

Subword-based Cross-lingual Transfer of Embeddings from Hindi to Marathi

Anonymous ACL submission

Abstract

Word embeddings are growing to be a crucial resource in the field of NLP for any language. This work focuses on static subword embeddings transfer for Indian languages from a relatively higher resource language to a genealogically related low resource language. We work with Hindi-Marathi as our language pair, simulating a low-resource scenario for Marathi. We demonstrate the consistent benefits of unsupervised morphemic segmentation on both source and target sides over the treatment performed by FastText. We show that a trivial “copy-and-paste” embeddings transfer based on even perfect bilingual lexicons is inadequate in capturing language-specific relationships. Our best-performing approach uses an EM-style approach to learning bilingual subword embeddings; the resulting embeddings are evaluated using the publicly available Marathi Word Similarity task as well as WordNet-Based Synonymy Tests. We find that our approach significantly outperforms the FastText baseline on both tasks; on the former task, its performance is close to that of pretrained FastText Marathi embeddings that use two orders of magnitude more Marathi data.

1 Introduction

Subword-level embeddings are useful for many tasks, but require large amounts of monolingual data to train. While about 14 Indian languages such as Hindi, Bengali, Tamil, and Marathi have the required magnitudes of data, most Indian languages are highly under-resourced; they have very little monolingual data and almost no parallel data, and not much digitization. For example, to the best of our knowledge, Marwadi, spoken by 14M people, has no available monolingual corpus; Konkani, spoken by about 3M people, has a monolingual corpus containing 3M tokens, and no parallel data.¹

¹The Opus Corpus (Tiedemann, 2012), one of the most popular collection of parallel texts, contains no parallel data for languages such as Konkani or Mundari.

However, many of these languages have very close syntactic, morphological, and lexical connections to surrounding languages including the mentioned high-resource languages. Our approach aims to leverage these connections in order to build embeddings for these low-resource languages, in the hope that this will aid further development of other NLP tools such as MT or speech tools for these languages.²

While there is a growing interest in shifting towards contextual embeddings with BERT, as well as extending them to low-resource languages, static embeddings retain value in being lightweight and less computationally expensive, especially as studies show that they can perform comparably to contextual embeddings in certain settings (Arora et al., 2020) and encode similar linguistic information (Miaschi and Dell’Orletta, 2020). Thus, an efficient method to develop static embeddings for languages with minimal or no NLP research remains a relevant step to building a basic range of resources in these languages. In this study, we work with Hindi-Marathi as our genealogically and culturally related language pair, and use asymmetric resources (large data for Hindi, artificially small monolingual data for Marathi). We are constrained by the necessity of evaluation datasets for the resulting embeddings.

Most Indian languages are morphologically rich, including Hindi and Marathi. This means that while related language pairs may have a high number of cognates, these may be “disguised” by surrounding inflectional or derivational morphemes. Therefore, even with an identical underlying syntactic structure, lexical correspondences between languages may be obscured or rendered incongruent. Further, when working with small data, the corpus frequencies of fully inflected surface forms would be much less reliable than those of stem

²While some languages may have a little parallel data, we assume none, so as to cater to languages that are just undergoing digitization.

078 and affix morphemes, intuitively resulting in a less
079 robust embeddings transfer. These factors add
080 weight to the intuition that many Indian languages
081 share morpheme-level correspondences with each
082 other. This motivated us to apply unsupervised
083 morphemic segmentation on both the source and
084 target language data; we demonstrate the benefits
085 of doing so in our evaluations. Note that this also
086 makes it natural to work with subword-level em-
087 beddings rather than word embeddings; studies
088 show that the former have an advantage over word
089 embeddings especially for morphologically rich
090 languages. (Chaudhary et al., 2018) (Zhu et al.,
091 2019b) (Li et al., 2018).

092 The idea of the transfer is to project the low-
093 resource language (LRL) subwords into a shared
094 bilingual space with the high-resource language
095 (HRL). We first attempt a trivial transfer that simply
096 finds the “closest” HRL subword for each LRL sub-
097 word, and copies its embedding. We demonstrate
098 that this approach, while tempting, is not enough
099 to capture the relationships between even identical
100 words in both languages; embeddings spaces ap-
101 pear to encode more complex information than this
102 approach would suggest. For our best performing
103 approach, we use the EM-style algorithm described
104 in Artetxe et al. (2017), which alternately optimizes
105 the distance between pairs belonging to a bilingual
106 mapping, and generates a bilingual mapping be-
107 tween words from the resulting bilingual embed-
108 dings. As far as we know, our work is the first to
109 apply this algorithm in the context of embeddings
110 transfer. We compare the resulting Marathi bilin-
111 gual embeddings to a FastText model trained on
112 the available data as well as pretrained models, on
113 the Word Similarity tasks and the WordNet-Based
114 Synonymy Tests.

