

Don't Forget Cheap Training Signals Before Building Unsupervised Bilingual Word Embeddings

Silvia Severini, Viktor Hangya, Masoud Jalili Sabet, Alexander Fraser, Hinrich Schütze

Center for Information and Language Processing

LMU Munich, Germany

{silvia, hangyav, masoud, fraser}@cis.uni-muenchen.de

Abstract

Bilingual Word Embeddings (BWEs) are one of the cornerstones of cross-lingual transfer of NLP models. They can be built using only monolingual corpora without supervision leading to numerous works focusing on unsupervised BWEs. However, most of the current approaches to build unsupervised BWEs do not compare their results with methods based on easy-to-access cross-lingual signals. In this paper, we argue that such signals should always be considered when developing unsupervised BWE methods. The two approaches we find most effective are: 1) using identical words as seed lexicons (which unsupervised approaches incorrectly assume are not available for orthographically distinct language pairs) and 2) combining such lexicons with pairs extracted by matching romanized versions of words with an edit distance threshold. We experiment on thirteen non-Latin languages (and English) and show that such cheap signals work well and that they outperform using more complex unsupervised methods on distant language pairs such as Chinese, Japanese, Kannada, Tamil, and Thai. In addition, they are even competitive with the use of high-quality lexicons in supervised approaches. Our results show that these training signals should not be neglected when building BWEs, even for distant languages.

Keywords: Bilingual Word Embeddings, Bilingual Dictionary Induction, Romanization

1. Introduction

Bilingual Word Embeddings (BWEs) are useful for many cross-lingual tasks. They can be built effectively even when only a small seed lexicon is available by mapping monolingual embeddings into a shared space. This makes them particularly valuable for low-resource settings (Mikolov et al., 2013). In addition, unsupervised mapping approaches can build BWEs for some languages when no seed lexicon is available. Various unsupervised methods have been proposed relying on the assumption that embedding spaces are isomorphic (Zhang et al., 2017; Lample et al., 2018; Artetxe et al., 2018; Alvarez-Melis and Jaakkola, 2018; Chen and Cardie, 2018; Hoshen and Wolf, 2018; Mohiuddin and Joty, 2019; Alaux et al., 2019; Dou et al., 2020; Grave et al., 2019; Li et al., 2020). However, with one exception, none of them compare their results with the widely available baseline of using identical words as seed lexicons.

It has been shown that identical word pairs of two languages can be used to build high quality BWEs (Smith et al., 2017; Artetxe et al., 2017). However, they were only tested on language pairs with similar scripts. The only exception is the work of Søgaard et al. (2018), who tested identical word pairs on English and Greek which use different alphabetical characters but the same numerals. Regardless of these experiments, recent works still propose novel unsupervised approaches without considering such cheap training signals, at least as baseline systems (Mohiuddin and Joty, 2019; Alaux et al., 2019; Dou et al., 2020; Grave et al., 2019; Li et al., 2020).

In this paper however, we argue that such signals should be used as a cheap and effective baseline in the devel-

opment of future unsupervised methods. We define them cheap as they require widely available monolingual corpora only, e.g., Wikipedia dumps, but no parallel data. We study two approaches for extracting the initial seed lexicons to build BWEs without relying on expensive dictionaries. (1) First, we leverage identical pairs as proposed by Smith et al. (2017; Artetxe et al. (2017)). Previous work assumed such pairs not to be available for language pairs with distinct scripts, hence the development of various unsupervised mapping approaches. We show that, surprisingly, they do appear in large quantities in the monolingual corpora that we use, even for distinct-script pairs. In contrast to Søgaard et al. (2018), we test identical word pairs on multiple language pairs with distinct scripts, including pairs using distinct numerals. In addition, we propose to (2) strengthen identical pairs by extending them with further easily accessible pairs based on romanization and edit distance, which exploits implicit links between languages in the form of approximate word transliteration pairs.

We focus on distant language pairs having distinct scripts for many of which unsupervised approaches have failed or had very poor performance so far. For instance, English to Chinese, Japanese, Kannada, Tamil, and Thai, which all obtain a score close to 0 on the Bilingual Dictionary Induction (BDI) task (Vulić et al., 2019). We evaluate the two approaches on thirteen different non-Latin¹ languages paired with English on BDI. We compare our lexicons' performance with unsupervised mapping and the frequently used MUSE training lexi-

¹We use (*non-*)Latin language here as a short form for language standardly written in a (*non-*)Latin script.

cons (Lample et al., 2018) and show that our noisy word pairs make it possible to build BWEs for language pairs where unsupervised approaches failed before and give accuracy scores similar to high quality lexicons.

