

Improving Word Translation via Two-Stage Contrastive Learning

Anonymous ACL submission

Abstract

Word translation or bilingual lexicon induction (BLI) is a key cross-lingual task, aiming to bridge the lexical gap between different languages. In this work, we propose a robust and effective two-stage contrastive learning framework for the BLI task. As Stage C1, we propose to refine standard cross-lingual linear maps between static word embeddings (WEs) via a contrastive learning objective; we also show how to integrate it into the self-learning procedure for even more refined cross-lingual maps. In Stage C2, we conduct BLI-oriented contrastive fine-tuning of mBERT, unlocking its word translation capability. We also show that static WEs induced from the ‘C2-tuned’ mBERT complement static WEs from Stage C1. Comprehensive experiments on standard BLI datasets for diverse languages and different experimental setups demonstrate substantial gains achieved by our framework. While the BLI method from Stage C1 already yields substantial gains over all state-of-the-art BLI methods in our comparison, even stronger improvements are met with the full two-stage framework: e.g., we report gains for 112/112 BLI setups, spanning 28 language pairs.

1 Introduction and Motivation

Bilingual lexicon induction (BLI) or word translation is one of the seminal and long-standing tasks in multilingual NLP (Rapp, 1995; Gaussier et al., 2004; Heyman et al., 2017; Shi et al., 2021, *inter alia*). Its main goal is learning translation correspondences across languages, with applications of BLI ranging from language learning and acquisition (Yuan et al., 2020; Akyurek and Andreas, 2021) to machine translation (Qi et al., 2018; Duan et al., 2020; Chronopoulou et al., 2021) and the development of language technology in low-resource languages and domains (Irvine and Callison-Burch, 2017; Heyman et al., 2018). A large body of recent BLI work has focused on the so-called *mapping-based* methods (Mikolov et al., 2013; Artetxe et al.,

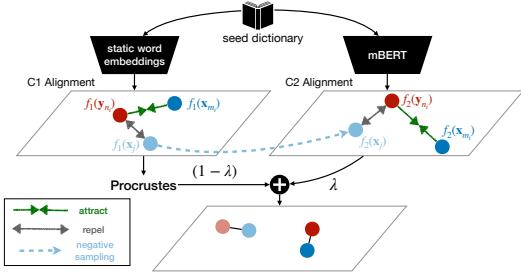


Figure 1: Illustration of the proposed two-stage BLI approach (see §2). It combines contrastive tuning on both static WEs (C1) and pretrained multilingual LMs (C2), where the static WEs are leveraged for selecting negative examples in contrastive tuning of the LM. Output of C1 and C2 is combined for the final BLI task.

2018; Ruder et al., 2019).¹ Such methods are particularly suitable for low-resource languages and weakly supervised learning setups: they support BLI with only as much as few thousand word translation pairs (e.g., 1k or at most 5k) as the only bilingual supervision (Ruder et al., 2019).²

Unlike for many other tasks in multilingual NLP (Doddapaneni et al., 2021; Chau and Smith, 2021; Ansell et al., 2021), state-of-the-art (SotA) BLI results are still achieved via static word embeddings (WEs) (Vulić et al., 2020b; Liu et al., 2021b). A typical *modus operandi* of mapping-based approaches is to first train monolingual WEs independently on monolingual corpora and then map them to a shared cross-lingual space via linear (Mikolov et al., 2013; Glavaš et al., 2019) or non-linear mapping func-

¹They are also referred to as *projection-based* or *alignment-based* methods (Glavaš et al., 2019; Ruder et al., 2019).

²In the extreme, *fully unsupervised* mapping-based BLI methods can leverage monolingual data only without any bilingual supervision (Lample et al., 2018; Artetxe et al., 2018; Hoshen and Wolf, 2018; Mohiuddin and Joty, 2019; Ren et al., 2020, *inter alia*). However, comparative empirical analyses (Vulić et al., 2019) show that, with all other components equal, using seed sets of only 500–1,000 translation pairs, always outperforms fully unsupervised BLI methods. Therefore, in this work we focus on this more pragmatic (weakly) supervised BLI setup (Artetxe et al., 2020); we assume the existence of at least 1,000 seed translations per each language pair.

tions (Mohiuddin et al., 2020). In order to achieve even better results, many BLI methods also apply a self-learning loop where training dictionaries are iteratively (and gradually) refined, and improved mappings are then learned in each iteration (Artetxe et al., 2018; Karan et al., 2020). However, there is still ample room for improvement, especially for lower-resource languages and dissimilar language pairs (Vulić et al., 2019; Nasution et al., 2021).

On the other hand, another line of recent research has demonstrated that a wealth of lexical semantic information is encoded in large multilingual pre-trained language models (LMs) such as mBERT (Devlin et al., 2019), but **1**) it is not straightforward to transform the LMs into multilingual lexical encoders (Liu et al., 2021b), **2**) extract word-level information from them (Vulić et al., 2020b, 2021), and **3**) word representations extracted from these LMs still cannot surpass static WEs in the BLI task (Vulić et al., 2020b; Zhang et al., 2021). Motivated by these insights, in this work we investigate following research questions:

(RQ1) Can we further improve (weakly supervised) mapping-based BLI methods based on static WEs?

(RQ2) How can we extract more useful cross-lingual word representations from pretrained multilingual LMs such as mBERT or XLM-R?

(RQ3) Is it possible to boost BLI by combining cross-lingual representations based on static WEs and the ones extracted from multilingual LMs?

Inspired by the wide success of contrastive learning techniques in *sentence-level* representation learning (Reimers and Gurevych, 2019; Carlsson et al., 2021; Gao et al., 2021), we propose a *two-stage contrastive learning framework for effective word translation* in (weakly) supervised setups; it leverages and combines multilingual knowledge from static WEs and pretrained multilingual LMs. **Stage C1** operates solely on static WEs: in short, it is a mapping-based approach with self-learning, where in each step we additionally fine-tune linear maps with contrastive learning that operates on gradually refined positive examples (i.e., true translation pairs), and hard negative samples. **Stage C2** fine-tunes a pretrained multilingual LM (e.g., mBERT), again with a contrastive learning objective, using positive examples as well as negative examples extracted from the output of C1. Finally, we extract word representations from the multilingual LM fine-tuned in Stage C2, and combine them with static cross-lingual WEs from Stage C1; the

combined representations are then used for BLI.

We run a comprehensive set of BLI experiments on the standard BLI benchmark (Glavaš et al., 2019), comprising 8 diverse languages, in several setups. Our results indicate large gains over state-of-the-art BLI models: e.g., $\approx +8$ Precision@1 points on average, $+10$ points for many language pairs, gains for 107/112 BLI setups already after Stage C1 (cf., RQ1), and for all 112/112 BLI setups after Stage C2 (cf., RQ2 and RQ3). Moreover, our findings also extend to BLI for lower-resource languages from another BLI benchmark (Vulić et al., 2019). Finally, as hinted in recent work (Zhang et al., 2021), our findings validate that multilingual lexical knowledge in LMs, when exposed and extracted as in our contrastive learning framework, can complement the knowledge in static cross-lingual WEs (RQ3), and benefit BLI. We release the code and share the data at: [\[URL\]](#).

2 Methodology

Preliminaries and Task Formulation. In BLI, we assume two vocabularies $\mathcal{X} = \{w_1^x, \dots, w_{|\mathcal{X}|}^x\}$ and $\mathcal{Y} = \{w_1^y, \dots, w_{|\mathcal{Y}|}^y\}$ associated with two respective languages L_x and L_y . We also assume that each vocabulary word is assigned its (static) type-level word embedding (WE), that is, the respective WE matrices for each vocabulary are $X \in \mathbb{R}^{|\mathcal{X}| \times d}$, $Y \in \mathbb{R}^{|\mathcal{Y}| \times d}$. d is WE dimensionality, with typical values $d = 300$ for static WEs (e.g., fastText) (Bojanowski et al., 2017), and $d = 768$ (mBERT and XLM-R WEs). We also assume a set of *seed* translation pairs $\mathcal{D}_0 = \{(w_{m_1}^x, w_{n_1}^y), \dots, (w_{m_{|\mathcal{D}_0|}}^x, w_{n_{|\mathcal{D}_0|}}^y)\}$ for training (Mikolov et al., 2013; Glavaš et al., 2019), where $1 \leq m_i \leq |\mathcal{X}|, 1 \leq n_i \leq |\mathcal{Y}|$. Typical values for the seed dictionary size $|\mathcal{D}_0|$ are 5k pairs and 1k pairs (Vulić et al., 2019), often referred to as *supervised* (5k) and semi-supervised or *weakly supervised* settings (1k) (Artetxe et al., 2018). Given another *test* lexicon $\mathcal{D}_T = \{(w_{t_1}^x, w_{g_1}^y), \dots, (w_{t_{|\mathcal{D}_T|}}^x, w_{g_{|\mathcal{D}_T|}}^y)\}$, where $\mathcal{D}_0 \cap \mathcal{D}_T = \emptyset$, for each L_x test word $w_{t_i}^x$ in \mathcal{D}_T the goal is to retrieve its correct translation from the L_y 's vocabulary \mathcal{Y} , and evaluate it against the gold L_y translation $w_{g_i}^y$ from the pair.

Method in a Nutshell. We propose a novel two-stage contrastive learning (CL) method, with both stages C1 and C2 realised via contrastive learning objectives, see Figure 1. Stage C1 (§2.1) operates solely on static WEs, and can be

seen as a contrastive extension of mapping-based BLI approaches with static WEs. In practice, we blend contrastive learning with the standard SotA mapping-based framework with self-learning: VecMap (Artetxe et al., 2018), with some modifications. Stage C1 operates solely on static WEs in exactly the same BLI setup as prior work, and thus it can be evaluated independently. In Stage C2 (§2.2), we propose to leverage pretrained multilingual LMs for BLI, contrastively fine-tuning them for BLI, and extracting static WEs from the tuned LMs. These LM-based WEs can be combined with WEs obtained in Stage C1 (§2.3).

2.1 Stage C1

Stage C1 is based on the VecMap framework (Artetxe et al., 2018) which features **1)** *dual linear mapping*, where two separate linear transformation matrices map respective source and target WEs to a shared cross-lingual space; and **2)** a *self-learning* procedure that, in each iteration i refines the training dictionary and iteratively improves the mapping. We extend and refine VecMap’s self-learning for supervised and semi-supervised settings via CL.

Initial Advanced Mapping. After ℓ_2 -normalising word embeddings,³ the two mapping matrices, denoted as \mathbf{W}_x for the source language L_x and \mathbf{W}_y for L_y , are computed via the Advanced Mapping (AM) procedure based on the training dictionary, as fully described in Appendix A.1; while VecMap leverages whitening, orthogonal mapping, re-weighting and de-whitening operations to derive mapped WEs, we compute \mathbf{W}_x and \mathbf{W}_y such that a one-off matrix multiplication produces the same result, see Appendix A.1 for the details.

Contrastive Fine-Tuning. At each iteration i , after the initial AM step, the two mapping matrices \mathbf{W}_x and \mathbf{W}_y are then further contrastively fine-tuned via the InfoNCE loss (Oord et al., 2018), a standard and robust choice of a loss function in CL research (Musgrave et al., 2020; Liu et al., 2021c,b). The core idea is to ‘attract’ aligned WEs of positive examples (i.e., true translation pairs) coming from the dictionary \mathcal{D}_{i-1} , and ‘repel’ *hard negative samples*, that is, words which are semantically similar but do not constitute a word translation pair.

