# Exposing Cross-Lingual Lexical Knowledge
# from Multilingual Sentence Encoders

**Ivan Vulić**[1]    **Goran Glavaš**[2]    **Fangyu Liu**[1]
**Nigel Collier**[1]    **Edoardo Maria Ponti**[3,4]    **Anna Korhonen**[1]

[1]Language Technology Lab, TAL, University of Cambridge
[2]CAIDAS, Julius Maximilian University of Würzburg
[3]Mila – Quebec Artificial Intelligence Institute   [4]McGill University
{iv250, alk23}@cam.ac.uk

## Abstract

Pretrained multilingual language models (LMs) can be successfully transformed into multilingual sentence encoders (SEs; e.g., LABSE, XMPNET) via additional fine-tuning or model distillation on parallel data. However, it remains uncertain how to best leverage their knowledge to represent sub-sentence lexical items (i.e., words and phrases) in cross-lingual lexical tasks. In this work, we probe these SEs for the amount of cross-lingual lexical knowledge stored in their parameters, and compare them against the original multilingual LMs. We also devise a novel method to expose this knowledge by additionally fine-tuning multilingual models through inexpensive contrastive learning procedure, requiring only a small amount of word translation pairs. We evaluate our method on bilingual lexical induction (BLI), cross-lingual lexical semantic similarity, and cross-lingual entity linking, and report substantial gains on standard benchmarks (e.g., +10 Precision@1 points in BLI), validating that the SEs such as LABSE can be 'rewired' into effective cross-lingual lexical encoders. Moreover, we show that resulting representations can be successfully interpolated with static embeddings from cross-lingual word embedding spaces to further boost the performance in lexical tasks. In sum, our approach provides an effective tool for exposing and harnessing multilingual lexical knowledge 'hidden' in multilingual sentence encoders.

## 1 Introduction

Transfer learning with pretrained Masked Language Models (LMs) such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) offers unmatched performance in a large number of NLP tasks (Wang et al., 2019; Raffel et al., 2020). However, despite the wealth of semantic knowledge stored in the pretrained LMs (Rogers et al., 2020; Vulić et al., 2020b), they do not produce coherent and effective sentence encodings when used off-the-shelf (Liu et al., 2021c). Their further adaptation is required, similar to standard task fine-tuning (Reimers and Gurevych, 2019; Li et al., 2020; Yan et al., 2021, *inter alia*). They get *transformed* into sentence encoders (SEs) via dual-encoder frameworks relying on contrastive learning objectives (van den Oord et al., 2018; Musgrave et al., 2020). This transformation can be achieved via supervised (i.e., leveraging labeled external data such as NLI or sentence similarity annotations) (Reimers and Gurevych, 2019; Vulić et al., 2021b; Liu et al., 2021a) or, more recently, even fully unsupervised fine-tuning (Liu et al., 2021c; Gao et al., 2021).

Following the procedures from monolingual setups, another line of research has been transforming multilingual LMs into *multilingual sentence encoders* (Feng et al., 2022; Reimers and Gurevych, 2020). These multilingual SEs enable effective sentence matching and ranking in multiple languages as well as cross-lingually (Litschko et al., 2022). The transformation is typically done by combining **1)** LM objectives on monolingual data available in multiple languages with **2)** cross-lingual objectives such as Translation Language Modeling (TLM) (Conneau and Lample, 2019) and/or cross-lingual contrastive ranking (Yang et al., 2020). Such multilingual SEs consume a large number of parallel sentences for the latter objective. Consequently, they outperform multilingual off-the-shelf MLMs in cross-lingual sentence similarity and ranking applications (Liu et al., 2021d; Litschko et al., 2022). However, as we show in this work, off-the-shelf multilingual SEs might still lag behind traditional *static* cross-lingual word embeddings (CLWEs) when encoding sub-sentence *cross-lingual lexical items* (e.g., words or phrases) (Liu et al., 2021c) for tasks such as bilingual lexicon induction (BLI).

In this work, we probe multilingual sentence encoders for *cross-lingual lexical knowledge*. We demonstrate that multilingual SEs, due to their fine-
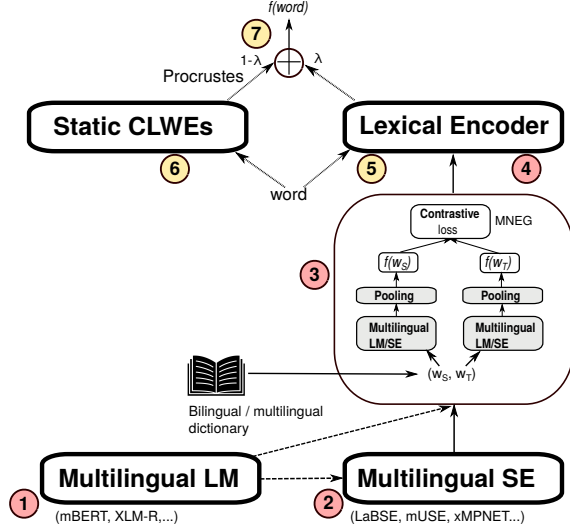
Figure 1: An illustration of the full pipeline of exposing cross-lingual lexical knowledge from multilingual language models (LMs) and sentence encoders (SEs), described later in §2. Multilingual LMs (①) can be transformed into multilingual SEs (②) as done in previous work (Reimers and Gurevych, 2020; Feng et al., 2022). A contrastive fine-tuning procedure (③) can be applied on any multilingual LM (①) or SE (②), leveraging an external bilingual (seed) dictionary: this yields a cross-lingual lexical encoder (④). At inference, a word/phrase is encoded by the fine-tuned lexical encoder (⑤). In addition, its encoding can be interpolated with the corresponding static (cross-lingual) word embedding (⑥), producing the final lexical representation/vector (⑦). Before the interpolation, static CLWEs must be mapped into the vector space of the lexical encoder (③): to this end, we use the orthogonal Procrustes projection. The interpolation can also done using word vectors obtained directly from multilingual LMs (①) and SEs (②), not shown above for clarity.

tuning on multilingual and parallel data, indeed store a wealth of such knowledge, much more than what 'meets the eye' when they are used as off-the-shelf encoders. However, this lexical knowledge has to be *exposed* from the original multilingual SEs, (again) through additional fine-tuning. From another perspective, we show that the input multilingual sentence encoder can be 'rewired' into an effective cross-lingual *lexical encoder*, see Figure 1. This rewiring transformation is again done via a quick and inexpensive contrastive learning procedure: turning the multilingual SE into an enhanced bilingual lexical encoder requires only up to 1k-5k word translation pairs for any language pair.[1]

---

[1]Note that this is a typical requirement of standard mapping-based approaches for learning static cross-lingual word embeddings which excel in the BLI task (Mikolov et al., 2013; Conneau et al., 2018; Glavaš and Vulić, 2020).

We demonstrate the usefulness of such 'contrastively exposed' cross-lingual lexical knowledge on three standard cross-lingual lexical tasks using standard evaluation data and protocols: BLI, cross-lingual lexical semantic similarity, and cross-lingual entity linking. Our main results indicate substantial gains from the contrastive lexical knowledge exposure: e.g., ≈+10 Precision@1 BLI points on standard BLI benchmarks (Glavaš et al., 2019). Furthermore, while cross-lingual lexical knowledge can also be exposed via contrastive learning on mBERT and XLM-R directly (Li et al., 2022) (see again Figure 1), we empirically validate that multilingual SEs such as LABSE (Feng et al., 2022) and multilingual xMPNET (Song et al., 2020; Reimers and Gurevych, 2020) indeed contain much richer lexical knowledge than mBERT or XLM-R, likely owing to their additional exposure to parallel data. This is evidenced by large performance gains over mBERT and XLM-R, both before *and* after additional contrastive fine-tuning of the input models with word translation pairs (see again Figure 1).

We also demonstrate that word encodings from the lexical encoder obtained by applying contrastive tuning on input multilingual SEs such as LABSE and xMPNET can be easily and effectively interpolated with standard static CLWEs (Artetxe et al., 2018) (see Figure 1) for even stronger performance in lexical tasks such as BLI. Our results further indicate the importance of exposing cross-lingual lexical knowledge from multilingual SEs especially in the case of low-resource languages, as demonstrated on a low-resource BLI benchmark (Vulić et al., 2019).

## 2 From Multilingual Sentence Encoders to Cross-Lingual Lexical Encoders

**Motivation.** In a nutshell, the motivation for this work largely stems from the research on *probing and analyzing* pretrained language models and encoders for various types of knowledge they might implicitly store in their parameters (Ethayarajh, 2019; Jawahar et al., 2019; Rogers et al., 2020). Here, we focus on a *particular knowledge type*: cross-lingual lexical knowledge, and its extraction from multilingual LMs and especially from SEs.

Previous work tried to prompt off-the-shelf multilingual LMs for word translation knowledge via masked natural language templates (Gonen et al., 2020), extracting type-level word embeddings from

the LMs directly without context (Vulić et al., 2020a, 2021a), or averaging over their contextual encodings using a large auxiliary corpus in the target language (Bommasani et al., 2020; Litschko et al., 2022). However, previous work **1)** demonstrated that even sophisticated templates and extraction strategies still cannot match static cross-lingual word embeddings (e.g., based on fastText) in cross-lingual lexical tasks such as BLI (Vulić et al., 2020b), and **2)** did not attempt to expose such knowledge from multilingual SEs, nor **3)** compared the lexical knowledge extracted from multilingual LMs to the knowledge from multilingual SEs.

**Multilingual Sentence Encoders.** A large body of work focuses on inducing multilingual encoders that capture meaning of sentences across multiple languages (Artetxe and Schwenk, 2019; Feng et al., 2022; Yang et al., 2020, *inter alia*). The most popular approach fine-tunes multilingual MLMs on cross-lingual sentence matching datasets (Reimers and Gurevych, 2020). Such SEs specialized for semantic similarity are supposed to encode sentence meaning more accurately, supporting tasks that require unsupervised semantic text similarity and ranking, monolingually and across languages.

