

WECHSEL: Effective initialization of subword embeddings for cross-lingual transfer of monolingual language models

Anonymous ACL submission

Abstract

001 Recently, large pretrained language models
002 (LMs) have gained popularity. Training these
003 models requires ever more computational re-
004 sources and most of the existing models are
005 trained on English text only. It is exceedingly
006 expensive to train these models in other lan-
007 guages. To alleviate this problem, we intro-
008 duce a method – called WECHSEL – to trans-
009 fer English models to new languages. We ex-
010 change the tokenizer of the English model with
011 a tokenizer in the target language and initialize
012 token embeddings such that they are close to
013 semantically similar English tokens by utiliz-
014 ing multilingual static word embeddings cover-
015 ing English and the target language. We use
016 WECHSEL to transfer GPT-2 and RoBERTa
017 models to 4 other languages (French, Ger-
018 man, Chinese and Swahili). WECHSEL im-
019 proves over a previously proposed method for
020 cross-lingual parameter transfer and outper-
021 forms models of comparable size trained from
022 scratch in the target language with up to 64x
023 less training effort. Our method makes train-
024 ing large language models for new languages
025 more accessible and less damaging to the envi-
026 ronment. We make our code and models pub-
027 licly available.

028 1 Introduction

029 Large LMs based on the Transformer architec-
030 ture (Vaswani et al., 2017) have become increas-
031 ingly popular since GPT (Radford et al., 2018)
032 and BERT (Devlin et al., 2019) were introduced,
033 prompting the creation of many large LMs pre-
034 trained on English text (Yang et al., 2019; Clark
035 et al., 2020; Lewis et al., 2020; Joshi et al., 2020;
036 Ram et al., 2021). There is a tendency towards
037 training larger and larger models (Brown et al.,
038 2020; Fedus et al., 2021) while restricting focus
039 to the English language. Recent work has called
040 attention to the costs associated with training in-
041 creasingly large LMs, including environmental cost

and financial cost (Bender et al., 2021). If train- 042
ing large LMs for English is already costly, it is 043
prohibitively expensive to train new, similarly pow- 044
erful models to cover all other relevant languages. 045

046 One approach to address this issue is creating
047 massively multilingual models (Devlin et al., 2019;
048 Conneau et al., 2020; Xue et al., 2021) which are
049 trained on a concatenation of text in many different
050 languages. These models exhibit natural language
051 understanding capabilities in a wide variety of lan-
052 guages, but suffer from what Conneau et al. (2020)
053 call the *curse of multilinguality*: beyond a certain
054 number of languages in the training data, overall
055 performance decreases on monolingual as well as
056 cross-lingual tasks. Consistent with this finding,
057 Nozza et al. (2020) observe that monolingual LMs
058 often outperform massively multilingual models.
059 It is thus desirable to train monolingual models
060 in more languages. Training monolingual models
061 in non-English languages is commonly done by
062 training a new model with randomly initialized pa-
063 rameters (Antoun et al., 2020; Louis, 2020; Chan
064 et al., 2020; Martin et al., 2020). But to train a
065 model with capabilities comparable to that of an
066 English model in this way, presumably a similar
067 amount of compute to what was used to train the
068 English model would be required.

069 To address this issue, we introduce WECHSEL¹,
070 a novel method to transfer monolingual language
071 models to a new language. WECHSEL uses multi-
072 lingual static word embeddings between the source
073 language and the target language to initialize model
074 parameters. We copy all inner (non-embedding)
075 parameters of the English model, exchange the tok-
076 enizer with a tokenizer for the target language and
077 instead of randomly initializing the token embed-
078 dings as done in prior work (de Vries and Nissim,
079 2021), we initialize token embeddings in the target
080 language such that they are close to semantically

¹Word Embeddings Can Help initialize Subword Embed-
dings in a new Language.

similar English tokens by mapping multilingual static word embeddings to subword embeddings. Embeddings take up roughly 31% of the parameters of RoBERTa (Liu et al., 2019) and roughly 33% of the parameters of GPT2 (Radford et al., 2019). Intuitively, semantically transferring embeddings instead of randomly initializing one third of the model should result in improved performance. Our parameter transfer aims to provide an effective initialization in the target language, requiring significantly fewer training steps to reach high performance. As multilingual static word embeddings are available for many languages (Bojanowski et al., 2017), WECHSEL is widely applicable.

We evaluate our method by transferring English RoBERTa and GPT-2 – as representative models of encoder and decoder language models respectively – to 4 new languages (French, German, Chinese and Swahili). We evaluate our RoBERTa models by fine-tuning on Neural Entity Recognition (NER) and Natural Language Inference (NLI) tasks in the respective languages. Our GPT-2 models are evaluated by computing Language Modelling Perplexity (PPL) on a hold-out set. We compare WECHSEL initialization with randomly initialized models (denoted as FullRand) as well as a recently proposed method which only transfers the inner (non-embedding) parameters (denoted as TransInner, de Vries and Nissim (2021)) under the same training conditions (around 4 days on a TPUv3-8). We also compare our model with models of comparable size trained from scratch under significantly larger training regimes, in particular CamemBERT (Martin et al., 2020) (French), GBERT_{Base} (Chan et al., 2020) (German), and BERT_{Base}-Chinese (Devlin et al. (2019) (Chinese). These models are trained on 6.4, 3.9, and 2 times more tokens, respectively. Results show that RoBERTa models initialized with WECHSEL outperform randomly initialized models by an average of 6.17% accuracy on NLI and 1.37% micro F1 score on NER, and models initialized with TransInner by an average of 0.9% accuracy and 0.5% micro F1 score on NLI and NER, respectively. GPT-2 models initialized with WECHSEL outperform randomly initialized models by an average 0.75 PPL and models initialized with TransInner by an average 1.42 PPL. Our models already outperform GBERT_{Base} and CamemBERT on average on downstream tasks after 10% of training steps. Our contribution is summarized as follows.

- We propose WECHSEL, a novel method for transferring monolingual language models to a new language by utilizing multilingual static word embeddings between the source and the target language.
- We show effective transfer of RoBERTa and GPT-2 using WECHSEL to 4 different languages and high performance after minimal training effort.
- We train more effective GPT-2 and RoBERTa models for German, French, Chinese and Swahili than previously published models under a more efficient training setting. Our code and models are publicly available at <https://github.com/anonymized>.

In the following, we review related work in Section 2. We then introduce the WECHSEL method in Section 3, followed by explaining the experiment setup in Section 4. We show and discuss results in Section 5.

2 Related Work

Large Language Models. Training Language Models is usually done in a self-supervised manner i.e. deriving labels from the training text instead of needing explicit annotations. One widely-used optimization objective is Masked Language Modelling (Devlin et al., 2019, MLM), where random tokens in the input are masked (replaced by a special [MASK] token), and the task is to predict the original tokens. Another common objective is Causal Language Modelling (CLM), where the task is to predict the next token. These two objectives highlight a fundamental distinction between language models: models can be trained as encoders (e.g. with MLM) or as decoders (e.g. with CLM).

