

# Only for the Unseen Languages, Say the Llamas: On the Efficacy of Language Adapters for Cross-lingual Transfer in English-centric LLMs

Anonymous ACL submission

## Abstract

Most state-of-the-art large language models (LLMs) are trained mainly on English data, limiting their effectiveness on non-English, especially low-resource, languages. This study investigates whether language adapters can facilitate cross-lingual transfer in English-centric LLMs. We train language adapters for 13 languages using Llama 2 (7B) and Llama 3.1 (8B) as base models, and evaluate their effectiveness on two downstream tasks (MLQA and SIB-200) using either task adapters or in-context learning. Our results reveal that language adapters improve performance for languages not seen during pretraining, but provide negligible benefit for seen languages. These findings highlight the limitations of language adapters as a general solution for multilingual adaptation in English-centric LLMs.

## 1 Introduction

Most state-of-the-art LLMs are English-centric (Touvron et al., 2023; Jiang et al., 2023). To illustrate, in Llama 2 (Touvron et al., 2023), English constitutes 90% of the pre-training data. Despite this data imbalance, recent English-centric LLMs exhibit some multilingual capabilities (Kew et al., 2024; Ye et al., 2023). However, these capabilities are inconsistent across languages and tasks, with low-resource languages being particularly affected (Razumovskaia et al., 2024).

To endow LLMs with more profound multilingual capabilities, cross-lingual transfer (XLT) has emerged as a prevalent paradigm, aiming to transfer task-specific knowledge from a high-resource source language to a lower-resource target language, thereby alleviating the constraint of having supervised task data (Philippy et al., 2023).

As LLMs grow larger and full fine-tuning becomes less feasible, parameter-efficient fine-tuning (PEFT) methods have been explored for XLT (Houlsby et al., 2019; Hu et al., 2021). One com-

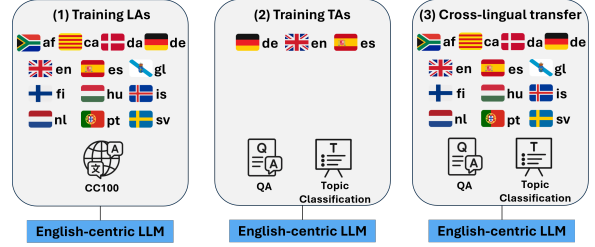


Figure 1: To evaluate cross-lingual transfer, language adapters (for 13 languages) and task adapters (for 3 high-resource source languages) are trained on top of a frozen English-centric LLM. Task adapters are evaluated on all languages of interest on two selected tasks.

mon setup for enhancing XLT abilities is to combine small language and task adaptation modules, as introduced by Pfeiffer et al. (2020b). The authors propose language adapters (LAs) and task adapters (TAs), parameter-efficient modules that are trained on top of a frozen base LLM and capture language- and task-specific representations, respectively.

While LAs have been extensively evaluated for small-scale multilingual LLMs (Pfeiffer et al., 2020b; Parović et al., 2022; Rathore et al., 2023; Yong et al., 2023), there is only a paucity of work that assesses its applicability to large-scale English-centric LLMs (Lin et al., 2024; Razumovskaia et al., 2024). Our work closes this gap by making the following contributions:

1. We evaluate in a systematic manner whether LAs help enhance XLT abilities of English-centric LLMs across 13 linguistically diverse languages and two tasks (one QA and one NLU task) to inspect the impact of typological relatedness and task-related intricacies.
2. We conduct a detailed analysis of the variables critical for successful XLT in English-centric LLMs by comparing different task adaptation methods (TAs vs. in-context learning (ICL)) and base LLMs (Llama 2 vs. Llama 3.1).

Our main findings on English-centric LLMs uncover that (1) surprisingly, **LAs are beneficial exclusively for languages that are unseen** during pretraining, while (2) they are **at best redundant for rarely seen languages**; and (3) that - in contrast to previous findings on multilingual models - the typological relatedness of languages for language transfer has only **a minimal effect**.

## 2 Related Work

**Language Adapters.** LAs represent a parameter-efficient and modular method for language adaptation (Poth et al., 2023). They are added to a frozen base LLM and typically trained on monolingual, unsupervised data using a language modeling objective in order to learn language-specific representations (Pfeiffer et al., 2020a). In general, any adapter architecture can be utilized for LA training: Prior work on small-scale, multilingual base LLMs has primarily employed *bottleneck adapters* (Houlsby et al., 2019) for LA training (Pfeiffer et al., 2020b; Parović et al., 2022; Faisal and Anastasopoulos, 2022; Yong et al., 2023; Gurgurov et al., 2024). They observed enhanced XLT, particularly for lower-resource languages. However, Kunz and Holmström (2024) find that the effect of LAs varies considerably across target languages and omitting LAs is beneficial in some cases. More recent work that employs large-scale, English-centric base LLMs prefers *LoRA adapters* (Hu et al., 2021) for LA training (Lin et al., 2024; Razumovskaia et al., 2024), arguably due to the inference latency that bottleneck adapters introduce, which LoRA helps mitigate by merging its weights with the base LLM’s weights (Hu et al., 2021). An alternative strand of work made use of other PEFT methods such as soft prompts for XLT (Philippy et al., 2024; Vykopal et al., 2025)

### Cross-lingual transfer in English-centric LLMs.

Previous work evaluating XLT in English-centric LLMs can be roughly divided into two approaches: *one-stage* XLT, which omits LAs entirely and applies task adaptation only, and *two-stage* XLT, in which LAs are trained prior to task adaptation.

*One-stage* XLT. Three task adaptation methods can be distinguished: In (1), single-task TAs are trained followed by an ICL<sup>1</sup> evaluation at inference. Ye et al. (2023) show that minimal pre-training data for a given target language suffices to enable

successful zero-shot XLT. In (2), ICL is applied exclusively. Asai et al. (2024) and Ahuja et al. (2024) establish XLT ICL benchmarks, revealing that English-centric LLMs perform well in high-resource languages but struggle with low-resource languages. Finally, in (3), multi-task instruction tuning (IT) is employed to fine-tune a base LLM, followed by ICL at inference. Previous work finds that multilingual IT with only a few languages (Aggarwal et al., 2024; Kew et al., 2024; Chen et al., 2024), or even monolingual IT in English (Chirkova and Nikoulina, 2024), suffices to elicit robust XLT abilities. In this study, we omit multi-task IT and focus on a comparison between single-task TAs and ICL.

*Two-stage* XLT. Lin et al. (2024) train a single LA covering 534 languages. They report performance gains for languages with low-resource scripts while performance drops for high-resource languages. Razumovskaia et al. (2024) train language-specific LAs and emphasize that performance improvements over setups without LAs are limited to NLG tasks. Kunz (2025) conducts a case study on Icelandic summarization, comparing several PEFT methods for language adaptation. It is shown that LoRAs situated in the feed-forward layers and bottleneck adapters yield the largest performance improvements.

## 3 Experimental Setup

Unlike most previous work that assessed the XLT abilities of English-centric LLMs, we begin by adapting the XLT setup as commonly employed for *multilingual* LLMs, i.e., we train LAs and TAs. Figure 1 illustrates our training and evaluation pipeline, including the selected languages and tasks. Subsequently, we study the effect of the task adaptation method and the base LLM, resulting in three different XLT configurations.

### 3.1 Models

The open-weights LLMs Llama 2 7B (Touvron et al., 2023) and its successor Llama 3.1 8B (Dubey et al., 2024) are selected as base LLMs. Both models are decoder-only, autoregressive LLMs. Despite the limited non-English pre-training data (2% in Llama 2 and 5% in Llama 3.1<sup>2</sup>), the models have demonstrated certain XLT abilities when fine-tuned for specific tasks (Ye et al., 2023) or evaluated using ICL (Asai et al., 2024; Ahuja et al., 2024).

<sup>1</sup>Following Li (2023), ICL encompasses any learning without parameter updates, including zero-shot evaluation.

<sup>2</sup>See Appendix B for a detailed language distribution.

### 3.2 Adapter Method

In this study, we use *bottleneck adapters*<sup>3</sup> as proposed by Pfeiffer et al. (2020b) to train LAs and TAs (see Appendix A for details). This method injects trainable adapter layers into the frozen base LLM, consisting of a down- and an up-projection which are situated after the feed-forward block of each transformer layer. Crucially, this architecture allows composition, i.e., multiple bottleneck adapters can be easily stacked on top of each other.

### 3.3 Data

**Language Data** Following previous work (Pfeiffer et al., 2022; Kunz, 2025), this work trains LAs on monolingual, unlabeled data extracted from CC-100, a multilingual, web-crawled corpus created by Conneau et al. (2020) for XLM-R pre-training. All LAs are trained on the first 200k<sup>4</sup> CC-100 samples of the respective language. While not explicitly stated, it is likely that CC-100 was seen during Llama 2 and 3.1 pre-training. Thus, the models are not necessarily trained on new data but rather *primed* towards the respective target languages.

**Task Data** We evaluate the effect of LAs based on model performance on one Question Answering (QA) and one NLU downstream task. For QA, we use *MLQA-en* (*T*) (henceforth *MLQA*), an extractive QA dataset from the Aya Collection (Singh et al., 2024), that extends the English subset of MLQA (Lewis et al., 2020) with translations into 100 languages. F1 as implemented for SQuAD (Rajpurkar et al., 2018) is used as evaluation metric.

For NLU, *SIB-200* (Adelani et al., 2024) is selected, a topic classification dataset with seven labels. Exact Match (EM) is used as evaluation metric.<sup>5</sup> These datasets were chosen primarily for their extensive language coverage and availability of parallel data. Given the use of autoregressive LLMs, both tasks - though not inherently generative - are framed as generation problems; that is, we generate targets (see Appendix D for task templates).

### 3.4 Languages

The set of languages comprises 13 Latin-script languages from three language groups. We exam-

ine seven Germanic languages (English, German, Dutch, Swedish, Danish, Icelandic, Afrikaans), four Romance languages (Spanish, Portuguese, Catalan, Galician), and two Finno-Ugric languages (Finnish, Hungarian). In each XLT setup, one language is selected as the source language, with the remaining ones as target languages.

All experiments use English, German, and Spanish as source languages. English serves as a reference, given its frequent use as source language (Pfeiffer et al., 2020b; Parović et al., 2022). Due to data availability and based on the assumption that higher-resource languages transfer more effectively than lower-resource languages (Senel et al., 2024), German and Spanish are chosen as non-English source languages. Each source language is evaluated on all 13 target languages.

