

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

Anonymous ACL submission

## Abstract

Large pretrained language models (LMs) have become the central building block of many NLP applications. Training these models requires ever more computational resources and most of the existing models are trained on English text only. It is exceedingly expensive to train these models in other languages. To alleviate this problem, we introduce a novel method – called WECHSEL – to efficiently and effectively transfer pretrained LMs to new languages. WECHSEL can be applied to any model which uses subword-based tokenization and learns an embedding for each subword. The tokenizer of the source model (in English) is replaced with a tokenizer in the target language and token embeddings are initialized such that they are semantically similar to the English tokens by utilizing multilingual static word embeddings covering English and the target language. We use WECHSEL to transfer the English RoBERTa and GPT-2 models to four languages (French, German, Chinese and Swahili). We also study the benefits of our method on very low-resource languages. WECHSEL improves over proposed methods for cross-lingual parameter transfer and outperforms models of comparable size trained from scratch with up to 64x less training effort. Our method makes training large language models for new languages more accessible and less damaging to the environment. We make our code and models publicly available.

## 1 Introduction

Large LMs based on the Transformer architecture (Vaswani et al., 2017) have become increasingly popular since GPT (Radford et al., 2018) and BERT (Devlin et al., 2019) were introduced, prompting the creation of many large LMs pretrained on English text (Yang et al., 2019; Clark et al., 2020; Lewis et al., 2020; Joshi et al., 2020; Ram et al., 2021). There is a tendency towards training larger and larger models (Brown et al.,

2020; Fedus et al., 2021) while the main focus is on the English language. Recent work has called attention to the costs associated with training increasingly large LMs, including environmental and financial cost (Bender et al., 2021). If training large LMs for English is already costly, it is prohibitively expensive to train new, similarly powerful models to cover other languages.

One approach to address this issue is creating massively multilingual models (Devlin et al., 2019; Conneau et al., 2020; Xue et al., 2021) trained on a concatenation of texts in many different languages. These models show strong natural language understanding capabilities in a wide variety of languages, but suffer from what Conneau et al. (2020) call the *curse of multilinguality*: beyond a certain number of languages, overall performance decreases on monolingual as well as cross-lingual tasks. Consistent with this finding, Nozza et al. (2020) observe that monolingual LMs often outperform massively multilingual models. It is thus desirable to train monolingual models in more languages. Training monolingual models in non-English languages is commonly done by training a new model with randomly initialized parameters (Antoun et al., 2020; Louis, 2020; Chan et al., 2020; Martin et al., 2020). However, to train a model with capabilities comparable to that of an English model in this way, presumably a similar amount of compute to what was used to train the English model would be required.

To address this issue, we introduce WECHSEL<sup>1</sup>, a novel method to transfer monolingual language models to a new language. WECHSEL uses multilingual static word embeddings between the source language and the target language to initialize model parameters. WECHSEL first copies all inner (non-embedding) parameters of the English model, and exchanges the tokenizer with a tokenizer for the tar-

<sup>1</sup>Word Embeddings Can Help initialize Subword Embeddings in a new Language.

get language. Next, in contrast to prior work doing random initialization (de Vries and Nissim, 2021), the token embeddings in the target language are initialized such that they are close to semantically similar English tokens by mapping multilingual static word embeddings to subword embeddings. The latter step is particularly important considering that token 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 provides an effective initialization in the target language, requiring significantly fewer training steps to reach high performance than training from scratch. As multilingual static word embeddings are available for many languages (Bojanowski et al., 2017), WECHSEL is widely applicable.

We conduct our experiments on RoBERTa and GPT-2 as representative models of encoder and decoder language models, respectively. We transfer the English RoBERTa model to four languages (French, German, Chinese and Swahili), and the English GPT-2 model to the same four plus another four very low-resource languages (Sundanese, Scottish Gaelic, Uyghur and Malagasy). We evaluate the transferred RoBERTa models on Neural Entity Recognition (NER), and Natural Language Inference (NLI) tasks in the respective languages. The transferred GPT-2 models are evaluated in terms of Language Modelling Perplexity (PPL) on a held-out set. We compare WECHSEL with randomly initialized models (denoted as FullRand), as well as the recently proposed TransInner method which only transfers the inner (non-embedding) parameters (de Vries and Nissim, 2021). All mentioned models are trained under the same 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<sub>Base</sub> (Chan et al., 2020) (German), and BERT<sub>Base</sub>-Chinese (Devlin et al., 2019).

Results show that models initialized with WECHSEL outperform randomly initialized models and models initialized with TransInner across all languages and all tasks, for both RoBERTa and GPT-2. In addition, strong performance is reached

at a fraction of the training steps of other methods. 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 four and eight languages, respectively, achieved after minimal training effort.
- We release more effective GPT-2 and RoBERTa models than previously published non-English models, achieved under our more efficient training setting. Our code and models are publicly available at [github.com/anonymized](https://github.com/anonymized).

In the following, we review related work in Section 2. We 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 optimization objective is Masked Language Modelling (Devlin et al., 2019, MLM), where randomly selected tokens in the input are 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 recently proposed language models commonly consists of subwords (Clark et al., 2020; Liu et al., 2019; Devlin et al., 2019).

**Multilingual representations.** There has been a significant amount of work in creating multilingual static word embeddings. A common method is learning embeddings from scratch using data in multiple languages (Luong et al., 2015; Duong et al., 2016). Alternatively, multilinguality can be

180 achieved by aligning existing monolingual word  
181 embeddings using a bilingual dictionary, so that  
182 the resulting embeddings share the same seman-  
183 tic space (Xing et al., 2015; Joulin et al., 2018).  
184 Recent studies improve on this by reducing (or  
185 even completely removing) the need for bilingual  
186 data (Artetxe et al., 2017, 2018; Lample et al.,  
187 2018).

188 Beside static word embeddings, multilinguality  
189 is also well studied in the area of contextualized  
190 representations. One approach to learn multilingual  
191 contextualized representations is through training  
192 a model on a concatenation of corpora in differ-  
193 ent languages. Some models created based on  
194 this approach are mBERT (Devlin et al., 2019),  
195 XLM-R (Conneau et al., 2020) and mT5 (Xue  
196 et al., 2021), trained on text in 104, 100, and 101  
197 languages, respectively. As shown by Pires et al.  
198 (2019), a multilingual model such as mBERT can  
199 enable cross-lingual transfer by using task-specific  
200 annotations in one language to fine-tune the model  
201 for evaluation in another language. Despite the ben-  
202 efits, recent studies outline a number of limitations  
203 of massively multilingual LMs. Wu and Dredze  
204 (2020) empirically show that in mBERT “the 30%  
205 languages with least pretraining resources perform  
206 worse than using no pretrained language model at  
207 all”. Conneau et al. (2020) report that beyond a  
208 certain number of languages in the training data,  
209 the overall performance decreases on monolingual  
210 as well as cross-lingual tasks. These studies moti-  
211 vate our work on introducing an efficient approach  
212 for creating effective monolingual LMs for more  
213 languages.