Instead of words, the vocabulary of language models usually consists of subwords. A subword is a combination of characters below or at word-level. Many recently proposed language models use subword tokenization (Clark et al., 2020; Liu et al., 2019; Devlin et al., 2019). WECHSEL can be used for any model which (1) uses subword-based tokenization and (2) learns an embedding for each token.

Multilingual representations. There has been a significant amount of work in creating multilingual static word embeddings. Multilingual static word embeddings can be created by learning static word

180 embeddings from scratch using data in multiple
181 languages (Luong et al., 2015; Duong et al., 2016).
182 Alternatively, multilinguality can be achieved by
183 aligning existing monolingual word embeddings
184 using a bilingual dictionary, so that the resulting
185 embeddings share the same semantic space (Xing
186 et al., 2015; Joulin et al., 2018). Recent studies
187 improve this by reducing the need for bilingual data
188 (or even requiring no bilingual data at all) (Artetxe
189 et al., 2017, 2018; Lample et al., 2018).

190 Besides multilingual static word embeddings,
191 multilinguality is also relevant to contextualized
192 representations. Multilingual contextualized repre-
193 sentations can be learned through training a model
194 on a concatenation of corpora in different lan-
195 guages. Among such model are mBERT (Devlin
196 et al., 2019), XLM-R (Conneau et al., 2020) and
197 mT5 (Xue et al., 2021), which are trained on text
198 in 104, 100 and 101 languages, respectively. As
199 shown by Pires et al. (2019), a multilingual model
200 such as mBERT can enable cross-lingual transfer
201 by using task-specific annotations in one language
202 to fine-tune the model for evaluation in another lan-
203 guage. However, recent studies outline a number
204 of problems with massively multilingual models.
205 Wu and Dredze (2020) empirically show that in
206 mBERT “the 30% languages with least pretraining
207 resources perform worse than using no pretrained
208 language model at all”. Conneau et al. (2020) re-
209 port that beyond a certain number of languages in
210 the training data, the overall performance decreases
211 both on monolingual as well as cross-lingual tasks.
212 These studies motivate our work on creating mono-
213 lingual LMs for more languages.

214 **Cross-lingual transfer of monolingual LMs.**

215 Studies related to the cross-lingual transfer of
216 monolingual language models can be divided into
217 two categories:

- 218 • **Bilingualization of a monolingual LM** is
219 concerned with transferring a model to a
220 new language while preserving capabilities
221 in the original language. Artetxe et al. (2020)
222 achieve this goal by replacing the tokenizer
223 and relearning the token embeddings, while
224 freezing other parameters. Such a model be-
225 comes bilingual, since the initial tokenizer
226 and embeddings can be used for tasks in the
227 source language while the new tokenizer and
228 embeddings can be used for tasks in the tar-
229 get language. Thus, a model can be finetuned
230 on annotated task data in the source language,

231 then zero-shot transferred to the target lan-
232 guage. Tran (2020) follow a similar approach,
233 while instead of randomly initializing embed-
234 dings, they utilize static word embeddings to
235 initialize embeddings in the target language
236 close to semantically similar English tokens.
237 They then continue training the model on an
238 English text corpus as well as on the target lan-
239 guage in order to preserve model capabilities
240 in English.

- 241 • **Creating a new monolingual LM in the tar-**
242 **get language** is, in contrast, concerned with
243 creating a model in the target language with-
244 out the necessity to preserve its capabilities in
245 the source language. Zoph et al. (2016) and
246 Nguyen and Chiang (2017) show that cross-
247 linguistically transferring a machine translation
248 model can improve performance, especially
249 for low-resource languages. They replace the
250 model tokenizer with a tokenizer for the target
251 language. Zoph et al. (2016) then use embed-
252 dings of random tokens in the original vocabu-
253 lary to initialize token embeddings in the new
254 vocabulary, while Nguyen and Chiang (2017)
255 improve on this by utilizing vocabulary over-
256 lap between the source and target language.
257 More recently, de Vries and Nissim (2021) fol-
258 low a similar approach to the one of Artetxe
259 et al. (2020) by transferring a GPT-2 model to
260 a new language. de Vries and Nissim (2021)
261 add an additional step, where they train the
262 entire model for some amount of steps to al-
263 low adapting to the target language beyond
264 the lexical level. We refer to the method of
265 de Vries and Nissim (2021) as TransInner and
266 consider it as a baseline in our experiments.

267 Our WECHSEL method belongs to the second
268 category. WECHSEL is an extension to the method
269 proposed by Tran (2020) with the goal of creating a
270 new monolingual LM instead of bilingualizing the
271 LM. This allows removing the constraints imposed
272 by the need to preserve capabilities in the source
273 language. In addition, we generalize the semantic
274 subword mapping done by Tran (2020) to consider
275 an arbitrary number of neighbors with an arbitrary
276 temperature. We are the first to show that a cross-
277 linguistically transferred model can outperform mono-
278 lingual models which have been trained extensively
279 from scratch in the target language, while requiring
280 substantially less computational resources.

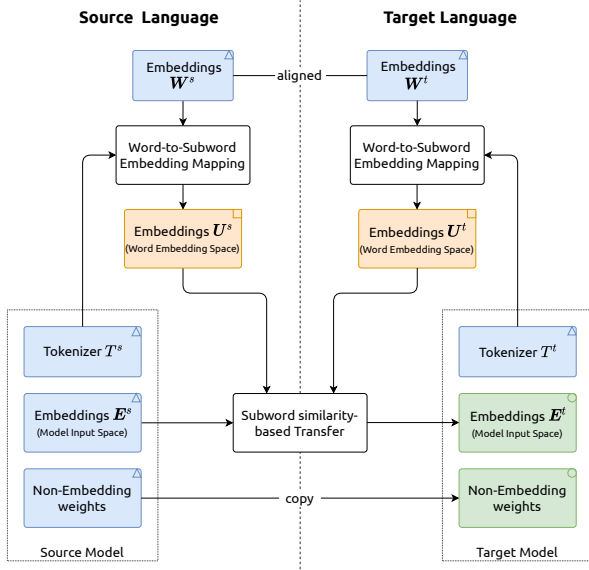


Figure 1: Summary of our WECHSEL method. We show **inputs** \triangle , **intermediate results** \square and **outputs** \circ .

3 Methodology

Given the tokenizer T^s in the source language (with vocabulary \mathbb{U}^s), the corresponding token embeddings E^s and a tokenizer T^t in the target language (with vocabulary \mathbb{U}^t), our goal is to find a good initialization of the embeddings E^t by using E^s . To this end, we use existing bilingual word embeddings, containing a set of words in the source and target language and their aligned vectors, as well as word frequency statistics. We denote the set of words in the source and target language as \mathbb{V}^s and \mathbb{V}^t respectively, and the aligned static word embeddings as W^s and W^t .

