

An Empirical Study on Cross-lingual Vocabulary Adaptation for Efficient Generative LLM Inference

Anonymous ACL submission

Abstract

The development of state-of-the-art generative large language models (LLMs) disproportionately relies on English-centric tokenizers, vocabulary and pre-training data. Despite the fact that some LLMs have multilingual capabilities, recent studies have shown that their inference efficiency deteriorates when generating text in languages other than English. This results in increased inference time and costs. Cross-lingual vocabulary adaptation methods have been proposed for adapting models to a target language aiming to improve downstream performance. However, the effectiveness of these methods on increasing inference efficiency of generative LLMs has yet to be explored. In this paper, we perform an empirical study of various cross-lingual vocabulary adaptation methods on five generative LLMs (including monolingual and multilingual models) across four typologically-diverse languages and four natural language understanding tasks. We find that cross-lingual vocabulary adaptation substantially contributes to LLM inference speedups of up to 271.5%. We also show that adapting LLMs that have been pre-trained on more balanced multilingual data results in downstream performance comparable to the original models.¹

1 Introduction

Generative large language models (LLMs) obtain strong generalization performance in many downstream natural language processing (NLP) tasks (OpenAI, 2023; Touvron et al., 2023a; Jiang et al., 2023) across various languages. For example, BLOOM (Scao et al., 2022) supports 46 languages while Open AI’s ChatGPT reportedly supports 90 languages (Ahuja et al., 2023).

Despite the multilingual capabilities of state-of-the-art LLMs, their development disproportionately relies on English-oriented tokenizers, vocabulary

¹Our code and models will be made publicly available on GitHub and Hugging Face Hub.

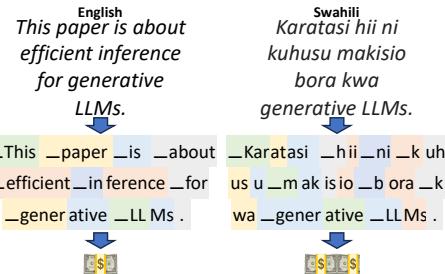


Figure 1: Example of overfragmentation when applying the Mistral-7B tokenizer to non-English text.

and pre-training data. For example, around 30% of the training data in BLOOM (i.e. a multilingual LLM) is English. This negatively affects the efficiency and downstream performance of LLMs in other languages. It has been demonstrated that LLMs overfragment text in underrepresented languages with different writing systems (Rust et al., 2021; Muller et al., 2021), resulting in increased processing time, latency and costs for non-English speakers (Ahia et al., 2023; Petrov et al., 2023). Moreover, recent studies (Lin et al., 2022; Ahuja et al., 2023; Muennighoff et al., 2023) found that LLMs often perform better in a given language other than English when prompted in English instead of prompting directly in the other language. This is an unrealistic setting for non-English speakers that introduces extra disadvantages. Figure 1 shows an illustrative example of overfragmentation in non-English text generation.

Cross-lingual vocabulary adaptation (Tran, 2019; Wang et al., 2020; Chau et al., 2020) is a resource-efficient method for cross-lingual transfer. The vocabulary of a source model is first updated (or replaced) with tokens from a target language, followed by fine-tuning the embedding matrix on data from the target language. Previous work on cross-lingual vocabulary adaptation methods primarily aims to improve downstream performance such as natural language inference and named-entity recog-

nition (Minixhofer et al., 2022; Dobler and de Melo, 2023). However, the effectiveness of these methods on improving inference efficiency of generative LLMs has yet to be explored. We hypothesize that LLM inference in a target language can be improved by adapting the vocabulary of the source model to reduce text overfragmentation.

To test our hypothesis, we perform an empirical study of various cross-lingual vocabulary adaptation methods, on four generative LLMs, including a wide range of downstream tasks, from text classification, and span prediction, to summarization in both zero-shot and few-shot settings across four languages (i.e. German, Japanese, Arabic, and Swahili). Our contributions are as follows:

- We demonstrate that cross-lingual vocabulary adaptation accelerates inference by up to 271.5% in 99% of cases (§5.1).
- We show that multilingual LLM vocabulary adaptation leads to comparable downstream performance to source models pre-trained on more balanced multilingual data (§5.2).
- We conduct an analysis to shed light on different design choices regarding the practical application of cross-lingual vocabulary adaptation in generative LLMs (§6).

2 Related Work

2.1 Impact of Tokenization on LLMs

Subword tokenization splits text into subword units and is the standard approach for tokenization in LLMs (Scao et al., 2022; Touvron et al., 2023a; Jiang et al., 2023). It includes methods such as WordPiece (Schuster and Nakajima, 2012), Byte Pair Encoding (BPE) (Sennrich et al., 2016), and Unigram (Kudo, 2018). Other approaches include word- (Bengio et al., 2000; Mikolov et al., 2013), character- (Al-Rfou et al., 2019) and byte-level (Xue et al., 2022) tokenization.

The impact of tokenization on LLMs has been actively studied including model performance (Bostrom and Durrett, 2020; Rust et al., 2021; Gow-Smith et al., 2022; Toraman et al., 2023; Fujii et al., 2023), inference speed (Hofmann et al., 2022; Sun et al., 2023; Petrov et al., 2023), memory usage (Sun et al., 2023), training (Ali et al., 2023) and API costs (Ahia et al., 2023; Petrov et al., 2023). It is acknowledged that tokenizers lead to disproportionate fragmentation for different languages and scripts in multi- and cross-lingual settings (Rust et al., 2021; Muller et al., 2021).

2.2 Cross-lingual Vocabulary Adaptation

Tran (2019) use English BERT as a source LM. They initialized target language token representations as a weighted sum of the source embeddings followed by fine-tuning both the source and target models. Wang et al. (2020) and Chau et al. (2020) added a fixed number of new target language tokens to the source vocabulary, expanding the source embedding matrix and output projection layers accordingly. The embeddings of the new tokens are randomly initialized over the expanded elements. Both studies performed additional pre-training on a target language corpus, often called language adaptive pre-training, i.e. LAPT (Chau et al., 2020), after the target vocabulary initialization. LAPT enables learning a target language model more efficiently than training it from scratch which is prohibitive with the size of current LLMs. It has become standard practice in more recent cross-lingual vocabulary adaptation studies (Minixhofer et al., 2022; Dobler and de Melo, 2023; Downey et al., 2023; Ostendorff and Rehm, 2023). More recently, state-of-the-art methods completely replace the source embeddings with target language embeddings instead of expanding the source vocabulary (Minixhofer et al., 2022; Dobler and de Melo, 2023; Ostendorff and Rehm, 2023; Downey et al., 2023). The aim is to utilize overlapping tokens between the source and target vocabularies for efficiency.

Cross-lingual vocabulary adaptation has been extensively used to adapt generative LLMs to specific target languages (Cui et al., 2023; Balachandran, 2023; Larcher et al., 2023). However, the majority of these approaches simply expand the source embedding matrix followed by LAPT, while vocabulary replacement approaches have not been explored. To the best of our knowledge, this is the first systematic study on the efficacy of various cross-lingual vocabulary adaptation methods for improving the inference efficiency of generative LLMs across languages.

3 Cross-lingual Vocabulary Adaptation

3.1 Problem Setting

Let \mathcal{M}_s be a source pre-trained LLM with \mathcal{T}_s and \mathcal{V}_s its corresponding tokenizer and vocabulary. The aim is to learn a model \mathcal{M}_t with the same architecture as \mathcal{M}_s for a target language that supports a target vocabulary \mathcal{V}_t given a tokenizer \mathcal{T}_t .

\mathcal{M}_t is first initialized with the weights of \mathcal{M}_s . Subsequently, its input embedding and output layer

matrices are replaced such that the former is of dimensionality $|\mathcal{V}_t| \times H_t$ and the latter $H_t \times |\mathcal{V}_t|$, where H_t is the hidden dimensionality of \mathcal{M}_t . The weights of both matrices are tied and they can be initialized by applying one of the target vocabulary initialization methods in §3.2. Finally, \mathcal{M}_t is adapted to the target language (i.e. with LAPT) by training it on target language data \mathcal{D} using a causal language modeling objective.

3.2 Target Vocabulary Initialization Methods

Random. The simplest approach is to randomly initialize the embeddings of \mathcal{M}_t (de Vries and Nissim, 2021; Downey et al., 2023).

Cross-lingual and Progressive Initialization (CLP). CLP (Ostendorff and Rehm, 2023) first finds overlapping tokens between \mathcal{V}_t and \mathcal{V}_s , i.e. $\mathcal{V}_t \cap \mathcal{V}_s$, and simply copies their weights from \mathcal{M}_s to \mathcal{M}_t . Each target token that does not overlap with any source token, i.e. $\mathcal{V}_t \setminus (\mathcal{V}_t \cap \mathcal{V}_s)$ is initialized by its weighted average across all embeddings in $\mathcal{V}_t \cap \mathcal{V}_s$, i.e. common tokens in the source and target vocabularies. The weight of each embedding in $\mathcal{V}_t \cap \mathcal{V}_s$ is computed as the cosine similarity score between the respective overlapping token and the target non-overlapping token. Since there is no common representation between the non-overlapping token and the overlapping tokens, CLP employs vector representations from an auxiliary target language-specific pre-trained LM with the same tokenizer and vocabulary as \mathcal{M}_t so that both tokens are mapped.

Heuristics. Downey et al. (2023) proposed a heuristic-based initialization that consists of the following three rules according to IDENTITY, SCRIPT, and POSITION of a token. First, embeddings are initialized according to their identity, in the same way that overlapping tokens are initialized in CLP, i.e. by copying them from \mathcal{M}_s . For all the remaining tokens in \mathcal{V}_t , their embeddings are initialized based on the type of SCRIPT identified by the Unicode block. Each token that belongs to a particular script (e.g. Hebrew) is represented by a vector sampled from a Normal distribution with the same mean and standard deviation. The mean and standard deviation of each group are computed using the embeddings of \mathcal{M}_s that belong to the same group. In conjunction with SCRIPT, a group can further be divided into two according to the POSITION of each subword token in a word, i.e. whether it is placed at the beginning or in the middle (e.g. “_the” vs.

“the”). Finally, the embeddings of any remaining tokens are randomly initialized.

FOCUS. Dobler and de Melo (2023) proposed fast overlapping token combinations using sparsemax (FOCUS) initialization. It is an approach similar to CLP that reuses the embeddings of \mathcal{M}_s in \mathcal{M}_t for tokens in $\mathcal{V}_t \cap \mathcal{V}_s$. For non-overlapping tokens $\mathcal{V}_t \setminus (\mathcal{V}_t \cap \mathcal{V}_s)$, FOCUS uses fastText (Bojanowski et al., 2017) trained on target specific data \mathcal{D} tokenized by \mathcal{T}_t and computes the cosine similarities between tokens in $\mathcal{V}_t \cap \mathcal{V}_s$ and $\mathcal{V}_t \setminus (\mathcal{V}_t \cap \mathcal{V}_s)$ in the fastText model. It then applies sparsemax (Martins and Astudillo, 2016), which is a sparse variant of softmax that assigns zero to any low-probability elements, over the similarity scores. The embeddings of tokens in $\mathcal{V}_t \setminus (\mathcal{V}_t \cap \mathcal{V}_s)$ are finally initialized by taking the weighted sum of the source embeddings of tokens in $\mathcal{V}_t \cap \mathcal{V}_s$, whose weights are the similarity scores with sparsemax applied.

CLP+. Finally, we propose *CLP+*, a modification to CLP motivated by the use of sparsemax in FOCUS. The aim is to dynamically select semantically similar tokens from $\mathcal{V}_t \cap \mathcal{V}_s$ to initialize a target embedding for a token in $\mathcal{V}_t \setminus (\mathcal{V}_t \cap \mathcal{V}_s)$, leading to a better initialization of the embeddings (Tran, 2019). We follow the same process as CLP for tokens in $\mathcal{V}_t \cap \mathcal{V}_s$. For non-overlapping tokens in $\mathcal{V}_t \setminus (\mathcal{V}_t \cap \mathcal{V}_s)$, instead of taking the weighted average of *all* overlapping source embeddings of $\mathcal{V}_t \cap \mathcal{V}_s$ as in CLP, we use the weighted sum of embeddings whose weight is calculated with sparsemax.

4 Experimental Setup

4.1 Source Models

We first use **BLOOM-1B** and **BLOOM-7B** (Scao et al., 2022) as source models, which are trained on data from 46 languages including Arabic (4.6%) and Swahili (0.02%). We also use **TigerBot-7B** (Chen et al., 2023), which is based on LLaMA 2 (Touvron et al., 2023b) adapted using data from East Asian languages, i.e. Chinese (54%), Korean (0.001%), and Japanese (0.01%). Finally, we experiment with **Mistral-7B** (Jiang et al., 2023) which is an English-centric model. Table 1 shows the tokenizer and vocabulary size of each source model.

4.2 Target Languages and Adaptation Data

We experiment with a typologically diverse set of target languages including German (Indo-European), Japanese (Japonic), Arabic

Source (\mathcal{M}_s)	Tokenizer (\mathcal{T}_s)	$ \mathcal{V}_s $
BLOOM	Byte-level BPE	250,680
TigerBot	Byte-level BPE	60,512
Mistral	Byte-level BPE	32,000
Target (\mathcal{M}_t)	Tokenizer (\mathcal{T}_t)	$ \mathcal{V}_t $
German	Byte-level BPE	50,257
Japanese	Unigram	32,000
Arabic	Byte-level BPE	64,000
Swahili	Byte-level BPE	50,257

Table 1: Tokenizers and vocabulary size for source and target models.

(Afro–Asiatic), and Swahili (Niger–Congo). We opted to use these languages because of the availability of language-specific (1) tokenizers; and (2) downstream task datasets with the same task formulation across languages.²

For adapting the source models, we use the OSCAR language-specific subcorpus (Jansen et al., 2022) for German, Arabic, and Japanese (January 2023 version). For Swahili, we use the Swahili subset of CC-100 (Conneau et al., 2020) following Minixhofer et al. (2022). We use publicly available existing tokenizers and vocabularies for each target language. Table 1 shows the type of tokenizer and vocabulary size for the target models. More details are available in Table 3 in the Appendix.

