

---

# HyperOFA: Expanding LLM Vocabulary to New Languages via Hypernetwork Based Embedding Initialization

---

Enes Özeren

Department of Statistics, LMU Munich  
enes.oezeren@campus.lmu.de

## Abstract

Most pre-trained language models are trained primarily for high-resource languages, limiting their usability in mid and low-resource languages. A common approach to adapt these models for other languages involves introducing new tokens specific to these languages and continuing pre-training. However, the method used to initialize these newly introduced tokens significantly impacts the duration and efficiency of continued pre-training. Poor initialization can lead to longer training times and increased computational costs. OFA [Liu et al., 2023a] method provides an effective initialization strategy. Building on this, HyperOFA introduces a hypernetwork-based approach to initialize new tokens. Experimental results show that HyperOFA consistently outperforms the random initialization baseline and performs on par with OFA, sometimes surpassing it and sometimes falling behind.

## 1 Introduction

Current language models (LMs) require pre-training on vast amounts of text data. As a result, most LMs are trained on high-resource languages, such as English and Chinese, which have a lot of available linguistic data. High-resource languages are those with extensive digital text corpora, while mid- and low-resource languages lack sufficient data for effective model training. Consequently, many LMs are either monolingual or multilingual, but for limited high-resource languages. This limits the usability and accessibility of LMs for speakers of mid or low-resource languages, hindering their potential benefits and widespread adoption. Training a language model from scratch for such languages is costly and often impractical due to limited text data and the large number of model parameters—hence the term Large Language Models (LLMs).

Most LMs consist of two primary components: transformer block parameters, which capture general language reasoning capabilities, and embedding parameters, which encode token representations for specific languages. Extending pre-trained LMs (PLMs) to support additional languages often requires modifying and initializing the embedding parameters for the new languages. Therefore, some research focuses on extending the language capabilities of PLMs by initializing new token embeddings from new languages. Some of these approaches aim to transfer the capabilities of a monolingual model to another language, resulting in a new monolingual model [de Vries and Nissim, 2020, Minixhofer et al., 2021], while others explore expanding existing monolingual or multilingual models to include a broader range of languages [Artetxe et al., 2019, Tran, 2020, Dobler and De Melo, 2023, Liu et al., 2023a].

To initialize the token embeddings from new languages, methods like WECHSEL [Minixhofer et al., 2021] and OFA [Liu et al., 2023a] utilize external information sources such as static multilingual word vectors. These methods identify similar token embeddings from LM’s supported languages and generate new token embeddings by combining the similar embeddings linearly. However, combining

LM’s supported language tokens can be limiting, especially when introducing new tokens for a new language with very different semantic meanings. Additionally, the use of linear combinations may restrict the ability of these methods.

To address these challenges, this study introduces HyperOFA. Like previous methods, HyperOFA utilizes static multilingual word vectors as an external information source. However, it trains a hypernetwork to learn the relationship between the word vectors and the LM’s existing token embeddings. During initialization, the hypernetwork generates new token embeddings by taking external word vectors as input. To empirically test the method, the vocabulary size of RoBERTa (50K tokens) ([Liu, 2019]) and XLM-RoBERTa (250K tokens) ([Conneau, 2019]) models were extended to 401K tokens. The results were compared against a random initialization baseline and the OFA method. The HyperOFA code-base is available at <https://github.com/enesozeren/hyper-ofa>.

## 2 Related Work

**Tokenization** Language models typically use embedding vectors that encode the semantic information of basic units of text, known as tokens. Tokens can be character level, word level, or more commonly, sub-word level. The process of generating tokens from given text is performed by a tokenizer which is trained on a dataset similar to LM’s pre-training dataset but a smaller one ([Radford et al., 2019]). The token set is referred as the LM’s vocabulary. This vocabulary consists of units from the languages in the smaller dataset, potentially, limiting the LMs ability to learn new languages through continued pre-training. To address this, Glot500-m ([Imani et al., 2023]) expands the vocabulary of XLM-RoBERTa ([Conneau, 2019]) from 250K tokens to 401K tokens, performs continued pre-training and demonstrates superior performance on both high- and low-resource languages.

**Introducing New Tokens** To improve the ability of LLMs to handle languages beyond those included in their original pre-training dataset, recent research has proposed various methods. One promising approach is to introduce new tokens for the new languages and then performing continued pre-training. These new tokens’ embeddings are initialized in different ways in the literature, which this paper categorizes as either random or wise initialization.

**Random Token Initialization** The simplest approach for initializing new token embeddings is random initialization. de Vries and Nissim [2020] and Artetxe et al. [2019] randomly initialize the new token embeddings and perform continued pre-training by freezing model weights and training only the embedding layer. Similarly, ([Imani et al., 2023]) employ random initialization but perform continued pre-training across all model parameters. However, random initialization requires more computational resources and training data to achieve good performance compared to wise initialization methods ([Liu et al., 2023a]).

**Wise Token Initialization** Wise initialization methods generally uses a external source of information to initialize the new embeddings better than random. An early work, Tran [2020] created a bilingual LM by introducing new tokens from a new language by using a parallel or non-parallel corpus. The method learns a sparse word translation matrix to initialize the new token embeddings.

The WECHSEL method ([Minixhofer et al., 2021]) proposes a wise initialization technique by using an external source of information to transform a monolingual model to another monolingual model in a different language. The method utilizes multilingual static word vectors which contain both the source and target language. First, the method matches the words in the external source to the tokens of source and target language if they are contained in the word. Then, using this matching information, both the source and target tokens are projected into the external word vector space. By considering the cosine similarity of the tokens in this shared space, the target token embeddings are initialized as the convex combination of similar source token embeddings. The paper applies this method for nine different languages.

Liu et al. [2023a] introduced the OFA framework that builds upon the token initialization method from the WECHSEL approach. However, OFA expands the vocabulary of both monolingual and multilingual models to cover significantly larger number of languages. To prevent the embedding matrix from becoming disproportionately large compared to the rest of the model, OFA employs a factorization technique which was proposed in ALBERT ([Lan, 2019]). The method leverages

→  
ColexNet+ word vectors [Liu et al., 2023b] as external word vectors and utilizes Glot500-m tokenizer for introducing the new tokens.

Similar to WECHSEL and OFA methods, HyperOFA also utilizes an external source of information, specifically ColexNet+ word vectors, as used in OFA. However, the key distinction of HyperOFA lies in its use of a hypernetwork to directly predict embeddings for new token initialization, whereas the previous two methods initialize new tokens by combining embeddings of existing ones. This approach can be particularly advantageous for initializing tokens from different languages which have very different semantics than the model’s original languages. This paper directly compares HyperOFA with OFA method under experiment settings similar with those in the OFA paper to evaluate the impact of this different approach.