### 3.5 Training and Evaluation Settings

We include three main experiments, each of which essentially compares two XLT setups:

- (1) *noLA* employs one-stage XLT, i.e., omits LAs entirely and relies only on task adaptation. Thus, this setup relies on cross-lingual representations that emerge during pre-training.
- (2) *LA* employs two-stage XLT, i.e., trains LAs prior to task adaptation. Thus, this setup relies on strengthening cross-lingual representations after pre-training through LAs.

We hypothesize that if LAs show a positive effect, *LA* should outperform *noLA* which serves as a baseline. We define a base configuration in Experiment 1 and modify one variable at a time, resulting in Experiment 2 and 3.

**Experiment 1: Llama-2/TA** We adapt the MAD-X framework (Pfeiffer et al., 2020b) to English-centric LLMs (see Appendix E for a detailed walk-through example): As for the *LA* setup, language-specific LAs for all relevant languages are trained on top of frozen Llama 2. Next, a TA in the selected source language is trained on top of the frozen source LA. At inference, XLT is evaluated zero-shot by replacing the source LA with the target LA while retaining the source TA. As for the *noLA* setup, only a TA is trained in the source language, then evaluated zero-shot in the target languages.

**Experiment 2: Llama-2/ICL** We modify the task adaptation method: Instead of task-specific TAs, we use ICL and craft a prompt, consisting

<sup>3</sup>In preliminary experiments, we observed that *prompt tuning* (Lester et al., 2021) and *LoRA* (Hu et al., 2021) underperform.

<sup>4</sup>Doubling the number of LA training samples to 400k did not yield any performance gains.

<sup>5</sup>We cut off generations after the first word to account for verbose model outputs.

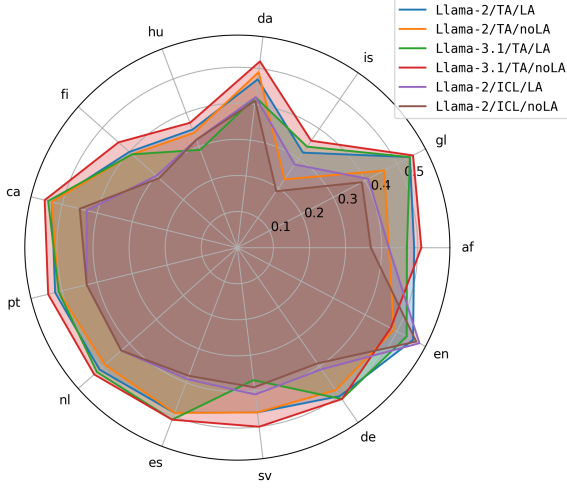


Figure 2: MLQA F1 scores for all target languages averaged across the three source languages *en*, *de*, *es* for all experiments.

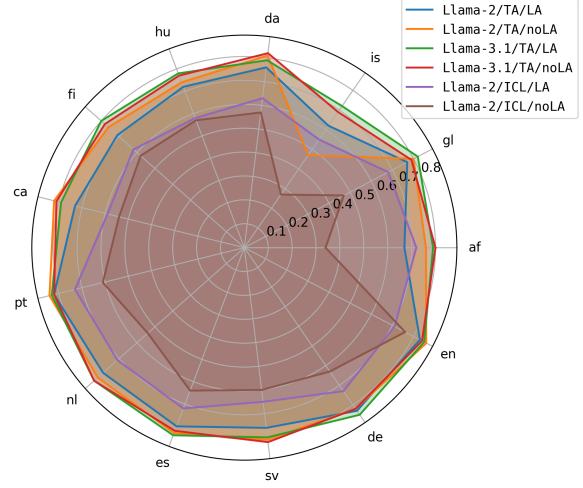


Figure 3: SIB-200 EM scores for all target languages averaged across the three source languages *en*, *de*, *es* for all experiments.

of five and ten randomly sampled source language demonstrations for MLQA and SIB-200, respectively,<sup>6</sup> followed by the test instance in the respective target language (see Appendix D.2 for the full prompt templates). Hence, we reduce the required computational cost, as only LAs need to be trained. We also address issues that may arise from stacking adapters. Again, we utilize Llama 2 as base LLM.

**Experiment 3: Llama-3.1/TA** We modify the base LLM and replace Llama 2 by Llama 3.1, potentially benefiting from more multilingual corpora. We train TAs for task adaptation. LAs and TAs are trained similar to Experiment 1.

## 4 Results and Analysis

In the following section, the findings of the three experiments are presented and discussed. Full scores are reported in Tables 5 to 10 in Appendix F. We use italic *en*, *de*, *es* to denote the source language of a specific setup, i.e., ‘with *en*’ means ‘with English as source language’.

### 4.1 General Findings

LAs do not consistently enhance XLT across target languages and tasks; they are often redundant or harm performance. Table 5 and 6 show that even for the source languages themselves, *noLA* outperforms or is on par with *LA*. This aligns with prior work (Kunz and Holmström, 2024; Oji and Kunz,

2025), which reports inconsistencies across languages and tasks in multilingual LLMs, as well as performance degradation with LAs in some cases.

As a topic classification task, SIB-200 requires less language-specific knowledge than the extractive QA task MLQA, where more fine-grained language understanding is necessary. This is reflected in Table 1 which shows that models generally achieve substantially better performance on SIB-200 than on MLQA with a less pronounced gap between English and non-English languages.

Regarding target-language related differences, Figures 2 and 3 show that Finnish Hungarian and Icelandic (summarized as *IsFiHu*) perform the worst across tasks. We attribute the poor performance of *IsFiHu* to a misaligned vocabulary. Due to their typological distance from English, languages like *IsFiHu* may lack language-specific tokens in the English-centric vocabulary. This leads to a less efficient tokenization<sup>7</sup> which in turn results in a suboptimal flow of input through the model and a decreased downstream task performance as similarly shown by Ali et al. (2024).

### 4.2 Experiment 1: Base

**MLQA.** As Table 5 illustrates, target languages unseen during Llama 2 pre-training (i.e., Afrikaans, Galician and Icelandic) benefit most from the usage of LAs across all source languages. Regarding seen languages, LAs do not reveal a discernible pattern. As Figure 2 shows, with *en* and *de*, LAs tend to

<sup>6</sup>First experiments revealed that for SIB-200, five demonstrations result in an overreliance on the label *geography*.

<sup>7</sup>Indicated by higher fertility (token/word ratio) scores in Table 4 in Appendix C.



Setup	Llama-2/TA						Llama-2/ICL						Llama-3.1/TA					
	MLQA			SIB-200			MLQA			SIB-200			MLQA			SIB-200		
	non-en	en	avg.	non-en	en	avg.	non-en	en	avg.	non-en	en	avg.	non-en	en	avg.	non-en	en	avg.
$LA_{en}$	0.44	<b>0.78</b>	<b>0.47</b>	0.67	0.85	0.68	<b>0.40</b>	0.57	<b>0.42</b>	0.59	0.75	0.60	0.42	<b>0.80</b>	0.46	0.79	<b>0.88</b>	0.80
$LA_{de}$	<b>0.47</b>	0.44	<b>0.47</b>	0.80	0.83	0.80	0.39	0.57	0.40	<b>0.68</b>	0.73	<b>0.68</b>	0.46	0.37	0.46	0.80	0.80	0.80
$LA_{es}$	0.44	0.43	0.44	0.76	0.82	0.76	0.35	0.49	0.36	<b>0.68</b>	0.65	0.67	0.42	0.43	0.42	<b>0.81</b>	0.86	<b>0.81</b>
$noLA_{en}$	0.45	<b>0.78</b>	<b>0.47</b>	0.74	<b>0.86</b>	0.75	0.37	<b>0.64</b>	0.40	0.37	<b>0.78</b>	0.41	0.46	0.79	0.48	0.79	0.83	0.79
$noLA_{de}$	0.44	0.38	0.44	<b>0.81</b>	<b>0.86</b>	<b>0.81</b>	0.36	0.56	0.38	0.61	<b>0.78</b>	0.63	<b>0.51</b>	0.35	<b>0.50</b>	0.79	0.84	0.79
$noLA_{es}$	0.38	0.32	0.38	0.76	0.85	0.77	0.35	0.47	0.36	0.54	0.73	0.56	0.46	0.31	0.45	<b>0.81</b>	0.84	<b>0.81</b>

Table 1: Average scores across experiments and tasks, containing scores across non-English languages (non-en), English (en) and all languages (avg.). For MLQA, F1 scores and for SIB-200, EM scores are reported. Bold numbers indicate best scores per experiment and dataset.



Figure 4: Heatmap comparing MLQA F1  $LA$  and  $noLA$  scores across source and target languages for Experiment 1. Positive scores mean  $LA$  is superior.

show negligible or detrimental effects (with  $LA_{en}$ : -0.04 for Swedish, Catalan and Danish compared to  $noLA_{en}$ ). All non-English *seen* target languages are *rarely seen*, thus, possess minimal pre-training data compared to English. We hypothesize that LAs might interfere with language-specific representations, existing in the base LLM for the respective target language, resulting in reduced downstream task performance. For unseen languages, this interference is reduced, which facilitates learning more meaningful language-specific representations.

As for the impact of the source language, we find that *en* and *de* generally yield similar results while *es* falls behind. German can be leveraged effectively as a source language despite constituting only 0.17% of Llama 2’s pre-training data. Notably, as Table 1 shows, performance drops drastically for English as target language when transferring from German or Spanish under both  $noLA$  and  $LA$ . We conjecture that training TAs reinforces a source language bias, and that using non-English source languages introduces noise, as all training data is *translated* from English, leading to lower-quality data and hindering generalization to English.

**SIB-200.** Figure 17 illustrates that the benefit of LAs vanishes for SIB-200. This aligns with previous work (Kew et al., 2024; Razumovskaia et al., 2024). A topic classification task such as SIB-200

probably requires less language-specific knowledge and rather relies on high-level, language-agnostic semantic features that are already well-encoded in the base LLM. Adding LAs may disrupt existing task-relevant features.

We notice other differences to MLQA: LAs are less harmful for *de* (−0.04) and *es* (−0.02) than for *en* (−0.09)<sup>8</sup>. We assume that while source languages with a weaker pre-training bias are beneficial, they cannot fully mitigate the disruptions induced by the LAs. As for English as target language, in both  $LA$  and  $noLA$ , *de* and *es* are competitive with *en*, suggesting effective cross-lingual generalization to English on SIB-200.