While the primary goal of off-the-shelf mBERT and XLM-R models is contextualizing (sub)word representations, multilingual SEs produce a semantic encoding of input text (e.g., a sentence in any of the supported languages) directly. Their primary purpose is encoding sentences, but they can be equally applied to sub-sentential text: words and phrases. Our working hypothesis is that the pretrained multilingual SEs can be turned into effective static/decontextualized cross-lingual word encoders. This can be done via additional lexical fine-tuning (Vulić et al., 2021a) on a small set of word translation pairs from an external seed bilingual dictionary.

## 2.1 Methodology

**Text Input** (e.g., words or phrases), in any supported language, is provided to each multilingual encoder fully "in isolation" (ISO), without any surrounding context. The ISO approach has been verified in previous work (Vulić et al., 2021a; Litschko et al., 2022; Li et al., 2022), and is lightweight: **1)** it disposes of any external text corpora and is not impacted by the external data; **2)** it encodes words more efficiently due to the absence of context. Moreover, it allows us to directly study the

richness of cross-lingual information stored in the encoders' parameters, and its interaction with additional cross-lingual signal from bilingual lexicons.

The input word $w$ gets tokenized via the encoder's dedicated tokenizer into the sequence $[SPEC1][sw_1] \ldots [sw_M][SPEC2]$, where $[sw_1] \ldots [sw_M]$, $M \geq 1$, refers to the sequence of $M$ constituent subwords/WordPieces of $w$, and $[SPEC1]$ and $[SPEC2]$ are placeholders for the encoder's special tokens.[2]

For simplicity, we assume that each language $L_i$ is associated with its word vocabulary $\mathcal{V}_i$, spanning a finite amount of word types.

**Lexical Encoding.** We treat LMs and SEs exactly the same, and denote the encoding function of any multilingual encoder as $f_\theta(\cdot)$, where $\theta$ denotes the encoder's parameters. The final $d$-dimensional encoding of the input word $w$, $f_\theta(w)$, is the mean average of the representations of each subword from $w$'s corresponding sequence $[SPEC1][sw_1] \ldots [sw_M][SPEC2]$, where the representations are obtained from the last Transformer layer.[3]

**Contrastive Fine-Tuning.** Let us assume that we have a bilingual dictionary at our disposal, spanning $N$ (typically $N$=5,000) word translation pairs $\mathcal{D} = \{(w_{1,s}, v_{1,t}), \ldots, (w_{N,s}, v_{N,t})\}$ for a language pair $L_s$-$L_t$.[4] We treat the dictionary entries as *positive examples* for the contrastive fine-tuning procedure, which should be 'attracted' (i.e., pulled closer together) in the contrastively fine-tuned space.

For each of the $N$ source language words in the dictionary ($w_{1,s}$), we also precompute a set of $N_n$ hard negative samples: these words from $\mathcal{V}_t$ are highly similar or related to $w_{1,s}$ but are not direct word translations of $w_{1,s}$, and should be 'repelled' (i.e., pushed further away) from the word $w_{1,s}$ in the fine-tuned semantic space. These hard negatives are obtained based on the cosine similarity scores of the encodings from the input multilingual

---

[2]For instance, in case of multilingual BERT $[SPEC1]$ is the special [CLS] token, while $[SPEC2]$ is the [SEP] token.

[3]Put simply, we process sub-sentential text input in the same way as multilingual SEs process sentence-level input. We also experimented with other approaches to obtaining the final encoding from prior work, such as taking the representation of the [CLS] token from the last layer (Liu et al., 2021b; Li et al., 2022); however, in our preliminary experiments, we obtained slightly better results with mean-averaging.

[4]Such bilingual dictionaries are one of the most widespread and cheapest-to-build resources in multilingual NLP (Ruder et al., 2019; Wang et al., 2022).

encoder prior to its contrastive fine-tuning.

We then rely on the variant of the standard and widely used multiple negatives ranking loss (MNEG) (Cer et al., 2018; Henderson et al., 2019, 2020), which combines **1)** the $N_n$ hard negatives per each positive example with **2)** random in-batch negatives. The aim of the loss, when adapted to word-level inputs (Vulić et al., 2021a), is to rank true $L_t$ word translations from $\mathcal{D}$ over hard negatives and randomly paired $L_t$ words. The similarity between any two words $w_i$ and $w_j$ is quantified via the similarity function $S$ operating on their encodings $S(f_\theta(w_i), f_\theta(w_j))$. We use the scaled cosine similarity following (Henderson et al., 2019; Vulić et al., 2021a): $S(f_\theta(w_i), f_\theta(w_j)) = C \cdot cos(f_\theta(w_i), f_\theta(w_j))$, where $C$ is the scaling constant. Contrastive fine-tuning with MNEG then proceeds in batches of $B$ positive pairs $(w_i, v_i), \ldots, (w_B, v_B)$ from $\mathcal{D}$. The MNEG loss for a single batch, when also taking into account $N_n$ hard negatives per each positive example, is computed as follows:

$$\mathcal{L} = -\sum_{i=1}^{B} S(f_\theta(w_i), f_\theta(v_i)) \qquad \text{(positives)}$$
$$+ \sum_{i=1}^{B} \log \sum_{j=1, j\neq i}^{B} e^{S(f_\theta(w_i), f_\theta(v_j))} \quad \text{(in-batch negatives)}$$
$$+ \sum_{i=1}^{B} \log \sum_{k=1}^{N_n} e^{S(f_\theta(w_i), f_\theta(v_{k,i}))} \quad \text{(hard negatives)}$$

where $v_{k,i}$ denotes the $k$-th hard negative from the language $L_t$ for the $L_s$ word $w_i$.[5] The fine-tuned input encoder, now with parameters $\theta'$, can again be used to obtain the type-level representation of any input word $w$: $f_{\theta'}(w)$, see Figure 1.

**Interpolation with Static CLWEs.** Li et al. (2022) recently demonstrated that further performance gains in the BLI task might be achieved by combining the type-level output of the encoding function $f$ with static CLWEs, but they experimented only with multilingual LMs. Static CLWEs and encoder-based representations of the same set of words can be seen as two different views of the same data point. Following Li et al. (2022), we learn an additional linear orthogonal mapping from the static cross-lingual WE space – e.g., a CLWE space induced from monolingual fastText embeddings (Bojanowski et al., 2017) using VECMAP

(Artetxe et al., 2018) – into the cross-lingual space spanned by the encoder function $f$. The mapping transforms $\ell_2$-normed $d_1$-dimensional static cross-lingual WEs into $d_2$-dimensional cross-lingual WEs obtained through the multilingual encoder (tuned $f_{\theta'}$ or original $f_\theta$).

Learning the linear map $\boldsymbol{W} \in \mathbb{R}^{d_1 \times d_2}$, when $d_1 < d_2$,[6] is formulated as a Generalised Procrustes problem (Schönemann, 1966; Viklands, 2006). It operates on all (i.e., both $L_s$ and $L_t$) words from the seed translation dictionary $\mathcal{D}$.[7]

Unless noted otherwise, a final representation of an input word $w$ is then computed as follows:

$$(1 - \lambda) \frac{\mathbf{s}_w \boldsymbol{W}}{\|\mathbf{s}_w \boldsymbol{W}\|_2} + \lambda \frac{f_\theta(w)}{\|f_\theta(w)\|_2}, \qquad (1)$$

where $\lambda$ is a tunable interpolation hyper-parameter, $\mathbf{s}_w$ denotes the static CLWE of $w$, and $f_\theta(w)$ can be replaced by $f_{\theta'}(w)$. This simple procedure yields an 'interpolated' shared cross-lingual WE space.

## 3 Experimental Setup

**Multilingual Sentence Encoders.** We probe two widely used multilingual SEs: **1)** Language-agnostic BERT Sentence Embedding (**LaBSE**) (Feng et al., 2022) which adapts pretrained multilingual BERT (**mBERT**) (Devlin et al., 2019) into a multilingual SE; **2)** Multilingual **xMP-NET** is a distillation-based adaptation (Reimers and Gurevych, 2020) of **XLM-R** (Conneau et al., 2020) as the student model into a multilingual SE, based on monolingual MPNet encoder (Song et al., 2020) as the teacher model. LaBSE is the current state-of-the-art multilingual SE and supports 109 languages, while xMPNET is the best-performing multilingual SE in the Sentence-BERT repository (Reimers and Gurevych, 2019).[8] Along with multilingual SEs LaBSE and xMPNET, we also experiment with the original multilingual LMs – mBERT and XLM-R – in the same training and evaluation protocols (see Figure 1 and §2), aiming to quantify: (i) the extent to which cross-lingual lexical knowledge can be exposed from LMs that have not been

---

specialized for sentence-level semantics, as well as (ii) the increase in quality of lexical knowledge when multilingual LMs get transformed into SEs.

**Evaluation Tasks.** We evaluate on the standard and diverse cross-lingual lexical semantic tasks:

**Task 1: Bilingual Lexicon Induction (BLI)**, a standard task to assess the "semantic quality" of static cross-lingual word embeddings (CLWEs) (Ruder et al., 2019), allows us to **1)** directly assess the extent to which cross-lingual word translation knowledge can be exposed from multilingual LMs and SEs; **2)** immediately test the ability to transform multilingual sentence encoders into bilingual lexical encoders; **3)** investigate different model variants. We run a series of BLI evaluations on two standard BLI evaluation benchmarks. First, GT-BLI (Glavaš et al., 2019), constructed semi-automatically from Google Translate, comprises 28 language pairs with a good balance of typologically similar and distant languages (Croatian: HR, English: EN, Finnish: FI, French: FR, German: DE, Italian: IT, Russian: RU, Turkish: TR). Second, PanLex-BLI (Vulić et al., 2019) focuses on BLI evaluation for lower-resource languages, deriving training and test data from PanLex (Kamholz et al., 2014). We evaluate on 10 pairs comprising the following languages: Bulgarian (BG), Catalan (CA), Estonian (ET), Hebrew (HE), and Georgian (KA).