214 **Cross-lingual transfer of monolingual LMs.**  
215 Studies in this area can be divided into two cat-  
216 egories:

- 217 • **Bilingualization of a monolingual LM** is  
218 concerned with extending a model to a new  
219 language while preserving its capabilities in  
220 the original language. Artetxe et al. (2020)  
221 approach this problem by replacing the to-  
222 kenizer and relearning the subword embed-  
223 dings, while freezing other (non-embedding)  
224 parameters. Such a model becomes bilingual,  
225 since the initial tokenizer and embeddings can  
226 be used for tasks in the source language, while  
227 the new tokenizer and embeddings can be used  
228 for tasks in the target language. Thus, a model  
229 can be finetuned on annotated task data in

230 the source language, and then zero-shot trans-  
231 ferred to the target language. Tran (2020)  
232 follow a similar approach, while instead of  
233 randomly initializing embeddings, they utilize  
234 static word embeddings to initialize embed-  
235 dings in the target language close to semanti-  
236 cally similar English tokens. They then contin-  
237 ue training the model on an English text  
238 corpus as well as on the target language in or-  
239 der to preserve model capabilities in English.

- 240 • **Creating a new monolingual LM in the tar-**  
241 **get language** is, in contrast, concerned with  
242 transferring a model from a source to a tar-  
243 get language without the necessity to preserve  
244 its capabilities in the source language. Zoph  
245 et al. (2016) and Nguyen and Chiang (2017)  
246 show that cross-lingually transferring a ma-  
247 chine translation model can improve perfor-  
248 mance, especially for low-resource languages.  
249 Zoph et al. (2016) use embeddings of random  
250 tokens in the original vocabulary to initial-  
251 ize token embeddings in the new vocabulary,  
252 while Nguyen and Chiang (2017) utilize vo-  
253 cabulary overlap between the source and tar-  
254 get language. More recently, de Vries and Nis-  
255 sim (2021) follow a similar approach to the  
256 one of Artetxe et al. (2020) for transferring a  
257 GPT-2 model to a new language. de Vries and  
258 Nissim (2021) add an additional step, where  
259 they train the entire model for some amount  
260 of steps to allow adapting to the target lan-  
261 guage beyond the lexical level. We refer to  
262 the method of de Vries and Nissim (2021) as  
263 TransInner and consider it as a baseline in our  
264 experiments.

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

### 3 Methodology

To initialize the model in the target language, we copy the inner (non-embedding) parameters from the source model. Our goal, then, is given the tokenizer  $T^s$  in the source language with vocabulary  $\mathbb{U}^s$ , the corresponding token embeddings  $\mathbf{E}^s$ , and a tokenizer  $T^t$  in the target language with vocabulary  $\mathbb{U}^t$ , to find a good initialization of the embeddings  $\mathbf{E}^t$  by using  $\mathbf{E}^s$ . To this end, we use existing bilingual word embeddings enriched with subword information, containing a set of words and subword n-grams in the source and target language and their aligned vectors. We denote the set of words and n-grams in the source and target language as  $\mathbb{V}^s$  and  $\mathbb{V}^t$  respectively, and the aligned static embeddings as  $\mathbf{W}^s$  and  $\mathbf{W}^t$ . In Appendix D we consider an alternative method if no subword information is available in the bilingual word embeddings.

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  $\mathbf{U}^s$  and  $\mathbf{U}^t$  for the source and target language, respectively. Next, we use  $\mathbf{U}^s$  and  $\mathbf{U}^t$  to compute the 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  $\mathbf{E}^t$  through a convex combination of embeddings in  $\mathbf{E}^s$  (Section 3.2). By applying WECHSEL, the vectors of  $\mathbf{E}^t$  are in the same semantic space as  $\mathbf{E}^s$ , 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 more detail in the following.

#### 3.1 Subword Embedding Computation

The process of mapping word embeddings to subword embeddings is done individually for the source and the target language. Given a tokenizer  $T$  with vocabulary  $\mathbb{U}$  and embeddings  $\mathbf{W}$ , the goal is to find subword embeddings  $\mathbf{U}$  for subwords in  $\mathbb{U}$  in the same semantic space as  $\mathbf{W}$ . To this end, we decompose subwords in  $\mathbb{U}$  into n-grams and compute the embedding by taking the sum of the embeddings of all occurring n-grams, equivalent to how embeddings for out-of-vocabulary words are computed in fastText (Bojanowski et al., 2017).

$$\mathbf{u}_x = \sum_{g \in \mathbb{G}^{(x)}} \mathbf{w}_g$$

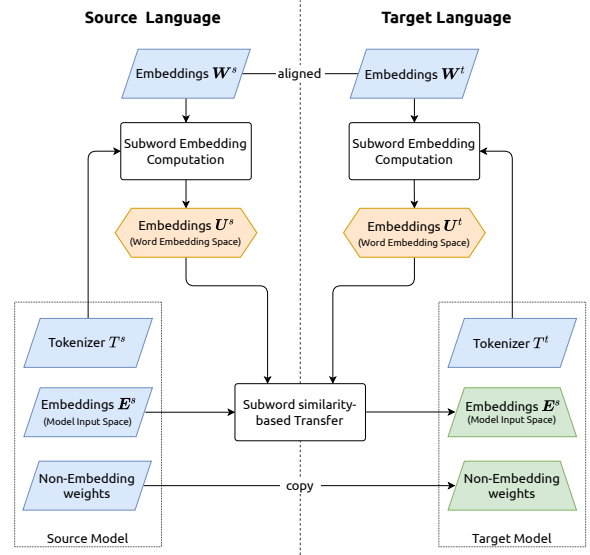


Figure 1: Summary of our WECHSEL method. We show **inputs**, **intermediate results** and **outputs**.

where  $\mathbb{G}^{(x)}$  is the set of n-grams occurring in the subword  $x$  and  $\mathbf{w}_g$  is the embedding of the n-gram  $g$ . Subwords in which no known n-gram occurs are initialized to zero.

#### 3.2 Subword similarity-based Transfer

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

$$s_{x,y} = \frac{\mathbf{u}_x^t \mathbf{u}_y^{sT}}{\|\mathbf{u}_x^t\| \|\mathbf{u}_y^s\|}$$

We now exploit these similarities to initialize embeddings in  $\mathbf{E}^t$  by a convex combination of embeddings in  $\mathbf{E}^s$ . In particular, each subword embedding in  $\mathbf{E}^t$  is defined as the weighted mean of the  $k$  nearest embeddings in  $\mathbf{E}^s$  according to the similarity values. The weighting is done by a softmax of the similarities with temperature  $\tau$ .

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

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

## 4 Experiment Design

We evaluate our method by transferring the English RoBERTa (Liu et al., 2019) and the English GPT-2 model (Radford et al., 2019) to French, German, Chinese and Swahili. We refer to these languages as *medium-resource languages*. In addition, we study the benefits of our method on four *low-resource languages*, namely Sundanese, Scottish Gaelic, Uyghur and Malagasy.

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. GPT-2 is evaluated by Perplexity (PPL) on a held-out set from the same corpus on which the model was trained on.

