', 'PROPN', 'PUNCT', 'SCONJ', 'SYM', 'VERB', 'X']. The input sentence will be a list of words in the sentence. The output format should be a list of tuples, where each tuple consists of a word from the input text and its corresponding POS tag label from the POS tag label set provided.

Your response should include only a list of tuples, in the order that the words appear in the input sentence, including punctuations, with each tuple containing the corresponding POS tag label for a word.

Sentence: {{text}}  
Output:

##### Listing 5: MMasakhaPOS Prompt 5

Given the following sentence, identify the part of speech (POS) for each word. Use the following POS tag set:

NOUN: Noun (person, place, thing),  
VERB: Verb (action, state),  
ADJ: Adjective (describes a noun),  
ADV: Adverb (modifies a verb, adjective, or adverb),  
PRON: Pronoun (replaces a noun),  
DET: Determiner (introduces a noun),  
ADP: Adposition (preposition or postposition),  
CCONJ: Conjunction (connects words, phrases, clauses)  
PUNCT: Punctuation,  
PROPN: Proper Noun,  
AUX: Auxiliary verb (helper verb), \nSCONJ: Subordinating conjunction  
PART: Particle,  
SYM: Symbol,  
INTJ: Interjection,  
NUM: Numeral,  
X: others. The output format should be a list of tuples, where each tuple consists of a word from the input text and its corresponding POS tag label key only from the POS tag set provided.

Your response should include only a list of tuples, in the order that the words appear in the input sentence, including punctuations, with each tuple containing the corresponding POS tag label for a word.

Sentence: {{text}}  
Output:

#### NER prompts:

##### Listing 1: MasakhaNER Prompt 1

Named entities refers to names of location, organisation and personal name.

For example, 'David is an employee of Amazon and he is visiting New York next week to see Esther' will be

PERSON: David \$ ORGANIZATION: Amazon \$ LOCATION: New York \$ PERSON: Esther

Ensure the output strictly follows the format: label : entity \$ label: entity, with each unique entity on a separate label line, avoiding grouped entities (e.g., avoid LOC: entity, entity) or irrelevant entries like none.

Text: {{text}}  
Return only the output

##### Listing 2: MasakhaNER Prompt 2

1498	You are working as a named entity recognition expert	1560
1499	and your task is to label a given text with	1561
1500	named entity labels. Your task is to identify	1562
1501	and label any named entities present in the	1563
1502	text. The named entity labels that you will be	1564
1503	using are PER (person), LOC (location), ORG (	1565
1504	organization) and DATE (date). Label multi-word	1566
1505	entities as a single named entity. For words	1567
1506	which are not part of any named entity, do not	1568
1507	return any value for it.	1569
1508	Ensure the output strictly follows the format: label	
1509	: entity \$\$ label: entity, with each unique	
1510	entity on a separate label line, avoiding	
1511	grouped entities (e.g., avoid LOC: entity,	
1512	entity) or irrelevant entries like none. Return	
1513	only the output	
1514		
1515	Text: {{text}}	

### Listing 3: MasakhaNER Prompt 3

1516	You are a Named Entity Recognition expert in {{	
1517	language}} language.	
1518	Extract all named entities from the following {{	
1519	language}} text and categorize them into PERSON	
1520	, LOCATION, ORGANIZATION, or DATE.	
1521	Ensure the output strictly follows the format; label	
1522	: entity \$\$ label: entity, with each unique	
1523	entity on a separate label line, avoiding	
1524	grouped entities (e.g., avoid LOC: entity,	
1525	entity) or irrelevant entries like none. Return	
1526	only the output	
1527		
1528	Text: {{text}}	
1529	Return only the output	

### Listing 4: MasakhaNER Prompt 4

1530	As a {{language}} linguist, label all named entities	
1531	in the {{language}} text below with the	
1532	categories: PERSON, LOCATION, ORGANIZATION, and	
1533	DATE. Ensure the output strictly follows the	
1534	format; label: entity \$\$ label: entity, with	
1535	each unique entity on a separate label line,	
1536	avoiding grouped entities (e.g., avoid LOC:	
1537	entity, entity) or irrelevant entries like none	
1538	. Return only the output.	
1539		
1540	Text: {{text}}	
1541	Return only the output	

### Listing 5: MasakhaNER Prompt 5

1542	Provide a concise list of named entities in the text	
1543	below. Use the following labels: PERSON,	
1544	LOCATION, ORGANIZATION, and DATE. Ensure the	
1545	output strictly follows the format; label:	
1546	entity \$\$ label: entity, with each unique	
1547	entity on a separate label line, avoiding	
1548	grouped entities (e.g., avoid LOC: entity,	
1549	entity) or irrelevant entries like none. Return	
1550	only the output.	
1551		
1552	Text: {{text}}	
1553	Return only the output	

## Sentiment prompts:

### Listing 1: AfriSenti Prompt 1

1555	Does this statement; "{{tweet}}" have a Neutral,	
1556	Positive or Negative sentiment? Labels only	

### Listing 2: AfriSenti Prompt 2

1557	Does this {{language}} statement; "{{tweet}}" have a	
1558	Neutral, Positive or Negative sentiment?	
1559	Labels only	

### Listing 3: AfriSenti Prompt 3

You are an assistant able to detect sentiments in	1560
tweets.	1561
Given the sentiment labels Neutral, Positive or	1562
Negative; what is the sentiment of the {{	1563
language}} statement below? Return only the	1564
labels.	1565
	1566
text: {{tweet}}	1567
label:	1568
	1569

### Listing 4: AfriSenti Prompt 4

Label the following text as Neutral, Positive, or	1570
Negative. Provide only the label as your	1571
response.	1572
	1573
text: {{tweet}}	1574
label:	1575

### Listing 5: AfriSenti Prompt 5

You are tasked with performing sentiment	1576
classification on the following {{language}}	1577
text. For each input, classify the sentiment as	1578
positive, negative, or neutral. Use the	1579
following guidelines:	1580
	1581
Positive: The text expresses happiness, satisfaction	1582
, or optimism.	1583
Negative: The text conveys disappointment,	1584
dissatisfaction, or pessimism.	1585
Neutral: The text is factual, objective, or without	1586
strong emotional undertones.	1587
If the text contains both positive and negative	1588
sentiments, choose the dominant sentiment. For	1589
ambiguous or unclear sentiments, select the	1590
label that best reflects the overall tone.	1591
Please provide a single classification for each	1592
input.	1593
	1594
text: {{tweet}}	1595
label:	1596

### Listing 6: NollySenti Prompt 1

Does this movie description "{{review}}" have a	1597
Positive or Negative sentiment? Labels only	1598

### Listing 7: NollySenti Prompt 2

Does this {{language}} movie description; "{{review	1599
}}" have a Positive or Negative sentiment?	1600
Labels only	1601

### Listing 8: NollySenti Prompt 3

You are an assistant able to detect sentiment in	1602
movie reviews.	1603
Given the sentiment labels Positive or Negative;	1604
what is the sentiment of the English statement	1605
below? Return only the labels	1606
	1607
Review: {{review}}"	1608
	1609

### Listing 9: NollySenti Prompt 4

Label the following text as Positive, or Negative.	1610
Provide only the label as your response.	1611
	1612
text: {{review}}	1613
label:	1614

### Listing 10: NollySenti Prompt 5

You are tasked with performing sentiment	1615
classification on the following English text.	1616
For each input, classify the sentiment as	1617
positive, negative. Use the following	1618
guidelines:	1619



1748	Given the text: '{{text}}', classify it into one of	The following text is in {{language}}: '{{text}}'.	1818
1749	these intents: [alarm, balance, bill_balance,	Given the list of intents: [alarm, balance,	1819
1750	book_flight, book_hotel, calendar_update,	bill_balance, book_flight, book_hotel,	1820
1751	cancel_reservation, car_rental,	calendar_update, cancel_reservation, car_rental	1821
1752	confirm_reservation, cook_time, exchange_rate,	, confirm_reservation, cook_time, exchange_rate	1822
1753	food_last, freeze_account, ingredients_list,	, food_last, freeze_account, ingredients_list,	1823
1754	interest_rate, international_visa, make_call,	interest_rate, international_visa, make_call,	1824
1755	meal_suggestion, min_payment, pay_bill,	meal_suggestion, min_payment, pay_bill,	1825
1756	pin_change, play_music, plug_type, recipe,	pin_change, play_music, plug_type, recipe,	1826
1757	restaurant_reservation, restaurant_reviews,	restaurant_reservation, restaurant_reviews,	1827
1758	restaurant_suggestion, share_location,	restaurant_suggestion, share_location,	1828
1759	shopping_list_update, spending_history, text,	shopping_list_update, spending_history, text,	1829
1760	time, timezone, transactions, transfer,	time, timezone, transactions, transfer,	1830
1761	translate, travel_notification,	translate, travel_notification,	1831
1762	travel_suggestion, update_playlist, weather].	travel_suggestion, update_playlist, weather],	1832
1763	Only output one intent from the list.	identify the intent expressed in the text.	1833
		Return only the identified intent.	1834

## Listing 2: IngongoIntent Prompt 2

1764	Analyze the text: '{{text}}'. Choose the most
1765	appropriate intent from these options: [alarm,
1766	balance, bill_balance, book_flight, book_hotel,
1767	calendar_update, cancel_reservation,
1768	car_rental, confirm_reservation, cook_time,
1769	exchange_rate, food_last, freeze_account,
1770	ingredients_list, interest_rate,
1771	international_visa, make_call, meal_suggestion,
1772	min_payment, pay_bill, pin_change, play_music,
1773	plug_type, recipe, restaurant_reservation,
1774	restaurant_reviews, restaurant_suggestion,
1775	share_location, shopping_list_update,
1776	spending_history, text, time, timezone,
1777	transactions, transfer, translate,
1778	travel_notification, travel_suggestion,
1779	update_playlist, weather]. Respond with only
1780	the selected intent.