Standard BLI setups and data are adopted: 5k training word pairs are used as seed dictionary $\mathcal{D}$, and another 2k pairs as test data. Note that $\mathcal{D}$ is used to contrastively fine-tune multilingual encoders as well as to learn the mapping function of VECMAP. The evaluation metric is standard Precision@1 (P@1).[9] For PanLex-BLI, we also run additional tests with a smaller dictionary $\mathcal{D}$ spanning 1k translation pairs.

**Task 2: Cross-Lingual Lexical Semantic Similarity (XLSIM)** is another established cross-lingual lexical task. We use the recent comprehensive XLSIM benchmark Multi-SimLex (Vulić et al., 2020a), which comprises cross-lingual datasets of 2k-4k scored word pairs over 66 language pairs. We evaluate on a subset of language pairs shared with the GT-BLI dataset: EN, FI, RU, FR.[10] To avoid any test data leakage, we remove all XLSIM

test pairs from the bilingual dictionary $\mathcal{D}$ prior to fine-tuning.

**Task 3: Cross-Lingual Entity Linking (XEL)** is a standard task in knowledge base (KB) construction (Zhou et al., 2022), where the goal is to link an entity mention in any language to a corresponding entity in an English KB or in a language-agnostic KB.[11] We evaluate on the cross-lingual biomedical entity linking (XL-BEL) benchmark of Liu et al. (2021d): it requires the model to link an entity mention to entries in UMLS (Bodenreider, 2004), a language-agnostic medical knowledge base. We largely follow the XL-BEL experimental setup of Liu et al. (2021d) and probe the encoder models first *without* any additional fine-tuning on UMLS data, and then *with* subsequent UMLS fine-tuning (i) only on the EN UMLS data, (ii) on all the UMLS data in 10 languages of the XL-BEL dataset.[12] Due to a large number of experiments, we again focus on the subset of languages in XL-BEL shared with GT-BLI: EN, DE, FI, RU, TR. Standard P@1 and P@5 scores are reported.

**Static CLWEs and Word Vocabularies.** As monolingual static WEs, we select CommonCrawl fastText vectors (Bojanowski et al., 2017) of the top 200k most frequent word types in the training data, following prior work on learning static CLWEs (Conneau et al., 2018; Artetxe et al., 2018; Heyman et al., 2019).[13] Static CLWEs are then induced via the standard and popular supervised mapping-based VECMAP method (Artetxe et al., 2018), leveraging the seed dictionary $\mathcal{D}$. These CLWEs are used for interpolation with encoder-based WEs (see §3) but also as the baseline approaches for BLI and XLSIM tasks. We compute the type-level (ISO) WEs from multilingual LMs and SEs for the same 200K most frequent words of each language.

**Technical Details and Hyperparameters.** The implementation is based on the SBERT framework (Reimers and Gurevych, 2019), using the suggested settings: AdamW (Loshchilov and Hutter, 2018); learning rate of $2e-5$; weight decay rate of 0.01. We run contrastive fine-tuning for 5 epochs with all

---

[9]We observed very similar performance trends for P@5 and Mean Reciprocal Rank (MRR) as BLI measures.

[10]The evaluation metric is the standard Spearman's rank correlation between the average of gold human-elicited XL-SIM scores for word pairs and the cosine similarity between their respective word representations.

[11]Following prior work (Liu et al., 2021b; Zhou et al., 2022), XEL in this work also refers only to entity mention *disambiguation*; it does not cover the mention detection subtask.

[12]See (Liu et al., 2021d) for additional details.

[13]CommonCrawl-based fastText WEs typically outperform other popular choice for monolingual WEs: Wikipedia-based fastText (Glavaš et al., 2019; Li et al., 2022). We note that the main trends in our results also extend to the Wiki-based WEs.

| Multilingual LMs | | мBERT | | | | XLM-R | | | |
|---|---|---|---|---|---|---|---|---|---|
| Config → | VecMap | noCL (1.0) | noCL (λ) | +CL (1.0) | +CL (λ) | noCL (1.0) | noCL (λ) | +CL (1.0) | +CL (λ) |
| [BLI] λ=0.3 | 42.7 | 9.0 | 39.2 | 22.3 | 44.3 | 6.4 | 33.7 | 21.2 | 43.8 |
| [XLSIM] λ=0.5 | 45.8 | 5.7 | 35.4 | 38.4 | 48.1 | 1.7 | 23.5 | 46.1 | 51.8 |

| Multilingual SEs | | LaBSE | | | | xMPNET | | | |
|---|---|---|---|---|---|---|---|---|---|
| Config → | VecMap | noCL (1.0) | noCL (λ) | +CL (1.0) | +CL (λ) | noCL (1.0) | noCL (λ) | +CL (1.0) | +CL (λ) |
| [BLI] λ=0.3 | 42.7 | 21.4 | 45.7 | 30.8 | 49.1 | 17.0 | 41.7 | 28.6 | 47.9 |
| [XLSIM] λ=0.5 | 45.8 | 50.4 | 54.9 | 48.8 | 54.1 | 51.3 | 56.6 | 49.6 | 54.5 |

Table 1: (a) P@1 scores (×100%) averaged across all 28 language pairs in the GT-BLI dataset ([BLI] rows); (b) Spearman's $\rho$ correlation scores (×100) averaged across a subset of 6 language pairs from Multi-SimLex ([XLSIM rows]). See §3 for the description of different model configurations/variants. $|\mathcal{D}| = 5k$. The number in the parentheses denotes the value for $\lambda$ (see §3), which differs between the two tasks (0.3 for BLI and 0.5 for XLSIM). The $\lambda$ value of 1.0 effectively means 'no interpolation' with static VecMap CLWEs. Individual results per each language pair in both tasks and with other $\lambda$s are in Appendix B and Appendix C.

| | | мBERT | | XLM-R | | LaBSE | | xMPNET | |
|---|---|---|---|---|---|---|---|---|---|
| Pair ↓ / Config → | VecMap | +CL (1.0) | +CL (0.4) | +CL (1.0) | +CL (0.4) | +CL (1.0) | +CL (0.4) | +CL (1.0) | +CL (0.4) |
| BG–CA | 34.4 | 9.6 | 31.9 | 13.2 | 33.3 | 17.9 | **38.0** | 15.9 | 35.7 |
| BG–ET | 30.0 | 17.1 | 32.6 | 21.3 | 34.1 | 29.9 | **42.7** | 26.1 | 38.9 |
| BG–HE | 26.1 | 9.9 | 21.1 | 10.5 | 26.3 | 23.7 | **37.2** | 10.9 | 27.2 |
| BG–KA | 26.8 | 16.0 | 29.8 | 15.9 | 30.5 | 27.2 | **37.4** | 18.7 | 32.4 |
| CA–ET | 26.3 | 26.8 | 32.9 | 23.5 | 34.1 | 28.8 | **38.6** | 29.0 | 38.9 |
| CA–HE | 23.3 | 2.3 | 12.5 | 4.9 | 18.5 | 12.7 | **28.9** | 8.5 | 22.7 |
| CA–KA | 20.7 | 1.5 | 10.3 | 4.7 | 20.0 | 9.6 | **26.1** | 6.4 | 21.8 |
| ET–HE | 18.6 | 15.0 | 21.9 | 17.7 | 26.0 | 31.0 | **37.8** | 18.5 | 27.0 |
| ET–KA | 16.5 | 7.2 | 18.2 | 12.7 | 24.3 | 19.3 | **30.3** | 12.5 | 25.8 |
| HE–KA | 12.7 | 15.6 | 23.8 | 13.3 | 23.1 | 25.3 | **30.2** | 15.1 | 24.4 |
| **Average** | 23.5 | 12.1 | 23.5 | 13.8 | 27.0 | 22.5 | **34.7** | 16.2 | 29.5 |

Table 2: P@1 scores over a representative subset of 10 language pairs from the PanLex-BLI dataset of Vulić et al. (2019). See §3 for the description of different model configurations/variants. $|\mathcal{D}| = 5k$. Highest scores per row are in **bold**. Respective average scores for the *noCL (1.0)* config (i.e., without contrastive learning and without interpolation with static VecMap CLWEs) are: 4.2 (мBERT), 3.1 (XLM-R), 17.0 (LaBSE), 8.3 (xMPNET).

the models, with the batch size of $B = 128$ positive examples for MNEG. The number of hard negatives per each positive example is set to $N_n = 10$ (see §2.1).[14] Since standard BLI and XLSIM datasets lack a validation portion (Ruder et al., 2019), we follow prior work (Glavaš et al., 2019) and tune hyperparameters on a *single, randomly selected* language pair from each dataset, and use those hyperparameter values in all other runs.

**Model Configurations.** They are labelled as follows: **ENC-{noCL,+CL}** ($\lambda$), where (i) ENC denotes the input multilingual Transformer, which can be a multilingual LM (мBERT, XLM-R), or a multilingual SE (LaBSE, xMPNET), (ii) 'noCL' refers to using the input model 'off-the-shelf' without any contrastive fine-tuning, while '+CL' variants apply contrastive fine-tuning, and (iii) $\lambda$ is the interpolation with the static CLWE space, obtained with VecMap (see again Figure 1).[15]

**Important Disclaimer.** We note that the main pur-

---

[14]We also tested larger values: $N_n = \{20, 30, 50\}$. They slow down fine-tuning while offering very small and negligible gains in performance.