Due to the difficulty of extrinsic evaluation on low-resource languages, we only train GPT-2 models in these languages, and evaluate their performance intrinsically via Language Modelling Perplexity on a held-out set. We use the pretrained models RoBERTa<sub>Base</sub> 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 training data in the target language, we use a subset of 4GiB from the OSCAR corpus (Ortiz Suárez et al., 2019) for German, French and Chinese. For the other languages, we use data from the CC-100 corpus (Conneau et al., 2020) which contains 1.6GiB, 0.1GiB, 0.1GiB, 0.4GiB and 0.2GiB for Swahili, Sundanese, Scottish Gaelic, Uyghur and Malagasy, respectively. To obtain aligned word embeddings between the source and the target language we use monolingual fastText word embeddings<sup>2</sup> (Bojanowski et al., 2017). We align these embeddings using the Orthogonal Procrustes method (Schönemann, 1966; Artetxe et al., 2016) with bilingual dictionaries from MUSE<sup>3</sup> (Conneau et al., 2017) for French, German and Chinese and a bilingual dictionary from FreeDict<sup>4</sup> (Bański and Wójtowicz, 2009) for Swahili. For the low-resource languages, we use bilingual dictionaries scraped from Wiktionary.<sup>5</sup>

<sup>2</sup><https://fasttext.cc>

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

<sup>4</sup><https://freedict.org>

<sup>5</sup>available at [github.com/anonymized](https://github.com/anonymized)

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 <sub>Base</sub>	255.6B	3.9x
BERT <sub>Base</sub> -Chinese	131.1B	2.0x

Table 1: Tokens trained on in the target language between our models and previous monolingual models.

We choose temperature  $\tau = 0.1$  and neighbors  $k = 10$  for WECHSEL by conducting a parameter search over a grid with varying values for  $k$  and  $\tau$  using linear probes (Appendix A). 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.

- **TransInner:** 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. In preliminary experiments reported in Appendix E, we compare the method by Zoph et al. (2016) with TransInner, observing superior performance of TransInner, so we choose TransInner as the baseline for cross-lingual transfer in all our experiments.
- **FullRand:** 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).

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 the peak LR for the first 10% of steps, followed by a 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 model.

In addition to the mentioned baselines trained under this setting, we compare the results of

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>86.93</u>	<u>82.59</u>	<u>81.63</u>	<u>90.26</u>	<u>85.95</u>	<b>82.43</b>	<b>90.88</b>	<b>86.65</b>	-	-	-
	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 <sub>Base</sub>	-	-	-	-	-	-	-	-	-	79.25	89.48	84.37
German	WECHSEL-RoBERTa	<u>75.64</u>	<u>84.53</u>	<u>80.08</u>	<u>81.11</u>	<u>89.05</u>	<u>85.08</u>	<b>81.79</b>	<b>89.72</b>	<b>85.76</b>	-	-	-
	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 <sub>Base</sub>	-	-	-	-	-	-	-	-	-	78.64	89.46	84.05
	XLM-R <sub>Base</sub>	-	-	-	-	-	-	-	-	-	78.58	88.76	83.67
Chinese	WECHSEL-RoBERTa	<u>63.23</u>	<u>72.79</u>	<u>68.01</u>	<u>77.19</u>	<u>79.07</u>	<u>78.13</u>	<b>78.32</b>	<b>80.55</b>	<b>79.44</b>	-	-	-
	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 <sub>Base</sub> -Chinese	-	-	-	-	-	-	-	-	-	76.55	<b>82.05</b>	79.30
	XLM-R <sub>Base</sub>	-	-	-	-	-	-	-	-	-	76.41	78.36	77.38
Swahili	WECHSEL-RoBERTa	<u>60.28</u>	<u>74.38</u>	<u>67.33</u>	<u>73.87</u>	<u>87.63</u>	<u>80.75</u>	<b>75.05</b>	<b>87.39</b>	<b>81.22</b>	-	-	-
	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 <sub>Base</sub>	-	-	-	-	-	-	-	-	-	69.18	87.37	78.28

Table 2: 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)**), all trained on more tokens than our models. For each language, the best results in every column are indicated with underlines. The overall best results including the comparison with existing monolingual/multilingual models of comparable size are shown in bold.

442 RoBERTa models with existing comparable models  
443 trained from scratch with more training effort.  
444 We consider the total number of tokens the  
445 model has encountered in the target language,  
446 computed as the product of batch size  $\times$  sequence  
447 length  $\times$  train steps (shown in Table 1) as a proxy  
448 for training effort. We evaluate the performance  
449 of CamemBERT (Martin et al., 2020) (French),  
450 GBERT<sub>Base</sub> (Chan et al., 2020) (German), and  
451 BERT<sub>Base</sub>-Chinese (Devlin et al., 2019) as existing  
452 monolingual LMs,<sup>6</sup> as well as XLM-R<sub>Base</sub> (Artetxe  
453 et al., 2020) as a high-performing multilingual LM.

## 5 Results

454 We present our results on transferring RoBERTa  
455 and GPT-2 from English to other languages, fol-  
456 lowed by analyzing training behavior. In Ap-  
457 pendix C, we provide a qualitative assessment of  
458 how well subword tokens are mapped between the  
459 source and the target languages.

### 5.1 Transferring RoBERTa

460 Table 2 reports the evaluation results of RoBERTa.  
461 As shown, models initialized with WECHSEL out-  
462 perform models trained from scratch and models  
463 initialized with TransInner across all languages.  
464 Surprisingly, close relatedness of the source and  
465

<sup>6</sup>To the best of our knowledge there is no monolingual model available for Swahili.

467 target language is not necessary to achieve effective  
468 transfer, as e. g. on NLI WECHSEL improves abso-  
469 lute accuracy by 7.15%, 6.31%, 6.94% and 4.71%  
470 over models trained from scratch for French, Ger-  
471 man, Chinese and Swahili, respectively.

472 We observe that our parameter transfer-based  
473 model consistently outperforms the previously re-  
474 leased LMs on both monolingual and multilingual  
475 settings, while these models benefit from much  
476 larger training resources in terms of computation  
477 time and corpus size. In particular, the results  
478 show an improvement over XLM-R<sub>Base</sub> by an av-  
479 erage 3.54% accuracy for NLI and 1.14% micro  
480 F1 score for NER. For NLI, we improve over the  
481 prior monolingual models by 1.55%, 3.15% and  
482 1.77% absolute accuracy for French, German and  
483 Chinese, respectively. For NER, we observe im-  
484 provements over monolingual models with 0.62%  
485 and 0.26% absolute micro F1 score improvement  
486 for French and German, respectively. For Chinese,  
487 the monolingual model BERT<sub>Base</sub>-Chinese still out-  
488 performs our method by 1.5% absolute micro F1  
489 score. We suspect that the discrepancy between  
490 NLI and NER is due to the limited training cor-  
491 pus size (max. 4GiB), while a larger corpus can  
492 potentially improve NER as more named entities  
493 appear (Martin et al., 2020).