First, independently for both languages, we compute static subword embeddings for tokens in the tokenizer vocabulary in the same semantic space as the static word embeddings (Section 3.1). This results in subword embeddings U^s and U^t for the source and target language, respectively. Next, we use U^s and U^t to compute semantic similarity of every subword in \mathbb{U}^s to every subword in \mathbb{U}^t . Using these semantic similarities, we initialize the embeddings in E^t through an affine combination of embeddings in E^s (Section 3.2). By applying WECHSEL, the vectors of E^t are in the same semantic space as E^s , where a subword in the target language is semantically similar to its counterpart(s) in the source language. These steps are summarized in Figure 1 and explained in detail in the following.

3.1 Word-to-Subword Embedding Mapping

The process of mapping word embeddings to subword embeddings is done separately for source and target languages. Given a tokenizer T (with vocabulary \mathbb{U}), word embeddings W , and word frequencies f , the goal is to find subword embeddings U for subwords in \mathbb{U} using W . To this end, we apply the tokenizer T to every word v in \mathbb{V} resulting in a set of subwords for each word. We define $\mathbb{V}^{(x)}$ as the set of words containing the subword x when tokenized. The embedding u_x of the subword x is then defined as the average of the embeddings of words in $\mathbb{V}^{(x)}$, weighted by the word frequencies.

$$u_x = \frac{\sum_{v \in \mathbb{V}^{(x)}} w_v \cdot f_v}{\sum_{v \in \mathbb{V}^{(x)}} f_v}$$

where w_v is the embedding and f_v is the frequency of word v . Subwords which do not occur in any word are initialized to zero. We implement this method of word-to-subword mapping using a procedure we refer to as *tokenize-flatten-reduce* as depicted in Figure 2.

3.2 Subword similarity-based Transfer

Applying the previous step to both source and target language results in the subword embeddings U^s and U^t over the subword vocabularies \mathbb{U}^s and \mathbb{U}^t , respectively. Our aim is now to use these embeddings to find an effective transformation from E^s to E^t . We first compute the cosine similarity of every subword $x \in \mathbb{U}^t$ to every subword $y \in \mathbb{U}^s$, denoted as $s_{x,y}$.

$$s_{x,y} = \frac{u_x^t u_y^{sT}}{\|u_x^t\| \|u_y^s\|}$$

We now exploit these similarities to initialize embeddings in E^t by an affine combination of embeddings in E^s . Each subword embedding in E^t is defined as the weighted mean of the k nearest embeddings in E^s according to the similarity values. The weighting is done by a softmax of the similarities with temperature τ .

$$e_x^t = \frac{\sum_{y \in \mathcal{J}_x} \exp(s_{x,y}/\tau) \cdot e_y^s}{\sum_{y' \in \mathcal{J}_x} \exp(s_{x,y'}/\tau)}$$

where \mathcal{J}_x is the set of k neighbouring subwords in the source language. Subword embeddings for which U^t is zero are initialized from a random normal distribution $\mathcal{N}(\mathbb{E}[E^s], \text{Var}[E^s])$. The inner (non-embedding) parameters are simply copied from the source model.

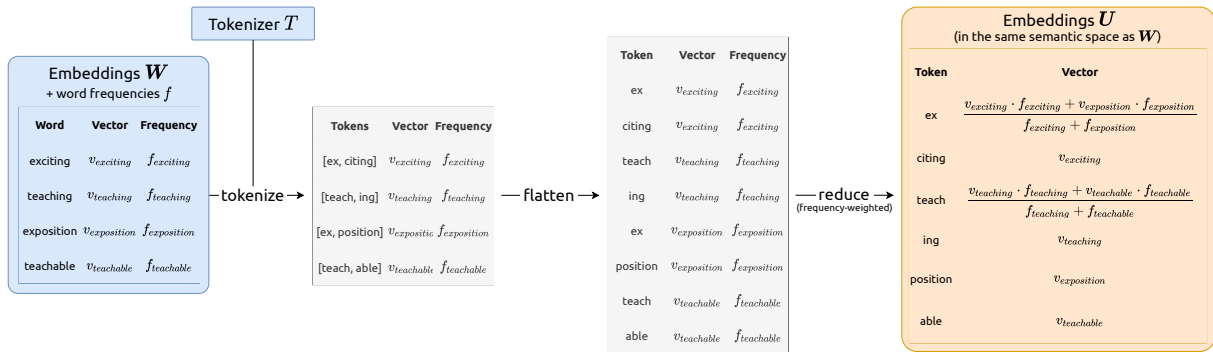


Figure 2: Word-to-subword embedding mapping. First, **tokenize** all words in the word embeddings. Then **flatten** the result by assigning the embeddings of the words in which it occurred and their word frequencies to each subword. Finally, **reduce** the embeddings assigned to each subword by taking their mean, weighted by word frequency.

4 Experiment Design

We evaluate our method by transferring the English RoBERTa model (Liu et al., 2019) and the English GPT-2 model (Radford et al., 2019) to a subset of the 7 languages proposed in Conneau et al. (2020) (French, German, Chinese and Swahili). Chinese allows evaluation of cross-lingual transfer to a strongly dissimilar language. Swahili serves to evaluate our method on a low-resource language. We use the pretrained models RoBERTa_{Base} with 125M parameters, and the small GPT-2 variant with 117M parameters provided by HuggingFace’s Transformers (Wolf et al., 2020) in all experiments.

To ensure our method does not depend on excessive amounts of data in the target language we restrict the amount of training data to subsets of 4GiB from the OSCAR corpus (Ortiz Suárez et al., 2019) for all experiments for all languages, except Swahili. For Swahili, we use the 1.6GiB Swahili subset of the CC-100 corpus (Conneau et al., 2020). To obtain aligned word embeddings between the source and the target language we use monolingual fastText word embeddings² (Bojanowski et al., 2017) and align them using the Orthogonal Procrustes method (Schönemann, 1966; Artetxe et al., 2016) with bilingual dictionaries from MUSE³ (Conneau et al., 2017) for French, German and Chinese and a bilingual dictionary from FreeDict⁴ (Bański and Wójtowicz, 2009) for Swahili. We use word frequency information provided as part of the fastText word vectors. We choose temperature $\tau = 0.1$ and neighbors $k = 10$ for WECHSEL by conducting a grid search over initializations with varying k and τ using linear probes (Appendix A).

²<https://fasttext.cc>

³<https://github.com/facebookresearch/MUSE>

⁴<https://freedict.org>

We train tokenizers in the target languages using a vocabulary size of 50k tokens and byte-level BPE (Radford et al., 2019). After applying WECHSEL, we continue training RoBERTa on the MLM objective and GPT-2 on the CLM objective. We compare against two baseline methods.

- Randomly initializing E^t while transferring all other parameters from the English model as in de Vries and Nissim (2021). After training only embeddings for a fixed amount of steps while freezing other parameters, the entire model is trained for the remaining steps. We refer to this method as TransInner.
- Training from scratch in the target language, as is commonly done when training BERT-like or GPT-like models in a new language (Antoun et al., 2020; Louis, 2020; Chan et al., 2020; Martin et al., 2020). We refer to this method as FullRand.

