

# Tokenizer-Aware Cross-Lingual Adaptation of Decoder-Only LLMs through Embedding Relearning and Swapping

Anonymous ACL submission

## Abstract

Extending Large Language Models (LLMs) to support new languages is a challenging problem, with most methods proposed suffering from high computational cost and catastrophic forgetting of original model capabilities. Embedding relearning (Artetxe et al., 2020), a technique that creates new tokenizers and tunes embeddings on fixed model weights for target language adaptation, is both lightweight and performant. However, it has only been demonstrated to work for older generation encoder-only models and for high resource languages. In this paper, we extend this framework to decoder-only LLMs focusing on joint adaptation to many languages, including low-resource languages. We experiment in three language groups over 100 languages each. Our approach adapts a pre-trained model via switching to a customized tokenizer, and relearning the embedding layer. Across 96 diverse languages spanning both classification and generation tasks, we demonstrate embedding relearning improves Gemma2 models (with up to 27B parameters) by up to 20%, being more effective than, or on par with, full-weight updating baselines while effectively mitigating English forgetting (1-3% regressions). Analysis reveals the critical role of customizing tokenizers in achieving effective language transfer, particularly for non-Latin script languages. We further show embedding relearning helps transfer reasoning abilities across languages, achieving a 14% improvement over a math-optimized LLM across 20 languages.

## 1 Introduction

Large Language Models (LLMs) have transformed the field of natural language processing through pre-training on extensive web-scale corpora (Brown et al., 2020; Anil et al., 2024). Despite these advancements, their success has been primarily centered on English, leaving the multilingual ability less explored. While the multilingual potential of

LLMs has been demonstrated across multiple languages (Shi et al., 2023), their practical applications remain largely confined to a limited set of high-resource languages. This limitation reduces their utility for users speaking under-represented languages (Ahia et al., 2023).

Recently, many works focus on increasing the language support of LLMs. For instance, continued pre-training approaches further train LLMs on additional multilingual data by using either the original vocabulary (Zheng et al., 2024; Üstün et al., 2024) or introducing language-specific tokens (Fujii et al., 2024; Lu et al., 2024) (Figure 1 (b)). Despite the effectiveness in enhancing multilingual support, this paradigm typically requires full-parameter tuning on vast data, making adapting an LLM to accommodate new languages expensive. Moreover, such adaptation poses a significant risk of catastrophic forgetting, whereby the LLM loses previously acquired knowledge from the initial pre-training phase (Luo et al., 2024; Shi et al., 2024). Additionally, these methods usually focus on a smaller number of high-resource languages (Aryabumi et al., 2024; Alves et al., 2024; Xu et al., 2024), neglecting mid- and low-resource languages. Given these challenges, it is crucial to explore efficient and scalable methods for developing multilingual LLMs that can support diverse languages across varying resource levels.

Alternatively, *embedding relearning* – a technique that adapts models to new languages by learning new tokenizers and embeddings while keeping the transformer layers fixed – has shown to be effective in improving performance in target languages (Artetxe et al., 2020). However, existing embedding relearning approaches face several limitations: 1) the major focus on ‘outdated’ encoder-only PLMs limits the applicability; 2) the monolingual transfer strategy (i.e., one embedding per language) reduces efficiency; 3) limited linguistic coverage and exclusive evaluation on classification

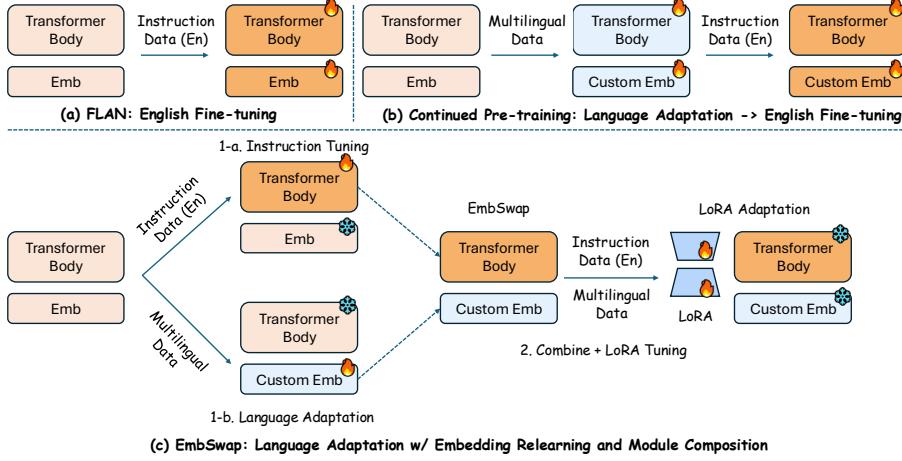


Figure 1: Methods for cross-lingual transfer of LLMs. (a): Use English instruction data for fine-tuning for vanilla cross-lingual transfer. (b): Continued pre-training: create customized tokenizers and do full-parameter tuning on multilingual data followed by English instruction-tuning. (c): EMBSWAP: 1-a: freeze the original embeddings of LLMs and instruction-tune the transformer body using English alignment data; 1-b: learn new multilingual embeddings by freezing the transformer body for target language adaptation of LLMs; 2) combine new embeddings with instruction-tuned transformer body as the EMBSWAP and further perform LoRA tuning to connect the combined components for enhanced cross-lingual transfer.

tasks leave the diverse linguistic and task settings underexplored (Liang et al., 2023). These limitations prompt important questions: 1) Can this approach be extended to modern decoder-only LLMs and how well it performs on generation tasks as compared to well-studied classification ones? 2) Can we achieve efficient cross-lingual transfer to diverse languages with different resource levels altogether? 3) What are the effective ways of building tokenizers for a large group of target languages that also help preserve LLM’s existing abilities?

We revisit *embedding relearning* and extend it to decoder-only LLMs for cross-lingual transfer across many languages. Figure 1 (c) shows that the pipeline starts with two parallel training from LLMs primarily pre-trained on English data: 1) We adapt it to a target language group<sup>1</sup> by learning multilingual embeddings on new tokenizers while fixing the transformer body; 2) We employ English alignment data to instruction-tune the transformer body, keeping the original embeddings fixed. After this, the relearned multilingual embeddings are combined with the instruction-tuned transformer body for efficient zero-shot cross-lingual transfer, which we denote as EMBSWAP. We can further enhance the compatibility of the composed modules via a cost-effective LoRA-based adaptation.

We conduct experiments on diverse multilingual NLP tasks spanning classification, generation and

mathematical reasoning. Results show that EMBSWAP surpasses the original Gemma2 models at all scales and achieves results comparable to full-weight updating baselines, while effectively mitigating knowledge forgetting. Our findings indicate that carefully-curated tokenizer construction and embedding initialization methods are crucial, and non-Latin script languages benefit the most from customized tokenizers. We believe these insights equip practitioners with a recipe for achieving effective and efficient cross-lingual transfer in LLMs through tokenizer and embedding relearning.

## 2 Methodology

Prior embedding relearning methods create dedicated tokenizer and embedding layer for each target language, which is practical as the involved language model is monolingual. By contrast, most recent LLMs exhibit a degree of multilingualism in both their representations and tokenizers (Blevins and Zettlemoyer, 2022). A naive solution to applying embedding relearning to LLMs is to continue training the embedding layer on additional multilingual data. However, the tokenizers employed in these LLMs are biased towards English and several high-resource languages. This bias results in the over-fragmentation of text from long-tail languages (Figure 2), thereby degrading the performance and efficiency of processing such languages (Ahuja et al., 2023). In this paper, we demonstrate that having tokenizers that provide fairer representation

<sup>1</sup>We focus on three language groups: South East Asia (SEA), Indic (IND), and African (AFR).

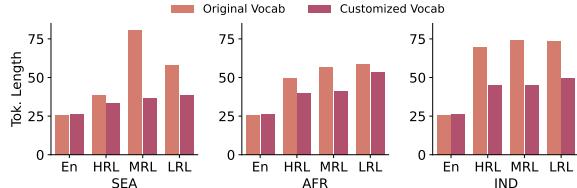


Figure 2: The tokenization comparison between using the vanilla and customized multilingual tokenizers on Gemma2. Tok. Length refers to the average number of tokens required to represent the same amount of texts.

across languages is critical to achieving effective cross-lingual transfer.

## 2.1 Customized Vocabulary Construction

Our strategy involves constructing distinct tokenizers for each language group (§3.2). Tailoring tokenizers to specific language groups enhances cross-lingual transfer among geographically related long-tail languages compared to using monolingual tokenizers. Moreover, this approach avoids the shortcomings of a universal tokenizer that treats all low-resource languages uniformly poorly. Based on this, we propose a *Prune-with-Extension* approach for developing tokenizers optimized for language adaptation while maintaining English ability. First, we prune the tokenizer of target LLMs by removing non-English tokens. Then the pruned tokenizer is extended through adding new tokens, which are obtained by training tokenizers for target languages using BPE (Gage, 1994; Sennrich et al., 2016).

**Pruning the Tokenizer** To preserve the pre-trained knowledge embedded in the language model, current approaches often expand the vocabulary by adding new tokens (Fujii et al., 2024; Cui et al., 2024). This, however, can substantially increase pre-training time due to the extra computational cost of the output embedding layer (Liang et al., 2023). To avoid this, we first prune the existing tokenizers by retaining only English tokens before adding those from low-resource languages. Given the predominant English training data for LLMs, we hypothesize that a significant portion of their knowledge is associated with English tokens, and reusing English tokens can effectively retain this knowledge (Garcia et al., 2021). In our implementation, for a given LLM, we identify English tokens by tokenizing a set of 20 million English sentences using its tokenizer, with further filtering through removing non-Latin script tokens.<sup>2</sup>

**Training Multilingual Tokenizers** To get the

<sup>2</sup>40% of tokens are discarded: non-Latin scripts tokens from high-resource languages and very rare English tokens.

data for building a multilingual vocabulary for long-tail languages, we sample from the Next Thousand Languages (NTL) corpus (Caswell et al., 2020; Bapna et al., 2022). Our empirical analysis reveals that sampling data for each language up to a maximum of 500K lines from NTL effectively addresses the imbalance between high- and low-resource languages, outperforming temperature-sampling techniques. Subsequently, we train a BPE tokenizer using the sampled data to generate a vocabulary whose size aligns with that of the target LLMs.

**Extending the Pruned Tokenizer** We sequentially add the identified English tokens followed by tokens from the newly built multilingual tokenizer. Both types of tokens are added in the same preference order as in their respective tokenizers. The final vocabulary maintains the same size as the original tokenizer, with over 60% token overlap, resulting in negligible variations in English tokenizations (see Table 6 in Appendix). Figure 2 shows our final vocabulary achieves significant compression rate improvements by consistently producing shorter sequences across languages of a spectrum of resource levels while barely affecting English.

## 2.2 Training Recipe

**Embedding Initialization** To maximally inherit the pre-trained knowledge embedded in the target LLM’s embedding layer, we adopt a strategy inspired by Gee et al. (2022). For tokens that overlap between the target LLM’s vocabulary and our multilingual vocabulary, we directly copy the corresponding embeddings. For new tokens, we employ the LLM’s original tokenizer to decompose them into subtokens and initialize their embeddings using the average of their subtoken embeddings.<sup>3</sup>

**Language Adaptation** We fine-tune the customized embeddings on 200B curated multilingual tokens  $\mathcal{D}_{la}$  while keeping the transformer body frozen (Figure 1 (c) step 1-b), with the same training objective used for pre-training the base LLM. This is based on the assumption that the pre-trained transformer body encapsulates universal cross-lingual knowledge (Zhao et al., 2024b; Wendler et al., 2024; Tang et al., 2024), while the embedding layer encodes language-specific information, which suggests embedding tuning should be effective for language adaptation.

**EMBSWAP** We instruction-tune the transformer

<sup>3</sup>We find initialization from original embeddings are crucial, while the exact method (averaging or max-pooling) makes minimal difference (See Appendix Figure 16).

body of base LLMs on a diverse range of English tasks  $\mathcal{D}_{it}$  (Wei et al., 2022) with 20B tokens (Figure 1 (c) step 1-a). Notably, we employ the LLM’s *original* embeddings and keep them frozen in this step.<sup>4</sup> We then integrate the customized embeddings obtained from the language adaptation stage into the instruction-tuned transformer body.

**LORA-Adaptation** Since the transformer body and customized embeddings are independently trained, EMBSWAP may suffer from incompatible parameters. Our empirical findings indicate that EMBSWAP is effective for discriminative tasks but sometimes underperforms an instruction-tuned baseline on generative tasks. To mitigate this and ensure the assembled model’s effectiveness across various tasks, we insert LoRA weights into the self-attention layer of the tuned transformer body (Figure 1 (c) step 2). These weights are then fine-tuned on a sub-sampled joint corpus  $\mathcal{D}_{mix} = \mathcal{D}_{la} \cup \mathcal{D}_{it}$  with 4B tokens, while keeping both the body and embeddings frozen.

### 3 Experiment setup

#### 3.1 Pre-training Data

The data  $\mathcal{D}_{la}$  for embedding training is a mixed corpus with 65% sentence-level and 35% document-level texts. The sentence-level data is exclusively from the Next Thousand Languages (NTL) corpora (Caswell et al., 2020; Bapna et al., 2022), which provides web-crawled monolingual sentences and translation pairs for over 1000 languages. For document-level texts, we sample data from multilingual Wikipedia and mC4 (Xue et al., 2021). We use UniMax sampling (Chung et al., 2023) with  $N = 5$  to up-sample low-resource languages. Please refer to Appendix Figure 14 for pre-training data ablations.

We take FLAN (Wei et al., 2022) as the instruction tuning data  $\mathcal{D}_{it}$ . The data  $\mathcal{D}_{mix}$  used in LoRA-Adaptation consists of a 10% sample of  $\mathcal{D}_{it}$ , combined with an equal number of instances from  $\mathcal{D}_{la}$ .<sup>5</sup>

#### 3.2 Languages

We select languages from three families based on geographic relations: South East Asian (SEA), African (AFR), and Indic (IND). This results in 212 languages from SEA, 392 from AFR, and 170

<sup>4</sup>Freezing embeddings makes the instruction tuning and language adaptation processes symmetric. This enhances modularity and improves the parameter compatibility.

<sup>5</sup>This ensures the instruction-following ability isn’t forgotten. Please refer to Appendix Figure 15 for detailed analysis.

from IND. Each regional dataset is processed separately, with a tailored tokenizer, language-adapted embeddings and LoRA update parameters.

### 3.3 Models

Our evaluation is focused on Gemma2 (2B, 9B, 27B) (Riviere et al., 2024). We also test the generalization ability of embedding relearning in two LLMs with varying degrees of multilinguality: Aya23 (8B, 35B) (Aryabumi et al., 2024) and PaLM2 (XXS, S) (Anil et al., 2023).

As shown in Figure 1, we end up having four types of models: (i) FLAN (step 1-a): models that undergo instruction tuning. (ii) Lang-Adapt (step 1-b): the LLMs after language adaptation with embedding tuning. (iii) EMBSWAP (step-2 Combine): model denoted as EMBSWAP is constructed by combining the transformer body of FLAN with the embeddings from Lang-Adapt. (iv): LoRA-Adapt (step-2 LoRA Tuning): EMBSWAP models with the LoRA-Adaptation process. Detailed training procedures for each model type are in Appendix A.1.

We compare EMBSWAP with three variants: 1)

**Continued Pre-training (CPT)**: the most costly method which first fully tunes the base model on multilingual data  $\mathcal{D}_{la}$  to obtain language-adapted model  $\theta_{la}$ , then instruction tunes  $\theta_{la}$  on English data  $\mathcal{D}_{it}$ . 2) **ChatVector** (Huang et al., 2024): which avoids repeated instruction tuning, through obtaining a task vector (Ilharco et al., 2023) from a base model and IT model:  $\theta_{task} = \theta_{it} - \theta_{base}$ , which is then added to  $\theta_{la}$ . 3) **IT+Lang-Adapt**: an embedding relearning variant that performs one-step embedding tuning on Gemma2-IT model with the combination of  $\mathcal{D}_{la}$  and  $\mathcal{D}_{it}$ . For all three variants, we adopt the same customized tokenizer, initialized embeddings, and training data conditions (i.e., 200B  $\mathcal{D}_{la}$  and 20B  $\mathcal{D}_{it}$ ) for fair comparison.

### 3.4 Evaluation Tasks

We use *five-shot* prompting for evaluating language-adapted LLMs. In contrast, EMBSWAP is evaluated in a *zero-shot* setting, given it has been instruction tuned. We also evaluate EMBSWAP using a compiled English benchmark (Appendix B) to examine potential regressions in general English ability.

We perform evaluation with various multilingual tasks including machine translation (MT), multiple choice question answering (MCQA), topic classification, mathematical reasoning, cross-lingual question answering (QA) and summarization. For MT, we focus on EN-XX and select 23 SEA, 42 AFR,

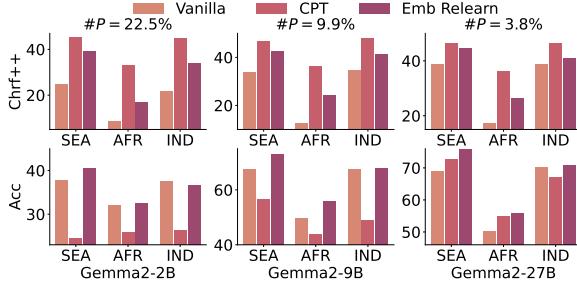


Figure 3: FLORES-200 EN-XX and BELEBELE Language Adaptation results across all sizes of Gemma2 models with 5-shot prompting.  $\#P$ : fraction of tuned parameters in embedding relearning.

Model	SEA	AFR	IND	Avg.
CALM (PaLM2-XXS-NTL+S) <sup>†</sup>	25.3	17.8	17.9	19.8
PaLM2-S-NTL	25.2	17.4	15.1	18.2
PaLM2-S	22.0	15.3	14.2	16.4
★ Lang-Adapt	<b>25.6</b>	<b>18.8</b>	<b>22.3</b>	<b>22.3</b>
Gemma2-9B	19.9	18.3	13.6	16.4
★ Lang-Adapt	<b>36.4</b>	<b>25.8</b>	<b>29.6</b>	<b>30.4</b>
Gemma2-27B	30.6	22.2	18.4	22.4
★ Lang-Adapt	<b>41.9</b>	<b>31.8</b>	<b>27.5</b>	<b>32.2</b>

Table 1: Language adaptation results on GSM8K-NTL, with the best in **bold**. <sup>†</sup>: from Bansal et al. (2024).

and 21 IND languages from FLORES-200 (Goyal et al., 2022). For MCQA, we use BELEBELE (Bopardkar et al., 2024) and evaluate on 15 SEA, 25 AFR, and 19 IND languages. For multilingual mathematical reasoning, we use GSM8K-NTL (Shi et al., 2023; Bansal et al., 2024) translated from GSM8K and select 5 SEA, 5 AFR, and 10 IND languages. For cross-lingual QA and summarization, the XORQA-IN and XSUM-IN datasets from INDICGENBENCH (Singh et al., 2024) are used and 29 IND languages are included for evaluation.

### 3.5 Main Results

We first present the results of Language Adaptation through embedding tuning on its own (§2.2 – Language Adaptation) in §3.5.1. Then we evaluate EMBSWAP alongside LoRA adaptation (§2.2 – EMBSWAP & LoRA-Adaptation) in §3.5.2.

#### 3.5.1 Language Adaptation Results

**Embedding relearning improves cross-lingual transfer.** Figure 3 shows embedding relearning consistently outperforms vanilla Gemma2 across three language groups and achieves better cross-lingual transfer than CPT, as evident by BELEBELE results, a task that relies on skills learnt from English. While CPT excels on FLORES (likely due to its full model adaptation for MT with training on parallel data), it does not generalize well to other

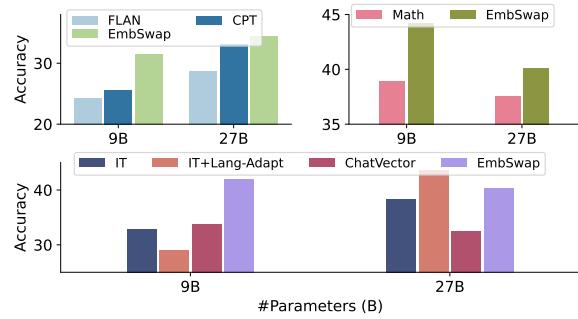


Figure 4: Zero-shot multilingual mathematical reasoning results on GSM8K-NTL. IT: LLMs are instruction-tuned without fixing embedding.

task types. See Appendix C.1 for more language adaptation results on Aya23 and PaLM2.

**Embedding relearning helps preserve LLMs’ reasoning knowledge.** We evaluate how language adaptation affects English reasoning transfer to other languages. Table 1 shows embedding relearning consistently improves the mathematical reasoning ability across various low-resource languages. Compared to CALM (Bansal et al., 2024), a form of adapter enabling model composition, embedding relearning achieves more substantial improvements over PaLM2-S with only embedding tuning and incurring no extra inference costs.<sup>6</sup> Moreover, it shows superior performance (+4%) compared to PaLM2-S-NTL that was full-parameter tuned on NTL.<sup>7</sup> Overall, our results suggest that embedding relearning offers an effective alternative to CALM and full-parameter tuning. We observe performance improvements with more modern models, with particularly pronounced gains in Gemma2.

#### 3.5.2 EMBSWAP Results

**EMBSWAP transfers mathematical reasoning abilities across languages.** Figure 4 shows that EMBSWAP outperforms the instruction-tuned baseline for zero-shot evaluation on mathematical reasoning tasks, with up to 8% gains across 20 low-resource languages. Moreover, EMBSWAP further advances performance by incorporating LLMs with enhanced mathematical reasoning ability. These include the IT models aligned via reinforcement learning and math models instruction-tuned on reasoning intensive data (Yue et al., 2024).

<sup>6</sup>Embedding relearning requires more training as it must be repeated for different linguistic regions.

<sup>7</sup>Both methods use the same training dataset. However, NTL training did not create region-specific models by splitting training data, potentially diminishing its effectiveness due to the curse of multilinguality (Conneau et al., 2020), which arises when a single model is trained on too many languages.