### 4.3 Experiment 2: ICL

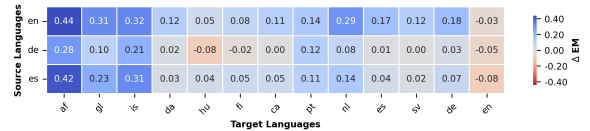


Figure 5: Heatmap comparing MLQA F1  $LA$  and  $noLA$  scores across source and target languages for Experiment 2 with ICL. Positive scores mean  $LA$  is superior.

**MLQA.** Figure 2 illustrates that performance generally drops only moderately when using ICL instead of TAs. This suggests robust ICL capabilities of the base LLM for even more complex tasks. Similar to Experiment 1, with ICL, LAs are most effective for the unseen languages Afrikaans, Galician and Icelandic across source languages (see Figure 5). *en* and *de* yield absolute performance gains of +0.08 and +0.07 on average over the  $noLA$  setup, respectively.

Regarding seen languages, Figure 5 shows mostly minimal performance differences between

<sup>8</sup>All numbers are averaged over five random seeds.

*LA* and *noLA* across source languages. Considering that ICL disentangles the *LA* effect from the task adaptation stage as the latter does not involve any parameter updates, results with ICL indicate that *LAs* may rather add *redundant* than *interfering* representations, as observed for Llama-2/TA.

**SIB-200.** Unlike Experiment 1 with TAs, Figure 18 shows that *LA* consistently outperforms *noLA* with ICL. However, Figure 3 illustrates that a single TA, a computationally cheaper setup, suffices to surpass *LA* with ICL across target languages, again making *LAs* an inefficient choice. Similar to MLQA, *LAs* provoke particularly pronounced performance improvements for unseen languages.

In line with Llama-2/TA, in any Llama-2/ICL setup examined, *de* and *es* considerably outperform *en*, suggesting that the heavy English pre-training bias may hinder the transfer of task-relevant knowledge stored in pre-trained representations.

#### 4.4 Experiment 3: Llama 3.1



Figure 6: Heatmap comparing MLQA F1 *LA* and *noLA* scores across source and target languages for Experiment 3 with Llama 3.1. Positive scores mean *LA* is superior.

**MLQA.** Figure 2 shows that Llama-3.1/TA surpasses Llama-2/TA. When comparing the overall best scores across experiments, there is no language where Llama 2 surpasses Llama 3.1. However, performance gains are only marginally across most non-English target languages, highlighting that simply switching to a stronger, more multilingual base LLM does not bridge the performance gap in English-centric LLMs.

Figure 6 shows that the positive effect of *LAs* for unseen languages vanishes with Llama 3.1. Moreover, across source languages, for unseen languages, Llama 3.1 under *noLA* is on par with or outperforms Llama 2 under *LA*. Considering the amplified pre-training data size in Llama 3.1 (15T tokens vs. 2T tokens in Llama 2), we hypothesize that previously *unseen* languages Afrikaans, Galician and Icelandic in Llama 2 effectively turn into *rarely seen* languages in Llama 3.1 and benefit from larger language-specific pre-training corpora.

Thus, *LAs* for these languages may be prone to the same interference as discussed for seen languages in Llama 2. These findings further suggest that adding language-specific representations *during* pre-training may be more effective for XLT than *after* pre-training through *LAs*, as highlighted by Pfeiffer et al. (2022).

Regarding seen languages, *LAs* with Llama 3.1 induce more severe deterioration than with Llama 2. This results in Llama 3.1 building both the performance-wise head (*noLA* setup) and the tail (*LA* setup) across base LLMs. More language-specific pre-training data seems to be generally beneficial for XLT in the *noLA* setup, while stacking *LAs* in the target language and a TA trained in the source language may be more susceptible to interference.

**SIB-200.** As Table 10 shows, performance is similar across source languages and within each source language, only marginal differences exist between *noLA* and *LA*. This is dissimilar to findings for Llama 2 where *de* and *es* outperformed *en* and *LAs* produced performance deterioration across the board.

Table 10 shows that *es* consistently yields the best EM scores across target languages in both XLT setups. *LA* outperforms *noLA* only marginally, with a maximum absolute performance improvement of +0.03 for Galician. Considering the generally high performance on SIB-200 (with *es*: avg. of 0.81 across target languages for both XLT setups), we do not assume that *LAs* add meaningful, language-specific representations, leading to better performance.

## 5 Qualitative Analysis

Based on the results of Experiment 1-3, we conduct a qualitative analysis and analyze model representations of intermediate layers. The goal here is to better understand the representation shifts that *LAs* do or do not induce.

**Method.** We utilize Logit Lens (nostalgebraist, 2020) to visualize next-token distributions. Logit Lens applies the unembedding matrix of the model to project hidden states of intermediate layers into the space of the output vocabulary. Thus, Logit Lens yields next-token distributions across different input positions and layers.

LogitLens: Llama 2 | setup: LA | source: English | target: German

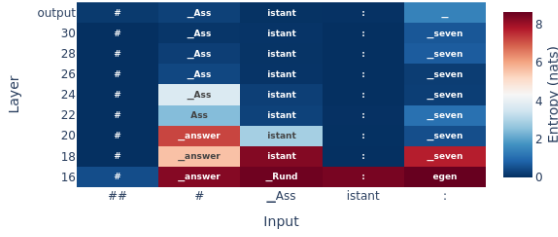


Figure 7: Logit Lens for MLQA test instance with *en* and German as target language. Target: *sieben* (seven). Base LLM: Llama 2. Setup: *LA*.

LogitLens: Llama 2 | setup: LA | source: English | target: Icelandic

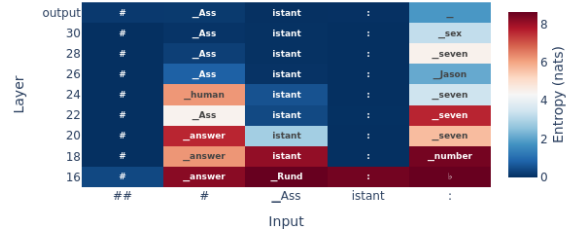


Figure 9: Logit Lens for MLQA test instance with *en* and Icelandic as target language. Target: *sjö* (seven). Base LLM: Llama 2. Setup: *LA*.

LogitLens: Llama 2 | setup: noLA | source: English | target: German

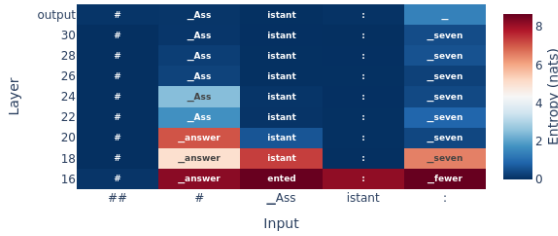


Figure 8: Logit Lens for MLQA test instance with *en* and German as target language. Target: *sieben* (seven). Base LLM: Llama 2. Setup: *noLA*.

LogitLens: Llama 2 | setup: noLA | source: English | target: Icelandic

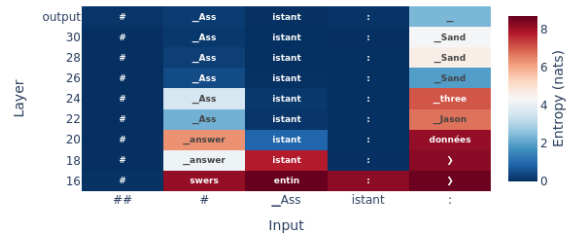


Figure 10: Logit Lens for MLQA test instance with *en* and Icelandic as target language. Target: *sjö* (seven). Base LLM: Llama 2. Setup: *noLA*.

**Setup.** Logit Lens<sup>9</sup> is used to investigate whether LAs introduce shifts in the next-token distributions. Given the observed interferences with Llama-2/TA, we stick to Llama-2/ICL for Logit Lens experiments. Again, we use 5 and 10 source language demonstrations for MLQA and SIB-200, respectively. We aim for test instances with *single-token*, *language-specific* targets, given that Logit Lens visualizes only the first token of the output by default and to assess the promotion of language-specific tokens through LAs, respectively.<sup>10</sup>

We select German and Icelandic as target languages to represent the two extremes of LA impact, with LAs being consistently redundant for German and beneficial for Icelandic. We discuss all examples with *en*. As LAs showed larger effects on MLQA, we focus on MLQA and present Logit Lens visualizations for SIB-200 in Appendix G.2.

**MLQA.** Figures 7 to 10 show the Logit Lens visualizations for German and Icelandic with *en*

under *LA* and *noLA*. The Figures show the final five input positions from layer 16 onward.<sup>11</sup> The token in the upper-right corner corresponds to the token being predicted, i.e., the target.<sup>12</sup>

Regarding German, LAs had no impact on MLQA. This is reflected in the Logit Lens analysis by negligible differences between *LA* (Figure 7) and *noLA* (Figure 8) across layers and positions, suggesting that next-token distributions are mainly preserved. Moreover, in both XLT setups, intermediate layers at the final position are dominated by English tokens. This aligns with findings by Wendler et al. (2024) and Zhang et al. (2024a), who made the identical observation for Chinese.

Regarding Icelandic, Figures 9 and 10 show that differences in the next-token distributions between *LA* and *noLA* are most salient at the final position. While similar to German, *LA* ranks the English

<sup>11</sup>Earlier layers mostly contain tokens without meaningful signal.

<sup>12</sup>Note that the underscore represents a whitespace. Models often predicted the digit 7 with a leading whitespace instead of the written-out variant.

variant of the correct token highest in intermediate layers, *noLA* fails to extract the correct token.<sup>13</sup> Thus, LAs may assist in steering the base LLM towards the correct token by upweighing contextually related English tokens.

If these observations can be verified to be a trend among more German and Icelandic MLQA test instances, Logit Lens provides valuable insights into why performance for German is unchanged and improved for Icelandic, and further strengthens the hypothesis that LAs provoke only marginal transformations to the base LLM.

**SIB-200.** As Figures 20 to 23 illustrate, the correct label *politics* emerges in intermediate layers and is predicted confidently in both XLT setups across target languages. This suggests that for SIB-200, ten task demonstrations suffice to elicit robust ICL abilities and establish a solid understanding useful for XLT. Furthermore, negligible differences between *LA* and *noLA* next-token distributions highlight that LAs are at best redundant for SIB-200 across target languages.

## 6 Main Take-Aways

We take the findings from the three experiments the qualitative analysis and summarize them as follows.