115 2 Previous Work

116 2.1 Subwords in Embedding Spaces

117 In a seminal work, Bojanowski et al. (2017) present
118 FastText embeddings, that treat morphology by rep-
119 resenting words as bags of chagrams. Kudo and
120 Richardson (2018) present a subword tokenizer for
121 neural text processing, and Kudo (2018) shows the
122 benefits of using multiple subword segmentations
123 in neural machine translation, especially in low-
124 resource settings. Zhu et al. (2019b) look at the
125 segmentation of a word, such as using chagrams,
126 Byte Pair Encoding (BPE) (Gage, 1994; Sennrich
127 et al., 2016), Morfessor, as well as the composi-

128 tion of the subword embeddings (addition, averag-
129 ing, etc.) to construct the final word vector, and
130 conclude that the best performing configuration is
131 highly language and task dependent. A subsequent
132 work (Zhu et al., 2019a) focuses on LRLs and finds
133 the combination of BPE and addition largely robust,
134 although they once again note language-dependent
135 variability. They also find that encoding “affix” in-
136 formation with positional embeddings is beneficial,
137 hinting that the embedding space may distinguish
138 the importance of different kinds of subwords.

139 2.2 Cross-lingual embeddings

140 The problem of learning bilingual embeddings has
141 usually been studied in a symmetric resources sce-
142 nario. Xu et al. (2018) propose an unsupervised
143 method of mapping two sets of monolingual static
144 embeddings into a shared space; they present re-
145 sults for English paired with Spanish, Chinese, and
146 French, evaluated on the bilingual lexicon induc-
147 tion and Word Similarity tasks. Chaudhary et al.
148 (2018) experiment with joint and transfer learn-
149 ing for training bilingual subword embeddings for
150 pairs of Indian LRLs from scratch, by projecting
151 different scripts into the International Phonetic Al-
152 phabet (IPA). Kayi et al. (2020) present an exten-
153 sion of the BiSkip cross-lingual learning objec-
154 tive that leverages subword information to train
155 English-paired bilingual embeddings for LRLs, us-
156 ing around 30K parallel sentences. We describe
157 Artetxe et al. (2017) in some detail below, since we
158 use this algorithm in our approach. There is also
159 growing interest in multilingual contextual embed-
160 dings (Devlin et al., 2018) (Kakwani et al., 2020)
161 (Ruder et al., 2019) such as multilingual BERT;
162 Wang et al. (2020) propose an approach to extend
163 multilingual BERT to low-resource languages with-
164 out retraining it, Pfeiffer et al. (2020) suggest an
165 approach towards incorporating previously unseen
166 scripts into a multilingual BERT model.

167 2.3 Bilingual Lexicon Induction

168 This task is closely related to that of embeddings
169 transfer; we see that these two tasks leverage each
170 other in the literature. Older works such as Koehn
171 and Knight (2002) and Haghghi et al. (2008) use
172 monolingual features such as frequency heuristics,
173 orthographic features, tags, and context vectors in
174 order to find bilingual mappings for mainly Eu-
175 ropean language pairs. Hauer et al. (2017) use
176 word2vec embeddings (Mikolov et al., 2013) in
177 order to iteratively train a translation matrix.

2.4 Summarizing Artetxe et al. (2017)

Artetxe et al. (2017) present an EM-style approach to training bilingual embeddings from monolingual embeddings without parallel data; however, it assumes high quality monolingual embeddings for both languages trained on at least 1 billion word corpora each. Given the two sets of word embeddings, they find a bilingual dictionary D by choosing the closest target word for each source word with respect to the cosine distance between source and target word embeddings. In the next step, they use the dictionary D to calculate a linear transformation matrix that minimizes the sum of cosine distances of the embeddings of all word pairs in D . They apply an orthogonality constraint on the transformation matrix in order to preserve monolingual invariance i.e. to prevent the degradation of the monolingual relationships in the resulting embedding space. These steps are repeated until convergence.

3 Note on languages

Hindi, spoken by about 340M people, is related to other large Indian languages such as Marathi, Punjabi, and Bangla, and has 48 recognized “dialects” over India, which makes it a good choice for the HRL in this project. Hindi is written in the Devanagari script, which is also used for over 120 other (often related) languages, including Marathi. Both Hindi and Marathi are largely free word order; nouns inflect for case and number, verbs inflects for tense, number, gender, and adjectives inflect for gender, case, and number. Some differences are that Marathi exhibits more agglutinative tendencies than Hindi; Marathi is conventionally written in a manner that allows suffix stacking with certain boundary changes. For example, a Marathi token may be a sequence of verb+nominalizing-morpheme+case-marker or noun+postposition+genitive, whereas Hindi separates these morphemes into tokens in many cases (while still exhibiting inflectional and some derivational morphology). See Figure 1.

4 Data and Resources

4.1 Training Data

For Hindi, we used 1M sentences containing roughly 18M tokens from the HindMonoCorp 0.5 (Bojar et al., 2014). For Marathi, we used 50K sentences containing 0.8M tokens from the Indic-



Figure 1: Tokens (with transliterations) in Marathi and Hindi. The stem for “do” is the same (i.e. “kar”) in both languages; Marathi uses one token whereas Hindi uses three.