These hard negative samples are extracted as follows. Let us suppose that $(w_{m_i}^x, w_{n_i}^y)$ is a translation pair in the current dictionary \mathcal{D}_{i-1} , with its

Algorithm 1 Stage C1: Self-Learning

```

1: Require:  $X, Y, \mathcal{D}_0, \mathcal{D}_{\text{add}} = \emptyset$ 
2: for  $i=1:N_{\text{iter}}$  do
3:    $\mathbf{W}_x, \mathbf{W}_y \leftarrow$  Initial AM using  $\mathcal{D}_{i-1}$ ;
4:    $\mathcal{D}_{\text{CL}} \leftarrow \mathcal{D}_0$  (supervised) or  $\mathcal{D}_{i-1}$  (semi-super);
5:   for  $i=1:N_{\text{CL}}$  do
6:     Retrieve  $\bar{\mathcal{D}}$  for the pairs from  $\mathcal{D}_{\text{CL}}$ ;
7:      $\mathbf{W}_x, \mathbf{W}_y \leftarrow$  Optimize Contrastive Loss;
8:     Compute new  $\mathcal{D}_{\text{add}}$ ;
9:     Update  $\mathcal{D}_i = \mathcal{D}_0 \cup \mathcal{D}_{\text{add}}$ ;
return  $\mathbf{W}_x, \mathbf{W}_y$ 

```

constituent words associated with d -dim static WEs \mathbf{x}_{m_i} and \mathbf{y}_{n_i} . We then retrieve the nearest neighbours of $\mathbf{y}_{n_i} \mathbf{W}_y$ from $X \mathbf{W}_x$ and derive $\bar{w}_{m_i}^x \subset \mathcal{X}$ ($w_{m_i}^x$ excluded), a set of hard negative samples of size N_{neg} . In a similar (symmetric) manner, we also derive the set of negatives $\bar{w}_{n_i}^y \subset \mathcal{Y}$ ($w_{n_i}^y$ excluded). We use $\bar{\mathcal{D}}_i$ to denote a collection of all hard negative set pairs over all training pairs in the current iteration i . We then fine-tune \mathbf{W}_x and \mathbf{W}_y by optimizing the following contrastive objective:

$$s_{i,j} = \exp(\cos(\mathbf{x}_i \mathbf{W}_x, \mathbf{y}_j \mathbf{W}_y) / \tau), \quad (1)$$

$$p_i = \frac{s_{m_i, n_i}}{\sum_{w_j^y \in \{w_{n_i}^y\} \cup \bar{w}_{n_i}^y} s_{m_i, j} + \sum_{w_j^x \in \bar{w}_{m_i}^x} s_{j, n_i}}, \quad (2)$$

$$\min_{\mathbf{W}_x, \mathbf{W}_y} - \mathbb{E}_{(w_{m_i}^x, w_{n_i}^y) \in \mathcal{D}_{\text{CL}}} \log(p_i). \quad (3)$$

τ denotes a standard temperature parameter. The objective, formulated here for a single positive example, spans all positive examples from the current dictionary, along with the respective sets of negative examples computed as described above.

Self-Learning. The application of (a) initial mapping via AM and (b) contrastive fine-tuning can be repeated iteratively. Such self-learning loops typically yield more robust and better performing BLI methods (Artetxe et al., 2018; Vulic et al., 2019). At each iteration i , a set of automatically extracted high-confidence translation pairs \mathcal{D}_{add} are added to the seed dictionary \mathcal{D}_0 , and this dictionary $\mathcal{D}_i = \mathcal{D}_0 \cup \mathcal{D}_{\text{add}}$ is then used in the next iteration $i+1$.

Our dictionary augmentation method slightly deviates from the one used by VecMap. We leverage the most frequent N_{freq} source and target vocabulary words, and conduct forward and backward dictionary induction (Artetxe et al., 2018). Unlike VecMap, we do not add stochasticity to the process, and simply select the top N_{aug} high-confidence word pairs from forward (i.e., source-to-target) induction and another N_{aug} pairs from the backward induction. In practice, we first retrieve the $2 \times N_{\text{aug}}$ pairs with the highest Cross-domain Similarity Lo-

³Unlike VecMap, we do not mean-center WEs as this yielded slightly better results in our preliminary experiments.

cal Scaling (CSLS) scores (Lample et al., 2018),⁴ remove duplicate pairs and those that contradict with ground truth in \mathcal{D}_0 , and add the rest into \mathcal{D}_{add} .

For the initial AM step, we always use the augmented dictionary $\mathcal{D}_0 \cup \mathcal{D}_{\text{add}}$; the same augmented dictionary is used for contrastive fine-tuning in weakly supervised setups.⁵ We repeat the self-learning loop for N_{iter} times; in each iteration, we optimise the contrastive loss N_{CL} times, that is, we go N_{CL} times over all the positive pairs from the training dictionary (at this iteration). N_{iter} and N_{CL} are tunable hyper-parameters. Self-learning in Stage C1 is summarised in Algorithm 1.

2.2 Stage C2

Previous work tried to prompt off-the-shelf multilingual LMs for word translation knowledge via masked natural language templates (Gonen et al., 2020), averaging over their contextual encodings in a large corpus (Vulić et al., 2020b; Zhang et al., 2021), or extracting type-level WEs from the LMs directly without context (Vulić et al., 2020a, 2021). However, even sophisticated templates and WE extraction strategies still typically result in BLI performance inferior to fastText (Vulić et al., 2021).

(BLI-Oriented) Contrastive Fine-Tuning. Here, we propose to fine-tune off-the-shelf multilingual LMs relying on the supervised BLI signal: the aim is to expose type-level word translation knowledge directly from the LM, without any external corpora. In practice, we first prepare a dictionary of positive examples for contrastive fine-tuning: (a) $\mathcal{D}_{\text{CL}} = \mathcal{D}_0$ when $|\mathcal{D}_0|$ spans 5k pairs, or (b) when $|\mathcal{D}_0| = 1k$, we add the $N_{\text{aug}} = 4k$ automatically extracted highest-confidence pairs from Stage C1 (based on their CSLS scores, not present in \mathcal{D}_0) to \mathcal{D}_0 (i.e., \mathcal{D}_{CL} spans $1k + 4k$ word pairs). We then extract N_{neg} hard negatives in the same way as in §2.1, relying on the shared cross-lingual space derived as the output of Stage C1. Our hypothesis is that a difficult task of discerning between true translation pairs and highly similar non-translations as hard negatives, formulated within a contrastive learning objective, will enable mBERT to expose its word translation knowledge, and complement the knowledge already available after Stage C1.

Throughout this work, we assume the use

⁴Further details on the CSLS similarity and its relationship to cosine similarity are available in Appendix A.2.

⁵When starting with 5k pairs, we leverage only \mathcal{D}_0 for contrastive fine-tuning, as \mathcal{D}_{add} might deteriorate the quality of the 5k-pairs seed dictionary due to potentially noisy input.

of pretrained mBERT_{base} model with 12 Transformer layers and 768-dim embeddings.⁶ Each raw word input w is tokenised, via mBERT’s dedicated tokeniser, into the following sequence: $[CLS][sw_1] \dots [sw_M][SEP]$, $M \geq 1$, where $[sw_1] \dots [sw_M]$ refers to the sequence of M constituent subwords/WordPieces of w , and $[CLS]$ and $[SEP]$ are special tokens (Vulić et al., 2020b).

The sequence is then passed through mBERT as the encoder, its encoding function denoted as $f_\theta(\cdot)$: it extracts the representation of the $[CLS]$ token in the last Transformer layer as the representation of the input word w . The full set of mBERT’s parameters θ then gets contrastively fine-tuned in Stage C2, again relying on the InfoNCE CL loss:

$$s'_{i,j} = \exp(\cos(f_\theta(w_i^x), f_\theta(w_j^y))/\tau), \quad (4)$$

$$p'_i = \frac{s'_{m_i, n_i}}{\sum_{w_j^y \in \{w_{n_i}^y\} \cup \bar{w}_{n_i}^y} s'_{m_i, j} + \sum_{w_j^x \in \bar{w}_{m_i}^x} s'_{j, n_i}}, \quad (5)$$

$$\min_{\theta} - \mathbb{E}_{(w_{m_i}^x, w_{n_i}^y) \in \mathcal{D}_{\text{CL}}} \log(p'_i). \quad (6)$$

Type-level WE for each input word w is then obtained simply as $f_{\theta'}(w)$, where θ' refers to the parameters of the ‘BLI-tuned’ mBERT model.

2.3 Combining Output of C1 and C2

In order to combine the output WEs from Stage C1 and the mBERT-based WEs from Stage C2, we also need to map them into a ‘shared’ space: in other words, for each word w , its C1 WE and its C2 WE can be seen as two different views of the same data point. We thus learn an additional linear orthogonal mapping from the C1-induced cross-lingual WE space into the C2-induced cross-lingual WE space. It transforms ℓ_2 -normed 300-dim C1-induced cross-lingual WEs into 768-dim cross-lingual WEs. Learning of the linear map $\mathbf{W} \in \mathbb{R}^{d_1 \times d_2}$, where in our case $d_1 = 300$ and $d_2 = 768$, is formulated as a Generalised Procrustes problem (Schönemann, 1966; Viklands, 2006) operating on all (i.e., both L_x and L_y) words from the seed translation dictionary \mathcal{D}_0 .⁷

⁶We also experimented with XLM-R_{base}, but substantially higher overall results were obtained with mBERT as the underlying/input multilingual LM. We plan to analyse these implications in more detail in future work.

⁷Technical details of the learning procedure are described in Appendix A.3. It is important to note that in this case we do not use word translation pairs $(w_{m_i}^x, w_{n_i}^y)$ directly to learn the mapping, but rather each word $w_{m_i}^x$ and $w_{n_i}^y$ is duplicated to create training pairs $(w_{m_i}^x, w_{m_i}^x)$ and $(w_{n_i}^y, w_{n_i}^y)$, where the left word/item in each pair is assigned its WE from C1, and the right word/item is assigned its WE after C2.

330 Unless noted otherwise, a final representation
 331 of an input word w is then a linear combination
 332 of (a) its C1-based vector \mathbf{v}_w mapped to a 768-
 333 dim representation via \mathbf{W} , and (b) its 768-dim
 334 encoding $f_{\theta'}(w)$ from BLI-tuned mBERT:

$$335 \quad (1 - \lambda) \frac{\mathbf{v}_w \mathbf{W}}{\|\mathbf{v}_w \mathbf{W}\|_2} + \lambda \frac{f_{\theta'}(w)}{\|f_{\theta'}(w)\|_2}, \quad (7)$$

336 where λ is a tunable interpolation hyper-parameter.

337 3 Experimental Setup

338 **Monolingual WEs and BLI Setup.** We largely
 339 follow the standard BLI setup from prior work
 340 (Artetxe et al., 2018; Joulin et al., 2018; Glavaš
 341 et al., 2019; Karan et al., 2020, *inter alia*). The
 342 main evaluation is based on the standard BLI
 343 dataset from Glavaš et al. (2019): it comprises
 344 28 language pairs with a good balance of typologi-
 345 cally similar and distant languages: English (EN),
 346 German (DE), Italian (IT), French (FR), Russian
 347 (RU), Croatian (HR), Turkish (TR), and Finnish (FI).
 348 Again following prior work, we rely on monolin-
 349 gual fastText vectors trained on full Wikipedias for
 350 each language (Bojanowski et al., 2017), where
 351 vocabularies in each language are trimmed to the
 352 200K most frequent words (i.e., $|\mathcal{X}| = 200k$ and
 353 $|\mathcal{Y}| = 200k$). The same fastText WEs are used
 354 for our Stage C1 and in all baseline BLI models.
 355 mBERT in Stage C2 operates over the same vocab-
 356 uaries spanning 200k word types in each language.

357 We use 1k translation pairs (semi-supervised
 358 BLI mode) or 5k pairs (supervised) as seed dictio-
 359 nary \mathcal{D}_0 ; test sets span 2k pairs (Glavaš et al.,
 360 2019). With 56 BLI directions in total,⁸ this yields
 361 a total of 112 BLI setups for each model in our
 362 comparison. The standard *Precision@1* (P@1) BLI
 363 measure is reported, and we rely on CSLS ($k=10$)
 364 to score word similarity (Lample et al., 2018).⁹

365 **Training Setup and Hyperparameters.** Since
 366 standard BLI datasets typically lack a validation set
 367 (Ruder et al., 2019), following prior work (Glavaš
 368 et al., 2019; Karan et al., 2020) we conduct hyper-
 369 parameter tuning on a *single, randomly selected*
 370 language pair EN → TR, and apply those hyperpar-
 371 rameter values in all other BLI runs.

⁸For any two languages L_i and L_j , we run experiments both for $L_i \rightarrow L_j$ and $L_j \rightarrow L_i$ directions.

⁹The same trends in results are observed with Mean Reciprocal Rank (MRR) as another BLI evaluation measure (Glavaš et al., 2019); we omit MRR scores for clarity. Moreover, similar relative trends, but with slightly lower absolute BLI scores, are observed when replacing CSLS with the simpler cosine similarity measure: the results are available in the Appendix.