[15]Note that $\lambda = 1.0$ implies no interpolation with static CLWE space, i.e., WEs come directly from the LM/SE.

| | | LaBSE | | |
|---|---|---|---|---|
| Pair ↓ | VecMap | noCL (1.0) | +CL (1.0) | +CL (0.5) |
| BG–CA | 15.2 | 14.0 | 18.0 | **28.9** |
| BG–ET | 12.5 | 20.3 | 25.5 | **35.1** |
| BG–HE | 5.6 | 18.3 | 20.8 | **24.7** |
| BG–KA | 9.1 | 16.0 | 21.6 | **29.7** |
| CA–ET | 9.8 | 24.8 | 25.6 | **31.2** |
| CA–HE | 5.0 | 10.6 | 12.2 | **15.6** |
| CA–KA | 5.5 | 5.7 | 8.3 | **14.8** |
| ET–HE | 3.1 | **27.7** | 25.1 | 25.4 |
| ET–KA | 4.6 | 13.2 | 16.0 | **21.1** |
| HE–KA | 3.2 | 19.0 | 22.0 | **25.4** |
| **Average** | 7.4 | 17.0 | 19.5 | **25.2** |

Table 3: P@1 scores over 10 language pairs from the PanLex-BLI dataset of Vulić et al. (2019) when $|\mathcal{D}| = 1k$, with different model variants based on LaBSE (see §3). Highest scores per row are in **bold**.

pose of the chosen evaluation tasks and experimental protocols is not necessarily achieving state-of-the-art performance, but rather probing different model variants in different cross-lingual lexical tasks, and offering fair and insightful comparisons.

## 4 Results and Discussion

**Bilingual Lexicon Induction (BLI).** Table 1 displays our main BLI results, aggregated over all 28 language pairs of GT-BLI (Glavaš et al., 2019), for two multilingual LMs (мBERT and XLM-R) and

their SE counterparts (LABSE and XMPNET).

Two trends hold across the board. First, multilingual SEs, LABSE and XMPNET, outperform their multilingual LM counterparts, MBERT and XLM-R, across the board. The gains are visible in all four experimental configurations (with/without contrastive cross-lingual lexical specialisation × with/without interpolation with the VECMAP CLWE space). This confirms our intuition that multilingual SEs, having been (additionally) trained on parallel data (Feng et al., 2022; Reimers and Gurevych, 2020), should better reflect the cross-lingual alignments at the lexical level than the off-the-shelf multilingual LMs, which have not been exposed to any cross-lingual signal at pretraining. The poor cross-lingual lexical alignment in the representation spaces of MBERT and XLM-R is also reflected by the fact that for those encoders, we only surpass the baseline VECMAP performance by a small margin (+1.1 for XLM-R, +1.6 for MBERT) when we apply the contrastive lexical fine-tuning *and* interpolate word vectors obtained with fine-tuned encoders with VECMAP vectors.

The behaviour of SEs, on the other hand, is much more favourable. For example, LABSE surpasses baseline VECMAP performance with interpolation alone, even without the additional lexical contrastive fine-tuning procedure. When contrastively fine-tuned (and then interpolated with VECMAP) both LABSE and XMPNET surpass the baseline VECMAP performance by a wide margin (+5.2 for XMPNET and +6.4 for LABSE). This also confirms that more cross-lingual lexical knowledge can be exposed from the multilingual SEs via contrastive fine-tuning.

Another finding is that both (i) contrastive cross-lingual lexical learning (+CL) and (ii) interpolation with VECMAP consistently improve the performance for all four encoders. We reach the best performance for all four encoders when combining the contrastive cross-lingual lexical fine-tuning with the interpolation with static CLWEs (+CL ($\lambda$)). The contrastive fine-tuning step crucially contributes to the overall performance: compared to interpolation alone (noCL ($\lambda$)), contrastive cross-lingual lexical fine-tuning (+CL ($\lambda$)) brings an average performance gains of over 6 BLI points.

Table 2 shows the BLI results for all four encoders on 10 low(er)-resource language pairs from PanLex-BLI (Vulić et al., 2019). While overall relative trends are similar to those observed

for high(er)-resource languages in GT-BLI (Table 3), the gains stemming from cross-lingual contrastive lexical fine-tuning are substantially larger for this set of languages. The best-performing configuration – contrastive fine-tuning and interpolation (+CL (0.4)) applied on LaBSE – surpasses VECMAP by 11 BLI points on average (compared to 6 points on GT-BLI), with gains for some language pairs (e.g., HE-KA, ET-HE) approaching the impressive margin of 20 BLI points. This finding indicates that cross-lingual lexical knowledge stored in multilingual SEs is even more crucial when dealing with lower-resource languages.

In Table 3 we compare the results of LaBSE (as the best-performing multilingual SE) against VECMAP on PanLex-BLI in a scenario with less external bilingual supervision: $|\mathcal{D}| = 1k$. Interestingly, in this setup LABSE already substantially outperforms VECMAP out of the box (in its noCL (1.0) configuration); contrastive lexical fine-tuning (+CL (1.0)) and additional interpolation with VECMAP embeddings (+CL (0.5)) again bring further substantial gains, and we again observe a strong synergistic effect of the two components: +CL (1.0) yields gains over noCL (1.0) for 9/10 language pairs, and +CL (0.5) yields further gains for all 10/10 pairs. Furthermore, the contrastively fine-tuned LABSE is much more resilient to training data scarcity than VECMAP: a five-fold reduction of the training dictionary size (from 5k to 1k) reduces the performance of LABSE +CL ($\lambda$) by 27% (from 34.7 to 25.2 P@1 points) compared to a massive drop of almost 70% for VECMAP (from 23.5 to mere 7.4 P@1).

**Cross-Lingual Lexical Semantic Similarity (XLSIM).** The average XLSIM results are summarised in Table 1. They again corroborate one of the main findings from BLI experiments: multilingual SEs store more cross-lingual lexical knowledge than multilingual LMs. This is validated by substantial gains of SEs over corresponding LMs across all configurations in Table 1. Interestingly, due to their contrastive learning objectives on sentence-level parallel data (Feng et al., 2022), LABSE and XMPNET provide very strong XLSIM results when used off-the-shelf: their noCL (1.0) configuration already outperforms the VECMAP vectors. Additional lexical fine-tuning with 5k word translation pairs does not bring any performance benefits. However, the opposite is true for multilingual LMs: even fine-tuning with
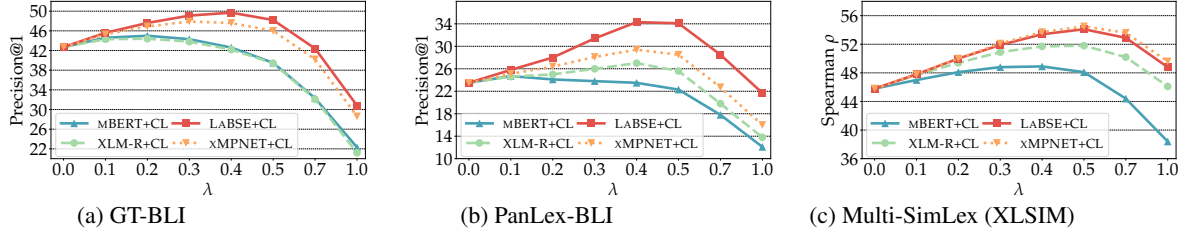
(a) GT-BLI  (b) PanLex-BLI  (c) Multi-SimLex (XLSIM)

Figure 2: Average scores across different interpolation values $\lambda$ for the BLI task on **(a)** GT-BLI and **(b)** PanLex-BLI, and **(c)** the XLSIM task on Multi-SimLex. $|\mathcal{D}| = 5k$. Additional results are in Appendix B and C.

| Config ↓ / Language $L_t \rightarrow$ | DE | | FI | | RU | | TR | |
|---|---|---|---|---|---|---|---|---|
| | P@1 | P@5 | P@1 | P@5 | P@1 | P@5 | P@1 | P@5 |
| XLM-R +noCL | 0.0 | 0.1 | 0.1 | 0.2 | 0.1 | 0.2 | 0.4 | 0.5 |
| XLM-R +noCL+UMLS$_{EN}$ | 27.6 | 32.0 | 12.2 | 14.7 | 21.8 | 25.9 | 29.3 | 35.9 |
| XLM-R +noCL+UMLS$_{all}$ | 31.8 | 37.3 | 18.6 | 22.2 | 35.4 | 41.2 | 42.8 | 48.9 |
| XLM-R +CL | 14.1 | 17.1 | 5.0 | 6.5 | 8.7 | 11.2 | 21.6 | 27.1 |
| XLM-R +CL+UMLS$_{EN}$ | 25.2 | 29.0 | 12.1 | 14.1 | 19.8 | 25.0 | 31.1 | 36.1 |
| XLM-R +CL+UMLS$_{all}$ | 32.1 | 36.7 | 19.1 | 23.8 | 34.9 | 42.4 | 43.4 | 49.0 |
| xMPNET +noCL | 19.5 | 25.9 | 12.2 | 14.8 | 19.2 | 24.3 | 28.9 | 36.3 |
| xMPNET +noCL+UMLS$_{EN}$ | 25.1 | 29.2 | 17.8 | 21.5 | 21.9 | 26.9 | 30.0 | 36.5 |
| xMPNET +noCL+UMLS$_{all}$ | **33.4** | 37.8 | **23.6** | **27.7** | **39.8** | 45.4 | **44.6** | **51.4** |
| xMPNET +CL | 20.8 | 26.5 | 9.1 | 12.5 | 12.8 | 17.1 | 30.4 | 36.5 |
| xMPNET +CL+UMLS$_{EN}$ | 25.1 | 28.7 | 11.4 | 14.0 | 21.8 | 27.2 | 31.0 | 37.5 |
| xMPNET +CL+UMLS$_{all}$ | 32.0 | **38.7** | 22.9 | 27.5 | 39.2 | **45.7** | 44.3 | 51.0 |

Table 4: A summary of results in the XEL task on the biomedical XL-BEL benchmark of Liu et al. (2021d). We show the results of the better-performing LM (XLM-R), and the more lightweight multilingual SE (xMPNET).

the small set of word translation pairs brings large benefits in the XLSIM task (e.g., compare noCL (1.0) and +CL (1.0) configurations for MBERT and XLM-R), and turns them into more effective lexical encoders. This result corroborates a similar finding from prior work in monolingual setups (Vulić et al., 2021a). Finally, interpolation with static CLWEs is again beneficial for the final XL-SIM performance with all four underlying multilingual models, with highest scores achieved with interpolated vectors ($\lambda = 0.5$) across the board, and substantial gains both over VECMAP and over configurations with $\lambda = 1.0$.