Corp Marathi monolingual dataset (Kakwani et al., 2020).³ The latter number was chosen because it seems to be the ballpark of the amount of monolingual data collected for newly digitized Indian languages.⁴

4.2 Pretrained Embeddings

We use pretrained FastText embeddings for Hindi, presented by Grave et al. (2018), in line with the assumption that we have good quality resources for the HRL. These embeddings (HIN-PRETR-2G⁵) are trained on the *Wikipedia* corpus as well as *Common Crawl*, containing a total of about 2G tokens. We also use the pretrained Marathi FastText embeddings (MAR-PRETR-60M) presented in the same work, solely for the purpose of evaluation; these embeddings are trained on 50M tokens and 85K Wikipedia articles.

4.3 Evaluation datasets

4.3.1 Word Similarity Dataset

A Word Similarity dataset is a set of word pairs, each annotated by humans according to the degree of similarity (integers ranging from 1 to 10) between the two words. Evaluation is usually performed by finding the cosine similarity between the two words vectors, and calculating the Spearman’s Rank Correlation between the human and model “similarity” judgments for all word pairs. We report this correlation multiplied by 100.

We present results on the Marathi Word Similarity dataset presented by Akhtar et al. (2017), con-

³Note that we do not lemmatize our data; good-quality lemmatizers are a scarce resource that we cannot assume for the LRL.

⁴See <https://www.ldcil.org/resourcesTextCorp.aspx> for efforts on collecting data on under-resourced languages such as Bodo, Dogri, Santhali, etc.

⁵We use the following shorthand to refer to our models unless otherwise specified: <language>-<method_label>-<tokens_of_training_data>. There may be two data slots in the case of bilingual embeddings, containing amount of Marathi and Hindi data respectively.

taining 104 word pairs, as our primary evaluation. This dataset is created by translating a subset of the WordSimilarity-353 English evaluation dataset into Marathi by native Marathi speakers fluent in English, and re-evaluating the similarity scores by 8 native speaker annotators.

4.3.2 WordNet-Based Synonymy Tests

Since the WordSim dataset is rather small, we also perform WordNet-Based Synonymy Tests (WBST) (Piasecki et al., 2018). A WBST consists of a set of “questions” consisting of one “query word”, and N options. One of the options is a synonym or closely related to the query word, while the rest are “detractors”, or randomly selected words. The task is to identify the synonym; we do this by calculating the cosine distances between the query word vector and each of the options and selecting the closest. The reported score is the percentage of correctly answered questions. We use the Marathi WordNet,⁶ containing 32K words, built by Sinha et al. (2006); Debasri et al. (2002), for generating the WBST, and the Python application interface given by Panjwani et al. (2018). Note that we use the WordNet solely for evaluation purposes.

5 Segmentation

Due to the fusional/agglutinative nature of the languages, as well as the morphological and tokenization differences as discussed in Section 3, we apply unsupervised morphemic segmentation to both source and target side data. This is motivated by the need to handle data scarcity on the LRL side, since fully inflected tokens are much rarer than their constituent subwords; we see that the unsegmented Marathi data has 100K distinct tokens, but only 20K distinct “morphemes” post-segmentation. The morphemic segmentation is also an attempt to isolate the morphs in the language data since, according to our hypothesis, it is easier to find correspondences between the two languages at this level rather than at the token level. This is clear in the fact that 50% of the “morphs” in the Marathi segmented data also occur in the Hindi corpus, whereas for the unsegmented data, this is only 20% of tokens. We experimented with BPE and Morfessor and decided to use the latter, since BPE seemed unable to preserve longer morphs regardless of parameter settings. However, this decision may vary according to language type.

⁶See <http://www.cfilt.iitb.ac.in/WordNet/webmwn/>

6 Approach

Our experiments test different intuitions about the cross-lingual interactions between the languages in question. As a baseline, we train a FastText model on the tokenized Marathi data with 0.8M tokens (MAR-BASE-0.8M). We work with 300-dimensional embeddings for all experiments.⁷

6.1 Normalized Edit Distance (NED) Approach

The NED approach is based on finding a bilingual subword-level mapping; it takes advantage of the high number of cognates between related languages as well as the common script. Its primary intuition is that since the languages share not only cognates but also syntactic and morphological properties, embedding vectors can essentially be “copied” over to the LRL from the HRL.