Eval. Metric	% Tuned Params (Lang Adapt)	Task Type	English	Classification						Generation						Avg.	
				BELEBELE Accuracy			SIB-200 Accuracy			FLORES-200 ChrF++			XORQA-IN Token-F1		XSUM-IN ChrF		
				SEA	AFR	IND	SEA	AFR	IND	SEA	AFR	IND	IND	EN	IND	EN	
<b>Gemma2-2B</b>																	
FLAN	-	73.0	52.0	36.0	50.2	65.9	47.3	67.8	27.8	9.6	20.7	7.6	47.8	3.7	<b>36.6</b>	31.7	
IT+Lang-Adapt	22.5%	54.1	51.2	36.7	42.9	67.0	<b>48.7</b>	67.1	<u>35.4</u>	<u>14.2</u>	<u>29.2</u>	<b>12.3</b>	40.2	<u>7.1</u>	27.5	33.0	
EMBSWAP	22.5%	<u>69.8</u>	<u>62.0</u>	<u>40.2</u>	<u>54.4</u>	<u>72.9</u>	<b>57.9</b>	<u>70.5</u>	<u>27.0</u>	11.0	<u>17.2</u>	8.6	<u>54.1</u>	6.6	<u>34.1</u>	<u>35.3</u>	
★ LoRA-Adapt	22.5%	69.7	<b>63.3</b>	<b>43.7</b>	<b>55.8</b>	<u>73.5</u>	47.8	<b>70.7</b>	<b>37.1</b>	<b>18.3</b>	<b>32.5</b>	9.8	<b>60.5</b>	<b>9.9</b>	<u>31.9</u>	<b>38.7</b>	
ChatVector	100%	39.1	28.8	32.5	30.6	24.5	29.3	52.2	42.1	32.7	36.1	19.4	16.1	19.3	21.9	28.0	
CPT	100%	64.3	55.4	49.0	58.9	66.9	65.9	78.4	44.2	24.4	36.0	11.2	68.2	14.8	34.0	43.0	
<b>Gemma2-9B</b>																	
FLAN	-	<b>83.5</b>	70.6	49.9	68.9	74.4	61.0	79.1	32.0	12.6	27.8	9.8	60.3	15.1	<b>37.5</b>	42.0	
IT+Lang-Adapt	9.9%	73.9	75.4	<u>57.9</u>	<u>72.7</u>	77.2	<b>69.2</b>	79.7	<u>39.3</u>	<u>22.2</u>	<u>38.7</u>	<b>23.1</b>	<b>69.6</b>	<u>17.0</u>	33.9	<b>48.3</b>	
EMBSWAP	9.9%	82.3	<u>78.1</u>	<u>57.5</u>	71.2	<b>78.6</b>	<u>68.8</u>	<u>79.8</u>	35.9	16.5	31.4	<u>12.5</u>	<u>65.2</u>	15.2	<u>37.1</u>	45.5	
★ LoRA-Adapt	9.9%	<u>82.5</u>	<b>78.4</b>	<b>60.2</b>	<b>73.1</b>	<u>78.5</u>	<b>69.2</b>	<b>80.1</b>	<b>40.0</b>	<b>25.6</b>	<b>40.5</b>	12.1	64.1	<b>17.4</b>	<u>37.0</u>	<u>47.3</u>	
ChatVector	100%	58.4	69.0	52.5	59.3	82.1	66.1	77.2	44.5	26.6	40.5	17.3	29.8	22.1	26.9	43.7	
CPT	100%	74.6	76.3	65.4	74.1	83.0	64.8	80.0	45.2	34.1	44.6	19.5	72.7	22.6	37.0	51.7	
<b>Gemma2-27B</b>																	
FLAN	-	<b>84.3</b>	71.9	52.6	72.9	72.6	60.2	76.4	33.2	15.2	29.6	<u>20.4</u>	61.7	15.0	<b>37.6</b>	43.7	
IT+Lang-Adapt	3.8%	81.9	<b>81.7</b>	<b>62.7</b>	<b>76.0</b>	<b>82.1</b>	<b>70.0</b>	<b>81.9</b>	<b>42.5</b>	<u>24.8</u>	<u>39.3</u>	<b>26.4</b>	<b>73.0</b>	<b>18.3</b>	35.8	<b>51.1</b>	
EMBSWAP	3.8%	<u>83.6</u>	78.8	56.3	73.1	78.1	66.1	78.5	<u>33.7</u>	<u>13.2</u>	<u>23.9</u>	<u>15.2</u>	<u>59.3</u>	<u>12.4</u>	<u>36.8</u>	<u>43.5</u>	
★ LoRA-Adapt	3.8%	83.4	<u>79.5</u>	<u>60.4</u>	74.3	<u>79.5</u>	<u>68.3</u>	<u>80.0</u>	<u>42.5</u>	<u>25.7</u>	<u>40.6</u>	<u>16.4</u>	<u>70.1</u>	18.0	<u>36.9</u>	49.1	
ChatVector	100%	62.9	78.5	55.7	71.7	81.5	55.9	77.1	43.4	31.4	37.8	16.2	18.2	20.0	22.8	43.0	
CPT	100%	79.5	81.2	68.6	77.6	72.6	72.3	78.2	45.4	34.4	45.0	21.1	72.4	23.7	38.7	52.5	

Table 2: Zero-shot cross-lingual transfer results. **Bold** and underlined: best and second-best results for embedding relearning methods. **Red** values indicate instances where EMBSWAP hurts FLAN. English results are excluded from the average. **% Tuned Params (Lang Adapt)**: proportions of trainable parameters during language adaptation.

**EMBSWAP is more effective for classification tasks.** Table 2 shows that EMBSWAP improves the performance on classification tasks across all model sizes with up to 10% gains. By contrast, for generation tasks, the method’s behavior is inconsistent, often leading to performance degradation. We attribute this phenomenon to intrinsic difficulty: classification tasks are generally easier as the solution space is typically small compared to generation tasks, making them more robust to embedding changes. However, the embedding layer is used for both text encoding and decoding in generation, and the auto-regressive generation paradigm makes the model sensitive to embedding changes due to error propagation accumulating over time steps.

**LoRA-Adaptation connects both worlds.** In Table 2, we demonstrate that the two components within EMBSWAP can cooperate better after LoRA-Adaptation, leading to significant gains on generation tasks. This results in an average improvement of 5.4% on the largest 27B variant. The results demonstrate the practical use of EMBSWAP for efficient zero-shot cross-lingual transfer of instruction-tuned LLMs. See Appendix C.2 and Table 4 for additional results on Aya23 and PaLM2.

**Embedding relearning on IT LLMs represents a viable alternative to EMBSWAP.** Table 2 reveals that for smaller model sizes, EMBSWAP

outperforms IT+Lang-Adapt in downstream tasks. However, as model size increases, IT+Lang-Adapt becomes the more effective solution. This suggests that IT+Lang-Adapt provides a highly effective alternative to EMBSWAP for zero-shot cross-lingual transfer involving a single model. However, when handling multiple IT models derived from a common base model, EMBSWAP offers greater modularity by eliminating the need for costly embedding relearning across individual models.

**EMBSWAP mitigates catastrophic forgetting.** CPT with full-weight updating boosts cross-lingual transfer but suffers from forgetting of English capabilities. ChatVector shows large drops both in English and multilingual tasks, demonstrating the technique is not well suited to our problem. In contrast, EMBSWAP shows only minor regressions in English, indicating the benefits of its modular composition which reuses existing knowledge. This advantage is clear in multilingual mathematical reasoning (Figure 4), where EMBSWAP, by preserving English knowledge, significantly outperforms full-weight updating methods. Similar findings are observed for Aya23 etc. (See Table 4).

### 3.6 Ablation Analysis

**Customized vocabulary amplifies the benefits of training on multilingual data.** We decouple

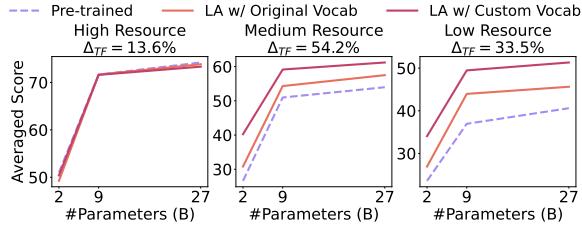


Figure 5: Tokenizer ablations for language adaptation. SEA averaged scores on FLORES-200 and BELEBELE are reported.  $\Delta_{TF}$ : % tokenizer fertility reduction.

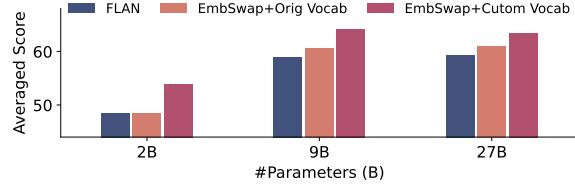


Figure 6: Tokenizer ablations under EMBSWAP. The macro-averaged scores on SEA subset of BELEBELE, SIB-200 and FLORES-200 are reported.

the effects of multilingual data in language adaptation from the change in vocabularies. Figure 5 shows that simply tuning embeddings on multilingual data can significantly improve the performance in medium and low-resource languages, indicating LLMs are under-fitting to these languages. Based on this, employing a new vocabulary customized for these languages can amplify the benefits of embedding tuning on multilingual data. This enhanced learning process facilitates better knowledge acquisition (Zhang et al., 2022; Hofmann et al., 2022). Moreover, the improvements are more pronounced in smaller models, highlighting the importance of effective tokenization for these models to adapt well to low-resource languages. Figure 6 indicates that the importance of customized vocabulary is also apparent in the EMBSWAP setting.

**Tokenization fertility is inversely correlated to downstream performance.** We study how tokenization fertility (the average number of tokens produced per word) affects the LLM’s performance across languages. To ablate this effect, we generate several re-sampled replicates of our tokenizer training datasets with different levels of priority given to high v.s. lower resource languages. Specifically, we use temperature sampling to manipulate the training sentences of each language for building different tokenizer training data, and train tokenizers with varying fertilities. We then relearn embeddings for each of these tokenizers and analyze the downstream performance. As shown in

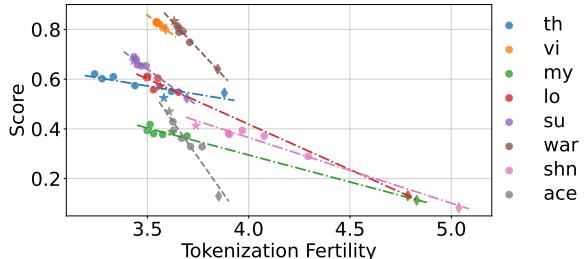


Figure 7: Correlation between the results of language adaptation on Gemma2-2B with tokenization fertility. Normalized ChrF++ on FLORES-SEA is reported. ◆ and ★ indicate the original and our customized tokenizers used in the other settings, respectively.

Figure 7, we find that the performance is inversely correlated to tokenization fertility, but the correlation is not uniform across languages. Notably, slight reductions in fertility lead to significant performance improvements in low-resource languages (e.g., ACE) while high-resource languages are less sensitive to fertility changes. Furthermore, Latin-script languages generally benefit more from fertility reductions compared to those in non-Latin scripts. Please refer to Appendix Figure 13 which reports similar findings for PaLM2-XXS.

**Pruning with extension outperforms other tokenizer construction methods.** We study how different ways of creating tokenizers affect how well LLMs adapt to new languages. We compared *Prune+Extension* to two variants: 1) *Scratch* trains tokenizers completely new for both English and target languages; and 2) *Extension* appends target language tokens to an existing tokenizer without removing any old ones, which adds 34% extra embedding parameters. Figure 8 shows *Prune+Extension* works the best overall. On a simpler task (FLORES), all methods perform similarly. We believe this is because FLORES is not as complex as BELEBELE, which needs reasoning skills likely learned from English. Supporting this, *Prune+Extension* greatly improves English performance on BELEBELE by keeping the original English tokens, and transfers this advantage to medium and low resource languages. By contrast, *Scratch* has far fewer original English tokens, suggesting that keeping these English tokens is important for retaining knowledge from pre-training.<sup>8</sup> *Prune+Extension* also outperforms *Extension*, showing that removing tokens

<sup>8</sup>It’s an open question of whether other facets of the tokenizer need to be retained to preserve other behaviours, e.g., markup tokens to facilitate code understanding.

471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504

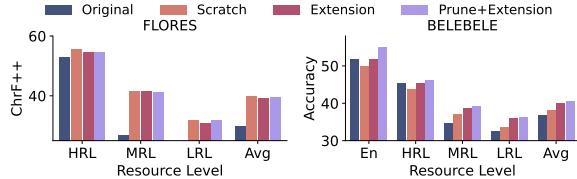


Figure 8: Ablations on tokenizer building methods. We report SEA language adaptation on Gemma2-2B.

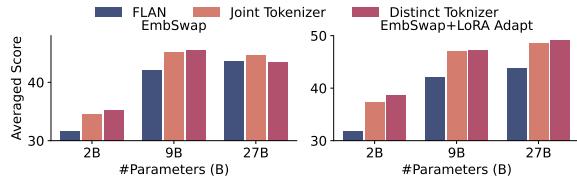


Figure 9: Distinct tokenizers per language group v.s. Joint tokenizer for all languages. Averaged results of five tasks included in Table 2 are reported.

irrelevant to target languages is beneficial.<sup>9</sup>

**A joint tokenizer for all languages is inferior to language-group specific tokenizers.** We test if using one joint tokenizer for SEA, AFR, and IND languages is better than using distinct tokenizers for each group. The joint tokenizer is trained on the combined tokenizer training corpora and has the same vocabulary size as the distinct ones. We relearn a single embedding layer with the joint tokenizer, followed by applying EMBSWAP with LoRA-Adaptation. As shown in Figure 9, although the joint tokenizer improves FLAN across all model scales, it is outperformed by distinct tokenizers. Smaller models show the biggest performance drop with the joint tokenizer. This happens because the shared vocabulary limits how much capacity is dedicated to each language, resulting in about 10% more tokens and making learning less effective.

## 4 Related Work

**Language Adaptation of LLMs.** Traditional ways of adapting LLMs to new languages involve continued pre-training on monolingual data (Cui et al., 2024; Zhao et al., 2024a) or multilingual instruction-tuning on synthetic data (Chen et al., 2024; Üstün et al., 2024; Aryabumi et al., 2024). These methods usually focus on monolingual transfer and a few high-resource languages. Moreover, their full-parameter tuning strategy is expensive, being prone to catastrophic forgetting. We achieve efficient adaptation to hundreds of languages through

<sup>9</sup>We suspect large vocabulary could increase ambiguity in output embeddings. Evidence also reveals that large vocabularies are not optimal for smaller LLMs (Tao et al., 2024).

embedding relearning and prevent knowledge forgetting with reusing parts of existing models.

**Tokenizer Adaptation and Embedding Relearning.** Previous tokenizer adaptation methods introduce languages into models through embedding relearning on new vocabularies (Artetxe et al., 2020; de Vries and Nissim, 2021; Marchisio et al., 2023; Chen et al., 2023). This increases flexibility in handling linguistic diversity and avoids over-segmentation that impairs task performance (Ahuja et al., 2023). However, they only consider encoder-only PLMs and monolingual transfer. We reevaluate this method for large-scale decoder-only LLMs with extension to hundreds of languages.

Many works focus on how to initialize the new embeddings of new vocabularies. They involve complex combinations (Dobler and de Melo, 2023; Liu et al., 2024) of overlapping tokens using auxiliary embeddings and external resources (Minixhofer et al., 2022; Remy et al., 2024), which limits their use across many languages. We avoid external resources by initializing the new embeddings with averaging the corresponding original embeddings (Gee et al., 2022; Mosin et al., 2023), making it scalable for hundreds of languages.

**Model Merging.** Model merging combines different specialized LLMs to create new ones with all their abilities (Wortsman et al., 2022). This can be done with simple arithmetic on existing parameters (Ilharco et al., 2023). We compared with one such method for instruction following across languages, ChatVector (Huang et al., 2024). Layer Swapping (Bandarkar et al., 2025) retains specific layers from language-adapted models while merging others from task fine-tuned models. Similar to Layer Swapping, EMBSWAP merges embeddings relearned on new vocabulary into instruction-tuned transformer body for cross-lingual transfer.

## 5 Conclusion

We tackled the challenge of adapting LLMs to diverse languages through embedding relearning. Our empirical findings reveal that embedding tuning with customized tokenizers contributes to effective language adaptation of LLMs, especially for low-resource languages with severe fragmentation. We also demonstrate that these embeddings can be integrated into any instruction-tuned LLMs to enable cross-lingual transfer with minimal training costs, outperforming or matching other strong baselines with high computational cost.

embedding relearning and prevent knowledge forgetting with reusing parts of existing models.	535
	536
<b>Tokenizer Adaptation and Embedding Relearning.</b> Previous tokenizer adaptation methods introduce languages into models through embedding relearning on new vocabularies (Artetxe et al., 2020; de Vries and Nissim, 2021; Marchisio et al., 2023; Chen et al., 2023). This increases flexibility in handling linguistic diversity and avoids over-segmentation that impairs task performance (Ahuja et al., 2023). However, they only consider encoder-only PLMs and monolingual transfer. We reevaluate this method for large-scale decoder-only LLMs with extension to hundreds of languages.	537
	538
	539
	540
	541
	542
	543
	544
	545
	546
	547
	548
Many works focus on how to initialize the new embeddings of new vocabularies. They involve complex combinations (Dobler and de Melo, 2023; Liu et al., 2024) of overlapping tokens using auxiliary embeddings and external resources (Minixhofer et al., 2022; Remy et al., 2024), which limits their use across many languages. We avoid external resources by initializing the new embeddings with averaging the corresponding original embeddings (Gee et al., 2022; Mosin et al., 2023), making it scalable for hundreds of languages.	549
	550
	551
	552
	553
	554
	555
	556
	557
	558
	559
<b>Model Merging.</b> Model merging combines different specialized LLMs to create new ones with all their abilities (Wortsman et al., 2022). This can be done with simple arithmetic on existing parameters (Ilharco et al., 2023). We compared with one such method for instruction following across languages, ChatVector (Huang et al., 2024). Layer Swapping (Bandarkar et al., 2025) retains specific layers from language-adapted models while merging others from task fine-tuned models. Similar to Layer Swapping, EMBSWAP merges embeddings relearned on new vocabulary into instruction-tuned transformer body for cross-lingual transfer.	560
	561
	562
	563
	564
	565
	566
	567
	568
	569
	570
	571
	572
We tackled the challenge of adapting LLMs to diverse languages through embedding relearning. Our empirical findings reveal that embedding tuning with customized tokenizers contributes to effective language adaptation of LLMs, especially for low-resource languages with severe fragmentation. We also demonstrate that these embeddings can be integrated into any instruction-tuned LLMs to enable cross-lingual transfer with minimal training costs, outperforming or matching other strong baselines with high computational cost.	573
	574
	575
	576
	577
	578
	579
	580
	581
	582
	583
	584

## 585 Limitations

586 Although embedding relearning involves only em-  
587 bedding tuning while keeping the transformer body  
588 frozen, it still incurs a computational cost of per-  
589 forming a full forward pass and back-propagation  
590 through all transformer layers. To address this,  
591 more efficient approaches, such as embedding tun-  
592 ing on a subset of transformer layers (Marchisio  
593 et al., 2023), could be employed to accelerate the  
594 training process.

595 We achieve substantial language adaptation by  
596 learning shared embeddings targeted at groups of  
597 geographically related languages, thereby avoiding  
598 the monolingual adaptation requiring the creation  
599 of language-specific embeddings. In addition, we  
600 demonstrate that constructing distinct embeddings  
601 for each language group outperforms the approach  
602 of treating all languages uniformly. However, this  
603 group-specific strategy may block certain forms of  
604 cross-lingual transfer and raises the critical ques-  
605 tion of how to best define ‘regions’ to facilitate the  
606 best transfer.

607 Although embedding relearning effectively ad-  
608 dresses catastrophic forgetting, we still observe  
609 a slight regression in English and several other  
610 high-resource languages (HRLs). This decline  
611 is primarily attributable to the fact that the orig-  
612 inal LLMs were already extensively pre-trained  
613 on these languages. Furthermore, the English and  
614 HRL data included in our multilingual dataset  $\mathcal{D}_{la}$   
615 is likely of lower quality compared to the orig-  
616 inal pre-training data. This limitation becomes par-  
617 ticularly pronounced in Gemma2, where the more  
618 advanced Gemini models (Anil et al., 2024) are em-  
619 ployed for knowledge distillation. We believe that  
620 such regressions in English and HRL performance  
621 could be alleviated with access to higher-quality  
622 data.

623 While we provide empirical results to demon-  
624 strate the effectiveness of embedding relearning  
625 and the importance of customized vocabularies tai-  
626 lored to target languages for cross-lingual transfer,  
627 future research should explore the conditions under  
628 which performing embedding relearning on LLMs  
629 is advantageous for cross-lingual transfer and the  
630 scenarios where it might be less effective.

631 There might also be safety concerns arising from  
632 the disruption of alignment abilities due to the  
633 changes to the embeddings, and how to align the  
634 model effectively and efficiently to reject unsafe  
635 queries after the EMBSWAP pipeline which would

need to be carefully addressed.

## 636 References

- Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo Kasai, David Mortensen, Noah Smith, and Yulia Tsvetkov. 2023. *Do all languages cost the same? tokenization in the era of commercial language models*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9904–9923, Singapore. Association for Computational Linguistics.
- Kabir Ahuja, Harshita Diddee, Rishav Hada, Milliecent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. 2023. *MEGA: Multilingual evaluation of generative AI*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4232–4267, Singapore. Association for Computational Linguistics.
- Duarte Miguel Alves, José Pombal, Nuno M Guerreiro, Pedro Henrique Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and Andre Martins. 2024. *Tower: An open multilingual large language model for translation-related tasks*. In *First Conference on Language Modeling*.
- Rohan Anil et al. 2023. *Palm 2 technical report*. Preprint, arXiv:2305.10403.
- Rohan Anil et al. 2024. *Gemini: A family of highly capable multimodal models*. Preprint, arXiv:2312.11805.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. *On the cross-lingual transferability of monolingual representations*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.
- Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. 2024. *Aya 23: Open weight releases to further multilingual progress*. Preprint, arXiv:2405.15032.
- Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. 2024. *The bele-bele benchmark: a parallel reading comprehension dataset in 122 language variants*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 749–775, Bangkok, Thailand. Association for Computational Linguistics.