## Listing 3: IngongoIntent Prompt 3

1781	You are a linguistic analyst trained to understand
1782	user intent. Based on the text: '{{text}}',
1783	choose the intent that best matches from this
1784	list: [alarm, balance, bill_balance,
1785	book_flight, book_hotel, calendar_update,
1786	cancel_reservation, car_rental,
1787	confirm_reservation, cook_time, exchange_rate,
1788	food_last, freeze_account, ingredients_list,
1789	interest_rate, international_visa, make_call,
1790	meal_suggestion, min_payment, pay_bill,
1791	pin_change, play_music, plug_type, recipe,
1792	restaurant_reservation, restaurant_reviews,
1793	restaurant_suggestion, share_location,
1794	shopping_list_update, spending_history, text,
1795	time, timezone, transactions, transfer,
1796	translate, travel_notification,
1797	travel_suggestion, update_playlist, weather].
1798	Return only the intent.

## Listing 4: IngongoIntent Prompt 4

1799	You are a English linguistic analyst trained to
1800	understand {{language}} user intent. Based on
1801	the {{language}} text: '{{text}}', choose the
1802	intent that best matches from this list: [alarm
1803	, balance, bill_balance, book_flight,
1804	book_hotel, calendar_update, cancel_reservation
1805	, car_rental, confirm_reservation, cook_time,
1806	exchange_rate, food_last, freeze_account,
1807	ingredients_list, interest_rate,
1808	international_visa, make_call, meal_suggestion,
1809	min_payment, pay_bill, pin_change, play_music,
1810	plug_type, recipe, restaurant_reservation,
1811	restaurant_reviews, restaurant_suggestion,
1812	share_location, shopping_list_update,
1813	spending_history, text, time, timezone,
1814	transactions, transfer, translate,
1815	travel_notification, travel_suggestion,
1816	update_playlist, weather]. Return only the
1817	intent.

## Listing 5: IngongoIntent Prompt 5

## Hate Speech prompts:

### Listing 1: AfriHate Prompt 1

I am providing you with the definition Hate speech,	1836
Abusive language and Normal tweets.	1837
Hate speech is a language content that expresses	1838
hatred towards a particular group or individual	1839
based on their political affiliation, race,	1840
ethnicity, religion, gender, sexual orientation	1841
, or other characteristics. It also includes	1842
threats of violence	1843
Abusive language is any form of bad language	1844
expressions including rude, impolite, insulting	1845
or belittling utterance intended to offend or	1846
harm an individual.	1847
Normal does not contain any bad language.	1848
	1849
Tweet: {{tweet}}	1850
	1851
Which category does the tweet above belong to: 'Hate	1852
', 'Abuse' or 'Normal'. Pick exactly one	1853
category. Return only the label	1854

### Listing 2: AfriHate Prompt 2

Read the following label definitions and provide a	1855
label without any explanations.	1856
	1857
Hate: Hate speech is public speech that expresses	1858
hate or encourages violence towards a person or	1859
group based on something such as race,	1860
religion, gender, ethnicity, sexual orientation	1861
or other characteristics.	1862
	1863
Abusive: Abusive and offensive language means verbal	1864
messages that use words in an inappropriate	1865
way and may include but is not limited to	1866
swearing, name-calling, or profanity. Offensive	1867
language may upset or embarrass people because	1868
it is rude or insulting.	1869
	1870
Normal: Normal language is neither hateful nor	1871
abusive or offensive. It does not contain any	1872
bad language.	1873
	1874
Text: {{tweet}}	1875
Label:	1876

### Listing 3: AfriHate Prompt 3

Read the following text and definitions:	1877
	1878
Text: {{tweet}}.	1879
	1880
Definitions:	1881
Hate: Hate speech is public speech that expresses	1882
hate or encourages violence towards a person or	1883
group based on something such as race,	1884
religion, gender, ethnicity, sexual orientation	1885
or other characteristics.	1886
	1887
Abuse: Abusive and offensive language means verbal	1888
messages that use words in an inappropriate way	1889
and may include but is not limited to swearing	1890
, name-calling, or profanity. Offensive	1891
language may upset or embarrass people because	1892
it is rude or insulting.	1893

1894  
1895  
1896  
1897  
1898  
1899  
1900

Normal: Normal language is neither hateful nor abusive or offensive. It does not contain any bad language.

Which of these definitions (hate, abuse, normal) apply to this tweet?, return only the label

#### Listing 4: AfriHate Prompt 4

1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925

Read the following definitions and text to categorize:

Definitions:

Hate: Hate speech is public speech that expresses hate or encourages violence towards a person or group based on something such as race, religion, gender, ethnicity, sexual orientation or other characteristics.

Abuse: Abusive and offensive language means verbal messages that use words in an inappropriate way and may include but is not limited to swearing, name-calling, or profanity. Offensive language may upset or embarrass people because it is rude or insulting.

Normal: Normal language is neither hateful nor abusive or offensive. It does not contain any bad language.

Text: {{tweet}}.

Which of these definitions (hate, abuse, normal) apply to this tweet? Return only the label

#### Listing 5: AfriHate Prompt 5

1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959

You will be given a text snippet and 3 category definitions.

Your task is to choose which category applies to this text.

Your text snippet is: {{tweet}}.

Your category definitions are:

HATE category definition: Hate speech is public speech that expresses hate or encourages violence towards a person or group based on something such as race, religion, gender, ethnicity, sexual orientation or other characteristics.

ABUSE category definition: Abusive and offensive language means verbal messages that use words in an inappropriate way and may include but is not limited to swearing, name-calling, or profanity. Offensive language may upset or embarrass people because it is rude or insulting.

NORMAL category definition: Normal language is neither hateful nor abusive or offensive. It does not contain any bad language.

Does the text snippet belong to the HATE, ABUSIVE, or the NORMAL category? Thinking step by step answer HATE, ABUSIVE, or NORMAL capitalizing all the letters.

Explain your reasoning FIRST, then output HATE, ABUSIVE, or NORMAL. Clearly return the label in capital letters.

### Natural Language Inference prompts:

#### Listing 1: AfriXNLI Prompt 1

1961  
1962  
1963  
1964  
1965  
1966  
1967

Please identify whether the premise entails or contradicts the hypothesis in the following premise and hypothesis. The answer should be exact entailment, contradiction, or neutral.

Premise: {{premise}}

Hypothesis: {{hypothesis}}.

Is it entailment, contradiction, or neutral?

#### Listing 2: AfriXNLI Prompt 2

{{premise}}

Question: {{hypothesis}} True, False, or Neither?

Answer:

#### Listing 3: AfriXNLI Prompt 3

Given the following premise and hypothesis in {{language}}, identify if the premise entails, contradicts, or is neutral towards the hypothesis. Please respond with exact 'entailment', 'contradiction', or 'neutral'.

Premise: {{premise}}

Hypothesis: {{hypothesis}}

#### Listing 4: AfriXNLI Prompt 4

You are an expert in Natural Language Inference (NLI) specializing in {{language}} language. Analyze the premise and hypothesis given in {{language}}, and determine the relationship between them.

Respond with one of the following options: 'entailment', 'contradiction', or 'neutral'.

Premise: {{premise}}

Hypothesis: {{hypothesis}}

#### Listing 5: AfriXNLI Prompt 5

Based on the given statement, is the following claim 'true', 'false', or 'inconclusive'.

Statement: {{premise}}

Claim: {{hypothesis}}

### D.2 Question Answering

#### CrosslingualQA prompts:

#### Listing 1: AfriQA Prompt 1

Your task is to answer a question given a context. Make sure you respond with the shortest span containing the answer in the context.

Question: {{question\_lang}}

Context: {{context}}

Answer:

#### Listing 2: AfriQA Prompt 2

Your task is to answer a question given a context. The question is in {{language}}, while the context is in English or French.

Make sure you respond with the shortest span in the context that contains the answer.

Question: {{question\_lang}}

Context: {{context}}

Answer:

#### Listing 3: AfriQA Prompt 3

Given the context, provide the answer to the following question.

Ensure your response is concise and directly from the context.

Question: {{question\_lang}}

Context: {{context}}

Answer:



2019  
2020  
2021  
2022  
2023  
2024  
2025

Listing 4: AfriQA Prompt 4

You are an AI assistant and your task is to answer the question based on the provided context. Your answer should be the shortest span that contains the answer within the context.  
Question: {{question\_lang}}  
Context: {{context}}  
Answer:

2026  
2027  
2028  
2029  
2030  
2031

Listing 5: AfriQA Prompt 5

Using the context, find the answer to the question. Respond with the briefest span that includes the answer from the context.  
Question: {{question\_lang}}  
Context: {{context}}  
Answer:

2032

Reading Comprehension prompts:

2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040

Listing 1: Belebele Prompt 1

P: {{passage}}  
Q: {{question}}  
A: {{option\_1}}  
B: {{option\_2}}  
C: {{option\_3}}  
D: {{option\_4}}  
Please choose the correct answer from the options above:

2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048

Listing 2: Belebele Prompt 2

Passage: {{passage}}  
Question: {{question}}  
1: {{option\_1}}  
2: {{option\_2}}  
3: {{option\_3}}  
4: {{option\_4}}  
Please select the correct answer from the given choices

2049  
2050  
2051  
2052  
2053  
2054  
2055  
2056

Listing 3: Belebele Prompt 3

Context: {{passage}}  
Query: {{question}}  
Option A: {{option\_1}}  
Option B: {{option\_2}}  
Option C: {{option\_3}}  
Option D: {{option\_4}}  
Please indicate the correct option from the list above:

2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067

Listing 4: Belebele Prompt 4

{{passage}}  
Based on the above passage, answer the following question:  
{{question}}  
Choices:  
A) {{option\_1}}  
B) {{option\_2}}  
C) {{option\_3}}  
D) {{option\_4}}  
Please provide the correct answer from the choices given