**Interpolation with Static CLWEs.** A more detailed analysis over different $\lambda$ values for BLI and XLSIM, summarised in Figure 2, reveals that interpolation can offer large performance benefits, especially in the $\lambda$ interval of $[0.3, 0.5]$, but the optimal interpolation values are task- and even dataset-dependent. For instance, for low-resource BLI on PanLex BLI more knowledge comes from the multilingual encoders as VECMAP CLWEs are of lower quality for such languages: in consequence, the optimal $\lambda$ value 'moves away' from the static CLWEs towards encodings obtained by fine-tuned multilingual SE. We also note that larger benefits of interpolation are observed when VECMAP CLWEs are combined with contrastively fine-tuned multilingual SEs: cf., the large gains in Figure 2b and in

Table 3 for the LABSE +CL model variant.

**Cross-Lingual Entity Linking (XEL).** Experiments on XL-BEL (Liu et al., 2021d), summarised in Table 4, demonstrate that additional contrastive tuning with word or phrase pairs can greatly boost performance of multilingual LMs: even fine-tuning with 5k word translation pairs without any domain-specific knowledge yields strong benefits for XLM-R. As expected, using a much larger and domain-specific external database UMLS yields much higher scores and is more crucial for performance. In fact, contrastively fine-tuning on UMLS generally improves XEL performance with all four underlying models. Again, we observe that xMPNET-based configurations outperform XLM-R-based configurations across the board. This finding again indicates that multilingual SEs store more cross-lingual lexical knowledge than multilingual LMs: this difference is particularly salient when the models are used off-the-shelf without any additional contrastive fine-tuning, and xMPNET retains higher performance even after fine-tuning with the UMLS-based domain-specific knowledge.

Further, we note that xMPNET-based XEL models outperform similar XEL models based on language-specific LMs, evaluated previously by Liu et al. (2021d), for FI and RU, and score on-par for TR. This is consistent with our hypothesis that multilingual SEs store rich multilingual lex-

ical knowledge, which then also enables them to benefit from the multilingual UMLS knowledge and to expose it further via contrastive UMLS-based fine-tuning. As mentioned, using UMLS synonyms yields higher gains than using a smaller set of generic word translations (+CL configurations). This suggests that in specialised domains such as biomedicine, in-domain multilingual lexical pairs are extremely valuable and should be preferred for exposing the lexical knowledge from multilingual SEs.

## 5 Conclusion and Future Work

In this work, we investigated strategies to expose cross-lingual lexical knowledge from pretrained models, including multilingual language models (LMs) and multilingual sentence encoders (SEs). Based on an extensive evaluation on a suite of cross-lingual lexical tasks spanning word translation, semantic similarity, and entity linking, we verified that multilingual sentence encoders (e.g., LABSE and XMPNET) are superior to multilingual language models (MBERT and XLM-R) in terms of stored cross-lingual lexical knowledge. Moreover, we proposed new methods to further fine-tune their representations based on contrastive learning to 'rewire' the models' parameters and transform them from LMs and SEs into more effective multilingual lexical encoders. We also empirically demonstrated that for some lexical tasks performance can be further boosted through interpolation of type-level lexical encodings with static cross-lingual word embeddings. These yield gains for all underlying models, which are especially significant for resource-poor languages and in low-data learning regimes.

While our work has focused on two widely used state-of-the-art multilingual SEs, the contrastive fine-tuning framework is versatile and model-independent and can be directly applied on top of other multilingual SEs in future work. We will also investigate other more sophisticated contrastive learning strategies, look into ensembling of knowledge extracted from different SEs, and expand our evaluation to more tasks and languages.

## Acknowledgements

## References

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 ACL 2018*, pages 789–798.

Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the ACL*, 7:597–610.

Olivier Bodenreider. 2004. The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the ACL*, 5:135–146.

Rishi Bommasani, Kelly Davis, and Claire Cardie. 2020. Interpreting Pretrained Contextualized Representations via Reductions to Static Embeddings. In *Proceedings of ACL 2020*, pages 4758–4781.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Universal sentence encoder for English. In *Proceedings of EMNLP 2018*, pages 169–174.

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 ACL 2020*, pages 8440–8451.

Alexis Conneau and Guillaume Lample. 2019. Cross-lingual language model pretraining. In *Proceedings of NeurIPS 2019*, pages 7057–7067.

Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. Word translation without parallel data. In *Proceedings of ICLR 2018*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of

deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT 2019*, pages 4171–4186.

Kawin Ethayarajh. 2019. How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings. In *Proceedings of EMNLP 2019*, pages 55–65.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of ACL 2022*.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of EMNLP 2021*, pages 6894–6910.

Goran Glavaš and Ivan Vulić. 2020. Non-linear instance-based cross-lingual mapping for non-isomorphic embedding spaces. In *Proceedings of ACL 2020*, pages 7548–7555.

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 *Proceedings of ACL 2019*, pages 710–721.

Hila Gonen, Shauli Ravfogel, Yanai Elazar, and Yoav Goldberg. 2020. It's not Greek to mBERT: Inducing word-level translations from multilingual BERT. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 45–56.

Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić, Pei-Hao Su, Tsung-Hsien Wen, and Ivan Vulić. 2020. ConveRT: Efficient and accurate conversational representations from transformers. In *Findings of EMNLP 2020*, pages 2161–2174.

Matthew Henderson, Ivan Vulić, Daniela Gerz, Iñigo Casanueva, Paweł Budzianowski, Sam Coope, Georgios Spithourakis, Tsung-Hsien Wen, Nikola Mrkšić, and Pei-Hao Su. 2019. Training neural response selection for task-oriented dialogue systems. In *Proceedings of ACL 2019*, pages 5392–5404.

Geert Heyman, Bregt Verreet, Ivan Vulić, and Marie-Francine Moens. 2019. Learning unsupervised multilingual word embeddings with incremental multilingual hubs. In *Proceedings of NAACL-HLT 2019*, pages 1890–1902.

Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does BERT learn about the structure of language? In *Proceedings of ACL 2019*, pages 3651–3657.

David Kamholz, Jonathan Pool, and Susan M. Colowick. 2014. PanLex: Building a resource for panlingual lexical translation. In *Proceedings of LREC 2014*, pages 3145–3150.

Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. On the sentence embeddings from pre-trained language models. In *Proceedings of EMNLP 2020*, pages 9119–9130.

Yaoyiran Li, Fangyu Liu, Nigel Collier, Anna Korhonen, and Ivan Vulić. 2022. Improving word translation via two-stage contrastive learning. In *Proceedings of ACL 2022*.

Robert Litschko, Ivan Vulić, Simone Paolo Ponzetto, and Goran Glavaš. 2022. On cross-lingual retrieval with multilingual text encoders. *Information Retrieval*.

Che Liu, Rui Wang, Jinghua Liu, Jian Sun, Fei Huang, and Luo Si. 2021a. DialogueCSE: Dialogue-based contrastive learning of sentence embeddings. In *Proceedings of EMNLP 2021*, pages 2396–2406.

Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. 2021b. Self-alignment pretraining for biomedical entity representations. In *Proceedings of NAACL-HLT 2021*, pages 4228–4238.

Fangyu Liu, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2021c. Fast, effective, and self-supervised: Transforming masked language models into universal lexical and sentence encoders. In *Proceedings of EMNLP 2021*, pages 1442–1459.

Fangyu Liu, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2021d. Learning domain-specialised representations for cross-lingual biomedical entity linking. In *Proceedings of ACL-IJCNLP 2021*, pages 565–574.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.

Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *Proceedings of ICLR 2018*.

Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. 2013. Exploiting similarities among languages for machine translation. *arXiv preprint, CoRR*, abs/1309.4168.

Kevin Musgrave, Serge J. Belongie, and Ser-Nam Lim. 2020. A metric learning reality check. In *Proceedings of ECCV 2020*, volume 12370, pages 681–699.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:140:1–140:67.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of EMNLP 2019*, pages 3982–3992.

Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In *Proceedings of EMNLP 2020*, pages 4512–4525.

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: what we know about how BERT works. *Transactions of the ACL*, 8:842–866.

Sebastian Ruder, Ivan Vulić, and Anders Søgaard. 2019. A survey of cross-lingual embedding models. *Journal of Artificial Intelligence Research*, 65:569–631.

Peter H Schönemann. 1966. A generalized solution of the orthogonal Procrustes problem. *Psychometrika*, 31(1):1–10.

Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. MPNet: Masked and permuted pretraining for language understanding. In *Proceedings of NeurIPS 2020*.

Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748.

Thomas Viklands. 2006. *Algorithms for the weighted orthogonal Procrustes problem and other least squares problems*. Ph.D. thesis, Datavetenskap.

Ivan Vulić, Simon Baker, Edoardo Maria Ponti, Ulla Petti, Ira Leviant, Kelly Wing, Olga Majewska, Eden Bar, Matt Malone, Thierry Poibeau, Roi Reichart, and Anna Korhonen. 2020a. Multi-SimLex: A large-scale evaluation of multilingual and crosslingual lexical semantic similarity. *Computational Linguistics*, 46(4):847–897.

Ivan Vulić, Goran Glavaš, Roi Reichart, and Anna Korhonen. 2019. Do we really need fully unsupervised cross-lingual embeddings? In *Proceedings of EMNLP-IJCNLP 2019*, pages 4407–4418.

Ivan Vulić, Edoardo Maria Ponti, Anna Korhonen, and Goran Glavaš. 2021a. LexFit: Lexical fine-tuning of pretrained language models. In *Proceedings of ACL-IJCNLP 2021*, pages 5269–5283.

Ivan Vulić, Edoardo Maria Ponti, Robert Litschko, Goran Glavaš, and Anna Korhonen. 2020b. Probing pretrained language models for lexical semantics. In *Proceedings of EMNLP 2020*, pages 7222–7240.