692	Lucas Bandarkar, Benjamin Muller, Pritish Yuvraj, Rui Hou, Nayan Singhal, Hongjiang Lv, and Bing Liu.	751
693	2025. Layer swapping for zero-shot cross-lingual transfer in large language models.	752
694	In <i>The Thirteenth International Conference on Learning Representations</i> .	753
695		
696		
697		
698	Rachit Bansal, Bidisha Samanta, Siddharth Dalmia, Nitish Gupta, Sriram Ganapathy, Abhishek Bapna, Prateek Jain, and Partha Talukdar.	754
699	2024. LLM augmented LLMs: Expanding capabilities through composition.	755
700	In <i>The Twelfth International Conference on Learning Representations</i> .	756
701		757
702		758
703		759
704	Ankur Bapna, Isaac Caswell, Julia Kreutzer, Orhan Firat, Daan van Esch, Aditya Siddhant, Mengmeng Niu, Pallavi Baljekar, Xavier Garcia, Wolfgang Macherey, Theresa Breiner, Vera Axelrod, Jason Riesa, Yuan Cao, Mia Xu Chen, Klaus Macherey, Maxim Krikun, Pidong Wang, Alexander Gutkin, Apurva Shah, Yanping Huang, Zhifeng Chen, Yonghui Wu, and Macduff Hughes.	760
705	2022. Building machine translation systems for the next thousand languages.	761
706	<i>Preprint</i> , arXiv:2205.03983.	762
707		763
708		764
709		765
710	Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi.	766
711	2019. Piqa: Reasoning about physical commonsense in natural language.	767
712	<i>Preprint</i> , arXiv:1911.11641.	768
713		769
714	Terra Blevins and Luke Zettlemoyer.	770
715	2022. Language contamination helps explains the cross-lingual capabilities of English pretrained models.	771
716	In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3563–3574,	772
717	Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	773
718		774
719		775
720		
721	Tom Brown, Benjamin Mann, Nick Ryder, Melanie 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, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei.	776
722	2020. Language models are few-shot learners.	777
723	In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.	778
724		779
725		780
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739	Isaac Caswell, Theresa Breiner, Daan van Esch, and Ankur Bapna.	781
740	2020. Language ID in the wild: Unexpected challenges on the path to a thousand-language web text corpus.	782
741	In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6588–6608, Barcelona, Spain (Online). International Committee on Computational Linguistics.	783
742		784
743		785
744		786
745		
746		
747	Pinzhen Chen, Shaoxiong Ji, Nikolay Bogoychev, Andrey Kutuzov, Barry Haddow, and Kenneth Heafield.	787
748	2024. Monolingual or multilingual instruction tuning: Which makes a better alpaca.	788
749	In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 836–846, Online. Association for Computational Linguistics.	789
750		790
751		791
752		792
753		793
754		794
755		795
756		
757		
758		
759		
760		
761		
762		
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		

807	specialization of multilingual models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13440–13454, Singapore. Association for Computational Linguistics.	865
808		866
809		867
810		868
811		869
812	Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. <b>DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs.</b> In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.	870
813		871
814		872
815		873
816		
817		
818		
819		
820		
821		
822	Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. 2024. <b>Continual pre-training for cross-lingual LLM adaptation: Enhancing Japanese language capabilities.</b> In <i>First Conference on Language Modeling</i> .	874
823		881
824		
825		
826		
827		
828	Philip Gage. 1994. A new algorithm for data compression. <i>C Users J.</i> , 12(2):23–38.	879
829		880
830		881
831	Xavier Garcia, Noah Constant, Ankur Parikh, and Orhan Firat. 2021. <b>Towards continual learning for multilingual machine translation via vocabulary substitution.</b> In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1184–1192, Online. Association for Computational Linguistics.	882
832		883
833		884
834		885
835		886
836		887
837		888
838		889
839		
840		
841		
842		
843		
844		
845	Leonidas Gee, Andrea Zugarini, Leonardo Rigutini, and Paolo Torroni. 2022. <b>Fast vocabulary transfer for language model compression.</b> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 409–416, Abu Dhabi, UAE. Association for Computational Linguistics.	890
846		891
847		892
848		893
849		894
850		895
851		896
852	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. <b>The Flores-101 evaluation benchmark for low-resource and multilingual machine translation.</b> <i>Transactions of the Association for Computational Linguistics</i> , 10:522–538.	897
853		898
854		899
855		900
856		901
857	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. <b>Measuring massive multitask language understanding.</b> In <i>International Conference on Learning Representations</i> .	902
858		903
859		
860		
861		
862		
863		
864		
865	Shih-Cheng Huang, Pin-Zu Li, Yu-chi Hsu, Kuang-Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard Tsai, and Hung-yi Lee. 2024. <b>Chat vector: A simple approach to equip LLMs with instruction following and model alignment in new languages.</b> In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10943–10959, Bangkok, Thailand. Association for Computational Linguistics.	909
866		910
867		911
868		912
869		913
870		914
871		915
872		
873		
874	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. <b>Editing models with task arithmetic.</b> In <i>The Eleventh International Conference on Learning Representations</i> .	874
875		875
876		876
877		877
878		878
879	Diederik P. Kingma and Jimmy Ba. 2014. <b>Adam: A method for stochastic optimization.</b> <i>Preprint</i> , arXiv:1412.6980.	879
880		880
881		881
882	Davis Liang, Hila Gonen, Yuning Mao, Rui Hou, Naman Goyal, Marjan Ghazvininejad, Luke Zettlemoyer, and Madian Khabsa. 2023. <b>XLM-V: Overcoming the vocabulary bottleneck in multilingual masked language models.</b> In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13142–13152, Singapore. Association for Computational Linguistics.	882
883		883
884		884
885		885
886		886
887		887
888		888
889		889
890	Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schuetze. 2024. <b>OFA: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining.</b> In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1067–1097, Mexico City, Mexico. Association for Computational Linguistics.	890
891		891
892		892
893		893
894		894
895		895
896		896
897	Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. <b>LLaMAX: Scaling linguistic horizons of LLM by enhancing translation capabilities beyond 100 languages.</b> In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10748–10772, Miami, Florida, USA. Association for Computational Linguistics.	897
898		898
899		899
900		900
901		901
902		902
903		903
904	Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2024. <b>An empirical study of catastrophic forgetting in large language models during continual fine-tuning.</b> <i>Preprint</i> , arXiv:2308.08747.	904
905		905
906		906
907		907
908		908
909	Kelly Marchisio, Patrick Lewis, Yihong Chen, and Mikel Artetxe. 2023. <b>Mini-model adaptation: Efficiently extending pretrained models to new languages via aligned shallow training.</b> In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5474–5490, Toronto, Canada. Association for Computational Linguistics.	909
910		910
911		911
912		912
913		913
914		914
915		915
916	Benjamin Minixhofer, Fabian Paischer, and Navid Rekabsaz. 2022. <b>WECHSEL: Effective initialization of subword embeddings for cross-lingual transfer of monolingual language models.</b> In <i>Proceedings of</i>	916
917		917
918		918
919		919

920	<i>the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3992–4006, Seattle, United States. Association for Computational Linguistics.	978
921		979
922		980
923		981
924		982
925		983
926	Vladislav Mosin, Igor Samenko, Borislav Kozlovskii, Alexey Tikhonov, and Ivan P. Yamshchikov. 2023. <i>Fine-tuning transformers: Vocabulary transfer</i> . <i>Artificial Intelligence</i> , 317:103860.	
927		
928		
929	François Remy, Pieter Delobelle, Hayastan Avetisyan, Alfiya Khabibullina, Miryam de Lhoneux, and Thomas Demeester. 2024. Trans-tokenization and cross-lingual vocabulary transfers: Language adaptation of LLMs for low-resource NLP. In <i>First Conference on Language Modeling</i> .	
930		
931		
932		
933		
934		
935	Morgane Riviere et al. 2024. Gemma 2: Improving open language models at a practical size. <i>Preprint</i> , arXiv:2408.00118.	
936		
937		
938	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: an adversarial winograd schema challenge at scale. <i>Commun. ACM</i> , 64(9):99–106.	
939		
940		
941		
942	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	
943		
944		
945		
946		
947		
948		
949	Freida Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multilingual chain-of-thought reasoners. In <i>The Eleventh International Conference on Learning Representations</i> .	
950		
951		
952		
953		
954		
955		
956	Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyuan Wang, Yibin Wang, Zifeng Wang, Sayna Ebrahimi, and Hao Wang. 2024. Continual learning of large language models: A comprehensive survey. <i>Preprint</i> , arXiv:2404.16789.	
957		
958		
959		
960		
961	Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Dinenesh Tewari, and Partha Talukdar. 2024. IndicGen-Bench: A multilingual benchmark to evaluate generation capabilities of LLMs on Indic languages. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11047–11073, Bangkok, Thailand. Association for Computational Linguistics.	
962		
963		
964		
965		
966		
967		
968		
969	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dong-dong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. Language-specific neurons: The key to multilingual capabilities in large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5701–5715, Bangkok, Thailand. Association for Computational Linguistics.	
970		
971		
972		
973		
974		
975		
976		
977		
510	Chaofan Tao, Qian Liu, Longxu Dou, Niklas Muenighoff, Zhongwei Wan, Ping Luo, Min Lin, and Ngai Wong. 2024. Scaling laws with vocabulary: Larger models deserve larger vocabularies. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
511		
512	Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Dara Bahri, Tal Schuster, Steven Zheng, Denny Zhou, Neil Houlsby, and Donald Metzler. 2023. UL2: Unifying language learning paradigms. In <i>The Eleventh International Conference on Learning Representations</i> .	
513		
514	Ahmet Üstün, Viraat Aryabumi, Zheng 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. Aya model: An instruction finetuned open-access multilingual language model. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15894–15939, Bangkok, Thailand. Association for Computational Linguistics.	
515		
516	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. Finetuned language models are zero-shot learners. In <i>International Conference on Learning Representations</i> .	
517		
518	Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Do llamas work in English? on the latent language of multilingual transformers. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15366–15394, Bangkok, Thailand. Association for Computational Linguistics.	
519		
520	Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmom, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 23965–23998. PMLR.	
521		
522	Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024. A paradigm shift in machine translation: Boosting translation performance of large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	
523		
524	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In	
525		
526		
527		
528		
529		
530		
531		
532		
533		
534		
535		
536		
537		
538		
539		
540		
541		
542		
543		
544		
545		
546		
547		
548		
549		
550		
551		
552		
553		
554		
555		
556		
557		
558		
559		
560		
561		
562		
563		
564		
565		
566		
567		
568		
569		
570		
571		
572		
573		
574		
575		
576		
577		
578		
579		
580		
581		
582		
583		
584		
585		
586		
587		
588		
589		
590		
591		
592		
593		
594		
595		
596		
597		
598		
599		
600		
601		
602		
603		
604		
605		
606		
607		
608		
609		
610		
611		
612		
613		
614		
615		
616		
617		
618		
619		
620		
621		
622		
623		
624		
625		
626		
627		
628		
629		
630		
631		
632		
633		
634		
635		
636		
637		
638		
639		
640		
641		
642		
643		
644		
645		
646		
647		
648		
649		
650		
651		
652		
653		
654		
655		
656		
657		
658		
659		
660		
661		
662		
663		
664		
665		
666		
667		
668		
669		
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		
756		
757		
758		
759		
760		
761		
762		
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		
810		
811		
812		
813		
814		
815		
816		
817		
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828		
829		
830		
831		
832		
833		
834		
835		
836		
837		
838		
839		
840		
841		
842		
843		
844		
845		
846		
847		
848		
849		
850		
851		
852		
853		
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		
864		
865		
866		
867		
868		
869		
870		
871		
872		
873		
874		
875		
876		
877		
878		
879		
880		
881		
882		
883		
884		
885		
886		
887		
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		

- 1036        *Proceedings of the 2021 Conference of the North*  
1037        *American Chapter of the Association for Computa-*  
1038        *tional Linguistics: Human Language Technologies*,  
1039        pages 483–498, Online. Association for Compu-  
1040        tational Linguistics.
- 1041        Xiang Yue, Tianyu Zheng, Ge Zhang, and Wenhui Chen.  
1042        2024. **MAmmoTH2: Scaling instructions from the**  
1043        **web**. In *The Thirty-eighth Annual Conference on*  
1044        *Neural Information Processing Systems*.
- 1045        Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali  
1046        Farhadi, and Yejin Choi. 2019. **HellaSwag: Can**  
1047        **a machine really finish your sentence?** In *Pro-*  
1048        *ceedings of the 57th Annual Meeting of the Asso-*  
1049        *ciation for Computational Linguistics*, pages 4791–  
1050        4800, Florence, Italy. Association for Computational  
1051        Linguistics.
- 1052        Shiyue Zhang, Vishrav Chaudhary, Naman Goyal,  
1053        James Cross, Guillaume Wenzek, Mohit Bansal, and  
1054        Francisco Guzman. 2022. **How robust is neural ma-**  
1055        **chine translation to language imbalance in multilin-**  
1056        **gual tokenizer training?** In *Proceedings of the 15th*  
1057        *biennial conference of the Association for Machine*  
1058        *Translation in the Americas (Volume 1: Research*  
1059        *Track)*, pages 97–116, Orlando, USA. Association  
1060        for Machine Translation in the Americas.
- 1061        Jun Zhao, Zhihao Zhang, Luhui Gao, Qi Zhang, Tao  
1062        Gui, and Xuanjing Huang. 2024a. **Llama beyond**  
1063        **english: An empirical study on language capability**  
1064        **transfer**. *Preprint*, arXiv:2401.01055.
- 1065        Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji  
1066        Kawaguchi, and Lidong Bing. 2024b. **How do large**  
1067        **language models handle multilingualism?** In *The*  
1068        *Thirty-eighth Annual Conference on Neural Informa-*  
1069        *tion Processing Systems*.
- 1070        Wenzhen Zheng, Wenbo Pan, Xu Xu, Libo Qin, Li Yue,  
1071        and Ming Zhou. 2024. **Breaking language barriers:**  
1072        **Cross-lingual continual pre-training at scale**. In *Pro-*  
1073        *ceedings of the 2024 Conference on Empirical Meth-*  
1074        *ods in Natural Language Processing*, pages 7725–  
1075        7738, Miami, Florida, USA. Association for Com-  
1076        putational Linguistics.

## 1077 Overview of Appendix

1078 Our supplementary includes the following sections:  
 1079 • Section **A**: Experimental Settings, including  
 1080 implementation details for training and evalua-  
 1081 tion.  
 1082 • Section **B**: Detailed result comparison be-  
 1083 tween instruction-tuned models and EMB-  
 1084 SWAP in English tasks.  
 1085 • Section **C**: Additional Language Adaptation  
 1086 and EMBSWAP results on PaLM2, Aya23 and  
 1087 Gemma2-IT models.  
 1088 • Section **D**: Supplementary analysis including  
 1089 additional data ablations.  
 1090 • Section **E**: Results by language for both lan-  
 1091 guage adaptation and EMBSWAP.

## 1092 A Experimental Settings

### 1093 A.1 Training Details

1094 For language adaptation, we use a constant learning  
 1095 rate of  $1 \times 10^{-4}$  for PaLM2 and  $1 \times 10^{-5}$  for Gemma2  
 1096 and Aya23 with the Adam optimizer ([Kingma and](#)  
 1097 [Ba, 2014](#)). The embeddings are trained on a total of  
 1098 200B tokens,<sup>10</sup> with each batch consisting of exam-  
 1099 ples packed to a maximum sequence length of 2K  
 1100 for PaLM2 and 8K for Gemma2 and Aya23. We pre-  
 1101 train the model using the UL2 objectives ([Tay et al.,](#)  
 1102 [2023](#)) for PaLM2 and causal language modeling ob-  
 1103 jectives for Gemma2 and Aya23. Language adap-  
 1104 tation consumes up to 256 TPU-v5 chips for the  
 1105 largest Gemma2-27B and Aya23-35B models. The  
 1106 batch size is selected based on the model size and  
 1107 computing resources we have, with batch sizes of  
 1108 256, 128, and 64 assigned to the Gemma2-2B, 9B,  
 1109 and 27B models, respectively. A similar strategy  
 1110 is applied to the Aya23 models. For PaLM2 mod-  
 1111 els, we use the same batch size of 2048 for all  
 1112 variants. We choose the best checkpoint based on  
 1113 the performance of FLORES-200 development sets  
 1114 corresponding to each target language group. The  
 1115 training time varies with model size, where smaller  
 1116 models complete training within 24 hours while  
 1117 larger models require up to 1 week to finish.

1118 We instruction-tune the transformer body of  
 1119 LLMs on  $\mathcal{D}_{it}$  using the same hyper-parameter  
 1120 setting to obtain XX-FLAN, where we sample up  
 1121 to 200M instances from the FLAN mixture to  
 1122 construct  $\mathcal{D}_{it}$ . We use early stopping to select  
 1123 the best model based on the performance on  
 1124 MMLU ([Hendrycks et al., 2021](#)) and assemble it

Dataset	Prompt
BELEBELE	The following are multiple choice questions (with answers).  Passage: [Target Language Passage] Question: [Target Language Question] (A) [Choice A] (B) [Choice B] (C) [Choice C] (D) [Choice D] Answer:
SIB-200	[News Article in Target Language] Question: What label best describes this news article? (A) science/technology (B) travel (C) politics (D) sports (E) health (F) entertainment (G) geography Answer:
FLORES-200	Translate this from English to [Target Language Name]:  English: [Sentence in English] [Target Language Name]:
XSUM-IN	I will first show a news article in English and then provide a one sentence summary of it in [Target Language Name].  Summarize the following article: [Article in English] Summary in [Target Language Name]:
XORQA-IN	Generate an answer in [Target Language Name/English] for the question based on the given passage:  [Passage in English] Q: [Question in Target Language] A:
GSM8K-NTL	Q: [Question in Target Language] A: [Let's think step by step.]

1095 Table 3: Prompt templates used in each of the evalua-  
 1096 tion dataset. For few-shot evaluation, n-shot examples  
 1097 have the same format as the last test instance, which  
 1098 comes after the preamble but before the test instances.

1100 with customized embedding to obtain EMBSWAP  
 1101 models. We observe all LLMs converge fast, as in-  
 1102 dicated by the average performance on MMLU. In  
 1103 most cases, training is completed within 24 hours.

1105 For LoRA-Adaptation, we add LoRA weights  
 1106 to the self-attention module of all transformer lay-  
 1107 ers with a LoRA rank of 64 and exclusively op-  
 1108 timize these weights.<sup>11</sup> We use a learning rate of  
 1109  $5 \times 10^{-6}$  for all models with 10% steps of warm-up.  
 1110 Analogous to embedding tuning, the FLORES-200  
 1111 development set is used for model selection. The  
 1112 training process is computationally efficient, com-  
 1113 pleting within 12 hours even for the largest model.

<sup>10</sup>English accounts for 30.8% of the tokens.

<sup>11</sup>This adds less than 1% parameters.

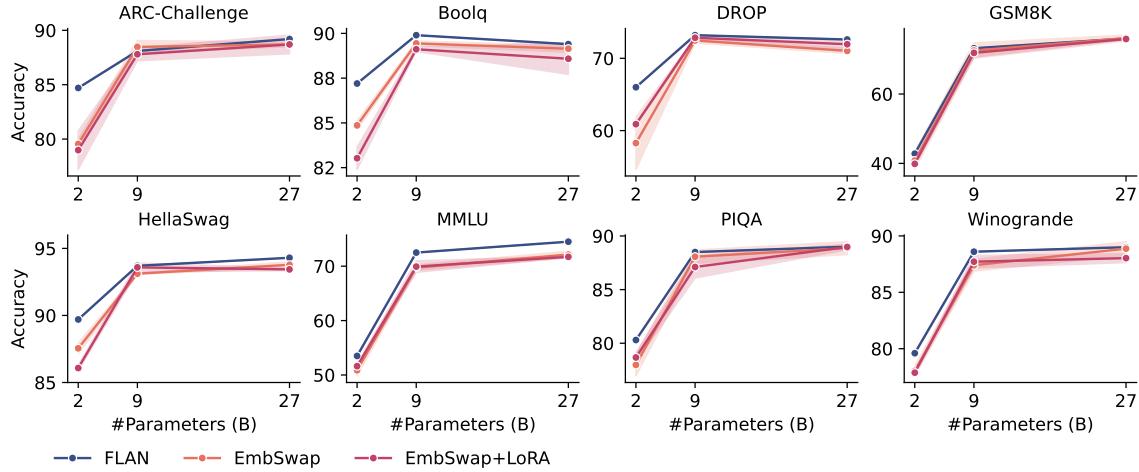


Figure 10: Performance on English tasks. For EMBSWAP methods, we present the averaged performance with variances across three embeddings (i.e., SEA, AFR, IND). Shading indicates the standard deviations measured over three embeddings.

## A.2 Evaluation Details

We use the prompt formats listed in Table 3 for evaluation. For generation tasks, greedy decoding is employed with a maximum sequence length of 256 tokens. For classification tasks, we calculate the logits of each available option (e.g., (A), (B)) and select the option with the highest score as the predicted answer.

## B Details on English Benchmarking

For benchmarking performance in English, we choose ARC-Challenge (Clark et al., 2018), Boolq (Clark et al., 2019), DROP (Dua et al., 2019), GSM8K (Cobbe et al., 2021), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), PIQA (Bisk et al., 2019), and Winogrande (Sakaguchi et al., 2021). We use different number of shots for evaluation following Riviere et al. (2024).

Figure 10 shows the comparison between the instruction-tuned FLAN model and EMBSWAP across various sizes of Gemma2. Both EMBSWAP variants exhibit performance regressions across all tasks compared to the FLAN model, although the performance gap closes as model capacity scales up. We believe that these minor regressions are justifiable in light of the substantial gains achieved in zero-shot cross-lingual transfer.

## C Additional Main Results

### C.1 Language Adaptation on PaLM2 and Aya23

We evaluate the generalization ability of language adaptation on two LLMs with varying levels of multilingualism. Among them, PaLM2-S exhibits the

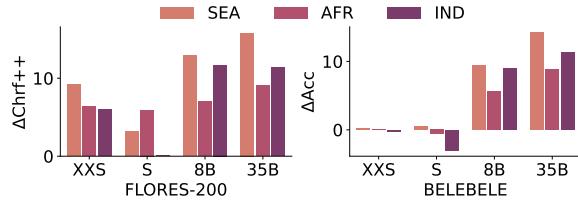


Figure 11: Language Adaptation on PaLM2 (XXS, S) and Aya23 (8B, 35B). Absolute gains over the pre-trained models are reported.

strongest multilingual abilities while Aya23 models demonstrate limited multilingual performance. Figure 11 shows that the performance gains from language adaptation decrease as the original multilingual scope of the LLMs expand (Aya23 → PaLM2-XXS → PaLM2-S). Moreover, larger performance improvements are observed on Aya23 models when scaling up their size, suggesting that language adaptation may be particularly effective for models with stronger English proficiency.

### C.2 EMBSWAP Results on PaLM2, Aya23 and Gemma2-IT

Employing the same EMBSWAP pipeline, we observe similar patterns of performance improvement across all sizes of the PaLM2 and Aya23 models, as shown in Table 4. This demonstrates the broad applicability of EMBSWAP to various types of LLMs. In addition, we also show the detailed results of Gemma2-IT models. EMBSWAP also performs effectively with off-the-shelf IT models that have undergone complex supervised fine-tuning and reinforcement learning. This further underscores the versatility of EMBSWAP in integrating pre-trained multilingual embeddings into LLMs that have been instruction-tuned using diverse methodologies for