372 In Stage C1, when $|\mathcal{D}_0|=5k$, the hyperpa-
 373 rameter values are $N_{\text{iter}}=2$, $N_{\text{CL}}=200$, $N_{\text{neg}}=150$,
 374 $N_{\text{freq}}=60k$, $N_{\text{aug}}=10k$. SGD optimiser is used,
 375 with a learning rate of 1.5 and $\gamma=0.99$. When
 376 $|\mathcal{D}_0|=1k$, the values are $N_{\text{iter}}=3$, $N_{\text{CL}}=50$, $N_{\text{neg}}=60$,
 377 $N_{\text{freq}}=20k$, and $N_{\text{aug}}=6k$; SGD with a learning rate
 378 of 2.0, $\gamma=1.0$. $\tau=1.0$ and dropout is 0 in both cases,
 379 and the batch size for contrastive learning is always
 380 equal to the size of the current dictionary $|\mathcal{D}_{\text{CL}}|$
 381 (i.e., $|\mathcal{D}_0|$ (5k case), or $|\mathcal{D}_0 \cup \mathcal{D}_{\text{add}}|$ which varies
 382 over iterations; see §2.1). In Stage C2, $N_{\text{neg}}=28$
 383 and the maximum sequence length is 6. We use
 384 AdamW (Loshchilov and Hutter, 2019) as the opti-
 385 miser with learning rate of $2e-5$ and weight decay
 386 of 0.01. We fine-tune mBERT for 5 epochs, with
 387 a batch size of 100; dropout rate is 0.1 and $\tau=0.1$.
 388 Unless noted otherwise, λ is fixed to 0.2.

389 **Baseline Models.** Our BLI method is evaluated
 390 against four strong SotA BLI models from recent
 391 literature, all of them with publicly available imple-
 392 mentations. Here, we provide brief summaries:¹⁰

393 **RCSLS** (Joulin et al., 2018) optimises a relaxed
 394 CSLS loss, learns a non-orthogonal mapping, and
 395 has been established as a strong BLI model in em-
 396 pirical comparative analyses as its objective func-
 397 tion is directly ‘BLI-oriented’ (Glavaš et al., 2019).

398 **VecMap**’s core components (Artetxe et al., 2018)
 399 have been outlined in §2.1.

400 **LNMap** (Mohiuddin et al., 2020) non-linearly
 401 maps the original static WEs into two latent seman-
 402 tic spaces learned via non-linear autoencoders,¹¹
 403 and then learns another non-linear mapping be-
 404 tween the latent autoencoder-based spaces.

405 **FIPP** (Sachidananda et al., 2021), in brief, first
 406 finds common (i.e., isomorphic) geometric struc-
 407 tures in monolingual WE spaces of both languages,
 408 and then aligns the Gram matrices of the WEs
 409 found in those common structures.

410 For all baselines, we have verified that the hy-
 411 perparameter values suggested in their respective
 412 repositories yield (near-)optimal BLI performance.
 413 Unless noted otherwise, we run VecMap, LNMap,
 414 and FIPP with their own self-learning procedures.¹²

¹⁰For further technical details and descriptions of each BLI model, we refer to their respective publications. We used publicly available implementations of all the baseline models.

¹¹This step is directed towards mitigating anisomorphism (Søgaard et al., 2018; Dubossarsky et al., 2020) between the original WE spaces, which should facilitate their alignment.

¹²RCSLS is packaged without self-learning; extending it to support self-learning is non-trivial and goes beyond the scope of this work.

415
416 **Model Variants.** We denote the full two-stage BLI
417 model as **C2 (Mod)**, where **Mod** refers to the ac-
418 tual model/method used to derive the shared cross-
419 lingual space used by Stage C2. For instance, **C2**
420 (**C1**) refers to the model variant which relies on
421 our Stage C1, while **C2 (RCSLS)** relies on RC-
422 SLS as the base method. We also evaluate BLI
423 performance of our Stage **C1** BLI method alone.

4 Results and Discussion

424 The main results are provided in Table 1, while
425 the full results per each individual language pair,
426 and also with cosine similarity as the word retrieval
427 function, are provided in Appendix D. The main
428 findings are discussed in what follows.

429 **Stage C1 versus Baselines.** First, we note that
430 there is not a single strongest baseline among the
431 four SotA BLI methods. For instance, RCSLS and
432 VecMap are slightly better than LNMap and FIPP
433 with 5k supervision pairs, while FIPP and VecMap
434 come forth as the stronger baselines with 1k su-
435 pervision. There are some score fluctuations over
436 individual language pairs, but the average perfor-
437 mance of all baseline models is within a relatively
438 narrow interval: the average performance of all
439 four baselines is within 3 P@1 points with 5k pairs
440 (i.e., ranging from 38.22 to 41.22), and VecMap,
441 FIPP, and LNMap are within 2 points with 1k pairs.

442 Strikingly, contrastive learning in Stage C1 al-
443 ready yields substantial gains over all four SotA
444 BLI models, which is typically much higher than
445 the detected variations between the baselines. We
446 mark that C1 improves over all baselines in 51/56
447 BLI setups (in the 5k case), and in all 56/56 BLI
448 setups when D_0 spans 1k pairs. The average gains
449 with the C1 variant are ≈ 5 P@1 points over the
450 strongest baseline in both cases. Note that all the
451 models in comparison, all currently considered
452 SotA in the BLI task, use exactly the same monolin-
453 gual WEs and leverage exactly the same amount of
454 bilingual supervision. The gains achieved with our
455 Stage C1 thus strongly indicate the potential and
456 usefulness of word-level contrastive fine-tuning
457 when learning linear cross-lingual maps with static
458 WEs (see RQ1 from §1).

459 **Stage C1 + Stage C2.** The scores improve further
460 with the full two-stage procedure. The **C2 (C1)**
461 BLI variant increases average P@1 for another 3.3
462 (5k) and 3 P@1 points (1k), and we observe gains
463 for all language pairs in both translation directions,
464 rendering Stage C2 universally useful. These gains

[5k] Pairs	RCSLS ⁺	VecMap ^x	LNMap	FIPP	C1	C2 (C1)
DE→*	43.77	40.49	40.35	40.95	<u>46.14</u>	48.86
*→DE	44.74	42.18	39.55	41.66	<u>46.39</u>	50.12
EN→*	50.94	45.43	44.74	45.76	<u>51.31</u>	54.31
*→EN	49.17	50.19	44.32	47.96	<u>52.61</u>	55.47
FI→*	35.11	36.29	33.18	34.83	<u>39.80</u>	43.44
*→FI	33.49	33.40	34.15	33.00	<u>38.82</u>	41.97
FR→*	47.02	44.67	42.80	44.03	<u>49.12</u>	51.91
*→FR	49.42	48.86	46.25	48.08	<u>51.84</u>	54.53
HR→*	34.06	36.26	33.41	33.52	<u>40.22</u>	45.53
*→HR	32.80	32.96	31.34	31.52	<u>37.82</u>	42.65
IT→*	46.59	44.77	43.23	44.11	<u>48.92</u>	51.91
*→IT	48.41	47.85	45.53	46.64	<u>50.99</u>	53.85
RU→*	40.99	41.01	37.94	39.72	<u>44.17</u>	47.24
*→RU	40.10	35.62	35.66	36.03	<u>42.15</u>	45.20
TR→*	31.29	31.54	30.14	30.34	<u>36.61</u>	39.86
*→TR	31.66	29.42	28.99	28.37	<u>35.67</u>	39.26
Avg.	41.22	40.06	38.22	39.16	<u>44.54</u>	47.88

[1k] Pairs	RCSLS ⁺	VecMap ^x	LNMap	FIPP	C1	C2 (C1)
DE→*	33.43	36.69	37.28	37.70	<u>43.94</u>	46.61
*→DE	32.23	38.63	36.74	39.47	<u>43.15</u>	46.01
EN→*	38.16	38.63	40.44	42.26	<u>47.16</u>	49.84
*→EN	38.57	48.39	43.61	46.68	<u>51.59</u>	54.03
FI→*	22.49	33.08	30.00	32.11	<u>36.81</u>	40.28
*→FI	22.29	27.40	29.95	29.88	<u>36.61</u>	39.63
FR→*	34.98	38.65	39.77	41.08	<u>46.23</u>	48.57
*→FR	36.83	46.61	43.81	46.26	<u>49.75</u>	52.17
HR→*	21.59	33.22	30.05	30.93	<u>37.28</u>	42.16
*→HR	20.87	28.15	27.67	28.15	<u>34.00</u>	38.77
IT→*	36.67	39.45	39.93	42.20	<u>46.55</u>	49.22
*→IT	38.33	45.49	43.47	45.17	<u>48.50</u>	50.94
RU→*	28.45	37.75	35.13	38.24	<u>42.21</u>	44.61
*→RU	27.78	26.16	29.71	31.28	<u>38.02</u>	41.04
TR→*	18.72	26.97	26.63	27.05	<u>33.77</u>	36.89
*→TR	17.59	23.63	24.26	24.68	<u>32.34</u>	35.57
Avg.	29.31	35.56	34.90	36.45	<u>41.74</u>	44.77

Table 1: P@1 scores on the BLI benchmark of Glavaš et al. (2019) with bilingual supervision (i.e., D_0 size) of 5k (upper half) and 1k translation pairs (bottom half). $L \rightarrow *$ and $* \rightarrow L$ denote the average BLI scores of BLI setups where L is the source and the target language, respectively. The word similarity measure is CSLS (see §3). Underlined scores are the peak scores among methods that rely solely on static fastText WEs; Bold scores denote the highest scores overall (i.e., the use of word translation knowledge exposed from mBERT is allowed). ⁺RCSLS is always used without self learning (see the footnote in 3); ^xWe report VecMap with self-learning in the 1k-pairs scenario, and its variant without self-learning when using supervision of 5k pairs as it performs better than the variant with self-learning.

indicate that mBERT does contain word translation knowledge in its parameters. However, the model must be fine-tuned (i.e., transformed) to ‘unlock’ the knowledge from its parameters: this is done through a BLI-guided contrastive fine-tuning procedure (see §2.2). Our findings thus further confirm the ‘rewiring hypothesis’ from prior work (Vulić et al., 2021; Liu et al., 2021b; Gao et al., 2021), here validated for the BLI task (see RQ2 from §1), which states that task-relevant knowledge at sentence- and word-level can be ‘rewired’/exposed from the off-the-shelf LMs, even when leveraging very limited task supervision, e.g., with only 1k or 5k word translation pairs as in our experiments.

[1k] Pairs	BG→CA	CA→HE	HE→BG
VecMap	39.43	24.64	31.55
FIPP	34.29	20.63	26.38
C1	41.88	30.56	33.49
C2 (C1)	44.28	33.99	37.78

Table 2: BLI scores on the Panlex-BLI sets.

[5k] Pairs	DE→TR	TR→HR	HR→RU
RCSLS	30.99	24.60	37.19
C2 (RCSLS)	36.52	33.17	44.77
VecMap	27.18	25.99	37.98
C2 (VecMap)	34.95	34.29	44.98
C1	<u>34.69</u>	<u>32.37</u>	<u>41.66</u>
C2 (C1)	38.86	36.32	46.40

[1k] Pairs	DE→TR	TR→HR	HR→RU
RCSLS	18.21	13.84	24.72
C2 (RCSLS)	25.40	22.52	33.88
VecMap	23.37	20.50	36.09
C2 (VecMap)	27.91	26.84	40.45
C1	<u>32.03</u>	<u>27.00</u>	<u>39.40</u>
C2 (C1)	34.85	32.16	42.14

Table 3: Stage C2 with different ‘support’ methods: RCSLS, VecMap, and C1. P@1×100% scores.

Performance over Languages. The absolute BLI scores naturally depend on the actual source and target languages: e.g., the lowest absolute performance is observed for morphologically rich (HR, RU, FI, TR) and non-Indo-European languages (FI, TR). However, both C1 and C2 (C1) mode variants offer wide and substantial gains in performance for *all* language pairs, irrespective of the starting absolute score. This result further suggests wide applicability and robustness of our BLI method.

4.1 Further Discussion

Evaluation on Lower-Resource Languages. The robustness of our BLI method is further tested on another BLI evaluation set: PanLex-BLI (Vulić et al., 2019), which focuses on BLI evaluation for lower-resource language; 1k training pairs and 2k test pairs are derived from PanLex (Kamholz et al., 2014). The results for a subset of three languages (Bulgarian: BG, Catalan: CA, Hebrew: HE) are presented in Table 2, with more results available in Appendix E. Overall, the results further confirm the efficacy of the C2 (C1), with gains observed even with typologically distant language pairs (e.g., CA→HE and HE→BG).