Our work calls into question – at least for BDI – the strong trend toward unsupervised approaches in recent literature, similarly to Vulić et al. (2019), given that cheap signals are (i) available and easy to exploit, (ii) sufficient to obtain performance similar to dictionaries based on parallel resources like MUSE and (iii) able to make up for the failure of unsupervised methods. Finally, we analyze which lexicon properties impact performance and show that our lexicon outperform unsupervised methods also for non-English language pairs. Our paper calls for the need to use easily accessible bilingual signals, such as identical and/or transliteration word pairs, as baselines when developing unsupervised BWE approaches.

2. Unsupervised pair extraction

We show that we can extract the seed lexicon needed for mapping systems without the need for labeled data, making up for the failure of unsupervised methods. First, we show that identical pairs do appear in corpora of distant languages and can be exploited. Secondly, we propose a novel method to boost the identical pairs sets by extracting the initial seed lexicon without the need for any bilingual knowledge, starting from monolingual corpora, and using romanization and edit distance.

2.1. Identical pair approach

When dealing with languages with different scripts, identical pairs would seem to be unlikely to occur, which is assumed by unsupervised mapping methods. Smith et al. (2017; Artetxe et al. (2017) form dictionaries from identical strings which appear in both languages but limit their approach to similar languages sharing a common alphabet, such as European ones. Similarly, (Lample et al., 2018) refrain from using such identical word pairs, assuming they are not available for distant languages. An exception is the work of Sjøgaard et al. (2018) which shows the presence of identical pairs between English and Greek, which share numerals only but not alphabetical characters.

However, we show that there are domains where these pairs are actually available in large quantity even for pairs with different scripts, including the use of different numerals; an example is Wikipedia: see the statistics of fastText Wikipedia embeddings (Bojanowski et al., 2017) in Table 1. Most of these identical pairs are punctuation marks and digits, non-transliterated named entities written in the Latin script, or English words (assumably words of a title) which were not translated in the non-English languages. This is also true for language pairs not including English. In this paper, we build BWEs based on these pairs and show that they are sufficient for good BDI results on distant language pairs with distinct scripts.

| Lang | ID | Lang | ID | Lang | ID |
|--------|-----|--------|-----|--------|-----|
| ko-th* | 17K | ko-he* | 11K | he-th* | 15K |
| en-zh* | 62K | en-bn* | 31K | en-ar* | 19K |
| en-th | 46K | en-hi* | 30K | en-ru | 18K |
| en-ja | 43K | en-ta* | 23K | en-he* | 17K |
| en-el | 35K | en-kn* | 21K | en-ko* | 15K |
| en-fa* | 32K | | | | |

Table 1: Number of identical pairs per language pair. Language pairs using different digits as their official numerals, on top of different alphabetical characters, are indicated with *.

2.2. Romanization based augmentation (ID++)

Identical pairs are noisy and may appear in smaller quantities for certain corpora and language pairs (e.g., he-ko). We propose our romanization approach that builds the seed lexicon completely automatically and can augment the identical pairs set. We exploit the concept of transliteration and orthographic similarity to find a cheap signal between languages (cf. (Riley and Gildea, 2018; Severini et al., 2020a; Severini et al., 2020b; Severini et al., 2022)) and to take advantage of cognates (Chakravarthi et al., 2019; Laville et al., 2020). It consists of 3 steps at the end of which we add the identical pairs and run VecMap in a semi-supervised setting.

1. Source candidates First, we generate a list of source language words, which are the candidates to be matched with a word on the target side. We use the English Wikipedia dumps² as our monolingual corpus and apply Flair (Akbi et al., 2018) to extract Universal Part-of-Speech (UPOS) tags. We collect all English proper nouns (PROPN), since names are often transliterated between languages. The resulting English proper noun set consists of $\approx 800K$ words.

2. Target candidates The language-specific target data is extracted from the vocabulary of the pre-trained Wikipedia fastText embeddings (Bojanowski et al., 2017). The sets are not pre-processed with a POS tagger assuming that such a tool is missing or perform poorly for low-resource languages. Compared to the English proper noun set, the vocabularies are smaller: between 40K and 500K. Then, we romanize the corpora to obtain equivalent words but with only Latin characters – this supports the distance-based metrics in step (3). We use Uroman (Hermjakob et al., 2018) for romanization. Examples of romanization are $\kappa\alpha\rho\lambda$ (Russian) \rightarrow carl and $\beta\alpha\beta\upsilon\lambda\acute{\omega}\nu$ (Greek) \rightarrow babylon. Uroman mainly covers 1-1 character correspondences and does not vocalize words for Arabic and Hebrew. In general, its romanization is not as accurate as the transliteration of a neural model. However, neural models need a training corpus of labeled pairs to work well, while Uroman only