Task Type	ENGLISH	CLASSIFICATION						GENERATION							
Eval. Metric	ENGLISH	BELEBELE Accuracy			SIB-200 Accuracy			FLORES-200 ChrF++			XORQA-IN Token-F1		XSUM-IN ChrF		Avg.
Model		SEA	AFR	IND	SEA	AFR	IND	SEA	AFR	IND	IND	EN	IND	EN	
PaLM2-XXS-FLAN	59.7	54.5	40.0	49.4	67.2	51.3	70.9	28.5	13.9	21.0	9.7	42.1	<b>4.8</b>	33.7	32.8
PaLM2-XXS-FA	<b>62.6</b>	58.0	40.2	50.9	<b>76.8</b>	<u>62.2</u>	<b>77.2</b>	<u>27.7</u>	<u>12.2</u>	<u>17.9</u>	10.4	59.4	<u>2.1</u>	<b>34.2</b>	35.9
★ LoRA-Adapt	58.9	<b>59.7</b>	<b>43.4</b>	<b>53.8</b>	<u>73.7</u>	<b>63.5</b>	<u>76.9</u>	<b>32.3</b>	<b>16.0</b>	<b>26.7</b>	11.7	<b>64.2</b>	<u>1.0</u>	<b>30.7</b>	<b>37.7</b>
PaLM2-S-FLAN	<b>86.3</b>	80.2	67.4	<b>82.1</b>	70.6	60.9	74.5	37.6	19.3	36.7	<b>21.3</b>	46.9	15.8	41.1	45.4
PaLM2-S-FA	85.1	83.2	<u>67.1</u>	<b>79.9</b>	77.0	66.9	74.6	<u>36.1</u>	<u>23.1</u>	<u>34.6</u>	<b>21.1</b>	51.8	<b>17.3</b>	<b>40.1</b>	47.1
★ LoRA-Adapt	85.7	<b>83.6</b>	<b>68.4</b>	<u>81.8</u>	<u>77.7</u>	<b>68.7</b>	<u>77.8</u>	<u>39.3</u>	<b>24.1</b>	<u>38.7</u>	18.9	<b>57.3</b>	<u>15.2</u>	<b>41.2</b>	<b>48.2</b>
Aya23-8B-FLAN	<b>74.3</b>	47.8	34.1	42.7	64.3	48.5	60.6	22.9	6.0	16.1	<u>10.5</u>	<u>58.5</u>	<b>11.7</b>	<b>34.5</b>	32.3
Aya23-8B-FA	71.1	<b>56.6</b>	<u>38.4</u>	<u>48.0</u>	<b>72.4</b>	<b>58.7</b>	<b>71.3</b>	<u>28.8</u>	<u>9.1</u>	<u>11.5</u>	<u>7.8</u>	<b>61.6</b>	<u>3.9</u>	<b>25.7</b>	34.2
★ LoRA-Adapt	71.3	<b>59.1</b>	<b>39.0</b>	<b>50.8</b>	65.6	55.3	65.2	<u>37.1</u>	<b>16.9</b>	33.4	<u>9.6</u>	<u>56.5</u>	<u>9.4</u>	<u>30.7</u>	<b>36.8</b>
Aya23-35B-FLAN	<b>83.6</b>	56.9	<b>40.3</b>	<u>54.6</u>	67.0	52.9	67.2	27.1	9.1	23.6	<u>11.4</u>	<u>64.7</u>	<u>10.4</u>	<b>37.9</b>	36.5
Aya23-35B-FA	79.2	57.6	<b>40.2</b>	<u>53.0</u>	70.8	<u>53.8</u>	68.5	<u>27.7</u>	<u>9.5</u>	<u>13.9</u>	<u>6.8</u>	<u>62.3</u>	<u>7.3</u>	<u>22.1</u>	<b>35.2</b>
★ LoRA-Adapt	82.0	<b>72.6</b>	<b>50.1</b>	<b>65.2</b>	<b>75.2</b>	<b>64.7</b>	<b>74.0</b>	<u>41.6</u>	<u>23.4</u>	<u>40.1</u>	<u>12.6</u>	<u>67.8</u>	<u>15.7</u>	<u>37.1</u>	<b>45.1</b>
Gemma2-2B-IT + Lang-Adapt	<u>55.5</u>	47.5	35.4	45.8	57.4	42.8	61.2	24.2	7.3	17.7	<u>14.9</u>	53.0	<u>10.1</u>	<u>31.9</u>	31.6
Gemma2-2B-FA	54.1	51.2	36.7	42.9	<u>67.0</u>	<u>48.7</u>	<u>67.1</u>	<u>35.4</u>	<u>14.2</u>	<u>29.2</u>	12.3	40.2	<u>7.1</u>	<u>27.5</u>	33.0
★ LoRA Adapt	55.1	<u>53.7</u>	<u>37.4</u>	<u>46.0</u>	66.7	<u>49.5</u>	64.0	27.4	11.7	27.0	<b>16.7</b>	<u>58.7</u>	<b>13.4</b>	<b>32.0</b>	<u>35.9</u>
Gemma2-9B-IT + Lang-Adapt	<b>79.9</b>	71.5	51.4	72.5	72.5	57.8	79.3	32.2	14.0	31.1	<b>24.0</b>	64.4	<u>16.6</u>	<u>34.3</u>	44.5
Gemma2-9B-FA	73.9	75.4	<b>57.9</b>	<b>72.7</b>	77.2	<u>69.2</u>	<u>79.7</u>	<u>39.3</u>	<u>22.2</u>	<u>38.7</u>	23.1	<b>69.6</b>	<u>17.0</u>	<u>33.9</u>	48.3
★ LoRA-Adapt	78.5	<b>77.4</b>	54.5	<u>72.6</u>	<u>78.2</u>	66.1	<u>78.5</u>	33.8	16.6	<u>18.8</u>	21.4	65.2	<u>9.6</u>	<u>34.2</u>	<u>43.8</u>
Gemma2-27B-IT + Lang-Adapt	79.5	76.5	<b>58.1</b>	<u>70.8</u>	<b>81.3</b>	<b>71.2</b>	<b>82.8</b>	<u>41.5</u>	<u>25.1</u>	<u>39.8</u>	18.2	<u>68.2</u>	<u>15.7</u>	<u>33.0</u>	<b>48.4</b>
Gemma2-27B-FA	82.1	74.9	54.6	<u>75.9</u>	76.2	63.1	<b>83.0</b>	35.2	20.1	34.5	<b>27.7</b>	68.2	<b>19.0</b>	<u>34.7</u>	48.0
★ LoRA-Adapt	81.9	<b>81.7</b>	<b>62.7</b>	<b>76.0</b>	<b>82.1</b>	<u>70.0</u>	81.9	<b>42.5</b>	<u>24.8</u>	<u>39.3</u>	<u>26.4</u>	<b>73.0</b>	<u>18.3</u>	<b>35.8</b>	<b>51.1</b>
Gemma2-27B-FA	79.9	80.8	59.5	<u>75.9</u>	81.4	68.9	<b>83.0</b>	<u>35.1</u>	<u>6.2</u>	<u>29.9</u>	20.0	<u>68.8</u>	<u>15.0</u>	<u>34.6</u>	<u>46.7</u>
★ LoRA-Adapt	<b>82.2</b>	78.8	60.1	<u>74.3</u>	81.7	<b>71.7</b>	<u>82.6</u>	42.3	<b>25.6</b>	<b>40.6</b>	24.2	<u>67.9</u>	<u>15.6</u>	<u>34.0</u>	49.6

Table 4: Additional EMBSWAP results on PaLM2, Aya23 and Gemma2-IT models. The best and second-best results are marked in **bold** and underlined. Red values indicate EMBSWAP hurts the performance.

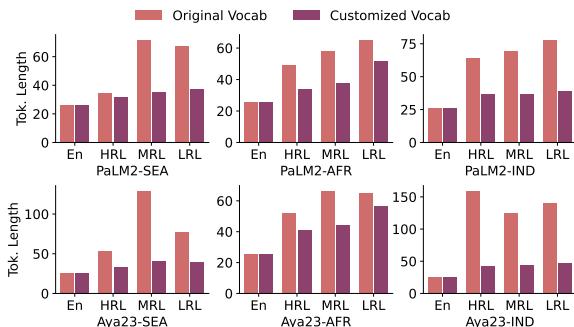


Figure 12: The tokenization comparison between using the vanilla and customized multilingual tokenizers on Gemma2. Tok. Length refers to the average number of tokens required to represent the same amount of texts.

efficient zero-shot cross-lingual transfer.

## D Supplementary Analysis

**Additional tokenization results on the customized vocabularies.** We present additional results on tokenization fertility using customized vocabularies developed for PaLM2 and Aya23. Figure 12 shows similar patterns as those reported for Gemma2, where the fertility tokenization for both MRLs and LRLs shows a substantial decrease, while the English tokenization remains roughly unchanged. This effect is particularly pronounced

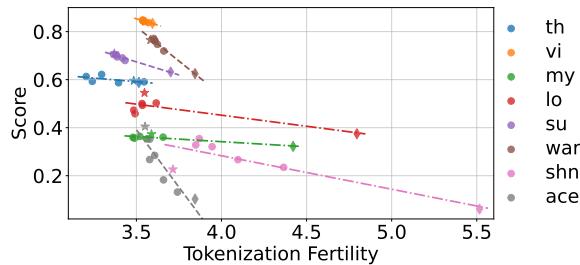


Figure 13: Correlation between the performance of language adaptation on PaLM2-XXS with tokenizer fertility. Normalized ChrF++ on FLORES-SEA are reported. ♦ and ★ indicate the original and customized tokenizers in PaLM2.

in Aya23, where we observe over  $\times 3$  reduction in fertility for SEA and IND languages. In Table 6, we also show a few tokenized examples for low-resource languages. We find that the customized tokenizer produces more meaningful tokens and avoids overtokenization.

**Tokenizer fertility is also inversely correlated to downstream performance in PaLM2.** We replicate the analysis shown in Figure 7 using PaLM2-XXS, with results shown in Figure 13. We observe patterns analogous to those reported for Gemma2-2B, wherein reduced tokenizer fertility is generally associated with improved downstream performance

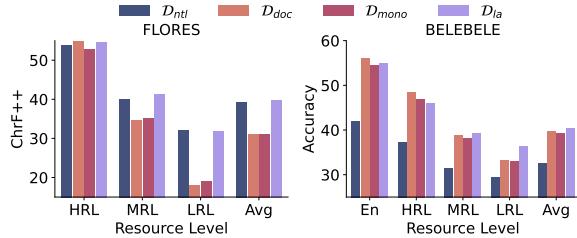


Figure 14: Pre-training data ablation for language adaptation.  $\mathcal{D}_{ntl}$ : multilingual data sampled from the NTL corpus;  $\mathcal{D}_{doc}$ : multilingual data sampled Wikipedia and mC4;  $\mathcal{D}_{la}$ : our final data mixture for language adaptation, i.e.,  $\mathcal{D}_{la} = \mathcal{D}_{ntl} \cup \mathcal{D}_{doc}$ ;  $\mathcal{D}_{mono}$ : monolingual data by excluding parallel sentences from  $\mathcal{D}_{la}$ . SEA languages results on Gemma2-2B are reported.

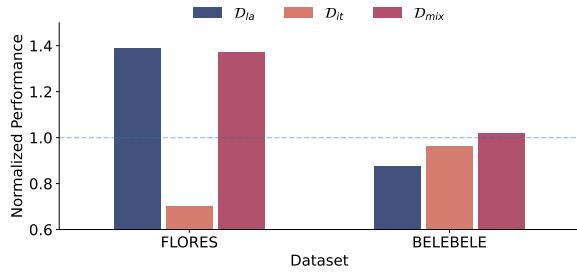


Figure 15: Training data ablation for Lora-Adaptation.  $\mathcal{D}_{la}$ : multilingual data for embedding tuning;  $\mathcal{D}_{it}$ : the FLAN mixture for instruction-tuning;  $\mathcal{D}_{mix}$ : a combination of equally sub-sampled  $\mathcal{D}_{la}$  and  $\mathcal{D}_{it}$ . Averaged SEA language results (normalized score v.s. EMBSWAP) on Gemma2-2B are reported.

and this relationship is particularly pronounced in LRLs and languages written in non-Latin script. **Long-tail NTL and document-level data are both important for language adaptation.** We perform language adaptation on Gemma2-2B in SEA languages with different data mixture. Figure 14 reveals that improved performance in long-tail languages primarily stems from the inclusion of NTL data, while document-level data plays a crucial role in preserving knowledge for high-resource languages. The significance of incorporating document-level data is further underscored by the results on BELEBELE, where the removal of  $\mathcal{D}_{doc}$  leads to a substantial performance decline across all resource levels. This finding highlights the importance of  $\mathcal{D}_{doc}$  in maintaining the ability of LLMs in processing various types of texts. In addition, excluding parallel data from the training mixture (i.e.,  $\mathcal{D}_{mono}$ ) leads to a significant decline in performance on translation tasks. Moreover, the performance of LRLs on BELEBELE also declines substantially, indicating the critical role of parallel data in enhancing tasks beyond translation through facilitated cross-lingual transfer (Anil et al., 2023).

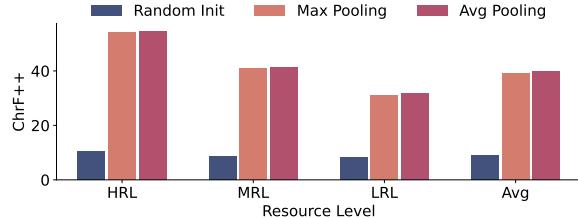


Figure 16: Ablations on embedding initialization methods. FLORES-SEA language performance of language adaptation on Gemma2-2B is reported. Max Pooling: for each new token in the customized vocabulary, we use the original tokenizer to tokenize it and apply max pooling over the embeddings of the corresponding subtokens as the initialization.

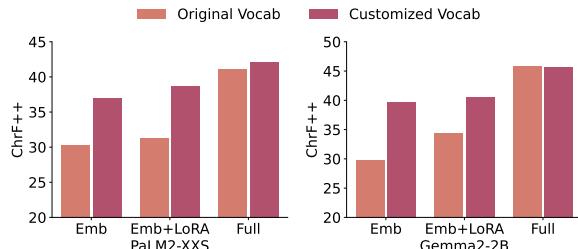


Figure 17: The effects of using customized vocabularies with different proportions of tuned parameters. The averaged score on SEA languages of FLORES-200 is reported.

### Both multilingual and instruction data are important for LoRA-Adaptation

We investigate the impact of employing multilingual  $\mathcal{D}_{la}$  and instruction-tuning data  $\mathcal{D}_{it}$  in LoRA-Adaptation. As shown in Figure 15, the removal of either  $\mathcal{D}_{la}$  or  $\mathcal{D}_{it}$  harms the performance compared to the vanilla EMBSWAP model on FLORES-200 or BELEBELE, while combining both datasets enables the adapted LLM to achieve the best overall results.

### Employing the original embeddings for initialization is essential to language adaptation.

In Figure 16, we show that without initializing the customized embeddings using the original embeddings from the LLM, the language adaptation does not perform well at all, yielding ChrF++ scores below 10 even for HRLs. In contrast, initializing with the original embeddings significantly improves the effectiveness of language adaptation, where the average pooling method slightly outperforms the max pooling variant.

### The benefits of employing customized vocabulary decrease with more tuned parameters.

We study the effects of employing customized embeddings with varying ratios of tuned parameters. As shown in Figure 17, the benefits of using customized embeddings diminish as the number of tuned parameters increases (Emb → Emb+LoRA

1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269

Size	Vocab	HRL	MRL	LRL
2B	Original	59.1	55.6	57.1
	Custom	61.1 $\uparrow$ 3.38%	61.2 $\uparrow$ 10.07%	60.3 $\uparrow$ 5.60%
9B	Original	30.2	27.2	28.6
	Custom	33.8 $\uparrow$ 11.92%	33.3 $\uparrow$ 22.43%	32.1 $\uparrow$ 12.24%
27B	Original	17.1	15.0	16.0
	Custom	19.6 $\uparrow$ 14.62%	19.2 $\uparrow$ 28.0%	18.4 $\uparrow$ 15.0%

Table 5: Comparing latency for using original and customized vocabularies in EMBSWAP. The number of instances processed per second (i.e., prefilling) by Gemma2 are reported. We use the passages in BELEBELE for all SEA, AFR, and IND languages.

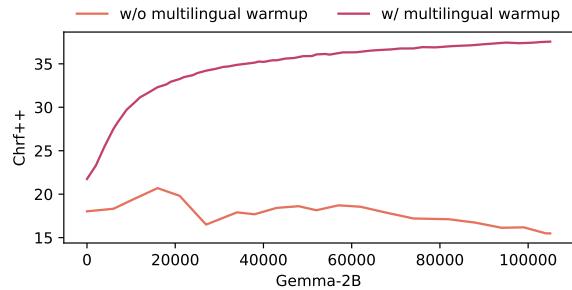


Figure 18: Learning curves of language adaptation on LLMs with limited multilingual abilities. Averaged results on SEA languages of FLORES-200 are reported.

→ Full). Specifically, on Gemma2-2B model, customized embeddings show no advantage when full-parameter tuning is employed. This phenomenon arises because increasing the number of tuned parameters allocates greater model capacity for language adaptation, which simplifies the adaptation process compared to relying solely on embedding tuning. Nonetheless, full-parameter tuning could exacerbate the problem of catastrophic forgetting, while embedding tuning provides a safer alternative. Furthermore, the use of customized embeddings amplifies the advantages of embedding tuning, making it a promising technique.

**Customized vocabulary improves latency.** We evaluate latency by measuring the number of texts processed per second by LLMs. We use the passages from all SEA, AFR, and IND languages in BELEBELE as test texts. For comparison, we test EMBSWAP models integrated with embeddings trained using either the original Gemma2 tokenizer or our customized one. Table 5 shows that employing customized tokenizer consistently improves latency, particularly in MRLs and LRLs. This trend becomes increasingly pronounced as model size scales, highlighting the importance of customized tokenizers in achieving low-latency processing for long-tail languages in larger LLMs.

**Multilingual warmup is necessary for embed-**

**ding relearning on LLMs with limited multilingual abilities.** We investigate whether embedding tuning with customized embeddings is a universal technique for the language adaptation of any types of LLMs. As shown in Figure 18, simply performing embedding tuning on Gemma-2B, a model with very limited multilingual capabilities, does not successfully adapt it to various languages. In contrast, when a multilingual continued pre-training process is conducted prior to embedding tuning, where the document-level data  $\mathcal{D}_{doc}$  is used to warm up the LLM, we observe consistent improvements throughout the training process. This suggests that, for LLMs, a good initial multilingual ability is essential for the success of embedding relearning.

## E Results by Language

Table 7 presents an overview of languages available across our evaluation benchmarks. We show the per-language results on each task for both language adaptation (Tables 8 – 10) and EMBSWAP (Tables 11 – 16).

1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313

1314  
1315  
1316  
1317  
1318  
1319

English:	
Sentence	It is the biggest acquisition in eBay's history
Gemma	It is the biggest acquisition in 'eBay's history
Custom	It is the biggest acquisition in eBay's history
<hr/>	
Achinese:	
Sentence	Biasajih tv dipeugot deungön cara keu peusenang masyarakat umum
Gemma	Bias aj ih tv di pe ug ot de ung ön cara ke u pe usen ang masyarakat umum
Custom	Bias aj ih tv dipe ug ot deung ön cara keu pe usen ang masyarakat umum
<hr/>	
Shan:	
Sentence	၂၁၁။
Gemma	၂၁၁။
Custom	၂၁၁။

Table 6: Qualitative examples for comparing tokenization using the customized vocabulary against Gemma’s original one. We find that the customized tokenizer reduces overtokenization without affecting English tokenization. Sentences have been word-segmented to simplify presentation.

Language name	ISO code	BELE	BELE	SIB-200	FLORES-200	XORQA-IN	XSUM-IN	GSM8K-NTL
English	eng_Latn	✓	✓			✓	✓	
SEA								
Achinese (Arabic)	ace_Arab <sup>L</sup>			✓	✓			
Achinese	ace_Latn <sup>L</sup>			✓	✓			
Balinese	ban_Latn <sup>L</sup>			✓	✓			
Betawi	bew_Latn <sup>M</sup>							✓
Banjar (Arabic)	bjn_Arab <sup>L</sup>			✓	✓			
Banjar	bjn_Latn <sup>L</sup>			✓	✓			
Buginese	bug_Latn <sup>L</sup>			✓	✓			
Cebuano	ceb_Latn <sup>M</sup>	✓		✓	✓			
Ilocano	ilo_Latn <sup>M</sup>	✓		✓	✓			✓
Indonesian	ind_Latn <sup>H</sup>	✓		✓	✓			
Javanese	jav_Latn <sup>M</sup>	✓		✓	✓			
Kachin	kac_Latn <sup>L</sup>	✓		✓				
Khmer	khm_Khmr <sup>M</sup>	✓		✓		✓		
Lao	lao_Laoo <sup>M</sup>	✓		✓		✓		
Kedah Malay	meo_Latn <sup>M</sup>							✓
Pattani Malay	mfa_Arab <sup>L</sup>							✓
Minangkabau (Arabic)	min_Arab <sup>L</sup>			✓	✓			
Minangkabau	min_Latn <sup>L</sup>			✓	✓			✓
Myanmar (Burmese)	mya_Mymr <sup>M</sup>	✓		✓	✓			
Pangasinan	pag_Latn <sup>L</sup>			✓	✓			
Shan	shn_Mymr <sup>L</sup>	✓		✓	✓			
Sundanese	sun_Latn <sup>M</sup>	✓		✓	✓			
Tagalog	tgl_Latn <sup>H</sup>			✓	✓			
Thai	tha_Thai <sup>H</sup>	✓		✓	✓			
Vietnamese	vie_Latn <sup>H</sup>	✓		✓	✓			
Waray (Philippines)	war_Latn <sup>L</sup>	✓		✓	✓			
Standard Malay	zsm_Latn <sup>H</sup>			✓	✓			
AFR								
Afrikaans	afr_Latn <sup>H</sup>	✓		✓				
Twi	aka_Latn <sup>L</sup>			✓				
Amharic	amh_Ethi <sup>H</sup>	✓		✓				
Bambara	bam_Latn <sup>L</sup>	✓		✓				
Bemba (Zambia)	bem_Latn <sup>L</sup>			✓		✓		
Chokwe	cjk_Latn <sup>L</sup>			✓		✓		

Language name	ISO code	BELE	BELE	SIB-200	FLORES-200	XORQA-IN	XSUM-IN	GSM8K-NTL
Dinka	din_Latn <sup>L</sup>				✓			
Dyula	dyu_Latn <sup>L</sup>			✓				
Efik	efi_Latn <sup>L</sup>							✓
Ewe	ewe_Latn <sup>L</sup>			✓	✓			
Fon	fon_Latn <sup>L</sup>			✓	✓			
Fulfulde	ful_Latn <sup>L</sup>				✓			
Nigerian Fulfulde	fuv_Latn <sup>L</sup>	✓		✓				
Hausa	hau_Latn <sup>H</sup>	✓		✓				
Igbo	ibo_Latn <sup>H</sup>	✓		✓		✓		
Kamba (Kenya)	kam_Latn <sup>L</sup>			✓	✓			
Kabiyè	kbp_Latn <sup>L</sup>			✓		✓		
Kabuverdianu	kea_Latn <sup>L</sup>	✓		✓				
Kikuyu	kik_Latn <sup>L</sup>			✓		✓		
Kinyarwanda	kin_Latn <sup>M</sup>	✓		✓		✓		
Kimbundu	kmb_Latn <sup>L</sup>			✓		✓		
Central Kanuri (Arabic)	knc_Arab <sup>L</sup>				✓			
Central Kanuri	knc_Latn <sup>L</sup>					✓		
Kongo	kon_Latn <sup>L</sup>			✓		✓		
Krio	kri_Latn <sup>L</sup>							✓
Lingala	lin_Latn <sup>L</sup>	✓		✓		✓		
Tshiluba (Luba-Lulua)	lua_Latn <sup>L</sup>			✓		✓		
Luganda	lug_Latn <sup>M</sup>	✓		✓		✓		
Luo	luo_Latn <sup>L</sup>	✓		✓		✓		
Mossi	mos_Latn <sup>L</sup>			✓		✓		
Sepedi	nso_Latn <sup>L</sup>	✓		✓		✓		
Nuer	nus_Latn <sup>L</sup>			✓		✓		
Chicewa (Zambia)	nya_Latn <sup>L</sup>	✓		✓		✓		
Oromo	orm_Latn <sup>M</sup>				✓			
Nigerian Pidgin	pcm_Latn <sup>L</sup>							✓
Rundi	run_Latn <sup>L</sup>			✓		✓		
Sango	sag_Latn <sup>L</sup>			✓		✓		
Shona	sna_Latn <sup>M</sup>	✓		✓		✓		
Somali	som_Latn <sup>H</sup>	✓		✓		✓		
Sesotho	sot_Latn <sup>H</sup>	✓		✓		✓		
Swati	ssw_Latn <sup>L</sup>	✓		✓		✓		
Swahili	swa_Latn <sup>H</sup>				✓			
Tamasheq (Latin)	taq_Latn <sup>L</sup>			✓		✓		
Tamasheq (Tifinagh)	taq_Tfng <sup>L</sup>			✓		✓		
Tigrinya	tir_Ethi <sup>M</sup>	✓		✓		✓		
Tswana	tsn_Latn <sup>M</sup>	✓		✓		✓		
Tsonga	tso_Latn <sup>L</sup>	✓		✓		✓		
Tumbuka	tum_Latn <sup>L</sup>			✓		✓		
Umbundu	umb_Latn <sup>L</sup>			✓		✓		
Wolof	wol_Latn <sup>L</sup>	✓		✓		✓		
Xhosa	xho_Latn <sup>M</sup>	✓		✓		✓		
Yoruba	yor_Latn <sup>H</sup>	✓		✓		✓		
Zulu	zul_Latn <sup>H</sup>	✓		✓		✓		

#### IND

Assamese	asm_Beng <sup>M</sup>	✓	✓	✓	✓	✓	✓	
Awadhi	awa_Deva <sup>L</sup>		✓	✓	✓	✓	✓	
Bengali	ben_Beng <sup>H</sup>	✓	✓	✓	✓	✓	✓	
Bengali (Latin)	ben_Latn <sup>M</sup>	✓						
Haryanvi	bgc_Deva <sup>L</sup>				✓		✓	
Bhojpuri	bho_Deva <sup>L</sup>		✓	✓	✓	✓	✓	
Tibetan	bod_Tibt <sup>M</sup>	✓		✓	✓	✓	✓	
Bodo (India)	brx_Deva <sup>L</sup>				✓		✓	
Dhivehi	div_Thaa <sup>M</sup>							✓
Dogri	doi_Deva <sup>L</sup>							✓
Dzongkha	dzo_Tibt <sup>M</sup>		✓					✓
Garhwali	gbm_Deva <sup>L</sup>				✓		✓	
Goan Konkani	gom_Deva <sup>L</sup>				✓		✓	
Gujarati	guj_Gujr <sup>H</sup>		✓	✓	✓	✓	✓	
Hindi	hin_Deva <sup>H</sup>	✓	✓	✓	✓	✓	✓	

Language name	ISO code	BELEBELE	SIB-200	FLORES-200	XORQA-IN	XSUM-IN	GSM8K-NTL
Hindi (Latin)	hin_Latin <sup>M</sup>	✓					
Chhattisgarhi	hne_Deva <sup>L</sup>		✓	✓	✓		✓
Hadathi	hoj_Deva <sup>L</sup>				✓		✓
Kannada	kan_Knda <sup>H</sup>	✓	✓	✓	✓		✓
Kashmiri	kas_Arab <sup>L</sup>		✓				
Kashmiri (Devanagari)	kas_Deva <sup>L</sup>		✓				
Mizo	lus_Latin <sup>M</sup>		✓				
Magahi	mag_Deva <sup>L</sup>		✓				
Maithili	mai_Deva <sup>L</sup>		✓	✓	✓	✓	✓
Malayalam	mal_Mlym <sup>H</sup>	✓	✓	✓	✓	✓	✓
Marathi	mar_Deva <sup>H</sup>	✓	✓	✓	✓	✓	✓
Meitei (Bengali)	mni_Beng <sup>L</sup>		✓	✓	✓	✓	✓
Malvi	mup_Deva <sup>L</sup>				✓	✓	
Marwari	mwr_Deva <sup>L</sup>				✓	✓	
Mazanderani	mzn_Arab <sup>M</sup>						✓
Nepali	npi_Deva <sup>M</sup>			✓		✓	
Odia	ory_Orya <sup>M</sup>	✓		✓	✓	✓	✓
Punjabi	pan_Guru <sup>M</sup>	✓	✓	✓	✓	✓	
Southern Pashto	pbt_Arab <sup>M</sup>			✓			
Pashto	pbu_Arab <sup>M</sup>				✓	✓	
Sanskrit	san_Deva <sup>M</sup>		✓	✓	✓	✓	✓
Santali (Ol Chiki)	sat_Olck <sup>L</sup>		✓		✓	✓	✓
Sinhala	sin_Latin <sup>L</sup>	✓					
Sinhala	sin_Sinh <sup>H</sup>	✓		✓			
Sindhi (Perso-Arabic)	snd_Arab <sup>M</sup>	✓		✓			
Tamil	tam_Taml <sup>H</sup>	✓	✓	✓	✓	✓	✓
Telugu	tel_Telu <sup>H</sup>	✓	✓	✓	✓	✓	✓
Urdu	urd_Arab <sup>H</sup>	✓	✓	✓	✓	✓	
Urdu	urd_Latin <sup>L</sup>	✓					

Table 7: Overview of languages included in each of our evaluation benchmark. H, M, L indicate high, medium, and low-resource languages, respectively.