494 The first two columns of Figure 2 show the  
495 performance of RoBERTa models on downstream

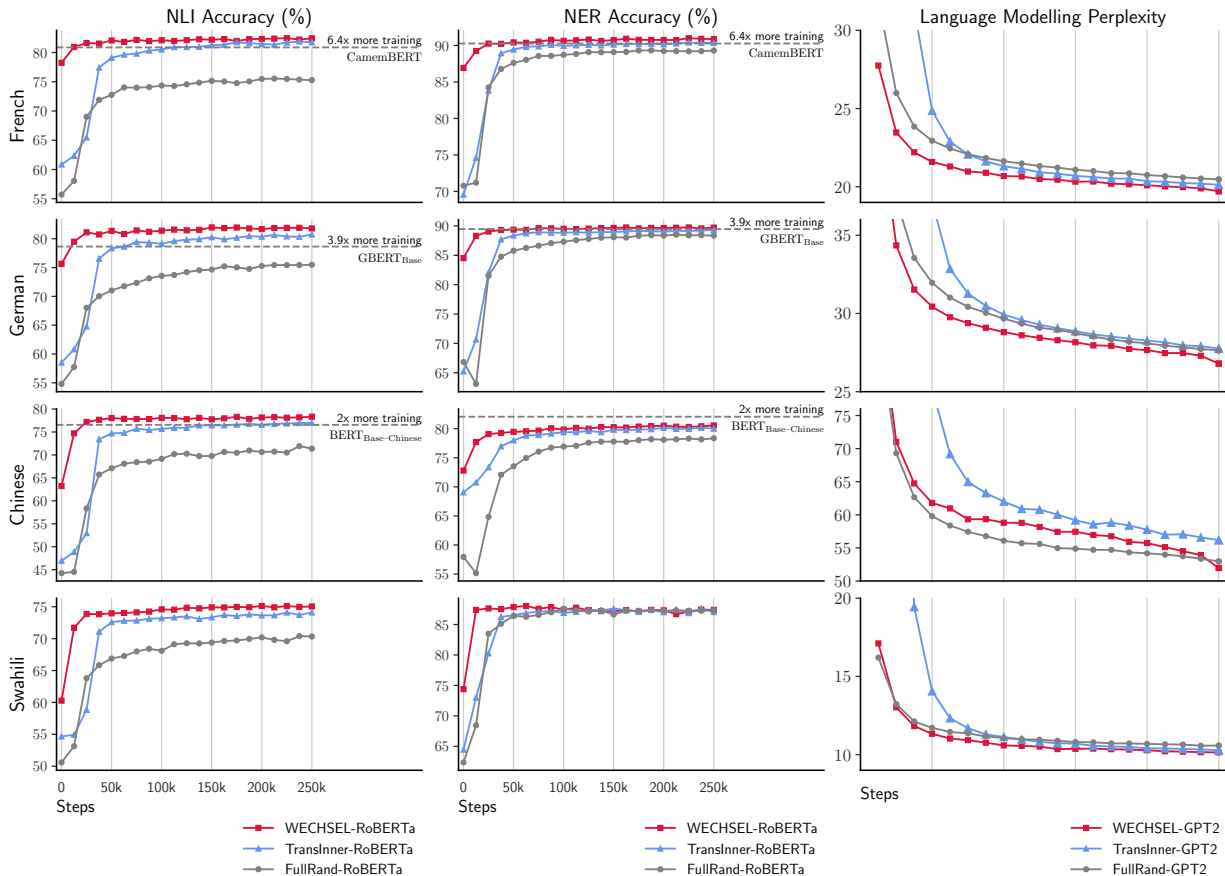


Figure 2: Test scores over training steps from fine-tuning RoBERTa models on NLI (using XNLI) and NER (using WikiANN). Perplexity on the held-out set over training steps of GPT-2 models. We evaluate every 12.5k steps.

Lang	Model	PPL@0	PPL@25k	PPL@250k
French	WECHSEL-GPT2	1.7e+3	23.47	<b>19.71</b>
	TransInner-GPT2	1.4e+5	67.97	20.13
	FullRand-GPT2	5.9e+4	25.99	20.47
German	WECHSEL-GPT2	3.7e+3	34.35	<b>26.80</b>
	TransInner-GPT2	1.5e+5	121.67	27.76
	FullRand-GPT2	5.8e+4	37.29	27.63
Chinese	WECHSEL-GPT2	2.4e+4	71.02	<b>51.97</b>
	TransInner-GPT2	1.5e+5	231.05	56.17
	FullRand-GPT2	5.8e+4	<u>69.29</u>	52.98
Swahili	WECHSEL-GPT2	1.4e+5	<b>13.02</b>	<b>10.14</b>
	TransInner-GPT2	1.4e+5	42.95	10.28
	FullRand-GPT2	<u>5.8e+4</u>	13.22	10.58

Table 3: Results of training GPT2 models. We report Perplexity before training (**PPL@0**), after 10% of steps (**PPL@25k**) and after training (**PPL@250k**).

Lang	Model	Best PPL
Sundanese	WECHSEL-GPT2	<b>111.72</b>
	TransInner-GPT2	151.86
	FullRand-GPT2	149.46
Scottish Gaelic	WECHSEL-GPT2	<b>16.43</b>
	TransInner-GPT2	18.62
	FullRand-GPT2	19.53
Uyghur	WECHSEL-GPT2	<b>34.33</b>
	TransInner-GPT2	39.06
	FullRand-GPT2	42.82
Malagasy	WECHSEL-GPT2	<b>14.01</b>
	TransInner-GPT2	14.85
	FullRand-GPT2	15.93

Table 4: Results of training GPT2 models on low-resource languages. We report the best Perplexity on the held-out set, evaluated every 2.5k steps.

tasks after each 12.5k training steps. Models initialized with WECHSEL reach high performance in significantly fewer steps than models initialized with FullRand or TransInner.

We expect FullRand-RoBERTa to approach performance of the respective prior monolingual mod-

els when trained on the same amount of tokens<sup>7</sup>. For French, WECHSEL-RoBERTa outperforms CamemBERT after 10% of training steps, reducing training effort by 64x. For German, WECHSEL-

<sup>7</sup>It would presumably be slightly worse because we restrict training corpus size to 4GiB.

506 RoBERTa outperforms GBERT<sub>Base</sub> after 10% of  
 507 training steps, reducing training effort by 39x.  
 508 For Chinese, WECHSEL-RoBERTa outperforms  
 509 BERT<sub>Base</sub>-Chinese on NLI, but does not outper-  
 510 form BERT<sub>Base</sub>-Chinese on NER.

## 5.2 Transferring GPT-2

### 5.2.1 To Medium-Resource Languages

513 Results on medium-resource languages are shown  
 514 in Table 3. Similar to the results for WECHSEL-  
 515 RoBERTa, the GPT-2 models trained with WECH-  
 516 SEL consistently outperform the models trained  
 517 from scratch and with TransInner across all lan-  
 518 guages.

519 The last column of Figure 2 depicts the perfor-  
 520 mance of GPT-2 models after each 12.5k train-  
 521 ing steps. Comparing the results across all lan-  
 522 guages throughout training, we observe a stronger  
 523 dependence on similarity of the source to the tar-  
 524 get language than for downstream tasks such as  
 525 NLI or NER. In particular, for French and German,  
 526 WECHSEL is consistently better than TransInner  
 527 and FullRand throughout the entire training, while  
 528 for Chinese, a decrease in perplexity towards the  
 529 end of training causes WECHSEL to surpass train-  
 530 ing from scratch.