²<https://dumps.wikimedia.org/> (01.04.2020)

| | en-th | en-ja | en-kn | en-ta | en-zh |
|---|-------|-------|-------|-------------------|-------|
| Unsupervised | | | | | |
| 1. | 0.00 | 0.96 | 0.00 | 0.07 | 0.07 |
| 2. | 0.00 | 0.48 | 0.00 | 0.07 | 0.00 |
| 3. | 0.00 | 0.00 | 0.00 | 0.00 [◊] | 0.00 |
| Semi-supervised (Artetxe et al., 2018) | | | | | |
| ID | 24.40 | 48.87 | 22.03 | 17.93 | 37.00 |
| Rom. | 23.33 | 48.46 | 22.90 | 18.00 | 0.27 |
| ID++ | 23.47 | 49.14 | 24.23 | 18.20 | 35.00 |
| MUSE | 24.33 | 48.73 | 23.78 | 18.80 | 36.53 |

Table 2: acc@1 on BDI for unsupervised (1: Artetxe et al. (2018), 2: Grave et al. (2019), 3: Mohiuddin and Joty (2019)) and semi-supervised approaches for 5 languages for which unsupervised methods fail. The semi-supervised results are obtained using VecMap with three different initial lexicons: the identical pair set (ID), ID extended with romanization based pairs (ID++) and the MUSE dictionary. We show an ablation study as well, i.e., the romanized pairs only (Rom.). Scores from Mohiuddin et al. (2020) are marked with [◊].

uses the character descriptions from the Unicode table,³ manually created tables and some heuristics, supporting a large number of languages.

3. Candidate matching To find the corresponding target word for an English noun, the noun is compared with each (romanized) target word based on their orthography. The similarity of two words w_1 and w_2 is defined as $1 - \text{NL}(w_1, w_2)$, where NL is the Levenshtein distance (Levenshtein, 1966) divided by the length of the longer string. We select a pair of words if the similarity is ≥ 0.8 ; this ensures a trade off between number of pairs and quality, based on manual investigation. We use the Symmetric Delete algorithm to speed up computation, similarly to (Riley and Gildea, 2018). It takes the lists of source and target words, and a constant k and identifies all the source-target pairs that are identical after k insertion or deletions.⁴ The final step is to look up, for each romanized target word, its original non-romanized form.

3. Evaluation

We evaluate our seed lexicons on BDI to show the quality of the BWEs obtained with them. Recent papers (Marchisio et al., 2020) show that there is a direct relationship between BDI accuracy and downstream BLEU for machine translation. Moreover, Sabet et al. (2020) show that good-quality word embeddings directly reflect the performance also for extrinsic tasks like word alignment. We use the VecMap tool to build BWEs since it supports both unsupervised, semi-supervised and supervised techniques (Artetxe et al., 2018). The

semi-supervised approach is of particular interest to us since it performs well with small and noisy seed lexicons by iteratively refining them. VecMap iterates over two steps: embedding mapping and dictionary induction. The process starts from an initial dictionary that is iteratively augmented and refined by extracting probable word pairs from the BWEs built in the current iteration with BDI. The method is repeated until the improvement on the average dot product for the induced dictionary stays above a given threshold. We use pre-trained Wikipedia fastText embeddings (Bojanowski et al., 2017) as the input monolingual vectors, taking only the 200K most frequent words and using default parameters otherwise. We compare the performance of VecMap using our lexicons with MUSE. MUSE contains dictionaries for many languages and it was created using a Facebook internal translation tool (Lample et al., 2018), thus it can be considered as a higher quality cross-lingual resource based on parallel data. Since Kannada is not supported by MUSE, we use the dictionary provided by Anzer et al. (2020). We show *acc@1* scores based on CSLS vector similarity calculated by the MUSE evaluation tool (Lample et al., 2018).⁵

Tables 2 and 3 show accuracy for all language pairs considering English as the source; see Table 7 in Appendix B for the full table containing results in both directions. Table 2 gives scores for language pairs for which unsupervised methods completely diverge ($\text{acc@1} < 1$). We report results for three unsupervised methods (Artetxe et al., 2018; Mohiuddin and Joty, 2019; Grave et al., 2019). In contrast, using identical word pairs as lexicon (ID) or its extension with the romanization based pairs (ID++) with VecMap leads to successful BWEs without any parallel data or manually created lexicons. In addition, scores are even comparable to high-quality dictionaries like MUSE. Looking at results for all language pairs in Table 2 and 3, our sets always obtain results comparable to MUSE (baseline dictionaries), with improvements for Arabic, Chinese, Russian and Greek. In the unsupervised cases (Table 2), both ID and ID++ pair sets lead to an accuracy improvement of at least 17 points. ID++ outperform ID for three of the five low-resource pairs and five out of eight high-resource pairs proving that the romanized pairs can indeed strengthen the identical pairs sets. These results show that good quality BWEs can be built by relying on implicit cross-lingual signals without expensive supervision or fragile unsupervised approaches.