**LAs are beneficial for *unseen* languages on tasks requiring more language-specific knowledge.** Unseen languages (Afrikaans, Galician and Icelandic in Llama 2) evaluated on MLQA are the only languages that consistently benefit from the usage of LAs. This is corroborated by Experiment 2 with ICL which disentangles the effect of the LA from the task adaptation stage more explicitly.

**LAs are at best redundant for *rarely seen* languages and tasks requiring less language-specific knowledge.** Across all experiments, *noLA* is competitive with or surpasses *LA* for most task-language-combinations. Experiment 3 with Llama 3.1 as base LLM substantiates this finding, as the positive effect of LAs vanishes entirely; attributed to previously unseen languages in Llama 2 turning into rarely seen languages in Llama 3.1. Hence, in most cases, adding language-specific representations *during* pre-training appears performance-wise more effective and computationally more efficient than *after* pre-training via LAs.

<sup>13</sup>Tokens like *\_Sand* and *\_Jason* occur in the instance’s passage and denote names.

**The impact of the typological relatedness between source and target language is *minimal*.** Rather, the source language bias and task-specific requirements were found to be critical for the source language choice. English as source language consistently yielded the best performance across target languages on the QA task, whereas German and Spanish were superior on the NLU task.

**LAs and XLT to underrepresented target languages are constrained by the inherent English bias of the base LLM.** While the competitive results of the XLT setup without LAs across our experiments suggest that English-centric representations are able to generalize across non-English target languages, this generalization is severely limited, as evidenced by the performance gap between English and non-English languages on the QA task. Preliminary analyses using the Logit Lens, based on a limited number of test instances and languages, further suggest that LAs, as implemented in our work, may not be able to induce profound language-specific transformations and mitigate the strong English bias of the base LLM.

## 7 Conclusion

We comprehensively evaluated the efficacy of LAs for XLT in English-centric LLMs on 13 languages and 2 downstream tasks. We investigated several XLT configurations with varying task adaptation methods and base LLMs and found that the effect of LAs is largely inconsistent across target languages and tasks. Omitting LAs entirely and relying on a single TA often yielded superior results. A positive effect of LAs was mostly observed for unseen languages, while even minimal language-specific pre-training data tended to diminish this effect. We conclude that LAs do not consistently help enhance XLT and cannot fully mitigate the evident performance gap between English and non-English languages in English-centric LLMs.

From a broader perspective, our findings establish a solid foundation for future research to explore, in greater depth, the capabilities of LAs and the transformations they provoke within English-centric LLMs.

## Limitations

**Languages.** As we rely on automatic evaluation, data sparsity hinders the inclusion of truly low-resource languages. We focus on mainly mid-



to high-resource languages, underrepresented in English-centric LLMs. Future work is encouraged to include low-resource languages that are likely to have yet less pre-training data in the respective base LLMs to test the hypothesis that LAs can help enhance XLT to unseen languages in greater detail. Besides, all languages examined use the Latin script. It is, therefore, straightforward to include non-Latin script languages in future experiments.

**Tasks & Data.** This study is restricted to one QA and one NLU task. Naturally, this hinders us from asserting strong conclusions regarding XLT in English-centric LLMs and implications for real-world applications that rely on robust multilingual generation capabilities. Moreover, we note that automatic translations and metric flaws may confound the obtained results for non-English languages on MLQA.

**Base LLMs.** Our XLT evaluations are limited to two Llama variants. To account for potential Llama-specific biases and to strengthen our hypothesis that LAs primarily benefit unseen languages, a more diverse set of base LLMs is essential.

**Language Adapters.** We highlight four LA-related limitations: First, we did not conduct comprehensive LA hyperparameter tuning. While we briefly explored the number of training samples by doubling the default, we did not examine the reduction factor or potential domain mismatches in the LA data -factors that may be especially important for performance. Second, LAs, as utilized in this study, do not operate on vocabulary level. Thus, the English-centric vocabulary of the base LLM remains unchanged throughout LA training, potentially adversely affecting excessively tokenized languages. Third, we restricted the evaluation of the effect of LAs to an extrinsic evaluation based on downstream task performance. Finally, LAs, as trained in this work, follow a data-driven, post-hoc approach, meaning that we rely on the ability of the base LLM to learn language-specific representations after pre-training by simply feeding in unlabeled, language-specific data while freezing all parameters of the base LLM. Hence, we do not take into account language-specific neurons or regions of the base LLM that may impact performance, as shown by Tang et al., 2024; Zhang et al., 2024b, *inter alia*.

## References

- David Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba Alabi, Yanke Mao, Haonan Gao, and En-Shiun Lee. 2024. [SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian’s, Malta. Association for Computational Linguistics.
- Divyanshu Aggarwal, Ashutosh Sathe, Ishaan Watts, and Sunayana Sitaram. 2024. [MAPLE: Multilingual evaluation of parameter efficient finetuning of large language models](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 14824–14867, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma, Ishaan Watts, Ashutosh Sathe, Millicent Ochieng, Rishav Hada, Prachi Jain, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. 2024. [MEGAVERSE: Benchmarking large language models across languages, modalities, models and tasks](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2598–2637, Mexico City, Mexico. Association for Computational Linguistics.
- Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leueling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, Charvi Jain, Alexander Weber, Lena Jurkschat, Hammam Abdelwahab, Chelsea John, Pedro Ortiz Suarez, Malte Ostendorff, Samuel Weinbach, Rafet Sifa, and 2 others. 2024. [Tokenizer choice for LLM training: Negligible or crucial?](#) In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3907–3924, Mexico City, Mexico. Association for Computational Linguistics.
- Akari Asai, Sneha Kudugunta, Xinyan Yu, Terra Blevins, Hila Gonen, Machel Reid, Yulia Tsvetkov, Sebastian Ruder, and Hannaneh Hajishirzi. 2024. [BUFFET: Benchmarking large language models for few-shot cross-lingual transfer](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1771–1800, Mexico City, Mexico. Association for Computational Linguistics.
- Pinzhen Chen, Shaoxiong Ji, Nikolay Bogoychev, Andrey Kutuzov, Barry Haddow, and Kenneth Heafield. 2024. [Monolingual or multilingual instruction tuning: Which makes a better alpaca](#). In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1347–1356, St. Julian’s, Malta. Association for Computational Linguistics.
- Nadezhda Chirkova and Vassilina Nikoulina. 2024. [Zero-shot cross-lingual transfer in instruction tuning](#)

695	of large language models. In <i>Proceedings of the 17th International Natural Language Generation Conference</i> , pages 695–708, Tokyo, Japan. Association for Computational Linguistics.	
696		
697		
698		
699	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. <i>Unsupervised cross-lingual representation learning at scale</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	
700		
701		
702		
703		
704		
705		
706		
707		
708	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 516 others. 2024. <i>The llama 3 herd of models</i> . <i>Preprint</i> , arXiv:2407.21783.	
709		
710		
711		
712		
713		
714		
715		
716	Fahim Faisal and Antonios Anastasopoulos. 2022. <i>Phylogeny-inspired adaptation of multilingual models to new languages</i> . In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 434–452, Online only. Association for Computational Linguistics.	
717		
718		
719		
720		
721		
722		
723		
724		
725	Daniil Gurgurov, Mareike Hartmann, and Simon Osermann. 2024. <i>Adapting multilingual LLMs to low-resource languages with knowledge graphs via adapters</i> . In <i>Proceedings of the 1st Workshop on Knowledge Graphs and Large Language Models (KaLLM 2024)</i> , pages 63–74, Bangkok, Thailand. Association for Computational Linguistics.	
726		
727		
728		
729		
730		
731		
732	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. <i>Parameter-efficient transfer learning for NLP</i> . In <i>Proceedings of the 36th International Conference on Machine Learning</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 2790–2799. PMLR.	
733		
734		
735		
736		
737		
738		
739		
740	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. <i>Lora: Low-rank adaptation of large language models</i> . <i>Preprint</i> , arXiv:2106.09685.	
741		
742		
743		
744	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. <i>Mistral 7b</i> . <i>Preprint</i> , arXiv:2310.06825.	
745		
746		
747		
748		
749		
750		
751		
	Tannon Kew, Florian Schottmann, and Rico Sennrich. 2024. <i>Turning English-centric LLMs into polyglots: How much multilinguality is needed?</i> In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 13097–13124, Miami, Florida, USA. Association for Computational Linguistics.	752
		753
		754
		755
		756
		757
	Jenny Kunz. 2025. <i>Train more parameters but mind their placement: Insights into language adaptation with PEFT</i> . In <i>Proceedings of the Joint 25th Nordic Conference on Computational Linguistics and 11th Baltic Conference on Human Language Technologies (NoDaLiDa/Baltic-HLT 2025)</i> , pages 323–330, Tallinn, Estonia. University of Tartu Library.	758
		759
		760
		761
		762
		763
		764
	Jenny Kunz and Oskar Holmström. 2024. <i>The impact of language adapters in cross-lingual transfer for NLU</i> . In <i>Proceedings of the 1st Workshop on Modular and Open Multilingual NLP (MOOMIN 2024)</i> , pages 24–43, St Julians, Malta. Association for Computational Linguistics.	765
		766
		767
		768
		769
		770
	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. <i>The power of scale for parameter-efficient prompt tuning</i> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	771
		772
		773
		774
		775
		776
		777
	Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. <i>MLQA: Evaluating cross-lingual extractive question answering</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7315–7330, Online. Association for Computational Linguistics.	778
		779
		780
		781
		782
		783
		784
	Yinheng Li. 2023. <i>A practical survey on zero-shot prompt design for in-context learning</i> . In <i>Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing</i> , pages 641–647, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.	785
		786
		787
		788
		789
		790
	Peiqin Lin, Shaoxiong Ji, Jörg Tiedemann, André F. T. Martins, and Hinrich Schütze. 2024. <i>Mala-500: Massive language adaptation of large language models</i> . <i>Preprint</i> , arXiv:2401.13303.	791
		792
		793
		794
	nostalgebraist. 2020. <i>interpreting gpt: the logit lens</i> . Accessed: 2024-12-17.	795
		796
	Romina Oji and Jenny Kunz. 2025. <i>How to tune a multilingual encoder model for Germanic languages: A study of PEFT, full fine-tuning, and language adapters</i> . In <i>Proceedings of the Joint 25th Nordic Conference on Computational Linguistics and 11th Baltic Conference on Human Language Technologies (NoDaLiDa/Baltic-HLT 2025)</i> , pages 433–439, Tallinn, Estonia. University of Tartu Library.	797
		798
		799
		800
		801
		802
		803
		804
	Marinela Parović, Goran Glavaš, Ivan Vulić, and Anna Korhonen. 2022. <i>BAD-X: Bilingual adapters improve zero-shot cross-lingual transfer</i> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational</i>	805
		806
		807
		808
		809