Ivan Vulić, Pei-Hao Su, Samuel Coope, Daniela Gerz, Paweł Budzianowski, Iñigo Casanueva, Nikola Mrkšić, and Tsung-Hsien Wen. 2021b. ConFiT: Conversational fine-tuning of pretrained language models. In *Proceedings of EMNLP 2021*, pages 1151–1168.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *Proceedings of NeurIPS 2019*, pages 3261–3275.

Xinyi Wang, Sebastian Ruder, and Graham Neubig. 2022. Expanding pretrained models to thousands more languages via lexicon-based adaptation. In *Proceedings of ACL 2022*.

Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. ConSERT: A contrastive framework for self-supervised sentence representation transfer. In *Proceedings of ACL-IJCNLP 2021*, pages 5065–5075.

Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo, Jax Law, Noah Constant, Gustavo Hernandez Abrego, Steve Yuan, Chris Tar, Yun-hsuan Sung, Brian Strope, and Ray Kurzweil. 2020. Multilingual universal sentence encoder for semantic retrieval. In *Proceedings of ACL 2020: System Demonstrations*, pages 87–94.

Wenxuan Zhou, Fangyu Liu, Ivan Vulić, Nigel Collier, and Muhao Chen. 2022. Prix-LM: Pretraining for multilingual knowledge base construction. In *Proceedings of ACL 2022*.

| Languages in: GT-BLI, Multi-SimLex, XL-BEL | |
| --- | --- |
| EN | English |
| DE | German |
| TR | Turkish |
| FI | Finnish |
| HR | Croatian |
| RU | Russian |
| IT | Italian |
| FR | French |
| **Languages in:** PanLex-BLI | |
| BG | Bulgarian |
| CA | Catalan |
| ET | Estonian |
| HE | Hebrew |
| KA | Georgian |

Table 5: Languages and their ISO 639-1 codes.

## A List of Languages

The list of languages used in this work, along with their ISO 639-1 codes, is available in Table 5.

## B BLI Results across Individual Language Pairs

Additional experiments and analyses over individual language pairs and other $\lambda$ values, which further support the main claims of the paper, have been relegated to the appendix for clarity and compactness of the presentation in the main paper:

**Table 6.** It provides results over all 28 language pairs in GT-BLI with 2 multilingual LMs and 2 multilingual SEs in the *noCL* variant without contrastive fine-tuning.

**Table 7.** It provides results over all 28 language pairs in GT-BLI with 2 multilingual LMs and 2 multilingual SEs in the *+CL* variant with contrastive fine-tuning.

**Table 8.** It provides results over all 28 language pairs in GT-BLI and across different $\lambda$ values with the LABSE +noCL variant.

**Table 9.** It provides results over all 28 language pairs in GT-BLI and across different $\lambda$ values with the LABSE +CL variant.

## C XLSIM Results across Individual Language Pairs

**Table 10.** It provides results over selected 6 language pairs from Multi-SimLex with 2 multilingual LMs and 2 multilingual SEs in the *noCL* variant without contrastive fine-tuning.

**Table 11.** It provides results over selected 6 language pairs from Multi-SimLex with 2 multilingual

LMs and 2 multilingual SEs in the *+CL* variant with contrastive fine-tuning.

**Table 12.** It provides results over selected 6 language pairs from Multi-SimLex and across different $\lambda$ values with the XMPNET +noCL variant.

**Table 13.** It provides results over selected 6 language pairs from Multi-SimLex and across different $\lambda$ values with the XMPNET +CL variant.

## D Models and Evaluation Data

URLs to the models used in this paper are provided in Table 14. Training and test data for all three tasks (BLI, XLSIM, XEL) is available online:

- GT-BLI is available here: `https://github.com/codogogo/xling-eval`

- PanLex-BLI: `https://github.com/cambridgeltl/panlex-bli`

- Multi-SimLex [XLSIM]: `https://multisimlex.com/`

- XL-BEL [XEL]: `https://github.com/cambridgeltl/sapbert`

Our code is based on PyTorch, and relies on the two following widely used repositories:

- `sentence-transformers` `www.sbert.net`

- `huggingface.co/transformers/`

| Pair ↓ / Config → | VECMAP | MBERT | | XLM-R | | LABSE | | XMPNET | |
|---|---|---|---|---|---|---|---|---|---|
| | | +noCL (1.0) | +noCL (0.3) | +noCL (1.0) | +noCL (0.3) | +noCL (1.0) | +noCL (0.3) | +noCL (1.0) | +noCL (0.3) |
| EN–DE | 55.6 | 15.6 | 50.7 | 12.7 | 45.5 | 25.4 | 54.1 | 22.0 | 47.7 |
| EN–TR | 40.4 | 7.2 | 34.9 | 6.0 | 28.9 | 23.6 | 42.1 | 16.1 | 33.7 |
| EN–FI | 45.6 | 7.9 | 38.7 | 6.6 | 33.4 | 19.3 | 45.1 | 14.8 | 39.5 |
| EN–HR | 37.5 | 8.9 | 31.8 | 6.8 | 25.6 | 24.7 | 45.3 | 18.5 | 38.9 |
| EN–RU | 45.6 | 3.2 | 40.7 | 0.9 | 34.7 | 24.7 | 49.9 | 17.6 | 41.7 |
| EN–IT | 60.2 | 12.3 | 57.1 | 9.3 | 53.0 | 26.4 | 62.3 | 23.5 | 58.8 |
| EN–FR | 64.1 | 25.2 | 62.5 | 19.5 | 56.6 | 34.2 | 67.5 | 29.2 | 61.7 |
| DE–TR | 32.5 | 9.1 | 28.9 | 6.9 | 24.4 | 17.7 | 33.0 | 13.1 | 29.7 |
| DE–FI | 39.7 | 9.2 | 34.1 | 7.3 | 30.1 | 16.0 | 37.4 | 13.3 | 34.7 |
| DE–HR | 33.3 | 11.5 | 31.0 | 9.7 | 25.7 | 19.2 | 38.5 | 14.9 | 34.1 |
| DE–RU | 40.0 | 4.2 | 36.4 | 0.9 | 32.4 | 14.3 | 41.5 | 9.5 | 37.7 |
| DE–IT | 49.5 | 10.9 | 45.9 | 8.1 | 42.7 | 19.4 | 51.4 | 18.3 | 49.1 |
| DE–FR | 50.0 | 15.8 | 49.7 | 10.5 | 42.3 | 22.5 | 53.2 | 20.2 | 49.8 |
| TR–FI | 31.3 | 6.7 | 26.2 | 5.2 | 22.0 | 15.5 | 31.7 | 11.3 | 30.9 |
| TR–HR | 25.4 | 10.8 | 24.3 | 8.3 | 20.5 | 18.7 | 33.1 | 14.4 | 28.2 |
| TR–RU | 32.9 | 2.6 | 29.1 | 0.8 | 25.5 | 14.1 | 36.9 | 11.3 | 33.5 |
| TR–IT | 37.1 | 7.9 | 34.7 | 5.5 | 27.4 | 17.0 | 38.9 | 14.8 | 38.3 |
| TR–FR | 39.4 | 7.8 | 37.3 | 5.7 | 30.9 | 20.9 | 43.1 | 16.6 | 39.4 |
| FI–HR | 30.4 | 7.2 | 26.6 | 5.5 | 22.8 | 17.4 | 36.4 | 12.2 | 32.7 |
| FI–RU | 38.2 | 2.6 | 34.0 | 0.9 | 30.5 | 15.1 | 41.0 | 9.0 | 37.1 |
| FI–IT | 39.9 | 7.9 | 36.7 | 6.8 | 30.4 | 18.1 | 43.4 | 16.4 | 41.8 |
| FI–FR | 42.8 | 7.5 | 38.9 | 5.9 | 32.2 | 18.6 | 45.9 | 16.2 | 42.2 |
| HR–RU | 40.6 | 6.0 | 35.8 | 1.6 | 30.4 | 24.5 | 45.7 | 16.3 | 41.4 |
| HR–IT | 40.4 | 11.2 | 39.0 | 8.4 | 31.3 | 24.5 | 47.9 | 22.4 | 44.5 |
| HR–FR | 43.6 | 9.7 | 42.3 | 6.1 | 30.6 | 25.7 | 50.0 | 19.8 | 43.9 |
| RU–IT | 46.6 | 3.1 | 42.2 | 1.5 | 36.4 | 21.8 | 47.6 | 18.4 | 46.5 |
| RU–FR | 48.7 | 4.1 | 44.3 | 1.5 | 38.4 | 26.6 | 50.9 | 18.9 | 47.5 |
| IT–FR | 64.1 | 16.6 | 62.8 | 9.3 | 58.6 | 33.9 | 65.6 | 27.5 | 63.1 |
| **Average** | 42.7 | 9.0 | 39.2 | 6.4 | 33.7 | 21.4 | 45.7 | 17.0 | 41.7 |

Table 6: Individual P@1 scores (×100%) for all 28 language pairs in the GT-BLI dataset of Glavaš et al. (2019), with multilingual LMs and SEs used 'off-the-shelf' *without contrastive fine-tuning* (§2). See §3 for the description of different model configurations/variants. $|\mathcal{D}| = 5k$. The number in the parentheses denotes the value for $\lambda$ (see §3): the value of 1.0 effectively means 'no interpolation' with static VECMAP CLWEs.