2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075

Listing 5: Belebele Prompt 5

Read the passage: {{passage}}  
Then answer the question: {{question}}  
Options:  
A. {{option\_1}}  
B. {{option\_2}}  
C. {{option\_3}}  
D. {{option\_4}}  
Please choose the correct option from the above list

P: {{story}}  
Q: {{question}}  
A: {{options\_A}}  
B: {{options\_B}}  
C: {{options\_C}}  
D: {{options\_D}}  
Please choose the correct answer from the options above

Listing 6: NaijaRC Prompt 1  
Listing 7: NaijaRC Prompt 2  
Passage: {{story}}  
Question: {{question}}  
1: {{options\_A}}  
2: {{options\_B}}  
3: {{options\_C}}  
4: {{options\_D}}  
Please select the correct answer from the given choices

Listing 8: NaijaRC Prompt 3  
Context: {{story}}  
Query: {{question}}  
Option A: {{options\_A}}  
Option B: {{options\_B}}  
Option C: {{options\_C}}  
Option D: {{options\_D}}  
Please indicate the correct option from the list above

Listing 9: NaijaRC Prompt 4  
{{story}}  
Based on the above passage, answer the following question  
{{question}}  
Choices:  
A) {{options\_A}}  
B) {{options\_B}}  
C) {{options\_C}}  
D) {{options\_D}}  
Please provide the correct answer from the choices given

Listing 10: NaijaRC Prompt 5  
Read the passage: {{story}}  
Then answer the question: {{question}}  
Options:  
A. {{options\_A}}  
B. {{options\_B}}  
C. {{options\_C}}  
D. {{options\_D}}  
Please choose the correct option from the above list

D.3 Knowledge  
Arc-E prompts:

Listing 1: UHURA Prompt 1  
You are a virtual assistant that answers multiple-choice questions with the correct option only.  
Question: {{question}}  
Choices:  
A. {{options\_A}}  
B. {{options\_B}}  
C. {{options\_C}}  
D. {{options\_D}}  
Answer:

Listing 2: UHURA Prompt 2  
Choose the correct option that answers the question below:

2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142

```
Question: {{question}}

Choices:
A. {{options_A}}
B. {{options_B}}
C. {{options_C}}
D. {{options_D}}
Answer: .
```

Listing 3: UHURA Prompt 3

2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153

```
Answer the following multiple-choice question by
picking 'A', 'B', 'C', or 'D'

Question: {{question}}

Options:
A. {{options_A}}
B. {{options_B}}
C. {{options_C}}
D. {{options_D}}
Answer:
```

Listing 4: UHURA Prompt 4

2154  
2155  
2156  
2157  
2158  
2159  
2160  
2161

```
Question: {{question}}

Options:
A. {{options_A}}
B. {{options_B}}
C. {{options_C}}
D. {{options_D}}
Answer:
```

Listing 5: UHURA Prompt 5

2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170

```
Which of the following options answers this question
: {{question}}

Options:
A. {{options_A}}
B. {{options_B}}
C. {{options_C}}
D. {{options_D}}
Answer:
```

2171

**MMLU prompts:**

Listing 1: OpenAIMMLU Prompt 1

2172  
2173  
2174  
2175  
2176  
2177  
2178

```
Q: {{Question}}
A: {{A}}
B: {{B}}
C: {{C}}
D: {{D}}
Please choose the correct answer from the options
above
```

Listing 2: OpenAIMMLU Prompt 2

2179  
2180  
2181  
2182  
2183  
2184  
2185

```
Question: {{Question}}
1: {{A}}
2: {{B}}
3: {{C}}
4: {{D}}
Please select the correct answer from the given
choices
```

Listing 3: OpenAIMMLU Prompt 3

2186  
2187  
2188  
2189  
2190  
2191  
2192

```
Input Question: {{Question}}
Option A: {{A}}
Option B: {{B}}
Option C: {{C}}
Option D: {{D}}
Please indicate the correct option from the list
above
```

Listing 4: OpenAIMMLU Prompt 4

```
Critically analyze the question and select the most
probable answer from the list:
{{Question}}
Choices:
A) {{A}}
B) {{B}}
C) {{C}}
D) {{D}}
```

2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200

Listing 5: OpenAIMMLU Prompt 5

```
Answer the question and pick the correct answer from
the options:
{{Question}}
Options:
A. {{A}}
B. {{B}}
C. {{C}}
D. {{D}}
Please choose the correct option from the above list
```

2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209

Listing 6: AfriMMLU Prompt 1

```
You are a highly knowledgeable and intelligent
artificial intelligence model answers multiple-
choice questions about {{subject}}.

Question: {{question}}
Choices:
A: {{options_A}}
B: {{options_B}}
C: {{options_C}}
D: {{options_D}}

Answer:
```

2210  
2211  
2212  
2213  
2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221

Listing 7: AfriMMLU Prompt 2

```
As an expert in {{subject}}, choose the most
accurate answer to the question below. Your
goal is to select the correct option 'A', 'B',
'C', or 'D' by understanding the nuances of the
topic.

Question: {{question}}
Choices:
A: {{options_A}}
B: {{options_B}}
C: {{options_C}}
D: {{options_D}}

Answer:
```

2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235

Listing 8: AfriMMLU Prompt 3

```
You are a subject matter expert in {{subject}}.
Utilizing your expertise in {{subject}}, answer
the following multiple-choice question by
picking 'A', 'B', 'C', or 'D'.

Question: {{question}}
Choices:
A: {{options_A}}
B: {{options_B}}
C: {{options_C}}
D: {{options_D}}

Answer:
```

2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248

Listing 9: AfriMMLU Prompt 4

```
Analyze each question critically and determine the
most correct option based on your understanding
of the subject matter
Question: {{question}}
Choices:
A: {{options_A}}
B: {{options_B}}
C: {{options_C}}
D: {{options_D}}

Answer:
```

2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259

#### Listing 10: AfriMMLU Prompt 5

Given your proficiency in {{subject}}, please answer the subsequent multiple-choice question  
Question: {{question}}  
Choices:  
A: {{options\_A}}  
B: {{options\_B}}  
C: {{options\_C}}  
D: {{options\_D}}  
Answer:

### D.4 Reasoning

#### Math prompts: from IROKOBENCH (Adelani et al., 2024b)

##### Listing 1: AfriMGSM Prompt 1

{{question}}  
Step-by-step Answer:

##### Listing 2: AfriMGSM Prompt 2

Give direct numerical answers for the question provided.  
Question: {{question}}  
Step-by-step Answer:

##### Listing 3: AfriMGSM Prompt 3

Solve the following math question  
Question: {{question}}  
Step-by-step Answer:

##### Listing 4: AfriMGSM Prompt 4

Answer the given question with the appropriate numerical value, ensuring that the response is clear and without any supplementary information.  
Question: {{question}}  
Step-by-step Answer:

##### Listing 5: AfriMGSM Prompt 5

For mathematical questions provided in {{language}} language. Supply the accurate numeric step by step answer to the provided question.  
Question: {{question}}  
Step-by-step Answer:

### D.5 Text Generation

#### Machine Translation prompts

##### Listing 1: Machine Translation Prompt 1

{{source\_lang}} sentence: {{source\_text}}  
{{target\_lang}} sentence:

##### Listing 2: Machine Translation Prompt 2

You are a translation expert. Translate the following {{source\_lang}} sentences to {{target\_lang}}  
{{source\_lang}} sentence: {{source\_text}}  
{{target\_lang}} sentence:

#### Listing 3: Machine Translation Prompt 3

As a {{source\_lang}} and {{target\_lang}} linguist, translate the following {{source\_lang}} sentences to {{target\_lang}}.  
{{source\_lang}} sentence: {{source\_text}}  
{{target\_lang}} sentence:

#### Summarization prompts

##### Listing 1: XL-SUM Prompt 1

Provide a summary of the document written in {{language}}. Ensure that you provide the summary in {{language}} and nothing else.  
Document in {{language}}: {{text}}  
Summary:

##### Listing 2: XL-SUM Prompt 2

Summarize the document below in triple backticks and return only the summary and nothing else.  
{{text}}

##### Listing 3: XL-SUM Prompt 3

You are an advanced Summarizer, a specialized assistant designed to summarize documents in {{language}}. Your main goal is to ensure summaries are concise and informative. Ensure you return the summary only and nothing else.  
Document: {{text}}  
Summary:

#### Diacritics Restoration prompts

##### Listing 1: AFRIADR Prompt 1

Please restore the missing diacritics in the following sentence: {{text}}.  
Return output sentence only

##### Listing 2: AFRIADR Prompt 2

Given a sentence without diacritics, add the appropriate diacritics to make it grammatically and semantically correct.  
Sentence: {{text}}.  
Return output sentence only

##### Listing 3: AFRIADR Prompt 3

This text is in {{language}}. Restore all diacritical marks to their proper places in the following sentence: {{text}}. Return output sentence only

##### Listing 4: AFRIADR Prompt 4

You are a linguist specializing in diacritical marks for {{language}}. Add the appropriate diacritics to this {{language}} sentence: {{text}}. Return output sentence only

##### Listing 5: AFRIADR Prompt 5

You are a linguist specializing in diacritical marks for {{language}}. Diacritics are essential for proper pronunciation and meaning in {{language}}. You are tasked with converting {{language}} sentences without diacritics into their correctly accented forms. Here's the input: {{text}}. Return output sentence only

## E Detailed Results Per Language

This appendix presents detailed per-language performance results for each dataset. We group them by the task category shown in Figure 2. Each figure shows the model performance on the best prompt per language.