For each Marathi morph, we choose the Hindi morph with the minimum normalized edit distance from it. NED is calculated in the following way:

$$NED(l, h) = \frac{edit_distance(l, h)}{max(length(l), length(h))}$$

To obtain the embedding of any Marathi word, we first segment it. For each subword, we

- Look for the closest Hindi morph by NED
- Retrieve the corresponding Hindi subword embedding

Finally, we compose the subword embeddings, using addition, to give the word embedding. See Algorithm 1 for a depiction.⁸

Algorithm 1: NED Approach

```

l_word ← LRL word;
H_EMB ← HRL embeddings;
l_morphs ← segment_lrl(l_word);
l_subwords_emb ← empty list;
for l_morph in l_morphs do
  | h_closest ← closest_HRL_morph(l_word);
  | append(l_subwords_emb, H_EMB(h_closest));
end
l_emb ← compose_subwords(l_subwords_emb);
return l_emb ;

```

⁷Repeating some experiments for 100 dimensional embeddings spaces, we observe similar trends, with a generally lower performance.

⁸Of course, an NED-based approach is highly limited to related words in the language. However, testing it out gives us an interesting insight about cognates and identical words (see Section 9.1)

6.2 Iterative approach

The iterative approach is based on an algorithm proposed by Artetxe et al. (2017), intended to generate bilingual *word* embeddings for equally well-resourced languages (See Section 2.4 for details). We hypothesize that the algorithm will maintain its quality at the subword level for morphologically rich languages; further, we hypothesize that in our data-asymmetry situation, this approach will serve to “transfer” some of the higher quality of the HRL embedding space to the LRL embeddings, by leveraging a bilingual mapping to induce the relationships already encoded in the HRL embeddings.

As the initial set of LRL embeddings, we use FastText vectors trained on available Marathi segmented data (MAR-SEGM-0.8M). For the HRL, we can use any available resource. We try using pre-trained FastText vectors (HIN-PRETR-2G); we also retrain FastText on the segmented Hindi data (HIN-SEGM-18M). For all runs, we set the initial seed dictionary as identical words⁹ in the source and target corpora.¹⁰ We use the MAR-SEGM-0.8M FastText model as a backoff for unseen morphs, as shown in Algorithm 2. For composing the subword embed-

Algorithm 2: Using bilingual embeddings with backoff

```

L_word ← LRL word;
L_EMB ← Bilinual LRL embeddings;
L_EMB_backup ← MAR-SEGM-0.8M;
l_morphs ← segment_lrl(l_word);
l_subwords_emb ← empty list;
for l_morph in l_morphs do
  l_morph_emb ← empty list ;
  if l_morph in L_EMB then
    | l_morph_emb ← L_EMB(l_word);
  end
  else
    | l_morph_emb ← L_EMB_backup(l_morph);
  end
  append(l_subwords_emb, l_morph_emb);
end
l_emb ← compose_subwords(l_subwords_emb);
return l_emb ;

```

dings of a word, we tried addition, averaging, and also simply picking the first subword embeddings and discarding the rest. The idea behind the last one is that this approximates the stem of the word,

⁹Of course, this is only possible when the languages share a script.

¹⁰Note that this approach does not use any parallel data or bilingual lexicons; this aligns with our assumptions about parallel data. However, in the case that parallel data does exist, it can be used to find a good quality bilingual seed lexicon in lieu of using identical words.; this has been shown to improve the quality of the resulting bilingual embeddings.

| Approach | Score |
|-------------------|--------------|
| MAR-BASE-0.8M | 24.64 |
| MAR-SEGM-0.8M | 43.23 |
| BI-JOINT-0.8M-18M | 35.48 |

Table 1: Marathi monolingual and Marathi-Hindi Joint results on Marathi WordSim task. Notation of models explained in Section 4.2.

| Embeddings | Score |
|----------------|--------------|
| MAR-PRETR-60M | 54.89 |
| MAR-SKIPGR-27M | 41.12 |
| HIN-PRETR-2G | 39.94 |

Table 2: Scores of high-resource Marathi and Hindi models on Marathi WordSim task for comparison.

and also reduces the noise created by summing different subword embeddings.

7 Results: Word Similarity Task

All results are evaluated on the Marathi Word Similarity dataset as explained in Section 4.3.1.

7.1 Baseline and Comparison Models

We show the performance of MAR-BASE-0.8M and MAR-SEGM-0.8M; taking motivation from Chaudhary et al. (2018), we also try a joint approach i.e. we train bilingual embeddings jointly on the segmented Hindi and Marathi data (BI-JOINT-0.8M-18M). See Table 1 for these scores. We observe that simple segmentation of the data causes an improvement of over 20 points, outdoing not only MAR-BASE-0.8M but MAR-SKIPGR-27M (See Table 2). Surprisingly, the joint model BI-JOINT-0.8M-18M dips in performance in comparison to the MAR-SEGM-0.8M. We discuss this effect of the Hindi data on the bilingual embeddings in Section 9.1.

In Table 2, we show the performance of pre-trained FastText Marathi embeddings mentioned in Section 4.2 (MAR-PRETR-60M), as well as the best performing model score from Akhtar et al. (2017) on this evaluation dataset. Akhtar et al. (2017) test different sets of embeddings including Skip-gram, CBOW (Mikolov et al., 2013) and FastText (Bojanowski et al., 2017) algorithms, all trained on a corpus with 27M tokens, of which the Skip-Gram (MAR-SKIPGR-27M) performed best.