Language Adaptation Results on FLORES-200 EN-XX																
Model	PaLM2				Gemma2				Aya23							
	Size		XXS		S		2B		9B		27B		8B		35B	
Variant	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA
SEA																
ind_Latn	65.9	65.0	70.7	69.1	64.4	64.6	68.9	68.3	69.9	68.5	73.2	69.5	71.9	68.6		
tha_Thai	42.9	40.9	50.3	49.8	35.6	35.3	44.7	43.8	47.4	45.6	20.3	33.9	28.4	37.0		
vie_Latn	56.4	56.4	61.7	61.6	52.1	54.2	58.7	58.3	60.0	59.8	65.7	63.5	64.9	62.3		
zsm_Latn	63.0	64.6	68.0	67.4	58.4	63.1	64.6	65.6	65.6	65.9	54.2	64.6	57.4	64.9		
tgl_Latn	54.2	55.1	62.1	61.4	48.4	55.2	57.8	58.3	60.4	59.7	36.7	55.3	47.4	57.9		
mya_Mymr	22.4	26.9	40.9	37.4	12.9	27.5	25.4	33.5	30.1	33.3	6.5	23.4	7.3	29.5		
lao_Laoo	30.6	37.8	48.8	47.9	9.5	39.0	26.7	42.9	33.7	42.1	4.0	27.0	10.2	41.1		
khm_Khmr	24.7	29.0	36.7	36.4	12.7	25.9	24.3	28.0	27.4	32.0	2.2	20.6	8.7	26.3		
ceb_Latn	41.7	56.3	58.5	60.6	31.8	54.0	50.6	58.2	54.9	59.1	25.8	55.6	35.1	57.2		
jav_Latn	43.9	50.1	53.9	53.3	26.4	48.7	43.5	51.6	47.0	52.6	25.8	51.4	30.1	52.6		
sun_Latn	38.9	48.1	51.1	52.0	21.8	45.1	39.4	44.5	41.7	47.8	29.2	43.7	29.2	48.7		
ilo_Latn	19.3	42.9	45.1	51.8	19.0	43.3	32.0	49.0	42.6	52.2	18.0	45.5	20.0	49.9		
war_Latn	29.9	51.6	54.6	59.2	33.7	54.1	49.0	58.2	52.4	59.3	25.6	52.5	35.2	55.5		
bug_Latn	7.9	15.9	16.5	25.0	14.1	22.2	18.1	28.9	21.7	28.5	15.4	21.9	18.8	27.2		
pag_Latn	14.5	32.8	28.2	43.6	25.1	38.6	28.2	43.3	34.2	43.8	17.9	39.2	20.5	43.0		
shn_Mymr	4.4	17.7	3.0	21.0	11.6	29.6	14.9	32.4	16.2	32.7	4.6	9.2	5.0	28.4		
min_Latn	25.0	46.9	49.8	53.6	23.8	47.7	38.4	51.3	45.9	52.8	35.1	47.6	36.0	50.8		
ace_Latn	8.2	29.0	25.7	40.0	9.4	32.5	16.1	37.0	26.8	40.3	16.0	32.6	15.2	38.3		
ban_Latn	13.7	41.3	39.4	42.4	21.7	44.4	30.1	46.8	37.1	47.5	21.8	45.1	24.7	47.2		
bjn_Latn	21.6	26.1	47.2	44.5	24.9	29.9	36.3	38.0	42.7	40.6	36.8	32.6	36.7	35.5		
ace_Arab	2.7	3.5	7.9	10.1	3.0	15.7	3.9	15.7	9.3	18.0	3.9	1.6	2.2	14.8		
bjn_Arab	2.7	4.5	9.8	11.3	4.7	11.0	6.5	11.9	12.0	20.6	3.9	7.3	4.1	19.2		
min_Arab	3.4	5.2	10.9	16.2	5.1	20.0	6.6	17.9	10.3	20.9	5.3	1.1	2.7	18.8		
Avg.	27.7	<b>36.9</b>	40.9	<b>44.1</b>	24.8	<b>39.2</b>	34.1	<b>42.8</b>	38.7	<b>44.5</b>	23.8	<b>36.7</b>	26.6	<b>42.4</b>		

Language Adaptation Results on FLORES-200 EN-XX														
Model	PaLM2				Gemma2						Aya23			
Size	XXS		S		2B		9B		27B		8B		35B	
Variant	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA
AFR														
swa_Latn	52.2	54.4	63.8	62.1	40.3	51.1	57.0	58.7	61.5	59.5	9.8	39.9	15.1	49.6
lin_Latn	9.0	25.7	20.6	40.8	7.9	30.2	14.7	42.2	21.0	42.1	9.9	30.8	12.1	35.1
yor_Latn	9.5	14.9	24.8	23.0	6.3	14.7	11.6	21.6	18.4	22.8	4.0	13.6	5.9	16.9
ful_Latn	4.0	6.2	5.5	13.4	8.0	10.0	5.3	11.0	6.0	16.0	6.4	11.2	7.8	11.0
ibo_Latn	23.2	30.1	39.0	38.9	8.7	23.8	20.5	34.4	29.6	35.9	5.0	22.6	7.0	27.3
orm_Latn	4.6	8.8	12.1	24.6	8.3	14.1	9.0	25.1	11.2	25.7	9.2	12.7	9.3	15.8
som_Latn	24.6	32.5	42.5	42.4	11.2	29.0	26.4	38.4	32.3	40.0	17.6	28.5	25.4	37.3
tso_Latn	9.8	20.3	16.4	38.4	9.3	20.7	16.0	37.0	26.3	40.7	7.6	17.3	9.2	18.4
nya_Latn	25.5	37.5	44.8	44.6	9.0	28.1	17.9	40.1	27.5	40.2	7.7	26.5	10.9	30.2
zul_Latn	27.2	39.0	49.0	47.9	10.3	25.2	24.4	39.5	32.4	40.6	8.3	22.1	7.8	26.1
kin_Latn	5.9	22.8	25.1	37.0	7.8	15.9	16.0	29.7	22.5	33.9	8.3	10.8	8.3	20.3
run_Latn	5.0	18.6	17.0	30.4	5.5	14.2	11.6	26.3	18.0	27.8	6.3	11.8	7.1	15.2
sna_Latn	25.3	33.5	40.3	39.9	8.3	22.8	18.0	33.6	26.5	35.5	6.9	18.2	9.4	22.3
xho_Latn	26.3	35.6	44.4	44.5	12.2	26.9	23.5	36.5	32.2	38.4	9.1	22.4	8.0	28.1
tsn_Latn	10.4	26.5	34.1	42.2	7.8	21.0	17.6	39.7	30.6	42.1	8.4	19.5	10.5	24.7
tir_Ethi	3.0	5.1	13.2	13.1	3.0	14.0	4.4	20.5	5.4	18.9	2.2	6.6	2.1	9.4
kik_Latn	5.9	9.0	9.6	12.6	9.1	11.0	9.5	9.0	12.0	11.2	7.1	7.8	7.0	8.2
kon_Latn	8.1	28.6	16.3	38.7	8.8	33.4	13.8	38.9	16.0	38.2	7.8	31.7	9.9	35.6
lua_Latn	10.2	19.1	7.3	20.9	8.4	19.7	8.8	27.5	11.8	30.5	7.1	11.8	9.4	18.7
umb_Latn	5.6	3.6	6.2	7.1	7.9	5.8	7.3	8.1	8.6	8.5	6.4	7.0	6.9	7.1
sot_Latn	21.1	35.6	50.0	50.5	7.9	30.2	21.0	42.5	34.3	42.9	8.1	29.3	10.9	31.7
mos_Latn	3.8	7.7	3.5	5.6	5.7	5.4	5.9	4.6	6.9	8.3	4.7	4.3	4.3	4.7
nso_Latn	10.0	29.9	33.5	46.2	7.5	20.5	16.7	42.3	31.3	45.0	8.4	15.8	10.8	25.9
knc_Latn	5.3	7.9	6.5	8.4	8.3	9.4	5.9	12.5	7.6	12.8	6.9	8.6	5.8	9.6
knc_Arab	7.9	4.5	6.1	5.5	6.2	5.9	4.4	3.9	7.8	8.9	8.2	3.5	9.0	5.4
luo_Latn	6.7	9.5	8.5	15.2	12.2	11.0	10.1	18.1	13.0	22.2	5.9	7.9	7.7	10.8
bem_Latn	6.2	8.6	19.3	23.4	7.0	13.0	11.3	17.5	17.3	22.2	7.7	12.1	9.8	10.0
lug_Latn	5.4	13.9	15.6	30.6	6.4	13.4	9.1	22.1	16.0	27.1	6.4	11.4	8.0	12.5
wol_Latn	3.3	7.4	11.9	18.9	5.5	9.9	7.3	15.0	9.2	17.0	6.9	6.1	8.0	8.9
kmb_Latn	7.7	7.4	7.5	11.6	8.0	10.8	8.8	14.9	10.8	21.3	6.9	9.4	8.2	11.3
kam_Latn	13.4	8.3	8.5	12.0	10.1	9.2	10.9	10.5	10.2	11.3	6.2	10.2	6.9	10.3
ewe_Latn	10.0	11.9	7.4	26.0	6.2	16.1	7.0	27.7	9.6	27.8	4.6	8.0	6.2	19.3
ssw_Latn	14.2	28.7	31.6	37.6	8.9	17.9	15.9	28.5	23.6	32.5	7.4	17.4	6.5	18.1
tum_Latn	9.0	24.1	23.4	33.2	7.7	18.8	11.9	27.8	16.9	31.5	6.6	17.4	8.9	18.5
fon_Latn	4.8	6.1	2.7	11.3	4.6	11.8	4.5	18.1	4.5	15.1	4.5	6.1	4.8	11.6
din_Latn	4.6	2.7	5.5	8.5	8.2	9.0	6.1	13.4	7.6	18.8	6.7	8.4	6.6	9.7
kbp_Latn	4.7	8.6	9.9	15.4	8.0	12.6	8.6	18.7	9.7	19.3	6.5	6.7	6.7	11.5
cjk_Latn	4.7	6.3	5.2	7.2	6.5	9.9	6.8	9.2	7.4	14.9	5.5	7.8	6.7	8.4
nus_Latn	3.4	5.2	4.3	8.7	8.6	13.1	6.0	17.1	6.6	16.4	5.1	7.8	6.8	11.4
taq_Latn	4.1	4.1	5.4	7.4	5.9	6.6	7.4	9.2	7.4	10.5	7.4	7.0	8.2	7.5
taq_Tfng	2.6	3.5	8.7	6.1	9.4	3.5	4.7	3.9	5.8	6.2	3.6	1.2	6.0	1.8
sag_Latn	8.0	11.3	10.4	15.1	10.2	15.5	7.9	23.0	11.3	27.8	8.0	11.5	10.4	20.8
Avg.	10.9	<b>17.3</b>	19.5	<b>25.4</b>	8.7	<b>16.8</b>	12.7	<b>24.2</b>	17.3	<b>26.4</b>	7.1	<b>14.1</b>	8.5	<b>17.6</b>
IND														
hin_Deva	50.2	50.3	60.2	58.3	47.0	50.2	54.8	54.1	55.8	54.4	60.2	48.7	59.9	53.1
ben_Beng	34.5	39.8	48.7	47.8	26.8	37.2	41.9	44.0	45.3	43.8	13.9	36.8	27.5	41.4
urd_Arab	31.8	37.9	47.8	48.3	26.8	37.0	40.8	45.5	44.8	45.5	11.3	38.7	19.5	42.4
tel_Telu	33.0	36.7	54.6	51.8	22.8	42.0	43.3	49.1	48.6	50.2	9.0	29.4	13.9	42.3
tam_Taml	33.7	37.0	54.4	52.4	27.0	38.3	45.7	48.9	49.8	48.9	19.0	31.7	32.6	41.6
mar_Deva	31.2	35.1	47.7	46.3	22.9	34.7	38.3	42.0	42.9	42.5	20.0	33.3	24.5	35.9
mai_Deva	24.8	38.1	47.2	49.4	23.4	38.5	35.7	44.6	41.8	46.2	32.3	35.5	37.7	41.8
bho_Deva	26.7	33.7	39.9	41.7	25.5	34.9	34.3	40.8	37.6	40.1	32.5	37.0	36.2	34.8
pbt_Arab	15.6	24.2	32.6	33.2	6.0	18.2	14.6	24.9	15.9	27.4	2.9	24.5	4.6	25.9
guj_Gujr	27.6	36.8	49.6	49.2	22.5	40.2	41.8	46.4	46.1	47.0	11.9	37.4	23.0	41.6
kan_Knda	30.6	34.0	50.8	49.2	21.7	40.4	40.9	47.2	45.6	46.2	9.1	32.7	18.9	39.9
awa_Deva	33.4	35.9	45.6	46.5	32.2	36.4	39.6	43.5	41.2	44.4	39.0	39.3	42.4	41.1
ory_Orya	14.8	19.8	46.3	44.5	10.1	24.9	21.3	39.6	27.7	35.7	4.6	5.8	9.1	14.3
mal_Mlym	27.3	29.9	52.4	49.1	24.1	35.2	42.1	45.5	47.1	45.4	3.3	23.8	5.1	34.4
pan_Guru	26.1	36.3	48.2	48.6	21.0	37.5	41.4	45.7	46.1	45.7	7.2	33.7	11.9	38.7
hne_Deva	30.2	40.5	48.2	50.4	29.8	40.5	40.8	47.9	44.0	46.9	39.8	39.5	42.4	43.7
npi_Deva	34.6	43.1	52.3	51.3	27.1	41.3	44.0	46.9	47.5	46.8	25.4	38.5	30.0	40.2

Language Adaptation Results on FLORES-200 EN-XX													
Model	PaLM2				Gemma2				Aya23				
Size	XXS		S		2B		9B		27B		8B		35B
Variant	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT
asm_Beng	13.5	26.3	36.9	36.4	13.3	25.0	27.8	31.7	33.3	33.4	8.4	18.2	15.4
mni_Beng	3.9	6.8	8.1	20.0	7.3	19.4	8.8	31.7	10.7	27.3	7.0	14.4	8.6
bod_Tibt	12.7	12.1	26.8	24.6	9.0	13.8	14.4	18.8	16.3	17.0	9.0	6.9	8.3
san_Deva	11.2	20.1	27.6	28.8	11.3	20.8	21.5	26.5	26.2	27.3	14.0	20.0	19.4
Avg.	26.1	<b>32.1</b>	44.1	<b>44.2</b>	21.8	<b>33.6</b>	34.9	<b>41.2</b>	38.8	<b>41.0</b>	18.1	<b>29.8</b>	23.4
													<b>34.8</b>

Table 8: Language Adaptation results on FLORES-200 EN-XX with 5-shot in-context learning. PT: Pre-training; LA: Language Adaptation.

Language Adaptation Results on GSM8K-NTL													
Model	PaLM2				Gemma2				Aya23				
Size	S				9B		27B		8B		35B		
Variant	PT	NTL	CALM	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT
asm_Beng	5.2	4.0	9.2	17.2	14.8	29.2	15.6	29.6	0.8	3.6	2.0	10.4	
bew_Latn	33.6	33.6	34.8	33.6	33.6	42.8	48.4	50.4	33.6	31.6	45.6	42.4	
bho_Deva	23.6	22.8	29.2	26.0	19.6	35.6	29.2	40.0	14.8	12.0	21.6	24.0	
doi_Deva	17.2	21.6	22.4	26.8	12.0	33.2	17.2	30.0	9.6	9.6	18.8	24.0	
div_Thaa	11.2	13.2	14.8	19.2	6.8	26.8	8.0	21.2	1.6	3.2	1.2	14.0	
dzo_Tibt	0.8	0.0	0.4	7.6	1.2	15.2	2.8	8.8	0.4	0.0	0.0	1.2	
efi_Latn	14.8	14.0	18.0	22.0	10.8	26.0	16.4	31.2	4.0	5.2	4.0	18.4	
gom_Deva	22.4	22.8	25.2	30.0	18.8	36.4	26.8	33.6	8.4	8.0	9.2	20.4	
ilo_Latn	14.8	14.0	16.8	18.8	18.0	27.6	23.6	34.8	5.6	9.6	9.2	16.0	
kri_Latn	12.4	20.0	18.8	18.8	20.4	25.2	21.6	34.4	8.0	8.8	13.2	11.6	
mai_Deva	22.8	21.2	24.8	25.6	19.6	29.2	25.2	30.4	15.6	12.8	29.2	26.4	
meo_Latn	28.8	33.2	34.0	28.4	28.4	42.0	45.6	46.4	28.8	29.2	41.2	39.6	
mfa_Arab	14.0	20.4	17.6	22.4	6.0	34.8	5.2	38.8	4.0	16.8	4.4	27.2	
min_Latn	25.2	24.8	23.6	24.8	13.6	34.8	30.0	39.2	9.2	19.2	17.2	32.0	
mni_Beng	2.8	6.0	4.4	9.6	3.2	19.6	2.8	15.6	1.6	2.8	1.2	7.6	
mzn_Arab	31.6	27.6	36.4	36.8	33.2	40.8	44.0	38.4	30.4	18.8	41.2	34.0	
nso_Latn	8.4	9.6	8.4	13.2	6.4	18.4	14.0	22.0	2.8	6.0	5.2	8.4	
ory_Orya	9.6	12.0	12.4	24.0	6.4	30.4	12.0	27.6	1.2	1.6	2.8	10.0	
pcm_Latn	34.4	31.6	33.6	30.0	43.2	44.0	47.6	51.2	28.8	26.0	41.2	38.8	
tso_Latn	7.2	11.6	10.0	10.0	10.8	15.2	11.2	20.0	3.2	5.2	4.8	8.8	
Avg.	17.0	18.2	19.7	<b>22.2</b>	16.3	<b>30.4</b>	22.4	<b>32.2</b>	10.6	<b>11.5</b>	15.7	<b>20.8</b>	

Table 9: Language Adaptation results on GSM8K-NTL with 5-shot in-context learning. PT: Pre-training; LA: Language Adaptation.