810	<i>Linguistics: Human Language Technologies</i> , pages	Vipul Rathore, Rajdeep Dhingra, Parag Singla, and	868
811	1791–1799, Seattle, United States. Association for	Mausam. 2023. <a href="#">ZGUL: Zero-shot generalization to</a>	869
812	Computational Linguistics.	<a href="#">unseen languages using multi-source ensembling of</a>	870
813	Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James	<a href="#">language adapters</a> . In <i>Proceedings of the 2023 Con-</i>	871
814	Cross, Sebastian Riedel, and Mikel Artetxe. 2022.	<i>ference on Empirical Methods in Natural Language</i>	872
815	<a href="#">Lifting the curse of multilinguality by pre-training</a>	<i>Processing</i> , pages 6969–6987, Singapore. Associa-	873
816	<a href="#">modular transformers</a> . In <i>Proceedings of the 2022</i>	tion for Computational Linguistics.	874
817	<i>Conference of the North American Chapter of the</i>		
818	<i>Association for Computational Linguistics: Human</i>	Evgeniia Razumovskaia, Ivan Vulić, and Anna Korho-	875
819	<i>Language Technologies</i> , pages 3479–3495, Seattle,	nen. 2024. <a href="#">Analyzing and adapting large language</a>	876
820	United States. Association for Computational Lin-	<a href="#">models for few-shot multilingual nlu: Are we there</a>	877
821	guistics.	<a href="#">yet?</a> <i>Preprint</i> , arXiv:2403.01929.	878
822	Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya	Lütfi Kerem Senel, Benedikt Ebing, Konul Baghirova,	879
823	Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun	Hinrich Schuetze, and Goran Glavaš. 2024. <a href="#">Kardeş-</a>	880
824	Cho, and Iryna Gurevych. 2020a. <a href="#">AdapterHub: A</a>	<a href="#">NLU: Transfer to low-resource languages with the</a>	881
825	<a href="#">framework for adapting transformers</a> . In <i>Proceedings</i>	<a href="#">help of a high-resource cousin – a benchmark and</a>	882
826	<i>of the 2020 Conference on Empirical Methods in Nat-</i>	<a href="#">evaluation for Turkic languages</a> . In <i>Proceedings of</i>	883
827	<i>ural Language Processing: System Demonstrations</i> ,	<i>the 18th Conference of the European Chapter of the</i>	884
828	pages 46–54, Online. Association for Computational	<i>Association for Computational Linguistics (Volume 1:</i>	885
829	Linguistics.	<i>Long Papers)</i> , pages 1672–1688, St. Julian’s, Malta.	886
830	Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Se-	Association for Computational Linguistics.	887
831	bastian Ruder. 2020b. <a href="#">MAD-X: An Adapter-Based</a>		
832	<a href="#">Framework for Multi-Task Cross-Lingual Transfer</a> .	Shivalika Singh, Freddie Vargus, Daniel D’souza,	888
833	In <i>Proceedings of the 2020 Conference on Empirical</i>	Börje Karlsson, Abinaya Mahendiran, Wei-Yin Ko,	889
834	<i>Methods in Natural Language Processing (EMNLP)</i> ,	Herumb Shandilya, Jay Patel, Deividas Mataciun-	890
835	pages 7654–7673, Online. Association for Computa-	nas, Laura O’Mahony, Mike Zhang, Ramith Het-	891
836	tional Linguistics.	tiarachchi, Joseph Wilson, Marina Machado, Luisa	892
837	Fred Philippy, Siwen Guo, and Shohreh Haddadan.	Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem	893
838	2023. <a href="#">Towards a common understanding of con-</a>	Ergun, Ifeoma Okoh, and 14 others. 2024. <a href="#">Aya</a>	894
839	<a href="#">tributing factors for cross-lingual transfer in multi-</a>	<a href="#">dataset: An open-access collection for multilingual</a>	895
840	<a href="#">lingual language models: A review</a> . In <i>Proceedings</i>	<a href="#">instruction tuning</a> . In <i>Proceedings of the 62nd An-</i>	896
841	<i>of the 61st Annual Meeting of the Association for</i>	<i>nual Meeting of the Association for Computational</i>	897
842	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	<i>Linguistics (Volume 1: Long Papers)</i> , pages 11521–	898
843	pages 5877–5891, Toronto, Canada. Association for	11567, Bangkok, Thailand. Association for Compu-	899
844	Computational Linguistics.	tational Linguistics.	900
845	Fred Philippy, Siwen Guo, Shohreh Haddadan, Cedric	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dong-	901
846	Lothritz, Jacques Klein, and Tegawendé F. Bissyandé.	dong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei,	902
847	2024. <a href="#">Soft prompt tuning for cross-lingual trans-</a>	and Ji-Rong Wen. 2024. <a href="#">Language-specific neurons:</a>	903
848	<a href="#">fer: When less is more</a> . In <i>Proceedings of the 1st</i>	<a href="#">The key to multilingual capabilities in large language</a>	904
849	<i>Workshop on Modular and Open Multilingual NLP</i>	<a href="#">models</a> . In <i>Proceedings of the 62nd Annual Meeting</i>	905
850	<i>(MOOMIN 2024)</i> , pages 7–15, St Julians, Malta. As-	<i>of the Association for Computational Linguistics (Vol-</i>	906
851	sociation for Computational Linguistics.	<i>ume 1: Long Papers)</i> , pages 5701–5715, Bangkok,	907
852	Clifton Poth, Hannah Sterz, Indraneil Paul, Sukannya	Thailand. Association for Computational Linguistics.	908
853	Purkayastha, Leon Engländer, Timo Imhof, Ivan		
854	Vulić, Sebastian Ruder, Iryna Gurevych, and Jonas	NLLB Team, Marta R. Costa-jussà, James Cross, Onur	909
855	Pfeiffer. 2023. <a href="#">Adapters: A unified library for</a>	Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hef-	910
856	<a href="#">parameter-efficient and modular transfer learning</a> . In	ernan, Elahe Kalbassi, Janice Lam, Daniel Licht,	911
857	<i>Proceedings of the 2023 Conference on Empirical</i>	Jean Maillard, Anna Sun, Skyler Wang, Guillaume	912
858	<i>Methods in Natural Language Processing: System</i>	Wenzek, Al Youngblood, Bapi Akula, Loic Barrault,	913
859	<i>Demonstrations</i> , pages 149–160, Singapore. Associa-	Gabriel Mejia Gonzalez, Prangthip Hansanti, and	914
860	tion for Computational Linguistics.	20 others. 2022. <a href="#">No language left behind: Scal-</a>	915
861	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.	<a href="#">ing human-centered machine translation</a> . <i>Preprint</i> ,	916
862	<a href="#">Know what you don’t know: Unanswerable ques-</a>	arXiv:2207.04672.	917
863	<a href="#">tions for SQuAD</a> . In <i>Proceedings of the 56th Annual</i>		
864	<i>Meeting of the Association for Computational Lin-</i>	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	918
865	<i>guistics (Volume 2: Short Papers)</i> , pages 784–789,	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	919
866	Melbourne, Australia. Association for Computational	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	920
867	Linguistics.	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	921
		Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	922
		Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 oth-	923
		ers. 2023. <a href="#">Llama 2: Open foundation and fine-tuned</a>	924
		<a href="#">chat models</a> . <i>Preprint</i> , arXiv:2307.09288.	925



- Ivan Vykopal, Simon Ostermann, and Marian Simko. 2025. [Soft language prompts for language transfer](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 10294–10313, Albuquerque, New Mexico. Association for Computational Linguistics.
- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. [Do llamas work in English? on the latent language of multilingual transformers](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15366–15394, Bangkok, Thailand. Association for Computational Linguistics.
- Jiacheng Ye, Xijia Tao, and Lingpeng Kong. 2023. [Language versatilists vs. specialists: An empirical revisiting on multilingual transfer ability](#). *ArXiv*, abs/2306.06688.
- Zheng Xin Yong, Hailey Schoelkopf, Niklas Muenighoff, Alham Fikri Aji, David Ifeoluwa Adelani, Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Kasai, Ahmed Baruwa, Genta Winata, Stella Biderman, Edward Raff, Dragomir Radev, and Vasilina Nikoulina. 2023. [BLOOM+1: Adding language support to BLOOM for zero-shot prompting](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11682–11703, Toronto, Canada. Association for Computational Linguistics.
- Shimao Zhang, Changjiang Gao, Wenhao Zhu, Jiajun Chen, Xin Huang, Xue Han, Junlan Feng, Chao Deng, and Shujian Huang. 2024a. [Getting more from less: Large language models are good spontaneous multilingual learners](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8037–8051, Miami, Florida, USA. Association for Computational Linguistics.
- Zhihao Zhang, Jun Zhao, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024b. [Unveiling linguistic regions in large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6228–6247, Bangkok, Thailand. Association for Computational Linguistics.



A Training Details

Hyperparameter	Value
<i>LAs</i>	
Reduction factor	16
Trainable parameters	67.1M
Batch size	4
Training steps	50k
Context length	1024
<i>MLQA TAs</i>	
Reduction factor	16
Trainable parameters	67.1M
Dropout	0.0
Batch size	4
Training epochs	3
<i>SIB-200 TAs</i>	
Reduction factor	32
Trainable parameters	33.6M
Dropout	0.1
Batch size	4
Training epochs	20

Table 2: Details for training LAs and TAs. These values apply to all languages. I.e., LAs are trained on 200k samples per language à 1024 tokens. Due to the same hidden dimension and the same number of hidden layers, the number of trainable parameters applies to both Llama 2 and Llama 3.1. Unspecified hyperparameters were set to the default values as provided in the adapters and transformers library.

B Llama 2 Language Distribution

Language	Data (in %)
en	90.00
de	0.17
sv	0.15
es	0.13
nl	0.12
pt	0.09
ca	0.04
fi	0.03
hu	0.03
da	0.02
is	0.00
gl	0.00
af	0.00

Table 3: Amounts of pre-training data in Llama 2 for languages relevant to this work. No detailed language distribution is available for Llama 3.1.

