

DM-BLI: Dynamic Multiple Subspaces Alignment for Unsupervised Bilingual Lexicon Induction

Anonymous ACL submission

Abstract

Unsupervised bilingual lexicon induction (BLI) task aims to find word translations between languages and has achieved great success in similar language pairs. However, related works mostly rely on a single linear mapping for language alignment and fail on distant or low-resource language pairs, achieving less than half the performance observed in rich-resource language pairs. In this paper, we introduce **DM-BLI**, a **D**ynamic **M**ultiple subspaces alignment framework for unsupervised **BLI**. DM-BLI improves language alignment by utilizing multiple subspace alignments instead of a single mapping. We begin via unsupervised clustering to discover these subspaces in source embedding space. Then we identify and align corresponding subspaces in the target space using a rough global alignment. DM-BLI further employs intra-cluster and inter-cluster contrastive learning to refine precise alignment for each subspace pair. Experiments conducted on standard BLI datasets for 12 language pairs (6 rich-resource and 6 low-resource) demonstrate substantial gains achieved by our framework. We release our code to facilitate the community.

1 Introduction

Unsupervised bilingual lexicon induction (BLI) has shown to be a key multilingual NLP task to align cross-lingual word embeddings (CLWE) (Mikolov et al., 2013a; Ruder et al., 2019) and bridge lexical gap between low-resource languages (Eder et al., 2021; Marchisio et al., 2022).

Existing BLI approaches can be roughly divided into two categories: mapping-based methods (Conneau et al., 2017; Artetxe et al., 2018; Ren et al., 2020; Li et al., 2022) and generation-based methods (Gonen et al., 2020; Ghazvininejad et al., 2023; Li et al., 2023). Mapping-based methods aim to align monolingual embeddings from various languages into a shared CLWE space via linear or non-linear projections. Generation-based methods

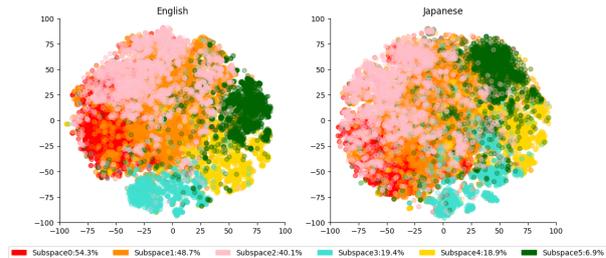


Figure 1. t-SNE visualization of the clustered monolingual word embedding in a distant language pair of English (left) and Japanese (right). Different colors represent different subspaces. With a global orthogonal mapping from English to Japanese, BLI accuracies for subspaces 0-5 are 54.3%, 48.7%, 40.1%, 19.4%, 18.9% and 6.9%, respectively.

leverage the machine translation capacities of large language models (LLMs) (Briakou et al., 2023) to directly generate word translations via zero-shot or few-shot prompting. Mapping-based methods are superior to generation-based methods in unsupervised settings, especially are far superior on low-resource languages (Li et al., 2023), primarily due to the unbalanced training corpus size of each language supported by LLMs (Zhu et al., 2023a).

The existing fully unsupervised mapping-based approaches still need to carefully address two issues. First, these approaches rely on the strong assumption that monolingual word embedding spaces are isomorphic and the mapping matrix should be under orthogonal constraint, but this assumption does not hold true for all languages (Søgaard et al., 2018; Glavaš et al., 2019), especially for distant language pairs (Ormazabal et al., 2019; Vulić et al., 2019). Therefore, weak orthogonal constraints have been proposed to tackle this issue (Mohiuddin et al., 2020; Glavaš and Vulić, 2020).

Second, a global mapping matrix does not consistently perform optimally across all subspaces (Nakashole, 2018; Wang et al., 2020). As shown in Figure 1, subspaces exhibit inconsistent structural similarity. With a global orthogonal mapping,

068 BLI accuracy varies among different subspaces:
069 the highest accuracy is 54.3% in subspace 0 and
070 the lowest accuracy is 6.89% in subspace 5. To al-
071 leviate the issue, recent research proposed a multi-
072 adversarial learning method (Wang et al., 2020)
073 and a graph-based paradigm (Ren et al., 2020) to
074 learn or refine a specific mapping for each subspace.
075 However, in these approaches, multiple subspaces
076 assigned by initial mappings are static. Once initial
077 solutions of these mappings are not good enough,
078 they may get stuck in poor local optima.