Language Adaptation Results on BELEBELE													
Model	PaLM2				Gemma2				Aya23				
Size	XXS		S		2B		9B		27B		8B		35B
Variant	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT
SEA													
eng_Latn	23.8	28.0	85.9	86.8	61.2	54.9	90.6	91.3	92.3	91.6	82.4	80.0	91.4
ind_Latn	25.9	26.4	80.7	80.0	49.3	45.6	84.2	84.1	87.8	87.0	73.7	71.7	84.2
tha_Thai	27.2	26.6	77.9	78.2	44.1	40.0	78.4	76.7	82.8	81.8	44.6	56.9	65.0
vie_Latn	25.9	25.8	82.1	80.6	46.3	43.9	84.3	84.7	88.1	86.6	73.3	71.3	84.2
mya_Mymr	25.8	25.9	75.7	71.2	32.6	38.1	68.0	72.6	65.7	75.2	27.9	48.8	36.2
lao_Laoo	29.1	27.0	70.3	67.7	28.7	36.7	60.6	69.6	59.1	73.4	27.4	47.1	29.4
khm_Khmr	27.1	26.9	77.4	70.1	30.6	38.6	64.7	72.1	65.2	75.6	25.9	47.0	33.8
ceb_Latn	25.2	25.8	74.7	75.1	38.2	41.2	76.0	79.4	78.7	81.6	44.1	55.8	55.0
jav_Latn	27.4	27.2	74.9	73.7	37.9	42.0	73.0	76.1	75.6	77.0	44.8	57.6	56.3
sun_Latn	26.3	26.2	71.0	70.1	35.9	38.8	68.6	73.2	70.1	76.0	37.0	59.2	47.6
ilo_Latn	27.9	28.0	64.6	67.3	34.7	39.0	60.9	69.9	63.8	74.6	36.9	45.2	39.9
war_Latn	26.4	27.2	77.0	77.3	36.9	41.4	70.7	79.3	73.3	81.4	43.9	56.3	50.1
shn_Mymr	27.4	26.2	28.9	40.6	25.2	31.2	32.0	47.9	33.0	49.8	23.6	32.4	26.4
													39.2

Language Adaptation Results on BELEBELE																
Model	PaLM2						Gemma2						Aya23			
	Size		XXS		S		2B		9B		27B		8B		35B	
Variant	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA	PT	LA
kac_Latn	24.9	27.1	33.4	43.2	27.6	35.3	33.1	48.1	29.3	48.7	31.0	33.9	28.6	43.4		
Avg.	26.5	<b>26.7</b>	69.6	<b>70.1</b>	37.8	<b>40.5</b>	67.5	<b>73.2</b>	68.9	<b>75.7</b>	44.0	<b>54.5</b>	52.0	<b>66.2</b>		
AFR																
eng_Latn	23.8	26.8	85.9	86.7	61.2	50.9	90.6	90.7	92.3	92.1	82.4	82.1	91.4	88.4		
afr_Latn	26.8	27.6	82.8	81.0	46.4	40.9	87.4	86.4	88.6	88.0	61.3	71.9	78.4	81.3		
hau_Latn	29.6	26.8	65.6	61.0	35.0	33.3	63.9	65.0	65.8	70.2	29.3	41.2	29.8	49.0		
lin_Latn	26.6	25.8	44.7	47.4	30.0	33.8	36.6	51.8	39.8	51.6	30.9	37.6	30.4	40.7		
yor_Latn	28.3	27.3	50.8	48.4	28.9	31.0	41.6	48.7	40.6	46.7	26.7	32.8	26.8	36.6		
ibo_Latn	25.2	27.9	54.4	49.0	29.8	30.4	48.1	51.2	49.4	54.2	30.0	36.8	30.1	41.8		
amh_Ethi	28.2	25.8	75.3	68.1	32.8	35.9	57.8	70.6	57.3	67.4	29.2	40.4	27.3	48.4		
som_Latn	27.3	28.7	60.8	56.3	30.8	30.3	49.6	56.3	51.8	58.0	30.8	36.1	32.9	42.4		
bam_Latn	27.9	27.4	35.7	38.6	27.4	28.6	34.9	40.4	33.0	39.2	32.4	33.7	30.3	36.2		
tso_Latn	24.7	26.8	44.4	51.4	29.7	32.1	43.7	56.1	44.9	57.9	32.0	39.2	33.2	41.9		
nya_Latn	28.2	26.8	54.7	53.4	29.9	30.2	47.0	54.1	50.1	54.9	28.6	33.2	31.4	40.2		
zul_Latn	27.2	27.1	54.6	56.9	29.9	32.2	55.3	61.8	55.4	61.2	30.4	37.4	31.3	44.3		
fuv_Latn	23.8	28.8	28.9	30.4	25.3	27.3	30.0	30.8	29.6	30.8	28.7	27.9	27.4	28.1		
kin_Latn	25.9	24.3	54.3	54.0	32.6	33.8	51.0	60.0	55.2	59.9	30.7	37.7	33.9	38.3		
sna_Latn	25.8	28.7	62.4	60.9	34.0	35.7	57.7	62.1	59.8	66.8	32.0	38.6	35.6	46.7		
xho_Latn	28.7	26.9	59.8	56.3	32.4	31.0	56.0	63.3	54.3	61.9	29.8	37.3	32.8	41.1		
tsn_Latn	28.2	27.4	54.6	56.8	31.0	32.7	46.9	58.3	46.8	54.4	28.2	35.6	33.3	42.3		
tir_Ethi	27.6	26.3	56.8	52.6	26.9	31.1	37.6	55.8	35.3	51.9	28.4	33.3	27.0	36.2		
sot_Latn	29.1	27.9	57.4	57.3	31.6	32.7	51.6	62.2	50.2	58.9	29.9	39.7	31.7	43.4		
nso_Latn	27.1	25.4	47.6	53.3	28.8	33.2	45.8	55.9	45.4	51.3	29.0	36.2	32.6	42.1		
luo_Latn	25.1	27.4	34.6	33.2	29.4	27.8	34.9	36.8	33.8	35.9	29.4	33.8	32.2	33.2		
lug_Latn	27.1	25.7	42.9	41.7	27.0	27.6	38.8	44.1	38.3	42.9	27.4	33.0	31.3	32.7		
wol_Latn	27.3	25.7	36.0	34.8	26.3	28.1	34.0	35.7	33.6	33.8	27.3	30.7	27.8	30.2		
ssw_Latn	26.1	26.4	48.2	48.2	29.0	29.9	46.2	51.0	43.0	51.3	28.4	34.6	32.4	35.6		
kea_Latn	27.3	27.1	65.1	63.4	36.4	31.3	54.8	47.3	58.6	52.4	43.6	40.0	50.2	41.9		
Avg.	26.9	26.9	<b>54.3</b>	53.7	32.1	<b>32.5</b>	49.7	<b>55.9</b>	50.1	<b>55.7</b>	33.5	<b>39.2</b>	36.1	<b>43.3</b>		
IND																
eng_Latn	23.8	27.6	85.9	86.0	61.2	49.8	90.3	90.8	92.3	92.0	82.4	84.2	91.4	89.9		
hin_Deva	26.3	25.0	75.7	73.2	40.8	40.0	77.0	72.3	80.1	76.3	60.6	64.8	73.8	70.9		
ben_Beng	27.1	28.1	76.0	73.6	41.1	39.2	76.7	73.6	78.6	77.6	36.1	54.1	48.7	67.7		
urd_Arab	27.3	27.8	77.0	73.8	41.4	39.3	76.6	75.2	79.8	78.9	39.2	53.8	49.2	68.3		
tel_Telu	27.2	27.7	71.7	69.6	38.1	36.1	71.3	69.4	75.2	73.0	31.6	43.1	34.6	58.8		
tam_Taml	28.4	26.1	75.2	72.6	42.3	40.6	75.7	73.6	77.9	76.3	40.3	51.2	61.0	64.7		
mar_Deva	25.6	24.8	77.1	73.3	38.1	36.0	76.7	73.1	79.2	75.7	39.1	51.8	52.3	65.4		
kan_Knda	28.6	23.6	78.1	73.0	40.3	41.1	74.7	73.9	77.1	78.2	28.9	52.2	41.1	68.7		
ory_Orya	26.9	27.0	77.6	72.3	30.9	36.6	60.9	70.4	63.1	74.0	30.6	30.0	40.8	49.6		
mal_Mlym	26.8	29.2	79.7	75.9	42.6	39.6	76.1	75.2	80.1	77.1	35.3	44.2	51.6	60.9		
pan_Guru	28.6	27.9	76.8	72.1	36.8	35.2	77.0	74.3	77.4	74.1	29.6	49.2	33.0	65.2		
snd_Arab	26.6	26.2	68.4	68.4	32.4	34.9	59.0	66.4	64.9	69.4	32.3	44.1	39.6	59.7		
sin_Latn	26.6	24.1	35.1	34.9	28.2	26.8	38.3	36.8	38.9	38.3	30.9	28.9	34.3	32.8		
sin_Sinh	26.2	28.1	76.9	75.7	35.0	39.9	70.9	74.3	71.4	73.2	30.2	50.3	38.0	61.4		
asm_Beng	28.7	24.4	75.9	72.3	36.0	32.8	67.9	70.8	73.3	71.1	34.0	41.3	40.9	58.8		
hin_Latn	26.7	29.1	68.3	63.9	37.9	38.0	69.4	69.9	74.2	74.3	45.3	48.9	57.6	60.8		
bod_Tibt	27.3	27.8	50.2	38.7	27.3	28.4	39.1	43.0	35.0	44.8	25.7	26.3	30.1	31.8		
ben_Latn	28.6	26.4	52.4	51.9	29.8	30.6	49.8	54.1	55.6	57.8	32.2	34.7	38.9	40.7		
urd_Latn	26.8	27.7	56.3	54.8	34.2	31.8	56.8	56.8	60.0	61.7	36.6	37.8	43.6	46.6		
Avg.	<b>27.0</b>	26.8	<b>70.2</b>	67.2	<b>37.6</b>	36.7	67.6	<b>68.1</b>	70.2	<b>70.7</b>	37.9	<b>46.9</b>	47.4	<b>59.1</b>		

Table 10: Language Adaptation results on BELEBELE with 5-shot in-context learning. PT: Pre-training; LA: Language Adaptation.

---

**EMBSWAP Results on BELEBELE**


---

Model	PaLM2						Gemma2						Aya23								
	XXS			S			2B			9B			27B			8B			35B		
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR

---

**EMBSWAP Results on BELEBELE**


---

Model	PaLM2						Gemma2						Aya23								
	XXS			S			2B			9B			27B			8B			35B		
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR

**SEA**

eng_Latn	80.9	76.0	78.9	95.7	95.9	95.8	87.7	85.9	87.4	93.9	93.9	93.6	95.4	95.0	94.0	88.9	83.4	86.4	94.2	87.0	93.2
ind_Latn	74.0	72.3	73.4	91.8	91.9	92.0	78.7	76.1	75.8	88.8	87.6	88.2	90.4	89.7	88.3	82.0	74.3	74.3	89.6	72.7	86.3
tha_Thai	66.8	65.2	65.8	88.3	88.3	88.1	70.0	67.7	68.3	83.6	83.3	82.7	85.7	85.1	85.3	57.2	59.3	63.7	73.9	60.7	79.4
vie_Latn	74.1	72.0	72.0	92.6	92.7	92.6	75.6	74.0	74.9	89.7	89.2	88.4	90.0	89.7	88.8	81.9	75.9	74.6	90.2	75.6	87.6
mya_Mymr	52.0	56.1	56.7	86.8	85.1	85.7	36.3	62.1	60.8	69.7	76.1	77.9	65.2	78.2	78.2	34.0	51.1	56.2	38.1	49.2	69.0
lao_Laoo	44.9	55.9	57.3	81.2	82.3	83.1	43.3	60.9	61.4	61.4	77.3	75.9	59.7	77.1	77.3	31.8	50.6	51.9	38.8	49.2	66.4
khm_Khmr	51.9	52.6	55.8	86.3	84.6	86.7	45.1	65.0	62.6	69.2	78.6	78.4	68.6	81.1	79.4	26.8	46.9	51.4	43.6	51.7	71.2
ceb_Latn	52.2	61.8	62.6	88.9	90.6	90.0	52.9	64.2	65.8	79.0	82.4	84.4	82.2	84.1	85.9	45.1	59.8	60.9	57.9	59.1	77.8
jav_Latn	61.0	61.8	62.2	86.7	86.7	87.3	48.1	59.8	64.7	78.7	80.8	80.9	79.0	81.7	81.4	48.1	62.9	65.3	62.0	58.7	77.8
sun_Latn	55.2	56.4	58.3	86.3	84.3	85.4	41.1	56.8	61.7	70.9	79.6	78.8	71.2	79.6	79.9	41.0	57.7	61.1	51.2	53.7	75.4
ilo_Latn	41.9	50.7	53.4	79.2	84.0	84.3	43.1	53.0	56.8	63.8	75.9	76.9	68.6	74.6	78.2	34.4	45.7	45.3	42.1	47.6	65.0
war_Latn	50.6	60.8	63.4	89.3	90.6	90.8	49.3	65.3	64.8	75.1	84.7	83.7	79.1	84.1	87.9	42.4	58.8	64.2	53.9	59.0	77.2
shn_Mymr	26.1	33.9	36.8	28.8	52.1	50.8	28.0	39.8	40.0	33.4	53.3	53.8	35.7	54.4	54.1	26.8	33.7	34.6	26.4	30.1	46.2
kac_Latn	31.7	36.8	39.4	40.8	55.4	58.3	28.3	36.0	38.9	31.6	50.1	54.1	35.3	48.6	53.1	28.1	32.0	37.8	34.1	33.0	43.9
Avg.	54.5	58.0	<b>59.7</b>	80.2	83.2	<b>83.6</b>	52.0	61.9	<b>63.1</b>	70.6	78.1	<b>78.4</b>	71.9	78.8	<b>79.4</b>	47.8	56.6	<b>59.1</b>	56.9	56.2	<b>72.6</b>

**AFR**

eng_Latn	80.9	79.1	80.7	95.7	95.4	95.7	87.7	85.2	87.7	93.9	94.3	93.6	95.4	94.9	93.1	88.9	84.2	86.1	94.2	92.2	91.8
afr_Latn	70.4	69.8	71.6	93.3	93.8	93.1	72.4	74.9	75.3	90.1	90.6	89.4	92.1	91.1	89.4	69.6	74.9	74.4	84.0	81.9	87.3
hau_Latn	40.4	38.4	43.8	80.2	75.7	77.9	35.0	42.9	45.4	67.4	68.1	71.2	68.1	68.1	74.1	28.0	40.7	40.2	33.8	43.1	58.3
lin_Latn	34.2	36.7	40.7	57.1	62.9	65.7	30.2	34.4	42.3	36.6	55.0	57.3	42.0	47.6	58.8	28.1	34.0	37.2	35.7	35.7	47.0
yor_Latn	32.7	29.7	32.7	62.7	58.9	59.9	29.7	31.9	34.4	40.9	46.4	51.8	41.6	45.3	51.9	27.3	30.4	33.4	34.4	31.2	44.6
ibo_Latn	34.7	35.6	37.2	68.6	64.7	65.8	31.4	35.1	44.2	47.6	54.8	55.7	48.8	51.7	57.2	29.7	35.6	37.1	35.0	37.0	49.1
amh_Ethi	48.7	43.4	50.4	86.9	82.2	83.0	35.7	56.6	55.9	60.3	74.2	75.6	60.3	72.3	73.7	25.2	41.9	40.8	32.8	43.6	57.8
som_Latn	39.9	37.2	39.9	73.3	70.9	72.7	33.9	38.3	41.8	50.0	59.3	60.9	52.4	58.2	63.0	29.8	34.8	34.3	36.1	37.6	49.6
bam_Latn	34.7	34.3	37.4	43.7	45.7	46.2	30.2	30.6	35.6	34.2	40.8	41.1	36.8	40.7	43.7	32.6	32.6	30.8	36.2	32.6	38.3
tso_Latn	36.3	39.6	44.3	57.9	69.7	70.8	34.7	41.1	45.8	43.7	54.7	62.6	48.7	58.3	62.3	31.2	36.4	37.0	40.1	37.0	52.6
nya_Latn	38.4	38.8	40.7	70.8	67.1	70.1	30.1	34.3	39.2	46.7	55.3	58.1	50.6	53.4	59.6	27.1	34.7	34.3	32.1	35.6	46.3
zul_Latn	41.2	41.8	46.7	76.2	72.0	75.0	33.7	41.3	42.9	57.4	66.0	69.6	58.2	62.7	68.9	33.0	36.9	37.6	38.8	40.8	52.3
fuv_Latn	30.7	28.0	27.9	31.0	31.9	32.1	29.0	27.3	28.0	28.2	29.8	31.0	29.7	31.2	29.9	27.4	25.7	26.7	31.0	28.1	29.3
kin_Latn	37.1	41.3	44.3	68.9	71.8	73.8	35.3	37.2	41.0	51.6	64.2	67.2	58.6	59.6	68.4	32.0	35.3	34.9	37.3	36.0	46.8
sna_Latn	42.9	44.1	50.3	77.6	75.7	77.9	32.4	42.6	46.4	58.7	67.0	69.4	62.3	65.0	68.9	30.3	39.0	42.2	38.2	42.7	56.0
xho_Latn	42.9	42.4	47.0	76.6	74.1	76.2	35.4	40.3	44.1	58.4	64.7	68.3	60.2	65.3	70.2	30.3	36.9	38.2	38.2	39.9	52.0
tsn_Latn	35.4	38.8	44.6	68.3	69.9	70.8	31.8	37.1	43.1	48.8	57.4	63.4	53.3	55.8	61.4	32.3	35.4	36.6	36.6	35.7	49.2
tir_Ethi	32.2	38.1	40.9	70.2	71.9	71.6	29.3	41.6	44.3	33.2	59.1	60.4	37.6	51.9	57.7	27.4	33.3	32.9	31.6	29.7	45.1
sot_Latn	38.1	40.4	42.2	76.9	76.0	77.9	29.4	38.7	42.9	52.7	64.6	66.9	53.4	60.6	66.8	30.1	37.4	37.1	34.4	39.1	53.3
nso_Latn	35.4	38.3	41.7	65.3	69.8	70.7	31.3	35.7	42.4	44.0	57.6	63.1	51.4	58.1	63.3	30.0	34.2	35.7	35.8	35.8	47.2
luo_Latn	33.6	31.1	34.4	40.8	42.2	42.7	30.0	29.7	32.7	34.1	34.6	36.0	36.2	37.0	36.7	27.9	32.1	32.0	34.0	33.0	36.9
lug_Latn	31.6	32.0	34.8	51.8	54.9	57.9	29.9	29.7	33.0	38.8	45.0	47.0	40.0	43.9	47.4	29.8	32.2	32.1	35.2	28.9	39.3
wol_Latn	31.4	32.3	31.8	44.0	38.7	39.2	29.2	26.4	29.4	30.2	30.4	35.3	31.9	34.9	34.2	28.2	28.9	29.6	32.4	30.3	32.0
ssw_Latn	33.3	36.3	41.0	63.2	61.3	66.2	29.2	33.9	38.0	42.7	51.7	57.1	46.8	51.9	56.1	28.7	33.1	32.7	32.9	31.6	41.7
kea_Latn	41.9	37.7	38.8	84.9	79.3	80.6	42.4	37.1	37.4	56.3	51.3	52.3	59.4	47.0	52.8	46.6	40.1	40.2	57.4	44.9	48.4
Avg.	40.0	40.2	<b																		

---

**EMBSWAP Results on BELEBELE**


---

Model	PaLM2						Gemma2						Aya23											
	Size			XXS			S			2B			9B			27B			8B			35B		
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR
bod_Tibt	27.9	29.4	31.9	61.0	48.1	54.2	27.1	32.4	35.6	40.6	45.4	47.4	37.6	41.9	47.7	27.4	25.8	27.3	30.1	24.6	31.4			
ben_Latn	35.3	38.3	38.6	68.7	66.9	70.7	33.6	34.9	38.7	44.6	53.1	58.6	56.9	60.7	63.1	30.6	31.7	34.9	38.9	37.2	44.9			
urd_Latn	33.2	35.6	39.0	72.0	69.9	73.3	35.6	38.8	39.3	56.9	56.1	60.7	64.2	64.0	63.4	35.9	36.0	37.9	49.1	41.4	50.9			
Avg.	49.4	50.9	<b>53.8</b>	<b>82.1</b>	79.9	81.8	50.2	54.5	<b>55.8</b>	68.9	71.2	<b>73.1</b>	72.9	73.1	<b>74.3</b>	42.7	48.0	<b>50.8</b>	54.6	53.0	<b>55.2</b>			

Table 11: EMBSWAP results on BELEBELE with zero-shot prompting. FL: LLMs instruction-tuned on FLAN mixture; ES: EMBSWAP; +LR: EMBSWAP with LoRA Adaptation.

---

**EMBSWAP Results on SIB-200**


---

Model	PaLM2						Gemma2						Aya23											
	Size			XXS			S			2B			9B			27B			8B			35B		
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR
<b>SEA</b>																								
eng_Latn	85.3	83.8	82.4	76.0	76.5	77.9	82.4	82.4	79.9	83.3	81.4	80.4	83.3	81.9	82.4	82.8	80.8	79.8	82.8	80.8	79.8			
ace_Arab	27.9	51.5	47.1	53.4	70.1	67.6	37.7	45.6	50.5	53.4	64.2	65.2	46.6	62.7	67.6	37.4	46.5	38.4	39.4	51.5	65.7			
ace_Latn	63.2	78.4	71.6	71.1	80.9	77.0	62.3	68.6	74.0	73.5	76.5	79.9	71.1	73.0	81.9	63.6	74.7	61.6	68.7	73.7	76.8			
ban_Latn	75.0	80.9	76.5	73.0	77.0	77.9	68.6	74.0	76.0	81.4	78.4	78.9	81.9	77.5	81.4	70.7	75.8	66.7	71.7	70.7	75.8			
bjn_Arab	17.2	43.6	47.1	55.4	63.7	62.3	31.4	42.2	48.0	47.1	59.8	62.3	44.1	59.3	64.7	36.4	42.4	39.4	40.4	48.5	58.6			
bjn_Latn	73.0	78.9	75.0	76.5	79.4	81.4	67.6	71.6	71.1	80.4	78.9	79.9	77.9	78.9	79.9	72.7	73.7	60.6	73.7	70.7	71.7			
bug_Latn	57.8	67.2	71.1	63.7	74.0	72.5	57.4	62.3	64.2	62.7	69.6	69.1	63.2	68.6	70.1	55.6	66.7	54.5	65.7	69.7	75.8			
ceb_Latn	79.4	84.3	79.4	79.9	77.9	77.0	75.5	79.4	77.9	85.3	84.3	79.9	83.8	84.3	85.3	73.7	77.8	75.8	77.8	76.8	79.8			
ilo_Latn	68.6	81.9	80.4	72.5	80.9	81.9	73.5	79.4	81.4	78.9	83.8	82.4	79.9	85.3	82.4	68.7	78.8	73.7	66.7	76.8	81.8			
ind_Latn	85.8	83.8	80.4	77.5	80.4	80.9	84.8	83.3	80.9	83.8	83.3	83.3	82.8	80.9	83.3	84.8	78.8	74.7	80.8	77.8	78.8			
jav_Latn	75.5	81.4	77.9	75.5	76.5	76.5	73.0	76.0	75.5	82.4	81.4	79.9	79.1	81.4	81.4	76.8	76.8	67.7	79.8	73.7	77.8			
kac_Latn	41.2	64.7	66.7	42.2	70.6	70.1	35.8	63.7	67.6	50.0	71.6	76.5	49.0	70.6	74.5	48.5	62.6	60.6	45.5	60.6	72.7			
khm_Khmr	82.4	85.3	79.9	80.9	77.9	77.9	71.6	79.9	76.0	83.8	81.4	83.3	77.0	82.8	83.8	55.6	73.7	70.7	62.6	74.7	74.7			
lao_Laoo	79.9	86.8	80.4	78.4	82.8	83.8	69.6	81.4	78.4	76.5	81.4	81.4	73.0	82.8	82.8	55.6	75.8	66.7	58.6	74.7	77.8			
min_Arab	17.6	40.7	45.1	57.8	68.1	68.6	32.4	44.1	51.0	44.6	60.3	59.3	40.2	59.8	62.3	27.3	43.4	30.3	40.4	41.4	54.5			
min_Latn	76.0	81.9	77.0	76.5	78.9	81.9	70.1	75.5	77.0	81.9	85.3	83.8	77.5	83.3	82.4	66.7	78.8	71.7	73.7	77.8	79.8			
mya_Mymr	78.9	83.3	78.4	77.9	79.1	81.9	65.2	79.9	77.5	80.4	84.3	83.3	78.9	81.9	80.9	40.4	72.7	66.7	56.6	63.6	74.7			
pag_Latn	67.6	83.3	78.9	76.0	80.4	82.4	67.6	76.0	77.5	77.0	80.4	81.9	72.5	85.3	83.8	65.7	73.7	69.7	68.7	71.7	75.8			
shn_Mymr	39.7	76.5	79.4	42.2	78.9	80.4	42.6	76.0	80.9	56.9	79.9	81.4	56.4	78.4	78.9	53.5	70.7	66.7	51.5	71.7	74.7			
sun_Latn	77.0	82.8	79.4	76.5	78.9	81.4	74.0	78.4	79.4	80.4	81.4	80.9	80.4	82.4	72.7	77.8	72.7	71.7	73.7	82.8				
tgl_Latn	83.3	81.9	77.9	78.4	77.0	80.4	77.9	77.0	76.5	85.3	83.3	81.9	83.3	80.4	82.8	79.8	82.8	77.8	81.8	79.8	77.8			
tha_Thai	84.8	85.8	79.4	76.5	78.4	79.9	82.8	81.4	77.0	80.4	83.3	82.4	81.9	82.4	81.9	72.7	73.7	66.7	75.8	75.8	75.8			
vie_Latn	82.4	85.8	78.9	77.0	77.9	79.9	81.9	81.4	80.9	81.9	82.8	82.4	82.4	82.4	81.9	85.9	83.8	72.7	80.8	76.8	79.8			
war_Latn	75.0	80.9	75.0	72.5	77.9	77.9	77.5	81.4	78.4	82.8	81.9	79.9	81.9	83.3	79.8	83.8	82.8	78.8	78.8	78.8	78.8			
zsm_Latn	85.3	83.3	77.0	77.0	79.4	81.9	83.8	80.9	80.9	85.3	85.8	82.4	86.8	84.3	85.3	80.8	82.8	70.7	79.8	79.8	78.8			
Avg.	67.2	<b>76.7</b>	73.7	70.6	77.0	<b>77.6</b>	65.9	72.9	<b>73.5</b>	74.4	<b>78.6</b>	78.5	72.6	78.1	<b>79.5</b>	64.3	72.4	<b>65.6</b>	66.9	70.8	<b>75.2</b>			