4.3 Tasks

Following Ahia et al. (2023), we evaluate all models including baselines across four tasks in each target language: (1) textual entailment (NLI) consisting of JNLI (Kurihara et al., 2022) for Japanese and XNLI (Conneau et al., 2018) for the rest; (2) X-CSQA (Lin et al., 2021) for multiple choice question-answering (MC); (3) summarization (SUM) including MLSUM (Scialom et al., 2020) for German and XL-Sum (Hasan et al., 2021) for the rest; and (4) span prediction (SPAN) consisting of XQuAD (Kurihara et al., 2022) for Arabic and German, JSQuAD (Kurihara et al., 2022) for Japanese and KenSwQuAD (Wanjawa et al., 2023) for Swahili. Similarly, we use 500 random samples.

Due to the computational constraints, we conduct zero-shot experiments on SUM. For the rest of the tasks, we evaluate models in zero- and few-shot settings. We use five demonstrations for NLI and MC and three for SPAN in the few-shot cases.

²Note that data for the same task across languages does not match. Model performance is not directly comparable.

4.4 Prompt Templates

We use the same English prompt templates as Ahia et al. (2023) for NLI and SUM. For MC and SPAN, we formulate a task-specific English prompt. We translate the English prompt templates into each corresponding target language using a machine translation API (i.e. Google Translate), following Yong et al. (2023). The prompt templates can be found in Appendix A.7.

4.5 Baselines

We compare the cross-lingual vocabulary adaptation methods against two baselines: (1) we use the source models directly on the target language tasks without any adaptation (**Source**); and (2) following Yong et al. (2023), we adapt the source models by continuing pre-training on data from a target language by keeping the source vocabulary (**LAPT**).

4.6 Evaluation Metrics

Inference Efficiency. We calculate the average number of prompt tokens per sample for each dataset and tokenizer, and use its relative ratio to each source tokenizer as a proxy for inference speedup following (Ahia et al., 2023; Petrov et al., 2023). We use the average number of prompt tokens rather than the actual inference time because commercial APIs (e.g. OpenAI) often charge users on the basis of the total number of prompt and generated tokens. Note that inference efficiency is independent of the model size.

Downstream Performance. For downstream performance evaluation, we use standard metrics for each dataset such as accuracy for NLI and MC, F1 for SPAN, and ROUGE-L for SUM.

4.7 Implementation Details

Efficient LLM Adaptation. We perform our experiments under resource-constrained settings due to limited access to computational resources. For computational efficiency, we use a low-rank adaptation approach with LAPT, i.e. LoRA (Hu et al., 2021) applied on all linear layers (setting rank $r = 8$), following (Yong et al., 2023; Cui et al., 2023; Balachandran, 2023; Larcher et al., 2023). We pre-train each model for a maximum of four days. We use a batch size of 8 for BLOOM-1B and 16 for the 7B models with gradient accumulation steps set to 4 and a sequence length of 1,024. We set the learning rate to 1e-4 and save checkpoints every 1,000 steps. For a fair comparison,

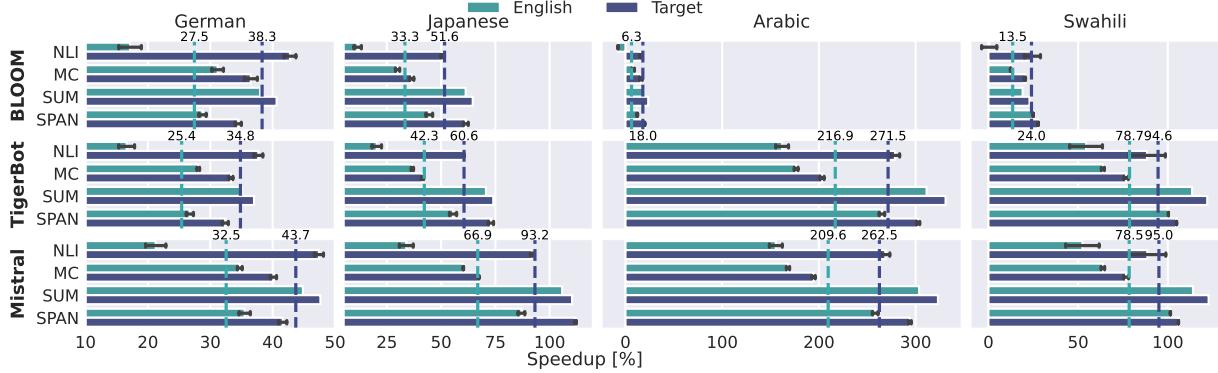


Figure 2: Relative speedup ratios to each base model/tokenizer when prompted in English and a target language. Dotted lines denote the average speedup ratio across tasks in each setting.

we use the checkpoints with the largest number of steps available across all vocabulary initialization approaches and the LAPT baseline for the same source model size and language.³

Libraries and Hardware. We implement our models using PyTorch (Paszke et al., 2019), HF Transformers (Wolf et al., 2020) and PEFT (Mangrulkar et al., 2022). We preprocess data with HF Datasets (Lhoest et al., 2021). We use a single NVIDIA A100 (80GB) GPU for all experiments.

5 Results

5.1 Inference Efficiency

Figure 2 shows the relative inference speedup ratio between the target and source model prompted in the target language and English. Overall, the results confirm our hypothesis that cross-lingual vocabulary adaptation accelerates inference in 95 out of 96 cases including zero- and few-shot settings.

Next, we examine the efficiency of target models in each language. We first observe that the models adapted in German show moderate average speedup ratios (25.4-43.7%) across different tasks, source models and prompting languages. This is possibly due to the close relationship between German and English (i.e. both Germanic and Indo-European languages). The Japanese target models also exhibit moderate but slightly greater average speedups compared to German of up to 60.6% using BLOOM and TigerBot as source models. In contrast, inference speedups are substantially greater using Mistral as the source model (66.9-93.2% on average). These differences can stem from the inclusion of Chinese pre-training data in

BLOOM,⁴ and Chinese and Japanese data in TigerBot. Target models in Arabic and Swahili obtain smaller speedups than the other languages (i.e. up to 24.0% on average) using BLOOM as the source model. This is likely due to the inclusion of Arabic and Swahili pre-training data in BLOOM. In contrast, target models in both languages obtain substantial gains using TigerBot and Mistral as the source models, i.e. up to an impressive 271.5% for Arabic and 95.0% for Swahili. We speculate that this happens because the two languages are not included in the training data of TigerBot and Mistral⁵, and the different Arabic script.

Looking into individual tasks, we observe that adapted models gain larger speedup ratios in SPAN and SUM compared to the other two tasks across source models and languages. In particular, we record a maximum speedup of 331% in Arabic SUM with in-language prompting using TigerBot as the source model. In contrast, speedup ratios tend to be smaller than average in NLI and MC across different target and source models, except for NLI with in-language prompting. Specifically, the Arabic model using BLOOM as source shows a slowdown of 7.63% when prompting in English. Our hypothesis is that this is due to the ratio of English-related words included in a prompt in each task, resulting to overfragmentation of such words by a *target-language* tokenizer, which is detrimental to inference speedup. Indeed, the number of tokens of the NLI English prompt template⁶ is ten when tokenized with the BLOOM source tokenizer,

⁴Note that the Japanese script includes Chinese characters.

⁵We do not have enough information about the training data of the model. Mistral.ai states that it is good at English tasks according to the blog post.

⁶Question: True, False, or Neither? Answer:

³For more details, refer to Appendix A.4.

Approach	German				Japanese				Arabic				Swahili			
	NLI	MC	SUM	SPAN	NLI	MC	SUM	SPAN	NLI	MC	SUM	SPAN	NLI	MC	SUM	SPAN
BLOOM-1B																
Source	.33	.21	17.8	.06	.29	.20	18.2	.22	.35	.20	12.0	.15	.32	.22	12.0	.03
LAPT	.36	.22	14.3	.09	.28	.20	20.7	.26	.36	.19	11.4	.13	.31	.18	7.7	.07
+ Random	.35_{.61}	.22_{.48}	15.3 _{.15.0}	.14_{.13.3}	.29_{.00}	.21_{.5.0}	19.0 _{.5.6}	.32_{.15.5}	.35_{.00}	.19_{.5.0}	11.5 _{.4.2}	.14_{.6.7}	.33_{.3.1}	.22_{.0.0}	10.2 _{.15.0}	.08_{.66.7}
+ Heuristics	.34_{.3.0}	.19_{.5.5}	15.3 _{.15.0}	.13_{.11.6}	.29_{.00}	.19_{.5.0}	19.2 _{.6.7}	.31_{.40.9}	.36_{.9}	.22_{.10.0}	11.3 _{.5.8}	.13_{.3.3}	.34_{.4.3}	.22_{.0.0}	11.9 _{.0.8}	.11_{.366.7}
+ CLP	.34_{.3.0}	.19_{.4.3}	14.6 _{.18.9}	.14_{.13.3}	.29_{.00}	.25_{.25.0}	18.8 _{.4.4}	.33_{.50.0}	.36_{.9}	.21_{.5.0}	11.2 _{.6.7}	.14_{.6.7}	.34_{.4.3}	.22_{.0.0}	11.5 _{.4.2}	.11_{.366.7}
+ FOCUS	.34_{.3.0}	.19_{.5.5}	16.1 _{.10.6}	.13_{.11.6}	.29_{.00}	.21_{.5.0}	19.2 _{.6.7}	.33_{.50.0}	.35_{.00}	.20_{.00}	11.2 _{.6.7}	.14_{.6.7}	.34_{.4.3}	.22_{.0.0}	11.2 _{.6.7}	.12_{.300.0}
+ CLP+	.31_{.6.1}	.15_{.28.6}	15.8 _{.12.2}	.13_{.11.6}	.29_{.00}	.19_{.5.0}	19.4 _{.7.8}	.33_{.50.0}	.35_{.00}	.17_{.15.0}	11.3 _{.5.8}	.15_{.0.0}	.34_{.6.3}	.20_{.9.1}	10.4 _{.13.3}	.10_{.233.3}
BLOOM-7B																
Source	.36	.21	23.1	.15	.28	.21	19.0	.33	.38	.17	11.5	.25	.36	.22	14.3	.22
LAPT	.37	.21	19.4	.14	.21	.21	21.6	.36	.36	.16	11.5	.21	.36	.20	13.0	.14
+ Heuristics	.32_{.11.1}	.22_{.4.8}	19.7 _{.14.3}	.21_{.40.0}	.30_{.7.1}	.23_{.9.5}	19.5 _{.2.6}	.38_{.15.2}	.37_{.2.6}	.19_{.11.8}	10.7 _{.2.7}	.21_{.16.0}	.35_{.8.3}	.22_{.0.0}	11.6 _{.17.1}	.16_{.27.3}
+ CLP+	.35_{.8.8}	.21_{.0.0}	18.7 _{.18.7}	.20_{.33.3}	.21_{.0.0}	.19_{.5.2}	.40_{.21.2}	.38_{.0.0}	.21_{.23.5}	.11.0.0	.21_{.16.0}	.35_{.8.3}	.23_{.4.5}	.10.9.21	.17_{.22.7}	
TigerBot-7B																
Source	.33	.24	23.9	.26	.17	.24	19.4	.57	.33	.21	9.0	.04	.29	.22	12.4	.03
LAPT	.32	.21	18.5	.18	.17	.21	21.6	.49	.33	.18	9.8	.13	.31	.21	15.9	.10
+ Heuristics	.35_{.6.1}	.20_{.16.7}	16.1 _{.32.9}	.18_{.30.8}	.29_{.7.6}	.22_{.8.3}	19.6 _{.3.2}	.40_{.28.9}	.36_{.1}	.17_{.19.0}	.10.3_{.14.4}	.08_{.00.0}	.36_{.4.1}	.22_{.0.0}	8.1 _{.3.2}	.05_{.66.7}
+ CLP+	.33_{.0.0}	.20_{.16.7}	14.1 _{.41.2}	.19_{.26.9}	.29_{.7.6}	.20_{.6.7}	19.8 _{.4.2}	.41_{.28.1}	.38_{.5.2}	.22_{.4.8}	.11.2_{.24.4}	.16_{.300.0}	.30_{.3.4}	.22_{.0.0}	8.6 _{.28.3}	.09_{.300.0}
Mistral-7B																
Source	.34	.25	24.1	.35	.17	.28	23.7	.60	.33	.20	11.2	.21	.35	.22	15.4	.07
LAPT	.33	.25	24.2	.28	.17	.20	23.4	.60	.33	.18	10.8	.14	.33	.22	16.2	.12
+ Heuristics	.40_{.17.6}	.26_{.4.0}	21.2 _{.11.7}	.22_{.37.1}	.29_{.7.6}	.20_{.28.6}	19.7 _{.7.9}	.43_{.28.3}	.39_{.8.2}	.19_{.5.0}	.10.7_{.2.7}	.15_{.38.1}	.34_{.2.9}	.22_{.0.0}	10.6 _{.29.3}	.14_{.100.0}
+ CLP+	.35_{.9}	.25_{.0}	20.2 _{.15.8}	.21_{.40.0}	.28_{.6.4}	.20_{.28.6}	19.9 _{.7.1}	.46_{.23.3}	.38_{.5.2}	.16_{.20.0}	.11.5_{.4.5}	.21_{.0.0}	.35_{.7}	.21_{.4.5}	10.2 _{.32.0}	.16_{.28.6}
BLOOM-1B																
Source	.38	.20	-	.10	.44	.19	-	.32	.35	.17	-	.20	.37	.23	-	.02
LAPT	.37	.17	-	.13	.26	.21	-	.34	.34	.16	-	.16	.34	.19	-	.02
+ Random	.34_{.10.5}	.21_{.5.0}	-	.16 _{.0.0}	.29_{.4.1}	.21_{.10.5}	-	.34_{.6.3}	.35_{.0}	.22_{.9.4}	-	.16 _{.0.0}	.34_{.8.1}	.20_{.13.0}	-	.06_{.90.0}
+ Heuristics	.33_{.13.2}	.22_{.15.0}	-	.17 _{.7.0}	.30_{.31.8}	.22_{.15.8}	-	.32_{.0}	.36_{.2.9}	.21_{.23.5}	-	.15_{.25.0}	.32_{.0.8}	.19_{.17.4}	-	.07_{.250.0}
+ CLP	.34_{.10.5}	.21_{.5.0}	-	.17 _{.7.0}	.30_{.31.8}	.20_{.25.3}	-	.33_{.1.1}	.35_{.0.0}	.21_{.23.5}	-	.15_{.25.0}	.32_{.10.8}	.19_{.17.4}	-	.06_{.300.0}
+ FOCUS	.34_{.10.5}	.18_{.10.0}	-	.17 _{.7.0}	.27_{.38.6}	.20_{.25.3}	-	.36_{.12.5}	.36_{.2.9}	.20_{.17.6}	-	.15_{.25.0}	.34_{.8.1}	.19_{.17.4}	-	.06_{.300.0}
+ CLP+	.34_{.10.5}	.20_{.0.0}	-	.19 _{.9.0}	.29_{.34.1}	.22_{.15.8}	-	.36_{.12.5}	.37_{.5.7}	.20_{.17.6}	-	.15_{.25.0}	.30_{.18.9}	.18_{.21.7}	-	.06_{.300.0}
BLOOM-7B																
Source	.38	.23	-	.29	.41	.19	-	.49	.37	.18	-	.29	.34	.18	-	.11
LAPT	.35	.24	-	.23	.34	.19	-	.53	.36	.18	-	.23	.37	.18	-	.07
+ Heuristics	.33_{.13.2}	.22_{.4.3}	-	.28 _{.3.4}	.28_{.31.7}	.21_{.10.5}	-	.46_{.6.1}	.36_{.7}	.21_{.16.7}	-	.24_{.7.2}	.36_{.9}	.19_{.5.6}	-	.13_{.18.2}
+ CLP+	.34_{.10.5}	.22_{.4.3}	-	.25 _{.13.8}	.30_{.26.8}	.20_{.5.3}	-	.46_{.6.1}	.36_{.7}	.22_{.22.2}	-	.25_{.13.8}	.36_{.5.9}	.18_{.0.0}	-	.13_{.18.2}
TigerBot-7B																
Source	.31	.37	-	.42	.16	.34	-	.65	.30	.19	-	.10	.36	.19	-	.03
LAPT	.30	.39	-	.27	.16	.34	-	.66	.30	.20	-	.17	.36	.20	-	.09
+ Heuristics	.33_{.6.5}	.26_{.29.7}	-	.21 _{.8.0}	.29_{.31.2}	.24_{.29.4}	-	.49_{.24.6}	.35_{.16.7}	.20_{.5.3}	-	.09_{.10.0}	.35_{.2.8}	.21_{.10.5}	-	.04_{.33.3}
+ CLP+	.36_{.16.1}	.31_{.16.2}	-	.31 _{.26.2}	.30_{.37.5}	.21_{.38.2}	-	.50_{.23.1}	.37_{.23.3}	.19_{.0.0}	-	.19_{.90.0}	.34_{.5.6}	.18_{.5.3}	-	.06_{.100.0}
Mistral-7B																
Source	.33	.53	-	.48	.16	.42	-	.69	.30	.32	-	.31	.40	.21	-	.12
LAPT	.33	.46	-	.27	.16	.37	-	.68	.30	.30	-	.26	.36	.34	-	.21
+ Heuristics	.45_{.36.4}	.41_{.22.6}	-	.24 _{.50.0}	.30_{.37.5}	.24_{.42.9}	-	.49_{.29.0}	.34_{.3.3}	.18_{.43.8}	-	.17_{.45.2}	.34_{.15.0}	.18_{.14.3}	-	.09_{.25.0}
+ CLP+	.37_{.12.1}	.47_{.11.3}	-	.25 _{.47.9}	.29_{.81.2}	.25_{.40.5}	-	.50_{.27.5}	.38_{.26.7}	.23_{.28.1}	-	.23_{.25.8}	.33_{.17.5}	.20_{.4.8}	-	.14_{.16.7}