**Hypernetworks** Hypernetworks are networks designed to generate the weights of another network. A recent survey by Chauhan et al. [2024] highlights their application across various domains of deep learning. One of the earlier works in initializing embeddings with hypernetworks is MIMICK ([Pinter et al., 2017]) which focuses on out-of-vocabulary (OOV) words for word vectors. The method uses an RNN architecture for the hypernetwork and predicts the OOV word embeddings. On the other hand, Schick and Schütze [2019] integrates a hypernetwork into BERT ([Devlin, 2018]) to generate embeddings for rare words. More recently, Minixhofer et al. [2024] proposed a hypernetwork-based method for zero-shot tokenizer transfer, enabling a language model to detach from its tokenizer and work with any tokenizer.

### 3 Methodology

The HyperOFA method builds upon certain aspects of the OFA method. To provide a clear understanding, similar notation used in OFA paper is used. The HyperOFA pipeline can be seen in the Figure 1.

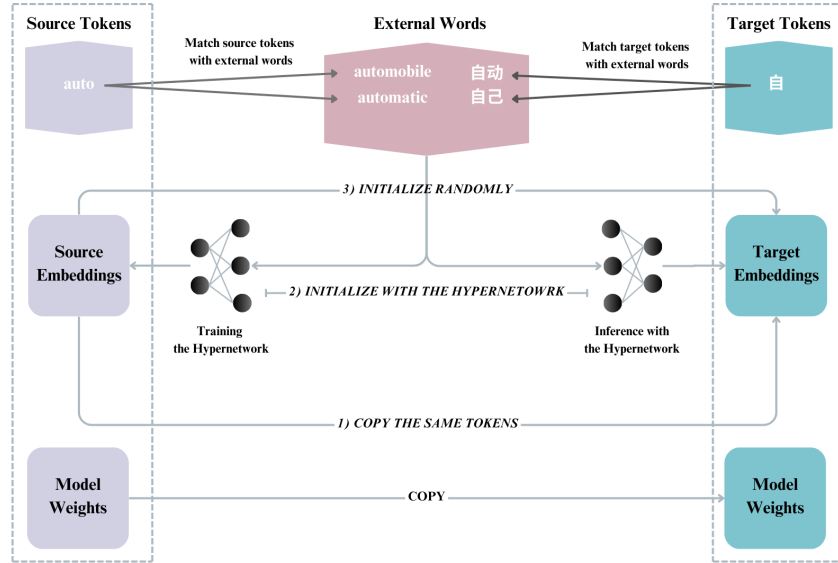


Figure 1: HyperOFA pipeline. The source model (left) transfers weights to the target model (right). The target embeddings are initialized by first copying embeddings for matching tokens, then generating embeddings via a hypernetwork for tokens with matching external words, and finally randomly initializing the rest.

#### 3.1 Problem Setting

Given a model with a source tokenizer  $TOK^s$ , the goal is to replace the source tokenizer with a target tokenizer  $TOK^t$  which supports broader range of tokens across various languages. The source vocabulary size  $V^s$  of  $TOK^s$  is smaller than target vocabulary size  $V^t$  of  $TOK^t$ . Starting from

the source embeddings  $E^s \in \mathbb{R}^{|V^s| \times D}$ , the HyperOFA method aims to initialize target embeddings  $E^t \in \mathbb{R}^{|V^t| \times D}$  wisely by keeping the same embedding dimension  $D$ .

### 3.2 Factorization

Since  $V^t > V^s$  the number of embedding parameters grows significantly from  $V^s \times D$  to  $V^t \times D$  in the target model. This can result having large ratio of model parameters in the embedding matrix. To address this, Lan [2019] and Liu et al. [2023a] use factorization method.

Factorization decomposes the  $E^s$  into two smaller matrices using the Singular Value Decomposition (SVD) method, such that  $E^s \approx F^s P$  where  $F^s \in \mathbb{R}^{|V^s| \times D'}$  is the coordinate matrix containing token-specific parameters, and  $P \in \mathbb{R}^{D' \times D}$  is the primitive embedding matrix capturing language-agnostic features. When  $D' < D$ , the total number of parameters of  $F^s$  and  $P$  is smaller than  $E^s$ . After factorizing  $E^s$ , the goal is to initialize the token-specific embeddings in  $F^t \in \mathbb{R}^{|V^t| \times D'}$  for the  $TOK^t$  while reusing the same  $P$ .

Liu et al. [2023a] experiments with different  $D'$  values to explore the trade-off between parameter saving and model performance.

### 3.3 External Word Vectors

As an external source of information for wise initialization,  $\xrightarrow{\text{CoxNet+}}$  ([Liu et al., 2023b]), static word vectors  $W \in \mathbb{R}^{4\text{million} \times 200}$ , is used in this method.  $\xrightarrow{\text{CoxNet+}}$  contains over four million words spanning more than one thousand languages. To utilize these words vectors, the tokens from  $TOK^s$  and  $TOK^t$  are matched with the words in  $W$ . Specifically, a token will be matched with a word if that word contains the token (see Figure 1). The resulting set of matched word vectors will be represented as  $W(s_i)$  for token  $i$  from  $TOK^s$  and  $W(t_j)$  for token  $j$  from  $TOK^t$ .

### 3.4 Hypernetwork

After performing factorization and creating the set of matched words and tokens, a hypernetwork  $HN_\theta$  with parameters  $\theta$  will be used. The aim of this hypernetwork is to generate the token embedding  $F_j$  by using the matched word vectors  $W(s_j)$  where  $j \in V^t$ . The training set for  $HN_\theta$  is  $(W(s_i), F_i)$  in which the inputs are the matched word vectors and the target is the token embedding in the source coordinate matrix  $F^s$ . Since the number of matched word vectors varies for each token, at least one vector is always present, though multiple word vectors are typically used to predict the target token embedding.

During training, a custom loss function is employed to ensure that the predicted token embeddings not only align semantically with the target embeddings but also maintain a similar norm. This loss combines contrastive loss, which enforces semantic similarity, and normalized L1 loss, which preserves magnitude consistency.

$$\mathcal{L}(\theta) = \lambda \cdot \mathcal{L}_{\text{contrastive}} + (1 - \lambda) \cdot \mathcal{L}_{\text{L1}}$$

$$\text{where } \mathcal{L}_{\text{contrastive}} = -\frac{1}{N} \sum_{k=1}^N \log \frac{\exp(\text{sim}(F_k^s, \hat{F}_k^s)/T)}{\sum_{l=1}^N \exp(\text{sim}(F_l^s, \hat{F}_l^s)/T)} \text{ and } \mathcal{L}_{\text{L1}} = \frac{1}{N} \sum_{k=1}^N \frac{\|F_k^s - HN_\theta(W(s_k))\|_1}{-\frac{1}{N} \sum_{l=1}^N \|F_l^s\|}$$

Here the  $N$  is batch size,  $\text{sim}$  is cosine similarity,  $T$  is temperature and  $\lambda$  is a parameter to balance loss components. The  $\mathcal{L}_{\text{L1}}$  component is a normalized by average embedding value to prevent very small loss values.