MUSE test w/o proper nouns The work of Kementchedjhieva et al. (2019) highlights that MUSE test sets contain a high number of proper nouns for German, Danish, Bulgarian, Arabic and Hindi. Since our romanization augmentation is based on such names, we evaluate their performance on the subsets of MUSE test

³<http://unicode.org/Public/UNIDATA/UnicodeData.txt>

⁴We used minimum frequency and minimum length equal to 1, k equals to 2.

⁵We follow Artetxe et al. (2018) work for comparison reasons and did not remove identical pairs from the test sets. However, overlaps between train romanized lexicons and test lexicons correspond to less than 1%.

| | Unsup. | ID | Rom. | ID++ | MUSE |
|-------|--------|-------|-------|-------|-------|
| en-ar | 36.30 | 40.27 | 39.33 | 40.20 | 39.87 |
| en-hi | 40.20 | 40.47 | 39.60 | 40.20 | 40.33 |
| en-ru | 44.80 | 49.13 | 48.87 | 49.53 | 48.80 |
| en-el | 47.90 | 47.87 | 48.00 | 48.27 | 48.00 |
| en-fa | 36.70 | 37.67 | 36.80 | 37.67 | 38.00 |
| en-he | 44.60 | 44.47 | 44.53 | 44.67 | 45.00 |
| en-bn | 18.20 | 19.87 | 19.80 | 20.13 | 21.60 |
| en-ko | 19.80 | 27.92 | 28.40 | 28.81 | 28.94 |

Table 3: acc@1 on BDI for (best) unsupervised method and semi-supervised VecMap with different initial lexicons. (full table in Appendix B, Table 7).

sets that don’t contain proper nouns. We remove proper nouns using the list of names obtained in Section 2.2 and evaluate the performance of all the approaches presented above. The new sets contains 10% less pairs on average. Results are shown in Table 8, Appendix C. The performance is similar to the one obtained on the original test sets, proving that our dictionaries and methods are not biased towards aligning word embeddings of proper nouns.

Non-English centric evaluation We analyze the performance of ID and ID++ for language pairs that do not include English. We use the test dictionaries from Vulić et al. (2019) that are derived from PanLex (Baldwin et al., 2010; Kamholz et al., 2014) by automatically translating each source language word into the target languages. We run VecMap for all combinations of Korean, Hebrew, and Thai. Romanized train lexicons are extracted by combining the languages through English (e.g., th-ko is obtained using en-th and en-ko), i.e., words are paired if their English translation is the same. Table 4 shows results. When Thai is involved, the unsupervised method fails as for English-Thai. Both ID and ID++ always outperform the respective unsupervised scores, and perform similar to higher-quality dictionaries. Additionally, ID++ outperforms ID in 3 out of 6 cases. These results demonstrate further the simplicity and high quality of our methods.

Romanized-only We analyze the performance of romanized pair lexicons on their own. Line Rom. in Table 2 and 3 shows that they obtain competitive results to the other two approaches, with improvements for Japanese, and perform similarly to MUSE dictionaries. The only failure is for Chinese (en-zh) – presumably because Chinese has a logographic script that does not represent phonemes directly, so romanization is less effective. These results show that the romanized pairs on their own also represent strong signals that shouldn’t be neglected. Moreover, they constitute a good alternative when identical pairs are not available in such quantities (e.g., corpora of religious domain, law field, or cultural-specific documents).

Impact of OOVs We analyze the pairs used for the various sets (Appendix A, Table 5). We define OOVs

| | Unsup. | ID | Rom. | ID++ | PanLex |
|-------|--------|--------------|--------------|--------------|--------|
| th-ko | 0.00 | 2.81 | <u>3.37</u> | 3.09 | 2.95 |
| th-he | 0.00 | <u>9.75</u> | 0.00 | 8.86 | 10.13 |
| ko-th | 0.00 | <u>15.90</u> | 14.23 | 15.26 | 14.36 |
| ko-he | 14.62 | 15.68 | <u>16.08</u> | 16.00 | 15.11 |
| he-th | 0.00 | 16.42 | 0.00 | <u>16.54</u> | 17.90 |
| he-ko | 14.30 | <u>15.39</u> | 15.15 | 15.09 | 16.06 |

Table 4: acc@1 on BDI for unsupervised and semi-supervised VecMap for all combinations of Korean, Hebrew, and Thai. PanLex are results obtained with training lexicons from Vulić et al. (2019) and semi-supervised VecMap.

as words for which there is no embedding available among the pre-trained Wikipedia fastText embeddings. Our romanized sets contain a substantial number of OOVs. (The identical pair sets do not contain OOVs because words are extracted from the top 200K most frequent.) The main reason for OOVs is that the selected English pair of a word is so rare that they do not have embeddings. On the other hand, the high number of OOVs (and resulting reduction of usable pairs) has only a limited negative impact on the performance.