All models are trained for 250k steps with the same hyperparameters across all languages (reported in Appendix B). Training one model takes around 4 days on a TPUv3-8. For WECHSEL and FullRand we use a learning rate (LR) schedule with linear warmup from zero to peak LR for the first 10% of steps, then linear decay to zero. For TransInner we perform two warmup phases from zero to peak LR, once for the first 10% of steps for training embeddings only, then again for the remaining steps while training the entire models.

5 Results

We show results for RoBERTa on two downstream tasks (NLI, NER) and for GPT-2 (CLM). The performance of our models throughout training is shown in Figure 3.

Lang	Model	Score@0			Score@25k			Score@250k			Score (more training)		
		NLI	NER	Avg	NLI	NER	Avg	NLI	NER	Avg	NLI	NER	Avg
French	WECHSEL-RoBERTa	<u>78.25</u>	<u>87.43</u>	<u>82.84</u>	<u>81.86</u>	<u>90.07</u>	<u>85.96</u>	<u>82.55</u>	<u>90.80</u>	<u>86.68</u>	-	-	-
	TransInner-RoBERTa	60.86	69.57	65.21	65.49	83.82	74.66	81.75	90.34	86.04	-	-	-
	FullRand-RoBERTa	55.71	70.79	63.25	69.02	84.24	76.63	75.28	89.30	82.29	-	-	-
	CamemBERT	-	-	-	-	-	-	-	-	-	80.88	90.26	85.57
	XLM-R _{Base}	-	-	-	-	-	-	-	-	-	79.25	89.48	84.37
German	WECHSEL-RoBERTa	<u>77.00</u>	<u>84.70</u>	<u>80.85</u>	<u>80.71</u>	<u>89.09</u>	<u>84.90</u>	<u>82.04</u>	<u>89.72</u>	<u>85.88</u>	-	-	-
	TransInner-RoBERTa	58.51	65.23	61.87	64.78	82.05	73.42	80.75	89.30	85.02	-	-	-
	FullRand-RoBERTa	54.82	66.84	60.83	68.02	81.53	74.77	75.48	88.36	81.92	-	-	-
	GBERT _{Base}	-	-	-	-	-	-	-	-	-	78.64	89.46	84.05
	XLM-R _{Base}	-	-	-	-	-	-	-	-	-	78.58	88.76	83.67
Chinese	WECHSEL-RoBERTa	<u>62.75</u>	<u>72.87</u>	<u>67.81</u>	<u>77.07</u>	<u>78.03</u>	<u>77.55</u>	<u>77.99</u>	<u>80.65</u>	<u>79.32</u>	-	-	-
	TransInner-RoBERTa	46.95	69.06	58.01	52.96	73.35	63.16	76.99	80.00	78.49	-	-	-
	FullRand-RoBERTa	44.24	57.95	51.09	58.34	64.84	61.59	71.38	78.35	74.86	-	-	-
	BERT _{Base} -Chinese	-	-	-	-	-	-	-	-	-	76.55	82.05	79.30
	XLM-R _{Base}	-	-	-	-	-	-	-	-	-	76.41	78.36	77.38
Swahili	WECHSEL-RoBERTa	<u>60.14</u>	<u>75.42</u>	<u>67.78</u>	<u>74.04</u>	<u>87.79</u>	<u>80.92</u>	<u>74.58</u>	<u>87.66</u>	<u>81.12</u>	-	-	-
	TransInner-RoBERTa	54.67	64.46	59.56	58.85	80.27	69.56	74.10	87.05	80.57	-	-	-
	FullRand-RoBERTa	50.59	62.35	56.47	63.79	83.49	73.64	70.34	87.34	78.84	-	-	-
	XLM-R _{Base}	-	-	-	-	-	-	-	-	-	69.18	87.37	78.28

Table 1: Results from fine-tuning RoBERTa models. We report accuracy for NLI on XNLI and micro F1 score for NER on WikiANN. Results are averaged over 3 runs. We report scores before training (**Score@0**), after 10% of steps (**Score@25k**) and after training (**Score@250k**). We also report results from fine-tuning prior monolingual models and XLM-R (**Score (more training)**) which are trained on more tokens than our models (c.f. Section 5.3). For each language, the best results in every column are indicated with underlines. The overall best results, including the comparison with prior monolingual models of comparable size, are shown in bold.

5.1 Transferring RoBERTa

We evaluate WECHSEL-RoBERTa by fine-tuning on XNLI (Conneau et al., 2018), and on the balanced train-dev-test split of WikiANN (Rahimi et al., 2019; Pan et al., 2017) to evaluate NLI and NER performance, respectively. The hyperparameters used for fine-tuning are reported in Appendix B.

Table 1 reports the evaluation results on RoBERTa. As shown, models initialized with WECHSEL outperform models trained from scratch and models initialized with TransInner across all languages. Surprisingly, close relatedness of the source and target language is not necessary to achieve effective transfer, as e.g. on NLI WECHSEL improves by 7.27%, 6.57%, 6.61% and 4.23% absolute accuracy over models trained from scratch for French, German, Chinese and Swahili, respectively.

Next, we compare WECHSEL-RoBERTa to monolingual models CamemBERT (Martin et al., 2020) (French), GBERT_{Base} (Chan et al., 2020) (German), and BERT_{Base}-Chinese (Devlin et al., 2019) (Chinese), and to fine-tuning XLM-R_{Base} (Artetxe et al., 2020) in the target language. To the best of our knowledge there is no monolingual model available for Swahili. We observe a consistent improvement over XLM-R_{Base} by an av-

erage 3.43% accuracy for NLI and 1.21% micro F1 score for NER. For NLI, we improve over the prior monolingual models by 1.67%, 3.4% and 1.44% absolute accuracy for French, German and Chinese, respectively and even outperform prior monolingual models after 10% of our training steps. For NER, we observe a less marked improvement over monolingual models with 0.54% and 0.26% absolute micro F1 score improvement for French and German, respectively. For Chinese, the monolingual model BERT_{Base}-Chinese still outperforms our method by 1.4% absolute micro F1 score. We suspect that the discrepancy between NLI and NER is due to the limited training corpus size (4GiB), while a larger corpus can potentially improve NER as more named entities appear (Martin et al., 2020).

5.2 Transferring GPT-2

GPT-2 is evaluated by Perplexity (PPL) on a hold-out set from the same corpus on which the model was trained on (OSCAR for French, German and Chinese; CC-100 for Swahili). Results are shown in Table 2. Consistent with results for WECHSEL-RoBERTa, the GPT-2 models trained with WECHSEL outperform the models trained from scratch and the models trained with TransInner across all languages.