Table 2: Mean performance over five runs with in-language prompting on 500 randomly selected test samples for each dataset. The baselines are in grey. **Bold** indicates comparable or better results than the baselines. **Green** indicates positive relative performance change over Source. **Red** denotes negative relative performance change.

increasing to 21 with the Arabic target tokenizer.

Finally, we investigate the inference efficiency for the target models by prompting language. We observe that the target models show greater inference speedup ratios with in-language prompts than English in all cases. The average differences between in-language and English prompts are 12.8%, 24.6%, and 26.8%, using BLOOM, TigerBot, and Mistral as source models, respectively. This suggests that the target models are susceptible to code-mixed text (i.e. including English prompts), leading to overfragmentation for words not in a target language. Furthermore, in-language prompting is a more realistic scenario for non-English speakers to use LLMs than English prompting. Therefore, these differences reflect the advantage of cross-lingual vocabulary adaptation and confirm the disparity of using a source tokenizer, reported by Ahia et al. (2023) and Petrov et al. (2023).

5.2 Downstream Performance

We compare the downstream performance of all cross-lingual vocabulary adaptation methods (§3.2) to the baselines, i.e. Source and LAPT (§4.5) using BLOOM-1B as the source model. Due to computa-

tional costs, we only evaluate the best two vocabulary adaptation approaches (Heuristics and CLP+) against the baselines with larger source models (BLOOM-7B, TigerBot-7B and Mistral-7B). Table 2 shows the zero- and few-shot performance of all models with in-language prompting. Results using English prompts are included in the Appendix (Table 9). We examine differences between English and target language prompting in §6.

Overall, adapted models show comparable or better performance than the baselines in the majority of the cases across tasks and languages using BLOOM-1B as source. Models adapted with simple Random target vocabulary initialization are competitive compared to more sophisticated approaches and the baselines in the majority of the cases – 18 for Source and 24 for LAPT out of 28 cases, respectively. However, they are not as robust with English prompting (see Table 9 in the Appendix). Models adapted with Heuristics also perform competitively with the semantic similarity-based methods (i.e. CLP, FOCUS and CLP+). They are similar to or better than Source in 17 out of 28 cases and LAPT in 19 out of 28 cases without a substantial drop in performance observed in Ran-

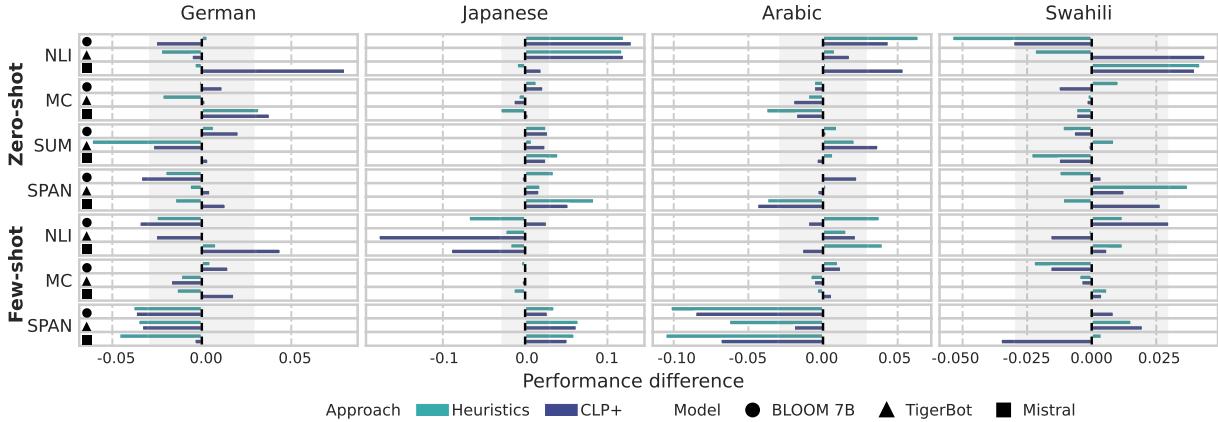


Figure 3: Performance difference between English and in-language prompts. Positive and negative values indicate better performance using in-language or English prompts respectively.

dom with English prompting. These results generally corroborate findings by Downey et al. (2023) that (1) Random initialization is not robust, and (2) Heuristics rivals the semantic similarity-based methods. CLP+ outperforms CLP and FOCUS in the few-shot cases excluding Random, while Heuristics is the best zero-shot approach.

Experimenting with larger source models, we observe that models adapted with Heuristics perform competitively with CLP+ with BLOOM-7B. They also show competitive performance with LAPT in 18 out of 28 cases across tasks and languages. However, they exhibit slightly worse performance than Source in 15 out of 28 cases. Models adapted with CLP+, in contrast, perform on par with Source (14 out of 28 cases) and slightly better than LAPT (17 out of 28 cases). When using TigerBot and Mistral as source, we find that adapted models fail to perform better than the baselines. The only exception is TigerBot, where models adapted with Heuristics and CLP+ show competitive performance with Source in 15 and 16 out of 28 cases across tasks and languages. This suggests that LLMs not as multilingual as BLOOM (i.e. TigerBot and Mistral) may not perform well with cross-lingual vocabulary adaptation, possibly due to less transferable cross-lingual knowledge, and the small amount of target language data included during pre-training.

6 Analysis

In-language vs. English Prompting. Recent studies (Lin et al., 2022; Ahuja et al., 2023; Muenighoff et al., 2023) report that in-language prompting yields lower downstream performance than prompting in English. Figure 3 shows the perfor-

mance difference between English and in-language prompting across models, languages and tasks for the best two performing cross-lingual vocabulary adaptation methods (Heuristics and CLP+).

Surprisingly, we observe no major performance drop with in-language prompting in the majority of the zero-shot settings across languages. We note similar or better performance with in-language prompting in 10 out of 16 settings. We also observe a drop of 0.03 or larger (non-shaded areas in Figure 3) in only 5 out of 16 settings. The few-shot settings also exhibited similar trends with substantial performance degradation of 0.03 or more in 5 out of 12 settings, and similar or better performance in 8 out of 12 settings. Some in-language prompting cases with lower performance than English, such as in German across tasks and zero-shot NLI in Arabic and Swahili, can be related to the tokenization effects discussed in §5.1. Previous studies (Rust et al., 2021; Bostrom and Durrett, 2020; Fujii et al., 2023) have also found a strong correlation between tokenization and performance. Further investigation is needed on how to mitigate performance degradation, especially in few-shot settings such as Japanese NLI and Arabic SPAN.

LAPT Steps. LAPT is an integral part of recent cross-lingual vocabulary adaptation methods (Minixhofer et al., 2022; Dobler and de Melo, 2023; Ostendorff and Rehm, 2023; Downey et al., 2023). However, it is a computationally intensive task that requires loading and updating models with billions of parameters over a large number of training steps. Therefore, we investigate the relationship between downstream performance and the number of LAPT steps, i.e. every 2k steps starting from

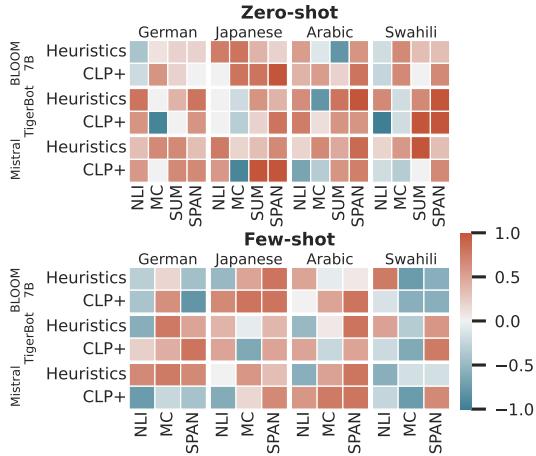


Figure 4: Kendall’s τ correlation between the number of LAPT steps and performance (in-language prompting).

1k and every 10k after 13k. Figure 4 shows the Kendall’s tau (Kendall, 1938) correlation coefficients between LAPT steps and performance.

Overall, LAPT helps improve downstream performance in both zero- and few-shot settings in 69.5% and 59.4% of the cases, respectively.⁷ In particular, both TigerBot and Mistral, which are not as multilingual as BLOOM, tend to benefit more from LAPT, especially in zero-shot SUM and SPAN across languages. This suggests that LAPT helps source LLMs to supplement target language knowledge to reach competitive performance with BLOOM 7B in a similar number of steps.

Next, we examine the correlation coefficients by task across zero- and few-shot settings. We often observe negative or no correlation in zero-shot MC across languages, zero-shot NLI in Japanese and Swahili, few-shot NLI in German, Japanese and Swahili, few-shot MC in Swahili, and few-shot SPAN in German and Swahili. In contrast, SPAN and SUM generally benefit well from LAPT in zero-shot settings across languages, in addition to few-shot SPAN in Japanese and Arabic. We hypothesize that because zero-shot tasks, especially SUM and SPAN, can be more challenging than the other text classification tasks (Davletov et al., 2021; Yamaguchi et al., 2022), they may require better target language representations to perform well.

LoRA Rank r . There is a trade-off between computational efficiency and performance when adapting LLMs with LoRA (Hu et al., 2021). We analyze how the LoRA rank affects performance

⁷We observe similar trends with English prompting (Figure 6 in the Appendix).

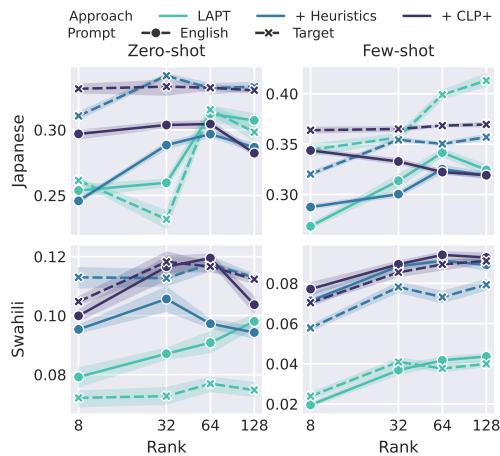


Figure 5: SPAN performance changes with respect to LoRA rank r .

in cross-lingual vocabulary adaptation. To keep computational costs low, we experiment by setting $r = \{8, 32, 64, 128\}$ using BLOOM-1B on SPAN in Japanese and Swahili where we observed large performance variations (Table 2). Figure 5 shows how performance changes with respect to r .

In general, the performance of target models does not increase with r in the zero-shot setting. The only exception is Heuristics in Japanese with English prompting, i.e. from 0.246 ($r = 8$) to 0.297 ($r = 64$). We observe that performance improves with r in the few-shot setting. CLP+ in Japanese with English prompting is an exception, i.e. from 0.344 ($r = 8$) to 0.319 ($r = 128$). This suggests that setting $r = 8$ is a good option in zero-shot settings. Increasing r to 32, 64 or 128 can yield better few-shot performance but results to higher computational costs. For instance, the best-performing Swahili model ($r = 64$) results in a 14% increase in the number of trainable parameters compared to $r = 8$. Careful cost-benefit consideration is needed to choose an optimal r .