### E.1 Natural Language Understanding (NLU)

#### E.1.1 POS

##### MasakhaPOS

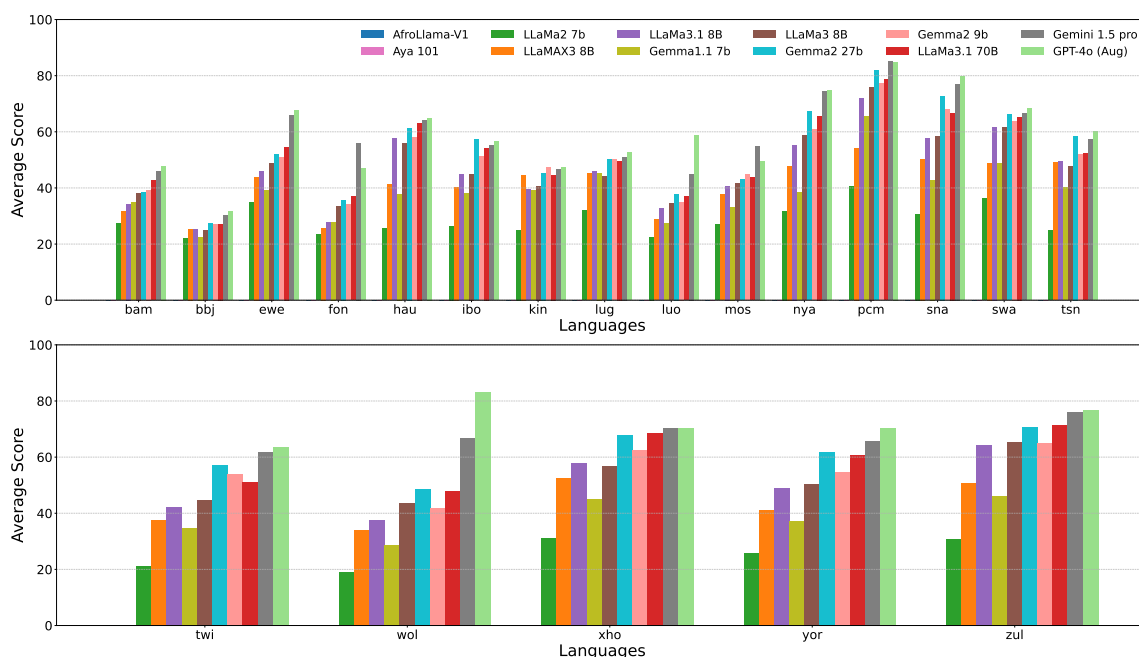


Figure 6: Per-language performance results for the MasakhaPOS dataset.

#### E.1.2 NER

##### MasakhaNER

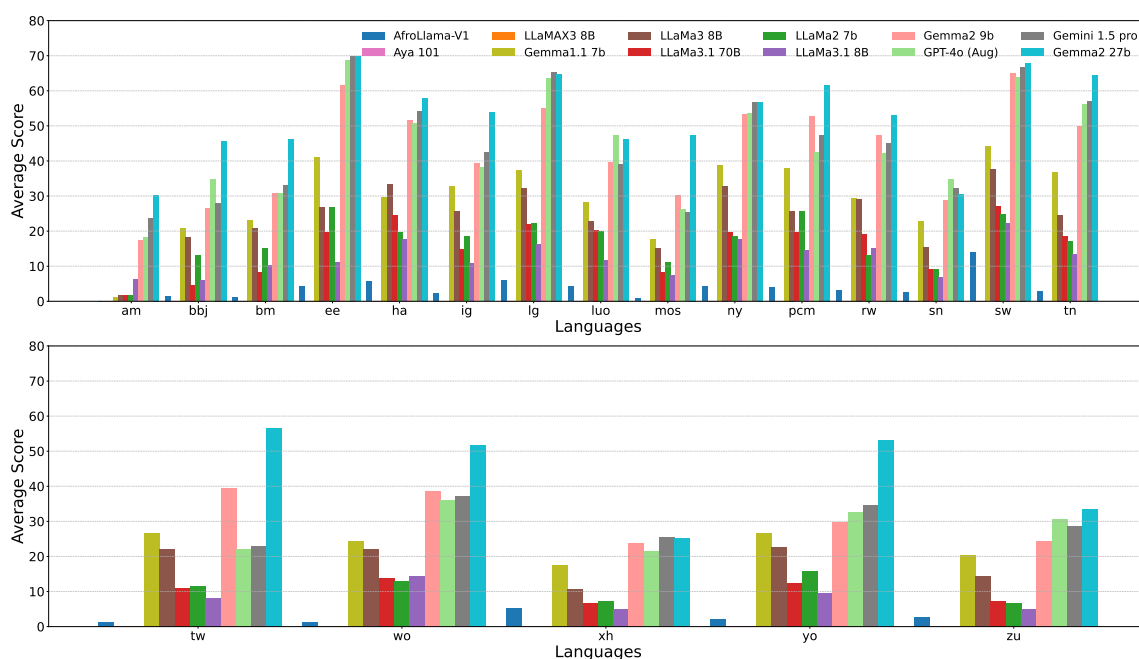


Figure 7: Per-language performance results for the MasakhaNER dataset.

E.1.3 Sentiment Analysis  
AfriSenti

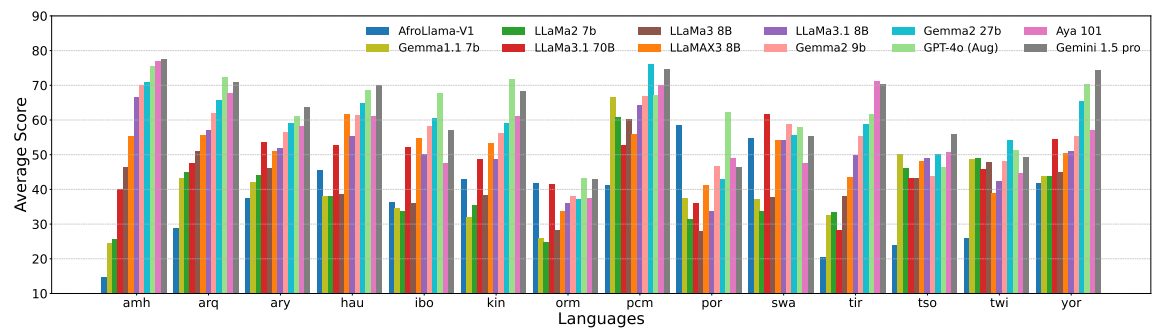


Figure 8: Per-language performance results for the AfriSenti dataset.



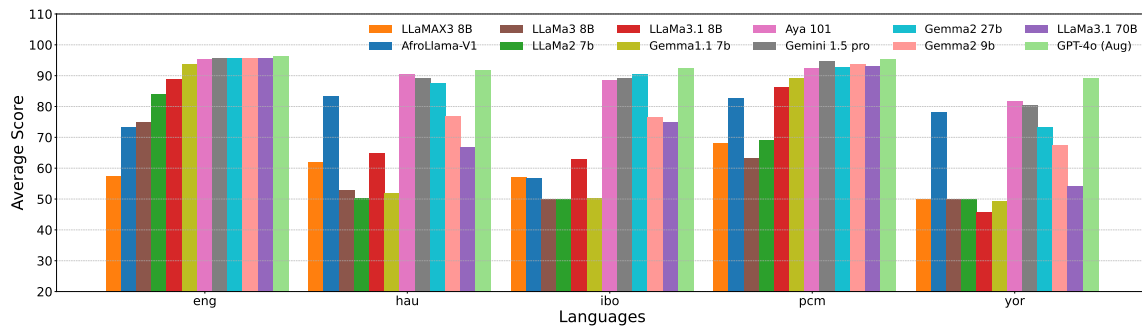


Figure 9: Per-language performance results for the NollySenti dataset.

### E.1.4 Intent Detection

#### Injongo Intent

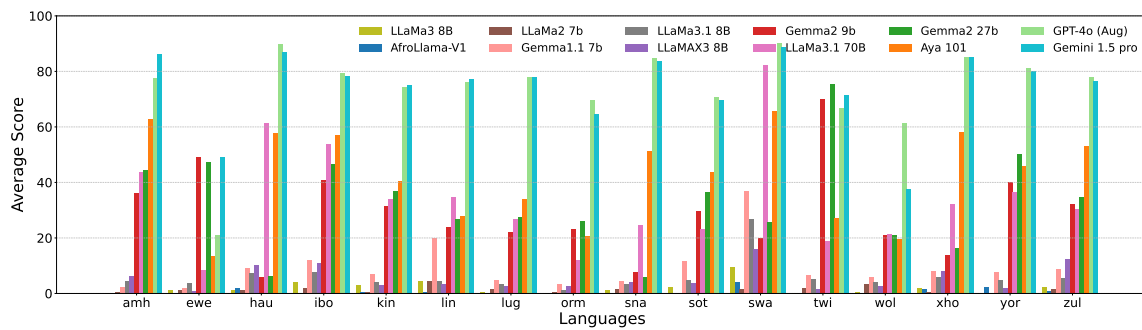


Figure 10: Per-language performance results for the InjongoIntent dataset.

### E.1.5 Topic Classification

#### MasakhaNEWS

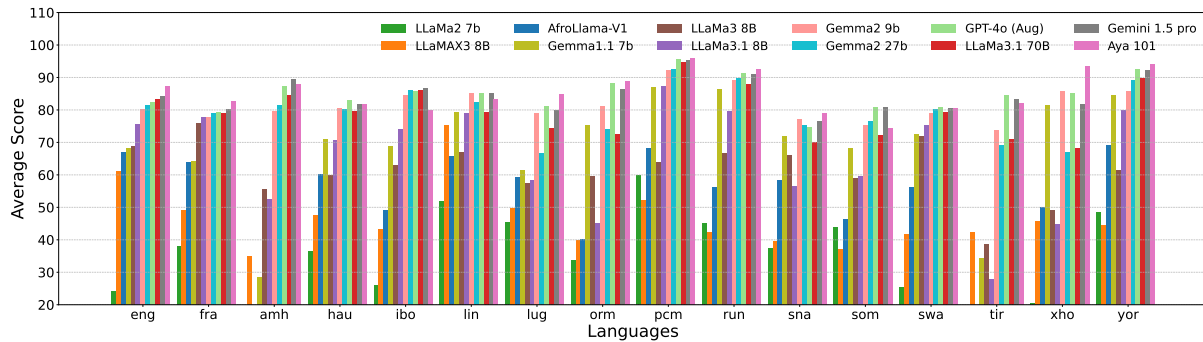


Figure 11: Per-language performance results for the MasakhaNEWS dataset.