For the language modeling task, we observe a stronger dependence on similarity of the source to

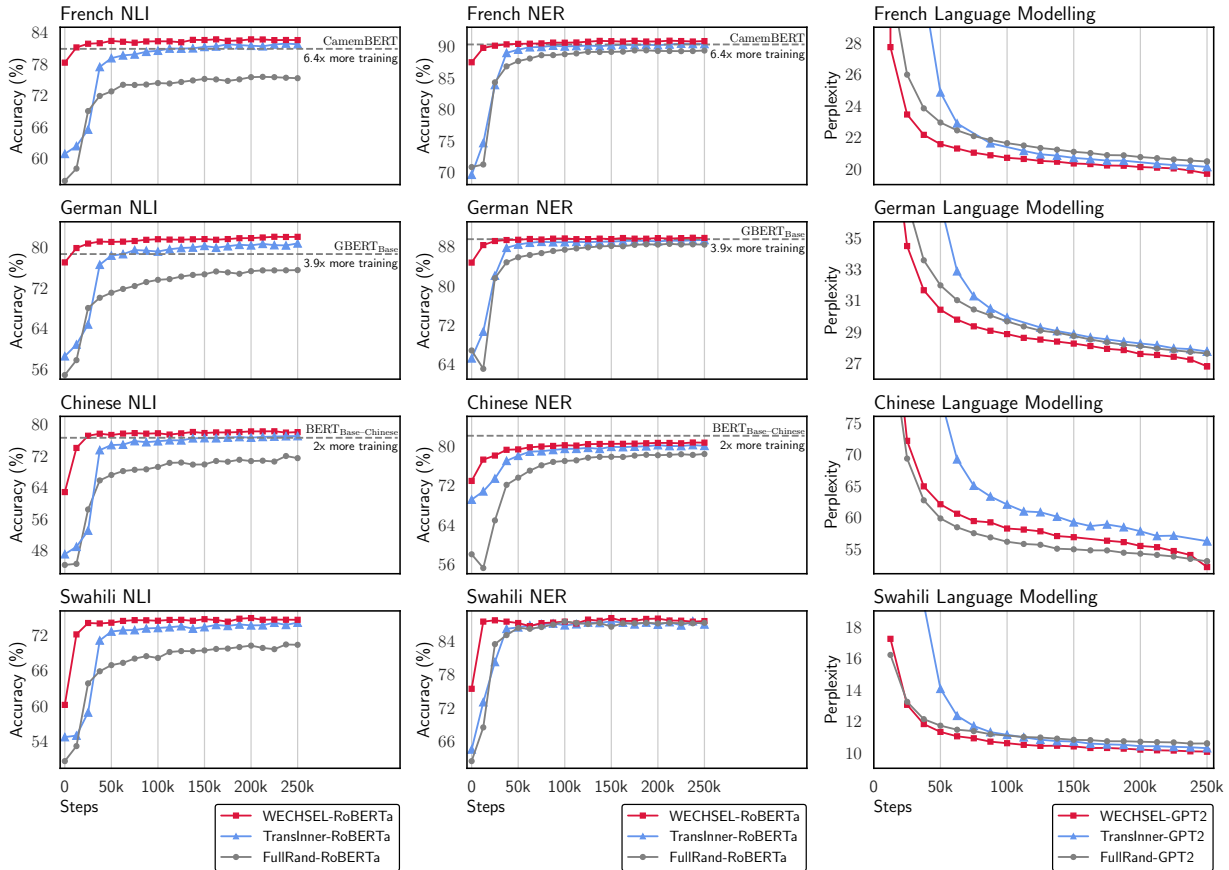


Figure 3: Test scores over training steps from fine-tuning RoBERTa models on NLI (using XNLI) and NER (using WikiANN). Perplexity on the hold-out set over training steps of GPT-2 models. We evaluate every 12.5k steps.

Lang	Model	PPL@0	PPL@25k	PPL@250k
French	WECHSEL-GPT2	2.3e+3	23.45	19.70
	TransInner-GPT2	1.4e+5	67.97	20.13
	FullRand-GPT2	5.9e+4	25.99	20.47
German	WECHSEL-GPT2	5.0e+3	34.46	26.82
	TransInner-GPT2	1.5e+5	121.67	27.76
	FullRand-GPT2	5.8e+4	37.29	27.63
Chinese	WECHSEL-GPT2	2.5e+4	72.11	52.07
	TransInner-GPT2	1.5e+5	231.05	56.17
	FullRand-GPT2	5.8e+4	69.29	52.98
Swahili	WECHSEL-GPT2	1.5e+5	13.03	10.06
	TransInner-GPT2	1.4e+5	42.95	10.28
	FullRand-GPT2	5.8e+4	13.22	10.58

Table 2: Results of training GPT2 models. We report Perplexity before training (**PPL@0**), after 10% of steps (**PPL@25k**) and after training (**PPL@250k**).

the target language than for downstream tasks such as NLI or NER. For French and German, WECHSEL is consistently better than TransInner and FullRand throughout the entire training. For Chinese, a decrease in perplexity towards the end of training causes WECHSEL to surpass training from scratch.

5.3 Effect of training effort

To highlight the improvement in training efficiency of WECHSEL as opposed to prior monolingual models, we consider the total number of tokens the model has encountered in the target language, computed as the product of batch size \times sequence length \times train steps (shown in Table 3). We expect FullRand-RoBERTa to approach performance of the respective prior monolingual models when trained on the same amount of tokens⁵. This allows quantifying the difference in training effort required to achieve good performance with WECHSEL as opposed to training from scratch. For French, WECHSEL-RoBERTa outperforms CamemBERT after 10% of training, reducing training effort by 64x. For German, WECHSEL-RoBERTa outperforms GBERT_{Base} after 10% of training steps, reducing training effort by 39x. For Chinese, WECHSEL-RoBERTa outperforms BERT_{Base}-Chinese on NLI, but does not outperform BERT_{Base}-Chinese on NER.

⁵It would presumably be slightly worse because we restrict training corpus size to 4GiB.

Model	Tokens trained on	Factor
WECHSEL-RoBERTa	65.5B	1.0x
TransInner-RoBERTa	65.5B	1.0x
FullRand-RoBERTa	65.5B	1.0x
CamemBERT	419.4B	6.4x
GBERT _{Base}	255.6B	3.9x
BERT _{Base} -Chinese	131.1B	2.0x

Table 3: Tokens trained on in the target language between our models and previous monolingual models.

5.4 Additional Analyses

To qualitatively assess how well subword tokens can be mapped between the source and the target language, we show a random sample of tokens in the target language and their most similar English tokens (according to WECHSEL) for each language (Appendix C). We also consider using fast-Text subword information as an alternative way to map word to subword embeddings (Appendix D).

5.4.1 Relearning embeddings

To quantitatively evaluate the mapping resulting from WECHSEL we conduct an additional experiment where we keep all non-embedding parameters frozen and only train the embeddings. The better the initialization, the less improvement would be possible with more training. We conduct this experiment for French as most similar language to English and Chinese as most dissimilar language and train GPT-2 models. Hyperparameters match the ones of our main experiments, except that we train for 75k steps only. We observe a strong decrease in Perplexity from training, indicating that our mapping is far from optimal. Especially early merges in the BPE vocabulary (i.e. common tokens) change compared to their initial value. Future work could investigate improving the mapping done by WECHSEL under this metric. Additional information is shown in Appendix E.