7 Conclusion

We investigated the effectiveness of cross-lingual vocabulary adaptation on LLM inference efficiency. Our extensive experiments in four diverse languages demonstrated that cross-lingual vocabulary adaptation substantially contributes to LLM inference speedups of up to 271.5% while maintaining comparable downstream performance to baselines when adapting multilingual LLMs. In future work, we will explore various inference-aware methods for cross-lingual transfer, such as cost-effective subword vocabulary selection (Gee et al., 2023).

599 Limitations

600 **Prompt Tuning.** We use a translated version of
601 in-language prompts from English. This may af-
602 fect the downstream performance due to machine
603 translation noise, underestimating the performance
604 of in-language prompting.

605 **Languages.** Although this study covers four lin-
606 guistically diverse languages (German, Arabic,
607 Japanese, and Swahili), it is an interesting study for
608 future work to assess more languages.

609 **Model Size.** This paper considers LLMs of var-
610 ious sizes ranging from 1B to 7B, which are far
611 larger than those tested in previous cross-lingual
612 vocabulary adaptation studies. Note that infer-
613 ence efficiency measured by the number of pro-
614 cessed/generated tokens is not affected by the
615 model size. However, investigating the perfor-
616 mance of cross-lingual vocabulary adaptation ap-
617 proaches with larger models would be valuable in
618 future studies.

619 References

620 Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo
621 Kasai, David Mortensen, Noah Smith, and Yulia
622 Tsvetkov. 2023. **Do all languages cost the same?**
623 **tokenization in the era of commercial language mod-**
624 **els.** In *Proceedings of the 2023 Conference on Em-*
625 *pirical Methods in Natural Language Processing*,
626 pages 9904–9923, Singapore. Association for Com-
627 putational Linguistics.

628 Kabir Ahuja, Harshita Diddee, Rishav Hada, Milli-
629 cent Ochieng, Krithika Ramesh, Prachi Jain, Ak-
630 shay Nambi, Tanuja Ganu, Sameer Segal, Mohamed
631 Ahmed, Kalika Bali, and Sunayana Sitaram. 2023.
632 **MEGA: Multilingual evaluation of generative AI.**
633 In *Proceedings of the 2023 Conference on Empirical*
634 *Methods in Natural Language Processing*, pages
635 4232–4267, Singapore. Association for Compu-
636 tational Linguistics.

637 Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy
638 Guo, and Llion Jones. 2019. **Character-level lan-**
639 **guage modeling with deeper self-attention.** *Proced-*
640 *ings of the AAAI Conference on Artificial Intelligence*,
641 33(01):3159–3166.

642 Mehdi Ali, Michael Fromm, Klaudia Thellmann,
643 Richard Rutmann, Max Lübbing, Johannes
644 Leveling, Katrin Klug, Jan Ebert, Niclas Doll,
645 Jasper Schulze Buschhoff, Charvi Jain, Alexan-
646 der Arno Weber, Lena Jurkschat, Hammam Abdel-
647 wahab, Chelsea John, Pedro Ortiz Suarez, Malte
648 Ostendorff, Samuel Weinbach, Rafet Sifa, Stefan
649 Kesselheim, and Nicolas Flores-Herr. 2023. **Tok-**
650 **enizer choice for LLM training: Negligible or cru-**
651 **cial?** *ArXiv*, abs/2310.08754.

Wissam Antoun, Fady Baly, and Hazem Hajj. 2021.	652
AraGPT2: Pre-trained transformer for Arabic lan-	653
guage generation. In <i>Proceedings of the Sixth Ara-</i>	654
<i>bic Natural Language Processing Workshop</i> , pages	655
196–207, Kyiv, Ukraine (Virtual). Association for	656
Computational Linguistics.	657
Abhinand Balachandran. 2023. Tamil-LLaMA: A new	658
Tamil language model based on LLaMA 2. <i>ArXiv</i> ,	659
abs/2311.05845.	660
Yoshua Bengio, Réjean Ducharme, and Pascal Vincent.	661
2000. A neural probabilistic language model. In	662
<i>Advances in Neural Information Processing Systems</i> ,	663
volume 13. MIT Press.	664
Piotr Bojanowski, Edouard Grave, Armand Joulin, and	665
Tomas Mikolov. 2017. Enriching word vectors with	666
subword information. <i>Transactions of the Associa-</i>	667
<i>tion for Computational Linguistics</i> , 5:135–146.	668
Kaj Bostrom and Greg Durrett. 2020. Byte pair encod-	669
ing is suboptimal for language model pretraining. In	670
<i>Findings of the Association for Computational Lin-</i>	671
<i>guistics: EMNLP 2020</i> , pages 4617–4624, Online.	672
Association for Computational Linguistics.	673
Ethan C. Chau, Lucy H. Lin, and Noah A. Smith. 2020.	674
Parsing with multilingual BERT, a small corpus, and	675
a small treebank. In <i>Findings of the Association</i>	676
<i>for Computational Linguistics: EMNLP 2020</i> , pages	677
1324–1334, Online. Association for Computational	678
Linguistics.	679
Ye Chen, Wei Cai, Liangmin Wu, Xiaowei Li, Zhanxuan	680
Xin, and Cong Fu. 2023. TigerBot: An open multi-	681
lingual multitask LLM. <i>ArXiv</i> , abs/2312.08688.	682
Alexis Conneau, Kartikay Khandelwal, Naman Goyal,	683
Vishrav Chaudhary, Guillaume Wenzek, Francisco	684
Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer,	685
and Veselin Stoyanov. 2020. Unsupervised	686
cross-lingual representation learning at scale. In <i>Pro-</i>	687
<i>ceedings of the 58th Annual Meeting of the Associa-</i>	688
<i>tion for Computational Linguistics</i> , pages 8440–	689
8451, Online. Association for Computational Lin- guistics.	690
Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina	691
Williams, Samuel Bowman, Holger Schwenk, and	692
Veselin Stoyanov. 2018. XNLI: Evaluating cross-	693
lingual sentence representations. In <i>Proceedings of</i>	694
<i>the 2018 Conference on Empirical Methods in Natu-</i>	695
<i>ral Language Processing</i> , pages 2475–2485, Brus- sels, Belgium. Association for Computational Lin- guistics.	696
Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient	697
and effective text encoding for Chinese LLaMA and	698
Alpaca. <i>ArXiv</i> , abs/2304.08177.	699
Adis Davletov, Nikolay Arefyev, Denis Gordeev, and	700
Alexey Rey. 2021. LIORI at SemEval-2021 task 2:	701
Span prediction and binary classification approaches	702
to word-in-context disambiguation. In <i>Proceedings</i>	703

707	<i>of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 780–786, Online. Association for Computational Linguistics.	765
708		766
709		767
710		768
711	Wietse de Vries and Malvina Nissim. 2021. As good as new: how to successfully recycle English GPT-2 to make models for other languages . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 836–846, Online. Association for Computational Linguistics.	769
712		770
713		771
714		772
715		773
716	Konstantin Dobler and Gerard de Melo. 2023. FOCUS: Effective embedding initialization for monolingual 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.	774
717		775
718		776
719		777
720		778
721		779
722		780
723	C.m. Downey, Terra Blevins, Nora Goldfine, and Shane Steinert-Threlkeld. 2023. Embedding structure matters: Comparing methods to adapt multilingual vocabularies to new languages . In <i>Proceedings of the 3rd Workshop on Multi-lingual Representation Learning (MRL)</i> , pages 268–281, Singapore. Association for Computational Linguistics.	781
724		782
725		783
726		784
727		785
728		786
729		787
730	Takuro Fujii, Koki Shibata, Atsuki Yamaguchi, Terufumi Morishita, and Yasuhiro Sogawa. 2023. How do different tokenizers perform on downstream tasks in scriptio continua languages?: A case study in Japanese . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)</i> , pages 39–49, Toronto, Canada. Association for Computational Linguistics.	788
731		789
732		790
733		791
734		792
735		793
736		794
737		795
738		796
739	Leonidas Gee, Leonardo Rigutini, Marco Ernandes, and Andrea Zugarini. 2023. Multi-word tokenization for sequence compression . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 612–621, Singapore. Association for Computational Linguistics.	797
740		798
741		799
742		800
743		801
744		802
745		803
746	Edward Gow-Smith, Harish Tayyar Madabushi, Carolina Scarton, and Aline Villavicencio. 2022. Improving tokenisation by alternative treatment of spaces . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11430–11443, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	804
747		805
748		806
749		807
750		808
751		809
752		810
753	Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubashir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. XL-sum: Large-scale multilingual abstractive summarization for 44 languages . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4693–4703, Online. Association for Computational Linguistics.	811
754		812
755		813
756		814
757		815
758		816
759		817
760		818
761	Valentin Hofmann, Hinrich Schuetze, and Janet Pierre-humbert. 2022. An embarrassingly simple method to mitigate undesirable properties of pretrained language model tokenizers . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 385–393, Dublin, Ireland. Association for Computational Linguistics.	819
762		820
763		821
764		821

822	Bill Yuchen Lin, Seyeon Lee, Xiaoyang Qiao, and Xiang Ren. 2021. Common sense beyond English: Evaluating and improving multilingual language models for commonsense reasoning. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1274–1287, Online. Association for Computational Linguistics.	881
823		882
824		883
825		884
826		885
827		886
828		887
829		888
830		889
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		
	Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, and Djamé Seddah. 2021. When being unseen from mBERT is just the beginning: Handling new languages with multilingual language models. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 448–462, Online. Association for Computational Linguistics.	881
		882
		883
		884
		885
		886
		887
		888
		889
	OpenAI. 2023. GPT-4 technical report. <i>ArXiv</i> , abs/2303.08774.	890
		891
	Malte Ostendorff and Georg Rehm. 2023. Efficient language model training through cross-lingual and progressive transfer learning. <i>ArXiv</i> , abs/2301.09626.	892
		893
		894
	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems 32</i> , pages 8024–8035. Curran Associates, Inc.	895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
	Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2023. Language model tokenizers introduce unfairness between languages. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	907
		908
		909
		910
		911
	Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. 2021. How good is your tokenizer? on the monolingual performance of multilingual language models. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3118–3135, Online. Association for Computational Linguistics.	912
		913
		914
		915
		916
		917
		918
		919
		920
	Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurencon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa Etxabe, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris C. Emezue, Christopher Klamm, Colin Leong, Daniel Alexander van Strien, David Ifeoluwa Adelani, Dragomir R. Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal	921
		922
		923
		924
		925
		926
		927
	Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.	928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939

940	Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady ElSahar, Hamza Benyamina, Hieu Trung Tran, Ian Yu, Idris Abdumumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jorg Frohberg, Josephine L. Tobiing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro von Werra, Leon Weber, Long Phan, Loubna Ben Allal, Ludovic Tangy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario vSavsko, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad Ali Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto L'opez, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma Sharma, S. Longpre, So maieh Nikpoor, S. Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debjyoti Datta, Eliza Szczeczla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-Shaibani, Matteo Manica, Nihal V. Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Févry, Trishala Neeraj, Urmiish Thakker, Vikas Raunak, Xiang Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Y Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre Francois Lavall'ee, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aur'elie N'ev'eol, Charles Lovering, Daniel H Garrette, Deepak R. Tunuguntla, Ehud Reiter, Ekaterina Taktashova, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Xiangru Tang, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, S. Osher Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Zdeněk Kasner, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ananda Santa Rosa Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajjHosseini, Bahareh Behroozi, Benjamin Ayoade Ajibade, Bharat Kumar Saxena, Carlos Muñoz Ferrandis, Danish Contractor, David M. Lansky, Davis David, Douwe Kiela, Duong Anh Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatim Tahira Mirza, Frankline Ononiwu, Habib Rezanejad, H.A. Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jan Passmore, Joshua Seltzer, Julio Bonis Sanz, Karen Fort, Lívia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolou, Michael McKenna, Mike Qiu, Mohammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nourhan Fahmy, Olanrewaju Samuel, Ran An, R. P. Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas L. Wang, Sourav Roy, Sylvain Viguer, Thanh-Cong Le, Tobi Oyebade, Trieu Nguyen Hai Le, Yoyo Yang, Zachary Kyle Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Kumar Singh, Benjamin Beilharz, Bo Wang, Caio Matheus Fonseca de Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel Le'on Perin'an, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Iman I.B. Bello, Isha Dash, Ji Soo Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthi Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, María Andrea Castillo, Marianna Nezhurina, Mario Sanger, Matthias Samwald, Michael Cullan, Michael Weinberg, M Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patricia Haller, R. Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo L. Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aroonsiri, Srishti Kumar, Stefan Schweter, Sushil Pratap Bharati, Tammay Laud, Th'eo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yashasvi Bajaj, Y. Venkatraman, Yifan Xu, Ying Xu, Yu Xu, Zhee Xao Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. 2022. <i>BLOOM: A 176b-parameter open-access multilingual language model</i> . <i>ArXiv</i> , abs/2211.05100.	1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1053 1054 1055 1056
941	Mike Schuster and Kaisuke Nakajima. 2012. <i>Japanese and Korean voice search</i> . In <i>2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5149–5152.	1057 1058 1059 1060
942	Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2020. <i>MLSUM: The multilingual summarization corpus</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> ,	1061 1062 1063 1064 1065
943		
944		
945		
946		
947		
948		
949		
950		
951		
952		
953		
954		
955		
956		
957		
958		
959		
960		
961		
962		
963		
964		
965		
966		
967		
968		
969		
970		
971		
972		
973		
974		
975		
976		
977		
978		
979		
980		
981		
982		
983		
984		
985		
986		
987		
988		
989		
990		
991		
992		
993		
994		
995		
996		
997		
998		
999		
1000		
1001		
1002		
1003		