C Fertility

Language	Fertility
en	1.45
de	2.04
sv	2.21
es	1.77
nl	2.00
pt	1.92
ca	1.96
fi	3.75
hu	3.00
da	2.22
is	3.03
gl	1.97
af	2.11

Table 4: Fertility (token/word ratio) as measured on the dev split of Flores-200 (Team et al., 2022) using the English-centric tokenizer of Llama 2.

D Task Templates  
D.1 Task Adapters

```
MLQA

### Human: Refer to the passage below and
then answer the question afterwards in the
same language as the passage:

Passage: {passage}

Question: {question}

### Assistant: {answer}
```

Figure 11: Prompt template used for MLQA during TA training and at inference for setups using TAs.

```
SIB-200

Classify the following sentence into one of
the following topics:
1. science/technology
2. travel
3. politics
4. sports
5. health
6. entertainment
7. geography

Sentence: {sentence}

Topic: {topic}
```

Figure 12: Prompt template used for SIB-200 during TA training and at inference for setups using TAs.

D.2 In-context Learning

```
MLQA

### Instruction: The task is to solve
reading comprehension problems. You will
be provided questions on a set of passages
and you will need to provide the answer
as it appears in the passage. The answer
should be in the same language as the
question and the passage. Provide nothing
else beyond the answer.

- n source language demonstrations -
### Human:
Passage: {passage}
Question: {question}

### Assistant: {answer}

### Human:
Passage: The aircraft involved in
the hijacking was a Boeing 757-222,
registration N591UA, delivered to the
airline in 1996. The airplane had a
capacity of 182 passengers; the September
11 flight carried 37 passengers and
seven crew, a load factor of 20 percent,
considerably below the 52 percent average
Tuesday load factor for Flight 93. The
seven crew members were Captain Jason Dahl,
First Officer LeRoy Homer Jr., and flight
attendants Lorraine Bay, Sandra Bradshaw,
Wanda Green, CeeCee Lyles, and Deborah
Welsh.
Question: How many crew members were there?

### Assistant: seven
```

Figure 13: ICL prompt template for MLQA. The string ‘- n source language demonstrations - -’ is not part of the prompt. This example is also the English test instance chosen for Logit Lens experiments on MLQA. Target is not provided. We set  $n = 5$ .

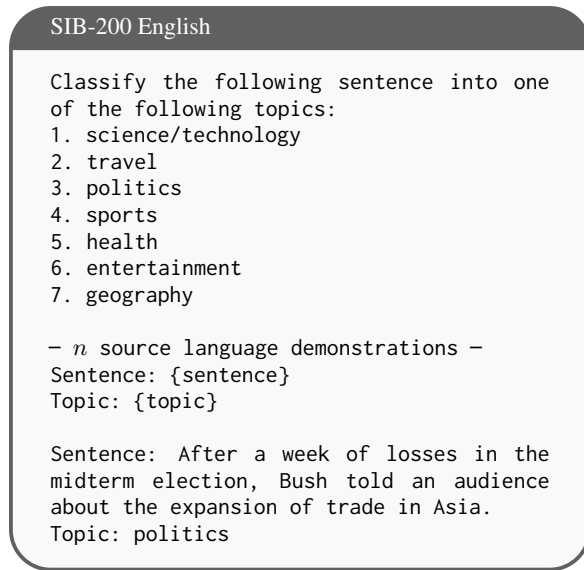


Figure 14: ICL prompt template for SIB-200. The string ‘-  $n$  source language demonstrations -’ is not part of the prompt. This example is also the English test instance chosen for Logit Lens experiments on SIB-200. Target is not provided. We set  $n = 10$ .

E Training & Evaluation Setups

E.1 LA Setup

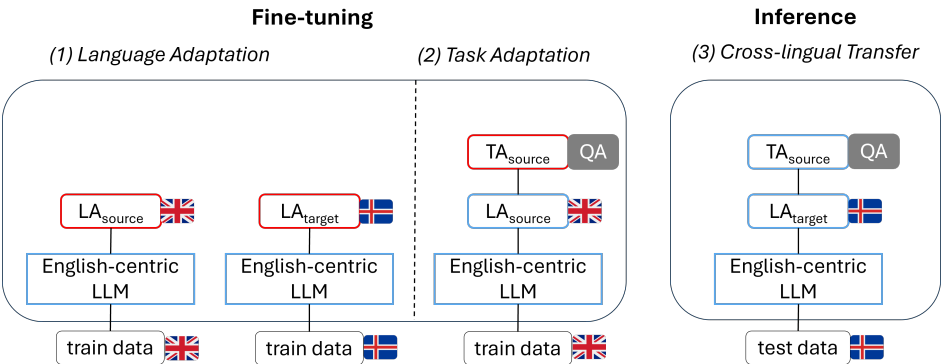


Figure 15: *LA* setup (blue and red edges indicate frozen and trainable parameters, respectively): (1) LAs are trained for each language of interest (here: English and Icelandic) on a frozen English-centric LLM (e.g., Llama 2 7B). (2) A TA (in this case, for a QA task) is trained in the source language (here: English) by stacking it on top of the frozen LA in the respective source language. (3) At inference, the source LA is replaced by the target LA (here: Icelandic) while retaining the TA in the source language. This setup is then evaluated zero-shot in the target language. Own illustration.

E.2 noLA Setup

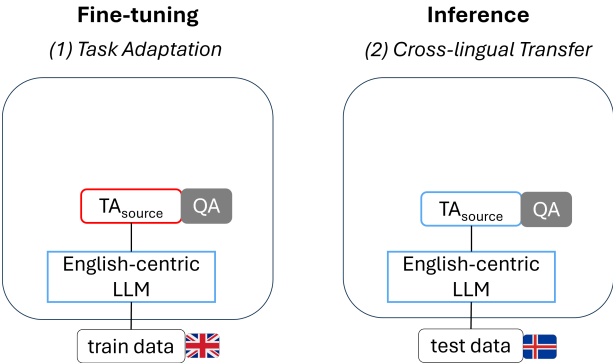


Figure 16: *noLA* setup (blue and red edges indicate frozen and trainable parameters, respectively): (1) A TA (in this case, for a QA task) is trained in the source language (here: English) on top of the frozen English-centric LLM. (2) At inference, the TA in the source language is retained and evaluated zero-shot in the target language (here: Icelandic). Own illustration.



## F Scores

### F.1 Experiment 1: Llama-2/TA

Setup	af	gl	is	da	fi	hu	ca	pt	nl	es	sv	de	en	avg.
$LA_{en}$	<b>0.51</b> ( $\pm 0.02$ )	<b>0.56</b> ( $\pm 0.01$ )	<b>0.32</b> ( $\pm 0.02$ )	0.49 ( $\pm 0.01$ )	0.33 ( $\pm 0.01$ )	0.39 ( $\pm 0.02$ )	0.53 ( $\pm 0.03$ )	0.53 ( $\pm 0.02$ )	0.53 ( $\pm 0.01$ )	0.47 ( $\pm 0.01$ )	0.46 ( $\pm 0.02$ )	0.51 ( $\pm 0.00$ )	<b>0.78</b> ( $\pm 0.00$ )	0.47
$LA_{de}$	0.50 ( $\pm 0.01$ )	0.54 ( $\pm 0.01$ )	<b>0.32</b> ( $\pm 0.01$ )	0.47 ( $\pm 0.01$ )	<b>0.37</b> ( $\pm 0.01$ )	0.42 ( $\pm 0.01$ )	0.54 ( $\pm 0.01$ )	0.52 ( $\pm 0.00$ )	0.52 ( $\pm 0.01$ )	0.47 ( $\pm 0.00$ )	0.47 ( $\pm 0.01$ )	<b>0.54</b> ( $\pm 0.00$ )	0.44 ( $\pm 0.09$ )	0.47
$LA_{es}$	0.45 ( $\pm 0.02$ )	0.51 ( $\pm 0.02$ )	0.31 ( $\pm 0.02$ )	0.45 ( $\pm 0.02$ )	0.34 ( $\pm 0.01$ )	0.39 ( $\pm 0.01$ )	0.52 ( $\pm 0.01$ )	0.51 ( $\pm 0.01$ )	0.48 ( $\pm 0.01$ )	<b>0.53</b> ( $\pm 0.01$ )	0.44 ( $\pm 0.01$ )	0.46 ( $\pm 0.01$ )	0.43 ( $\pm 0.05$ )	0.44
$noLA_{en}$	0.49 ( $\pm 0.01$ )	0.52 ( $\pm 0.01$ )	0.26 ( $\pm 0.01$ )	<b>0.53</b> ( $\pm 0.01$ )	0.34 ( $\pm 0.01$ )	0.39 ( $\pm 0.01$ )	<b>0.57</b> ( $\pm 0.01$ )	<b>0.55</b> ( $\pm 0.01$ )	<b>0.55</b> ( $\pm 0.01$ )	0.48 ( $\pm 0.01$ )	<b>0.50</b> ( $\pm 0.01$ )	0.51 ( $\pm 0.00$ )	<b>0.78</b> ( $\pm 0.00$ )	0.47
$noLA_{de}$	0.40 ( $\pm 0.01$ )	0.47 ( $\pm 0.01$ )	0.23 ( $\pm 0.00$ )	0.50 ( $\pm 0.01$ )	<b>0.37</b> ( $\pm 0.00$ )	<b>0.43</b> ( $\pm 0.01$ )	0.55 ( $\pm 0.01$ )	0.54 ( $\pm 0.01$ )	0.47 ( $\pm 0.02$ )	0.47 ( $\pm 0.01$ )	0.46 ( $\pm 0.00$ )	<b>0.54</b> ( $\pm 0.00$ )	0.38 ( $\pm 0.01$ )	0.44
$noLA_{es}$	0.38 ( $\pm 0.01$ )	0.38 ( $\pm 0.01$ )	0.20 ( $\pm 0.01$ )	0.44 ( $\pm 0.02$ )	0.31 ( $\pm 0.01$ )	0.34 ( $\pm 0.02$ )	0.46 ( $\pm 0.02$ )	0.45 ( $\pm 0.01$ )	0.45 ( $\pm 0.03$ )	<b>0.53</b> ( $\pm 0.01$ )	0.41 ( $\pm 0.03$ )	0.40 ( $\pm 0.03$ )	0.32 ( $\pm 0.04$ )	0.38

Table 5: MLQA-en F1 scores averaged over five random seeds for Experiment 1. We use the main variables *task adaptation method* = TA (*task-specific FT*), *adapter method* = *bottleneck adapter (BN)*, *base large language model (LLM)* = *Llama 2*, *LA coverage* = *monolingual*. Standard deviation in parentheses. Bold numbers indicate best scores between XLT setups ( $LA$ ,  $noLA$ ), underscored numbers indicate best scores within XLT setup between source languages ( $en$ ,  $de$ ,  $es$ ).