5.4.2 Is freezing necessary?

Previous work using the TransInner method freezes non-embedding parameters for a fixed amount of steps before training the entire model (de Vries and Nissim, 2021). This is done to prevent catastrophic forgetting at the beginning of training. To evaluate if freezing non-embedding parameters is still necessary with our method, we conduct an additional experiment. We train a model with WECHSEL and a model with TransInner without freezing any parameters, and the same models with freezing of non-embedding parameters for the first 10% of steps. We again match hyperparameters of

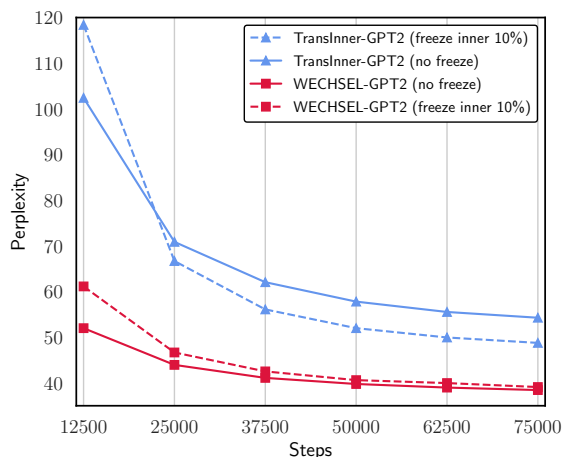


Figure 4: Comparison of German GPT-2 models trained with WECHSEL and TransInner between freezing non-embedding parameters at the start and not freezing any parameters.

the main experiments but train for 75k steps only. We train a German GPT-2 model and conjecture that the same result will hold for other languages and model types. Results are shown in Figure 4. We conclude that freezing is necessary when using TransInner, but there is no need for freezing when using WECHSEL (in fact, freezing slightly decreases performance).

6 Conclusion

We introduce WECHSEL, an effective method to transfer monolingual language models to new languages. WECHSEL exploits multilingual static word embeddings to compute an effective initialization of subword embeddings in the target language. Experiments on transferring representative transformers-based encoder and decoder language models from English to French, German, Chinese and Swahili show that the transferred RoBERTa and GPT-2 models are more efficient than strong baselines, and outperform prior monolingual models that have been trained for a significantly longer time. WECHSEL facilitates the creation of effective monolingual LMs in new languages in low resource and computationally-limited settings. Our work provides further evidence towards the hypothesis by Artetxe et al. (2020) that deep monolingual language models learn some abstractions that generalize across languages.

576
577
578
579
580
581
582
583

584
585
586
587
588
589
590

591
592
593
594
595
596
597

598
599
600
601
602
603
604

605
606
607
608
609
610

611
612
613

614
615
616
617
618
619

620
621
622
623

624
625
626
627
628

629
630
631
632

References

Wissam Antoun, Fady Baly, and Hazem Hajj. 2020. [AraBERT: Transformer-based model for Arabic language understanding](#). In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 9–15, Marseille, France. European Language Resource Association.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2016. [Learning principled bilingual mappings of word embeddings while preserving monolingual invariance](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2289–2294, Austin, Texas. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2017. [Learning bilingual word embeddings with \(almost\) no bilingual data](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 451–462, Vancouver, Canada. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. [A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 789–798, Melbourne, Australia. Association for Computational Linguistics.

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.

Piotr Bański and Beata Wójtowicz. 2009. [Freedict: an open source repository of tei-encoded bilingual dictionaries](#).

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. [On the dangers of stochastic parrots: Can language models be too big?](#). In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *arXiv preprint arXiv:2005.14165*.

Branden Chan, Stefan Schweter, and Timo Möller. 2020. [German’s next language model](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6788–6796, Barcelona,

Spain (Online). [International Committee on Computational Linguistics](#). 633
634

Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. [Electra: Pre-training text encoders as discriminators rather than generators](#). *arXiv preprint arXiv:2003.10555*. 635
636
637
638

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics. 639
640
641
642
643
644
645
646
647

Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. [Word translation without parallel data](#). *arXiv preprint arXiv:1710.04087*. 648
649
650
651

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. [Xnli: Evaluating cross-lingual sentence representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. 652
653
654
655
656
657
658

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 *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 836–846, Online. Association for Computational Linguistics. 659
660
661
662
663
664

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 665
666
667
668
669
670
671
672
673

Long Duong, Hiroshi Kanayama, Tengfei Ma, Steven Bird, and Trevor Cohn. 2016. [Learning crosslingual word embeddings without bilingual corpora](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1285–1295, Austin, Texas. Association for Computational Linguistics. 674
675
676
677
678
679
680

William Fedus, Barret Zoph, and Noam Shazeer. 2021. [Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity](#). *arXiv preprint arXiv:2101.03961*. 681
682
683
684

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. [SpanBERT: Improving pre-training by representing and predicting spans](#). *Transactions of the Association for Computational Linguistics*, 8:64–77. 685
686
687
688
689

690	Armand Joulin, Piotr Bojanowski, Tomas Mikolov, Hervé Jégou, and Edouard Grave. 2018. Loss in translation: Learning bilingual word mapping with a retrieval criterion . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2979–2984, Brussels, Belgium. Association for Computational Linguistics.	
691		
692		
693		
694		
695		
696		
697	Guillaume Lample, Alexis Conneau, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. Word translation without parallel data . In <i>International Conference on Learning Representations</i> .	
698		
699		
700		
701	Mike Lewis, Marjan Ghazvininejad, Gargi Ghosh, Armen Aghajanyan, Sida Wang, and Luke Zettlemoyer. 2020. Pre-training via paraphrasing . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 18470–18481. Curran Associates, Inc.	
702		
703		
704		
705		
706	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	
707		
708		
709		
710		
711	Antoine Louis. 2020. BelGPT-2: a GPT-2 model pre-trained on French corpora. https://github.com/antoiloui/belgpt2 .	
712		
713		
714	Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Bilingual word representations with monolingual quality in mind . In <i>Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing</i> , pages 151–159, Denver, Colorado. Association for Computational Linguistics.	
715		
716		
717		
718		
719		
720	Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamel Seddah, and Benoît Sagot. 2020. CamemBERT: a tasty French language model . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7203–7219, Online. Association for Computational Linguistics.	
721		
722		
723		
724		
725		
726		
727		
728	Toan Q. Nguyen and David Chiang. 2017. Transfer learning across low-resource, related languages for neural machine translation . In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 296–301, Taipei, Taiwan. Asian Federation of Natural Language Processing.	
729		
730		
731		
732		
733		
734		
735	Debora Nozza, Federico Bianchi, and Dirk Hovy. 2020. What the [mask]? making sense of language-specific bert models. <i>arXiv preprint arXiv:2003.02912</i> .	
736		
737		
738	Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. 2019. Asynchronous pipelines for processing huge corpora on medium to low resource infrastructures . Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019. Cardiff, 22nd July 2019, pages 9 – 16, Mannheim. Leibniz-Institut für Deutsche Sprache.	
739		
740		
741		
742		
743		
744		
745		
	Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1946–1958, Vancouver, Canada. Association for Computational Linguistics.	746
		747
		748
		749
		750
		751
		752
	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001, Florence, Italy. Association for Computational Linguistics.	753
		754
		755
		756
		757
		758
	Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning.	759
		760
		761
	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	762
		763
		764
	Afshin Rahimi, Yuan Li, and Trevor Cohn. 2019. Massively multilingual transfer for NER . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 151–164, Florence, Italy. Association for Computational Linguistics.	765
		766
		767
		768
		769
		770
	Ori Ram, Yuval Kirstain, Jonathan Berant, Amir Globerson, and Omer Levy. 2021. Few-shot question answering by pretraining span selection . 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 3066–3079, Online. Association for Computational Linguistics.	771
		772
		773
		774
		775
		776
		777
		778
		779
	Peter H Schönemann. 1966. A generalized solution of the orthogonal procrustes problem. <i>Psychometrika</i> , 31(1):1–10.	780
		781
		782
	Ke Tran. 2020. From english to foreign languages: Transferring pre-trained language models. <i>arXiv preprint arXiv:2002.07306</i> .	783
		784
		785
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	786
		787
		788
		789
		790
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi 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 M. 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.	791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802