Usefulness of Stage C2? The results in Table 1 have confirmed the effectiveness of our two-stage C2 (C1) BLI method (see RQ3 in §1). However, Stage C2 is in fact independent of our Stage C1, and thus can also be combined with other standard BLI methods. Therefore, we seek to validate whether combining exposed mBERT-based translation knowledge can also aid other BLI methods. In

[1k] Pairs	EN→*	DE→*	IT→*
C1 w/o CL	39.46	37.54	40.37
C1 w/o SL	39.31	32.59	36.45
C1	47.16	43.94	<u>46.55</u>
mBERT	9.55	9.39	8.13
mBERT (tuned)	17.29	20.92	23.29
C1 + mBERT	47.56	44.08	46.74
C2 (C1)	49.84	46.61	49.22

Table 4: Ablation study. CL = Contrastive Learning; SL = Self-Learning. ‘mBERT’ and ‘mBERT (tuned)’ refer to using word encodings from mBERT directly for BLI, before and after fine-tuning in Stage C2. Very similar trends are observed for all other language pairs (available in Appendix F).

other words, instead of drawing positive and negative samples from Stage C1 (§2.2) and combining C2 WEs with WEs from C1 (§2.3), we replace C1 with our baseline models. The results of these *C2 (RCSLS)* and *C2 (VecMap)* BLI variants for a selection of language pairs are provided in Table 3.

The gains achieved with all *C2 (.)* variants clearly indicate that Stage C2 produces WEs which aid all BLI methods. In fact, combining it with RCSLS and VecMap yields even larger relative gains over the base models than combining it with our Stage C1. However, since Stage C1 (as the base model) performs better than RCSLS and VecMap, the final absolute scores with *C2 (C1)* still outperform *C2 (RCSLS)* and *C2 (VecMap)*.

Combining C1 and C2? The usefulness of combining the representations from two stages is measured through varying the value of λ for several BLI setups. The plots are shown in Figure 2, and indicate that Stage C1 is more beneficial to the performance, with slight gains achieved when allowing the ‘influx’ of mBERT knowledge (e.g., λ in the [0.0 – 0.3] interval). While mBERT-based WEs are not sufficient as standalone representations for BLI, they seem to be even more useful in the combined model for lower-resource languages on PanLex-BLI, with steeper increase in performance, and peak scores achieved with larger λ -s.

Ablation Study, with results summarised in Table 4, displays several interesting trends. First, both CL and self-learning are key components in the 1k-setups: removing any of them yields substantial drops.¹³ Further, Table 4 complements the results from Figure 2 and again indicates that, while Stage C2 indeed boosts word translation capacity

¹³In 5k-setups, self-learning becomes less important, and removing it yields only negligible drops, while contrastive self-tuning remains a crucial component, see the Appendix.

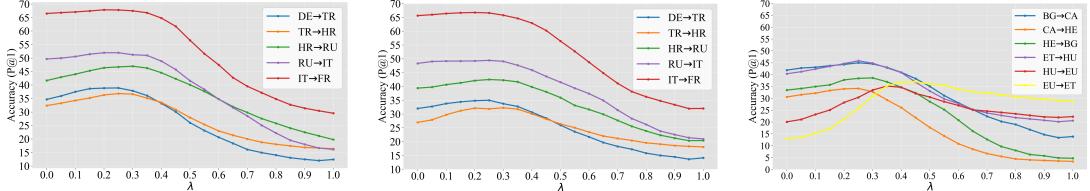


Figure 2: BLI scores with different λ values: (left) $|\mathcal{D}_0|=5k$; (middle) $|\mathcal{D}_0|=1k$; (right) PanLex-BLI, $|\mathcal{D}_0|=1k$.

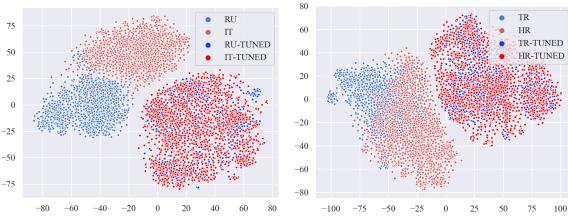


Figure 3: t-SNE visualisation (van der Maaten and Hinton, 2012) of mBERT encodings of words from BLI test sets for RU-IT (left) and TR-HR (right). Similar plots for more language pairs are in Appendix C.

of mBERT, using mBERT features alone is still not sufficient to achieve competitive BLI scores. Finally, Table 4 shows the importance of fine-tuning mBERT before combining it with C1-based WEs (§2.3): directly adding WEs extracted from the off-the-shelf mBERT does not yield any benefits (see the scores for the *C1+mBERT* variant).

The impact of contrastive fine-tuning on mBERT’s representation space for two language pairs is illustrated by a t-SNE plot in Figure 3. The semantic space of off-the-shelf mBERT displays a clear separation of language-specific subspaces (Libovický et al., 2020; Dufter and Schütze, 2020), which makes it unsuitable for the BLI task. On the other hand, contrastive fine-tuning reshapes the subspaces towards a shared (cross-lingual) space, the effects of which are then also reflected in mBERT’s improved BLI capability (see Table 4 again).

5 Related Work

This work is related to three topics, each with a large body of work; we can thus provide only a condensed summary of the most relevant research.

Mapping-based BLI. These BLI methods are highly popular due to reduced bilingual supervision requirements; consequently, they are applicable to low-resource languages and domains, learning linear (Lample et al., 2018; Artetxe et al., 2018; Joulin et al., 2018; Patra et al., 2019; Jawanpuria et al., 2019; Sachidananda et al., 2021) and non-linear maps (Mohiuddin et al., 2020; Glavaš and Vulić,

2020; Ganesan et al., 2021), typically using self-learning in weakly supervised setups.

Contrastive Learning in NLP aims to learn a semantic space such that embeddings of similar text inputs are close to each other, while ‘repelling’ dissimilar ones. It has shown promising performance on training generic sentence encoders (Giorgi et al., 2021; Carlsson et al., 2021; Liu et al., 2021a; Gao et al., 2021) and downstream tasks like summarisation (Liu and Liu, 2021) or NER (Das et al., 2021).

Exposing Lexical Knowledge from Pretrained LMs. Extracting lexical features from off-the-shelf multilingual LMs typically yields subpar performance in lexical tasks (Vulić et al., 2020b). To unlock the lexical knowledge encoded in PLMs, Liu et al. (2021a); Vulić et al. (2021) fine-tune LMs via contrastive learning with manually curated or automatically extracted phrase/word pairs to transform it into effective text encoders. Wang et al. (2021) and Liu et al. (2021c) apply similar techniques for phrase and word-in-context representation learning respectively. The success of these methods suggests that LMs store a wealth of lexical knowledge: yet, as we confirm here for BLI, fine-tuning is typically needed to expose this knowledge.

6 Conclusion

We have proposed a simple yet extremely effective and robust two-stage contrastive learning framework for improving bilingual lexicon induction (BLI). In Stage C1, we tune cross-lingual linear mappings between static word embeddings with a contrastive objective and achieve substantial gains in 107 out of 112 BLI setups on the standard BLI benchmark. In Stage C2, we further propose a contrastive fine-tuning procedure to harvest cross-lingual lexical knowledge from multilingual pre-trained language models. The representations from this process, when combined with Stage C1 embeddings, have resulted in further boosts in BLI performance, with large gains in all 112 setups. We have also conducted a series of finer-grained evaluations, analyses and ablation studies.

References

- Ekin Akyurek and Jacob Andreas. 2021. Lexicon learning for few shot sequence modeling. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP'21)*, pages 4934–4946, Online. Association for Computational Linguistics.
- Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. 2021. Composable sparse fine-tuning for cross-lingual transfer. *CoRR*, abs/2110.07560.
- 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 (ACL'18)*, pages 789–798, Melbourne, Australia. Association for Computational Linguistics.
- Mikel Artetxe, Sebastian Ruder, Dani Yogatama, Gorka Labaka, and Eneko Agirre. 2020. A call for more rigor in unsupervised cross-lingual learning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL'20)*, pages 7375–7388, Online. Association for Computational Linguistics.
- 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.
- Fredrik Carlsson, Amaru Cuba Gyllensten, Evangelia Gogoulou, Erik Ylipää Hellqvist, and Magnus Sahlgren. 2021. Semantic re-tuning with contrastive tension. In *Proceedings of the International Conference on Learning Representations (ICLR'21)*.
- Ethan C. Chau and Noah A. Smith. 2021. Specializing multilingual language models: An empirical study. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 51–61, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alexandra Chronopoulou, Dario Stojanovski, and Alexander Fraser. 2021. Improving the lexical ability of pretrained language models for unsupervised neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL'21)*, pages 173–180, Online. Association for Computational Linguistics.
- Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca J Passonneau, and Rui Zhang. 2021. Container: Few-shot named entity recognition via contrastive learning. *arXiv preprint arXiv:2109.07589*.
- 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 (NAACL'19)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sumanth Doddapaneni, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M. Khapra. 2021. A primer on pretrained multilingual language models. *CoRR*, abs/2107.00676.
- Xiangyu Duan, Baijun Ji, Hao Jia, Min Tan, Min Zhang, Boxing Chen, Weihua Luo, and Yue Zhang. 2020. Bilingual dictionary based neural machine translation without using parallel sentences. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL'20)*, pages 1570–1579, Online. Association for Computational Linguistics.
- Haim Dubossarsky, Ivan Vulić, Roi Reichart, and Anna Korhonen. 2020. The secret is in the spectra: Predicting cross-lingual task performance with spectral similarity measures. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP'20)*, pages 2377–2390, Online. Association for Computational Linguistics.
- Philipp Dufter and Hinrich Schütze. 2020. Identifying elements essential for BERT’s multilinguality. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP'20)*, pages 4423–4437, Online. Association for Computational Linguistics.
- Ashwinkumar Ganesan, Francis Ferraro, and Tim Oates. 2021. Learning a reversible embedding mapping using bi-directional manifold alignment. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3132–3139, Online. Association for Computational Linguistics.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP'21)*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Eric Gaussier, J.M. Renders, I. Matveeva, C. Goutte, and H. Dejean. 2004. A geometric view on bilingual lexicon extraction from comparable corpora. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL'04)*, pages 526–533, Barcelona, Spain.
- John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. 2021. DeCLUTR: Deep contrastive learning for unsupervised textual representations. In *Proceedings of the 59th Annual Meeting of the Association*

729	<i>for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP'21)</i> , pages 879–895, Online. Association for Computational Linguistics.	785
730		786
731		787
732		788
733	Goran Glavaš, Robert Litschko, Sebastian Ruder, and Ivan Vulić. 2019. How to (properly) evaluate cross-lingual word embeddings: On strong baselines, comparative analyses, and some misconceptions. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL'19)</i> , pages 710–721, Florence, Italy. Association for Computational Linguistics.	789
734		790
735		791
736		792
737		793
738		794
739		795
740		795
741	Goran Glavaš and Ivan Vulić. 2020. Non-linear instance-based cross-lingual mapping for non-isomorphic embedding spaces. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL'20)</i> , pages 7548–7555, Online. Association for Computational Linguistics.	796
742		797
743		798
744		799
745		800
746		801
747		802
748	Hila Gonen, Shauli Ravfogel, Yanai Elazar, and Yoav Goldberg. 2020. It's not Greek to mBERT: Inducing word-level translations from multilingual BERT. In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 45–56, Online. Association for Computational Linguistics.	803
749		804
750		805
751		806
752		807
753		807
754	Geert Heyman, Ivan Vulić, and Marie-Francine Moens. 2017. Bilingual lexicon induction by learning to combine word-level and character-level representations. In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL'17)</i> , pages 1085–1095, Valencia, Spain. Association for Computational Linguistics.	808
755		809
756		810
757		811
758		812
759		813
760		813
761		813
762	Geert Heyman, Ivan Vulić, and Marie-Francine Moens. 2018. A deep learning approach to bilingual lexicon induction in the biomedical domain. <i>BMC Bioinformatics</i> , 19(1):259:1–259:15.	814
763		815
764		816
765		817
766		818
767		819
768		820
769		821
770		821
771		821
772	Yedid Hoshen and Lior Wolf. 2018. Non-adversarial unsupervised word translation. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP'18)</i> , pages 469–478, Brussels, Belgium. Association for Computational Linguistics.	822
773		823
774		824
775		825
776		826
777		827
778		828
779		828
780		829
781	Ann Irvine and Chris Callison-Burch. 2017. A comprehensive analysis of bilingual lexicon induction. <i>Computational Linguistics</i> , 43(2):273–310.	830
782		831
783		832
784		833
785		834
786		835
787		836
788		836
789	Pratik Jawanpuria, Arjun Balgovind, Anoop Kunchukuttan, and Bamdev Mishra. 2019. Learning multilingual word embeddings in latent metric space: A geometric approach. <i>Transactions of the Association for Computational Linguistics</i> , 7:107–120.	837
790		838
791		839
792		840
793		840
794		841
795		841
796	Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. 2021a. Self-alignment pretraining for biomedical entity representations. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL'21)</i> , pages 4228–4238, Online. Association for Computational Linguistics.	842
797		843
798		844
799		845
800		846
801		847
802		847
803	Fangyu Liu, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2021b. Fast, effective, and self-supervised: Transforming masked language models into universal lexical and sentence encoders. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP'21)</i> , pages 1442–1459, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	848
804		849
805		850
806		851
807		852
808		852
809		853
810		854
811		855
812		856
813		856
814	Qianchu Liu, Fangyu Liu, Nigel Collier, Anna Korhonen, and Ivan Vulić. 2021c. MirrorWiC: On eliciting word-in-context representations from pretrained language models. In <i>Proceedings of the 25th Conference on Computational Natural Language Learning (CoNLL'21)</i> , pages 562–574, Online. Association for Computational Linguistics.	857
815		858
816		859
817		860
818		861
819		862
820		863
821		863
822		863
823		864
824		865
825		866
826		867
827		868
828		869
829		869
830		870
831		871
832		872
833		873
834		874
835		875
836		875
837		876
838		877
839		878
840		879
841		880