079 Different from previous methods, we propose a
080 Dynamic Multiple subspaces cross-lingual align-
081 ment framework for fully unsupervised Bilingual
082 Lexicon Induction, named DM-BLI. It leverages
083 intra-cluster and inter-cluster contrastive learning
084 to achieve precise alignment at subspace level for
085 both source and target languages, along with dy-
086 namically updating the subspace assignment of
087 each word. DM-BLI starts by clustering the embed-
088 dings of source language to establish multiple valid
089 subspaces. Then, we induce an initial solution to
090 discover corresponding multiple subspaces in the
091 target language. Finally, we iteratively refine a pair
092 of specific mappings for each subspace pair until
093 convergence is reached.

094 In summary, we make the following contribu-
095 tions:

- 096 • We propose a dynamic multiple subspaces
097 cross-lingual alignment framework for the
098 BLI task, which achieves customized map-
099 pings for each subspace pair.
- 100 • To boost the performance of our model, we de-
101 sign a contrastive learning framework includ-
102 ing intra-cluster and inter-cluster level based
103 on unsupervised clustering to dynamically up-
104 date the subspace assignment, avoiding falling
105 into local optima.
- 106 • We conduct extensive experiments to demon-
107 strate the effectiveness of our method on
108 twelve language pairs including six rich-
109 resource and six low-resource language pairs,
110 and DM-BLI achieves significant improve-
111 ments especially for distant and low-resource
112 language pairs.

113 2 Related Work

114 2.1 Cross-lingual Word Embedding

115 Bilingual lexicons can be induced via nearest neigh-
116 bour retrieval on CLWE, which represent lexical

117 words from two or more languages in a shared
118 space.

119 Based on whether parallel corpora are used
120 or not, CLWE approaches can be categorized
121 into three groups: supervised (Faruqui M, 2014;
122 Zou W Y, 2013; Vulić I, 2015), semi-supervised
123 (Artetxe M, 2017; Patra et al., 2019), and unsuper-
124 vised approaches (Conneau et al., 2017; Artetxe
125 et al., 2018). Because parallel corpora are not avail-
126 able for many languages, unsupervised approaches
127 gain much more attention.

128 But unsupervised methods do not require any
129 seed dictionary at all, it is more difficult to induce
130 a reliable initial solution which plays a crucial role
131 in alignment. Therefore, GAN-based adversarial
132 training (Zhang et al., 2017), optimal transport so-
133 lution (Alvarez-Melis and Jaakkola, 2018), Auto-
134 encoder (Mohiuddin and Joty, 2019), and graph-
135 based alignment (Ren et al., 2020) were utilized
136 to better match embedding distribution and find a
137 better initial solution in a fully unsupervised way.

138 Based on the type of pre-trained monolingual
139 embeddings, CLWE can be divided into two
140 groups: static CLWE and contextual CLWE. Most
141 works focused on static word embeddings (Ruder
142 et al., 2019), which can be derived by Word2Vec
143 (Mikolov et al., 2013b) or fastText (Bojanowski
144 et al., 2016). However, static embeddings lack con-
145 textual information to capture polysemy. Therefore,
146 contextual embeddings, generated from monolin-
147 gual and multilingual pre-trained language models
148 (Devlin et al., 2019; Lample and Conneau, 2019),
149 were utilized as input monolingual embeddings.
150 However, they cannot surpass static embedding
151 in the BLI task based on the same mapping tech-
152 nologies even with much more training time (Vulić
153 et al., 2020; Liu et al., 2021).

154 2.2 Bilingual Lexicon Induction

155 Bilingual lexicon induction is the task of inducing
156 word translations from monolingual corpora of two
157 languages.