| Pair ↓ / Config → | VECMAP | MBERT | | XLM-R | | LABSE | | XMPNET | |
|---|---|---|---|---|---|---|---|---|---|
| | | +CL (1.0) | +CL (0.3) | +CL (1.0) | +CL (0.3) | +CL (1.0) | +CL (0.3) | +CL (1.0) | +CL (0.3) |
| EN–DE | 55.6 | 26.4 | 59.2 | 24.5 | 56.3 | 31.6 | 61.1 | 30.2 | 58.7 |
| EN–TR | 40.4 | 17.4 | 39.3 | 19.2 | 42.0 | 33.1 | 50.1 | 30.4 | 45.7 |
| EN–FI | 45.6 | 18.6 | 45.6 | 20.2 | 45.7 | 30.8 | 53.3 | 28.7 | 50.4 |
| EN–HR | 37.5 | 23.6 | 44.3 | 21.8 | 44.3 | 36.9 | 53.9 | 31.8 | 49.6 |
| EN–RU | 45.6 | 23.9 | 48.9 | 28.3 | 48.8 | 46.4 | 55.8 | 37.9 | 53.1 |
| EN–IT | 60.2 | 26.8 | 64.0 | 25.4 | 61.7 | 33.3 | 66.9 | 33.0 | 65.3 |
| EN–FR | 64.1 | 34.4 | 67.7 | 33.0 | 65.4 | 42.1 | 71.2 | 39.5 | 68.5 |
| DE–TR | 32.5 | 19.5 | 32.8 | 15.5 | 33.5 | 24.4 | 37.6 | 22.7 | 36.2 |
| DE–FI | 39.7 | 20.3 | 38.5 | 19.0 | 37.5 | 25.0 | 43.3 | 24.4 | 42.1 |
| DE–HR | 33.3 | 24.2 | 37.3 | 22.4 | 36.9 | 27.9 | 42.9 | 27.2 | 41.8 |
| DE–RU | 40.0 | 21.2 | 42.0 | 19.1 | 41.5 | 27.5 | 45.6 | 26.8 | 44.1 |
| DE–IT | 49.5 | 23.0 | 49.9 | 19.0 | 48.6 | 24.6 | 52.5 | 25.6 | 51.4 |
| DE–FR | 50.0 | 27.4 | 52.5 | 24.3 | 51.5 | 31.3 | 54.5 | 30.0 | 53.9 |
| TR–FI | 31.3 | 18.0 | 29.6 | 15.9 | 31.7 | 22.2 | 34.9 | 22.4 | 35.9 |
| TR–HR | 25.4 | 21.8 | 30.2 | 19.0 | 30.8 | 27.4 | 36.8 | 26.0 | 36.4 |
| TR–RU | 32.9 | 16.2 | 33.9 | 15.7 | 33.5 | 24.7 | 37.4 | 24.7 | 37.8 |
| TR–IT | 37.1 | 17.4 | 37.8 | 15.4 | 36.8 | 22.8 | 41.6 | 22.2 | 41.0 |
| TR–FR | 39.4 | 18.4 | 40.4 | 17.4 | 38.9 | 29.4 | 43.9 | 26.3 | 43.6 |
| FI–HR | 30.4 | 20.2 | 32.7 | 17.4 | 32.5 | 27.4 | 39.7 | 24.3 | 38.4 |
| FI–RU | 38.2 | 17.6 | 37.3 | 18.0 | 38.7 | 27.6 | 42.4 | 27.6 | 41.3 |
| FI–IT | 39.9 | 18.4 | 40.5 | 17.6 | 40.0 | 25.2 | 45.3 | 24.0 | 44.9 |
| FI–FR | 42.8 | 17.5 | 43.0 | 18.2 | 41.9 | 26.4 | 47.6 | 26.1 | 45.8 |
| HR–RU | 40.6 | 26.4 | 41.4 | 29.5 | 43.9 | 36.1 | 46.4 | 33.3 | 46.4 |
| HR–IT | 40.4 | 24.6 | 44.0 | 22.8 | 42.8 | 32.3 | 49.7 | 30.2 | 49.2 |
| HR–FR | 43.6 | 24.1 | 46.5 | 23.4 | 44.1 | 35.0 | 51.8 | 29.3 | 50.7 |
| RU–IT | 46.6 | 22.0 | 46.9 | 19.4 | 45.0 | 32.0 | 50.2 | 28.4 | 51.1 |
| RU–FR | 48.7 | 22.1 | 49.1 | 20.5 | 47.6 | 37.0 | 53.5 | 28.9 | 51.2 |
| IT–FR | 64.1 | 33.4 | 65.5 | 32.9 | 64.3 | 42.2 | 66.0 | 39.0 | 66.4 |
| **Average** | 42.7 | 22.3 | 44.3 | 21.2 | 43.8 | 30.8 | 49.1 | 28.6 | 47.9 |

Table 7: Individual P@1 scores (×100%) for all 28 language pairs in the GT-BLI dataset of Glavaš et al. (2019), with model variants *with contrastive fine-tuning* (§2). $|\mathcal{D}| = 5k$.

| Pair ↓ / λ = → | VECMAP(0.0) | LABSE +noCL | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.7 | 1.0 |
| EN–DE | 55.6 | 57.3 | 56.3 | 54.1 | 50.9 | 47.8 | 40.6 | 25.4 |
| EN–TR | 40.4 | 42.0 | 41.9 | 42.1 | 41.6 | 40.2 | 35.9 | 23.6 |
| EN–FI | 45.6 | 46.1 | 46.1 | 45.1 | 44.4 | 42.8 | 34.4 | 19.3 |
| EN–HR | 37.5 | 39.8 | 43.0 | 45.3 | 46.5 | 46.5 | 41.3 | 24.7 |
| EN–RU | 45.6 | 46.5 | 48.0 | 49.9 | 50.5 | 50.2 | 46.1 | 24.7 |
| EN–IT | 60.2 | 61.6 | 63.0 | 62.3 | 61.3 | 60.1 | 49.9 | 26.4 |
| EN–FR | 64.1 | 66.0 | 67.1 | 67.5 | 65.9 | 64.8 | 55.3 | 34.2 |
| DE–TR | 32.5 | 33.6 | 33.4 | 33.0 | 32.3 | 30.2 | 24.4 | 17.7 |
| DE–FI | 39.7 | 39.7 | 39.3 | 37.4 | 35.9 | 34.3 | 25.8 | 16.0 |
| DE–HR | 33.3 | 35.8 | 37.0 | 38.5 | 38.8 | 36.1 | 29.2 | 19.2 |
| DE–RU | 40.0 | 40.9 | 41.2 | 41.5 | 41.1 | 38.4 | 29.6 | 14.3 |
| DE–IT | 49.5 | 51.3 | 51.3 | 51.4 | 49.4 | 46.3 | 36.4 | 19.4 |
| DE–FR | 50.0 | 52.2 | 53.2 | 53.2 | 51.8 | 49.1 | 37.8 | 22.5 |
| TR–FI | 31.3 | 32.3 | 31.6 | 31.7 | 31.2 | 29.8 | 24.8 | 15.5 |
| TR–HR | 25.4 | 28.3 | 31.2 | 33.1 | 34.1 | 33.6 | 28.9 | 18.7 |
| TR–RU | 32.9 | 34.9 | 35.5 | 36.9 | 36.5 | 34.0 | 29.0 | 14.1 |
| TR–IT | 37.1 | 38.7 | 39.6 | 38.9 | 38.8 | 37.7 | 31.5 | 17.0 |
| TR–FR | 39.4 | 41.4 | 42.1 | 43.1 | 43.4 | 41.9 | 34.6 | 20.9 |
| FI–HR | 30.4 | 32.3 | 34.7 | 36.4 | 36.9 | 36.3 | 28.6 | 17.4 |
| FI–RU | 38.2 | 39.5 | 40.0 | 41.0 | 40.5 | 37.6 | 28.9 | 15.1 |
| FI–IT | 39.9 | 42.9 | 42.9 | 43.4 | 44.2 | 41.5 | 34.0 | 18.1 |
| FI–FR | 42.8 | 44.7 | 45.9 | 45.9 | 46.1 | 43.6 | 34.6 | 18.6 |
| HR–RU | 40.6 | 41.8 | 43.9 | 45.7 | 45.8 | 45.0 | 39.0 | 24.5 |
| HR–IT | 40.4 | 43.5 | 46.0 | 47.9 | 48.6 | 47.6 | 41.8 | 24.5 |
| HR–FR | 43.6 | 46.8 | 48.6 | 50.0 | 50.1 | 47.9 | 42.0 | 25.7 |
| RU–IT | 46.6 | 48.1 | 48.1 | 47.6 | 46.8 | 44.5 | 38.7 | 21.8 |
| RU–FR | 48.7 | 50.2 | 51.0 | 50.9 | 50.0 | 49.0 | 43.6 | 26.6 |
| IT–FR | 64.1 | 64.9 | 65.8 | 65.6 | 64.9 | 63.0 | 54.0 | 33.9 |
| **Average** | 42.7 | 44.4 | 45.3 | 45.7 | 45.3 | 43.6 | 36.5 | 21.4 |

Table 8: Individual P@1 scores (×100%) for all 28 language pairs in the GT-BLI dataset of Glavaš et al. (2019), across different values for λ. The model variant is LABSE +noCL (see §3); similar patterns are observed with another multilingual SE in our evaluation (XMPNET). $|\mathcal{D}| = 5k$.