842	on Natural Language Processing (ACL-IJCNLP'21), pages 1065–1072, Online. Association for Computational Linguistics.	897
843		898
844		899
845	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In <i>Proceedings of the International Conference on Learning Representations (ICLR'19)</i> .	900
846		901
847		
848		
849	Tomás Mikolov, Quoc V. Le, and Ilya Sutskever. 2013. Exploiting similarities among languages for ma- chine translation. <i>ArXiv preprint</i> , abs/1309.4168.	902
850		903
851		904
852	Tasnim Mohiuddin, M Saiful Bari, and Shafiq Joty.	905
853	2020. LNMap: Departures from isomorphic as- sumption in bilingual lexicon induction through non- linear mapping in latent space. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natu- ral Language Processing (EMNLP'20)</i> , pages 2712– 2723, Online. Association for Computational Lin- guistics.	906
854		907
855		908
856		909
857		
858		
859		
860	Tasnim Mohiuddin and Shafiq Joty. 2019. Revisiting adversarial autoencoder for unsupervised word trans- lation with cycle consistency and improved training. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Compu- tational Linguistics: Human Language Technologies (NAACL'19)</i> , pages 3857–3867, Minneapolis, Min- nesota. Association for Computational Linguistics.	910
861		911
862		912
863		913
864		914
865		915
866		
867		
868	Kevin Musgrave, Serge J. Belongie, and Ser-Nam Lim. 2020. A metric learning reality check. In <i>Pro- ceedings of the European Conference on Computer Vi- sion (ECCV'20)</i> , pages 681–699.	916
869		917
870		918
871		919
872	Arbi Haza Nasution, Yohei Murakami, and Toru Ishida. 2021. Plan optimization to bilingual dictionary in- duction for low-resource language families. <i>ACM Transactions on Asian and Low-Resource Language Information Processing</i> , 20(2):29:1–29:28.	920
873		921
874		922
875		923
876		924
877	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive pre- dictive coding. <i>ArXiv preprint</i> , abs/1807.03748.	925
878		926
879		927
880	Barun Patra, Joel Ruben Antony Moniz, Sarthak Garg, Matthew R. Gormley, and Graham Neubig. 2019. Bilingual lexicon induction with semi-supervision in non-isometric embedding spaces. In <i>Proceed- ings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL'19)</i> , pages 184– 193, Florence, Italy. Association for Computational Linguistics.	928
881		929
882		930
883		931
884		932
885		933
886		934
887		935
888	Ye Qi, Devendra Sachan, Matthieu Felix, Sarguna Pad- manabhan, and Graham Neubig. 2018. When and why are pre-trained word embeddings useful for neu- ral machine translation? In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL'18)</i> , pages 529–535, New Orleans, Louisiana. Association for Computational Linguistics.	936
889		937
890		938
891		939
892		940
893		941
894		942
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		
918		
919		
920		
921		
922		
923		
924		
925		
926		
927		
928		
929		
930		
931		
932		
933		
934		
935		
936		
937		
938		
939		
940		
941		
942		
943		
944		
945		
946		
947		
948		
949		
950		
951		

952 and Anna Korhonen. 2020a. **Multi-SimLex: A large-**
953 **scale evaluation of multilingual and crosslingual lex-**
954 **ical semantic similarity.** *Computational Linguistics*,
955 46(4):847–897.

956 Ivan Vulić, Goran Glavaš, Roi Reichart, and Anna Ko-
957 rhonen. 2019. **Do we really need fully unsuper-**
958 **vised cross-lingual embeddings?** In *Proceedings*
959 *of the 2019 Conference on Empirical Methods in*
960 *Natural Language Processing and the 9th Interna-*
961 *tional Joint Conference on Natural Language Pro-*
962 *cessing (EMNLP-IJCNLP’19)*, pages 4407–4418,
963 Hong Kong, China. Association for Computational
964 Linguistics.

965 Ivan Vulić, Edoardo Maria Ponti, Anna Korhonen, and
966 Goran Glavaš. 2021. **LexFit: Lexical fine-tuning of**
967 **pretrained language models.** In *Proceedings of the*
968 *59th Annual Meeting of the Association for Compu-*
969 *tational Linguistics and the 11th International Joint*
970 *Conference on Natural Language Processing (ACL-*
971 *IJCNLP’21)*, pages 5269–5283, Online. Association
972 for Computational Linguistics.

973 Ivan Vulić, Edoardo Maria Ponti, Robert Litschko,
974 Goran Glavaš, and Anna Korhonen. 2020b. **Prob-**
975 **ing pretrained language models for lexical seman-**
976 **tics.** In *Proceedings of the 2020 Conference on*
977 *Empirical Methods in Natural Language Process-*
978 *ing (EMNLP’20)*, pages 7222–7240, Online. Asso-
979 ciation for Computational Linguistics.

980 Shufan Wang, Laure Thompson, and Mohit Iyyer. 2021.
981 **Phrase-BERT: Improved phrase embeddings from**
982 **BERT with an application to corpus exploration.** In
983 *Proceedings of the 2021 Conference on Empirical*
984 *Methods in Natural Language Processing*, pages
985 10837–10851, Online and Punta Cana, Dominican
986 Republic. Association for Computational Linguis-
987 tics.

988 Michelle Yuan, Mozhi Zhang, Benjamin Van Durme,
989 Leah Findlater, and Jordan Boyd-Graber. 2020. **In-**
990 **teractive refinement of cross-lingual word embed-**
991 **dings.** In *Proceedings of the 2020 Conference on*
992 *Empirical Methods in Natural Language Process-*
993 *ing (EMNLP’20)*, pages 5984–5996, Online. Asso-
994 ciation for Computational Linguistics.

995 Jinpeng Zhang, Baijun Ji, Nini Xiao, Xiangyu Duan,
996 Min Zhang, Yangbin Shi, and Weihua Luo. 2021.
997 **Combining static word embeddings and contextual**
998 **representations for bilingual lexicon induction.** In
999 *Findings of the Association for Computational Lin-*
1000 *guistics: ACL-IJCNLP 2021*, pages 2943–2955, On-
1001 line. Association for Computational Linguistics.

1002
1003

A Technical Details and Further Clarifications

1004

A.1 Advanced Mapping (AM) in Stage C1

1005 Suppose $\mathbf{X}_{\mathcal{D}}, \mathbf{Y}_{\mathcal{D}} \in \mathbb{R}^{|\mathcal{D}| \times d}$ are source and target
1006 embedding matrices corresponding to the training
1007 dictionary \mathcal{D} . Then $\mathbf{X}_{\mathcal{D}}^T$ and $\mathbf{Y}_{\mathcal{D}}^T$ are whitened, and
1008 singular value decomposition (SVD) is conducted
1009 on the whitened embeddings:

1010

$$\mathbf{X}'_{\mathcal{D}} = \mathbf{X}_{\mathcal{D}} (\mathbf{X}_{\mathcal{D}}^T \mathbf{X}_{\mathcal{D}})^{-\frac{1}{2}}, \quad (8)$$

1011

$$\mathbf{Y}'_{\mathcal{D}} = \mathbf{Y}_{\mathcal{D}} (\mathbf{Y}_{\mathcal{D}}^T \mathbf{Y}_{\mathcal{D}})^{-\frac{1}{2}}, \quad (9)$$

1012

$$\mathbf{U} \mathbf{S} \mathbf{V}^T = \mathbf{X}'_{\mathcal{D}} \mathbf{Y}'_{\mathcal{D}}. \quad (10)$$

1013
1014
1015
1016

\mathbf{W}_x and \mathbf{W}_y are then derived after re-weighting
and de-whitening as follows:

1017

$$\mathbf{W}_x = (\mathbf{X}_{\mathcal{D}}^T \mathbf{X}_{\mathcal{D}})^{-\frac{1}{2}} \mathbf{U} \mathbf{S}^{\frac{1}{2}} \mathbf{U}^T (\mathbf{X}_{\mathcal{D}}^T \mathbf{X}_{\mathcal{D}})^{\frac{1}{2}} \mathbf{U}, \quad (11)$$

1018

$$\mathbf{W}_y = (\mathbf{Y}_{\mathcal{D}}^T \mathbf{Y}_{\mathcal{D}})^{-\frac{1}{2}} \mathbf{V} \mathbf{S}^{\frac{1}{2}} \mathbf{V}^T (\mathbf{Y}_{\mathcal{D}}^T \mathbf{Y}_{\mathcal{D}})^{\frac{1}{2}} \mathbf{V}. \quad (12)$$

1019

A.2 Word Similarity/Retrieval Measures

1020 Given two word embeddings $\mathbf{x} \in \mathbf{X}$ and $\mathbf{y} \in \mathbf{Y}$, their similarity can be defined as their cosine
1021 similarity $m(\mathbf{x}, \mathbf{y}) = \text{cosine}(\mathbf{x}, \mathbf{y})$. In the FIPP
1022 model, we calculate dot product $m(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \cdot \mathbf{y}$
1023 between \mathbf{x} and \mathbf{y} instead without normalisation,
1024 as with FIPP this produces better BLI scores in
1025 general.¹⁴

1026 For the simple Nearest Neighbor (NN) BLI with
1027 cosine (or dot product), we retrieve the word from
1028 the entire target language vocabulary of size 200k
1029 with the highest similarity score and mark it as the
1030 translation of the input/query word in the source
1031 language.

1032 For the Cross-domain Similarity Local Scaling
1033 (CSLS) measure, a CSLS score is defined as
1034 $\text{CSLS}(\mathbf{x}, \mathbf{y}) = 2m(\mathbf{x}, \mathbf{y}) - r_{\mathbf{X}}(\mathbf{y}) - r_{\mathbf{Y}}(\mathbf{x})$. $r_{\mathbf{X}}(\mathbf{y})$
1035 is the average $m(\cdot, \cdot)$ score of \mathbf{y} and its k-NNs
1036 ($k = 10$) in \mathbf{X} ; $r_{\mathbf{Y}}(\mathbf{x})$ is the average $m(\cdot, \cdot)$ scores
1037 of \mathbf{x} and its k-NNs ($k = 10$) in \mathbf{Y} . Note that when
1038 using CSLS scores to retrieve the translation of \mathbf{x}
1039 in \mathbf{Y} , the term $r_{\mathbf{Y}}(\mathbf{x})$ can be ignored, as it is a
1040 constant for all \mathbf{y} , and we can similarly ignore $r_{\mathbf{X}}(\mathbf{y})$
1041 when doing BLI in the opposite direction.

1042
1043
1044

¹⁴<https://github.com/vinsachi/FIPPCLE/blob/main/xling-bli/code/eval.py>

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064

A.3 Generalised Procrustes in Stage C2

We consider the following Procrustes problem:

$$\underset{\mathbf{W}}{\operatorname{argmin}} \|\mathbf{XW} - \mathbf{Y}\|_F^2, \mathbf{WW}^T = \mathbf{I}, \quad (13)$$

1046 where $\mathbf{X} \in \mathbb{R}^{n \times d_1}$ is a C1-induced cross-lingual
1047 space spanning all source and target words in
1048 the training set \mathcal{D} , $\mathbf{Y} \in \mathbb{R}^{n \times d_2}$ is a C2-induced
1049 space representing all mBERT-encoded vectors
1050 corresponding to the same words from \mathbf{X} , and
1051 $\mathbf{W} \in \mathbb{R}^{d_1 \times d_2}$, $d_1 \leq d_2$. A classical Orthogonal
1052 Procrustes Problem assumes that $d_1 = d_2$ and \mathbf{W}
1053 is an orthogonal matrix (i.e., it should be a square
1054 matrix), where its optimal solution is given by
1055 \mathbf{UV}^T ; here, \mathbf{USV}^T is the full singular value de-
1056 composition (SVD) of $\mathbf{X}^T \mathbf{Y}$. In our experiments,
1057 we need to address the case $d_1 < d_2$ when mapping
1058 300-dimensional static fastText WEs to the 768-
1059 dimensional space of mBERT-based WEs. It is easy
1060 to show that when $d_1 < d_2$, $\mathbf{U}[\mathbf{S}, \mathbf{0}] \mathbf{V}^T = \mathbf{X}^T \mathbf{Y}$
1061 (again the full SVD decomposition), the optimal
1062 \mathbf{W} is then $\mathbf{U}[\mathbf{I}, \mathbf{0}] \mathbf{V}^T$ (it degrades to the Orthogo-
1063 nal Procrustes Problem when $d_1 = d_2$). Below, we
1064 provide a simple proof.