---

**EMBSWAP Results on SIB-200**


---

Model	PaLM2						Gemma2						Aya23								
	Size			XXS			S			2B			9B			27B			8B		
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR
nso_Latn	41.7	71.1	75.0	71.6	74.0	79.9	46.6	65.2	47.5	62.3	76.5	76.0	64.7	73.5	79.4	45.5	69.7	55.6	43.4	60.6	75.8
nus_Latn	24.0	41.2	44.6	34.8	54.9	56.9	34.3	41.7	36.8	42.2	60.3	62.7	43.6	50.5	54.4	44.4	42.4	42.4	52.5	36.4	55.6
nya_Latn	70.6	76.0	71.1	76.0	79.4	78.9	49.0	67.2	51.5	71.1	77.9	75.5	71.1	78.9	78.9	43.4	65.7	61.6	62.6	73.7	73.7
run_Latn	41.2	73.5	72.5	71.6	74.5	75.0	45.1	59.8	54.4	66.2	74.5	78.4	66.7	75.0	76.5	48.5	57.6	58.6	56.6	58.6	68.7
sag_Latn	44.6	54.4	60.3	46.6	62.3	62.7	41.7	53.4	40.7	54.4	66.2	68.6	50.0	61.8	66.7	50.5	56.6	58.6	53.5	48.5	60.6
sma_Latn	62.7	74.0	71.6	72.5	78.4	76.0	46.1	67.6	56.4	69.6	76.0	76.5	69.6	81.9	82.8	45.5	58.6	60.6	53.5	67.7	78.8
som_Latn	61.3	73.0	72.5	75.5	77.5	78.4	54.4	66.2	56.4	73.5	79.9	75.5	72.1	79.4	77.5	48.5	71.7	68.7	62.6	56.6	68.7
sot_Latn	57.4	72.5	74.5	73.0	74.0	76.5	49.5	69.1	52.9	68.6	77.5	77.0	67.2	78.9	79.9	51.5	60.6	61.6	44.4	63.6	69.7
ssw_Latn	52.5	74.0	73.5	70.1	76.0	77.0	48.0	67.2	52.0	68.6	77.5	76.5	63.7	72.1	75.5	42.4	64.6	59.6	48.5	56.6	69.7
taq_Latn	46.6	49.0	47.5	49.5	51.5	58.3	44.1	48.5	43.6	51.5	55.9	51.5	52.0	52.0	52.5	48.5	56.6	50.5	58.6	50.5	56.6
taq_Tfng	11.3	27.0	27.5	21.6	23.0	25.0	26.5	27.0	25.0	22.1	27.5	26.5	32.4	27.0	25.5	30.3	26.3	27.3	30.3	15.2	27.3
tir_Ethi	53.9	63.2	68.1	73.5	75.0	76.5	44.6	69.1	51.5	62.7	80.9	82.4	59.3	76.5	75.5	30.3	62.6	54.5	31.3	47.5	66.7
tsn_Latn	44.6	66.2	69.6	66.7	74.5	74.5	45.6	64.2	51.5	63.2	77.9	79.4	63.2	74.0	79.4	51.5	65.7	61.6	57.6	64.6	74.7
tso_Latn	39.7	61.8	70.1	64.2	76.0	77.0	46.1	59.8	47.5	56.9	74.0	75.5	58.3	72.5	77.9	45.5	60.6	60.6	53.5	60.6	72.7
tum_Latn	60.3	77.9	74.0	68.1	74.5	75.0	44.1	62.7	53.9	68.6	78.4	79.9	63.2	74.0	77.5	52.5	58.6	62.6	57.6	61.6	66.7
umb_Latn	38.7	48.0	47.1	40.2	45.6	46.1	37.3	38.7	37.3	50.5	49.0	50.0	50.0	47.1	45.1	44.4	45.5	45.5	51.5	40.4	48.5
wol_Latn	52.9	52.0	54.4	61.3	66.2	66.7	49.0	57.4	51.0	59.3	59.3	63.2	59.3	63.7	59.8	51.5	57.6	60.6	55.6	53.5	62.6
xho_Latn	69.6	78.4	77.9	76.0	77.0	80.4	54.9	71.1	54.4	73.5	81.9	79.4	77.0	79.4	83.3	47.5	67.7	58.6	57.6	62.6	74.7
yor_Latn	64.2	63.7	67.2	70.6	72.5	75.5	52.0	56.4	44.6	66.7	73.5	72.1	62.7	68.1	69.1	47.5	56.6	57.6	60.6	54.5	66.7
zul_Latn	63.2	78.9	73.0	75.5	78.9	79.9	52.5	73.0	53.9	70.6	81.9	78.4	70.6	82.4	81.4	53.5	66.7	61.6	63.6	68.7	72.7
Avg.	51.3	62.2	<b>63.5</b>	60.9	66.9	<b>68.7</b>	47.3	<b>57.9</b>	47.8	61.0	68.8	<b>69.2</b>	60.2	66.1	<b>68.3</b>	48.5	58.7	<b>55.3</b>	52.9	53.8	<b>64.7</b>

**IND**


---

eng_Latn	85.3	84.8	84.3	76.0	75.0	78.4	82.4	81.4	81.4	83.3	81.4	81.4	83.3	82.4	84.3	82.8	85.9	77.8	82.8	79.8	76.8
asm_Beng	73.0	80.9	79.4	77.0	78.4	79.9	69.6	73.0	68.6	82.8	83.8	82.4	79.4	82.4	79.4	56.6	65.7	62.6	61.6	68.7	72.7
awa_Deva	80.9	84.8	82.4	79.4	78.9	79.9	77.5	74.5	73.5	80.9	80.9	80.4	78.9	79.9	82.4	81.8	78.8	72.7	83.8	79.8	80.8
ben_Beng	81.9	84.8	79.9	80.9	78.9	81.4	77.9	80.4	76.0	83.8	84.3	82.4	84.3	82.4	84.3	61.6	80.8	72.7	70.7	76.8	77.8
bho_Deva	76.5	81.4	80.4	75.5	78.9	77.9	74.0	73.5	74.0	79.9	83.8	82.4	77.0	80.9	83.8	74.7	82.8	72.7	76.8	77.8	77.8
dzo_Tibet	40.2	71.6	68.1	63.7	75.5	77.0	23.5	60.8	56.9	60.8	74.5	77.5	49.0	60.8	67.2	26.3	33.3	32.3	28.3	38.4	59.6
guj_Gujr	77.5	82.4	82.4	78.4	79.9	81.9	74.5	79.4	77.5	81.4	85.8	84.8	81.9	81.4	85.3	55.6	80.8	69.7	67.7	74.7	77.8
hin_Deva	81.9	83.3	83.8	82.8	74.5	80.9	83.3	80.9	77.9	86.8	84.8	83.3	82.4	84.3	82.8	81.8	82.8	75.8	81.8	75.8	77.8
hne_Deva	78.9	84.3	83.8	77.9	75.0	80.4	74.5	73.0	77.5	82.4	81.4	79.9	78.4	81.4	80.4	77.8	83.8	72.7	78.8	79.8	79.8
kan_Knda	77.5	83.8	81.9	79.9	77.5	78.9	74.0	77.5	73.5	84.3	81.4	83.3	83.3	84.8	86.3	52.5	70.7	69.7	68.7	73.7	73.7
kas_Arab	64.7	69.6	70.1	77.0	79.4	82.4	63.7	63.7	66.2	74.5	75.5	74.0	70.6	75.0	77.5	52.5	68.7	64.6	63.6	67.7	70.7
kas_Deva	56.9	61.3	64.2	74.5	71.6	79.4	51.0	53.9	57.8	64.7	69.1	73.0	63.2	69.1	66.7	62.6	59.6	56.6	63.6	56.6	62.6
lus_Latn	58.3	78.4	75.5	59.8	72.1	77.9	61.8	71.1	70.6	77.9	80.9	80.9	78.4	82.4	81.9	54.5	65.7	67.7	54.5	66.7	76.8
mag_Deva	81.9	81.4	79.4	78.9	74.5	78.4	77.5	74.5	76.0	82.8	79.9	80.9	79.9	84.8	81.9	77.8	81.8	68.7	81.8	75.8	75.8
mai_Deva	83.8	83.8	82.8	77.0	77.5	81.9	76.0	74.5	76.0	80.4	83.8	82.8	78.9	82.4	81.4	76.8	81.8	71.7	79.8	78.8	79.8
mal_Mlym	77.0	79.4	83.8	80.9	77.9	80.9	73.5	72.1	76.5	83.8	81.4	80.9	80.4	83.3	84.8	66.7	70.7	66.7	73.7	71.7	76.8
mar_Deva	82.4	80.4	81.9	78.4	79.4	79.9	75.5	75.5	77.0	83.3	82.4	83.8	78.4	85.3	87.3	73.7	78.8	70.7	75.8	74.7	79.8
mmi_Beng	34.8	55.9	64.2	57.4	72.5	77.0	44.6	52.5	54.9	56.4	76.0	75.5	50.5	69.1	77.9	27.3	59.6	51.5	40.4	52.5	72.7
pan_Guru	75.5	80.4	78.9	76.5	75.5	80.9	72.1	78.4	74.5	85.3	84.3	84.3	83.3	81.9	83.8	44.4	67.7	61.6	63.6	70.7	76.8
san_Deva	71.6	76.0	75.5	76.0	75.5	77.9	66.2	65.2	67.2	78.4	77.9	79.4	69.6	73.5	77.5	68.7	68.7	60.6	75.8	67.7	71.7
sat_Olck	13.7	23.0	24.5	37.7	26.5	28.4	34.8	24.5	29.9	69.6	43.6	48.0	67.6	36.3	41.2	26.3	26.3	27.3	26.3	14.1	35.4
sin_Sinh	78.9	81.9	82.4	77.0	74.0	75.5	61.3	73.5	76.0	83.8	82.4	83.8	77.5	81.4	81.4	46.5	74.7	70.7	65.7	67.7	78.8
snd_Arab	72.1	81.4	80.9	77.9	76.5	79.9	63.7	70.6	70.1	77.0	82.8	83.3	80.4	80.9	84.8	54.5	74.7	72.7	63.6	74.7	80.8
tam_Taml	79.4	82.8	79.4	80.4	78.9	80.9	77.9	76.5	77.0	82.4	83.8	83.8	82.4	83.3	83.3	69.7	70.7	66.7	76.8	72.7	73.7
tel_Telu	77.9	85.3	84.3	79.4	78.4	82.8	78.9	77.0	77.0	85.8	84.3	85.8	83.3	86.8	87.3	60.6	77.8	67.7	63.6	71.7	77.8
urd_Arab	79.9	83.3	84.8	77.0	77.5	80.9	73.5	74.0	74.5	84.3	85.3	84.3	83.8	85.3	84.8	61.6	80.8	71.7	76.8	73.7	79.8
Avg.	70.9	<b>77.2</b>	76.9	74.5	74.6	<b>77.8</b>	67.8														

---

**EMBSWAP Results on FLORES-200**


---

Model	PaLM2												Gemma2												Aya23					
	Size			XXS			S			2B			9B			27B			8B			35B								
	Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR											
khm_Khmr	24.9	25.0	31.0	33.8	32.5	33.1	22.6	19.5	28.5	21.5	32.7	31.3	20.7	26.5	31.3	8.6	15.2	26.0	15.1	17.4	30.1	35.4	37.4	56.7	32.3	34.2	56.7			
ceb_Latn	35.4	37.4	56.5	57.0	52.8	58.8	34.1	34.1	53.7	37.0	48.0	57.7	46.9	52.4	58.0	27.0	46.0	55.9	32.3	34.2	56.7	30.3	19.3	46.7	29.0	37.4	50.4			
jav_Latn	30.3	19.3	46.7	49.8	48.3	49.7	30.4	30.8	46.7	36.2	33.1	50.6	37.7	38.9	51.0	19.4	45.2	48.5	29.0	37.4	50.4	32.7	42.4	42.2	34.9	43.6	51.4			
sun_Latn	32.7	42.4	42.2	45.2	43.9	46.1	30.4	31.5	41.5	33.8	42.1	43.9	36.3	40.9	43.4	26.8	24.1	43.8	31.0	34.9	43.6	19.6	18.7	22.1	43.5	23.8	51.4			
ilo_Latn	19.6	18.7	22.1	43.5	50.5	51.6	24.3	26.8	45.6	26.8	39.4	51.3	34.7	45.2	52.3	19.2	31.0	47.1	22.7	26.1	51.4	15.4	31.3	31.9	28.8	30.0	58.0			
war_Latn	15.4	31.3	31.9	38.3	43.2	54.3	28.8	30.5	52.9	32.8	37.8	57.0	40.2	51.5	58.5	23.8	34.1	55.3	27.6	27.5	58.0	21.7	21.8	24.7	23.9	23.8	26.7			
bug_Latn	21.7	21.8	24.7	23.4	25.3	27.5	13.5	17.3	23.9	25.7	25.7	26.6	25.5	25.2	27.1	9.1	16.0	23.0	23.8	17.4	26.7	25.7	23.0	23.9	23.8	34.3	31.0	34.9		
pag_Latn	25.7	23.0	23.9	26.3	28.1	28.6	25.8	26.7	40.7	27.9	28.5	42.5	26.5	36.4	43.6	21.0	26.6	42.4	24.7	21.9	43.8	4.5	5.8	1.6	3.7	5.0	6.4			
shn_Mymr	4.5	5.8	1.6	3.7	5.0	6.4	4.0	2.2	29.5	4.8	6.3	30.0	5.5	0.7	31.4	1.9	2.7	26.0	3.3	16.5	28.7	28.1	19.6	17.9	46.1	52.5	55.6	52.2		
min_Latn	28.1	19.6	17.9	46.1	48.9	52.5	35.1	34.4	47.1	37.3	43.0	50.7	37.6	45.2	52.5	34.2	43.1	50.0	35.6	35.5	52.2	16.8	22.9	28.7	28.8	28.9	39.0			
ace_Latn	16.8	22.9	28.7	28.8	28.5	36.7	11.1	7.8	33.1	28.6	32.9	37.5	28.1	32.9	37.9	2.3	19.2	35.5	1.0	24.7	39.0	26.3	29.6	31.8	34.3	34.2	44.9			
ban_Latn	26.3	29.6	31.8	33.2	34.4	38.3	29.8	29.0	42.0	32.6	38.4	44.2	32.0	40.9	45.5	21.9	38.1	42.5	30.7	34.2	44.9	20.4	9.4	16.0	36.5	30.3	32.0			
bjn_Latn	20.4	9.4	16.0	36.5	30.9	35.1	35.2	30.9	23.5	37.0	35.8	26.6	36.6	34.0	32.9	28.8	9.1	7.7	34.3	30.3	32.0	10.2	6.2	4.6	8.8	1.8	13.5			
ace_Arab	10.2	6.2	4.6	8.8	1.8	5.5	8.3	1.7	11.7	10.6	1.6	6.8	10.2	1.6	13.4	9.9	0.9	2.3	9.4	1.4	13.5	10.5	7.5	8.0	10.3	3.6	20.6			
bjn_Arab	10.5	7.5	8.0	5.0	4.6	10.3	6.6	1.7	10.2	10.9	9.8	8.3	10.3	1.6	17.9	10.1	1.3	8.4	10.0	3.6	20.6	9.9	5.5	7.5	1.1	11.2	11.2			
min_Arab	9.9	5.5	7.5	3.5	1.1	1.1	2.7	2.6	1.9	10.4	9.2	1.2	6.9	0.9	18.2	8.3	1.0	2.6	9.4	1.4	11.2	28.5	27.7	27.1	27.7	41.6	41.6			
Avg.	28.5	27.7	<b>32.3</b>	37.6	36.1	<b>39.3</b>	27.8	27.0	<b>37.1</b>	32.0	35.9	<b>40.0</b>	33.2	33.6	<b>42.5</b>	22.9	28.8	<b>37.1</b>	27.1	27.7	<b>41.6</b>									

AFR																												
swa_Latn	51.1	47.0	52.7	61.0	58.4	59.0	30.8	38.1	45.0	53.8	55.2	56.6	55.7	52.2	57.9	9.0	30.7	41.7	11.0	29.4	52.0	16.4	10.2	10.9	17.0	32.5	41.8	
lin_Latn	16.4	10.2	10.9	17.0	32.3	31.2	8.4	10.6	34.0	13.5	19.9	39.5	13.6	17.7	41.7	7.4	16.0	37.5	10.4	11.7	41.8	16.6	8.6	12.1	25.9	22.8	18.5	
yor_Latn	16.6	8.6	12.1	25.9	22.8	23.1	18.9	10.8	16.9	10.5	12.1	18.7	10.1	8.3	18.6	7.0	5.9	16.4	10.7	4.9	18.5	13.1	16.5	11.2	11.2	15.4	15.4	
ful_Latn	13.1	16.5	11.2	6.1	13.4	12.3	8.3	9.5	13.4	8.1	9.9	14.4	11.3	8.8	15.7	4.8	8.3	13.0	10.3	6.2	15.4	27.9	20.8	29.2	34.4	35.5	33.9	
ibo_Latn	27.9	20.8	29.2	34.4	35.5	36.9	21.8	20.4	27.9	23.2	29.5	33.7	23.1	13.0	34.0	8.5	17.0	27.0	15.2	15.3	33.9	1.7	5.7	7.5	11.7	22.1	23.3	
orm_Latn	1.7	5.7	7.5	11.7	16.0	22.1	9.1	13.1	15.6	9.8	9.1	24.8	10.7	11.1	24.6	9.2	4.1	13.4	10.8	10.0	23.3	27.0	30.6	34.4	34.4	35.3	35.3	
som_Latn	27.0	30.6	34.4	40.7	40.3	40.6	19.3	24.2	29.7	21.0	31.3	36.3	20.1	23.6	36.2	16.6	26.0	27.4	20.4	26.4	35.3	15.3	12.7	11.5	11.5	11.5	11.5	
tso_Latn	15.3	12.7	11.5	17.1	29.1	36.7	10.2	12.1	25.6	11.1	14.8	37.7	13.0	14.4	40.7	6.2	8.4	24.1	8.7	8.9	34.0	14.4	14.4	14.4	14.4	14.4	14.4	
nya_Latn	14.6	14.4	34.8	43.5	43.3	43.1	9.4	13.0	28.5	12.3	30.7	37.2	18.7	25.5	37.9	5.4	13.5	26.9	9.2	14.7	34.5	12.4	11.7	17.2	17.2	17.2	17.2	
zul_Latn	14.2	15.3	36.4	48.6	47.5	48.3	13.9	17.0	27.4	14.6	25.9	36.0	18.9	25.2	38.1	7.3	10.1	26.8	9.7	11.6	37.5	12.4	11.7	17.2	17.2	17.2	17.2	
kin_Latn	12.4	11.7	17.2	18.4	33.3	34.9	8.8	11.3	15.7	12.2	22.1	26.8	15.5	19.6	29.7	7.3	10.4	13.5	8.2	10.2	24.6	14.3	14.3	14.3	14.3	14.3	14.3	
run_Latn	14.3	14.0	15.6	16.5	28.0	30.1	9.9	11.9	15.9	10.0	19.2	24.6	12.8	17.6	26.5	6.8	11.0	15.9	8.8	10.2	21.7	11.4	11.4	11.4	11.4	11.4	11.4	
sna_Latn	11.4	15.4	29.6	39.4	39.6	39.9	7.9	14.9	26.7	12.6	26.9	34.3	19.8	26.2	35.7	4.8	14.9	24.7	7.9	11.5	32.4	14.4	14.4	14.4	14.4	14.4	14.4	
xho_Latn	25.3	14.8	30.5	43.3	43.2	44.4	16.3	17.0	29.9	14.9	27.0	35.8	20.7	25.4	37.4	7.5	12.0	27.1	9.7	11.8	35.8	10.1	17.4	20.9	20.9	20.9	20.9	
tsn_Latn	10.1	5.6	17.4	30.2	40.9	43.2	8.6	12.6	29.5	10.5	21.1	40.0	17.3	24.0	41.9	5.5	8.8	29.5	7.9	9.4	38.0	1.7	1.7	1.7	1.7	1.7	1.7	
tir_Ethi	0.9	0.6	1.0	7.1	7.8	10.8	0.4	0.4	13.5	4.2	7.6	19.4	1.9	2.9	19.1	1.3	3.7	10.2	1.2	2.7	15.8	13.3	13.3	13.3	13.3	13.3	13.3	
kik_Latn	13.3	11.1	16.0	15.8	17.5	13.7	13.6	11.7	7.8	18.7	18.2	14.5	15.7	12.0	15.2	4.9	7.5	4.2	7.1	9.7	15.4	14.0	14.0	14.0	14.0	14.0	14.0	
kon_Latn	14.0	9.6	9.9	14.6	20.5	28.7	8.1	10.0	33.5	10.9	19.4	41.8	15.5	11.5	41.3	7.2	11.0	37.0	8.5	10.5	42.0	15.2	15.2	15.2	15.2	15.2	15.2	
lua_Latn	15.2	8.6	19.9	10.4	25.1	25.3	6.8	10.3	22.7	8.9	14.9	26.7	16.0	7.3	30.8	5.8	11.3	21.0	7.3	7.5	24.8	12.2	12.2	12.2	12.2	12.2	12.2	
umb_Latn	10.4	12.2	14.3	12.8	16.1	16.0	9.4	11.5	7.3	12.7	8.6	8.2	15.7	3.4	11.2	6.4	7.6	5.6	12.7	8.0	10.5	12.4	12.4	12.4	12.4	12.4	12.4	12.4
sot_Latn	11.4	13.6	30.9	43.9	4																							

---

**EMBSWAP Results on FLORES-200**


---

Model	PaLM2									Gemma2									Aya23					
	Size			XXS			S			2B			9B			27B			8B			35B		
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR
bho_Deva	18.2	8.0	25.2	34.2	34.8	38.0	24.1	16.7	34.5	29.8	31.6	39.4	29.0	31.5	39.1	25.1	3.2	37.7	27.8	9.9	39.7			
pbt_Arab	8.1	2.6	18.3	29.1	28.3	31.0	9.5	10.0	22.5	12.4	21.4	29.3	10.9	16.4	27.3	6.5	13.1	17.7	9.2	3.2	29.8			
guj_Gujr	35.4	4.1	9.7	48.5	39.4	45.6	29.9	25.8	41.9	35.9	40.9	47.0	39.2	33.6	45.9	16.7	14.6	42.1	29.2	12.8	46.0			
kan_Knda	28.1	25.2	36.1	47.1	44.6	46.5	10.5	9.5	37.1	27.8	36.4	45.6	37.6	22.7	46.0	9.2	24.1	39.3	15.4	20.4	46.4			
awa_Deva	2.9	11.0	38.6	35.7	40.5	39.1	30.2	13.8	34.9	38.5	40.1	42.3	37.6	37.1	42.5	8.4	11.4	39.5	36.1	9.5	42.3			
ory_Orya	11.7	3.2	8.5	34.4	21.7	25.8	14.2	7.7	17.9	13.0	26.4	37.2	15.3	6.6	35.7	12.1	3.8	11.0	11.0	4.4	31.0			
mal_Mlym	25.6	21.1	28.4	40.8	38.0	40.1	23.8	15.9	30.7	29.4	30.8	41.6	31.0	7.5	44.3	18.0	3.6	34.3	25.0	16.8	44.0			
pan_Guru	20.0	14.3	28.0	46.0	43.8	45.8	22.3	19.4	37.7	32.9	38.2	44.7	37.9	10.4	44.6	13.5	23.5	38.0	27.2	22.7	44.3			
hme_Deva	2.5	21.0	36.9	33.4	40.1	43.4	16.2	15.9	37.8	34.6	29.7	47.2	37.2	33.3	46.5	18.4	20.4	41.4	33.7	13.2	46.3			
npi_Deva	33.4	27.4	37.1	49.3	46.1	48.4	26.7	27.1	40.1	36.1	42.7	46.0	37.6	38.5	46.5	19.7	3.6	40.3	28.2	21.6	45.8			
asm_Beng	2.3	6.3	21.7	32.9	31.7	35.0	11.0	8.5	23.6	15.1	24.5	32.6	21.7	17.0	33.7	2.5	6.1	26.7	11.2	10.0	33.4			
mni_Beng	0.4	1.2	6.6	0.8	7.2	17.9	0.5	1.5	21.7	4.6	3.8	31.0	0.9	0.1	30.6	0.4	10.1	22.6	0.2	2.4	29.1			
bod_Tibt	2.2	1.0	0.7	0.5	5.2	7.8	0.8	3.6	10.2	1.1	7.2	19.0	2.5	1.0	18.5	3.6	4.7	9.6	4.6	4.5	15.2			
san_Deva	13.8	0.6	4.2	20.0	23.4	24.1	16.5	5.3	16.7	18.1	19.0	26.0	18.0	17.0	25.7	15.8	0.4	20.3	17.7	5.4	25.6			
Avg.	21.0	17.9	<b>26.7</b>	36.7	34.6	<b>38.7</b>	20.7	17.2	<b>32.5</b>	27.8	31.4	<b>40.5</b>	29.6	23.9	<b>40.6</b>	16.1	11.5	<b>33.4</b>	23.6	13.9	<b>40.1</b>			

Table 13: EMBSWAP results on FLORES-200 with zero-shot prompting. FL: LLMs instruction-tuned on FLAN mixture; ES: EMBSWAP; +LR: EMBSWAP with LoRA Adaptation.