803 Shijie Wu and Mark Dredze. 2020. [Are all languages](#)
 804 [created equal in multilingual BERT?](#) In *Proceedings*
 805 *of the 5th Workshop on Representation Learning for*
 806 *NLP*, pages 120–130, Online. Association for Com-
 807 *putational Linguistics.*

808 Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. 2015.
 809 [Normalized word embedding and orthogonal trans-](#)
 810 [form for bilingual word translation.](#) In *Proceedings*
 811 *of the 2015 Conference of the North American Chap-*
 812 *ter of the Association for Computational Linguistics:*
 813 *Human Language Technologies*, pages 1006–1011,
 814 Denver, Colorado. Association for Computational
 815 *Linguistics.*

816 Linting Xue, Noah Constant, Adam Roberts, Mi-
 817 hira Kale, Rami Al-Rfou, Aditya Siddhant, Aditya
 818 Barua, and Colin Raffel. 2021. [mT5: A massively](#)
 819 [multilingual pre-trained text-to-text transformer.](#) In
 820 *Proceedings of the 2021 Conference of the North*
 821 *American Chapter of the Association for Computa-*
 822 *tional Linguistics: Human Language Technologies,*
 823 *pages 483–498, Online. Association for Computa-*
 824 *tional Linguistics.*

825 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-
 826 bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.
 827 [Xlnet: Generalized autoregressive pretraining for](#)
 828 [language understanding.](#) In *Advances in Neural In-*
 829 *formation Processing Systems*, volume 32. Curran
 830 *Associates, Inc.*

831 Barret Zoph, Deniz Yuret, Jonathan May, and Kevin
 832 Knight. 2016. [Transfer learning for low-resource](#)
 833 [neural machine translation.](#) In *Proceedings of the*
 834 *2016 Conference on Empirical Methods in Natu-*
 835 *ral Language Processing*, pages 1568–1575, Austin,
 836 Texas. Association for Computational Linguistics.

837 A Grid search over k and τ

838 To choose number of neighbors k and temperature
 839 τ for WECHSEL we conduct a grid search over
 840 linear probes of models with different initializa-
 841 tion shown in Table 7. For RoBERTa, we compute
 842 scores on NLI (using XNLI) and POS tagging (us-
 843 ing the French, German and Chinese GSD corpora
 844 in Universal Dependencies) using linear probes of
 845 the last hidden state. We probe on NLI by taking
 846 a concatenation of the mean of all token represen-
 847 tations in the premise with the mean of all token
 848 representations in the hypothesis. We probe on
 849 POS tagging by taking the mean of all token rep-
 850 resentations belonging to each word. For GPT2,
 851 we compute language modelling Perplexity on the
 852 hold-out set also used to evaluate performance of
 853 the trained models.

854 B Hyperparameters

855 Hyperparameters used to fine-tune RoBERTa on
 856 downstream tasks are shown in Table 4. Hyperpa-

rameters used to train models in our main experi- 857
 858 ments are shown in Table 5.

Parameter	NLI	NER
peak learning rate	2e-5	2e-5
batch size	128	32
sequence length	128	128
Adam ϵ	1e-8	1e-8
Adam β_1	0.9	0.9
Adam β_2	0.999	0.999
train epochs	2	10
warmup	10% of steps	10% of steps
warmup schedule	linear	linear
LR decay	linear to zero	linear to zero

Table 4: Hyperparameters used to fine-tune RoBERTa models on NLI (XNLI) and NER (WikiANN).

Parameter	RoBERTa	GPT2
peak learning rate	1e-4	5e-4
batch size	512	512
sequence length	512	512
weight decay	0.01	0.01
Adam ϵ	1e-6	1e-6
Adam β_1	0.9	0.9
Adam β_2	0.98	0.98
train steps	250k	250k

Table 5: Hyperparameters of the models transferred from RoBERTa and GPT2.

859 C Qualitative subword correspondence

860 We show a small random sample of tokens in the
 861 target language and their closest English token (ac-
 862 cording to WECHSEL) in Table 6.

863 D Using fastText subword information

864 As an alternative to our tokenize-flatten-reduce pro-
 865 cedure (Figure 2) for mapping word embeddings
 866 to subword embeddings, one could also use fast-
 867 Text vectors to generate embeddings for out-of-
 868 vocabulary words by decomposing them into n-
 869 grams (Bojanowski et al., 2017) and treating sub-
 870 words as out-of-vocabulary words. We compare
 871 this option to WECHSEL at varying values for tem-
 872 perature τ and neighbors k (Table 8). We find that
 873 this does not improve performance compared to
 874 our initial method, so we choose not to use sub-
 875 word information so as to not restrict applicability
 876 of WECHSEL to word embeddings where subword
 877 information is available.

Lang	Target Token	Closest English Token
French	héritage	legacy
	trep	soaked
	évêque	bishop
	scandaleux	udicrous
	vertig	astounding
	enregistrer	rec
	sucrés	sweets
	Emmanuel	Emmanuel
	entourage	confid
secrétariat	ariat	
German	machen	ize
	mit	with
	Spruchwort	proverb
	erischen	Austrian
	minuten	utes
	Haustechnik	umbing
	dringen	urgent
	verfeinern	refine
	umgebung	vironments
	ternehmen	irms
	Chinese	到处
巧合		coinc
第三		third
杂交		recomb
利来		chnology
政务		Govern
石		stone
喊麦		sing
中海		iterranean
张某		defendant
Swahili	shirikishe	ive
	Harusi	Marriage
	pesile	ery
	tihani	graduate
	changi	ool
	kuugua	ingestion
	kuzidi	acclaim
	vipigo	Trouble
	dhamiri	conscience
aliposimama	Slowly	

Table 6: Samples of tokens in each language and the corresponding closest tokens from the English vocabulary according to WECHSEL.