158 Existing BLI approaches achieved promising
159 performance on semantically similar and rich-
160 resource language pairs, but were still far from
161 satisfied on distant and low-resource language
162 pairs. For example, unsupervised BLI accuracy
163 on English-Spanish exceeded 80%, while under
164 40% on English-Chinese (e.g. Conneau et al.,
165 2017; Wang et al., 2020; Ren et al., 2020). In low-
166 resource language pairs like Bulgarian-Hungarian,

LLaMA_{13B} achieved 23.61% accuracy, whereas VecMap (Artetxe et al., 2018) achieved 39.24% (Li et al., 2023).

To address this issue, Taitelbaum et al. (2019) suggested leveraging auxiliary languages to bridge the gap between semantically distant and low-resource language pairs. Based on the observation that words are naturally grouped into different semantic subspaces and the BLI accuracies of different subspaces are not uniform, Wang et al. (2020) proposed a multi-adversarial learning method to learn a specific mapping for each subspace. However, this GAN-based method was less robust and its assignment of subspaces was fixed initially which would bring the noise of initial solution.

Different from previous work, we propose a dynamic multiple subspaces alignment framework for unsupervised BLI to achieve more robust and precise alignment at subspace level for both source and target languages, along with dynamically updating the subspace assignment of each word.

3 Methodology

3.1 Formulation

Given the source and target languages, let X and Y be the normalized pre-trained monolingual embeddings for source and target languages, respectively. Our goal is to find the optimal mapping matrices W_X^* and W_Y^* , with which XW_X^* and YW_Y^* are projected in a shared CLWE space, where semantically similar words across languages are close to each other.

Figure 2 illustrates the four procedural steps of our BLI method: multiple subspaces clustering on the source language, initial alignment, intra-cluster and inter-cluster contrastive refinement, and bilingual lexicon induction.

3.2 Multiple Subspaces Discovery

Multiple subspaces discovery contains the first two steps in Figure 2: multiple subspaces clustering and initial alignment. It aims to find pairs of subspaces $\{C_{s_i}, C_{t_i}\}$ from the source and target languages, where $i = 1, 2, \dots, K$ and K is the number of subspaces.

Multiple subspaces clustering is only carried on source language embedding to obtain K subspaces. Let $C_{s_i} = \{v_1^{s_i}, v_2^{s_i}, \dots, v_n^{s_i}\}$ be the i -th subspace, where $v_k^{s_i}$ is the k -th word in source subspace C_{s_i} and n is the number of words in C_{s_i} . A major challenge in multiple subspaces clustering is to

determine the optimal number of subspaces in advance. To tackle this issue, we use a parameter-free hierarchical clustering called First Integer Neighbor Clustering Hierarchy (FINCH) (Sarfraz et al., 2019) to provide a reference number K . Then, K-means algorithm (MacQueen et al., 1967) is used to cluster X into K subspaces.

Then, an initial alignment is conducted for identifying corresponding K subspaces in the target language, denoted as $C_{t_i} = \{v_1^{t_i}, v_2^{t_i}, \dots, v_m^{t_i}\}$, where $i = 1, 2, \dots, K$ and $v_j^{t_i}$ is the j -th word in target subspace C_{t_i} . Specifically, we operate the initial alignment following (Artetxe et al., 2018) to get a pair of global initial mapping matrices W_X and W_Y , with which we can retrieve the translation of each target word in the source language. Subsequently, the subspace index of the target word is set to be the subspace index of its translation.

3.3 Multiple Subspaces Contrastive Refinement

A single global mapping does not consistently perform optimally across all subspaces (Nakashole, 2018; Wang et al., 2020). Therefore, the proposed framework will dynamically refine matrices for each subspace pair. This framework contains both inter-cluster and intra-cluster contrastive learning. Inter-cluster contrastive learning ensures the distinguishability of features from different subspaces, thereby facilitating more effective customized mapping. Intra-cluster contrastive learning brings translation pairs within the subspace closer together, while push non-translation pairs further apart, thus achieving finer-grained alignment. The whole refinement process will be completed subspace by subspace.