When designing the model architecture for  $HN_\theta$ , the input requirements played a crucial role. First, the number of matched word vectors varies for across tokens, meaning the  $HN_\theta$  architecture must be capable of handling variable-length inputs. Secondly, since the order of the input matched word vectors should not influence the prediction, the model architecture should be permutation-invariant. Initially, an encoder only transformer model ([Vaswani, 2017] without positional encoding layers (called as Setformer in this study) that satisfies both requirements was tested. However, after observing poor performance, the approach shifted to a BiLSTM architecture ([Schuster and Paliwal,

1997]) despite it not inherently satisfying the permutation-invariance requirement. Experimental results demonstrated that BiLSTM works better for this task when compared to a transformer encoder model without positional encoding layer (see Appendix Table 4 for experiments). To address the BiLSTM’s sensitivity to input order, data augmentation was implemented by randomly shuffling the word vectors’ order during each training epoch, effectively preventing the model from overfitting to specific sequence arrangements.

The trained BiLSTM-based hypernetwork  $HN_\theta$  generates target embeddings by taking matched word vectors  $W(t_j)$  as inputs, with the network parameters optimized using the specified loss function.

### 3.5 New Token Initialization

The target embeddings,  $F^t$ , are initialized in three steps similar to the OFA method (see Figure 1).

1. For tokens that overlap between  $TOK^s$  and  $TOK^t$ , the token embeddings in  $F^s$  are directly copied to  $F^t$ .
2. For tokens with matched words in  $W$ , their embeddings are predicted by  $HN_\theta$  using the matched word vectors  $W(t_j)$  as input.
3. For the remaining tokens, the embeddings are randomly initialized from a normal distribution with the same mean and standard deviation as the source embeddings.

## 4 Hypernetwork Training

In this study, the vocabularies of RoBERTa-base (RoBERTa) [Liu, 2019] (50K tokens) and XLM-RoBERTa-base (XLM-R) [Conneau, 2019] (250K tokens) are extended to match the Glot500-m vocabulary [Imani et al., 2023], which contains 401K tokens.

Before initializing new token embeddings, the factorization method described in Section 3.2 is applied for different dimensions  $D'$ . While the original embedding size for both models is 768, factorization is used to create coordinate matrices  $F^s$  with dimensions of 100, 200, and 400. The models with reduced embedding sizes are referred to as RoBERTa-100, RoBERTa-200, RoBERTa-400, and similarly, XLM-R-100, XLM-R-200, XLM-R-400.

To compare the performance of HyperOFA with OFA, different hypernetworks are trained to initialize the new token embeddings for these RoBERTa-xxx and XLM-R-xxx models.

**Hypernetwork Training Dataset** For the RoBERTa model, the hypernetwork training dataset consists of 22K word vector and target embedding pairs, as only 22K out of RoBERTa’s 50K vocabulary tokens match with the ColexNet+ words. Similarly, for XLM-R, the training dataset contains 103K pairs, corresponding to the 103K matched tokens from its 250K vocabulary. To mitigate overfitting, data augmentation is applied by shuffling word vectors before each epoch. Additionally, with 50% probability, the number of word vectors is randomly limited to 50–100% of the available vectors. A custom sampler is also employed to ensure that batches contain input sequences of similar length, improving training stability.

**Hypernetwork Training Hyperparameters** As described in Section 3.4, the hypernetworks follow a BiLSTM architecture. All RoBERTa-xxx and XLM-R-xxx hypernetworks share the same configuration: a maximum context size of 256, a dropout rate of 0.4, and an Adam optimizer. The learning rate starts at  $1 \times 10^{-4}$  and decays linearly by a factor of 0.95 every 10 epochs. Training was conducted on two Nvidia A100 GPUs, with each model requiring approximately 1 to 1.5 hours.

To ensure healthy training, the hyperparameters in the loss function, as explained in Section 3.4, were set as follows:  $\lambda = 0.1$  for all hypernetworks, and  $T = 0.5$  for the hypernetworks of RoBERTa-xxx, and  $T = 0.25$  for the hypernetworks of XLM-R-xxx.

As the dimension of the predicted embedding increased, a hypernetwork with higher capacity was necessary. To address this, the hidden dimension of the BiLSTM was increased for embeddings with 400 dimensions. Further details are provided in Table 1. All models were trained until the validation loss converged.

Notably, the hypernetworks have a substantial number of parameters compared to the LLMs they initialize. Experiments showed that larger hypernetworks, when combined with strong regularization (dropout of 0.4 and the data augmentation methods mentioned earlier), were performing better than smaller hypernetworks. For a detailed comparison, see Appendix Figure 3.

LM	Hypernetwork	Training Data	Layers	Hid Dim	Param	Epoch
RoBERTa-100	HN-R-100	22K	2	800	22.1M	370
RoBERTa-200	HN-R-200	22K	2	800	22.5M	470
RoBERTa-400	HN-R-400	22K	2	1600	87.4M	400
XLM-R-100	HN-X-100	103K	4	800	52.9M	120
XLM-R-200	HN-X-200	103K	4	800	52.3M	230
XLM-R-400	HN-X-400	103K	4	1600	210.3M	80

Table 1: Hypernetwork model details for predicting the new tokens for RoBERTa and XLM-R language models with different factorized dimensions. Epochs column indicated the converged epoch number for the hypernetwork.

## 5 Results

With the hypernetworks explained in Section 4, RoBERTa and XLM-R models’ embeddings are initialized with the HyperOFA method. Also those 2 LMs are initialized with OFA and random initialization methods as well. The models are described as follows:

**HyperOFA-mono-xxx** These are RoBERTa models with an extended vocabulary, expanding from 50K tokens to 401K tokens. The "xxx" denotes the embedding dimension of the model, and the "mono" suffix indicates that the RoBERTa model is originally monolingual.

**HyperOFA-multi-xxx** These are XLM-R models with an extended vocabulary, expanding from 250K tokens to 401K tokens. The "xxx" denotes the embedding dimension of the model, and the "multi" suffix indicates that the XLM-R model is originally multilingual.

**OFA-mono-xxx** RoBERTa models with extened vocabulary with the OFA method by using Liu et al. [2023a] code-base.

**OFA-multi-xxx** XLM-R models with extened vocabulary with the OFA method by using Liu et al. [2023a] code-base.

**Random-mono-xxx** RoBERTa models with extened vocabulary with the random initialization method. In this approach all overlapping tokens are copied, while the remaining tokens are randomly initialized with old embedding’s mean and standard deviations are used.

**Random-multi-xxx** XLM-R models with extened vocabulary with the random initialization method. In this approach all overlapping tokens are copied, while the remaining tokens are randomly initialized with old embedding’s mean and standard deviations are used.