1066 1067	pages 8051–8067, Online. Association for Computational Linguistics.	1125 1126 1127
1068 1069 1070 1071 1072 1073 1074	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.	1128 1129 1130 1131 1132 1133
1075 1076 1077 1078 1079 1080	Jimin Sun, Patrick Fernandes, Xinyi Wang, and Graham Neubig. 2023. A multi-dimensional evaluation of tokenizer-free multilingual pretrained models . In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 1725–1735, Dubrovnik, Croatia. Association for Computational Linguistics.	1134 1135 1136 1137 1138 1139
1081 1082 1083 1084 1085 1086 1087 1088 1089	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . 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 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	1140 1141 1142 1143 1144 1145
1090 1091 1092 1093	Cagri Toraman, Eyup Halit Yilmaz, Furkan Şahinuç, and Oguzhan Ozcelik. 2023. Impact of tokenization on language models: An analysis for Turkish . <i>ACM Trans. Asian Low-Resour. Lang. Inf. Process.</i> , 22(4).	1146 1147 1148 1149
1094 1095 1096 1097 1098 1099 1100	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. LLaMA: Open and efficient foundation language models . <i>ArXiv</i> , abs/2302.13971.	1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160 1161
1101 1102 1103 1104 1105 1106 1107 1108 1109 1110 1111 1112 1113 1114 1115 1116 1117 1118 1119 1120 1121 1122 1123 1124	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharar Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. LLaMA 2: Open foundation and fine-tuned chat models . <i>ArXiv</i> , abs/2307.09288.	1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182
Ke M. Tran. 2019. From English to foreign languages: Transferring pre-trained language models . <i>ArXiv</i> , abs/2002.07306.	1125 1126 1127	
Asahi Ushio, Yi Zhou, and Jose Camacho-Collados. 2023. Efficient multilingual language model compression through vocabulary trimming . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 14725–14739, Singapore. Association for Computational Linguistics.	1128 1129 1130 1131 1132 1133	
Zihan Wang, Karthikeyan K, Stephen Mayhew, and Dan Roth. 2020. Extending multilingual BERT to low-resource languages . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2649–2656, Online. Association for Computational Linguistics.	1134 1135 1136 1137 1138 1139	
Barack W. Wanjava, Lilian D. A. Wanzare, Florence Indede, Owen Mconyango, Lawrence Muchemi, and Edward Ombui. 2023. KenSwQuAD—A question answering dataset for swahili low-resource language . <i>ACM Trans. Asian Low-Resour. Lang. Inf. Process.</i> , 22(4).	1140 1141 1142 1143 1144 1145	
Miles Williams and Nikolaos Aletras. 2023. Frustratingly simple memory efficiency for pre-trained language models via dynamic embedding pruning . <i>ArXiv</i> , abs/2309.08708.	1146 1147 1148 1149	
Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrette Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160 1161	
Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with pre-trained byte-to-byte models . <i>Transactions of the Association for Computational Linguistics</i> , 10:291–306.	1162 1163 1164 1165 1166 1167	
Atsuki Yamaguchi, Gaku Morio, Hiroaki Ozaki, and Yasuhiro Sogawa. 2022. Hitachi at SemEval-2022 task 2: On the effectiveness of span-based classification approaches for multilingual idiomticity detection . In <i>Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)</i> , pages 135–144, Seattle, United States. Association for Computational Linguistics.	1168 1169 1170 1171 1172 1173 1174 1175	
Zheng Xin Yong, Hailey Schoelkopf, Niklas Muenninghoff, 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 Vasiliina Nikoulina. 2023. BLOOM+1: Adding language support to BLOOM for zero-shot prompting .	1176 1177 1178 1179 1180 1181 1182	

1183
1184 In *Proceedings of the 61st Annual Meeting of the As-*
1185 *sociation for Computational Linguistics (Volume 1:*
1186 *Long Papers)*, pages 11682–11703, Toronto, Canada.
Association for Computational Linguistics.

1187 Appendix

1188 A Implementation Details

1189 A.1 Tokenizer

1190 To reduce the computational costs, we utilized publicly available existing tokenizers for each target language, which means we used them as \mathcal{T}_t . Table 1191 3 lists the tokenizers used in our experiments.

1194 A.2 Language-specific pre-trained LM

1195 For language-specific pre-trained LMs used in CLP and CLP+, we used the corresponding models to 1196 \mathcal{T}_t , which are listed in Table 3 and are all decoder-based models. Note that we only used the embedding 1197 of each language-specific pre-trained LM for 1198 vocabulary adaptation, and therefore, one can also 1199 use encoder and encoder-decoder based models.

1202 A.3 fastText in FOCUS

1203 For FOCUS, we trained a fastText model for each 1204 language on a corresponding CC-100 (Conneau 1205 et al., 2020) text with the same configuration as 1206 Dobler and de Melo (2023).

1207 A.4 Hyperparameters and Generation 1208 Configurations

1209 **LAPT** Table 4 shows the hyperparameters in 1210 LAPT for each model size. Note that due to the 1211 computational resource constraints and funds for 1212 running experiments, we could run pre-training 1213 of up to four days for each approach. Therefore, 1214 we picked up checkpoints with the largest number 1215 of steps available across models with the same 1216 base model (i.e. BLOOM-1B, BLOOM-7B, etc.) 1217 and language for evaluation to make a fair 1218 comparison. We also temporarily trimmed the unused 1219 embeddings of BLOOM models for LAPT, whose 1220 tokens did not appear in the training corpus during 1221 pre-training to save memory and for faster 1222 computation.⁸

1223 **Generation** Following Cui et al. (2023), we 1224 introduced a verbalizer for the classification tasks: NLI 1225 and MC, where we mapped the first generated token 1226 into a label to compute accuracy. For mapping, we 1227 simply picked up a token with the maximum log- 1228 likelihood among candidate tokenized words. The 1229 list of candidate label words for each task is shown 1230 in Table 5. Table 6 lists the parameters used during 1231 evaluation. To make a fair comparison, we did not

⁸We used the implementations by Ushio et al. (2023) and Williams and Aletras (2023).

conduct any generation parameter tuning and used the same ones across all approaches. For SUM and few-shot SPAN in Swahili, we truncated an article whenever it exceeded the maximum prompt length of 4,096 to avoid the CUDA out-of-memory error.

1237 A.5 Checkpoints

1238 As explained in A.4, we trained all models for up 1239 to four days each due to limited computational re- 1240 sources and funds for experiments. The only excep- 1241 tion was Swahili since the dataset is small enough 1242 to complete LAPT. To make a fair comparison, we 1243 used checkpoints with the largest number of steps 1244 available across models with the same target lan- 1245 guage and base model. Table 7 shows the list of 1246 checkpoints used for evaluation.

Model	Language			
	de	ja	ar	sw
BLOOM-1B	47k	48k	50k	9k
BLOOM-7B	8k	8k	8k	4k
TigerBot-7B	8k	8k	8k	4k
Mistral-7B	6k	6k	6k	4k

Table 7: List of checkpoints used for evaluation. We used checkpoints with the largest number of steps available across all models with the same base model and language.

1247 A.6 Libraries

1248 We implement our models using PyTorch (Paszke 1249 et al., 2019), Hugging Face Transformers (Wolf 1250 et al., 2020) and PEFT (Mangrulkar et al., 1251 2022). We preprocess data with Hugging Face 1252 Datasets (Lhoest et al., 2021). For evaluation, we 1253 use Hugging Face Evaluate⁹ to compute down- 1254 stream performance metrics.

1255 A.7 Prompt Templates

1256 Table 8 shows the prompt templates used in our 1257 evaluation.

1258 A.7.1 Code

1259 The anonymized code is available as supplementary 1260 material in the submission for reference.

1261 B Licenses

1262 This study used various publicly available mod- 1263 els and datasets with different licenses, as detailed 1264 below, all of which permit their use for academic 1265 research.

⁹<https://github.com/huggingface/evaluate>

Language	Tokenization Algorithm	Hugging Face Identifier	Citation	License
German	Byte-level BPE	malteos/gpt2-xl-wechsel-german		MIT
Japanese	Unigram	rinna/japanese-gpt-neox-3.6b-instruction-ppo		MIT
Arabic	Byte-level BPE	aubmindlab/aragpt2-base	(Antoun et al., 2021)	See here
Swahili	Byte-level BPE	benjamin/gpt2-wechsel-swahili	(Minixhofer et al., 2022)	MIT

Table 3: List of tokenizers used for each language-specific model with vocabulary adaptation.

Hyperparameters	1B	7B
Batch size	8	16
Gradient accumulation steps	4	4
Maximum number of training epochs	1	1
Maximum number of training days	4	4
Adam ϵ	1e-8	1e-8
Adam β_1	0.9	0.9
Adam β_2	0.999	0.999
Sequence length	1,024	1,024
Learning rate	1e-4	1e-4
Learning rate scheduler	cosine	cosine
Warmup steps	100	100
Weight decay	0.01	0.01
Attention dropout	0.0	0.0
Dropout	0.05	0.05
LoRA rank r	8	8
LoRA dropout	0.05	0.05
LoRA α	32	32
Training precision	FP16	FP16
Model quantization	int 8	int 8

Table 4: Hyperparameters for LAPT.

Task	Language	Label words
NLI	English	True, False, Neither
	German	Wahr, Falsch, Weder
	Japanese	真, 偽, どちらでもない
	Arabic	لا هنا ولا ذاك, خطأ, صحيح
	Swahili	Kweli, Uongo, Wala
MC	All	A, B, C, D, E

Table 5: List of candidate label words for each classification task.

B.1 Models

BLOOM is licensed under the BigScience RAIL License.¹⁰ TigerBot and Mistral are licensed under the Apache-2.0 License. The licenses of the helper models are listed in Table 3.

B.2 Datasets

XNLI is distributed under CC BY-NC 4.0. JNLI, XQuAD, and JSQuAD are distributed under CC BY-SA 4.0. XCSQA is a derivative of CommonsenseQA (Talmor et al., 2019), which is licensed

¹⁰<https://huggingface.co/spaces/bigscience/license>

Parameters	Values
Maximum prompt length	4,096
Temperature	0.8
Repetition penalty	1.1
Top k	40
Top p	0.9
Beam width	5
Sampling	True
Early stopping	True

Table 6: Parameters for generation.

under an MIT license. OSCAR and KenSwQuAD are licensed under CC0 – no rights reserved. XL-Sum is licensed under CC BY-NC-SA 4.0, while MLSUM is distributed under an MIT license.

C Results

C.1 Additional Results

Table 9 shows the results with standard deviations when prompted in English, and Table 10 shows the results with standard deviations when prompted in a target language.

C.2 English Downstream Performance

Table 11 shows the results on the English datasets. Despite the entire replacement of embeddings for cross-lingual vocabulary adaptation approaches, their adapted models exhibit comparable or better results in most of the tasks for BLOOM, except for SPAN, where Source showed the best result followed by LAPT. This can be ascribed to the following reasons: First, LAPT can retain more source model knowledge than cross-lingual vocabulary adaptation approaches, as their embeddings have not changed. Second, SPAN can be seen as a challenging task as it requires more linguistic understanding of a prompt than simply classifying a text as in NLI and MC. We, therefore, see such a huge performance difference in SPAN since cross-lingual vocabulary adaptation approaches lost more

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

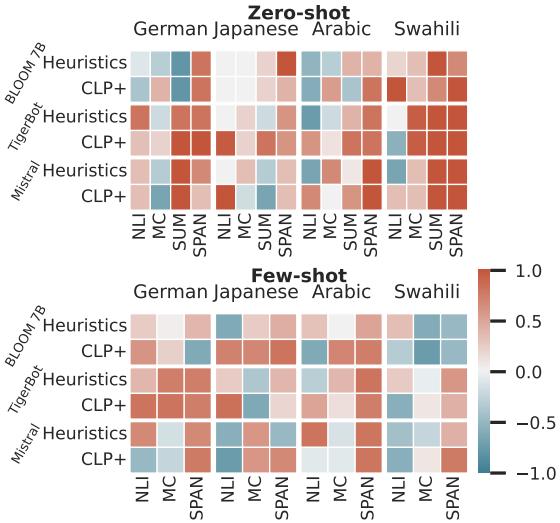


Figure 6: Kendall’s τ correlation between the number of LAPT steps and performance (English prompting).

source linguistic knowledge than LAPT counterparts in exchange for faster inference in a target language.

For TigerBot and Mistral, which are not as multilingual as BLOOM, we see quite similar trends observed in §5.2 in that (1) models with cross-lingual vocabulary adaptation fail to achieve competitive downstream performance to the baselines and (2) their few-shot performances are far lower than LAPT. These results suggest that there can be a relationship between downstream performances in a target language and those in English, and maintaining competitive downstream performance to the baselines in English might be a key to improving models with cross-lingual vocabulary adaptation in terms of their downstream performance.

C.3 How helpful is LAPT for LLMs with cross-lingual vocabulary adaptation?

Figure 6 visualizes Kendall’s tau correlation coefficient between the number of LAPT steps and downstream performance when prompted in English. Similar to Figure 4, we observe that LAPT helped improve downstream performance in both zero-shot and few-shot settings even when prompted in a target language in 65.6% and 63.5% of the cases, respectively.

C.4 Loss Curves

Figures 7 to 10 show the loss curves in LAPT for each model setting.

Approach	Japanese		Swahili	
	English	Target	English	Target
Zero-shot				
LAPT	0.720	0.333	0.788	0.187
+ Heuristics	0.453	0.173	0.173	0.160
+ CLP+	-0.160	-0.106	-0.226	0.066
Few-shot				
LAPT	0.626	1.00	0.453	0.906
+ Heuristics	0.706	0.600	0.591	0.701
+ CLP+	-0.946	0.626	0.886	0.756

Table 12: Kendall’s tau correlation coefficients corresponding to Figure 5. We include LAPT results for reference.

C.5 Kendall’s Tau Correlation Coefficient for Figure 5

Table 12 lists all Kendall’s tau correlation coefficients corresponding to Figure 5 in §6. Models using cross-lingual vocabulary adaptation do not exhibit a strong correlation in the zero-shot setting, ranging from -0.226 to 0.173. The only exception is Heuristics in Japanese with English prompting (0.45). We observe a positive correlation ranging (0.59-0.889) in the few-shot setting, except for CLP+ in Japanese with English prompting (-0.95).