Finally, Table 3 shows the performance of the MAR-PRETR-60M and HIN-PRETR-2G on certain word pairs in the Marathi WordSim dataset such that

| Embeddings | Identical Word Score |
|---------------|----------------------|
| HIN-PRETR-2G | 41.17 |
| MAR-PRETR-60M | 50.38 |

Table 3: Scores of pretrained embeddings on word pairs from the Marathi WordSim dataset that are identical in both languages

| Approach | Score |
|-----------------------|--------------|
| BI-SELF-SEGM-0.8M-18M | 43.62 |
| BI-SELF-PRETR-0.8M-2G | 42.72 |
| BI-NED-PRETR-0.8M-2G | 41.85 |
| BI-NED-SEGM-0.8M-18M | 39.37 |

Table 4: Scores on Marathi WordSim for self-mapping and NED strategies, using different Hindi embeddings. Notation: Bi-<mapping_method>-<hin_embs>-<mar_tokens>-<hin_tokens>.

both words are also used identically in Hindi.¹¹ These word pairs were manually identified from the Marathi evaluation dataset; we found that there were 64 such word pairs.¹² Surprisingly, we see a significant dip in the performance of HIN-PRETR-2G on these word pairs as compared to MAR-PRETR-60M, indicating that while the word pairs appear identical in both languages to a native speaker, their usage in the corpora or interaction with other words from the language is different.¹³

7.2 Normalized Edit Distance (NED)

Our NED models use only Hindi embeddings, and project Marathi morphs onto Hindi morphs as shown in Algorithm 1. For further simplicity, we also tried a self-mapping; i.e. we simply calculate the (Hindi) embeddings of the Marathi morphs obtained by segmentation, as they are. Note that this is only possible because Marathi and Hindi share a common script. The resulting embeddings are composed by addition unless otherwise mentioned. See Table 4 for the results on different combinations of embeddings and mappings.

¹¹That is, both of the words in the word pair must be both Hindi and Marathi words with the same spelling, and near-identical senses.

¹²Many of these are transliterations of English words. 24 of the total 135 unique words in the dataset are transliterations, and they occur 40 times i.e. 19.6% times in the 104 word pairs.

¹³Note that HIN-PRETR-2G performs very well on the Hindi WordSim dataset; its monolingual quality is not the problem.

| Approach | Comp. | Score |
|-----------------------|-------------|--------------|
| (MAR-BASE-0.8M | - | 24.64) |
| BI-ITER-PRETR-0.8M-2G | Sum | 44.28 |
| BI-ITER-SEGM-0.8M-9M | Sum | 49.49 |
| BI-ITER-SEGM-0.8M-18M | Sum | 49.21 |
| BI-ITER-SEGM-0.8M-18M | First morph | 50.06 |
| BI-ITER-SEGM-0.8M-36M | First morph | 50.10 |

Table 5: Iterative approach results on Marathi WordSim task using different sets of Hindi embeddings for the crosslingual transfer. Format of the approach name: Bi-Iter-<hin_embs>-<mar_tokens>-<hin_tokens>. **Comp.:** Composition function.

Firstly, we observe that the self-mapping performs better than NED in general.¹⁴ This is largely unsurprising; NED would only perform better for Marathi words that are cognates with Hindi words and show a slight difference in spelling; it will perform competitively with self-mapping for identical words in Hindi and Marathi. As we discuss in Section 7.1, such words form a large part of the evaluation dataset. As for the remaining words, it seems that the Hindi embeddings are able to capture the meaning of the unknown Marathi morphs, perhaps due to similarities at a subword level. Applying the NED mapping, however, can result in Marathi words being mapped to arbitrary Hindi words that may share no semantics with the original Marathi word.

Another interesting observation is that the BI-SELF-SEGM-0.8M-18M performs a little better than BI-SELF-PRETR-0.8M-2G. This affirms our intuition in Section 5 that segmentation on the Hindi side may indeed facilitate the correspondence between subwords common to Hindi and Marathi, leading to better performance on a Marathi evaluation set despite orders of magnitude less (Hindi) data.

7.3 Iterative Approach

This approach trains Marathi bilingual embeddings from Hindi and Marathi monolingual embeddings. The initial Marathi embeddings used are always the monolingual FastText MAR-SEGM-0.8M, whereas we try with some different Hindi embeddings. See results in Table 5.

There are three points of interest in the results:

¹⁴Note that there is a difference between the self-mapping model and directly applying HIN-PRETR-2G as in Table 2 In the former, we segment the Marathi word ourselves and apply Hindi embeddings to the resulting subwords; in the latter, we leave it up to FastText. We note that the former does better.