The details of the volume of tokens initialized with each step explained in Section 3.5 can be found in Table 2.

### 5.1 Evaluations

The performance of language models (LMs) with extended vocabularies is compared across three approaches: Random, OFA, and HyperOFA. Evaluations are conducted on four benchmarks under two conditions: without continued pre-training and with continued pre-training.

**Sentence Retrieval** Retrieval performance is assessed using the Sentence Retrieval Tatoeba (SR-T) [Artetxe and Schwenk, 2019] and Sentence Retrieval Bible (SR-B) benchmarks. Following Liu et al.

Method	Model	Copy	HyperOFA	OFA	Random	Total
HyperOFA	RoBERTa	27K	179K	0	195K	401K
	XLM-R	255K	84K	0	62K	401K
OFA	RoBERTa	27K	0	179K	195K	401K
	XLM-R	255K	0	84K	62K	401K
Random	RoBERTa	27K	0	0	374K	401K
	XLM-R	255K	0	0	146K	401K

Table 2: Distribution of tokens initialized using HyperOFA, OFA, and random initialization methods. The "Copy" column represents the number of tokens shared between the old and new vocabularies, which are directly copied. This distribution holds for all variants with different embedding factorization dimensions (100, 200, 400).

[2023a], Top-10 accuracy is used as the evaluation metric, where the correct translation must be among the ten nearest neighbors of an English sentence. Sentence representations are obtained by averaging contextualized word embeddings from the model’s 8th layer.

**Sequence Labeling** For sequence labeling, named entity recognition (NER) and part-of-speech tagging (POS) are evaluated using the WikiANN [Pan et al., 2017] and Universal Dependencies [De Marneffe et al., 2021] datasets, respectively. The methodology follows Liu et al. [2023a], where models are fine-tuned on the English training set. The best checkpoint, selected based on validation performance, is then used to report zero-shot performance on test sets in other languages. F1 scores are reported for both benchmarks.

#### 5.1.1 Performance before continued pre-training

The models’ embedding matrices are extended, and the newly added tokens are initialized using three different methods: Random, OFA, and HyperOFA. Test metrics for the benchmarks are then computed without any continued pre-training. The benchmark results are presented in Table 3.

The evaluation of vocabulary extension methods is conducted across all available languages in the benchmarks. SR-B covers 98 languages, SR-T covers 369 languages, NER covers 164 languages, and POS covers 91 languages.

The results without continued pre-training in Table 3 show that HyperOFA and OFA methods outperforms the Random initialization method in all cases. For the retrieval tasks (SR-B and SR-T), HyperOFA performs better than OFA on all cases except mono-400 models. This is because increasing the embedding dimension (from 100 to 200 and 400) leads to a decline in HyperOFA’s performance. With a fixed amount of training data, learning higher-dimensional embeddings becomes more challenging for the hypernetwork. In mono models, where the hypernetwork had only 20K training observations, HyperOFA-mono-400 underperforms compared to OFA-mono-400 across all four benchmarks. However, HyperOFA-multi-400 performs competitively with OFA-multi-400, as the hypernetwork had access to 100K training observations in the multi setting. For sequence labeling tasks (NER and POS), OFA generally outperforms HyperOFA, though HyperOFA occasionally achieves better results.

#### 5.1.2 Performance after continued pre-training

For continued pre-training, a subset of languages from Glot500-c ([Imani et al., 2023]), comprising 22 languages spanning high, mid, and low-resource categories is used. The full list of languages and their data volumes can be found in Appendix - Table 5. This dataset contains 1.1 billion tokens across 36 million sentences.

Continued pre-training is crucial because, even with carefully initialized embeddings, they must be fine-tuned on data from the new language tokens. Additionally, the model parameters need to adapt to these languages. To achieve this, six models were continued pre-training: RoBERTa-100 and XLM-R-400, each initialized with three different methods—HyperOFA, OFA, and Random initialization.

LM	SR-B	SR-T	NER	POS
Random-mono-100	3.5	4.6	23.4	22.5
OFA-mono-100	4.5	6.2	<b>25.0</b>	<b>23.5</b>
HyperOFA-mono-100	<b>5.0</b>	<b>6.4</b>	24.3	22.5
Random-mono-200	3.7	5.2	24.9	23.1
OFA-mono-200	4.5	7.2	<b>25.7</b>	<b>23.4</b>
HyperOFA-mono-200	<b>4.8</b>	<b>7.5</b>	25.3	<b>23.4</b>
Random-mono-400	4.1	5.3	25.8	23.0
OFA-mono-400	<b>4.8</b>	<b>7.2</b>	<b>26.1</b>	<b>24.5</b>
HyperOFA-mono-400	4.7	6.3	25.8	23.0
Random-multi-100	5.1	7.2	34.7	41.5
OFA-multi-100	5.1	7.5	36.3	<b>42.3</b>
HyperOFA-multi-100	<b>5.2</b>	<b>7.6</b>	<b>37.6</b>	<b>42.3</b>
Random-multi-200	5.7	10.0	38.1	47.3
OFA-multi-200	5.7	10.4	<b>40.2</b>	<b>48.6</b>
HyperOFA-multi-200	<b>6.0</b>	<b>10.6</b>	38.3	48.3
Random-multi-400	5.6	21.0	41.4	53.7
OFA-multi-400	5.9	<b>21.3</b>	<b>43.3</b>	<b>54.6</b>
HyperOFA-multi-400	<b>6.1</b>	<b>21.3</b>	42.6	53.8

Table 3: Performance comparison of OFA and HyperOFA in without continued pre-training setting. The values for OFA models are taken from Liu et al. [2023a] directly. SR-B covers 98 languages, SR-T covers 369 languages, NER covers 164 languages, and POS covers 91 languages. For SR-B and SR-T top 10 accuracy, for NER and POS f1 score is presented. All metrics are average across languages.

All six models were continued pre-trained using hyperparameters similar to those in Liu et al. [2023a], with three key differences: smaller dataset, an effective batch size of 512 instead of 384 and training on four Nvidia H100 GPUs. The continued pre-training was conducted for 4,000 steps (approx. 1 epoch) and required at most 1.5 hours to complete.

The results of training loss and retrieval benchmark evaluations are presented in Figure 2. The first observation is that all XLM-R-based models outperform the RoBERTa-based models. This aligns with expectations, as XLM-R is a multilingual model, whereas RoBERTa is English only. Hence, XLM-R is better at adaptation to new languages. Secondly, models with embeddings initialized using the HyperOFA and OFA methods achieve superior performance compared to those with randomly initialized embeddings. Finally, the performance trends between HyperOFA and OFA remain consistent across all cases, with neither method showing a decisive advantage overall. While OFA performs better on the SR-T benchmark, HyperOFA achieves slightly higher scores on the SR-B benchmark. This suggests that both approaches are effective, with their relative strengths depending on the specific evaluation metric.