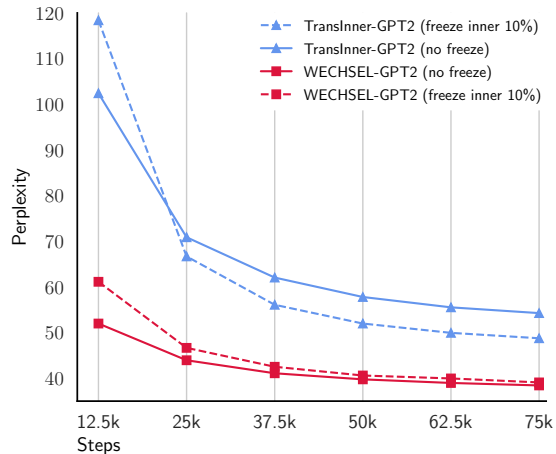
### 5.2.2 To Low-Resource Languages

532 Table 4 reports the perplexity of Language Mod-  
 533 elling on the low-resource languages. Again, we  
 534 observe consistent improvements using WECHSEL  
 535 on all languages. Furthermore, as discussed in Ap-  
 536 pendix G, we conduct a sensitivity analysis w. r. t.  
 537 the amount of available training data on French,  
 538 studying the relation of performance improvement  
 539 with training data size.

540 One difference of the low-resource models with  
 541 the ones trained on medium-resource languages is  
 542 that the low-resource LMs are prone to overfitting,  
 543 and require appropriate model selection even in  
 544 the early steps of training. Appendix F further  
 545 elaborates on this by showing the performance of  
 546 the low-resource LMs throughout training.

### 5.3 Is freezing necessary?

548 Previous work using the TransInner method freezes  
 549 non-embedding parameters for a fixed amount of  
 550 steps before training the entire model (de Vries  
 551 and Nissim, 2021). This is done to prevent cata-  
 552 strophic forgetting at the beginning of training. To  
 553 evaluate if freezing non-embedding parameters is  
 554 still necessary with our method, we conduct an



555 Figure 3: Comparison of German GPT-2 models  
 556 trained with WECHSEL and TransInner between freez-  
 557 ing non-embedding parameters at the start and not  
 558 freezing any parameters.  
 559  
 560  
 561  
 562  
 563  
 564

555 additional experiment. We train a German GPT-2  
 556 model with WECHSEL and a model with TransIn-  
 557 ner without freezing any parameters, and the same  
 558 models with freezing of non-embedding parameters  
 559 for the first 10% of steps. We match hyperparam-  
 560 eters of the main experiments except training for 75k  
 561 steps only. Based on the results shown in Figure 3,  
 562 we conclude that freezing is necessary when using  
 563 TransInner, but there is no need for freezing when  
 564 using WECHSEL.

## 6 Conclusion

566 We introduce WECHSEL, an effective method to  
 567 transfer monolingual language models to new lan-  
 568 guages. WECHSEL exploits multilingual static  
 569 word embeddings to compute an effective initializa-  
 570 tion of subword embeddings in the target language.  
 571 We conduct experiments by transferring RoBERTa  
 572 and GPT-2 models from English to French, Ger-  
 573 man, Chinese and Swahili, as well as English GPT-  
 574 2 to four low-resource languages. The evaluation  
 575 results show that the transferred RoBERTa and  
 576 GPT-2 models are more efficient and effective than  
 577 strong baselines, and consistently outperform prior  
 578 monolingual models that have been trained for a  
 579 significantly longer time. WECHSEL facilitates  
 580 the creation of effective monolingual LMs for new  
 581 languages with medium to low resources, particu-  
 582 larly in computationally-limited settings. In addi-  
 583 tion, our work provides strong evidence towards  
 584 the hypothesis by Artetxe et al. (2020) that deep  
 585 monolingual language models learn abstractions  
 586 that generalize across languages. We discuss limi-  
 587 tations and risks of our work in Appendix H.



588  
589  
590  
591  
592  
593  
594  
595  
  
596  
597  
598  
599  
600  
601  
602  
  
603  
604  
605  
606  
607  
608  
609  
  
610  
611  
612  
613  
614  
615  
616  
  
617  
618  
619  
620  
621  
622  
  
623  
624  
625  
  
626  
627  
628  
629  
630  
631  
  
632  
633  
634  
635  
  
636  
637  
638  
639  
640  
  
641  
642  
643  
644

## 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](#). 645  
646

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*. 647  
648  
649  
650

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. 651  
652  
653  
654  
655  
656  
657  
658  
659

Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. [Word translation without parallel data](#). *arXiv preprint arXiv:1710.04087*. 660  
661  
662  
663

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. 664  
665  
666  
667  
668  
669  
670

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. 671  
672  
673  
674  
675  
676

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. 677  
678  
679  
680  
681  
682  
683  
684  
685

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. 686  
687  
688  
689  
690  
691  
692

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*. 693  
694  
695  
696

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. 697  
698  
699  
700  
701