1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342

Task	Language	Template
NLI	English	{premise} Question: {hypothesis} True, False, or Neither? Answer:
	German	{premise} Frage: {hypothesis} Wahr, Falsch oder Weder? Antwort:
	Japanese	{premise} 質問: {hypothesis} 真、偽、どちらでもない? 答え:
	Arabic	{premise} سؤال: {hypothesis} صحيح، خطأ أو لا هذا ذاك؟ إجابة:
	Swahili	{premise} Swali: {hypothesis} Kweli, Uongo au Wala? Jibu:
MC	English	{question} A. {choice_1}, B. {choice_2}, C. {choice_3}, D. {choice_4}, E. {choice_5} Answer:
	German	{question} A. {choice_1}, B. {choice_2}, C. {choice_3}, D. {choice_4}, E. {choice_5} Antwort:
	Japanese	{question} A. {choice_1}, B. {choice_2}, C. {choice_3}, D. {choice_4}, E. {choice_5} 答え:
	Arabic	{question} A. {choice_1}, B. {choice_2}, C. {choice_3}, D. {choice_4}, E. {choice_5} إجابة:
	Swahili	{question} A. {choice_1}, B. {choice_2}, C. {choice_3}, D. {choice_4}, E. {choice_5} Jibu:
SUM	English	Write a short summary of the following text in {language}. Article: {text} Summary:
	German	Schreiben Sie eine kurze Zusammenfassung des folgenden Textes auf Deutsch. Artikel: {text} Zusammenfassung:
	Japanese	次の文章の要約を日本語で書きなさい。記事: {text} 要約:
	Arabic	{text} المختصر المقصود بالنص التالي باللغة العربية. المقالة:
	Swahili	Andika muhtasari mfupi wa maandishi yafuatayo kwa Kiswahili. Makala: {text} Muhtasari:
SPAN	English	Answer the following question. Context: {context} Question: {question} Answer:
	German	Beantworten Sie die folgende Frage. Artikel: {context} Frage: {question} Antwort:
	Japanese	次の文章の質間に答えなさい。文章: {context} 質問: {question} 答え:
	Arabic	الإجابة {question}: السؤال {context}: أجب على السؤال التالي. سياق:
	Swahili	Jibu swali lifuatalo. Makala: {context} Swali: {question} Jibu:

Table 8: Prompt template for each task and language.

Approach	German				Japanese				Arabic				Swahili			
	NLI	MC	SUM	SPAN	NLI	MC	SUM	SPAN	NLI	MC	SUM	SPAN	NLI	MC	SUM	SPAN
BLOOM-1B																
Source	.34.00	.20.00	14.9 _{0.7}	.08.00	.17.00	.25.00	3.7 _{0.1}	.23.00	.35.00	.19.00	9.9 _{0.1}	.17.01	.34.00	.21.00	10.6 _{0.4}	.08.00
LAPT	.33.00	.21.01	17.0 _{0.1}	.09.01	.17.00	.20.00	17.2 _{0.3}	.25.00	.34.00	.18.00	11.6 _{0.1}	.14.01	.34.00	.18.00	9.7 _{0.1}	.08.00
+ Random	.31.00	.22.00	19.0 _{0.2}	.12.00	.17 .00	.20.00	17.7 _{0.2}	.19.00	.36.00	.18.00	10.4 _{0.2}	.09.00	.32.00	.22.00	9.8 _{0.1}	.05.00
+ Heuristics	.32.00	.22.00	17.7 _{0.4}	.11.00	.17 .00	.20.00	17.9 _{0.1}	.25.00	.39.00	.18.00	10.9 _{0.1}	.13.00	.35.00	.22.00	11.9 _{0.1}	.10.00
+ CLP	.34 .00	.23.00	18.7 _{0.2}	.12.01	.17 .00	.21.00	17.6 _{0.1}	.27.01	.36.00	.18.00	11.2 _{0.1}	.14.00	.38.00	.22.00	11.6 _{0.1}	.11.00
+ FOCUS	.34 .00	.22.00	17.4 _{0.4}	.11.00	.17 .00	.21.00	16.5 _{0.1}	.28.00	.36.00	.18.00	11.2 _{0.1}	.13.00	.34.00	.22.00	12.0 _{0.1}	.11.00
+ CLP+	.32.00	.22.00	19.2 _{0.2}	.12.00	.17 .00	.20.00	18.6 _{0.1}	.30.00	.40.01	.18.00	11.2 _{0.1}	.14.00	.38.01	.22.00	11.4 _{0.1}	.10.00
BLOOM-7B																
Source	.33.00	.22.00	6.3 _{0.1}	.15.01	.17.00	.20.00	8.1 _{0.1}	.36.00	.33.00	.18.00	2.8 _{0.1}	.21.00	.34.00	.22.00	7.3 _{0.2}	.21.00
LAPT	.34.00	.19.01	18.7 _{0.5}	.17.01	.17.00	.19.01	18.2 _{0.2}	.36.01	.35.00	.17.00	9.7 _{0.2}	.22.00	.33.00	.20.00	13.2 _{0.1}	.14.00
+ Heuristics	.35 .00	.20.00	17.6 _{0.3}	.24.00	.17 .00	.21 .00	16.7 _{0.1}	.38.00	.33.00	.19 .00	10.5 _{0.1}	.19.00	.36.00	.23.00	12.3 _{0.1}	.16.00
+ CLP+	.35 .00	.21.00	18.0 _{0.3}	.22.00	.17 .00	.20 .00	16.9 _{0.1}	.37.01	.32.00	.21 .00	10.1 _{0.1}	.21.00	.39.01	.22.00	12.0 _{0.2}	.18.00
TigerBot-7B																
Source	.42.00	.22.00	5.4 _{0.2}	.32.03	.29.00	.28.01	1.9 _{0.1}	.51.01	.36.00	.20.00	2.4 _{0.1}	.05.00	.33.00	.22.00	9.0 _{0.2}	.06.00
LAPT	.37.00	.21.01	19.5 _{0.3}	.20.00	.32.01	.21.00	14.8 _{0.2}	.50.00	.41.00	.18.01	9.4 _{0.1}	.12.00	.43.00	.20.01	15.9 _{0.1}	.15.00
+ Heuristics	.35.00	.20.01	18.8 _{0.2}	.17.00	.17.00	.24.00	17.1 _{0.2}	.38.00	.34.00	.19.00	6.6 _{0.1}	.09.00	.31.00	.22.00	8.2 _{0.1}	.04.00
+ CLP+	.35.00	.22.00	20.2 _{0.1}	.19.00	.17.00	.21.00	18.9 _{0.1}	.40.00	.37.00	.23 .00	9.1 _{0.2}	.16.01	.32.00	.22.00	7.7 _{0.1}	.05.00
Mistral-7B																
Source	.36.00	.28.00	8.3 _{0.2}	.35.00	.21.00	.30.00	8.4 _{0.3}	.56.00	.41.00	.24.00	2.4 _{0.1}	.22.00	.34.00	.20.00	5.7 _{0.2}	.12.00
LAPT	.35.01	.25.01	25.2 _{0.3}	.30.00	.19.00	.21.01	23.1 _{0.1}	.48.00	.34.00	.18.00	10.9 _{0.1}	.12.00	.36.00	.22.01	16.5 _{0.1}	.18.00
+ Heuristics	.32.00	.22.00	20.8 _{0.3}	.21.00	.27 .00	.20.00	17.2 _{0.2}	.37.00	.33.00	.20.00	11.2 _{0.2}	.17.00	.30.00	.22.00	11.9 _{0.2}	.12.00
+ CLP+	.36 .00	.22.00	20.2 _{0.2}	.23.00	.29 .00	.23.00	15.9 _{0.2}	.38.00	.38.00	.20.00	10.8 _{0.1}	.24.00	.31.01	.21.00	12.6 _{0.1}	.17.00
BLOOM-1B																
Source	.35.00	.19.00	-	.11.00	.28.00	.18.00	-	.24.00	.34.00	.18.00	-	.21.01	.35.00	.22.00	-	.03.00
LAPT	.34.01	.18.01	-	.14.00	.30.01	.21.01	-	.27.00	.32.01	.17.01	-	.18.00	.36.00	.20.01	-	.02.00
+ Random	.35 .00	.21 .00	-	.16 .00	.28.00	.19.00	-	.24.00	.35 .00	.19 .01	-	.18.00	.34.00	.18.00	-	.03 .00
+ Heuristics	.35 .00	.22 .00	-	.18 .01	.34 .00	.16.00	-	.29 .00	.33.00	.21 .00	-	.19.00	.33.00	.18.00	-	.06 .00
+ CLP	.36 .00	.20 .00	-	.18 .00	.37 .00	.19.00	-	.33 .00	.35 .00	.20 .01	-	.20.00	.35.00	.19.00	-	.07 .00
+ FOCUS	.34.00	.19 .00	-	.19 .01	.54 .00	.21 .00	-	.33 .00	.34 .00	.19 .01	-	.20.00	.33.01	.20.00	-	.07 .00
+ CLP+	.34.00	.19 .00	-	.22.00	.51 .00	.20.00	-	.34.00	.35.01	.20 .01	-	.19.00	.35.01	.18.00	-	.07 .00
BLOOM-7B																
Source	.42.00	.20.01	-	.28.00	.23.00	.18.00	-	.33.00	.43.00	.20.01	-	.30.00	.38.00	.18.00	-	.11.00
LAPT	.34.01	.21.01	-	.28.00	.34.01	.20.01	-	.36.01	.38.01	.18.01	-	.28.00	.38.01	.20.01	-	.09.00
+ Heuristics	.36.00	.20.00	-	.31 .00	.26.00	.21 .00	-	.43 .00	.37.00	.20 .00	-	.32 .00	.33.00	.20 .00	-	.12 .00
+ CLP+	.37.00	.22.00	-	.29.00	.37.01	.20.00	-	.43.00	.32.00	.21 .00	-	.35 .00	.35.00	.20.00	-	.13 .00
TigerBot-7B																
Source	.43.00	.38.01	-	.38.00	.42.01	.33.00	-	.45.00	.39.00	.18.00	-	.13.00	.36.00	.23.00	-	.04.00
LAPT	.46.01	.39.00	-	.37.00	.29.00	.31.00	-	.47.00	.43.01	.19.01	-	.20.01	.44.01	.21.00	-	.15.00
+ Heuristics	.35.00	.27.00	-	.24.01	.47 .00	.24.00	-	.42.00	.33.00	.20 .00	-	.11.00	.36.00	.22.00	-	.02.00
+ CLP+	.36.00	.32.00	-	.34.00	.32.00	.21.00	-	.43.00	.35.00	.20 .00	-	.25 .00	.34.00	.19.00	-	.04.00
Mistral-7B																
Source	.54.00	.55.00	-	.45.00	.49.00	.42.00	-	.56.00	.46.00	.35.00	-	.32.00	.38.00	.21.00	-	.20.00
LAPT	.51.01	.47.00	-	.28.00	.43.01	.37.01	-	.55.00	.44.01	.30.01	-	.22.00	.47.01	.34.01	-	.25.00
+ Heuristics	.40.00	.40.00	-	.24.00	.39.00	.24.00	-	.44.00	.35.00	.18.00	-	.24.00	.33.00	.17.00	-	.13.00
+ CLP+	.37.00	.48.00	-	.30.00	.31.00	.26.00	-	.44.00	.34.00	.23.00	-	.33 .00	.32.00	.19.00	-	.14.00

Table 9: Mean performances over five runs with standard deviations when prompted in English on 500 randomly selected test samples for each dataset. The baselines are in grey. **Bold** indicates comparable or better results than the baselines.