1065 Let $\Omega = \mathbf{U}^T \mathbf{WV}$, then $\Omega \Omega^T = \mathbf{I}$. Therefore,
1066 each of its element $-1 \leq \Omega_{i,j} \leq 1$.

$$\begin{aligned} & \underset{\mathbf{W}}{\operatorname{argmin}} \|\mathbf{XW} - \mathbf{Y}\|_F^2 \\ &= \underset{\mathbf{W}}{\operatorname{argmin}} \langle \mathbf{XW} - \mathbf{Y}, \mathbf{XW} - \mathbf{Y} \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmin}} \|\mathbf{XW}\|_F^2 + \|\mathbf{Y}\|_F^2 - 2 \langle \mathbf{XW}, \mathbf{Y} \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} \langle \mathbf{XW}, \mathbf{Y} \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} \langle \mathbf{W}, \mathbf{X}^T \mathbf{Y} \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} \langle \mathbf{W}, \mathbf{X}^T \mathbf{Y} \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} \langle \mathbf{W}, \mathbf{U}[\mathbf{S}, \mathbf{0}] \mathbf{V}^T \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} \langle [\mathbf{S}, \mathbf{0}], \mathbf{U}^T \mathbf{WV} \rangle_F \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} \langle [\mathbf{S}, \mathbf{0}], \Omega \rangle_F \end{aligned} \quad (14)$$

1067 In the formula above, $\|\cdot\|_F$ and $\langle \cdot, \cdot \rangle_F$ are Frobe-
1068 nius norm and Frobenius inner product, and we
1069 leverage their properties throughout the proof. Note
1070 that \mathbf{S} is a diagonal matrix with non-negative el-
1071 ements and thus the maximum is achieved when
1072 $\Omega = [\mathbf{I}, \mathbf{0}]$ and $\mathbf{W} = \mathbf{U}[\mathbf{I}, \mathbf{0}] \mathbf{V}^T$.

1074 Note that the Procrustes mapping over word
 1075 embedding matrices keeps word similarities
 1076 on both sides intact. Since $WW^T = I$,
 1077 $\cos(\mathbf{x}_i W, \mathbf{x}_j W) = \cos(\mathbf{x}_i, \mathbf{x}_j)$.
 1078

1079 We would also like to add an additional note,
 1080 although irrelevant to our own experiments, that the
 1081 above derivation cannot address $d_1 > d_2$ scenarios:
 1082 in that case WW^T cannot be a full-rank matrix
 and thus $WW^T \neq I$.

1083 A.4 Languages in BLI Evaluation

	Language	Family	Code
XLING	Croatian	Slavic	HR
	English	Germanic	EN
	Finnish	Uralic	FI
	French	Romance	FR
	German	Germanic	DE
	Italian	Romance	IT
	Russian	Slavic	RU
	Turkish	Turkic	TR
PanLex-BLI	Basque	-(isolate)	EU
	Bulgarian	Slavic	BG
	Catalan	Romance	CA
	Estonian	Uralic	ET
	Hebrew	Afro-Asiatic	HE
	Hungarian	Uralic	HU

1084 Table 5: A list of languages in our experiments along
 1085 with their language family and ISO 639-1 code.
 1086

1087 B Reproducibility Checklist

- 1088 • **BLI Data:** The two BLI datasets are publicly
 1089 available.¹⁵ ¹⁶
- 1090 • **Static WEs:** We use the preprocessed fast-
 1091 Text WEs provided by Glavaš et al. (2019).
 1092 For PanLex-BLI, we follow the original paper’s setup (Vulić et al., 2019) and
 1093 adopt fastText WEs pretrained on both Common Crawl and Wikipedia(Bojanowski et al.,
 1094 2017).¹⁷ Following prior work, all static WEs
 1095 are trimmed to contain vectors for the top
 200k most frequent words in each language.
- 1096 • **Baseline BLI Models:** All models are ac-
 1097 cessible online as publicly available github repos-

¹⁵<https://github.com/vinsachi/FIPPCLE/blob/main/xling-bli/code/eval.py>

¹⁶<https://github.com/cambridgetl/panlex-bli>

¹⁷<https://fasttext.cc/docs/en/crawl-vectors.html>

1098 itories.

- 1099 • **Pretrained LM:** The used mBERT variant
 1100 is ‘bert-base-multilingual-uncased’, retrieved
 1101 from the huggingface.co model repository.
- 1102 • **Source Code:** Our code is available online at:
 1103 [\[URL-ANONYMOUS\]](https://URL-ANONYMOUS).
- 1104 • **Computing Infrastructure:** We run our
 1105 code on a machine with a 4.00GHz 4-core
 1106 i7-6700K CPU, 64GB RAM and two 12GB
 1107 NVIDIA TITAN X GPUs. We rely on Python
 1108 3.6.10, PyTorch 1.7.0 and huggingface.co
 1109 Transformers 4.4.2. Automatic Mixed Preci-
 1110 sion (AMP)¹⁸ is leveraged during C2 training.
- 1111 • **Runtime:** The training process (excluding
 1112 data loading and evaluation) typically takes
 1113 650 seconds for Stage C1 (seed dictionary of
 1114 5k, 2 self-learning iterations) and 200 seconds
 1115 for C1 (1k, 3 self-learning iterations) on a sin-
 1116 gle GPU. Stage C2 runs for ≈ 500 seconds on
 1117 two GPUs.

1118 C Visualisation of mBERT-Based Word 1119 Representations

1120 To illustrate the impact of the proposed BLI-
 1121 oriented fine-tuning of mBERT in Stage C2 on
 1122 its representation space, we visualise the 768-
 1123 dimensional mBERT word representations (i.e.,
 1124 mBERT-encoded word features alone, without the
 1125 infusion of C1-aligned static WEs). We encode
 1126 BLI test sets (i.e., these sets include 2k source-
 1127 target word pairs unseen during C2 fine-tuning),
 1128 before and after fine-tuning, relying on 1k training
 1129 samples as the seed dictionary D_0 .

1130 Here, we provide comparative t-SNE visuali-
 1131 sations between source and target word mBERT-
 1132 based decontextualised word representations (see
 1133 §2.2) for six language pairs from the BLI dataset of
 1134 Glavaš et al. (2019): EN-IT, FI-RU, EN-HR, HR-
 1135 RU, DE-TR, and IT-FR, while two additional vis-
 1136 1137 1138 1139 1140 1141 1142 1143
 1133 1134 1135 1136 1137 1138 1139 1140 1141 1142 1143
 1134 1135 1136 1137 1138 1139 1140 1141 1142 1143
 1135 1136 1137 1138 1139 1140 1141 1142 1143
 1136 1137 1138 1139 1140 1141 1142 1143
 1137 1138 1139 1140 1141 1142 1143
 1138 1139 1140 1141 1142 1143
 1139 1140 1141 1142 1143
 1140 1141 1142 1143
 1141 1142 1143
 1142 1143

¹⁸<https://pytorch.org/docs/stable/amp.html>

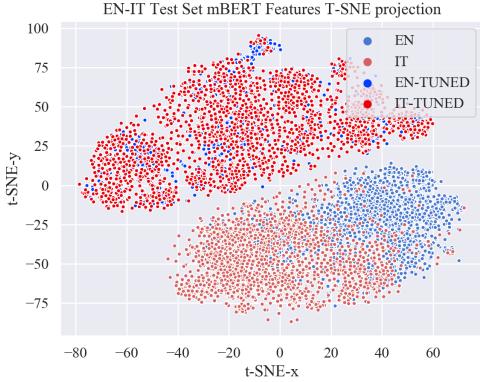


Figure 4: A t-SNE visualisation of mBERT-encoded representations of words from the EN-IT BLI test set. The representations before BLI-oriented fine-tuning of mBERT in Stage C2 are plotted in muted blue and red, and after fine-tuning in bright colours.

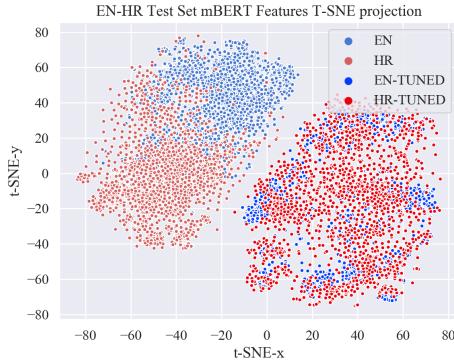


Figure 5: A t-SNE visualisation of mBERT-encoded representations of words from the EN-HR BLI test set. The representations before BLI-oriented fine-tuning of mBERT in Stage C2 are plotted in muted blue and red, and after fine-tuning in bright colours.

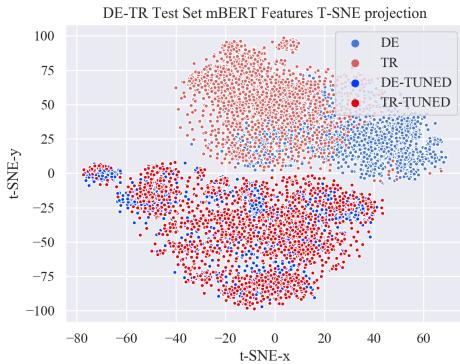


Figure 6: A t-SNE visualisation of mBERT-encoded representations of words from the DE-TR BLI test set. The representations before BLI-oriented fine-tuning of mBERT in Stage C2 are plotted in muted blue and red, and after fine-tuning in bright colours.

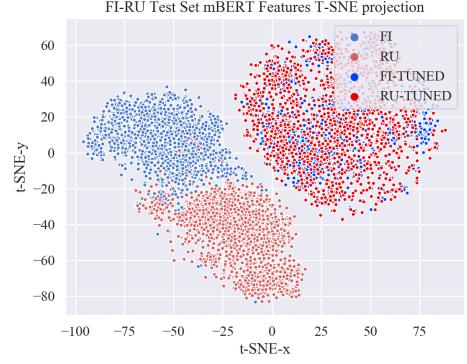


Figure 7: A t-SNE visualisation of mBERT-encoded representations of words from the FI-RU BLI test set. The representations before BLI-oriented fine-tuning of mBERT in Stage C2 are plotted in muted blue and red, and after fine-tuning in bright colours.

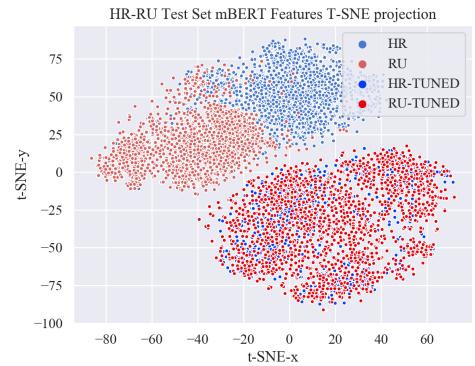


Figure 8: A t-SNE visualisation of mBERT-encoded representations of words from the HR-RU BLI test set. The representations before BLI-oriented fine-tuning of mBERT in Stage C2 are plotted in muted blue and red, and after fine-tuning in bright colours.

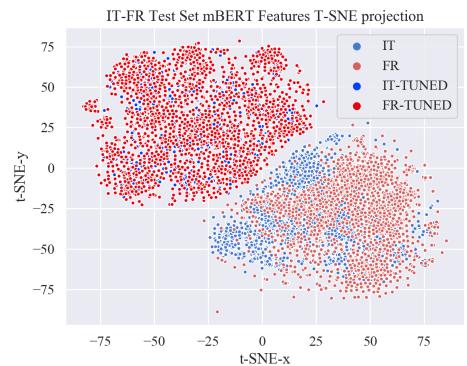


Figure 9: A t-SNE visualisation of mBERT-encoded representations of words from the IT-FR BLI test set. The representations before BLI-oriented fine-tuning of mBERT in Stage C2 are plotted in muted blue and red, and after fine-tuning in bright colours.