702	Armand Joulin, Piotr Bojanowski, Tomas Mikolov,	Xiaoman Pan, Boliang Zhang, Jonathan May, Joel	758
703	Hervé Jégou, and Edouard Grave. 2018. <a href="#">Loss in</a>	Nothman, Kevin Knight, and Heng Ji. 2017. <a href="#">Cross-</a>	759
704	<a href="#">translation: Learning bilingual word mapping with</a>	<a href="#">lingual name tagging and linking for 282 languages.</a>	760
705	<a href="#">a retrieval criterion.</a> In <i>Proceedings of the 2018</i>	In <i>Proceedings of the 55th Annual Meeting of the</i>	761
706	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>Association for Computational Linguistics (Volume</i>	762
707	<i>guage Processing</i> , pages 2979–2984, Brussels, Bel-	<i>1: Long Papers)</i> , pages 1946–1958, Vancouver,	763
708	gium. Association for Computational Linguistics.	Canada. Association for Computational Linguistics.	764
709	Guillaume Lample, Alexis Conneau, Marc’Aurelio	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019.	765
710	Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018.	<a href="#">How multilingual is multilingual BERT?</a> In <i>Pro-</i>	766
711	<a href="#">Word translation without parallel data.</a> In <i>Internat-</i>	<i>ceedings of the 57th Annual Meeting of the Asso-</i>	767
712	<i>ional Conference on Learning Representations.</i>	<i>ciation for Computational Linguistics</i> , pages 4996–	768
713	Mike Lewis, Marjan Ghazvininejad, Gargi Ghosh, Ar-	5001, Florence, Italy. Association for Computa-	769
714	men Aghajanyan, Sida Wang, and Luke Zettlemoyer.	tional Linguistics.	770
715	2020. <a href="#">Pre-training via paraphrasing.</a> In <i>Advances in</i>	Alec Radford, Karthik Narasimhan, Tim Salimans, and	771
716	<i>Neural Information Processing Systems</i> , volume 33,	Ilya Sutskever. 2018. Improving language under-	772
717	pages 18470–18481. Curran Associates, Inc.	standing with unsupervised learning.	773
718	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Alec Radford, Jeff Wu, Rewon Child, David Luan,	774
719	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Dario Amodei, and Ilya Sutskever. 2019. Language	775
720	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	models are unsupervised multitask learners.	776
721	Roberta: A robustly optimized bert pretraining ap-	Afshin Rahimi, Yuan Li, and Trevor Cohn. 2019. <a href="#">Mas-</a>	777
722	proach. <i>arXiv preprint arXiv:1907.11692.</i>	<a href="#">sively multilingual transfer for NER.</a> In <i>Proceed-</i>	778
723	Antoine Louis. 2020. BelGPT-2: a GPT-2 model pre-	<i>dings of the 57th Annual Meeting of the Association</i>	779
724	trained on French corpora. <a href="https://github.com/antoiloui/belgpt2">https://github.</a>	<i>for Computational Linguistics</i> , pages 151–164, Flo-	780
725	<a href="https://github.com/antoiloui/belgpt2">com/antoiloui/belgpt2.</a>	rence, Italy. Association for Computational Linguis-	781
726	Thang Luong, Hieu Pham, and Christopher D. Man-	tics.	782
727	ning. 2015. <a href="#">Bilingual word representations with</a>	Ori Ram, Yuval Kirstain, Jonathan Berant, Amir	783
728	<a href="#">monolingual quality in mind.</a> In <i>Proceedings of the</i>	Globerson, and Omer Levy. 2021. <a href="#">Few-shot ques-</a>	784
729	<i>1st Workshop on Vector Space Modeling for Natural</i>	<a href="#">tion answering by pretraining span selection.</a> In <i>Pro-</i>	785
730	<i>Language Processing</i> , pages 151–159, Denver, Col-	<i>ceedings of the 59th Annual Meeting of the Associa-</i>	786
731	orado. Association for Computational Linguistics.	<i>tion for Computational Linguistics and the 11th In-</i>	787
732	Louis Martin, Benjamin Muller, Pedro Javier Or-	<i>ternational Joint Conference on Natural Language</i>	788
733	tiz Suárez, Yoann Dupont, Laurent Romary, Éric	<i>Processing (Volume 1: Long Papers)</i> , pages 3066–	789
734	de la Clergerie, Djamel Seddah, and Benoît Sagot.	3079, Online. Association for Computational Lin-	790
735	2020. <a href="#">CamemBERT: a tasty French language model.</a>	guistics.	791
736	In <i>Proceedings of the 58th Annual Meeting of the</i>	Peter H Schönemann. 1966. A generalized solution of	792
737	<i>Association for Computational Linguistics</i> , pages	the orthogonal procrustes problem. <i>Psychometrika</i> ,	793
738	7203–7219, Online. Association for Computational	31(1):1–10.	794
739	Linguistics.	Ke Tran. 2020. From english to foreign languages:	795
740	Toan Q. Nguyen and David Chiang. 2017. <a href="#">Trans-</a>	Transferring pre-trained language models. <i>arXiv</i>	796
741	<a href="#">fer learning across low-resource, related languages</a>	<i>preprint arXiv:2002.07306.</i>	797
742	<a href="#">for neural machine translation.</a> In <i>Proceedings of</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	798
743	<i>the Eighth International Joint Conference on Natu-</i>	Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz	799
744	<i>ral Language Processing (Volume 2: Short Papers)</i> ,	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	800
745	pages 296–301, Taipei, Taiwan. Asian Federation of	<a href="#">you need.</a> In <i>Advances in Neural Information Pro-</i>	801
746	Natural Language Processing.	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	802
747	Debora Nozza, Federico Bianchi, and Dirk Hovy. 2020.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	803
748	What the [mask]? making sense of language-specific	Chaumond, Clement Delangue, Anthony Moi, Pier-	804
749	bert models. <i>arXiv preprint arXiv:2003.02912.</i>	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtow-	805
750	Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	806
751	Romary. 2019. <a href="#">Asynchronous pipelines for pro-</a>	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	807
752	<a href="#">cessing huge corpora on medium to low resource</a>	Teven Le Scao, Sylvain Gugger, Mariama Drame,	808
753	<a href="#">infrastructures.</a> Proceedings of the Workshop on	Quentin Lhoest, and Alexander M. Rush. 2020.	809
754	Challenges in the Management of Large Corpora	<a href="#">Transformers: State-of-the-art natural language pro-</a>	810
755	(CMLC-7) 2019. Cardiff, 22nd July 2019, pages	<a href="#">cessing.</a> In <i>Proceedings of the 2020 Conference on</i>	811
756	9 – 16, Mannheim. Leibniz-Institut für Deutsche	<i>Empirical Methods in Natural Language Processing:</i>	812
757	Sprache.	<i>System Demonstrations</i> , pages 38–45, Online. Asso-	813
		ciation for Computational Linguistics.	814

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

820 Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. 2015.  
 821 [Normalized word embedding and orthogonal trans-](#)  
 822 [form for bilingual word translation.](#) In *Proceedings*  
 823 *of the 2015 Conference of the North American Chap-*  
 824 *ter of the Association for Computational Linguistics:*  
 825 *Human Language Technologies*, pages 1006–1011,  
 826 Denver, Colorado. Association for Computational  
 827 *Linguistics.*

828 Linting Xue, Noah Constant, Adam Roberts, Mi-  
 829 hir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya  
 830 Barua, and Colin Raffel. 2021. [mT5: A massively](#)  
 831 [multilingual pre-trained text-to-text transformer.](#) In  
 832 *Proceedings of the 2021 Conference of the North*  
 833 *American Chapter of the Association for Computa-*  
 834 *tional Linguistics: Human Language Technologies,*  
 835 *pages 483–498, Online. Association for Computa-*  
 836 *tional Linguistics.*

837 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-  
 838 bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.  
 839 [Xlnet: Generalized autoregressive pretraining for](#)  
 840 [language understanding.](#) In *Advances in Neural In-*  
 841 *formation Processing Systems*, volume 32. Curran  
 842 *Associates, Inc.*

843 Barret Zoph, Deniz Yuret, Jonathan May, and Kevin  
 844 Knight. 2016. [Transfer learning for low-resource](#)  
 845 [neural machine translation.](#) In *Proceedings of the*  
 846 *2016 Conference on Empirical Methods in Natu-*  
 847 *ral Language Processing*, pages 1568–1575, Austin,  
 848 Texas. Association for Computational Linguistics.

## 849 A Grid search over $k$ and $\tau$

850 To choose number of neighbors  $k$  and temperature  
 851  $\tau$  for WECHSEL we conduct a grid search over  
 852 linear probes of models with different initializa-  
 853 tion shown in Table 8. For RoBERTa, we compute  
 854 scores on NLI (using XNLI) and POS tagging (us-  
 855 ing the French, German and Chinese GSD corpora  
 856 in Universal Dependencies) using linear probes of  
 857 the last hidden state. We probe on NLI by taking  
 858 a concatenation of the mean of all token represen-  
 859 tations in the premise with the mean of all token  
 860 representations in the hypothesis. We probe on  
 861 POS tagging by taking the mean of all token rep-  
 862 resentations belonging to each word. For GPT2,  
 863 we compute Language Modelling Perplexity on the  
 864 held-out set also used to evaluate performance of  
 865 the trained models.

## 866 B Hyperparameters

867 Hyperparameters used to fine-tune RoBERTa on  
 868 downstream tasks are shown in Table 5. Hyperpa-

rameters used to train models in our main experi- 869  
 870 ments are shown in Table 6.

Parameter	NLI	NER
peak learning rate	2e-5	2e-5
batch size	128	32
sequence length	128	128
Adam $\epsilon$	1e-8	1e-8
Adam $\beta_1$	0.9	0.9
Adam $\beta_2$	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 5: 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 $\epsilon$	1e-6	1e-6
Adam $\beta_1$	0.9	0.9
Adam $\beta_2$	0.98	0.98
train steps	250k	250k

Table 6: Hyperparameters of the models transferred from RoBERTa and GPT2.