Size of seed set and word frequency We analyze the impact of the size of the initial romanized seed set and of word frequency. Appendix A, Table 6, displays accuracy scores for MUSE and Romanized lexicons containing the $n \in \{25, 1000\}$ least and most frequent word pairs. Performance of VecMap applied to seed sets of size 25 is close to 0. The only exception is Russian, where the unsupervised approach already works well. Next, we investigate seed sets of size 1000 consisting of either the least frequent or the most frequent words. High-frequency seed sets give better results as expected. The effect is particularly strong for Tamil: the high-frequency set has performance close to the full set, whereas the low-frequency set is at ≤ 0.07 . The performance of MUSE seed sets of size 25 and romanized seed sets of size 1000 is similar, demonstrating the higher quality of MUSE. However, obtaining the romanized pairs is much cheaper.

4. Conclusion

We have analyzed two cheap resources for building BWEs which can alleviate the issues of unsupervised methods which fail on multiple language pairs. We focused on a wide range of non-Latin languages paired with English. (i) We exploited identical pairs that surprisingly appear in corpora of distinct scripts. We showed that they can be used even when numerals are distinct in contrast to previous work. (ii) We combined them with a simple method to extract the initial hypothesis set via romanization and edit distance. With both approaches, we obtained results that are competitive with high-quality dictionaries. Without using explicit cross-lingual signal, we outperformed previous unsupervised work for most languages and in particular for five

language pairs for which previous unsupervised work failed. Our results question the strong trend towards unsupervised mapping approaches, and show that cheap cross-lingual signals should always be considered for building BWEs, even for distant languages.

Acknowledgments

This work was funded by the European Research Council (grant #740516, #640550), the German Federal Ministry of Education and Research (BMBF, grant #01IS18036A), and the German Research Foundation (DFG; grant FR 2829/4-1).

5. Bibliographical References

- Akbik, A., Blythe, D., and Vollgraf, R. (2018). Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.
- Alaux, J., Grave, E., Cuturi, M., and Joulin, A. (2019). Unsupervised hyperalignment for multilingual word embeddings. In *Proceedings of the 7th International Conference on Learning Representations*.
- Alvarez-Melis, D. and Jaakkola, T. S. (2018). Gromov-wasserstein alignment of word embedding spaces. In *EMNLP*.
- Anzer, M., Chronopoulou, A., and Fraser, A. (2020). Comparing unsupervised and supervised approaches for kannada/english bilingual word embeddings. *Bachelor thesis at Ludwig Maximilians Universität München*.
- Artetxe, M., Labaka, G., and Agirre, E. (2017). Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 451–462.
- Artetxe, M., Labaka, G., and Agirre, E. (2018). A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 789–798.
- Baldwin, T., Pool, J., and Colowick, S. (2010). Panlex and lextract: Translating all words of all languages of the world. In *Coling 2010: Demonstrations*, pages 37–40.
- Chakravarthi, B. R., Arcan, M., and McCrae, J. P. (2019). Comparison of different orthographies for machine translation of under-resourced dravidian languages. In *2nd Conference on Language, Data and Knowledge (LDK 2019)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- Chen, X. and Cardie, C. (2018). Unsupervised multilingual word embeddings. *arXiv preprint arXiv:1808.08933*.
- Dou, Z. Y., Zhou, Z. H., and Huang, S. (2020). Unsupervised bilingual lexicon induction via latent variable models. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 621–626.
- Grave, E., Joulin, A., and Berthet, Q. (2019). Unsupervised alignment of embeddings with wasserstein procrustes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1880–1890. PMLR.
- Hoshen, Y. and Wolf, L. (2018). Non-adversarial unsupervised word translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 469–478.
- Kamholz, D., Pool, J., and Colowick, S. (2014). Panlex: Building a resource for panlingual lexical translation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 3145–3150.
- Kementchedjhieva, Y., Hartmann, M., and Søgaard, A. (2019). Lost in evaluation: Misleading benchmarks for bilingual dictionary induction. *arXiv preprint arXiv:1909.05708*.
- Laville, M., Hazem, A., and Morin, E. (2020). Taln/ls2n participation at the bucc shared task: bilingual dictionary induction from comparable corpora. In *Proceedings of the 13th Workshop on Building and Using Comparable Corpora*, pages 56–60.
- Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and reversals. *Soviet physics doklady*, 10(8):707–710.
- Li, Y., Luo, Y., Lin, Y., Du, Q., Wang, H., Huang, S., Xiao, T., and Zhu, J. (2020). A simple and effective approach to robust unsupervised bilingual dictionary induction. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5990–6001.
- Marchisio, K., Duh, K., and Koehn, P. (2020). When does unsupervised machine translation work? *arXiv preprint arXiv:2004.05516*.
- Mikolov, T., Le, Q. V., and Sutskever, I. (2013). Exploiting similarities among languages for machine translation. *arXiv preprint arXiv:1309.4168*.
- Mohiuddin, T. and Joty, S. (2019). Revisiting adversarial autoencoder for unsupervised word translation with cycle consistency and improved training. In *Proceedings of NAACL-HLT*, pages 3857–3867.
- Mohiuddin, M. T., Bari, M. S., and Joty, S. (2020). Lnmap: Departures from isomorphic assumption in bilingual lexicon induction through non-linear mapping in latent space. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2712–2723.
- Riley, P. and Gildea, D. (2018). Orthographic features for bilingual lexicon induction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 390–394.
- Sabet, M. J., Dufter, P., Yvon, F., and Schütze, H. (2020). Simalign: High quality word alignments without parallel training data using static and context-