1. We see that the BI-ITER-SEGM-0.8M-18M outperforms BI-ITER-PRETR-0.8M-2G; i.e. once again, we find that it is better to use embeddings trained on segmented Hindi data for the transfer, even though HIN-SEGM-18M is trained on two orders of magnitude fewer data than HIN-PRETR-2G. Since this approach is explicitly bilingual and attempts to project the Marathi and Hindi embeddings into a shared space, this is a much more direct affirmation that the similarities between Hindi and Marathi are best exploited at the subword level from *both* sides.
2. We see that the “first-morph” manner of composition does slightly better than summing or averaging¹⁵ the subword embeddings.¹⁶
3. Finally, we see that doubling the amount of Hindi data used to train the initial Hindi embeddings does not help. This indicates that the Hindi data is only useful up to a point; we discuss this further in Section 9.1.

8 Results: WordNet-Based Synonymy Tests

Due to the small size of the WordSim dataset, we also carried out an alternate form of evaluation. We generate the WBST questions ourselves from the Marathi WordNet¹⁷ and calculate scores as explained in Section 4.3.2. See Table 6 for the scores.¹⁸ These results confirm some of the findings from the WordSim results; here are some observations from Table 6.

1. Segmentation helps: MAR-SEGM-0.8M consistently outperforms the MAR-BASE-0.8M.
2. The iterative method is the best among the low-resource embeddings.
3. There is little or no difference between BI-ITER-SEGM-0.8M-18M and BI-ITER-SEGM-0.8M-36M: doubling the Hindi data for the bilingual

¹⁵We do not report averaging scores since they are almost identical to the summing scores

¹⁶This could be for several reasons; for example, if the first subword approximates the root of the word, then it may capture most of the meaning, whereas the remaining information may be irrelevant or add noise.

¹⁷For each query word, we randomly select one of its synonyms, and $N - 1$ non-synonym words from the WordNet, under the constraint that all words in the question occur in the corpus at least MIN times. We generate questions for each query word permitted by the value of MIN .

¹⁸Note that since a synonym as well as the detractors are selected randomly from the WordNet, the scores show some variation over different runs.

approach seems not to have much effect on the resulting embeddings.

4. The MAR-PRETR-60M still performs the best, with a seemingly larger margin than in the WordSim task.
5. As MIN increases, the performance of the low-resource methods generally increases; they naturally perform better on words they have seen more frequently in the corpus.

9 Discussion

Some of the clearer findings of our experiments are as regards segmentation and the benefits of a non-trivial bilingual embeddings transfer.

We see repeatedly that segmentation on both sides of the transfer helps the quality of the LRL embeddings. Segmenting the Marathi data causes a large boost in monolingual performance (Table 1); furthermore, when transferring from Hindi embeddings, BI-ITER-SEGM-0.8M-18M outperforms BI-ITER-PRETR-0.8M-2G (Table 5); the Hindi embeddings used in the latter are trained on 2 orders of magnitude higher (unsegmented) data.¹⁹ This suggests that the interaction between the two languages is indeed facilitated at a subword level, validating our bilingual native speaker intuition about the same. We also see that the iterative approach consistently outperforms both monolingual models MAR-BASE-0.8M and MAR-SEGM-0.8M, indicating that bilingual interaction between the related languages is indeed beneficial. This is a good sign for the project of building NLP tools for truly low-resource languages, although the impact of different typologies on this bilingual effect needs to be explored.

Finally, we find that, in agreement with the findings of the papers that investigate subword composition functions (Zhu et al., 2019a,b), the best-performing composition function for subword embeddings seems to be task and data dependent; counter-intuitively, even discarding everything except the first subword seems to work better in some cases than aggregating the embeddings of all parts of the token.

9.1 Using Hindi data

To the best of our knowledge, this is the first work that clearly demonstrates that a trivial “copy-and-paste” transfer approach, such as our NED models,

¹⁹Note that we are talking about performance in terms of the resultant Marathi bilingual embeddings rather than the direct evaluation of the Hindi embeddings.

| (MIN, N) | Test size | MAR-BASE -0.8M | MAR-SEGM -0.8M | BI-ITER-SEGM -0.8M-18M | BI-ITER-SEGM -0.8M-36M | MAR-PRETR -60M |
|----------|-----------|-------------------|-------------------|---------------------------|---------------------------|-------------------|
| (10,6) | 1183 | 51.23 | 58.92 | 61.62 | 57.06 | 84.70 |
| (10,5) | 1183 | 51.90 | 54.78 | 58.66 | 61.54 | 84.87 |
| (20,6) | 684 | 48.98 | 53.65 | 59.94 | 58.19 | 84.50 |
| (20,5) | 684 | 57.89 | 59.94 | 64.47 | 64.33 | 87.57 |
| (50,5) | 293 | 58.02 | 63.14 | 67.24 | 68.94 | 81.23 |

Table 6: WBST Results. *MIN*: minimum frequency of the question and options in the corpus, *N*: the number of total options, Test size: number of questions in the test. The two best-performing models have been bolded.

is not adequate, even when working with two culturally related languages that share a very high percentage of cognates as well as morphosyntactic properties. Our experiments with identical words pairs in Table 3 especially show that even identical words that are not false friends may behave differently depending on the language;²⁰ using Hindi embeddings *directly*, even for identical words, is problematic. We believe that this is an important insight into embeddings transfer that rejects relying on trivial or simplistic approaches.