Approach	German				Japanese				Arabic				Swahili			
	NLI	MC	SUM	SPAN												
BLOOM-1B																
Source	.33.00	.21.00	17.8 _{0.3}	.06.00	.29.00	.20.00	18.2 _{0.3}	.22.00	.35.00	.20.00	12.0 _{0.2}	.15.0 ₁	.32.00	.22.00	12.0 _{0.3}	.03.00
LAPT	.36. ₀₁	.22. ₀₁	14.3 _{0.2}	.09.00	.28.00	.20.00	20.7 _{0.2}	.26.00	.36.00	.19. ₀₁	11.4 _{0.1}	.13. ₀₁	.31. ₀₁	.18.00	7.7 _{0.1}	.07.00
+ Random	.35.00	.22.00	15.3 _{0.2}	.14.00	.29.00	.21.00	19.0 _{0.0}	.32.00	.35.00	.19.00	11.5 _{0.0}	.14.0 ₁	.33.00	.22.00	10.2 _{0.1}	.08.01
+ Heuristics	.34.00	.19.00	15.3 _{0.2}	.13.00	.29.00	.19.00	19.2 _{0.0}	.31.00	.36.00	.22.02	11.3 _{0.1}	.13.00	.34.00	.22.00	11.9 _{0.2}	.11.00
+ CLP	.34.00	.18.00	14.6 _{0.7}	.14.00	.29.00	.25.00	18.8 _{0.1}	.33.00	.36.00	.21.00	11.2 _{0.2}	.14.00	.34.00	.22.00	11.5 _{0.2}	.11.00
+ FOCUS	.34.00	.19.00	16.1 _{0.8}	.13.01	.29.00	.21.00	19.2 _{0.0}	.33.00	.35.00	.20.01	11.2 _{0.1}	.14.00	.34.00	.22.00	11.2 _{0.1}	.12.00
+ CLP+	.31.00	.15.00	15.8 _{0.6}	.13.00	.29.00	.19.00	19.4 _{0.0}	.33.00	.35.00	.17.00	11.3 _{0.1}	.15.01	.34.00	.20.00	10.4 _{0.1}	.10.00
BLOOM-7B																
Source	.36.00	.21.00	23.1 _{0.2}	.15.0 ₁	.28.00	.21.00	19.0 _{0.2}	.33.00	.38.00	.17.00	11.5 _{0.1}	.25.00	.36.0 ₁	.22.00	14.3 _{0.1}	.22.0 ₁
LAPT	.37. ₀₁	.21. ₀₁	19.4 _{0.3}	.14.00	.21.00	.21.00	21.6 _{0.1}	.36.00	.36.00	.16.00	11.5 _{0.2}	.21.00	.36. ₀₁	.20.00	13.0 _{0.1}	.14.00
+ Heuristics	.32.00	.22.00	19.7 _{0.3}	.21.00	.30.00	.23.00	19.5 _{0.1}	.38.00	.37.00	.19.00	10.7 _{0.2}	.21.00	.33.00	.22.00	11.6 _{0.1}	.16.00
+ CLP+	.35.00	.21.00	18.7 _{0.7}	.20.01	.29.00	.21.00	19.5 _{0.1}	.40.00	.38.00	.21.00	11.0 _{0.1}	.21.00	.33.00	.23.00	10.9 _{0.1}	.17.00
TigerBot-7B																
Source	.33.00	.24.00	23.9 _{0.2}	.26.01	.17.00	.24.00	19.4 _{0.3}	.57.0 ₁	.33.00	.21.00	9.0 _{0.1}	.04.00	.29.0 ₁	.22.00	12.4 _{0.2}	.03.00
LAPT	.32.00	.21.00	18.5 _{0.2}	.18.00	.17.00	.21.01	21.6 _{0.1}	.49.0 ₁	.33.00	.18.00	9.8 _{0.2}	.13.00	.31.00	.21.01	15.9 _{0.1}	.10.00
+ Heuristics	.35.00	.20.00	16.1 _{0.3}	.18.00	.29.00	.22.00	19.6 _{0.1}	.40.00	.36.00	.17.00	10.3.03	.08.00	.36.00	.22.00	8.1 _{0.1}	.05.00
+ CLP+	.33.00	.20.01	14.1 _{0.3}	.19.00	.29.00	.20.00	19.8 _{0.0}	.41.00	.38.00	.22.00	11.2.03	.16.00	.30.01	.22.00	8.6 _{0.1}	.09.00
Mistral-7B																
Source	.34.00	.25.00	24.1 _{0.2}	.35.01	.17.00	.28.00	23.7 _{0.1}	.60.00	.33.00	.20.00	11.2 _{0.1}	.21.00	.35.00	.22.00	15.4 _{0.1}	.07.00
LAPT	.33. ₀₁	.25.02	24.2 _{0.2}	.28.01	.17.00	.20.01	23.4 _{0.1}	.60.00	.33.00	.18.00	10.8 _{0.1}	.14.0 ₁	.33. ₀₁	.22.01	16.2 _{0.2}	.12.00
+ Heuristics	.40.00	.26.00	21.2 _{0.1}	.22.00	.29.00	.20.00	19.7 _{0.1}	.43.00	.39.00	.19.00	10.7 _{0.1}	.13.00	.34.00	.22.00	10.6 _{0.1}	.14.00
+ CLP+	.35.00	.25.00	20.2 _{0.3}	.21.00	.28.00	.20.00	19.9 _{0.0}	.46.00	.38.00	.16.00	11.5.02	.21.00	.33.00	.21.00	10.2 _{0.1}	.16.00
BLOOM-1B																
Source	.38.00	.20.00	-	.10.00	.44.00	.19.00	-	.32.00	.35.00	.17.00	-	.20.0 ₁	.37.00	.23.00	-	.02.00
LAPT	.37. ₀₁	.17. ₀₁	-	.13.01	.26.01	.21.01	-	.34.0 ₁	.34.00	.16.00	-	.16.00	.34.0 ₁	.19.0 ₁	-	.02.00
+ Random	.34.00	.21.00	-	.16.00	.29.00	.21.00	-	.34.00	.35.00	.22.00	-	.16.01	.34.00	.20.00	-	.06.00
+ Heuristics	.33.00	.23.00	-	.17.00	.30.00	.22.00	-	.32.00	.36.00	.21.01	-	.15.01	.33.00	.19.00	-	.07.00
+ CLP	.34.00	.21.00	-	.17.01	.30.00	.20.00	-	.33.00	.35.00	.21.00	-	.15.01	.33.00	.19.00	-	.08.00
+ FOCUS	.34.00	.18.00	-	.17.01	.27.00	.20.00	-	.36.00	.36.00	.20.00	-	.15.00	.34.00	.19.00	-	.08.00
+ CLP+	.34.00	.20.00	-	.19.00	.29.00	.22.00	-	.36.00	.37.01	.20.01	-	.15.01	.30.00	.18.00	-	.08.00
BLOOM-7B																
Source	.38.00	.23.00	-	.29.00	.41.0 ₁	.19.00	-	.49.0 ₁	.37.00	.18.00	-	.29.00	.34.00	.18.00	-	.11.00
LAPT	.35.00	.24.00	-	.23.00	.34.0 ₁	.19.0 ₁	-	.53.0 ₁	.36.0 ₁	.18.0 ₁	-	.23.00	.37.00	.18.0 ₁	-	.07.00
+ Heuristics	.33.00	.22.00	-	.28.00	.28.00	.21.00	-	.46.00	.36.00	.21.00	-	.24.00	.36.00	.19.00	-	.13.00
+ CLP+	.34.00	.22.00	-	.25.00	.30.00	.20.00	-	.46.00	.36.00	.22.00	-	.25.00	.36.00	.18.00	-	.13.00
TigerBot-7B																
Source	.31.00	.37.00	-	.42.00	.16.00	.34.00	-	.65.00	.30.00	.19.00	-	.10.00	.36.00	.19.0 ₁	-	.03.00
LAPT	.30.00	.39.01	-	.36.02	.16.00	.34.01	-	.66.00	.30.00	.20.0 ₁	-	.17.00	.36.00	.20.00	-	.09.00
+ Heuristics	.33.00	.26.01	-	.21.00	.29.00	.24.00	-	.49.00	.35.00	.20.00	-	.09.00	.35.00	.21.00	-	.04.00
+ CLP+	.36.01	.31.00	-	.31.00	.30.00	.21.00	-	.50.00	.37.01	.19.00	-	.19.00	.34.00	.18.00	-	.06.00
Mistral-7B																
Source	.33.00	.53.00	-	.48.00	.16.00	.42.00	-	.69.00	.30.00	.32.00	-	.31.00	.40.00	.21.00	-	.12.00
LAPT	.33.00	.46.00	-	.27.00	.16.00	.37.01	-	.68.00	.30.00	.30.01	-	.26.00	.36.01	.34.00	-	.21.00
+ Heuristics	.45.00	.41.00	-	.24.00	.30.00	.24.00	-	.49.00	.34.00	.18.00	-	.17.00	.34.00	.18.00	-	.09.00
+ CLP+	.37.00	.47.00	-	.25.00	.29.00	.25.00	-	.50.00	.38.00	.23.00	-	.23.00	.33.00	.20.00	-	.14.00

Table 10: Mean performances over five runs with standard deviations when prompted in a target language on 500 randomly selected test samples for each dataset. The baselines are in grey. **Bold** indicates comparable or better results than the baselines.

Approach	NLI				MC				SUM				SPAN			
	de	ja	ar	sw	de	ja	ar	sw	de	ja	ar	sw	de	ja	ar	sw
BLOOM-1B																
Zero-shot																
Source		.34.00				.18.00				11.2 _{0.1}				.21.0 ₁		
LAPT	.35. ₀₁	.33.00	.33.00	.33.00	.21. ₀₁	.20.00	.17. ₀₁	.18. ₀₁	9.7 _{0.1}	10.4 _{0.0}	10.6 _{0.1}	10.1 _{0.0}	.15. ₀₁	.18.00	.17.00	.15. ₀₁
+ Heuristics	.35.00	.34.00	.35.00	.36.00	.20.00	.20.00	.21.00	.20.00	11.2 _{0.1}	11.7 _{0.1}	12.3 _{0.1}	10.1 _{0.2}	.10. ₀₁	.06. ₀₁	.07.00	.09. ₀₁
+ CLP+	.35.00	.34.00	.32.00	.37.01	.19.00	.20.00	.20.00	.20.00	10.8 _{0.1}	12.1 _{0.1}	12.4 _{0.1}	10.4 _{0.2}	.12.00	.11.00	.07.00	.08.00
BLOOM-7B																
Source		.36.00				.17.00				11.1 _{0.1}				.31.00		
LAPT	.34.00	.34.00	.34.00	.36.00	.20. ₀₁	.20.01	.20.01	.18. ₀₁	10.8.00	11.0.0	10.9.0	10.6.0	.25.00	.27.01	.28.00	.23.00
+ Heuristics	.36.00	.34.00	.33.00	.36.00	.20.00	.20.00	.20.00	.19.00	11.2 _{0.1}	11.1 _{0.1}	12.2 _{0.1}	10.5 _{0.0}	.24.01	.19.00	.17.00	.17.00
+ CLP+	.36.00	.34.00	.33.00	.38.00	.20.00	.20.00	.20.00	.20.00	10.5 _{0.1}	11.8 _{0.1}	12.4 _{0.0}	11.1 _{0.1}	.21.00	.10.01	.22.00	.19.00
TigerBot-7B																
Source		.48.00				.29.00				12.7 _{0.1}				.42.0 ₁		
LAPT	.39. ₀₀	.49. ₀₁	.45. ₀₁	.45. ₀₁	.23.00	.25.00	.25.00	.28. ₀₁	11.6 _{0.1}	11.9 _{0.1}	12.2 _{0.1}	11.0 _{0.1}	.31.00	.34.01	.27.00	.35. ₀₁
+ Heuristics	.37.00	.34.00	.35.00	.31.00	.21.00	.24.00	.20.00	.20.00	10.2 _{0.1}	12.1 _{0.1}	9.8 _{0.1}	4.8 _{0.1}	.15.00	.24.01	.03.00	.02.00
+ CLP+	.39.00	.36.00	.36.00	.31.00	.24.00	.24.00	.21.00	.20.00	10.4 _{0.1}	12.5 _{0.1}	11.3 _{0.1}	6.3 _{0.1}	.20.00	.30.00	.20.00	.03.00
Mistral-7B																
Source		.42.00				.46.00				12.4 _{0.2}				.44.00		
LAPT	.36. ₀₁	.49. ₀₁	.45. ₀₁	.42.00	.34.01	.32.01	.28.01	.38. ₀₁	11.6.0	11.3.0	8.1.0	10.6.0	.39.00	.40.01	.28.00	.36.00
+ Heuristics	.33.00	.39.00	.36.00	.33.00	.21.00	.21.00	.20.00	.20.00	12.5 _{0.1}	12.9 _{0.1}	11.5 _{0.2}	8.5 _{0.2}	.23.00	.30.00	.19.00	.06.00
+ CLP+	.36.00	.37.00	.34.00	.31.01	.21.00	.25.00	.19.00	.22.00	12.2 _{0.1}	13.1 _{0.1}	12.9 _{0.1}	10.6 _{0.1}	.26.00	.27.01	.31.00	.18.00
BLOOM-1B																
Few-shot																
Source		.33.00				.20.00				-	-	-		.28.00		
LAPT	.34. ₀₁	.31.00	.32.01	.33.00	.17. ₀₁	.18.01	.19.01	.20. ₀₁	-	-	-	-	.20.01	.23.01	.25.00	.20. ₀₁
LoRA + Heuristics	.33.00	.34.00	.32.00	.34.00	.19.00	.20.00	.21.00	.19.00	-	-	-	-	.16.00	.11.00	.14.00	.17. ₀₁
LoRA + CLP+	.37.00	.34.00	.35.00	.35.00	.19.00	.21.00	.22.00	.20.00	-	-	-	-	.17.00	.10.01	.13.00	.18.00
BLOOM-7B																
Source		.43.00				.21.00				-	-	-		.39.00		
LAPT	.36. ₀₁	.38. ₀₁	.40.00	.39. ₀₁	.20.01	.21.01	.21.00	.19.00	-	-	-	-	.36.00	.38.00	.38.00	.37.00
+ Heuristics	.36.00	.35.00	.31.00	.32.00	.20.00	.19.00	.22.00	.21.00	-	-	-	-	.37.00	.33.00	.36.00	.34.00
+ CLP+	.36.00	.33.00	.33.00	.34.00	.18.00	.21.00	.20.00	.20.00	-	-	-	-	.35.00	.34.00	.36.03	.36.00
TigerBot-7B																
Source		.49.00				.58.00				-	-	-		.47.00		
LAPT	.47. ₀₁	.56. ₀₁	.45.00	.56. ₀₁	.57.00	.58.00	.52.00	.52. ₀₁	-	-	-	-	.47.00	.46.00	.44.00	.47.00
+ Heuristics	.36.00	.35.00	.37.00	.34.00	.20.00	.31.00	.23.00	.19.00	-	-	-	-	.24.00	.41.01	.03.00	.01.00
+ CLP+	.42.00	.34.00	.37.00	.36.00	.32.00	.24.00	.19.00	.19.00	-	-	-	-	.37.00	.43.00	.29.00	.02.00
Mistral-7B																
Source		.60.00				.66.00				-	-	-		.51.00		
LAPT	.55.00	.53.01	.49. ₀₁	.56.01	.62.00	.59.01	.57.01	.63.01	-	-	-	-	.38.00	.49.00	.31.00	.49.00
+ Heuristics	.43.00	.35.00	.36.00	.33.00	.33.00	.31.00	.20.00	.19.00	-	-	-	-	.27.00	.47.00	.32.00	.05.00
+ CLP+	.38.00	.34.00	.38.00	.37.00	.45.00	.31.00	.26.00	.21.00	-	-	-	-	.34.00	.45.00	.44.00	.28.01

Table 11: Mean performances over five runs with standard deviations on 500 randomly selected test samples for each English dataset. The baselines are in grey. **Bold** indicates comparable or better results than the baselines.