## 6 Limitations

In this study, as well as in Liu et al. [2023a], HyperOFA and OFA were explored for initializing new embeddings in encoder-only models. While both methods are theoretically applicable to decoder-only models like GPT [Radford et al., 2019] and encoder-decoder models like T5 [Raffel et al., 2020], their effectiveness in these settings remains untested, presenting an open research direction.

Another limitation concerns the embedding dimensions used in this study. Due to the embedding matrix factorization described in Section 3.2, the dimensions are relatively low compared to those in modern LLMs. While this approach reduces computational costs, it leaves open the question of how HyperOFA would perform with much higher-dimensional embeddings.

Finally, the continued pre-training dataset used in this study was relatively small compared to that of Liu et al. [2023a] due to disk quota limitations in this study. Exploring the impact of larger continued pre-training datasets, especially those having multiple languages, could provide deeper insights into the strengths and weaknesses of both methods in different settings.



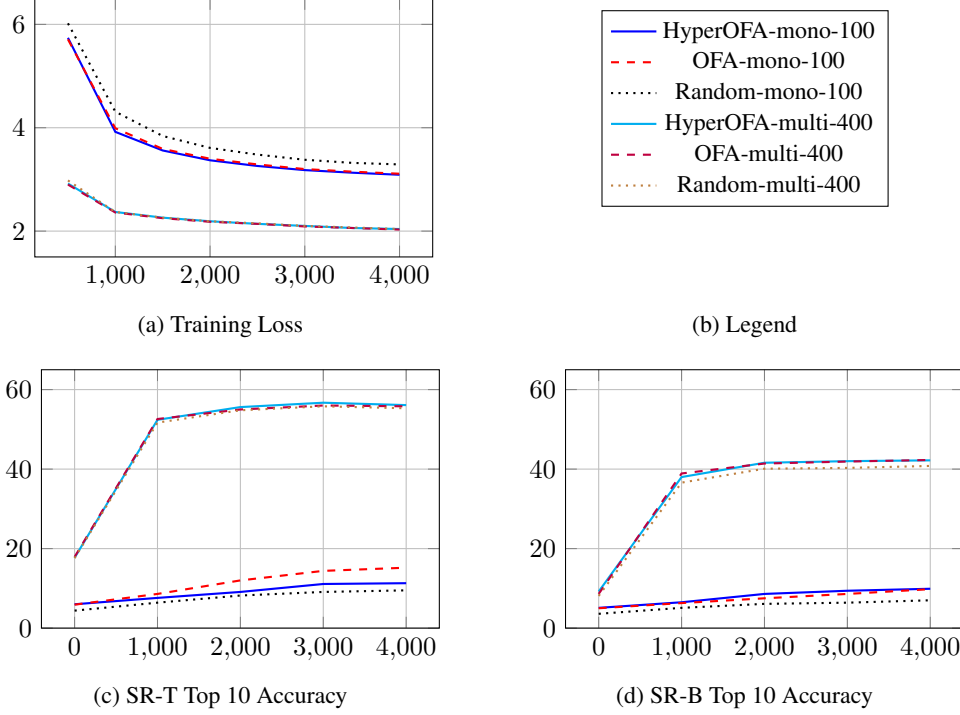


Figure 2: Training losses and evaluation metrics for models with initialized embeddings with HyperOFA, OFA, and random initialization during continued pre-training. (a) shows the training loss values, (b) provides the legend for all three plots, (c) presents SR-T average Top-10 accuracy for languages included in the continued pre-training set (20 languages out of 22 continued pre-training languages), and (d) shows SR-B average Top-10 accuracy for same languages as in SR-T.

## 7 Conclusion

This study introduced HyperOFA as a method for expanding the vocabulary of LLMs to new languages using a hypernetwork model. HyperOFA leverages external static multilingual word vectors as an external information source to initialize new token embeddings. Experiments across four different benchmarks demonstrated that HyperOFA consistently outperforms the random initialization baseline and performs competitively with the OFA method. These results highlight HyperOFA as a promising approach, alongside OFA, for efficient new token embedding initialization.

## Acknowledgements

I sincerely thank Mina Rezaei for generously providing access to computational resources and constructive feedback. I’m also deeply grateful to Yihong Liu for the thought-provoking discussions, invaluable insights, and assistance with the datasets.

## References

- Mikel Artetxe and Holger Schwenk. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the association for computational linguistics*, 7: 597–610, 2019.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolingual representations. *arXiv preprint arXiv:1910.11856*, 2019.
- Vinod Kumar Chauhan, Jiandong Zhou, Ping Lu, Soheila Molaei, and David A Clifton. A brief review of hypernetworks in deep learning. *Artificial Intelligence Review*, 57(9):250, 2024.

- A Conneau. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*, 2019.
- Marie-Catherine De Marneffe, Christopher D Manning, Joakim Nivre, and Daniel Zeman. Universal dependencies. *Computational linguistics*, 47(2):255–308, 2021.
- Wietse de Vries and Malvina Nissim. As good as new. how to successfully recycle english gpt-2 to make models for other languages. *arXiv preprint arXiv:2012.05628*, 2020.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Konstantin Dobler and Gerard De Melo. Focus: Effective embedding initialization for monolingual specialization of multilingual models. *arXiv preprint arXiv:2305.14481*, 2023.
- Ayyoob Imani, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André Martins, François Yvon, and Hinrich Schütze. Glot500: Scaling multilingual corpora and language models to 500 languages. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1082–1117, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.61. URL <https://aclanthology.org/2023.acl-long.61/>.
- Zhenzhong Lan. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schütze. Ofa: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining. *arXiv preprint arXiv:2311.08849*, 2023a.
- Yihong Liu, Haotian Ye, Leonie Weissweiler, and Hinrich Schuetze. Transfer learning for low-resource languages based on multilingual colexification graphs. *arxiv*, 2023b. URL <https://arxiv.org/abs/2305.12818>.
- Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364, 2019.
- Benjamin Minixhofer, Fabian Paischer, and Navid Rekabsaz. Wechsel: Effective initialization of subword embeddings for cross-lingual transfer of monolingual language models. *arXiv preprint arXiv:2112.06598*, 2021.
- Benjamin Minixhofer, Edoardo Maria Ponti, and Ivan Vulić. Zero-shot tokenizer transfer. *arXiv preprint arXiv:2405.07883*, 2024.
- Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. Cross-lingual name tagging and linking for 282 languages. In *Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: long papers)*, pages 1946–1958, 2017.
- Yuval Pinter, Robert Guthrie, and Jacob Eisenstein. Mimicking word embeddings using subword rnns. *arXiv preprint arXiv:1707.06961*, 2017.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Timo Schick and Hinrich Schütze. Bertram: Improved word embeddings have big impact on contextualized model performance. *arXiv preprint arXiv:1910.07181*, 2019.
- Mike Schuster and Kuldip K Paliwal. Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing*, 45(11):2673–2681, 1997.