## SIB

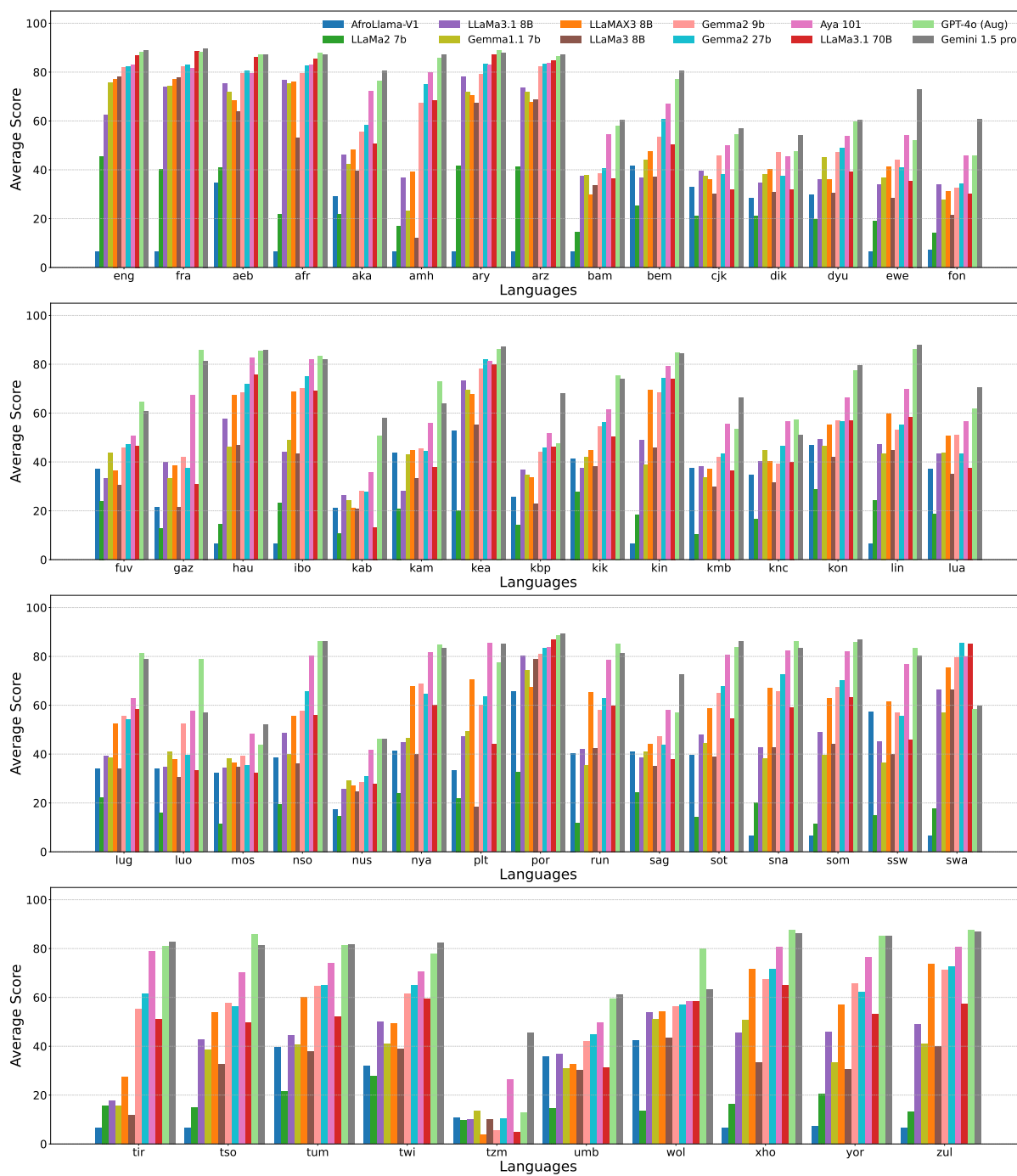


Figure 12: Per-language performance results for the SIB dataset.

## E.1.6 Hate Speech:

### AfriHate

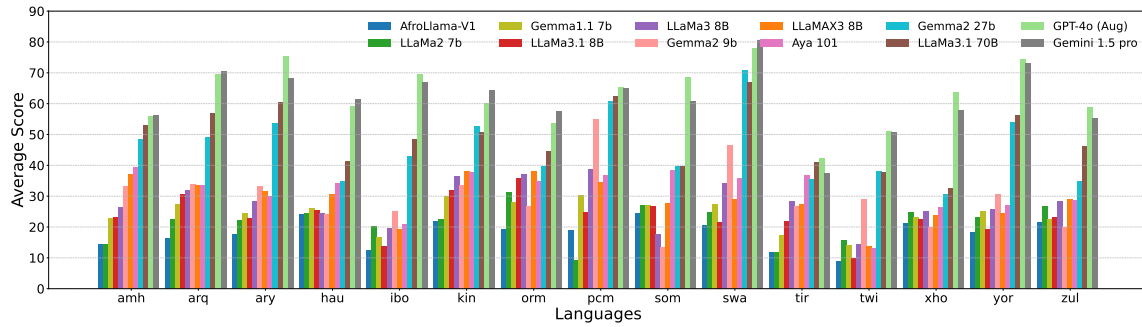


Figure 13: Per-language performance results for the AfriHate dataset.

## E.2 Natural Language Inference

### AfriXNLI

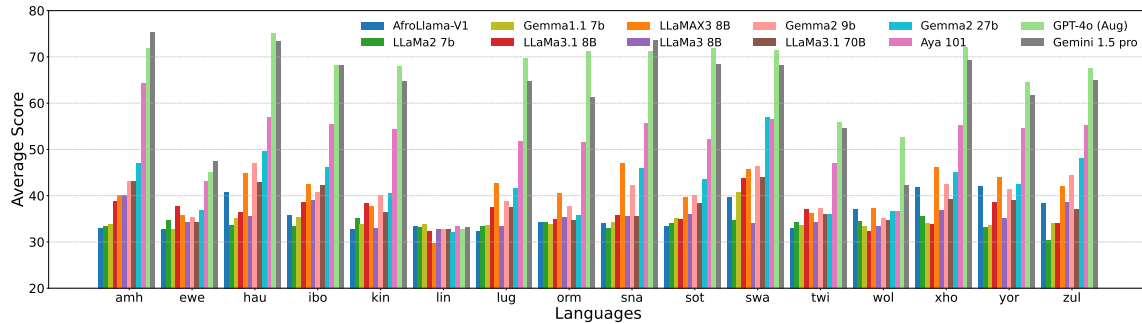


Figure 14: Per-language performance results for the AFRIXNLI dataset.

## E.3 Question Answering

### E.3.1 Cross-lingual Question Answering

#### AfriQA

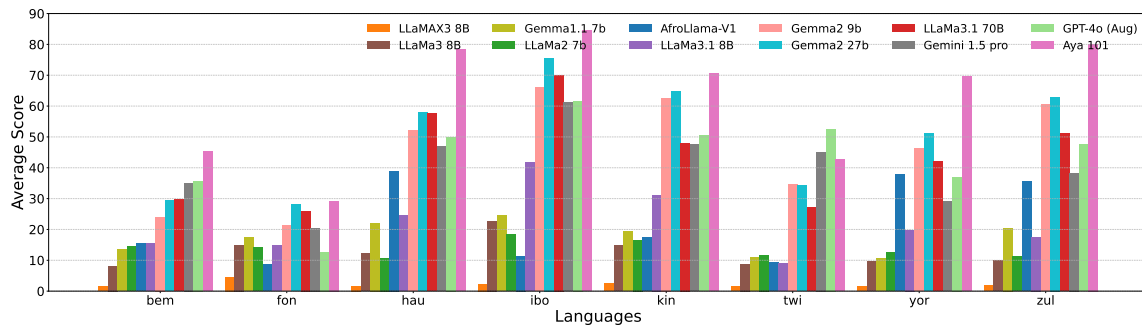
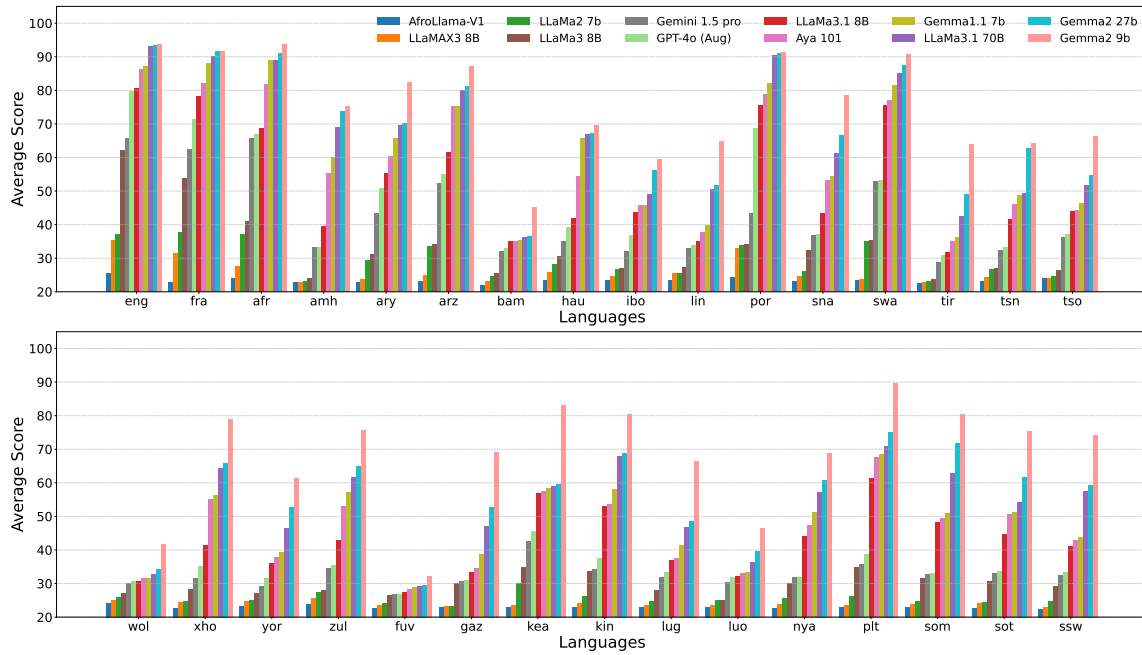


Figure 15: Per-language performance results for the AFRIQA dataset.