Many of our experiments are intended to indicate how useful the Hindi data and embeddings are to the Marathi tasks; e.g. we evaluate HIN-PRETR-2G directly on the Marathi WordSim task (Table 2), we experiment with different amounts of Hindi data for both tasks (Tables 5 and 6), and we try a self-mapping with the NED model (see Table 4). We see that doubling or halving the amount of Hindi data does not boost the results for either task and sometimes even harms performance. Similarly, we see that BI-JOINT-0.8M-18M performs worse than MAR-SEGM-0.8M (see Table 1). In conjunction, these results imply that under the current transfer paradigm, adding more Hindi data may sometimes hurt rather than benefit; too much Hindi data for the purpose of training bilingual embeddings may actually “conceal” Marathi word interactions. We invite further investigation of this effect.

10 Future Work

This work is intended to be the pilot in a series of similar studies. We hypothesize that we can obtain similar results for other genealogically related LRL-HRL pairs. We intend to repeat these experiments for language pairs (simulating LRL

environments) such as Punjabi-Hindi, Assamese-Bengali, Konkani-Marathi, and others. Some of the issues we will be working against are different scripts, morphemic segmentation of typologically different languages, and the lack of evaluation data. We would also like to experiment with the integration of parallel data into this approach. We mention one way of doing this in Section 6.2. Finally, we also think it would be interesting to extend our solution from a bilingual to a multilingual one, with multiple sources for a target language. This would be highly pertinent in the case of Indian languages, where even major Indian languages may be interconnected, and regional languages may benefit from the resources of more than one HRL.

11 Conclusion

Embeddings transfer from a high-resource language to a low-resource related language is an important task in today’s scenario of data inequality across languages. We take an Indian language pair, Hindi-Marathi, simulating a low-resource scenario for Marathi, and present an approach to embeddings transfer that uses very little monolingual data on the LRL side, and no parallel data. We believe that our work is the first to show that a “copy-and-paste” embeddings transfer fails even with a perfect bilingual dictionary for a closely related language pair. Our final approach improves significantly over a monolingual FastText baseline for both the WordSim and WBST tasks; its performance on the former task is close to that of high-resource pretrained FastText embeddings. We also demonstrate the benefits of unsupervised morphemic segmentation on both source and target sides for subword-level embeddings transfer.

²⁰This is to say even if words *a* and *b* occur identically and with the same senses in both languages, the word pair (*a*, *b*) may have a different relationship depending on the language.

600
601
602
603
604
605
606
607

608
609
610

611
612
613
614
615

616
617
618
619

620
621
622
623
624
625
626

627
628
629
630
631

632
633
634
635
636

637
638
639
640

641
642

643
644
645
646

647
648
649
650

651
652
653

References

Syed Sarfaraz Akhtar, Arihant Gupta, Avijit Vajpayee, Arjit Srivastava, and Manish Shrivastava. 2017. [Word similarity datasets for Indian languages: Annotation and baseline systems](#). In *Proceedings of the 11th Linguistic Annotation Workshop*, pages 91–94, Valencia, Spain. Association for Computational Linguistics.

Simran Arora, Avner May, Jian Zhang, and Christopher Ré. 2020. [Contextual embeddings: When are they worth it?](#) *CoRR*, abs/2005.09117.

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

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

Ondřej Bojar, Vojtěch Diatka, Pavel Rychlý, Pavel Straňák, Vít Suchomel, Aleš Tamchyna, and Daniel Zeman. 2014. [HindMonoCorp 0.5](#). LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

Aditi Chaudhary, Chunting Zhou, Lori Levin, Graham Neubig, David R Mortensen, and Jaime G Carbonell. 2018. Adapting word embeddings to new languages with morphological and phonological subword representations. *arXiv preprint arXiv:1808.09500*.

Chakrabarti Debasri, Narayan Dipak Kumar, Pandey Prabhakar, and Bhattacharyya Pushpak. 2002. Experiences in building the Indo-Wordnet: A Wordnet for Hindi. In *Proceedings of the First Global WordNet Conference*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Philip Gage. 1994. A new algorithm for data compression. *C Users Journal*, 12(2):23–38.

Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning word vectors for 157 languages. *arXiv preprint arXiv:1802.06893*.

Aria Haghighi, Percy Liang, Taylor Berg-Kirkpatrick, and Dan Klein. 2008. Learning bilingual lexicons from monolingual corpora. In *Proceedings of ACL-08: Hlt*, pages 771–779.