## 871 C Qualitative subword correspondence

872 We show a small random sample of tokens in the  
 873 target language and their closest English token (ac-  
 874 cording to WECHSEL) in Table 7.

## 875 D Using Word Embeddings without 876 subword information

877 As an alternative to n-gram decomposition, we in-  
 878 troduce a method for mapping word embeddings to  
 879 subword embeddings without using any subword  
 880 information (shown in Figure 4). For this method,  
 881 we require word frequency information in addition  
 882 to the word embeddings. We apply the tokenizer  $T$   
 883 to every word  $v$  in  $\mathbb{V}$  resulting in a set of subwords  
 884 for each word. We define  $\mathbb{V}^{(x)}$  as the set of words  
 885 containing the subword  $x$  when tokenized. The  
 886 embedding  $\mathbf{u}_x$  of the subword  $x$  is then defined as  
 887 the average of the embeddings of words in  $\mathbb{V}^{(x)}$ ,  
 888 weighted by the word frequencies.

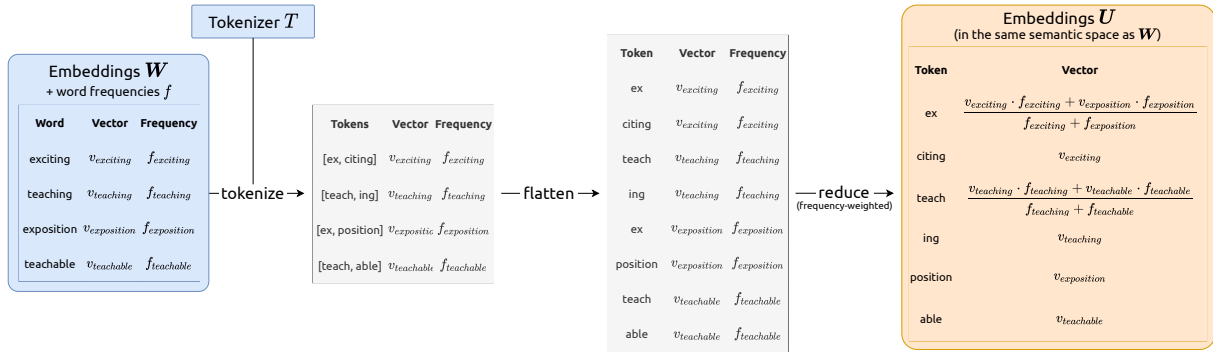


Figure 4: WECHSEL<sub>TFR</sub>, an alternative subword embedding computation method. First, **tokenize** all words in the word embeddings. Then **flatten** the result by assigning the embeddings of the words in which it occurred and their word frequencies to each subword. Finally, **reduce** the embeddings assigned to each subword by taking their mean, weighted by word frequency.

$$\mathbf{u}_x = \frac{\sum_{v \in \mathbb{V}(x)} \mathbf{w}_v \cdot f_v}{\sum_{v \in \mathbb{V}(x)} f_v}$$

where  $\mathbf{w}_v$  is the embedding and  $f_v$  is the frequency of word  $v$ . We call this variant of our method WECHSEL<sub>TFR</sub>. We evaluate WECHSEL<sub>TFR</sub> by training the same models as for WECHSEL. Results are shown in Table 9 for GPT2 and in Table 10 for RoBERTa. We find that, on average, performance is on par with WECHSEL.

## E Choosing a transfer baseline

We consider two baseline methods to transfer models to a new language without using any language-specific information. One method is copying non-embedding parameters to the target language and initializing embeddings from a random normal distribution as done by de Vries and Nissim (2021). We refer to this method as TransInner. Another option is copying non-embedding parameters and assigning the embedding of a random token in the source language to each embedding in the target language (effectively "shuffling" the embeddings) as done by Zoph et al. (2016) and Nguyen and Chiang (2017). We refer to this method as TransInnerShuffleEmb. We evaluate these two methods using a setup equivalent to the experiments in Section 5.3 and find that TransInner performs slightly better than TransInnerShuffleEmb (Figure 6), so we use TransInner for subsequent experiments.

## F Performance throughout training on low-resource languages

We show Language Modelling Perplexity on the held-out set throughout training in Figure 5.

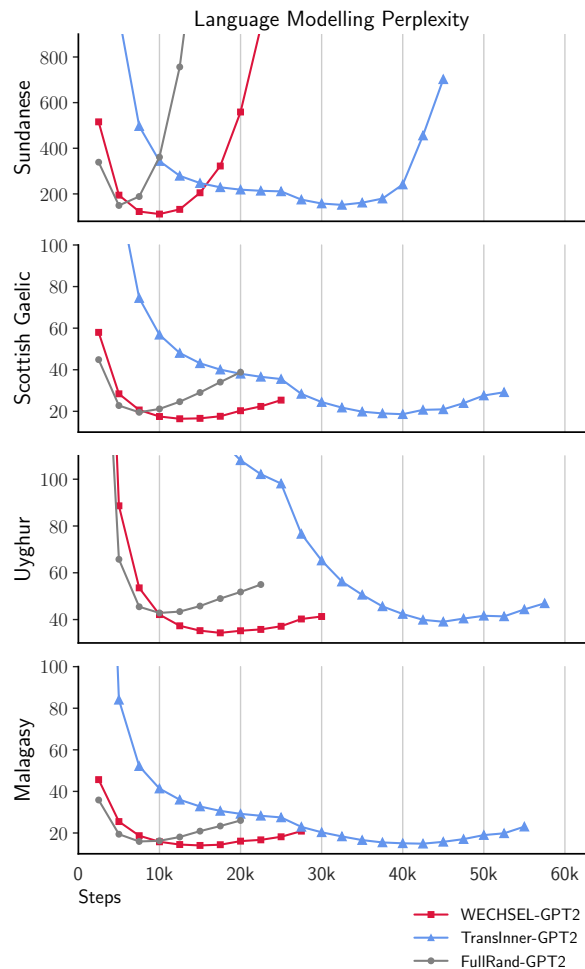


Figure 5: Perplexity throughout training on low-resource languages. We evaluate every 2.5k steps and stop training if Perplexity on the held-out set does not improve for 10k steps. TransInner-GPT2 takes more steps to overfit since all non-embedding parameters are frozen for the first 25k steps (c. f. Section 4).

Lang	Target Token	Closest English Token
French	héritage	legacy
	trempe	soaked
	épiscopat	bishop
	scandaleux	udicrous
	vertigineux	astounding
	enregistrer	rec
	sucrés	sweets
	Emmanuel	Emmanuel
	entourage	confidant
secrétariat	ariat	
German	machen	ize
	mit	with
	Sprichwort	proverb
	erischen	Austrian
	minuten	utes
	Haustechnik	umbing
	dringen	urgent
	verfeinern	refine
	umgebung	vironments
	ternehmen	irms
Chinese	到处	everywhere
	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 7: Samples of tokens in each language and the corresponding closest tokens from the English vocabulary according to WECHSEL.