- tualized embeddings. In *Findings of the Association for Computational Linguistics: EMNLP 2020*.
- Severini, S., Hangya, V., Fraser, A., and Schütze, H. (2020a). Combining word embeddings with bilingual orthography embeddings for bilingual dictionary induction. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6044–6055.
- Severini, S., Hangya, V., Fraser, A., and Schütze, H. (2020b). Lmu bilingual dictionary induction system with word surface similarity scores for bucc 2020. In *Proceedings of the 13th Workshop on Building and Using Comparable Corpora*, pages 49–55.
- Severini, S., Imani, A., Duffer, P., and Schütze, H. (2022). Towards a broad coverage named entity resource: A data-efficient approach for many diverse languages. *arXiv preprint arXiv:2201.12219*.
- Smith, S. L., Turban, D. H., Hamblin, S., and Hammerla, N. Y. (2017). Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *arXiv preprint arXiv:1702.03859*.
- Søgaard, A., Ruder, S., and Vulić, I. (2018). On the limitations of unsupervised bilingual dictionary induction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 778–788.
- Vulić, I., Glavaš, G., Reichart, R., and Korhonen, A. (2019). Do we really need fully unsupervised cross-lingual embeddings? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4398–4409.
- Zhang, M., Liu, Y., Luan, H., and Sun, M. (2017). Adversarial training for unsupervised bilingual lexicon induction. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1959–1970.

6. Language Resource References

- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Hermjakob, U., May, J., and Knight, K. (2018). Out-of-the-box universal romanization tool uroman. In *Proceedings of ACL 2018, System Demonstrations*, pages 13–18.
- Lample, G., Conneau, A., Ranzato, M., Denoyer, L., and Jégou, H. (2018). Word translation without parallel data. In *International Conference on Learning Representations*.

A. Statistics

In this section we show statistics on the language pairs analyzed and additional scores. Table 5 presents the number of pairs for each set that are not OOVs in the fastText wiki embeddings (Bojanowski et al., 2017).

| | MUSE | ID | Romanized | ID++ |
|-------|--------|--------|------------------|----------------|
| en-th | 6,799 | 46,653 | 10,721 / 53,804 | 58779 / 101066 |
| en-ja | 7,135 | 43,556 | 11,488 / 118,626 | 54970 / 161848 |
| en-kn | 1,552 | 21,090 | 12,888 / 59,207 | 33843 / 80032 |
| en-ta | 8,091 | 23,538 | 5,987 / 120,836 | 29472 / 143990 |
| en-zh | 8,728 | 62,289 | 6,360 / 41,829 | 68597 / 103971 |
| en-ar | 11,571 | 19,275 | 4,773 / 61,031 | 24019 / 80115 |
| en-hi | 8,704 | 30,502 | 16,180 / 73,553 | 46557 / 103791 |
| en-ru | 10,887 | 18,663 | 9,913 / 301,698 | 28520 / 319688 |
| en-el | 10,662 | 35,270 | 20,740 / 150,472 | 55841 / 185244 |
| en-fa | 8,869 | 32,866 | 10,226 / 85,210 | 43019 / 117817 |
| en-he | 9,634 | 17,012 | 4,005 / 40,258 | 20977 / 57059 |
| en-bn | 8,467 | 31,954 | 10,721 / 53,804 | 42573 / 85532 |
| en-ko | 7,999 | 15,518 | 9956 / 134156 | 25344 / 149031 |

Table 5: Number of pairs used that are not OOVs in the fastText wiki embeddings compared to the full size of the sets. For MUSE full and identical pairs sets there are no OOVs.

B. Main results

In Table 7 there are the accuracy scores based on CSLS vector similarity calculated by the MUSE evaluation tool (Lample et al., 2018). We show the scores for thirteen language pairs in both directions. The first five pairs are the ones for which unsupervised methods fail. We show both unsupervised and semi-supervised VecMap performance with baselines dictionaries and our three sets.

C. MUSE proper nouns removal

Table 8 shows results computed on the subsets of MUSE test sets that don’t contain proper nouns. We remove proper nouns using the list of names obtained in Section 2.2 The new sets contains 10% less pairs on average.