Ke Tran. From english to foreign languages: Transferring pre-trained language models. *arXiv preprint arXiv:2002.07306*, 2020.

A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.

## A Experiments for Hypernetwork

### A.1 Large vs Small Hypernetwork

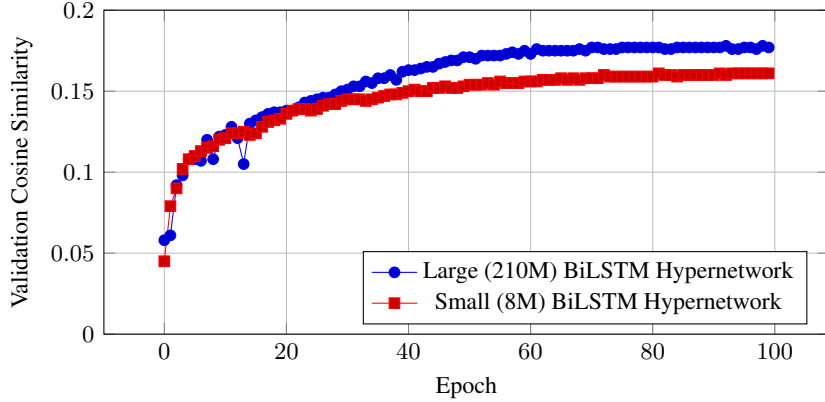


Figure 3: Comparison of large (210M parameters) and small (8M parameters) BiLSTM-based hypernetwork models in terms of validation cosine similarity between predicted and true embeddings over 100 epochs for the XLM-R embeddings.

### A.2 BiLSTM vs Setformer

Table 4 compares two candidate hypernetwork architectures for initializing token embeddings in the RoBERTa model with factorized embeddings (100 dimensional). Setformer (a Transformer encoder without positional encodings) and a BiLSTM (Bidirectional LSTM) are evaluated with two different parameter sizes.

Results show that LLMs initialized with smaller hypernetworks (8M parameters) perform similarly and better than a randomly initialized baseline. However, increasing the hypernetwork size to 22M parameters leads to a notable performance boost. Specifically, the LLM initialized with the larger BiLSTM hypernetwork achieves the highest SR-T Top 10 accuracy (6.4), outperforming both Setformer variants. This suggests that BiLSTM is more effective than Setformer as a hypernetwork, especially when given sufficient capacity.

LM	Hypernetwork	Hypernetwork Param.	LLM SR-T Top 10 Acc
HyperOFA-mono-100	Setformer	8M	5.7
HyperOFA-mono-100	BiLSTM	7M	5.6
HyperOFA-mono-100	Setformer	22M	5.2
HyperOFA-mono-100	BiLSTM	22M	<b>6.4</b>
Random-mono-100	Random Init	-	4.6

Table 4: Comparison of Setformer (Transformer encoder without positional encodings) and BiLSTM as hypernetworks for initializing token embeddings in HyperOFA-mono-100, a RoBERTa-based model with a new vocabulary and factorized embedding dimension of 100. The SR-T Top 10 Accuracy is reported without continued pre-training. Also random initialization baseline performance is given at the last row.

## B Continued Pre-training Dataset

The continued pre-training dataset was deliberately kept smaller than that used by Liu et al. [2023a] due to disk quota limitations in the HyperOFA study. The languages, their original sentence counts in Glot500-c [Imani et al., 2023] dataset and the sentence counts used in this study is listed in Table 5. For continued pre-training 36M sentences (approx. 1.1B tokens) across 22 languages are used. To categorize source category, thresholds used: high (>5M sentences), mid (>500K sentences), and low (<500K sentences).

Source Category	Language	Glott500-c Sentence Count	Subsampled Sentence Count
High	eng_Latn	36,121,560	5,000,000
	tur_Latn	29,182,577	5,000,000
	ell_Grek	22,031,905	5,000,000
	bul_Cyrl	21,822,051	5,000,000
	ces_Latn	20,374,860	5,000,000
	kor_Hang	6,348,091	5,000,000
Mid	kat_Geor	990,785	990,785
	fry_Latn	925,801	925,801
	zsm_Latn	849,033	849,033
	khm_Khmr	565,794	565,794
	jpn_Japn	507,538	507,538
Low	yue_Hani	483,750	483,750
	tuk_Latn	312,480	312,480
	uig_Arab	298,694	298,694
	pam_Latn	292,293	292,293
	kab_Latn	166,953	166,953
	gla_Latn	124,953	124,953
	mhr_Cyrl	91,557	91,557
	swh_Latn	43,876	43,876
	cmn_Hani	57,500	57,500
	pes_Arab	18,762	18,762
	dtp_Latn	1,355	1,355
<i>Total Sentence Count</i>		141,612,168	35,731,124

Table 5: Distribution of continued pre-training data across. The table shows the original Glot500-c volume and sub sampled volume for each language, grouped by their source category (High, Mid, Low).

## C Benchmark Language Coverage

### C.1 For Benchmark Performances in Table 3

SR-B Benchmark Languages:

mal\_Mlym, aze\_Latn, guj\_Gujr, ben\_Beng, kan\_Knda, tel\_Telu, mlt\_Latn, fra\_Latn, spa\_Latn, fil\_Latn, nob\_Latn, rus\_Cyrl, deu\_Latn, tur\_Latn, pan\_Guru, mar\_Deva, por\_Latn, nld\_Latn, zho\_Hani, ita\_Latn, ind\_Latn, ell\_Grek, bul\_Cyrl, swe\_Latn, ces\_Latn, isl\_Latn, pol\_Latn, ron\_Latn, dan\_Latn, hun\_Latn, tgk\_Cyrl, srp\_Latn, fas\_Arab, ceb\_Latn, heb\_Hebr, hrv\_Latn, fin\_Latn, slv\_Latn, vie\_Latn, mkd\_Cyrl, slk\_Latn, nor\_Latn, est\_Latn, ltz\_Latn, eus\_Latn, lit\_Latn, kaz\_Cyrl, lav\_Latn, epo\_Latn, cat\_Latn, tha\_Thai, ukr\_Cyrl, tgl\_Latn, sin\_Sinh, gle\_Latn, hin\_Deva, kor\_Hang, ory\_Orya, urd\_Arab, sqi\_Latn, bel\_Cyrl, afr\_Latn, nno\_Latn, tat\_Cyrl, hau\_Latn, sna\_Latn, msa\_Latn, som\_Latn, srp\_Cyrl, mlg\_Latn, zul\_Latn, arz\_Arab, nya\_Latn, tam\_Taml, hat\_Latn, uzb\_Latn, sot\_Latn, uzb\_Cyrl, als\_Latn, amh\_Ethi, sun\_Latn, war\_Latn, yor\_Latn, fao\_Latn, uzn\_Cyrl, smo\_Latn, bak\_Cyrl, ilo\_Latn, tso\_Latn, mri\_Latn, asm\_Beng, hil\_Latn, nso\_Latn, ibo\_Latn, kin\_Latn, hye\_Armn, lin\_Latn, tpi\_Latn, twi\_Latn, kir\_Cyrl, pap\_Latn, nep\_Deva, bcl\_Latn, xho\_Latn, cym\_Latn, gaa\_Latn, ton\_Latn, lat\_Latn, srn\_Latn, ewe\_Latn, bem\_Latn, efi\_Latn, bis\_Latn, haw\_Latn, hmo\_Latn, kat\_Geor, pag\_Latn, loz\_Latn, fry\_Latn, mya\_Mymr, nds\_Latn, run\_Latn, rar\_Latn, fij\_Latn, ckb\_Arab, ven\_Latn, zsm\_Latn, chv\_Cyrl, sag\_Latn, guw\_Latn, bre\_Latn, toi\_Latn, che\_Cyrl, pis\_Latn, oss\_Cyrl, nan\_Latn, tuk\_Latn, tir\_Ethi, yua\_Latn, min\_Latn, khm\_Khmr, tum\_Latn, lug\_Latn, tzo\_Latn, mah\_Latn, jav\_Latn, jpn\_Jpan, lus\_Latn, crs\_Latn, ndo\_Latn, snd\_Arab, yue\_Hani, kua\_Latn, hin\_Latn,

kal\_Latn, tdt\_Latn, mfe\_Latn, mos\_Latn, kik\_Latn, cnh\_Latn, gil\_Latn, pon\_Latn, ori\_Orya, luo\_Latn, nzi\_Latn, gug\_Latn, bar\_Latn, bci\_Latn, chk\_Latn, yap\_Latn, ssw\_Latn, quz\_Latn, sah\_Cyrl, tsn\_Latn, quy\_Latn, bbc\_Latn, wal\_Latn, uig\_Arab, pam\_Latn, seh\_Latn, zai\_Latn, gym\_Latn, bod\_Tibt, nde\_Latn, fon\_Latn, nbl\_Latn, kmr\_Latn, guc\_Latn, mam\_Latn, nia\_Latn, nyn\_Latn, cab\_Latn, top\_Latn, mco\_Latn, tzh\_Latn, plt\_Latn, iba\_Latn, kek\_Latn, sop\_Latn, kac\_Latn, qvi\_Latn, cak\_Latn, kbp\_Latn, ctu\_Latn, kri\_Latn, mau\_Latn, tyv\_Cyrl, btx\_Latn, nch\_Latn, ncj\_Latn, pau\_Latn, toj\_Latn, pcm\_Latn, dyu\_Latn, kss\_Latn, quc\_Latn, yao\_Latn, kab\_Latn, tuk\_Cyrl, ndc\_Latn, san\_Deva, qug\_Latn, arb\_Arab, mck\_Latn, arn\_Latn, pdt\_Latn, gla\_Latn, kmr\_Cyrl, nav\_Latn, ksw\_Mymr, mxv\_Latn, hif\_Latn, wol\_Latn, sme\_Latn, gom\_Latn, bum\_Latn, mgr\_Latn, ahk\_Latn, tsz\_Latn, bzj\_Latn, udm\_Cyrl, cce\_Latn, meu\_Latn, cbk\_Latn, bhw\_Latn, ngu\_Latn, nyy\_Latn, naq\_Latn, toh\_Latn, nse\_Latn, alz\_Latn, mhr\_Cyrl, djk\_Latn, gkn\_Latn, grc\_Grek, swl\_Latn, alt\_Cyrl, miq\_Latn, kaa\_Cyrl, lhu\_Latn, lzh\_Hani, cmn\_Hani, kjh\_Cyrl, mgh\_Latn, rmy\_Latn, srm\_Latn, gur\_Latn, yom\_Latn, cfm\_Latn, lao\_Lao, qub\_Latn, ote\_Latn, ldi\_Latn, ayr\_Latn, bba\_Latn, aln\_Latn, leh\_Latn, ban\_Latn, ace\_Latn, pes\_Arab, ary\_Arab, hus\_Latn, glv\_Latn, mai\_Deva, dzo\_Tibt, ctd\_Latn, nnb\_Latn, sxn\_Latn, mps\_Latn, gkp\_Latn, acr\_Latn, dtp\_Latn, lam\_Latn, poh\_Latn, quh\_Latn, tob\_Latn, ach\_Latn, npi\_Deva, myv\_Cyrl, tih\_Latn, gor\_Latn, ium\_Latn, teo\_Latn, kia\_Latn, crh\_Cyrl, enm\_Latn, mad\_Latn, cac\_Latn, hnj\_Latn, ikk\_Latn, sba\_Latn, zom\_Latn, bqz\_Latn, bim\_Latn, mdy\_Ethi, bts\_Latn, gya\_Latn, agw\_Latn, knv\_Latn, giz\_Latn, hui\_Latn, hif\_Deva

#### SR-T Benchmark Languages:

mal\_Mlym, aze\_Latn, ben\_Beng, tel\_Telu, fra\_Latn, spa\_Latn, nob\_Latn, rus\_Cyrl, deu\_Latn, tur\_Latn, mar\_Deva, por\_Latn, nld\_Latn, ara\_Arab, ita\_Latn, ind\_Latn, ell\_Grek, bul\_Cyrl, swe\_Latn, ces\_Latn, isl\_Latn, pol\_Latn, ron\_Latn, dan\_Latn, hun\_Latn, srp\_Latn, ceb\_Latn, heb\_Hebr, hrv\_Latn, glg\_Latn, fin\_Latn, slv\_Latn, vie\_Latn, mkd\_Cyrl, slk\_Latn, est\_Latn, eus\_Latn, lit\_Latn, kaz\_Cyrl, bos\_Latn, epo\_Latn, cat\_Latn, tha\_Thai, ukr\_Cyrl, tgl\_Latn, gle\_Latn, hin\_Deva, kor\_Hang, urd\_Arab, sqi\_Latn, bel\_Cyrl, afr\_Latn, nno\_Latn, tat\_Cyrl, ast\_Latn, mon\_Cyrl, arz\_Arab, tam\_Taml, uzb\_Cyrl, amh\_Ethi, war\_Latn, fao\_Latn, hye\_Armn, oci\_Latn, xho\_Latn, cym\_Latn, lat\_Latn, kat\_Geor, fry\_Latn, nds\_Latn, zsm\_Latn, bre\_Latn, tuk\_Latn, khm\_Khmr, jpn\_Jpan, yue\_Hani, gsw\_Latn, lvs\_Latn, kur\_Latn, ido\_Latn, uig\_Arab, pam\_Latn, pms\_Latn, wuu\_Hani, yid\_Hebr, ina\_Latn, kab\_Latn, gla\_Latn, cbk\_Latn, hsb\_Latn, mhr\_Cyrl, swl\_Latn, cmn\_Hani, pes\_Arab, dtp\_Latn, lfn\_Latn, ile\_Latn, csb\_Latn.