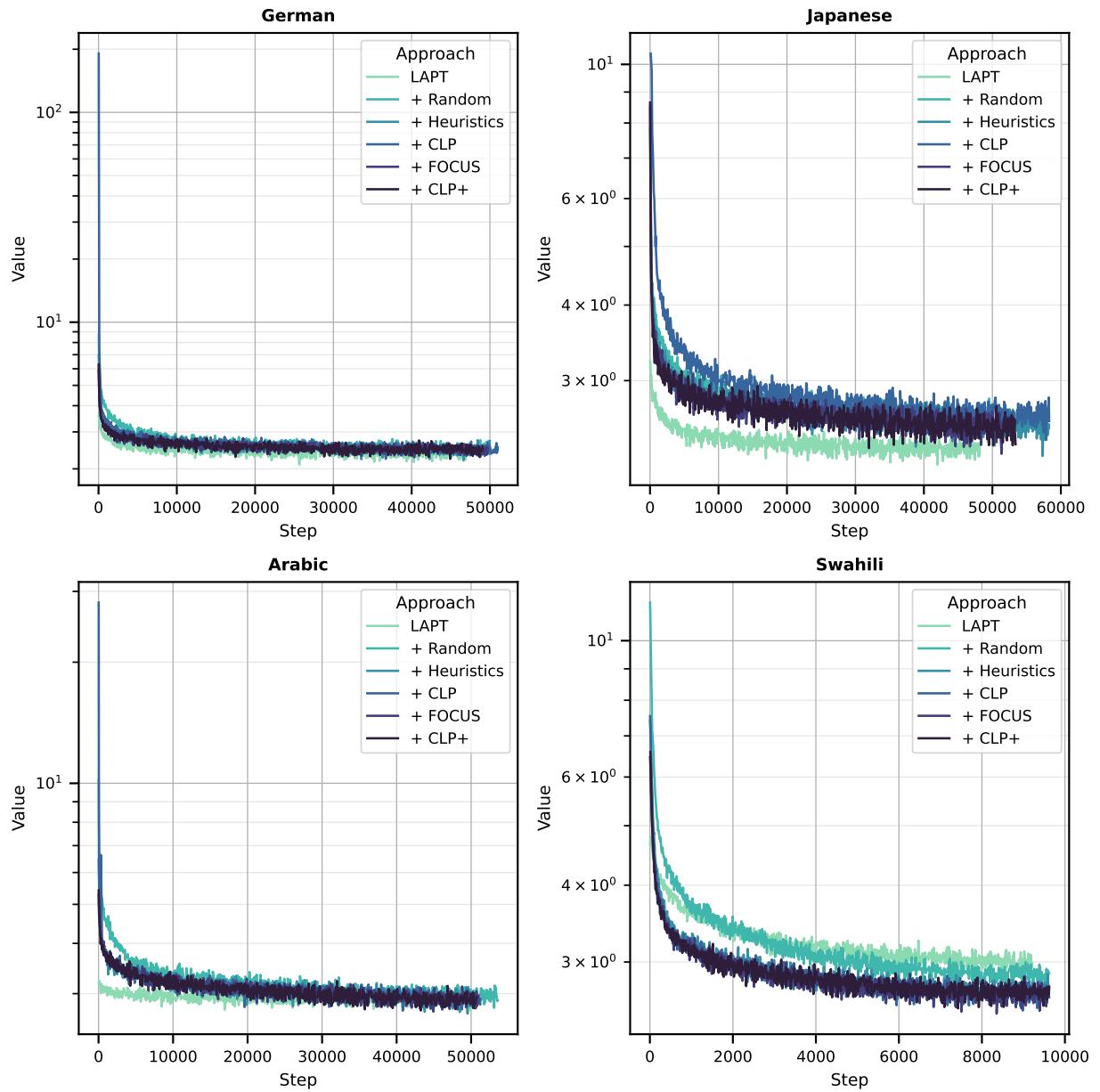


Figure 7: LAPT loss curves for BLOOM-1B

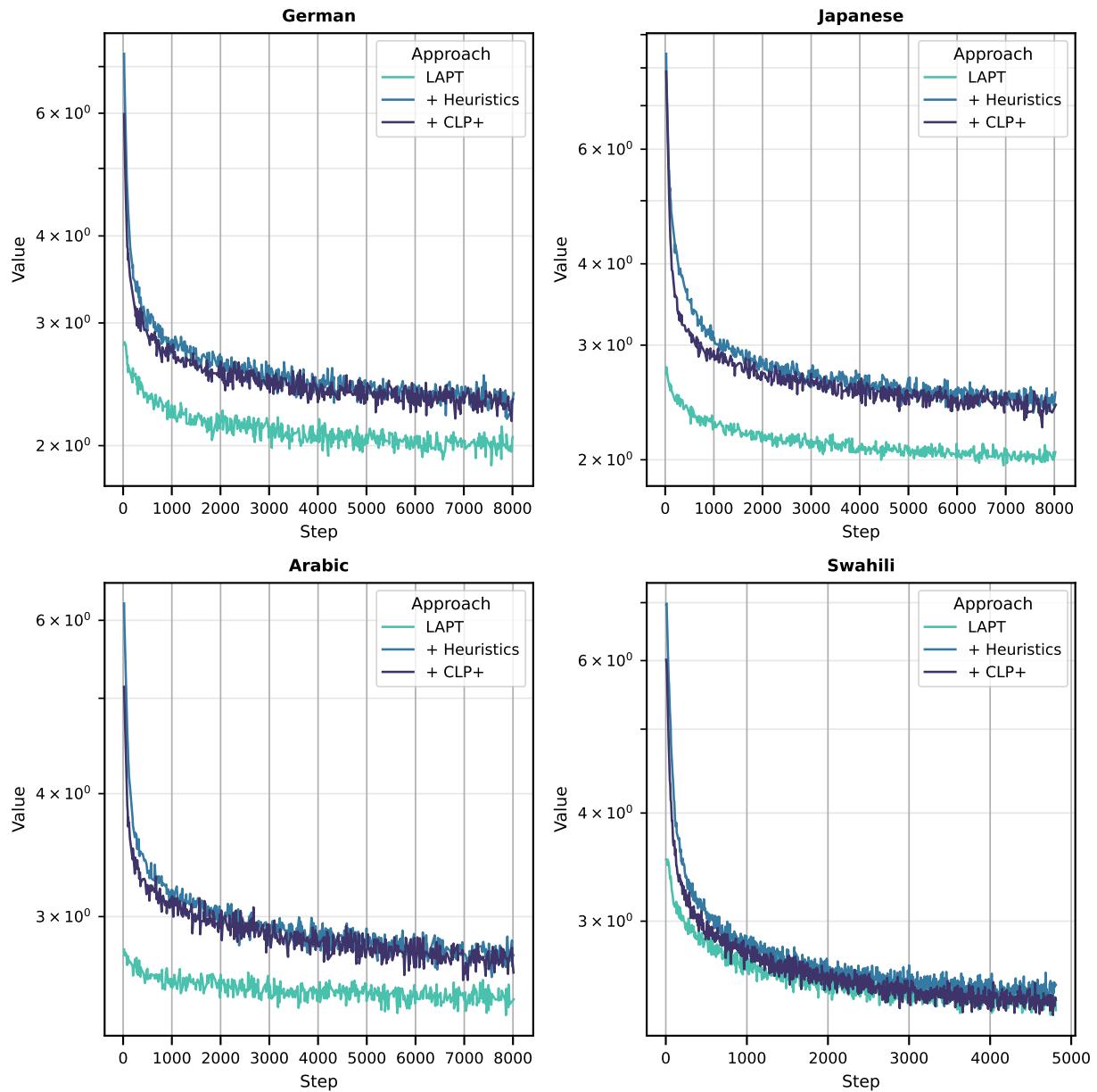


Figure 8: LAPT loss curves for BLOOM-7B

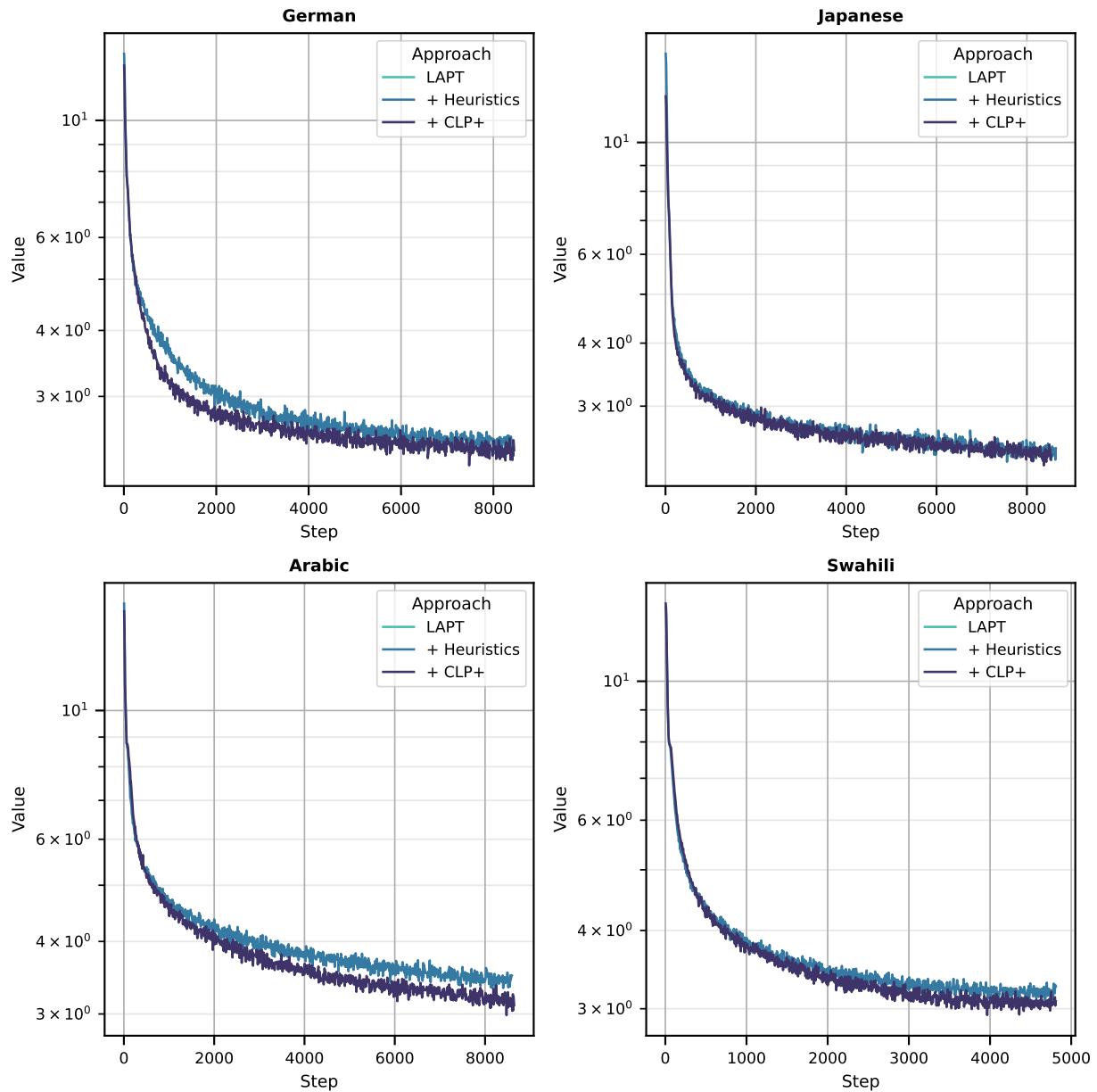


Figure 9: LAPT loss curves for TigerBot-7B

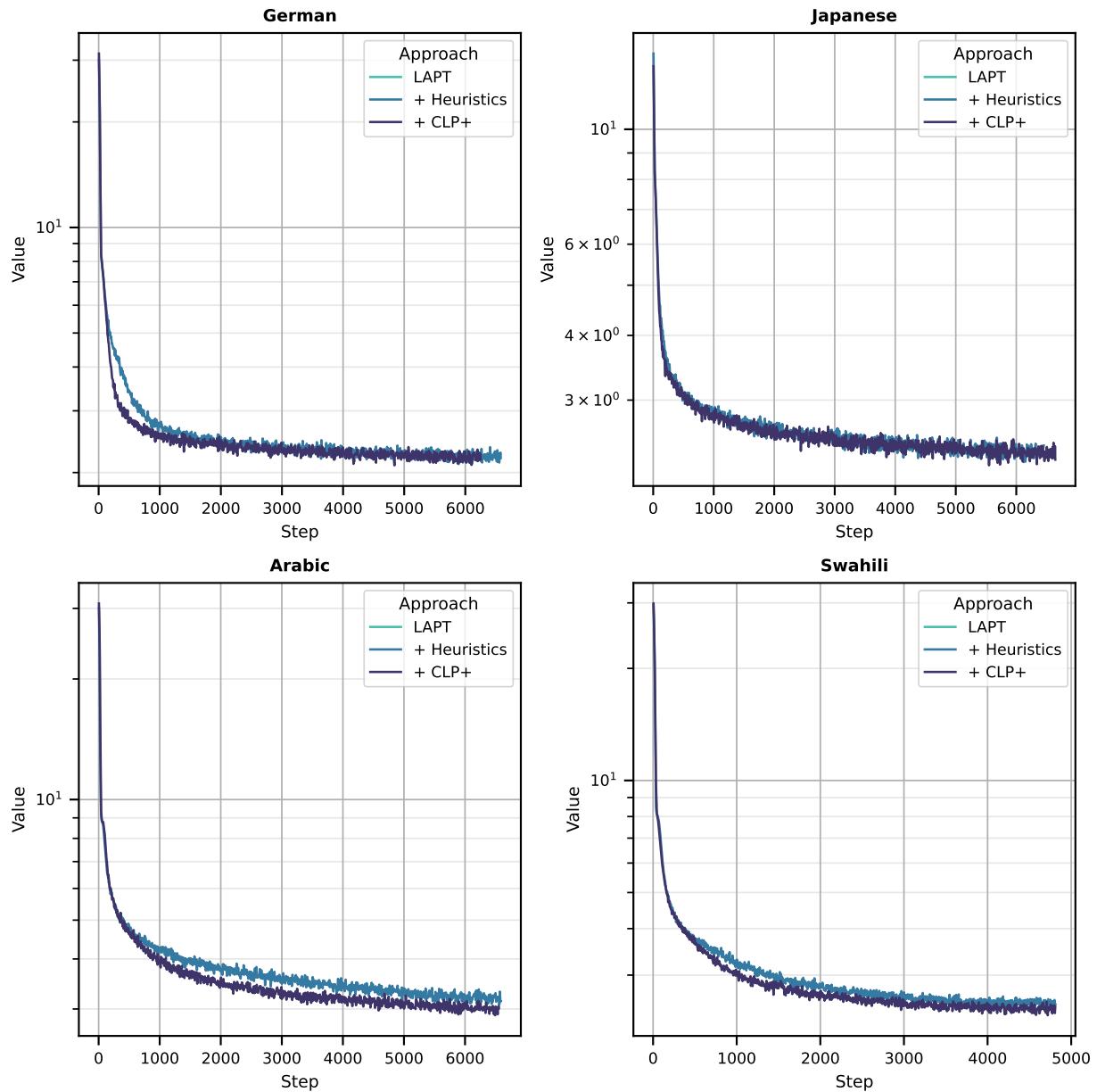


Figure 10: LAPT loss curves for Mistral-7B