D. Reproducibility

We run our method on up to 48 cores of Intel(R) Xeon(R) CPU E7-8857 v2 with 1TB memory and a single GeForce GTX 1080 GPU with 8GB memory. The training of semi-supervised BWEs using VecMap took approximately 1 hour per language pair. For VecMap, as well as for all others methods we analyzed, we used the latest code available in their git repositories with default parameters. ID++ is implemented in Python.

| | | MUSE | | | | Rom. | | | |
|-------|---|-------|-------|-------|-------|-------|-------|-------|-------|
| | | 25L | 25H | 1000L | 1000H | 25L | 25H | 1000L | 1000H |
| en-ta | → | 14.73 | 16.27 | 17.33 | 17.40 | 0.00 | 0.00 | 0.07 | 17.80 |
| | ← | 16.48 | 18.35 | 22.44 | 23.44 | 0.00 | 0.00 | 0.00 | 21.57 |
| en-fa | → | 35.33 | 34.20 | 38.07 | 37.20 | 0.00 | 0.20 | 37.47 | 37.47 |
| | ← | 41.73 | 42.60 | 44.14 | 44.21 | 0.07 | 0.13 | 42.40 | 43.40 |
| en-zh | → | 39.00 | 39.40 | 38.20 | 37.67 | 0.00 | 0.00 | 0.07 | 0.40 |
| | ← | 32.93 | 34.47 | 34.33 | 34.40 | 0.00 | 0.00 | 0.07 | 0.60 |
| en-ru | → | 49.07 | 43.07 | 49.07 | 49.27 | 49.33 | 47.73 | 49.40 | 49.00 |
| | ← | 65.93 | 60.60 | 65.93 | 66.13 | 65.80 | 64.47 | 65.60 | 66.40 |

Table 6: acc@1 using 25 or 1000 pairs lower-frequency (L) and higher-frequency (H) sets for MUSE and our romanized only (Rom.) set.

| | | Baselines | | | Semi-sup. MUSE | Our Semi-sup. | | | |
|----|-------|--------------|--------------|-------|---------------------------|------------------|--------------|--------------|--------------|
| | | Unsupervised | | | | ID | Rom. | ID++ | |
| | | 1 | 2 | 3 | | | | | |
| 1 | en-th | → | 0.00 | 0.00 | 0.00 | 24.33 | 24.40 | 23.33 | 23.47 |
| | | ← | 0.00 | 0.00 | 0.00 | 19.04 | 19.92 | 17.96 | 19.85 |
| 2 | en-ja | → | 0.96 | 0.48 | 0.00 | 48.73 | 48.87 | 48.46 | 49.14 |
| | | ← | 0.96 | 0.00 | 0.00 | 32.87 | 33.22 | 34.80 | 33.43 |
| 3 | en-kn | → | 0.00 | 0.00 | 0.00 | 23.78* | 22.03 | 22.90 | 24.23 |
| | | ← | 0.00 | 0.00 | 0.00 | 41.25* | 43.04 | 42.50 | 41.79 |
| 4 | en-ta | → | 0.07 | 0.07 | 0.00 [◊] | 18.80 | 17.93 | 18.00 | 18.20 |
| | | ← | 0.07 | 0.00 | 0.00 [◊] | 24.38 | 24.78 | 23.51 | 24.78 |
| 5 | en-zh | → | 0.07 | 0.00 | 0.00 | 36.53 | 37.00 | 0.27 | 35.00 |
| | | ← | 0.00 | 0.00 | 0.00 | 32.80 | 34.33 | 0.07 | 32.67 |
| 6 | en-ar | → | 33.60 | 7.67 | 36.30 [◊] | 39.87 | 40.27 | 39.33 | 40.20 |
| | | ← | 47.72 | 12.92 | 52.60 [◊] | 54.48 | 54.42 | 54.42 | 54.62 |
| 7 | en-hi | → | 40.20 | 0.00 | 0.00 [◊] | 40.33 | 40.47 | 39.60 | 40.20 |
| | | ← | 50.57 | 0.07 | 0.00 [◊] | 50.50 | 49.77 | 49.90 | 50.10 |
| 8 | en-ru | → | 48.80 | 37.33 | 46.90 [◊] | 48.80 | 49.13 | 48.87 | 49.53 |
| | | ← | 66.13 | 52.73 | 64.70 [◊] | 65.67 | 66.13 | 65.73 | 66.07 |
| 9 | en-el | → | 47.67 | 34.67 | 47.90 [◊] | 48.00 | 47.87 | 48.00 | 48.27 |
| | | ← | 63.40 | 49.20 | 63.50 [◊] | 63.33 | 63.27 | 64.40 | 63.47 |
| 10 | en-fa | → | 33.27 | 0.53 | 36.70 [◊] | 38.00 | 37.67 | 36.80 | 37.67 |
| | | ← | 39.99 | 0.40 | 44.50 [◊] | 43.47 | 43.67 | 42.93 | 43.60 |
| 11 | en-he | → | 44.60 | 37.13 | 44.00 [◊] | 45.00 | 44.47 | 44.53 | 44.67 |
| | | ← | 57.88 | 50.01 | 57.10 [◊] | 57.94 | 58.14 | 57.81 | 57.94 |
| 12 | en-bn | → | 18.20 | 0.00 | 0.00 [◊] | 21.60 | 19.87 | 19.80 | 20.13 |
| | | ← | 22.19 | 0.00 | 0.00 [◊] | 28.46 | 28.88 | 28.67 | 29.41 |
| 13 | en-ko | → | 19.80 | 9.62 | 0.00 | 28.94 | 27.92 | 28.40 | 28.81 |
| | | ← | 24.37 | 13.83 | 0.00 | 34.09 | 33.40 | 33.74 | 33.95 |