Bradley Hauer, Garrett Nicolai, and Grzegorz Kondrak. 2017. Bootstrapping unsupervised bilingual lexicon induction. In *Proceedings of the 15th Conference of*

the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, pages 619–624. 654
655
656

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, NC Gokul, Avik Bhattacharyya, Mitesh M Khapra, and Pratyush Kumar. 2020. [inlpsuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for Indian languages](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 4948–4961. 657
658
659
660
661
662
663
664

Efsun Sarioglu Kayi, Vishal Anand, and Smaranda Muresan. 2020. [Multiseg: Parallel data and subword information for learning bilingual embeddings in low resource scenarios](#). In *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)*, pages 97–105. 665
666
667
668
669
670
671
672

Philipp Koehn and Kevin Knight. 2002. Learning a translation lexicon from monolingual corpora. In *Proceedings of the ACL-02 workshop on Unsupervised lexical acquisition*, pages 9–16. 673
674
675
676

Taku Kudo. 2018. [Subword regularization: Improving neural network translation models with multiple subword candidates](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics. 677
678
679
680
681
682
683

Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics. 684
685
686
687
688
689
690

Bofang Li, Aleksandr Drozd, Tao Liu, and Xiaoyong Du. 2018. [Subword-level composition functions for learning word embeddings](#). In *Proceedings of the Second Workshop on Subword/Character Level Models*, pages 38–48, New Orleans. Association for Computational Linguistics. 691
692
693
694
695
696

Alessio Miaschi and Felice Dell’Orletta. 2020. [Contextual and non-contextual word embeddings: an in-depth linguistic investigation](#). In *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 110–119, Online. Association for Computational Linguistics. 697
698
699
700
701
702

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient Estimation of Word Representations in Vector Space](#). 703
704
705

Ritesh Panjwani, Diptesh Kanojia, and Pushpak Bhattacharyya. 2018. [pyiwn: A Python based API to access Indian Language WordNets](#). In *Proceedings of the 9th Global Wordnet Conference*, pages 378–383. 706
707
708
709

- 710 Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Se-
711 bastian Ruder. 2020. Unks everywhere: Adapting
712 multilingual language models to new scripts. *arXiv*
713 *preprint arXiv:2012.15562*.
- 714 Maciej Piasecki, Gabriela Czachor, Arkadiusz Janz,
715 Dominik Kaszewski, and Paweł Kędzia. 2018.
716 [Wordnet-based evaluation of large distributional](#)
717 [models for Polish](#). In *Proceedings of the 9th Global*
718 *Wordnet Conference*, pages 229–238, Nanyang
719 Technological University (NTU), Singapore. Global
720 Wordnet Association.
- 721 Sebastian Ruder, Ivan Vulić, and Anders Søgaard.
722 2019. A survey of cross-lingual word embedding
723 models. *Journal of Artificial Intelligence Research*,
724 65:569–631.
- 725 Rico Sennrich, Barry Haddow, and Alexandra Birch.
726 2016. Neural machine translation of rare words
727 with subword units. In *Proceedings of the 54th An-*
728 *annual Meeting of the Association for Computational*
729 *Linguistics (Volume 1: Long Papers)*, pages 1715–
730 1725.
- 731 Manish Sinha, Mahesh Reddy, and Pushpak Bhat-
732 tacharyya. 2006. An approach towards construction
733 and application of multilingual Indo-Wordnet. In
734 *3rd Global Wordnet Conference (GWC 06)*, Jeju Is-
735 land, Korea. Citeseer.
- 736 Jörg Tiedemann. 2012. Parallel Data, Tools and Inter-
737 faces in OPUS. In *Proceedings of the Eight Interna-*
738 *tional Conference on Language Resources and Eval-*
739 *uation (LREC’12)*, Istanbul, Turkey. European Lan-
740 guage Resources Association (ELRA).
- 741 Zihan Wang, Stephen Mayhew, Dan Roth, et al. 2020.
742 Extending multilingual BERT to low-resource lan-
743 guages. *arXiv preprint arXiv:2004.13640*.
- 744 Ruochen Xu, Yiming Yang, Naoki Otani, and Yuexin
745 Wu. 2018. Unsupervised cross-lingual transfer of
746 word embedding spaces. In *Proceedings of the 2018*
747 *Conference on Empirical Methods in Natural Lan-*
748 *guage Processing*, pages 2465–2474.
- 749 Yi Zhu, Benjamin Heinzerling, Ivan Vulić, Michael
750 Strube, Roi Reichart, and Anna Korhonen. 2019a.
751 On the importance of subword information for mor-
752 phological tasks in truly low-resource languages. In
753 *Proceedings of the 23rd Conference on Computa-*
754 *tional Natural Language Learning (CoNLL)*, pages
755 216–226.
- 756 Yi Zhu, Ivan Vulić, and Anna Korhonen. 2019b. A sys-
757 tematic study of leveraging subword information for
758 learning word representations. In *Proceedings of the*
759 *2019 Conference of the North American Chapter of*
760 *the Association for Computational Linguistics: Hu-*
761 *man Language Technologies, Volume 1 (Long and*
762 *Short Papers)*, pages 912–932.