2381  
2382

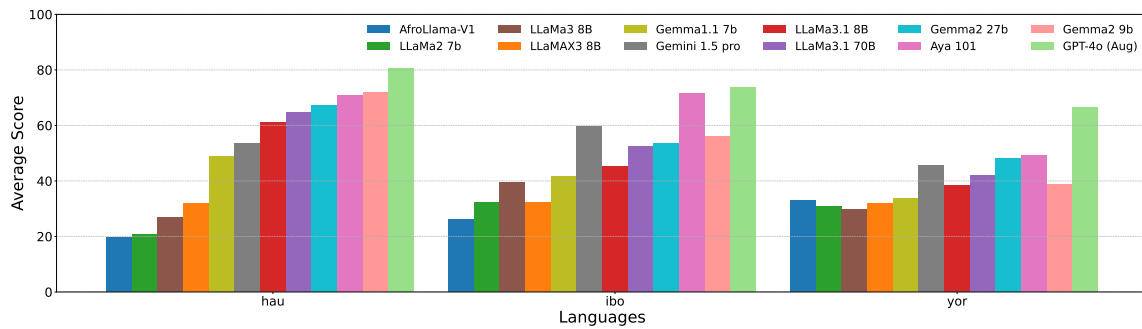
## E.3.2 Reading Comprehension

### Belebele



2383

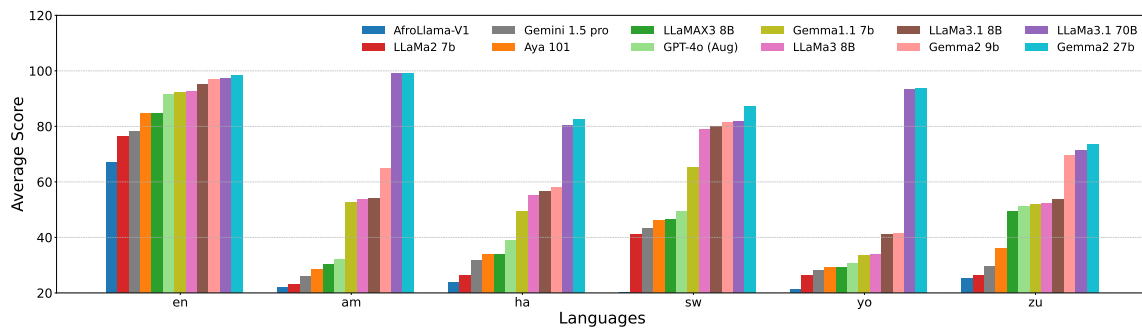
### NaijaRC



2384  
2385

## E.4 Knowledge

### Arc-E



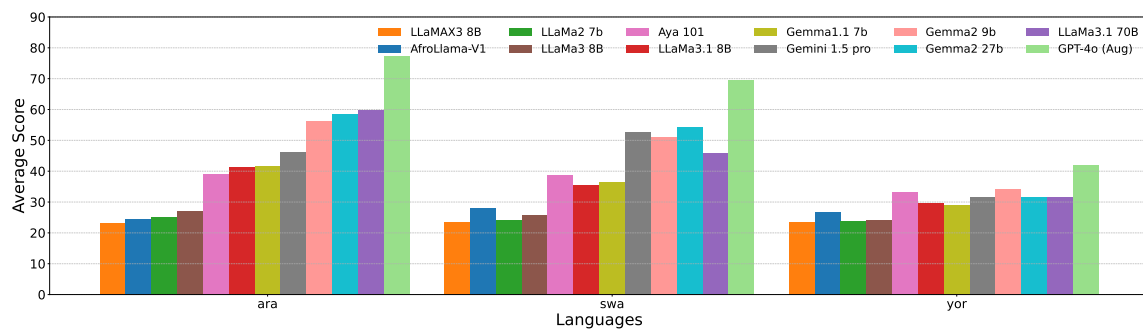


Figure 19: Per-language performance results for the OPENAI-MMLU dataset.

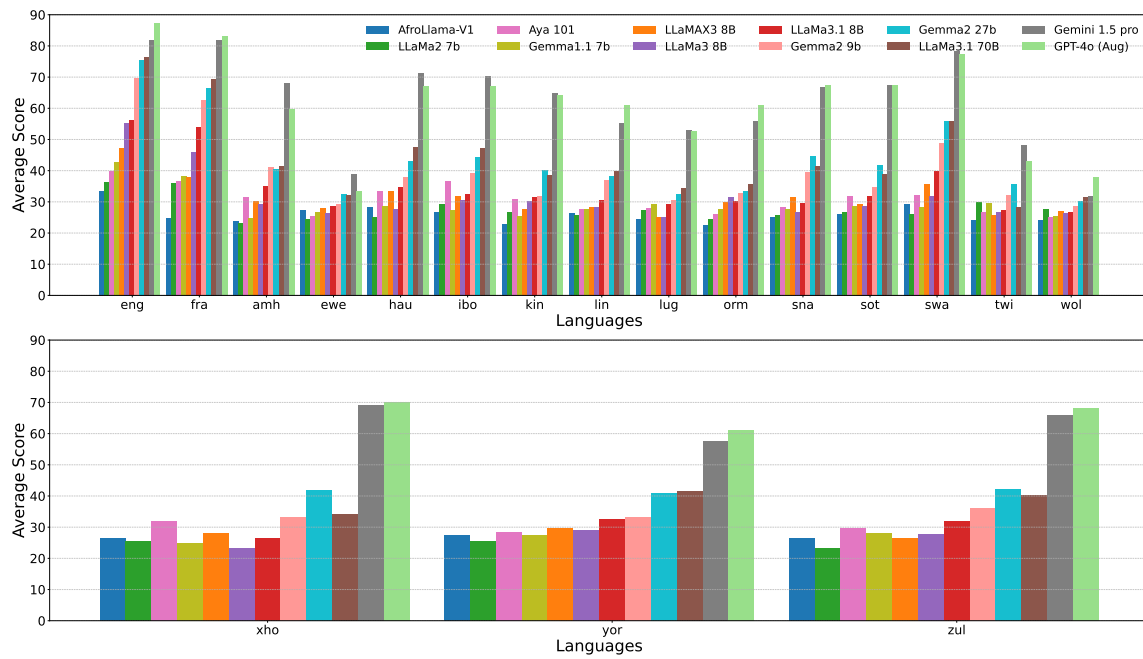


Figure 20: Per-language performance results for the AFRIMMLU dataset.



2388  
2389

## E.5 Reasoning

### AfriMGSM

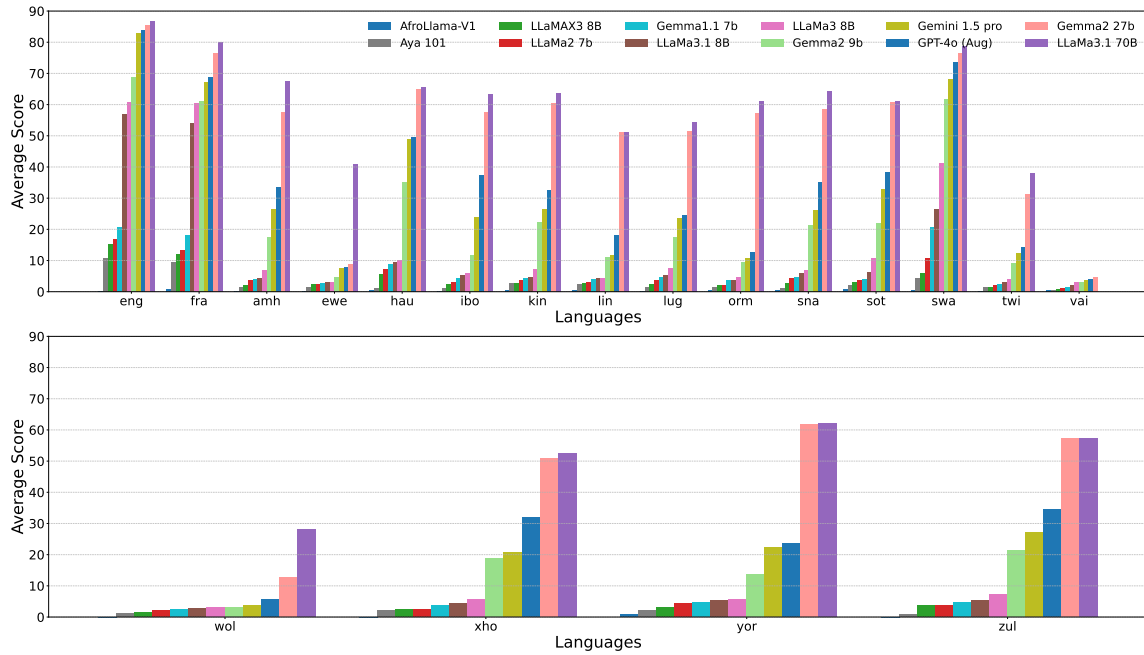


Figure 21: Per-language performance results for the AFRIMGSM dataset.