#### NER Benchmark Languages:

hbs\_Latn, mal\_Mlym, aze\_Latn, guj\_Gujr, ben\_Beng, kan\_Knda, tel\_Telu, mlt\_Latn, fra\_Latn, spa\_Latn, eng\_Latn, rus\_Cyrl, deu\_Latn, tur\_Latn, pan\_Guru, mar\_Deva, por\_Latn, nld\_Latn, ara\_Arab, zho\_Hani, ita\_Latn, ind\_Latn, ell\_Grek, bul\_Cyrl, swe\_Latn, ces\_Latn, isl\_Latn, pol\_Latn, ron\_Latn, dan\_Latn, hun\_Latn, tsk\_Cyrl, fas\_Arab, ceb\_Latn, heb\_Hebr, hrv\_Latn, glg\_Latn, fin\_Latn, slv\_Latn, vie\_Latn, mkd\_Cyrl, slk\_Latn, nor\_Latn, est\_Latn, ltz\_Latn, eus\_Latn, lit\_Latn, kaz\_Cyrl, lav\_Latn, bos\_Latn, epo\_Latn, cat\_Latn, tha\_Thai, ukr\_Cyrl, tgl\_Latn, sin\_Sinh, gle\_Latn, hin\_Deva, kor\_Hang, urd\_Arab, swa\_Latn, sqi\_Latn, bel\_Cyrl, afr\_Latn, nno\_Latn, tat\_Cyrl, ast\_Latn, mon\_Cyrl, msa\_Latn, som\_Latn, srp\_Cyrl, mlg\_Latn, arz\_Arab, tam\_Taml, uzb\_Latn, cos\_Latn, als\_Latn, amh\_Ethi, sun\_Latn, war\_Latn, div\_Thaa, yor\_Latn, fao\_Latn, bak\_Cyrl, ilo\_Latn, mri\_Latn, asm\_Beng, ibo\_Latn, kin\_Latn, hye\_Armn, oci\_Latn, lin\_Latn, kir\_Cyrl, nep\_Deva, cym\_Latn, lat\_Latn, kat\_Geor, fry\_Latn, mya\_Mymr, nds\_Latn, pnb\_Arab, ckb\_Arab, chv\_Cyrl, que\_Latn, bre\_Latn, pus\_Arab, che\_Cyrl, oss\_Cyrl, nan\_Latn, lim\_Latn, tuk\_Latn, min\_Latn, khm\_Khmr, jav\_Latn, vec\_Latn, jpn\_Jpan, snd\_Arab, yue\_Hani, sco\_Latn, ori\_Orya, arg\_Latn, kur\_Latn, bar\_Latn, roh\_Latn, aym\_Latn, sah\_Cyrl, lmo\_Latn, ido\_Latn, vol\_Latn, uig\_Arab, bod\_Tibt, pms\_Latn, wuu\_Hani, yid\_Hebr, scn\_Latn, ina\_Latn, xmf\_Geor, san\_Deva, gla\_Latn, mwl\_Latn, diq\_Latn, cbk\_Latn, szl\_Latn, hsb\_Latn, vls\_Latn, mhr\_Cyrl, grn\_Latn, lzh\_Hani, mzn\_Arab, nap\_Latn, ace\_Latn, frs\_Latn, eml\_Latn, vep\_Latn, sgs\_Latn, lij\_Latn, crh\_Latn, ksh\_Latn, zea\_Latn, csb\_Latn, jbo\_Latn, bih\_Deva, ext\_Latn, fur\_Latn.

#### POS Benchmark Languages:

mal\_Mlym, ben\_Beng, tel\_Telu, mlt\_Latn, fra\_Latn, spa\_Latn, eng\_Latn, rus\_Cyrl, deu\_Latn, tur\_Latn, mar\_Deva, por\_Latn, nld\_Latn, ara\_Arab, zho\_Hani, ita\_Latn, ind\_Latn, ell\_Grek, bul\_Cyrl, swe\_Latn, ces\_Latn, isl\_Latn, pol\_Latn, ron\_Latn, dan\_Latn, hun\_Latn, srp\_Latn, fas\_Arab, ceb\_Latn, heb\_Hebr, hrv\_Latn, glg\_Latn, fin\_Latn, slv\_Latn, vie\_Latn, slk\_Latn, nor\_Latn, est\_Latn, eus\_Latn, lit\_Latn, kaz\_Cyrl, lav\_Latn, cat\_Latn, tha\_Thai, ukr\_Cyrl, tgl\_Latn, sin\_Sinh, gle\_Latn, hin\_Deva, kor\_Hang, urd\_Arab, sqi\_Latn, bel\_Cyrl, afr\_Latn, tat\_Cyrl, tam\_Taml, amh\_Ethi, yor\_Latn, fao\_Latn, hye\_Armn, cym\_Latn, lat\_Latn, nds\_Latn, bre\_Latn, hyw\_Armn, jav\_Latn, jpn\_Jpan, yue\_Hani, gsw\_Latn, sah\_Cyrl, uig\_Arab, kmr\_Latn, pcm\_Latn, quc\_Latn, san\_Deva, gla\_Latn, wol\_Latn, sme\_Latn, hsb\_Latn, grc\_Grek, hbo\_Hebr, grn\_Latn, lzh\_Hani, ajp\_Arab, nap\_Latn, aln\_Latn, glv\_Latn, lij\_Latn, myv\_Cyrl, bam\_Latn, xav\_Latn.

## C.2 For Benchmark Performances in Figure 2

SR-T Benchmark Languages:

tur\_Latn, ell\_Grek, bul\_Cyrl, ces\_Latn, kor\_Hang, zsm\_Latn, kat\_Geor, fry\_Latn, khm\_Khmr, yue\_Hani, tuk\_Latn, uig\_Arab, pam\_Latn, kab\_Latn, gla\_Latn, mhr\_Cyrl, swl\_Latn, cmn\_Hani, pes\_Arab, dtp\_Latn

SR-B Benchmark Languages:

tur\_Latn, ell\_Grek, bul\_Cyrl, ces\_Latn, kor\_Hang, zsm\_Latn, kat\_Geor, fry\_Latn, khm\_Khmr, yue\_Hani, tuk\_Latn, uig\_Arab, pam\_Latn, kab\_Latn, gla\_Latn, mhr\_Cyrl, swl\_Latn, cmn\_Hani, pes\_Arab, dtp\_Latn