E Relearning embeddings

We show performance throughout training of GPT-2 models in French and Chinese where all non-embedding parameters are frozen in Figure 5. Similarity of token embeddings after training to their initialization is shown in Figure 6.

Lang	Model	k	τ	Scores		
				NLI	POS	LM
French	WECHSEL@0	1	1	57.0	84.5	4.1e+5
		10	0.1	59.2	86.2	2.6e+5
		10	1	58.1	84.7	5.8e+5
		50	0.1	55.9	84.8	3.0e+6
		50	1	53.7	80.5	1.0e+7
		-	-	46.3	60.6	5.7e+6
German	WECHSEL@0	1	1	54.1	71.7	9.2e+5
		10	0.1	57.8	76.5	5.6e+5
		10	1	56.8	75.0	1.1e+6
		50	0.1	54.2	75.1	1.9e+7
		50	1	51.8	70.8	7.6e+7
		-	-	44.5	49.1	6.2e+6
Chinese	WECHSEL@0	1	1	46.2	69.9	4.7e+6
		10	0.1	49.0	76.2	2.8e+6
		10	1	49.0	75.0	3.2e+6
		50	0.1	46.6	72.8	1.9e+7
		50	1	46.7	71.9	2.9e+7
		-	-	37.5	53.7	5.8e+6
Chinese	BERT _{Base} -Chinese	-	-	63.2	93.6	-
		-	-	63.2	93.6	-
		-	-	63.2	93.6	-
		-	-	63.2	93.6	-
		-	-	63.2	93.6	-
		-	-	63.2	93.6	-

Table 7: Grid search over the temperature τ and number of most similar tokens k parameters of WECHSEL.

Lang	Model	k	τ	Scores		
				NLI	POS	LM
French	WECHSEL-FSI@0	1	1	58.4	85.2	2.5e+5
		10	0.1	59.8	86.8	2.0e+5
		10	1	58.3	84.4	4.8e+5
		50	0.1	57.2	83.6	3.1e+6
		50	1	54.0	81.6	1.8e+7
		-	-	55.8	72.7	6e+5
German	WECHSEL-FSI@0	1	1	58.9	76.0	4.2e+5
		10	1	57.5	75.4	8.3e+6
		50	0.1	55.4	75.4	1.0e+7
		50	1	53.6	69.5	5.9e+7
		-	-	47.4	75.4	2.7e+6
		-	-	48.0	80.7	2.6e+6
Chinese	WECHSEL-FSI@0	1	1	48.3	80.3	3.1e+6
		10	1	48.3	77.8	3.7e+7
		50	0.1	48.3	77.8	3.7e+7
		50	1	47.9	76.5	8.6e+7
		-	-	47.9	76.5	8.6e+7
		-	-	47.9	76.5	8.6e+7

Table 8: The same grid search as in Table 7, but using subword information from fastText vectors instead of a tokenize-flatten-reduce procedure to map words to subwords.

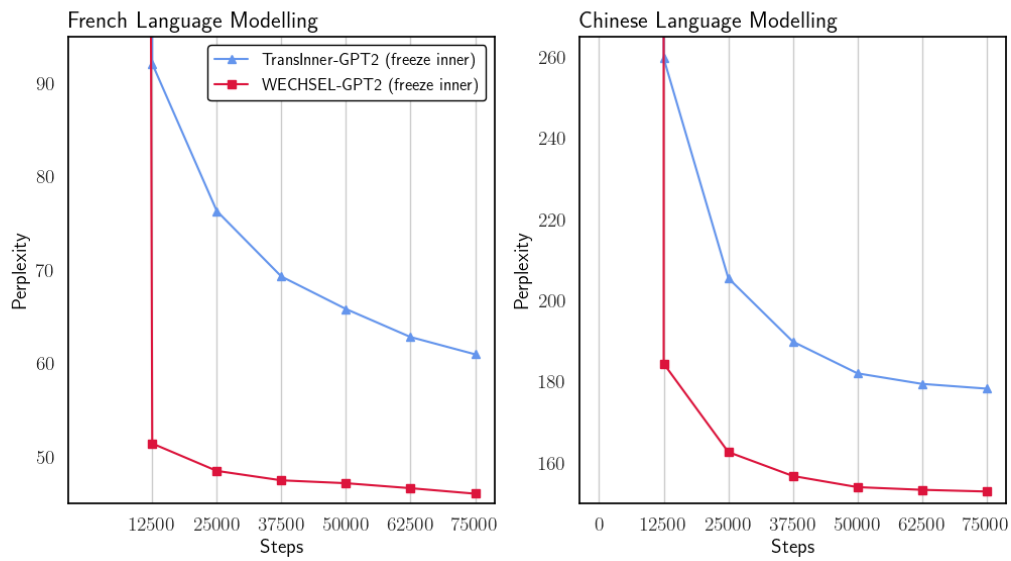


Figure 5: Perplexity over training steps from training French and Chinese GPT-2 models. We train models initialized with WECHSEL and with TransInner for 75k steps and freeze all non-embedding parameters for the entirety of training.

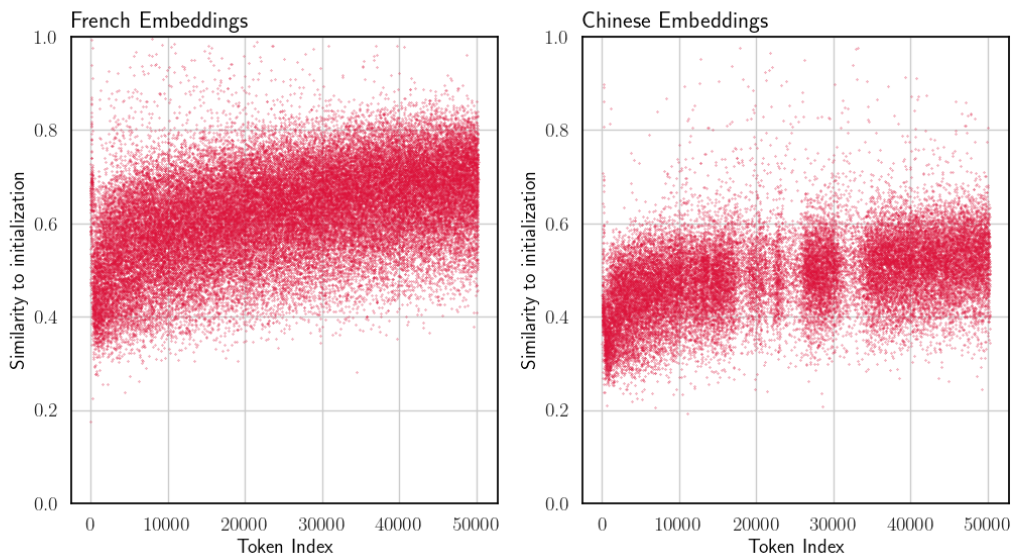


Figure 6: Cosine Similarity of embeddings to their initialization after training embeddings for 75k steps while freezing other parameters. We train French and German GPT-2 models initialized with WECHSEL.