Setup	af	gl	is	da	fi	hu	ca	pt	nl	es	sv	de	en	avg.
$LA_{en}$	0.50 ( $\pm 0.17$ )	0.74 ( $\pm 0.05$ )	0.55 ( $\pm 0.06$ )	0.71 ( $\pm 0.06$ )	0.66 ( $\pm 0.10$ )	0.59 ( $\pm 0.16$ )	0.66 ( $\pm 0.06$ )	0.79 ( $\pm 0.03$ )	0.71 ( $\pm 0.10$ )	0.78 ( $\pm 0.06$ )	0.68 ( $\pm 0.12$ )	0.82 ( $\pm 0.04$ )	0.85 ( $\pm 0.02$ )	0.68
$LA_{de}$	0.77 ( $\pm 0.09$ )	0.81 ( $\pm 0.04$ )	<b>0.70</b> ( $\pm 0.03$ )	0.78 ( $\pm 0.06$ )	<b>0.81</b> ( $\pm 0.04$ )	<b>0.82</b> ( $\pm 0.02$ )	0.77 ( $\pm 0.05$ )	<b>0.84</b> ( $\pm 0.06$ )	0.85 ( $\pm 0.03$ )	0.81 ( $\pm 0.04$ )	0.79 ( $\pm 0.07$ )	<b>0.87</b> ( $\pm 0.01$ )	0.83 ( $\pm 0.05$ )	0.80
$LA_{es}$	0.74 ( $\pm 0.06$ )	0.76 ( $\pm 0.02$ )	0.60 ( $\pm 0.11$ )	<u>0.80</u> ( $\pm 0.03$ )	0.69 ( $\pm 0.09$ )	0.71 ( $\pm 0.07$ )	0.76 ( $\pm 0.09$ )	0.82 ( $\pm 0.02$ )	0.82 ( $\pm 0.03$ )	<u>0.82</u> ( $\pm 0.02$ )	<u>0.81</u> ( $\pm 0.04$ )	0.81 ( $\pm 0.05$ )	0.82 ( $\pm 0.05$ )	0.76
$noLA_{en}$	0.72 ( $\pm 0.03$ )	0.79 ( $\pm 0.03$ )	0.40 ( $\pm 0.07$ )	0.79 ( $\pm 0.02$ )	0.68 ( $\pm 0.06$ )	0.73 ( $\pm 0.03$ )	0.80 ( $\pm 0.03$ )	<b>0.84</b> ( $\pm 0.03$ )	0.80 ( $\pm 0.03$ )	0.82 ( $\pm 0.03$ )	0.78 ( $\pm 0.02$ )	0.81 ( $\pm 0.02$ )	<b>0.86</b> ( $\pm 0.02$ )	0.75
$noLA_{de}$	<b>0.83</b> ( $\pm 0.02$ )	<b>0.83</b> ( $\pm 0.02$ )	0.56 ( $\pm 0.04$ )	<b>0.85</b> ( $\pm 0.01$ )	<b>0.81</b> ( $\pm 0.02$ )	<b>0.82</b> ( $\pm 0.02$ )	<b>0.84</b> ( $\pm 0.02$ )	<b>0.84</b> ( $\pm 0.03$ )	<b>0.86</b> ( $\pm 0.02$ )	<b>0.84</b> ( $\pm 0.02$ )	<b>0.84</b> ( $\pm 0.02$ )	0.85 ( $\pm 0.03$ )	<b>0.86</b> ( $\pm 0.02$ )	0.81
$noLA_{es}$	0.74 ( $\pm 0.05$ )	0.79 ( $\pm 0.02$ )	0.45 ( $\pm 0.05$ )	0.80 ( $\pm 0.03$ )	0.73 ( $\pm 0.06$ )	0.74 ( $\pm 0.04$ )	0.83 ( $\pm 0.03$ )	<b>0.84</b> ( $\pm 0.01$ )	0.81 ( $\pm 0.04$ )	0.83 ( $\pm 0.01$ )	0.81 ( $\pm 0.03$ )	0.81 ( $\pm 0.04$ )	0.85 ( $\pm 0.02$ )	0.77

Table 6: SIB-200 EM scores averaged over five random seeds for Experiment 1. We use the main variables *task adaptation method* = TA (*task-specific FT*), *adapter method* = *BN*, *base LLM* = *Llama 2*, *LA coverage* = *monolingual*. Standard deviation in parentheses. Bold numbers indicate best scores between XLT setups ( $LA$ ,  $noLA$ ), underscored numbers indicate best scores within XLT setup between source languages ( $en$ ,  $de$ ,  $es$ ).

## F.2 Experiment 2: Llama-2/ICL

Setup	af	gl	is	da	hu	fi	ca	pt	nl	es	sv	de	en	avg.
$LA_{en}$	<b>0.45</b>	<b>0.47</b>	<b>0.30</b>	<b>0.45</b>	<b>0.34</b>	<b>0.31</b>	<b>0.48</b>	<b>0.47</b>	<b>0.47</b>	<b>0.41</b>	<b>0.43</b>	<b>0.43</b>	<b>0.65</b>	0.42
$LA_{de}$	0.41	0.42	0.28	0.42	0.32	0.30	0.44	0.42	0.44	0.38	0.41	0.41	0.57	0.40
$LA_{es}$	0.40	0.35	0.25	0.38	0.29	0.28	0.37	0.41	0.39	0.37	0.38	0.38	0.49	0.36
$noLA_{en}$	<u>0.39</u>	<u>0.41</u>	<u>0.20</u>	<u>0.44</u>	<u>0.33</u>	<b>0.31</b>	<b>0.48</b>	<u>0.46</u>	<b>0.47</b>	0.40	<b>0.43</b>	<u>0.42</u>	0.64	0.40
$noLA_{de}$	0.36	0.37	0.18	0.40	0.32	0.29	0.44	0.42	0.42	0.37	0.38	0.40	0.56	0.38
$noLA_{es}$	0.35	0.38	0.18	0.40	0.30	0.28	0.43	0.40	0.39	0.36	0.36	0.35	0.47	0.36

Table 7: MLQA-en F1 scores for Experiment 2 which uses ICL instead of TAs for task adaptation. We use 5 source language task demonstrations. Bold numbers indicate best scores between XLT setups ( $LA$ ,  $noLA$ ), underscored numbers indicate best scores within XLT setup between source languages ( $en$ ,  $de$ ,  $es$ ).

Setup	af	gl	is	da	hu	fi	ca	pt	nl	es	sv	de	en	avg.
$LA_{en}$	0.67	0.61	0.55	0.57	0.53	0.55	0.50	0.64	0.66	0.64	0.59	0.69	<u>0.75</u>	0.60
$LA_{de}$	<b>0.74</b>	0.66	0.54	<b>0.70</b>	0.58	0.65	0.63	<b>0.80</b>	<b>0.74</b>	0.72	<b>0.70</b>	<b>0.81</b>	0.73	0.68
$LA_{es}$	<b>0.74</b>	<b>0.78</b>	<b>0.56</b>	0.62	<u>0.62</u>	<u>0.66</u>	<b>0.65</b>	0.75	0.73	<b>0.79</b>	0.66	0.68	0.65	0.67
$noLA_{en}$	0.23	0.30	0.23	0.45	0.48	0.47	0.39	0.50	0.37	0.47	0.47	0.51	<b>0.78</b>	0.41
$noLA_{de}$	<u>0.46</u>	<u>0.56</u>	<u>0.33</u>	0.68	<b>0.66</b>	<b>0.67</b>	<u>0.63</u>	<u>0.68</u>	<u>0.66</u>	0.71	<b>0.70</b>	<u>0.78</u>	<b>0.78</b>	0.63
$noLA_{es}$	0.32	0.55	0.25	0.59	0.58	0.61	0.60	0.64	0.59	<u>0.75</u>	0.64	0.61	0.73	0.56

Table 8: SIB-200 EM scores for Experiment 2 which uses ICL instead of TAs for task adaptation. We use 10 source language task demonstrations. Bold numbers indicate best scores between XLT setups ( $LA$ ,  $noLA$ ), underscored numbers indicate best scores within XLT setup between source languages ( $en$ ,  $de$ ,  $es$ ).