---

**EMBSWAP Results on GSM8K-NTL**


---

Model	Gemma2									Gemma2									Gemma2							
	2B			9B			27B			2B			9B			27B			2B			9B			27B	
Size	FL	Es	FL	Es	FL	Es	IT	Es	IT	Es	IT	Es	MH	Es	MH	Es	MH	Es	MH	Es	MH	Es	MH	Es	MH	Es
asm_Beng	11.2	14.4	30.8	26.4	34.8	35.2	15.6	12.4	48.0	42.4	53.2	48.0	20.4	29.2	46.0	28.0	43.2	36.0								
bew_Latn	28.4	24.8	48.4	46.4	48.8	50.4	41.6	30.0	58.0	57.6	62.4	53.2	51.6	42.0	70.0	68.8	71.6	68.8								
bho_Deva	16.4	17.2	29.2	35.6	38.0	35.6	23.2	19.2	48.4	48.0	50.8	51.6	29.2	36.4	51.2	48.8	49.6	49.6								
doi_Deva	11.2	14.0	19.6	34.0	26.4	33.6	13.6	19.6	32.4	44.0	37.6	46.0	17.6	25.2	36.4	36.4	29.2	37.6								
div_Thaa	6.8	10.4	9.6	23.2	18.8	26.8	2.8	12.4	15.6	36.8	22.0	32.8	2.8	15.6	17.2	26.4	13.2	12.0								
dzo_Tibt	2.0	2.8	6.8	12.8	10.4	8.4	0.4	3.6	10.0	14.0	14.0	20.4	0.8	3.6	11.6	6.4	8.8	4.0								
efi_Latn	8.0	10.4	13.2	27.6	16.0	30.8	5.2	5.6	17.2	31.6	24.8	20.8	6.8	20.0	21.6	48.0	19.6	38.0								
gom_Deva	15.2	15.6	25.2	35.6	28.4	37.2	13.6	14.4	37.2	48.4	44.0	49.2	11.6	30.4	38.4	47.6	36.4	41.6								
ilo_Latn	8.4	15.6	24.8	35.2	28.8	39.2	10.4	11.2	30.8	46.4	41.2	44.8	12.8	22.4	42.8	57.2	43.6	56.4								
kri_Latn	10.8	8.0	23.2	26.0	28.8	29.6	12.0	8.4	32.0	34.8	36.0	36.4	19.2	22.8	43.2	43.6	40.0	42.8								
mai_Deva	13.6	14.4	29.2	35.2	37.2	38.8	23.6	17.6	47.2	48.4	48.8	52.0	25.2	32.4	46.8	50.8	46.0	47.2								
meo_Latn	24.4	24.4	46.0	45.2	45.2	47.2	37.6	30.0	58.4	56.4	60.0	56.4	58.4	40.8	74.8	71.6	72.0	68.0								
mfa_Arab	6.4	10.8	6.8	37.2	9.2	42.8	3.6	18.0	12.0	52.0	8.8	49.2	3.2	25.2	7.2	54.8	7.2	48.0								
min_Latn	12.0	18.8	32.4	38.0	35.6	44.8	15.6	22.0	42.8	52.8	49.6	44.4	22.0	29.2	45.2	58.0	54.8	53.6								
mini_Beng	4.8	5.2	4.8	16.8	8.8	17.6	1.6	6.0	5.2	20.8	8.4	25.2	5.6	10.0	4.8	17.6	8.0	4.8								
mzn_Arab	21.6	20.0	40.0	42.8	46.8	48.4	29.2	19.2	56.4	56.4	60.0	59.6	33.6	32.8	65.6	59.2	61.2	52.4								
nso_Latn	4.4	10.8	14.8	26.0	21.2	24.0	4.4	7.2	19.2	26.8	25.6	13.6	7.2	10.8	24.4	36.4	22.4	23.2								
ory_Orya	10.0	12.8	19.2	26.4	24.8	27.2	2.8	11.6	29.6	34.8	37.2	42.4	6.4	21.6	33.6	26.4	27.2	23.6								
pem_Latn	28.0	23.6	46.8	43.2	48.4	50.0	43.6	36.0	60.0	60.0	61.6	51.6	61.2	62.8	77.6	72.4	77.6	74.4								
tso_Latn	8.8	9.6	12.4	16.4	17.2	20.4	4.0	5.2	14.4	24.8	22.4	10.0	7.6	8.8	19.6	24.8	17.6	20.0								
Avg.	12.6	<b>14.2</b>	24.2	<b>31.5</b>	28.7	<b>34.4</b>	15.2	<b>15.5</b>	33.7	<b>41.9</b>	38.4	<b>40.4</b>	20.2	<b>26.1</b>	38.9	<b>44.2</b>	37.5	<b>40.1</b>								

Table 14: EMBSWAP results on GSM8K-NTL with zero-shot prompting. FL: LLMs instruction-tuned on FLAN mixture; IT: LLMs aligned with supervised fine-tuning and reinforcement learning with human feedback; MH: LLMs instruction-tuned on the WebInstruct math dataset; ES: EMBSWAP;

---

**EMBSWAP Results on XSUM-IN**


---

Model	PaLM2									Gemma2									Aya23					
	Size			XXS			S			2B			9B											

EMBSWAP Results on XSUM-IN

Model	PaLM2						Gemma2						Aya23											
	XXS			S			2B			9B			27B			8B			35B					
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR
gom_Deva	0.2	0.2	0.2	0.2	7.3	3.7	0.4	5.0	10.3	14.2	14.9	17.5	12.1	17.4	18.5	13.2	0.3	10.5	0.4	6.4	15.2			
guj_Gujr	11.8	0.4	0.4	25.1	24.4	20.4	6.7	11.8	16.0	18.8	21.1	19.4	22.1	16.3	19.7	13.5	11.2	16.3	16.8	0.2	19.1			
hin_Deva	22.1	1.3	2.4	31.2	30.5	27.6	4.7	14.3	19.2	24.2	24.4	24.9	23.6	18.6	23.7	26.4	11.2	18.2	27.3	17.7	26.3			
hne_Deva	0.2	0.2	0.3	3.5	20.1	20.5	2.2	2.8	5.8	20.2	18.3	17.9	13.3	16.5	18.4	18.7	0.0	2.0	8.6	10.3	17.2			
hoj_Deva	0.2	0.2	0.1	0.6	0.5	0.5	1.1	1.0	0.5	1.9	14.9	15.7	3.1	4.6	12.8	0.8	0.1	9.8	0.3	0.1	0.7			
kan_Knda	4.4	2.7	0.7	28.6	28.3	28.9	7.9	4.9	14.9	17.9	21.1	23.4	24.0	8.5	24.3	8.2	7.2	15.8	13.4	12.4	22.1			
mai_Deva	0.1	0.1	0.1	13.7	19.1	17.3	2.4	1.6	4.0	19.0	18.2	17.6	19.1	15.2	17.1	14.5	0.1	11.8	3.4	5.3	15.7			
mal_Mlym	13.5	2.1	2.0	27.2	27.1	27.7	5.0	8.6	15.4	21.2	21.5	22.4	23.4	4.0	22.4	16.9	3.2	14.1	19.2	10.9	21.2			
mmi_Beng	0.2	0.2	0.2	0.2	8.3	4.9	0.1	1.7	3.1	4.9	1.4	8.7	0.1	0.0	10.6	0.8	1.8	3.0	1.3	0.3	8.6			
mar_Deva	16.6	10.9	1.0	30.5	28.9	27.4	13.7	14.1	16.7	19.7	19.1	20.7	21.2	21.8	19.9	17.6	5.1	16.1	21.5	11.2	21.0			
mup_Deva	0.3	0.3	0.3	0.4	0.4	0.5	0.5	0.4	2.8	2.5	9.6	14.9	9.3	14.5	16.0	8.5	0.0	0.3	2.3	3.0	0.9			
mwr_Deva	0.2	0.2	0.2	0.3	1.4	1.9	0.8	1.5	2.7	19.2	8.3	18.8	9.3	17.4	18.6	9.2	0.1	0.2	3.4	1.0	12.7			
npi_Deva	2.3	0.0	0.0	29.1	23.9	5.3	9.3	15.4	17.3	21.5	24.1	22.8	11.1	23.5	20.2	18.5	3.5	16.5	20.7	12.7	19.5			
ory_Orya	0.6	0.6	0.6	23.8	18.8	22.2	3.4	8.5	12.9	13.6	14.8	17.7	24.6	5.7	16.8	12.8	2.6	3.7	14.2	5.0	13.8			
pan_Guru	7.9	1.3	0.2	21.5	23.3	22.6	3.7	11.8	12.4	17.4	18.8	18.3	16.2	1.6	21.0	8.2	9.7	15.8	14.5	9.4	20.4			
pbu_Arab	0.2	0.2	0.2	6.2	13.5	4.3	1.8	7.2	10.1	11.4	15.3	18.0	16.4	11.0	18.3	7.6	11.7	15.2	5.7	3.1	18.2			
san_Deva	0.1	0.1	0.1	19.8	17.7	19.3	2.4	5.4	11.7	16.3	15.3	16.1	8.8	15.1	15.8	16.4	0.0	4.2	15.9	3.6	8.1			
sat_Olck	0.1	0.1	0.1	0.0	0.1	0.1	0.1	0.1	0.1	0.0	0.0	0.0	15.3	0.0	0.1	0.0	0.0	0.0	0.1	0.0	0.0			
tam_Taml	7.0	12.7	1.0	35.5	34.5	35.3	4.3	10.0	9.0	27.2	28.0	26.7	0.0	12.4	27.2	21.7	11.5	18.9	25.2	16.8	26.6			
tel_Telu	16.2	12.5	12.2	26.3	26.1	26.7	13.4	13.9	13.5	20.5	18.6	21.6	28.9	7.6	22.5	14.5	8.9	13.5	17.2	11.4	21.3			
urd_Arab	12.8	6.3	2.3	29.1	28.8	26.5	1.9	17.1	14.9	20.5	24.1	24.0	22.6	21.6	24.8	13.6	8.8	19.7	16.9	7.4	25.0			
Avg.	<b>5.8</b>	3.1	2.0	16.6	<b>18.1</b>	16.1	4.8	7.5	<b>10.6</b>	15.8	15.9	<b>18.1</b>	17.6	13.2	<b>18.6</b>	<b>12.5</b>	4.6	10.1	11.3	7.8	<b>16.4</b>			

Table 15: EMBSWAP results on XSUM-IN with zero-shot prompting. FL: LLMs instruction-tuned on FLAN mixture; ES: EMBSWAP; +LR: EMBSWAP with LoRA Adaptation.

Model	PaLM2						Gemma2						Aya23											
	XXS			S			2B			9B			27B			8B			35B					
Variant	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR
ANSWER IN EN																								
eng_Latin	70.9	75.5	74.8	75.1	82.1	82.6	71.2	66.6	73.3	75.7	74.6	74.4	75.5	78.5	75.9	75.6	72.6	71.0	78.4	65.2	74.9			
asm_Beng	47.9	63.0	68.4	52.0	51.1	58.3	54.3	61.3	67.7	71.6	72.4	71.2	65.7	64.3	74.1	56.7	67.2	57.9	64.8	69.2	73.9			
awa_Deva	48.4	64.0	64.6	49.7	52.9	57.7	49.6	54.2	57.9	63.2	66.3	68.1	63.3	58.8	70.7	68.3	64.9	63.2	77.1	67.6	72.9			
bgc_Deva	43.4	61.1	67.0	52.9	59.5	58.1	45.1	53.0	60.4	57.9	68.6	70.9	60.8	57.1	73.7	70.1	69.0	66.4	74.6	67.9	71.1			
bho_Deva	50.8	64.8	69.7	47.4	58.5	58.5	50.4	57.5	59.4	63.8	67.2	67.0	62.0	60.8	72.8	68.8	67.2	65.8	73.2	65.1	69.4			
ben_Beng	57.3	62.5	71.0	52.0	58.5	63.1	59.8	64.1	67.7	74.5	74.2	72.6	65.9	64.9	73.7	68.4	69.0	63.1	76.6	68.6	70.2			
bod_Tibt	3.2	2.6	17.5	23.5	12.4	38.5	11.4	36.0	40.9	27.5	41.7	33.8	43.5	44.7	60.8	32.0	30.1	23.5	30.3	36.4	30.6			
brx_Deva	5.3	45.6	58.1	25.6	57.6	61.7	12.4	40.9	49.5	19.6	56.5	61.2	25.2	48.1	60.5	28.2	47.4	52.4	28.0	46.3	63.4			
gbm_Deva	40.9	59.1	66.1	53.0	53.6	60.1	46.3	51.1	57.6	56.0	66.0	68.4	56.0	56.6	71.6	61.3	64.6	64.5	66.3	62.8	68.8			
gom_Deva	38.7	64.8	69.5	50.4	67.3	66.9	33.8	54.6	63.0	58.8	69.2	69.8	60.7	55.6	71.3	44.1	65.8	61.6	51.0	62.9	72.1			
guj_Gujr	45.1	70.3	72.7	44.0	43.5	54.8	60.4	61.6	67.4	67.9	70.3	59.8	72.4	66.8	75.2	57.1	68.2	64.4	73.4	70.2	75.6			
hin_Deva	61.8	72.5	73.2	55.0	47.5	52.1	68.8	66.7	68.7	72.6	71.8	64.9	65.5	57.7	75.4	74.7	70.4	67.0	76.8	71.5	76.8			
hne_Deva	49.0	61.8	66.9	47.3	58.5	58.3	45.2	51.6	59.8	58.3	65.8	67.8	65.8	59.3	71.4	65.6	67.6	61.2	73.8	64.2	70.4			
hoj_Deva	40.1	56.7	60.4	48.6	50.8	56.2	41.3	48.5	54.3	54.7	60.7	64.7	56.8	61.0	70.7	59.5	62.0	56.9	69.0	61.0	68.6			
kan_Knda	45.0	66.7	69.6	44.9	53.9	57.9	65.0	63.3	68.2	69.8	71.2	63.4	70.1	63.3	70.1	60.6	68.3	61.6	66.6	71.4	75.1			
mai_Deva	49.0	71.9	72.8	49.3	50.9	55.2	50.2	59.4	63.9	63.6	69.2	64.9	58.1	72.7	67.0	68.4	64.7	69.8	66.7	72.1				
mal_Mlym	51.8	69.9	71.0	56.8	57.1	63.7	68.1	69.2	72.5	74.7	71.2	67.7	70.3	63.3	74.4	71.8	66.9	61.2	78.0	74.0	76.9			
mmi_Beng	9.8	28.6	40.2	37.0	46.6	56.5	20.6	27.3	31.6	29.4	43.2	48.9	33.6	47.8	50.2	30.8	39.3	34.7	32.1	36.8	44.7			
mar_Deva	55.8	68.1	71.5	43.6	50.0	55.4	57.9	61.6	65.2	70.7	71.3	67.0	62.1	55.7	71.8	63.8	65.0	5						

		EMBSWAP Results on XORQA-IN																			
Model		PaLM2						Gemma2						Aya23							
Size		XXS			S			2B			9B			27B			8B				
Variant		FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR	FL	ES	+LR		
bho_Deva	5.8	3.4	5.3	20.4	22.1	17.7	1.6	3.7	2.8	1.5	2.0	2.3	16.8	13.4	9.0	5.5	2.6	2.6	7.3	2.1	8.8
ben_Beng	7.4	2.5	0.6	7.5	16.6	9.6	0.8	1.1	1.8	1.8	1.8	1.4	13.5	9.3	5.0	0.7	6.2	4.3	0.7	1.8	3.2
bod_Tibt	0.2	0.9	0.7	7.2	4.6	7.3	1.1	5.3	6.3	3.8	5.7	4.0	5.2	4.9	5.3	1.2	0.6	1.1	2.7	1.1	4.3
brx_Deva	0.8	13.8	17.7	15.0	22.5	22.7	3.1	15.4	17.9	6.5	18.0	19.9	9.3	15.8	20.0	6.4	11.0	16.9	6.6	7.0	18.5
gbm_Deva	10.4	0.9	2.5	20.2	8.3	8.5	11.5	12.9	13.3	12.3	13.4	13.3	18.1	15.1	17.4	16.3	6.7	13.1	16.2	7.6	16.2
gom_Deva	1.4	1.6	3.1	9.2	7.8	2.5	1.1	2.9	2.2	1.6	6.0	5.9	14.9	11.2	8.8	5.7	2.6	3.7	2.3	1.8	5.4
guj_Gujr	15.6	12.3	20.4	33.8	25.6	26.2	9.8	11.6	13.1	17.0	18.7	17.2	26.6	16.8	19.5	10.2	12.9	10.6	15.5	7.1	13.4
hin_Deva	40.7	28.1	35.2	59.5	51.4	51.0	23.0	25.3	25.6	25.3	27.3	28.5	45.0	37.6	37.9	21.4	22.1	21.8	22.7	18.1	25.4
hne_Deva	15.9	20.6	20.7	29.3	30.2	28.2	13.5	13.5	13.9	16.0	18.3	17.6	31.9	25.8	25.3	21.0	9.0	14.1	23.5	10.3	19.6
hoj_Deva	14.3	16.4	17.6	28.3	28.6	25.4	14.8	17.9	18.4	15.5	20.8	20.6	33.4	30.4	29.3	20.0	17.7	18.4	21.8	13.8	20.6
kan_Knda	10.8	11.6	10.4	29.4	23.1	18.2	11.3	8.8	12.3	11.8	14.2	18.1	32.2	15.6	25.0	8.4	10.3	10.9	10.6	11.9	20.6
mai_Deva	9.4	1.6	4.3	30.0	26.9	24.1	10.8	10.8	11.6	12.2	18.9	13.7	27.1	23.0	18.8	18.0	7.1	13.8	14.1	8.0	16.3
mal_Mlym	21.8	25.5	27.4	35.8	41.1	37.7	20.6	18.6	21.8	30.3	26.3	29.9	35.1	19.7	32.7	21.6	19.2	19.3	26.1	15.1	25.7
mni_Beng	0.8	0.3	0.2	1.1	10.3	4.2	0.6	0.1	0.3	0.7	0.6	0.6	1.5	0.4	2.0	1.6	1.1	2.3	0.4	0.2	3.6
mar_Deva	31.1	30.6	32.9	38.5	43.2	36.0	22.3	22.4	25.6	24.6	27.2	23.5	40.1	31.0	30.3	24.8	17.7	20.8	34.5	19.6	25.9
mup_Deva	4.6	16.9	16.4	12.4	21.5	23.8	5.3	5.8	5.4	5.9	6.4	7.3	14.8	11.6	10.0	6.9	3.9	5.0	5.8	3.6	6.2
mwr_Deva	2.4	1.6	3.3	10.5	12.1	9.7	1.6	1.7	2.4	4.2	8.7	5.8	14.8	11.7	10.8	6.7	2.2	5.8	6.0	2.0	4.0
ory_Orya	5.7	5.6	6.1	15.4	14.1	11.6	6.6	6.5	9.9	11.1	12.0	17.0	17.7	12.4	17.4	10.9	6.4	8.1	16.0	8.9	12.4
pan_Guru	8.6	7.1	10.8	20.7	11.5	11.5	4.7	5.8	7.0	7.4	10.6	12.8	21.9	7.2	16.4	5.3	12.3	11.0	7.0	5.6	12.7
pbu_Arab	3.8	6.9	6.8	8.2	8.1	7.5	1.2	2.9	3.6	3.2	5.1	3.5	7.4	8.8	6.0	3.2	4.5	4.9	3.9	3.5	8.1
san_Deva	3.1	4.9	5.1	29.7	15.9	17.4	1.3	2.7	3.5	7.6	12.0	11.4	19.4	11.1	11.8	10.3	1.7	4.4	7.8	3.0	10.3
sat_Olck	1.3	0.2	0.1	7.2	0.2	0.4	0.0	0.0	0.1	0.0	0.0	2.2	0.0	0.0	0.6	0.8	0.6	1.0	0.3	0.0	0.8
tam_Taml	16.8	20.5	19.0	23.8	32.7	27.1	13.9	12.8	15.0	17.2	24.6	17.7	29.8	16.8	22.6	14.5	14.1	12.3	17.5	11.8	17.3
tel_Telu	10.7	23.6	23.1	23.6	38.6	31.7	12.4	12.3	14.5	13.9	15.5	15.3	21.2	12.9	18.1	14.0	11.0	12.6	14.0	8.8	15.2
urd_Arab	8.4	6.7	9.3	17.6	22.7	17.4	1.7	2.5	3.5	4.1	11.3	4.1	13.4	16.5	9.3	4.3	2.9	3.3	1.7	1.9	4.0
Avg.	9.7	10.4	<b>11.7</b>	<b>21.3</b>	21.1	18.9	7.6	8.6	<b>9.8</b>	9.8	<b>12.5</b>	12.1	<b>20.4</b>	15.2	16.4	<b>10.5</b>	7.8	9.6	11.4	6.8	<b>12.6</b>

Table 16: EMBSWAP results on XORQA-IN with zero-shot prompting. FL: LLMs instruction-tuned on FLAN mixture; ES: EMBSWAP; +LR: EMBSWAP with LoRA Adaptation.