Table 7: acc@1 for unsupervised methods (1: Artetxe et al. (2018), 2: Grave et al. (2019), 3: Mohiuddin and Joty (2019)) and semi-supervised VecMap with different initial lexicons: MUSE set, identical pairs dataset (ID), our romanized only sets (Rom.), and the union of identical and romanized pairs (ID++). We show both forward (→) and backward (←) directions. In bold the best result for each pair of languages, for “Baselines” and “Our”. Scores from Mohiuddin et al. (2020) are marked with [◊].

*Kannada is not supported by MUSE, so we use the dictionary provided by (Anzer et al., 2020).

| | | | Baselines | | Our | | |
|----|-------|---|--------------|-------------------|-----------------|--------------|--------------|
| | | | Unsup | Semi-sup. MUSE | Semi-supervised | | |
| | | | | | ID | Rom. | ID++ |
| 1 | en-th | → | 0.00 | 27.21 | 27.13 | 26.35 | 26.11 |
| | | ← | 0.00 | 18.93 | 19.83 | 18.25 | 19.83 |
| 2 | en-ja | → | 0.71 | 46.15 | 45.04 | 46.31 | 46.39 |
| | | ← | 0.56 | 39.14 | 38.86 | 40.73 | 39.52 |
| 3 | en-kn | → | 0.00 | 23.78* | 22.03 | 22.90 | 24.23 |
| | | ← | 0.00 | 41.25* | 43.04 | 42.50 | 41.79 |
| 4 | en-ta | → | 0.08 | 20.12 | 19.35 | 18.97 | 19.43 |
| | | ← | 0.08 | 24.60 | 24.60 | 23.71 | 25.00 |
| 5 | en-zh | → | 0.07 | 37.34 | 38.14 | 0.07 | 35.74 |
| | | ← | 0.00 | 32.48 | 34.83 | 0.00 | 32.48 |
| 6 | en-ar | → | 35.44 | 39.70 | 40.23 | 39.24 | 40.15 |
| | | ← | 49.75 | 53.61 | 53.46 | 53.61 | 53.82 |
| 7 | en-hi | → | 42.49 | 42.42 | 42.79 | 42.11 | 42.57 |
| | | ← | 52.46 | 52.62 | 51.99 | 52.07 | 52.23 |
| 8 | en-ru | → | 45.64 | 45.64 | 46.40 | 45.64 | 46.70 |
| | | ← | 64.35 | 64.13 | 64.57 | 64.35 | 64.72 |
| 9 | en-el | → | 48.90 | 49.35 | 48.97 | 49.43 | 49.58 |
| | | ← | 63.87 | 63.80 | 63.87 | 64.56 | 63.72 |
| 10 | en-fa | → | 34.18 | 37.51 | 37.35 | 36.58 | 37.59 |
| | | ← | 41.78 | 43.59 | 44.06 | 43.35 | 43.82 |
| 11 | en-he | → | 42.22 | 42.60 | 42.29 | 42.14 | 42.29 |
| | | ← | 55.92 | 55.70 | 56.00 | 55.62 | 56.08 |
| 12 | en-bn | → | 20.44 | 22.74 | 21.59 | 20.52 | 20.98 |
| | | ← | 25.80 | 30.22 | 30.30 | 30.30 | 30.96 |
| 13 | en-ko | → | 20.30 | 26.57 | 25.63 | 26.02 | 26.49 |
| | | ← | 26.52 | 32.37 | 32.21 | 31.80 | 32.13 |

Table 8: acc@1 on MUSE test sets without proper nouns. Results are reported for unsupervised and semi-supervised Vecmap Artetxe et al. (2018) with different initial lexicons: MUSE set, identical pairs dataset (ID), our romanized only sets (Rom.), and the union of identical and romanized pairs (ID++). We show both forward (→) and backward (←) directions. In bold the best result for each pair of languages, for “Baselines” and “Our”.

*Kannada is not supported by MUSE, so we use the dictionary provided by (Anzer et al., 2020).