### F.3 Experiment 3: Llama-3.1/TA

Setup	af	gl	is	da	fi	hu	ca	pt	nl	es	sv	de	en	avg.
$LA_{en}$	<u>0.50</u> ( $\pm 0.01$ )	<u>0.56</u> ( $\pm 0.04$ )	0.34 ( $\pm 0.02$ )	<u>0.48</u> ( $\pm 0.05$ )	0.25 ( $\pm 0.05$ )	0.33 ( $\pm 0.05$ )	0.54 ( $\pm 0.05$ )	<u>0.54</u> ( $\pm 0.03$ )	0.54 ( $\pm 0.03$ )	0.49 ( $\pm 0.02$ )	0.38 ( $\pm 0.07$ )	0.51 ( $\pm 0.01$ )	<b>0.80</b> ( $\pm 0.00$ )	0.46
$LA_{de}$	0.47 ( $\pm 0.03$ )	0.53 ( $\pm 0.04$ )	<u>0.35</u> ( $\pm 0.02$ )	0.47 ( $\pm 0.04$ )	0.34 ( $\pm 0.02$ )	0.46 ( $\pm 0.01$ )	0.51 ( $\pm 0.06$ )	0.49 ( $\pm 0.06$ )	<u>0.55</u> ( $\pm 0.01$ )	0.49 ( $\pm 0.01$ )	0.44 ( $\pm 0.05$ )	<b>0.56</b> ( $\pm 0.00$ )	0.37 ( $\pm 0.11$ )	0.46
$LA_{es}$	0.44 ( $\pm 0.02$ )	0.52 ( $\pm 0.02$ )	0.32 ( $\pm 0.02$ )	0.32 ( $\pm 0.05$ )	0.28 ( $\pm 0.05$ )	0.39 ( $\pm 0.01$ )	<u>0.57</u> ( $\pm 0.01$ )	0.51 ( $\pm 0.01$ )	0.47 ( $\pm 0.03$ )	<b>0.56</b> ( $\pm 0.00$ )	0.30 ( $\pm 0.07$ )	0.47 ( $\pm 0.03$ )	0.43 ( $\pm 0.07$ )	0.42
$noLA_{en}$	0.51 ( $\pm 0.04$ )	0.56 ( $\pm 0.04$ )	0.37 ( $\pm 0.02$ )	0.52 ( $\pm 0.03$ )	0.34 ( $\pm 0.01$ )	0.42 ( $\pm 0.02$ )	0.55 ( $\pm 0.05$ )	0.53 ( $\pm 0.05$ )	0.54 ( $\pm 0.03$ )	0.47 ( $\pm 0.04$ )	0.50 ( $\pm 0.02$ )	0.50 ( $\pm 0.03$ )	<u>0.79</u> ( $\pm 0.00$ )	0.48
$noLA_{de}$	<b>0.54</b> ( $\pm 0.01$ )	<b>0.57</b> ( $\pm 0.01$ )	<b>0.38</b> ( $\pm 0.00$ )	<b>0.54</b> ( $\pm 0.01$ )	<b>0.40</b> ( $\pm 0.01$ )	<b>0.48</b> ( $\pm 0.00$ )	<b>0.59</b> ( $\pm 0.01$ )	<b>0.57</b> ( $\pm 0.01$ )	<b>0.56</b> ( $\pm 0.01$ )	0.50 ( $\pm 0.01$ )	<b>0.53</b> ( $\pm 0.01$ )	<b>0.56</b> ( $\pm 0.01$ )	0.35 ( $\pm 0.01$ )	0.50
$noLA_{es}$	0.48 ( $\pm 0.01$ )	0.51 ( $\pm 0.01$ )	0.34 ( $\pm 0.01$ )	0.49 ( $\pm 0.01$ )	0.36 ( $\pm 0.02$ )	0.42 ( $\pm 0.01$ )	0.51 ( $\pm 0.02$ )	0.51 ( $\pm 0.00$ )	0.50 ( $\pm 0.01$ )	<b>0.56</b> ( $\pm 0.00$ )	0.48 ( $\pm 0.01$ )	0.46 ( $\pm 0.01$ )	0.31 ( $\pm 0.08$ )	0.45

Table 9: MLQA-en F1 scores averaged over five random seeds for Experiment 4 with main variable *base LLM* = *Llama 3.1*. The remaining main variables are *task adaptation method* = *TA (task-specific FT)*, *adapter method* = *BN*, *LA coverage* = *monolingual*. Standard deviation in parentheses. Bold numbers indicate best scores between XLT setups (*LA*, *noLA*), underscored numbers indicate best scores within XLT setup between source languages (*en*, *de*, *es*).

Setup	af	gl	is	da	fi	hu	ca	pt	nl	es	sv	de	en	avg.
$LA_{en}$	0.78 ( $\pm 0.06$ )	0.81 ( $\pm 0.04$ )	0.71 ( $\pm 0.05$ )	0.78 ( $\pm 0.05$ )	0.78 ( $\pm 0.03$ )	0.80 ( $\pm 0.04$ )	0.78 ( $\pm 0.03$ )	0.82 ( $\pm 0.03$ )	0.85 ( $\pm 0.02$ )	0.85 ( $\pm 0.05$ )	0.80 ( $\pm 0.05$ )	<b>0.86</b> ( $\pm 0.02$ )	<b>0.88</b> ( $\pm 0.02$ )	0.80
$LA_{de}$	0.80 ( $\pm 0.03$ )	0.82 ( $\pm 0.03$ )	<b>0.72</b> ( $\pm 0.05$ )	0.81 ( $\pm 0.04$ )	0.78 ( $\pm 0.07$ )	0.80 ( $\pm 0.04$ )	0.80 ( $\pm 0.05$ )	0.81 ( $\pm 0.03$ )	0.82 ( $\pm 0.05$ )	0.81 ( $\pm 0.04$ )	0.79 ( $\pm 0.04$ )	0.84 ( $\pm 0.03$ )	0.80 ( $\pm 0.06$ )	0.80
$LA_{es}$	0.79 ( $\pm 0.04$ )	<b>0.84</b> ( $\pm 0.02$ )	<b>0.72</b> ( $\pm 0.07$ )	0.77 ( $\pm 0.08$ )	<b>0.79</b> ( $\pm 0.02$ )	<b>0.81</b> ( $\pm 0.03$ )	0.80 ( $\pm 0.04$ )	<b>0.85</b> ( $\pm 0.01$ )	<b>0.86</b> ( $\pm 0.02$ )	<b>0.86</b> ( $\pm 0.02$ )	0.82 ( $\pm 0.03$ )	<b>0.86</b> ( $\pm 0.01$ )	0.86 ( $\pm 0.01$ )	0.81
$noLA_{en}$	<b>0.81</b> ( $\pm 0.04$ )	0.79 ( $\pm 0.05$ )	0.69 ( $\pm 0.05$ )	<b>0.82</b> ( $\pm 0.07$ )	0.74 ( $\pm 0.05$ )	0.77 ( $\pm 0.05$ )	0.80 ( $\pm 0.05$ )	0.82 ( $\pm 0.05$ )	0.84 ( $\pm 0.06$ )	0.82 ( $\pm 0.05$ )	<b>0.83</b> ( $\pm 0.06$ )	0.80 ( $\pm 0.06$ )	0.83 ( $\pm 0.05$ )	0.79
$noLA_{de}$	0.79 ( $\pm 0.04$ )	0.78 ( $\pm 0.05$ )	0.68 ( $\pm 0.07$ )	0.81 ( $\pm 0.03$ )	0.78 ( $\pm 0.05$ )	0.76 ( $\pm 0.07$ )	0.80 ( $\pm 0.05$ )	0.80 ( $\pm 0.04$ )	0.84 ( $\pm 0.06$ )	0.81 ( $\pm 0.07$ )	0.82 ( $\pm 0.04$ )	0.83 ( $\pm 0.03$ )	0.84 ( $\pm 0.03$ )	0.79
$noLA_{es}$	0.79 ( $\pm 0.03$ )	0.81 ( $\pm 0.01$ )	0.70 ( $\pm 0.02$ )	<b>0.82</b> ( $\pm 0.02$ )	0.78 ( $\pm 0.03$ )	0.80 ( $\pm 0.01$ )	<b>0.84</b> ( $\pm 0.01$ )	0.83 ( $\pm 0.02$ )	0.84 ( $\pm 0.02$ )	0.83 ( $\pm 0.03$ )	0.82 ( $\pm 0.02$ )	0.83 ( $\pm 0.03$ )	0.84 ( $\pm 0.01$ )	0.81

Table 10: SIB-200 scores averaged over five random seeds for Experiment 4 with main variable *base LLM* = *Llama 3.1*. The remaining main variables are *task adaptation method* = *TA (task-specific FT)*, *adapter method* = *BN*, *LA coverage* = *monolingual*. Standard deviation in parentheses. Bold numbers indicate best scores between XLT setups (*LA*, *noLA*), underscored numbers indicate best scores within XLT setup between source languages (*en*, *de*, *es*).

## G Additional SIB-200 Results

### G.1 Heatmaps

#### G.1.1 Experiment 1: Llama-2/TA



Figure 17: Heatmap comparing SIB-200 EM *LA* and *noLA* scores across source and target languages for Experiment 1 with Llama 2 and TAs. Positive scores mean *LA* is superior.

#### G.1.2 Experiment 2: Llama-2/ICL



Figure 18: Heatmap comparing SIB-200 EM *LA* and *noLA* scores across source and target languages for Experiment 2 with Llama 2. Positive scores mean *LA* is superior.

#### G.1.3 Experiment 3: Llama-3.1/TA



Figure 19: Heatmap comparing SIB-200 EM *LA* and *noLA* scores across source and target languages for Experiment 3 with Llama 3.1. Positive scores mean *LA* is superior.



## G.2 Logit Lens Visualizations

LogitLens: Llama 2 | setup: LA | source: English | target: German

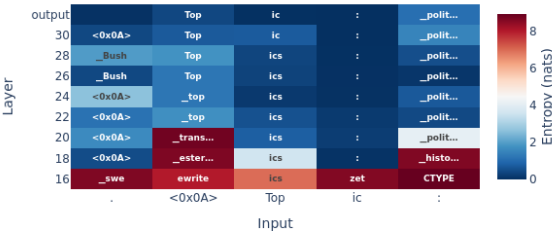


Figure 20: Logit Lens for SIB-200 test instance with *en* and German as target language. Base LLM: Llama 2. Setup: *LA*.

LogitLens: Llama 2 | setup: noLA | source: English | target: German

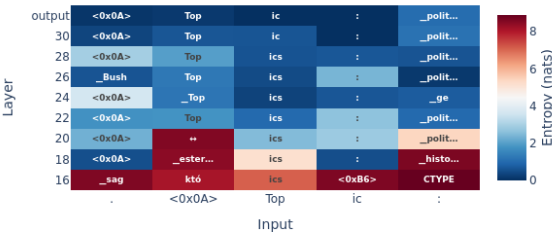


Figure 21: Logit Lens for SIB-200 test instance with *en* and German as target language. Base LLM: Llama 2. Setup: *noLA*.

LogitLens: Llama 2 | setup: LA | source: English | target: Icelandic

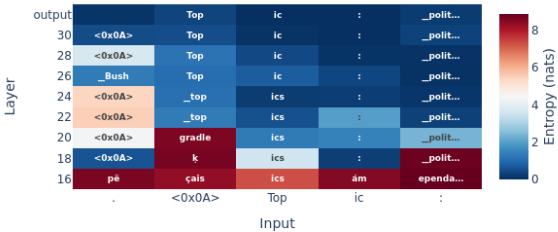


Figure 22: Logit Lens for SIB-200 test instance with *en* and Icelandic as target language. Base LLM: Llama 2. Setup: *LA*.

LogitLens: Llama 2 | setup: LA | source: English | target: Icelandic

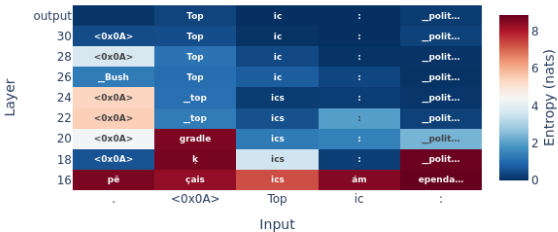


Figure 23: Logit Lens for SIB-200 test instance with *en* and Icelandic as target language. Base LLM: Llama 2. Setup: *noLA*.