| Pair ↓ / λ = → | VECMAP(0.0) | LABSE +CL | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.7 | 1.0 |
| EN–DE | 55.6 | 59.4 | 61.2 | 61.1 | 60.1 | 56.8 | 46.5 | 31.6 |
| EN–TR | 40.4 | 43.9 | 46.8 | 50.1 | 51.0 | 49.9 | 45.4 | 33.1 |
| EN–FI | 45.6 | 48.6 | 51.0 | 53.3 | 54.1 | 53.9 | 46.1 | 30.8 |
| EN–HR | 37.5 | 44.1 | 49.7 | 53.9 | 56.7 | 57.1 | 49.7 | 36.9 |
| EN–RU | 45.6 | 50.1 | 52.5 | 55.8 | 58.5 | 59.0 | 56.2 | 46.4 |
| EN–IT | 60.2 | 62.5 | 65.3 | 66.9 | 67.2 | 66.2 | 55.7 | 33.3 |
| EN–FR | 64.1 | 66.3 | 69.5 | 71.2 | 71.1 | 69.5 | 59.5 | 42.1 |
| DE–TR | 32.5 | 35.4 | 36.5 | 37.6 | 36.7 | 35.6 | 31.8 | 24.4 |
| DE–FI | 39.7 | 41.6 | 42.5 | 43.3 | 42.2 | 39.5 | 32.7 | 25.0 |
| DE–HR | 33.3 | 37.9 | 40.8 | 42.9 | 43.1 | 41.6 | 35.9 | 27.9 |
| DE–RU | 40.0 | 43.1 | 44.0 | 45.6 | 45.9 | 44.7 | 37.3 | 27.5 |
| DE–IT | 49.5 | 51.3 | 52.0 | 52.5 | 50.9 | 47.7 | 38.9 | 24.6 |
| DE–FR | 50.0 | 52.4 | 53.7 | 54.5 | 53.6 | 50.3 | 42.7 | 31.3 |
| TR–FI | 31.3 | 33.3 | 34.3 | 34.9 | 35.0 | 33.7 | 30.2 | 22.2 |
| TR–HR | 25.4 | 30.4 | 34.2 | 36.8 | 38.5 | 38.2 | 35.8 | 27.4 |
| TR–RU | 32.9 | 35.2 | 36.6 | 37.4 | 37.5 | 36.1 | 32.5 | 24.7 |
| TR–IT | 37.1 | 39.0 | 41.2 | 41.6 | 41.5 | 39.9 | 34.4 | 22.8 |
| TR–FR | 39.4 | 41.8 | 43.0 | 43.9 | 43.8 | 43.3 | 38.6 | 29.4 |
| FI–HR | 30.4 | 34.0 | 37.5 | 39.7 | 41.2 | 40.2 | 36.4 | 27.4 |
| FI–RU | 38.2 | 40.2 | 41.1 | 42.4 | 42.6 | 40.2 | 36.2 | 27.6 |
| FI–IT | 39.9 | 43.3 | 44.3 | 45.3 | 45.9 | 44.7 | 38.2 | 25.2 |
| FI–FR | 42.8 | 44.5 | 46.3 | 47.6 | 47.8 | 46.0 | 40.0 | 26.4 |
| HR–RU | 40.6 | 42.0 | 44.9 | 46.4 | 47.5 | 48.3 | 44.7 | 36.1 |
| HR–IT | 40.4 | 43.3 | 47.7 | 49.7 | 50.6 | 50.4 | 46.1 | 32.3 |
| HR–FR | 43.6 | 47.2 | 49.0 | 51.8 | 52.0 | 51.0 | 46.4 | 35.0 |
| RU–IT | 46.6 | 48.8 | 49.6 | 50.2 | 50.4 | 48.7 | 44.9 | 32.0 |
| RU–FR | 48.7 | 51.0 | 51.9 | 53.5 | 53.1 | 52.1 | 47.1 | 37.0 |
| IT–FR | 64.1 | 64.9 | 65.9 | 66.0 | 66.0 | 64.6 | 56.9 | 42.2 |
| **Average** | 42.7 | 45.6 | 47.6 | 49.1 | 49.4 | 48.2 | 42.4 | 30.8 |

Table 9: Individual P@1 scores (×100%) for all 28 language pairs in the GT-BLI dataset of Glavaš et al. (2019), across different values for λ. The model variant is LABSE +CL (see §3); similar patterns are observed with another multilingual SE in our evaluation (XMPNET). $|\mathcal{D}| = 5k$.

| Pair ↓ / Config → | VECMAP | MBERT | | XLM-R | | LABSE | | XMPNET | |
|---|---|---|---|---|---|---|---|---|---|
| | | +noCL (1.0) | +noCL (0.5) | +noCL (1.0) | +noCL (0.5) | +noCL (1.0) | +noCL (0.5) | +noCL (1.0) | +noCL (0.5) |
| EN–FI | 42.2 | 1.1 | 30.9 | 1.4 | 21.4 | 46.8 | 51.6 | 47.7 | 53.4 |
| EN–RU | 39.7 | 5.2 | 30.1 | 2.5 | 21.3 | 53.5 | 52.3 | 57.5 | 55.2 |
| EN–FR | 64.8 | 11.3 | 52.8 | 2.2 | 29.6 | 68.5 | 74.2 | 64.7 | 73.9 |
| FI–RU | 33.0 | 4.3 | 25.3 | 3.2 | 20.7 | 38.4 | 41.6 | 42.7 | 45.6 |
| FI–FR | 46.7 | 3.6 | 35.3 | 0.3 | 22.8 | 43.1 | 52.6 | 42.9 | 53.2 |
| RU–FR | 48.2 | 8.6 | 37.9 | 0.8 | 25.3 | 52.2 | 57 | 52.3 | 58.4 |
| Average | 45.8 | 5.7 | 35.4 | 1.7 | 23.5 | 50.4 | 54.9 | 51.3 | 56.6 |

Table 10: Individual Spearman's $\rho$ correlation scores ($\times 100$) on the XLSIM task (Multi-SimLex) for a subset of language pairs in our evaluation, with multilingual LMs and SEs used 'off-the-shelf' *without contrastive fine-tuning* (§2). See §3 for the description of different model configurations/variants. $|\mathcal{D}| = 5k$, with XLSIM test pairs removed from the dictionary. The number in the parentheses denotes the value for $\lambda$ (see §3).

| Pair ↓ / Config → | VECMAP | MBERT | | XLM-R | | LABSE | | XMPNET | |
|---|---|---|---|---|---|---|---|---|---|
| | | +noCL (1.0) | +noCL (0.5) | +noCL (1.0) | +noCL (0.5) | +noCL (1.0) | +noCL (0.5) | +noCL (1.0) | +noCL (0.5) |
| EN–FI | 42.2 | 32.4 | 43.5 | 42.3 | 48 | 45.6 | 50.6 | 48.1 | 51.5 |
| EN–RU | 39.7 | 34.8 | 42.1 | 45.8 | 47.3 | 47.1 | 49.2 | 50.5 | 50.4 |
| EN–FR | 64.8 | 56.5 | 67.4 | 57.9 | 69.2 | 64 | 72 | 64.1 | 71.9 |
| FI–RU | 33.0 | 28.4 | 35.9 | 38.3 | 40.8 | 38.3 | 42.3 | 40.1 | 43.1 |
| FI–FR | 46.7 | 34.7 | 47.3 | 41.7 | 50.6 | 45.9 | 53.5 | 46.6 | 54.2 |
| RU–FR | 48.2 | 43.6 | 52.1 | 50.4 | 55.1 | 51.8 | 56.8 | 48.5 | 55.6 |
| Average | 45.8 | 38.4 | 48.1 | 46.1 | 51.8 | 48.8 | 54.1 | 49.6 | 54.5 |

Table 11: Individual Spearman's $\rho$ correlation scores ($\times 100$) on the XLSIM task (Multi-SimLex) for a subset of language pairs in our evaluation, with model variants *with contrastive fine-tuning* (§2). $|\mathcal{D}| = 5k$.

| Pair ↓ / $\lambda$ = → | xMPNET +noCL | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 1.0 |
| EN–FI | 42.2 | 45.2 | 48.2 | 50.7 | 52.6 | 53.4 | 53 | 52.2 | 47.7 |
| EN–RU | 39.7 | 43.1 | 46.6 | 49.9 | 52.9 | 55.2 | 56.7 | 57.5 | 57.5 |
| EN–FR | 64.8 | 67.6 | 70.1 | 72.2 | 73.5 | 73.9 | 73.2 | 71.7 | 64.7 |
| FI–RU | 33 | 36.1 | 39.3 | 42.5 | 44.7 | 45.6 | 45.5 | 45.4 | 42.7 |
| FI–FR | 46.7 | 49.2 | 51.5 | 53.1 | 53.6 | 53.2 | 51.5 | 49.6 | 42.9 |
| RU–FR | 48.2 | 51.1 | 53.8 | 56.2 | 57.7 | 58.4 | 58.1 | 57.2 | 52.3 |
| Average | 45.8 | 48.7 | 51.6 | 54.1 | 55.8 | 56.6 | 56.3 | 55.6 | 51.3 |

Table 12: Individual Spearman's $\rho$ correlation scores ($\times 100$) on the XLSIM task (Multi-SimLex) for a subset of language pairs in our evaluation, across different values for $\lambda$. The model variant is XMPNET +noCL (see §3); similar patterns are observed with another multilingual SE in our evaluation (LABSE). $|\mathcal{D}| = 5k$.

| Pair ↓ / $\lambda$ = → | xMPNET +CL | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 1.0 |
| EN–FI | 42.2 | 44.2 | 46.4 | 48.6 | 50.4 | 51.5 | 51.8 | 51.3 | 48.1 |
| EN–RU | 39.7 | 41.6 | 43.9 | 46.3 | 48.6 | 50.4 | 51.6 | 52 | 50.5 |
| EN–FR | 64.8 | 66.9 | 69 | 70.7 | 71.8 | 71.9 | 71 | 69.5 | 64.1 |
| FI–RU | 33.0 | 35.1 | 37.5 | 39.9 | 41.9 | 43.1 | 43.5 | 43 | 40.1 |
| FI–FR | 46.7 | 48.9 | 51.2 | 53.1 | 54.2 | 54.2 | 53.3 | 51.7 | 46.6 |
| RU–FR | 48.2 | 50.1 | 52.1 | 53.8 | 55.2 | 55.6 | 55.1 | 53.8 | 48.5 |
| Average | 45.8 | 47.8 | 50 | 52.1 | 53.7 | 54.5 | 54.4 | 53.6 | 49.6 |

Table 13: Individual Spearman's $\rho$ correlation scores ($\times 100$) on the XLSIM task (Multi-SimLex) for a subset of language pairs in our evaluation, across different values for $\lambda$. The model variant is XMPNET +CL (see §3); similar patterns are observed with another multilingual SE in our evaluation (LABSE). $|\mathcal{D}| = 5k$.

| Name | URL |
|---|---|
| MBERT | huggingface.co/bert-base-multilingual-uncased |
| XLM-R | huggingface.co/xlm-roberta-base |
| LABSE | huggingface.co/sentence-transformers/LaBSE |
| XMPNET | huggingface.co/sentence-transformers/paraphrase-multilingual-mpnet-base-v2 |

Table 14: URLs of the multilingual Transformer models used in this work.