2390  
2391  
2392

## E.6 Text Generation

### E.6.1 Machine Translation

#### SALT (*en/fr-xx*)

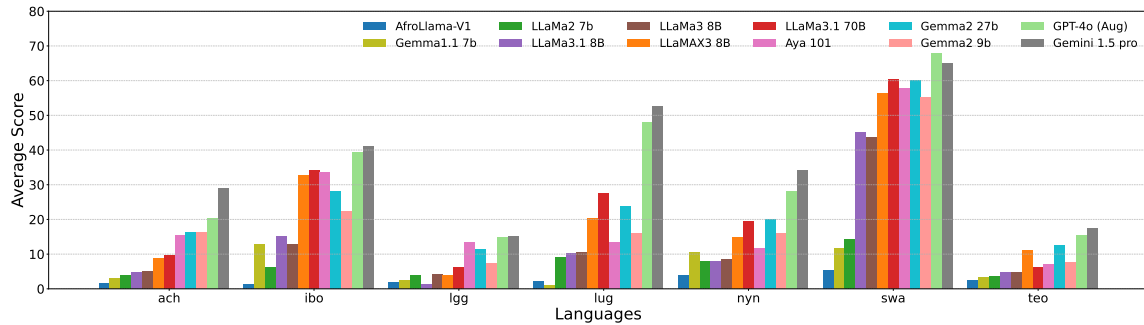


Figure 22: Per-language performance results for the SALT dataset (*en/fr-xx*).

2393

#### SALT (*xx-en/fr*)

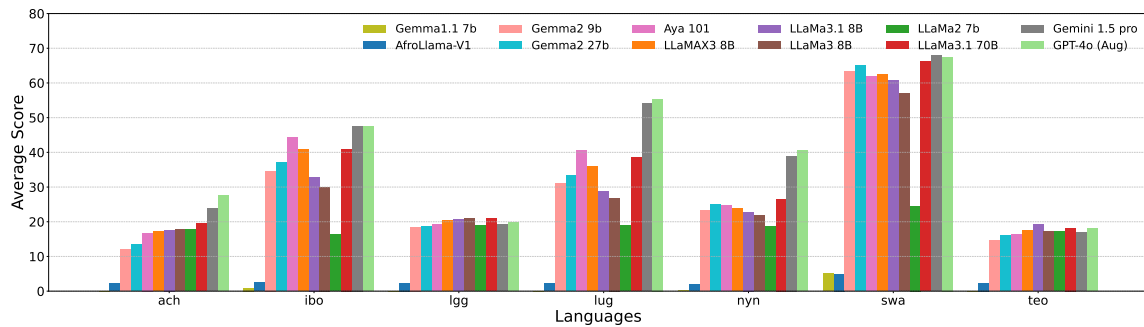


Figure 23: Per-language performance results for the SALT dataset (*xx-en/fr*).

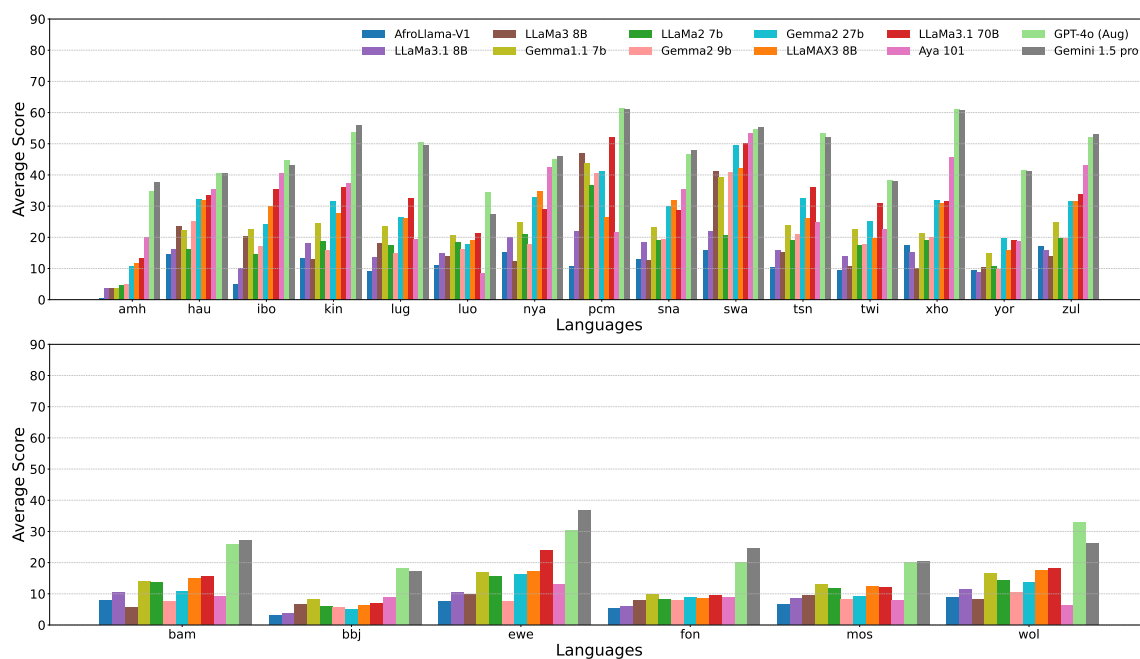


Figure 24: Per-language performance results for the MAFAND dataset.

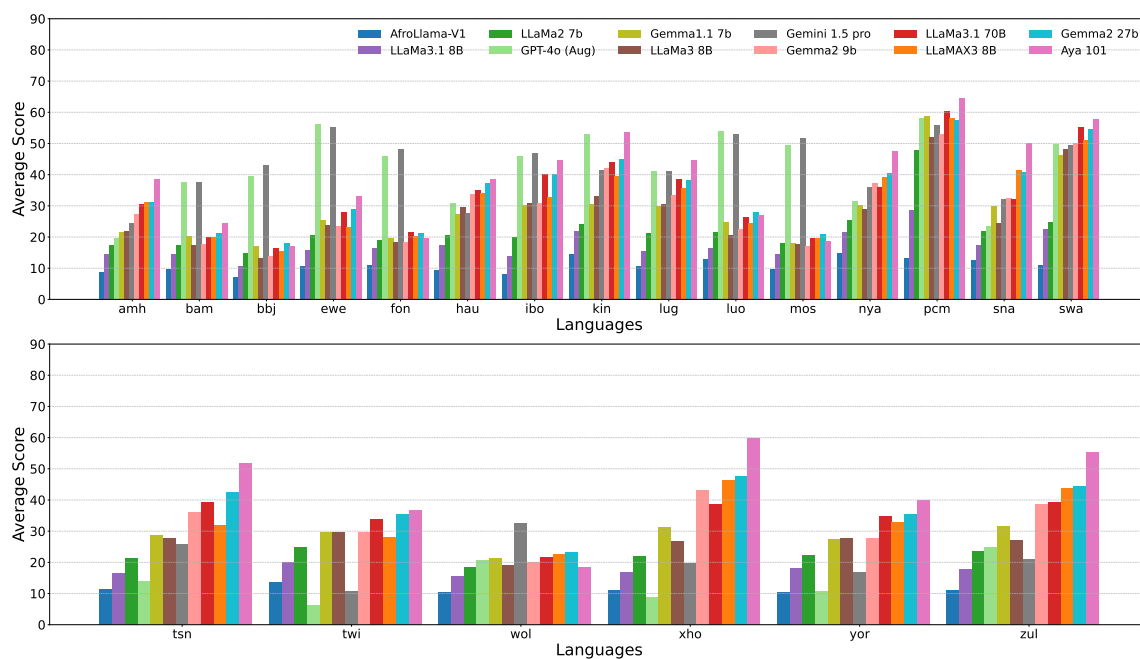


Figure 25: Per-language performance results for the MAFAND dataset.

NTREX (*en/fr-xx*)

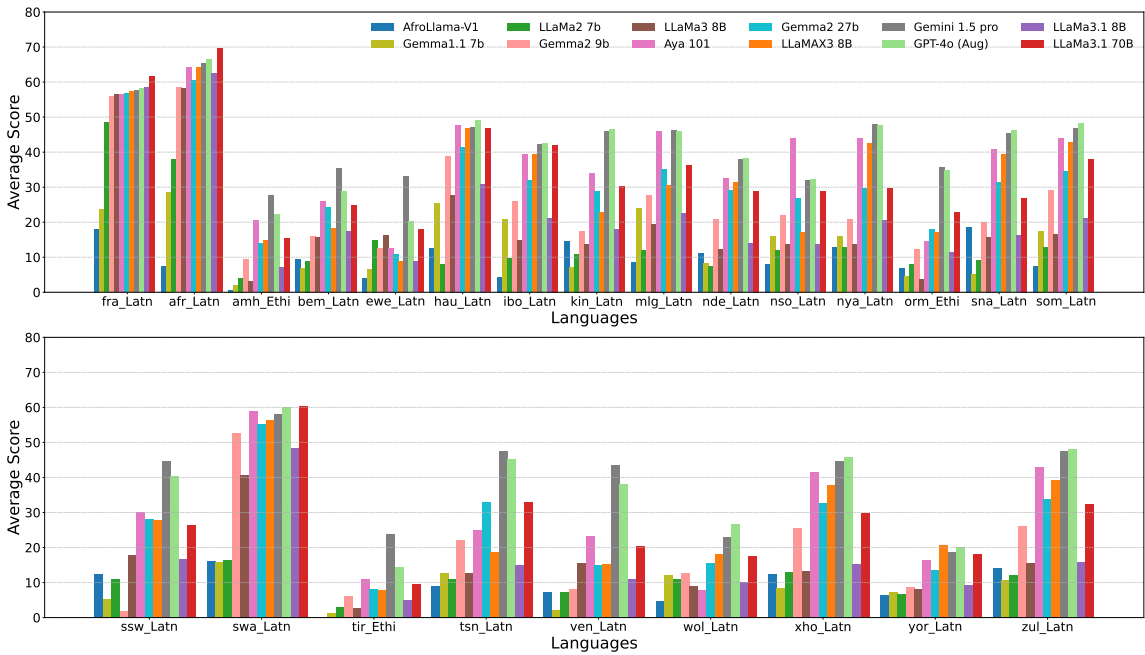


Figure 26: Per-language performance results for the NTREX-128 dataset (*en/fr-xx*).