1144

D Appendix: Full BLI Results

1145 Complete results on the BLI dataset of Glavaš et al.
1146 (2019), per each language pair and also including
1147 NN-based BLI scores, are provided in Tables 6-7.
1148 It can be seen as an expanded variant of the main
1149 Table 1 presented in the main paper.

1150

E Appendix: Additional Results on the 1151 PanLex-BLI Evaluation Set

1152 Additional results on the PanLex-BLI evaluation
1153 set, focused on typologically diverse and low-
1154 resource languages, for a subset of 3 more lan-
1155 guages are provided in Table 8. These results are
1156 related to the discussion in §4.1 in the main paper.

1157

F Appendix: Full Ablation Study

1158 Complete results of the ablation study, over all lan-
1159 guages in the evaluation set of Glavaš et al. (2019),
1160 are available in Table 9, and can be seen as addi-
1161 tional evidence which supports the claims from the
1162 main paper (see §4.1)

[5k] Pairs	RCSLS	VecMap-Sup	LNMap	FIPP	C1	C2 (C1)
DE→FI	30.62/37.35	29.21/33.59	31.35/36.10	30.93/35.37	38.97/42.10	41.47/44.65
FI→DE	32.48/39.36	35.42/38.73	31.32/36.73	36.05/39.41	39.83/42.46	44.30/47.03
DE→FR	47.63/52.74	46.64/50.44	44.91/48.46	47.89/50.44	51.49/53.78	54.09/55.56
FR→DE	47.23/51.22	45.37/47.75	41.65/44.80	45.73/47.85	50.13/51.37	53.23/53.29
DE→HR	29.26/33.75	27.07/32.08	27.65/32.34	27.65/31.09	34.17/37.66	39.07/42.41
HR→DE	30.30/36.35	32.98/37.24	28.98/33.72	31.51/34.30	39.14/41.35	45.03/48.29
DE→IT	47.68/52.63	47.78/50.55	44.91/47.94	46.90/49.97	50.65/52.79	52.48/54.77
IT→DE	46.51/51.01	44.96/47.29	42.58/45.53	44.86/46.67	49.97/51.21	53.90/53.80
DE→RU	37.87/42.41	31.98/34.38	35.21/37.92	36.57/37.09	42.67/44.29	44.71/46.79
RU→DE	40.54/45.78	40.65/43.32	36.72/40.28	40.18/42.38	46.05/46.73	48.51/49.71
DE→TR	24.93/30.99	23.84/27.18	25.46/29.16	23.94/27.65	31.30/34.69	35.84/38.86
TR→DE	27.00/31.84	26.46/29.93	24.92/27.85	26.09/29.18	33.33/36.74	38.50/40.95
EN→DE	52.95/ <u>57.60</u>	48.65/51.00	45.80/47.95	50.25/51.85	55.50/54.90	59.25/57.75
DE→EN	50.97/56.55	52.01/55.24	46.48/50.50	52.16/55.03	54.77/57.69	56.03/58.95
EN→FI	35.40/42.05	35.25/37.75	34.45/38.35	34.55/39.10	40.70/44.60	45.45/47.15
FI→EN	34.21/41.25	39.04/43.51	31.69/36.26	36.42/40.51	41.46/46.30	44.82/50.55
EN→FR	61.65/ <u>66.55</u>	60.65/63.10	57.75/62.10	61.15/63.25	64.35/65.05	68.45/67.20
FR→EN	59.23/63.11	59.60/62.75	54.53/58.72	59.03/61.87	62.23/63.84	64.30/65.49
EN→HR	31.40/37.90	29.70/34.05	28.40/31.75	28.50/31.95	37.50/40.70	43.60/47.20
HR→EN	28.51/35.67	35.24/39.08	27.83/32.61	31.93/34.72	38.66/42.40	42.61/49.08
EN→IT	58.85/ <u>64.05</u>	57.20/60.40	55.30/59.05	56.95/59.75	61.55/63.45	65.30/65.60
IT→EN	55.09/61.50	57.73/62.17	52.09/56.02	56.69/60.52	59.90/63.51	62.27/65.27
EN→RU	44.75/49.40	38.00/39.65	38.90/41.10	40.70/42.00	48.05/49.15	50.85/50.50
RU→EN	42.80/48.66	45.78/49.35	37.51/42.64	43.27/47.15	48.45/51.91	49.24/54.16
EN→TR	31.40/39.05	30.35/32.05	29.55/32.85	30.80/32.40	39.10/41.35	43.55/44.75
TR→EN	30.78/37.43	34.45/39.24	28.12/33.49	31.79/35.89	39.03/42.60	39.24/44.78
FI→FR	30.90/36.73	34.68/38.26	29.16/34.79	33.79/37.26	38.94/42.20	42.77/45.24
FR→FI	29.59/34.92	31.35/34.30	30.42/33.26	30.11/33.26	36.42/39.99	41.18/43.20
FI→HR	22.65/28.06	27.17/31.58	24.65/29.06	25.54/29.06	30.16/34.89	34.52/38.31
HR→FI	18.20/26.35	28.30/31.72	26.67/31.93	25.78/29.30	32.51/35.61	37.40/39.56
FI→IT	31.53/36.94	33.89/37.99	31.37/35.58	33.58/36.15	38.47/42.04	42.51/46.30
IT→FI	29.56/34.21	31.06/34.32	31.47/35.09	29.97/33.54	35.76/39.48	40.78/43.57
FI→RU	28.74/34.52	31.16/34.16	28.38/32.32	30.37/32.79	35.10/37.73	38.36/40.99
RU→FI	27.29/33.11	29.91/33.53	28.60/33.63	27.82/32.53	35.57/36.98	38.55/40.91
HR→FR	33.46/39.66	35.35/40.24	30.72/36.09	35.30/38.72	39.61/44.13	45.40/49.29
FR→HR	30.94/35.28	29.85/33.21	26.90/30.88	29.69/33.26	36.32/39.78	40.71/44.08
HR→IT	29.62/37.98	36.24/40.24	32.14/36.72	34.19/36.98	38.93/43.77	44.71/48.97
IT→HR	30.34/34.06	30.75/34.32	27.80/32.87	30.03/33.49	37.26/38.71	41.40/44.75
HR→RU	31.35/37.19	34.19/37.98	32.40/36.61	33.19/36.03	39.40/41.66	44.35/46.40
RU→HR	31.48/35.94	34.57/39.50	31.48/35.78	32.16/36.56	37.93/40.60	42.17/45.47
IT→FR	64.19/ <u>66.51</u>	64.03/65.89	62.12/64.60	63.57/65.32	65.37/66.51	66.82/67.86
FR→IT	62.96/66.11	62.70/64.72	61.05/63.68	62.18/64.30	64.25/66.27	66.79/67.20
RU→FR	44.00/47.67	43.58/47.51	38.82/43.64	42.90/47.15	48.04/50.55	50.13/52.70
FR→RU	41.02/ <u>45.01</u>	36.73/38.23	36.26/37.40	37.20/38.54	43.35/44.75	47.13/48.06
RU→IT	41.49/46.57	43.84/46.78	39.50/43.74	43.79/45.89	46.52/49.66	48.66/51.96
IT→RU	40.57/44.13	38.35/38.71	35.87/38.09	38.40/39.43	45.01/45.48	47.08/47.49
TR→FI	21.46/26.46	24.23/28.59	26.14/30.67	24.12/27.90	31.31/32.96	32.85/34.77
FI→TR	23.07/28.90	24.86/29.80	23.86/27.54	24.01/28.64	30.48/32.95	32.74/35.68
TR→FR	29.13/36.10	32.96/36.58	30.56/34.08	31.31/34.40	38.13/40.63	41.43/43.88
FR→TR	27.42/33.52	28.87/31.76	27.42/30.88	26.44/29.13	34.97/37.82	38.70/42.06
TR→HR	20.07/24.60	21.99/25.99	22.42/26.68	21.30/25.24	29.34/32.37	32.43/36.32
HR→TR	17.41/25.25	24.62/27.35	22.30/26.20	22.09/24.62	29.04/32.61	34.14/37.09
TR→IT	28.91/34.56	31.90/34.24	29.66/32.00	29.82/33.44	36.32/38.98	38.87/42.17
IT→TR	28.32/34.73	28.11/30.70	27.96/30.39	27.86/29.82	35.09/37.52	38.19/40.62
TR→RU	23.59/28.06	24.07/26.20	21.99/26.20	24.55/26.36	31.04/32.00	33.60/36.16
RU→TR	24.46/29.18	23.31/27.08	22.58/25.88	25.04/26.35	29.81/32.74	32.48/35.78
Avg.	35.78/41.22	36.76/40.06	34.37/38.22	36.22/39.16	41.95/44.54	45.41/47.88

Table 6: BLI results with 5k seed translation pairs. BLI prediction accuracy (P@1×100%) is reported in the NN/CSLS format (NN: Nearest Neighbor retrieval without CSLS adjustment; CSLS: CSLS retrieval). Underlined scores denote the highest scores among purely fastText-based methods; **bold** scores denote the highest scores in setups where both fastText and mBERT are allowed.