## G Sensitivity Analysis w. r. t. training data size

Evaluating on languages with different amounts of available data only indirectly measures the effect of training data size on WECHSEL since other factors (e.g. language similarity to English) are also involved. We conduct a sensitivity analysis to make the relation to the amount of training data explicit (Table 11). Due to computational constraints we only do this for French. We find that the improvement from WECHSEL increases as the amount of training data decreases. In addition, we find that using fastText embeddings trained on less data deteriorates performance, but still leaves a clear margin to TransInner and FullRand.

Lang	Model	$k$	$\tau$	Scores		
				NLI	POS	LM
French	WECHSEL@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	
	FullRand@0	-	-	46.3	60.6	5.7e+6
CamemBERT	-	-	63.5	93.6	-	
German	WECHSEL@0	1	1	55.8	72.7	6e+5
		10	0.1	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	
	FullRand@0	-	-	44.5	49.1	6.2e+6
GBERT <sub>Base</sub>	-	-	63.2	81.4	-	
Chinese	WECHSEL@0	1	1	47.4	75.4	2.7e+6
		10	0.1	48.0	80.7	2.6e+6
		10	1	48.3	80.3	3.1e+6
		50	0.1	48.3	77.8	3.7e+7
	50	1	47.9	76.5	8.6e+7	
	FullRand@0	-	-	37.5	53.7	5.8e+6
BERT <sub>Base</sub> -Chinese	-	-	61.9	91.9	-	

Table 8: Grid search over the temperature  $\tau$  and number of most similar tokens  $k$  parameters of WECHSEL.

Lang	Model	PPL@0	PPL@25k	PPL@250k
French	WECHSEL-GPT2	1.7e+3	23.47	19.71
	WECHSEL <sub>TRF</sub> -GPT2	2.3e+3	<u>23.45</u>	<u>19.70</u>
German	WECHSEL-GPT2	3.7e+3	34.35	<b>26.80</b>
	WECHSEL <sub>TRF</sub> -GPT2	5.0e+3	34.46	26.82
Chinese	WECHSEL-GPT2	2.4e+4	<u>71.02</u>	<b>51.97</b>
	WECHSEL <sub>TRF</sub> -GPT2	2.5e+4	72.11	52.07
Swahili	WECHSEL-GPT2	1.4e+5	<u>13.02</u>	10.14
	WECHSEL <sub>TRF</sub> -GPT2	1.5e+5	13.03	<b>10.06</b>

Table 9: Results of training WECHSEL<sub>TRF</sub> GPT2 models. We report Perplexity before training (PPL@0), after 10% of steps (PPL@25k) and after training (PPL@250k).

## H Limitations and Potential Risks

### H.1 Limitations

We conduct our experiments on up to eight languages, showing the benefits of our parameter transfer method to both medium- and low-resource languages. However, there are many more languages with diverse linguistic characteristics on which our WECHSEL method is not tested. This is a limitation forced by computational constraints, as we can not ascertain whether transfer to all other languages would result in similar improvements. In addition, our extrinsic evaluation is limited to two tasks (NLI and NER). While this choice is due to the limitations on the available collections in various languages, this evaluation does not necessarily provide a comprehensive view of language understanding tasks.

Lang	Model	Score@0			Score@25k			Score@250k		
		NLI	NER	Avg	NLI	NER	Avg	NLI	NER	Avg
French	WECHSEL-RoBERTa	<u>78.25</u>	86.93	82.59	81.63	<u>90.26</u>	85.95	82.43	<b>90.88</b>	86.65
	WECHSEL <sub>TFR</sub> -RoBERTa	<u>78.25</u>	<u>87.43</u>	<u>82.84</u>	<u>81.86</u>	90.07	<u>85.96</u>	<b>82.55</b>	90.80	<b>86.68</b>
German	WECHSEL-RoBERTa	75.64	84.53	80.08	<u>81.11</u>	89.05	85.08	81.79	<b>89.72</b>	85.76
	WECHSEL <sub>TFR</sub> -RoBERTa	<u>77.00</u>	<u>84.70</u>	<u>80.85</u>	80.71	<u>89.09</u>	84.90	<b>82.04</b>	<b>89.72</b>	<b>85.88</b>
Chinese	WECHSEL-RoBERTa	<u>63.23</u>	72.79	<u>68.01</u>	<u>77.19</u>	<u>79.07</u>	<u>78.13</u>	<b>78.32</b>	80.55	<b>79.44</b>
	WECHSEL <sub>TFR</sub> -RoBERTa	62.75	<u>72.87</u>	67.81	77.07	78.03	77.55	77.99	<b>80.65</b>	79.32
Swahili	WECHSEL-RoBERTa	<u>60.28</u>	74.38	67.33	73.87	87.63	80.75	<b>75.05</b>	87.39	<b>81.22</b>
	WECHSEL <sub>TFR</sub> -RoBERTa	60.14	<u>75.42</u>	<u>67.78</u>	<u>74.04</u>	<u>87.79</u>	<u>80.92</u>	74.58	<b>87.66</b>	81.12

Table 10: Results from fine-tuning WECHSEL<sub>TFR</sub>-RoBERTa models. Results shown the equivalently as in Table 2.

Model	Best PPL				
	Subsample Size	16MiB	64MiB	256MiB	1024MiB
WECHSEL-GPT2 (original fastText embeddings)		78.33	44.75	31.63	24.66
WECHSEL-GPT2 (fastText embeddings trained on subsample)		<u>97.42</u>	<u>49.50</u>	<u>32.88</u>	<u>24.75</u>
FullRand-GPT2		281.46	83.43	43.08	27.09
TransInner-GPT2		216.37	77.71	35.27	25.15

Table 11: Sensitivity Analysis w. r. t. the amount of training data on transfer to French. We train models on random subsamples of 16MiB, 64MiB, 256MiB and 1024MiB of the original training data, and evaluate on the same held-out set. For WECHSEL-GPT2, we train two models. One using the original, publicly available fastText embeddings trained on Common Crawl data. The other using fastText embeddings trained only on the corresponding subsample of text.

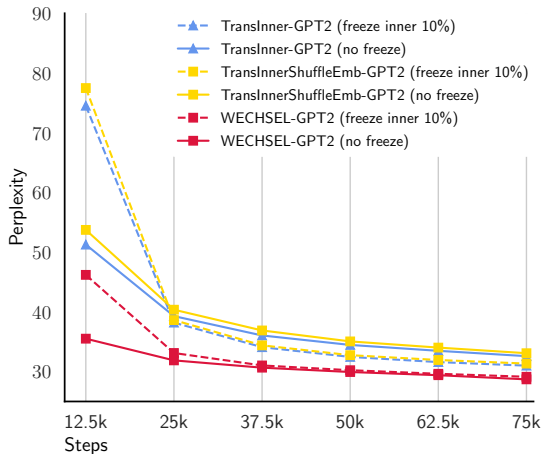


Figure 6: Comparison of German GPT-2 models trained with WECHSEL, TransInner and TransInner-ShuffleEmb between freezing non-embedding parameters at the start and not freezing any parameters.

## H.2 Risks

It is well-known that existing LMs trained on English text encode societal biases and stereotypes and using them in downstream tasks might lead to unfair treatment of various social groups. Since we propose a method to transfer the English LMs to new languages, it is highly probable that the existing biases are also transferred to the target LMs. We therefore advocate a conscious and responsible use of the transferred LMs in practice.

952  
953  
954  
955  
956  
957  
958  
959  
960  
961