NTREX (*xx-en/fr*)

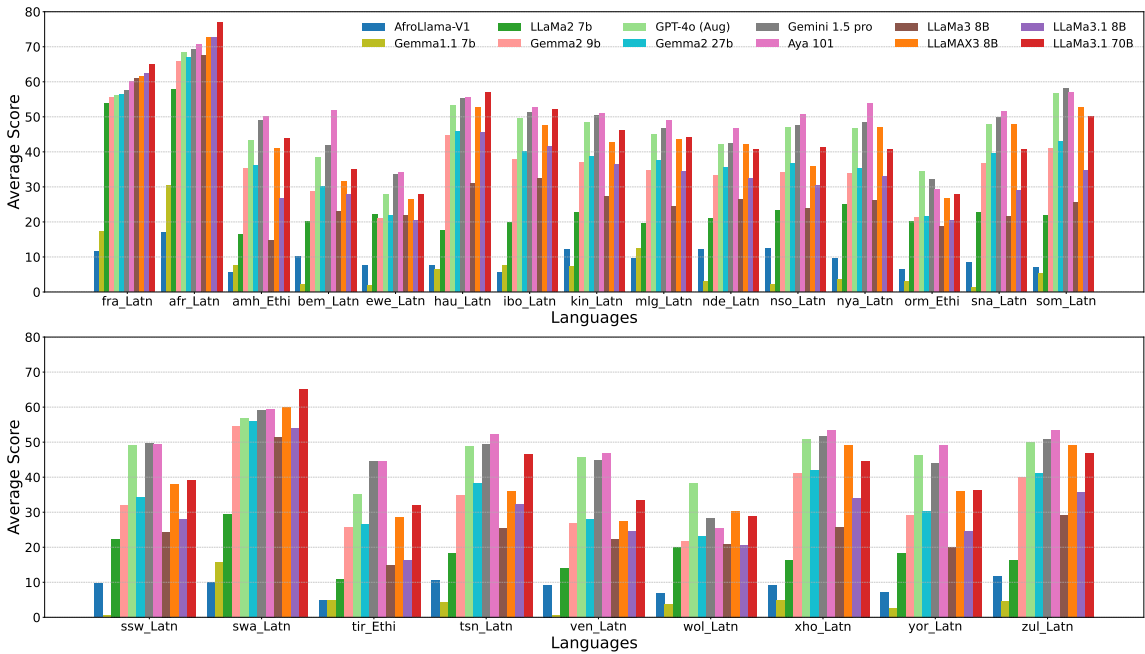


Figure 27: Per-language performance results for the NTREX-128 dataset (*xx-en/fr*).



Figure 28: Per-language performance results for the FLORES dataset (*en/fr-xx*).

Flores (African Languages only) (xx-en/fr)

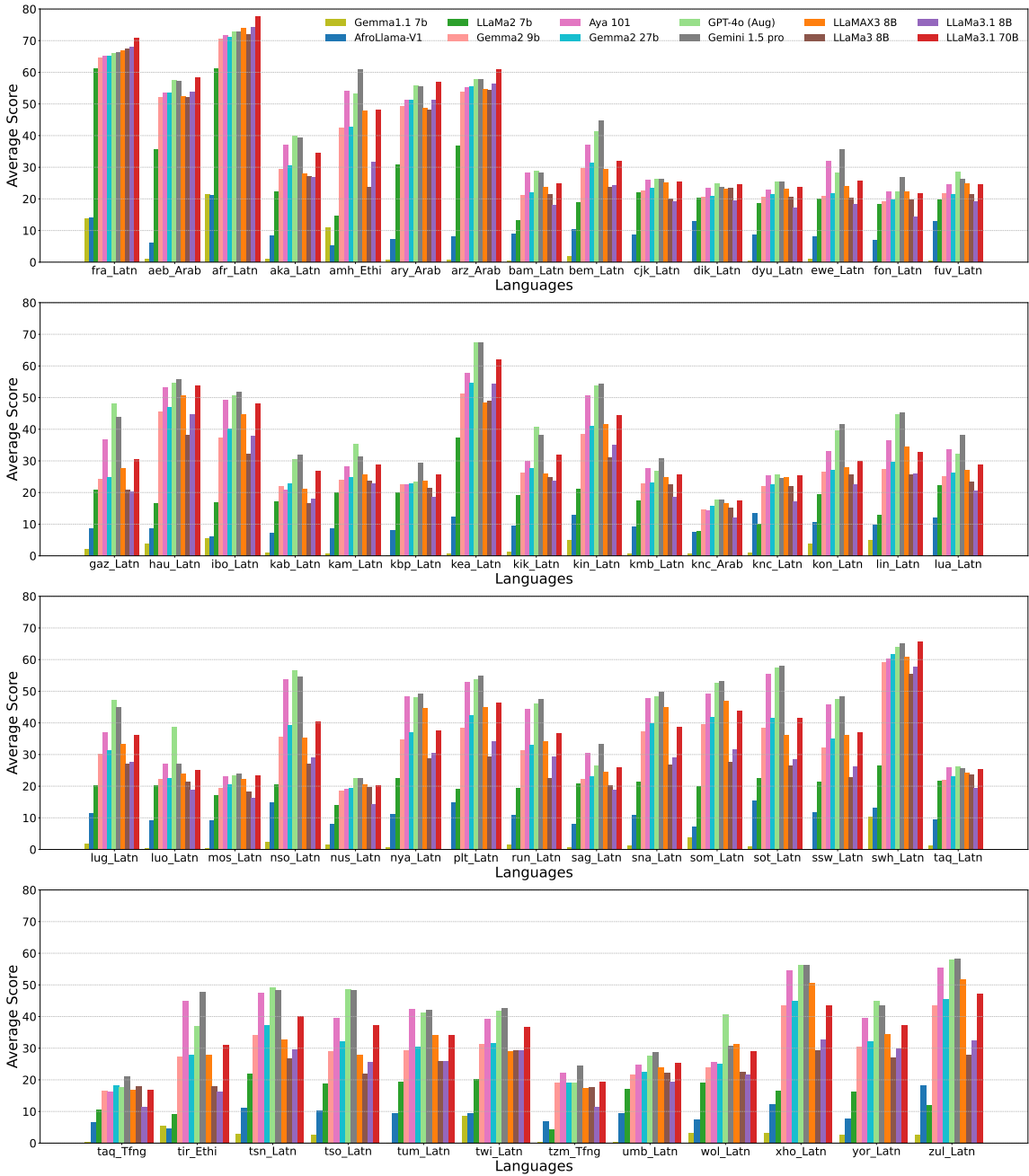


Figure 29: Per-language performance results for the FLORES dataset (xx-en/fr).



E.6.2 Summarization

2400

XL-SUM

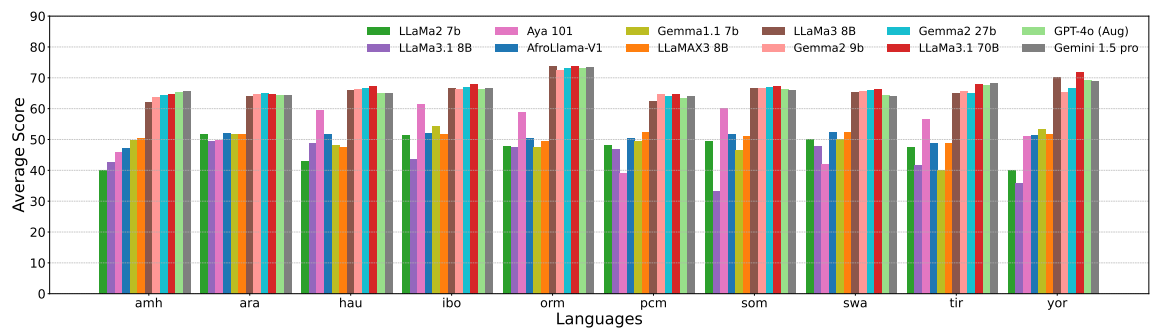


Figure 30: Per-language performance results for the XL-SUM dataset.

2401

E.6.3 Diacritics Restoration

2402

AFRIADR

2403

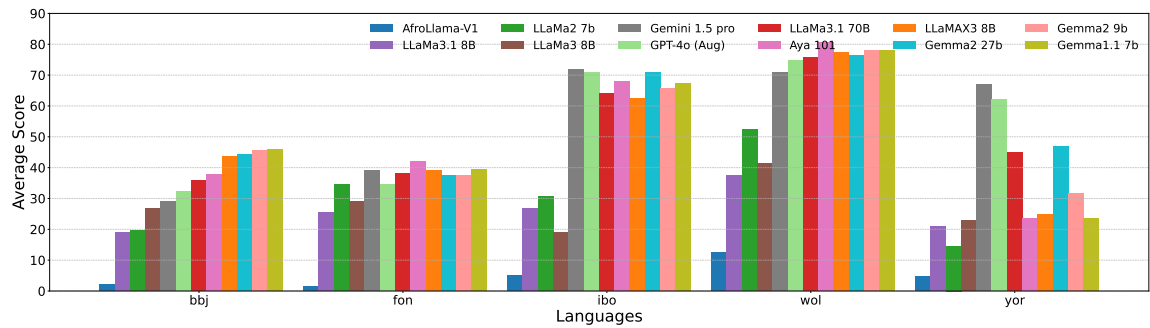


Figure 31: Per-language performance results for the AFRIADR dataset.