[1k] Pairs	RCSLS	VecMap-Semi	LNMap	FIPP	C1	C2 (C1)
DE→FI	20.97/26.34	23.68/28.33	29.47/32.24	25.56/30.26	<u>37.35/40.85</u>	40.79/43.77
FI→DE	21.18/27.01	32.05/35.00	27.64/34.47	31.79/36.73	<u>37.52/40.57</u>	42.83/44.93
DE→FR	34.06/41.94	46.17/49.03	43.82/47.21	46.48/50.18	<u>49.82/51.75</u>	52.11/54.04
FR→DE	33.89/37.92	42.11/44.34	39.63/42.99	43.30/46.51	<u>46.09/46.82</u>	48.01/48.16
DE→HR	19.25/22.59	22.64/27.39	24.26/28.64	21.91/27.18	<u>30.88/35.16</u>	36.46/40.48
HR→DE	19.10/23.04	30.98/32.82	25.25/29.46	28.77/31.56	<u>35.35/38.45</u>	41.19/44.35
DE→IT	38.81/44.03	46.58/48.72	43.82/47.52	46.01/48.98	<u>48.93/51.28</u>	50.39/52.53
IT→DE	36.64/40.83	41.91/44.39	39.69/42.58	42.95/45.94	<u>46.56/47.86</u>	49.41/49.66
DE→RU	27.80/32.66	20.97/25.46	27.86/30.73	26.03/30.05	<u>40.11/40.27</u>	42.15/42.83
RU→DE	27.82/32.58	36.46/39.08	33.84/37.30	37.98/40.65	<u>42.33/44.21</u>	45.00/46.99
DE→TR	14.03/18.21	20.40/23.37	21.39/24.36	18.94/22.85	<u>29.26/32.03</u>	32.24/34.85
TR→DE	14.43/18.10	23.22/26.57	20.13/24.55	21.67/25.24	<u>30.83/33.71</u>	34.45/37.11
EN→DE	43.00/46.10	46.40/48.20	43.05/45.80	47.95/49.65	<u>49.65/50.40</u>	51.75/50.85
DE→EN	43.14/48.25	51.90/54.56	47.16/50.23	50.97/54.41	<u>53.42/56.23</u>	55.24/57.75
EN→FI	22.40/28.35	24.30/27.95	29.50/33.60	30.40/34.50	<u>38.60/42.15</u>	43.75/45.00
FI→EN	22.70/28.38	37.41/41.15	29.01/35.47	33.68/37.10	<u>39.73/45.51</u>	42.93/48.77
EN→FR	49.00/56.50	57.90/60.00	56.85/60.50	59.65/61.60	<u>60.70/61.65</u>	63.65/62.50
FR→EN	49.46/55.56	58.35/61.41	54.32/58.41	58.72/61.61	<u>60.48/63.27</u>	62.65/64.05
EN→HR	18.65/22.50	21.95/24.95	21.30/25.55	21.70/26.65	<u>32.65/35.65</u>	39.20/42.35
HR→EN	16.57/22.88	34.61/37.45	26.35/30.72	29.77/32.93	<u>35.30/40.87</u>	40.35/47.55
EN→IT	48.65/55.20	55.15/57.55	54.70/57.60	56.00/58.30	<u>57.70/59.60</u>	60.70/61.05
IT→EN	48.22/53.64	56.85/60.78	52.61/56.69	56.59/60.78	<u>59.17/62.64</u>	61.40/63.67
EN→RU	31.50/35.50	21.10/25.05	28.50/32.25	32.75/35.15	<u>43.80/42.50</u>	46.55/46.05
RU→EN	32.37/36.62	44.37/46.20	36.46/41.17	43.27/46.20	<u>47.25/50.29</u>	48.35/53.17
EN→TR	19.35/23.00	24.45/26.70	25.15/27.75	26.40/29.95	<u>36.60/38.15</u>	39.05/41.05
TR→EN	19.81/24.65	33.49/37.17	26.94/32.59	29.98/33.76	<u>36.95/42.33</u>	37.86/43.24
FI→FR	16.13/22.49	31.84/34.79	25.70/30.01	29.58/33.74	<u>37.05/40.36</u>	40.67/43.30
FR→FI	17.69/21.73	21.11/23.95	25.14/28.50	26.49/29.49	<u>34.30/37.61</u>	37.09/40.56
FI→HR	15.24/17.24	25.22/29.90	21.86/26.33	23.49/26.90	<u>25.64/30.01</u>	30.74/34.26
HR→FI	14.05/18.52	25.04/27.62	23.57/27.83	23.99/27.41	<u>28.67/32.61</u>	33.46/36.14
FI→IT	20.13/25.33	32.11/34.68	28.38/31.84	30.27/34.21	<u>35.89/38.99</u>	40.04/42.88
IT→FI	19.07/24.60	22.84/26.10	27.80/30.13	27.96/31.01	<u>34.94/37.83</u>	38.71/41.65
FI→RU	18.44/21.91	26.69/30.27	23.33/27.69	26.48/30.43	<u>31.42/33.89</u>	34.73/37.15
RU→FI	15.72/20.48	29.02/33.11	25.93/31.01	25.93/30.28	<u>32.27/35.31</u>	34.94/37.35
HR→FR	17.99/23.04	35.61/39.14	28.35/32.93	30.19/34.67	<u>37.14/41.14</u>	43.08/45.71
FR→HR	16.76/20.54	23.80/27.52	24.00/28.45	25.50/28.56	<u>32.70/35.33</u>	36.26/39.68
HR→IT	20.52/26.20	36.40/38.77	29.46/33.09	31.93/35.03	<u>37.40/40.24</u>	42.40/46.19
IT→HR	18.81/23.72	23.88/28.68	24.81/28.63	26.10/30.44	<u>33.02/35.92</u>	37.62/41.29
HR→RU	20.99/24.72	32.40/36.09	29.35/34.30	30.30/34.09	<u>37.30/39.40</u>	40.72/42.14
RU→HR	20.32/25.67	34.10/38.08	29.70/33.94	30.91/36.14	<u>34.68/38.92</u>	38.03/41.17
IT→FR	55.25/59.95	<u>63.41/65.06</u>	60.93/63.93	63.05/65.22	<u>63.41/65.63</u>	65.27/66.77
FR→IT	55.25/59.91	62.13/63.58	60.37/62.80	61.98/64.15	<u>63.11/64.56</u>	64.46/65.49
RU→FR	26.72/33.68	42.33/45.42	36.04/40.54	41.91/46.57	<u>46.52/48.87</u>	48.87/51.28
FR→RU	27.06/30.83	20.33/24.57	27.57/31.92	29.69/32.90	<u>40.71/40.46</u>	43.66/43.61
RU→IT	30.59/35.36	41.91/43.74	38.92/41.80	42.54/44.94	<u>45.10/48.35</u>	46.46/49.24
IT→RU	29.82/32.97	22.89/26.10	29.20/31.47	33.49/35.76	<u>41.34/41.50</u>	43.41/43.57
TR→FI	13.31/16.03	19.81/24.76	21.73/26.36	21.73/26.20	<u>26.94/29.93</u>	30.35/32.96
FI→TR	11.77/15.08	21.97/25.80	19.71/24.17	21.49/25.64	<u>24.96/28.32</u>	27.80/30.64
TR→FR	16.67/20.23	30.46/32.85	26.57/31.52	28.27/31.84	<u>35.46/38.82</u>	38.92/41.59
FR→TR	14.43/18.37	22.19/25.19	23.02/25.30	21.83/24.37	<u>32.02/35.59</u>	35.70/38.44
TR→HR	11.66/13.84	16.19/20.50	19.01/22.15	17.15/21.19	<u>22.74/27.00</u>	27.85/32.16
HR→TR	10.10/12.73	19.57/20.67	18.57/21.99	18.36/20.83	<u>22.51/28.25</u>	28.88/33.04
TR→IT	17.15/22.31	29.29/31.42	26.94/29.66	26.62/30.56	<u>33.65/36.47</u>	36.42/39.19
IT→TR	16.12/20.98	22.22/25.06	23.93/26.10	23.62/26.25	<u>32.66/34.47</u>	35.50/37.93
TR→RU	12.94/15.87	13.05/15.55	15.87/19.60	17.04/20.55	<u>25.35/28.12</u>	29.82/31.95
RU→TR	11.42/14.77	16.61/18.60	17.02/20.12	20.90/22.89	<u>26.40/29.54</u>	30.07/33.05
Avg.	24.73/29.31	32.50/35.56	31.10/34.90	33.00/36.45	<u>38.97/41.74</u>	42.33/44.77

Table 7: BLI results with 1k seed translation pairs. BLI prediction accuracy (P@1 × 100%) is reported in the NN/CSLS format (NN: Nearest Neighbor retrieval without CSLS adjustment; CSLS: CSLS retrieval). Underlined scores denote the highest scores among purely fastText-based methods; **bold** scores denote the highest scores in setups where both fastText and mBERT are allowed.

[1k] Pairs	ET→HU	HU→EU	EU→ET
VecMap	35.55	20.03	9.83
FIPP	30.30	11.58	8.22
C1	40.35	20.09	13.00
C2	44.64	28.26	21.35
mBERT	15.40	16.97	23.70
mBERT(tuned)	20.59	22.30	28.62
C2($\lambda=0.4$)	-	34.62	36.70

Table 8: Additional BLI scores on the PanLex-BLI evaluation sets of [Vulić et al. \(2019\)](#); ‘mBERT’ and ‘mBERT (tuned)’ refer to using word encodings from mBERT directly for BLI, before and after fine-tuning in Stage C2.

[5k] Pairs	C1 w/o CL	C1 w/o SL	C1	mBERT	mBERT(tuned)	C1+mBERT	C2 (C1)
DE→*	35.16/39.30	41.70/45.07	<u>43.43/46.14</u>	8.90/9.39	17.70/18.66	43.13/46.25	46.24/48.86
*→DE	37.24/41.23	43.46/45.85	<u>44.85/46.39</u>	8.86/9.51	18.10/19.21	44.61/46.47	48.96/50.12
EN→*	37.99/41.58	48.41/50.99	<u>49.54/51.31</u>	9.29/9.55	15.08/15.87	49.44/51.55	53.78/54.31
*→EN	46.36/50.16	47.36/51.18	<u>49.21/52.61</u>	10.42/10.71	21.34/22.58	48.96/52.77	51.22/55.47
FI→*	31.92/36.78	33.62/38.21	<u>36.35/39.80</u>	5.73/5.93	12.23/13.23	35.97/40.00	40.00/43.44
*→FI	26.16/31.13	33.07/37.26	<u>35.89/38.82</u>	5.57/5.89	11.99/12.95	35.48/39.05	39.67/41.97
FR→*	38.60/42.41	45.27/48.40	<u>46.81/49.12</u>	9.65/10.18	18.37/19.70	46.65/49.29	50.29/51.91
*→FR	45.30/48.85	47.35/50.82	<u>49.42/51.84</u>	9.86/10.38	20.01/21.10	49.07/51.92	52.73/54.53
HR→*	30.88/35.52	33.95/38.51	<u>36.76/40.22</u>	7.11/7.72	17.52/18.57	36.13/40.40	41.95/45.53
*→HR	26.94/32.19	32.24/36.42	<u>34.67/37.82</u>	7.09/7.54	16.83/17.81	34.23/38.08	39.13/42.65
IT→*	39.06/42.67	45.55/48.39	<u>46.91/48.92</u>	7.47/8.13	18.64/20.18	46.35/48.91	50.06/51.91
*→IT	44.48/47.60	46.35/49.93	<u>48.10/50.99</u>	7.03/7.46	16.24/17.12	47.66/51.07	51.33/53.85
RU→*	37.46/40.84	39.30/42.81	<u>41.77/44.17</u>	1.95/2.29	14.50/15.74	41.56/44.38	44.25/47.24
*→RU	27.85/32.12	39.04/41.46	<u>40.66/42.15</u>	1.38/1.94	11.47/13.25	40.53/42.39	43.73/45.20
TR→*	26.14/30.92	31.12/35.08	<u>34.07/36.61</u>	6.18/6.53	12.10/12.87	33.41/36.81	36.70/39.86
*→TR	22.88/26.74	30.08/34.55	<u>32.83/35.67</u>	6.07/6.28	10.14/10.79	32.09/35.85	36.52/39.26
Avg.	34.65/38.75	39.87/43.43	<u>41.95/44.54</u>	7.04/7.46	15.77/16.85	41.58/44.70	45.41/47.88
[1k] Pairs	C1 w/o CL	C1 w/o SL	C1	mBERT	mBERT(tuned)	C1+mBERT	C2 (C1)
DE→*	33.39/37.54	24.74/32.59	<u>41.40/43.94</u>	8.90/9.39	20.26/20.92	41.46/44.08	44.20/46.61
*→DE	35.21/38.73	24.01/32.08	<u>41.19/43.15</u>	8.86/9.51	20.78/21.10	41.48/43.37	44.66/46.01
EN→*	35.65/39.46	33.21/39.31	<u>45.67/47.16</u>	9.29/9.55	16.92/17.29	46.05/47.56	49.24/49.84
*→EN	44.95/49.02	28.26/39.19	<u>47.47/51.59</u>	10.42/10.71	26.11/26.82	47.08/51.63	49.83/54.03
FI→*	29.34/33.91	13.17/21.10	<u>33.17/36.81</u>	5.73/5.93	15.66/16.13	33.15/36.90	37.11/40.28
*→FI	23.35/28.38	14.12/20.73	<u>33.30/36.61</u>	5.57/5.89	14.80/15.35	33.27/36.83	37.01/39.63
FR→*	36.34/39.49	27.86/34.51	<u>44.20/46.23</u>	9.65/10.18	20.74/21.59	44.15/46.52	46.83/48.57
*→FR	44.06/47.64	28.73/36.32	<u>47.16/49.75</u>	9.86/10.38	23.03/23.59	47.24/49.88	50.37/52.17
HR→*	28.42/33.07	12.40/20.76	<u>33.38/37.28</u>	7.11/7.72	20.41/20.97	33.01/37.38	38.58/42.16
*→HR	24.15/28.84	14.61/20.67	<u>30.33/34.00</u>	7.09/7.54	19.18/19.74	30.49/34.30	35.17/38.77
IT→*	36.71/40.37	29.04/36.45	<u>44.44/46.55</u>	7.47/8.13	22.25/23.29	44.42/46.74	47.33/49.22
*→IT	43.02/46.05	29.42/37.68	<u>45.97/48.50</u>	7.03/7.46	19.27/19.86	45.75/48.54	48.70/50.94
RU→*	35.36/38.69	18.95/27.72	<u>39.22/42.21</u>	1.95/2.29	18.86/19.12	39.09/42.27	41.67/44.61
*→RU	24.33/28.77	20.82/26.62	<u>37.15/38.02</u>	1.38/1.94	14.57/15.74	37.41/38.37	40.15/41.04
TR→*	24.06/28.62	11.39/18.22	<u>30.27/33.77</u>	6.18/6.53	14.80/15.28	30.07/33.92	33.67/36.89
*→TR	20.19/23.73	10.80/17.36	<u>29.20/32.34</u>	6.07/6.28	12.14/12.40	28.69/32.44	32.75/35.57
Avg.	32.41/36.39	21.35/28.83	<u>38.97/41.74</u>	7.04/7.46	18.74/19.32	38.93/41.92	42.33/44.77

Table 9: Full ablation study on 8 languages, 28 language pairs in both directions with training dictionary sizes of 5k and 1k respectively, that is, 112 BLI setups for each method. $L \rightarrow *$ and $* \rightarrow L$ denote the average BLI scores of BLI setups where L is the source and the target language, respectively. BLI prediction accuracy (P@1×100%) is reported in the NN/CSLS format (NN: Nearest Neighbor retrieval without CSLS adjustment; CSLS: CSLS retrieval). Underlined scores denote the highest scores among purely fastText-based methods; **bold** scores denote the highest scores in setups where both fastText and mBERT are allowed.