3.3.1 Inter-cluster Contrastive Learning

Given the subspace pair $\{C_{s_i}, C_{t_i}\}$, inter-cluster contrastive learning aims to bring the whole subspaces C_{s_i} closer to C_{t_i} , while pushing it away from other non-corresponding subspaces $C_{t_j, i \neq j}$.

We introduce optimal transport distance as the metric to evaluate distance of two subspaces distribution, in our work Wasserstein distance (Han et al., 2022) has been applied. The Wasserstein distance between the distributions of two subspaces can be calculated as:

$$D_w(C_{s_i}, C_{t_i}) = \min_{T \in \pi(C_{s_i}, C_{t_i})} \sum_{j=1}^n \sum_{k=1}^m T_{jk} c(v_j^{s_i}, v_k^{t_i}) \quad (1)$$

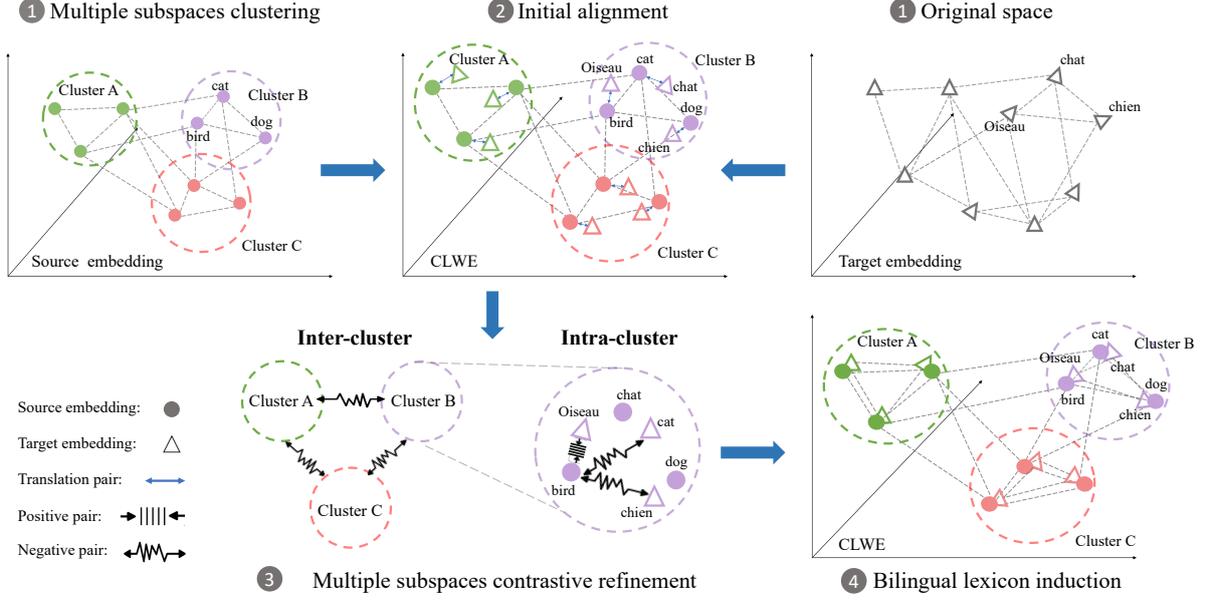


Figure 2. An illustration of the proposed DM-BLI framework. ❶ represents the monolingual word embedding spaces of source and target language, where English is the source language denoted by circles while French is target language denoted by triangles. Multiple subspaces clustering is only applied to source language(English) and different colors represent different subspaces. ❷ represents a cross-lingual word embedding space via an initial alignment. ❸ is a multiple subspaces contrastive learning refinement block aiming to push away words from different clusters and pull closer the words being translation for each other closer within the cluster. ❹ represents refined cross-lingual word embedding space, where words being translation for each other stay closer.

where $c(v_j^{t_i}, v_k^{s_i})$ is the transport cost between words $v_j^{t_i} \in C_{t_i}$ and $v_k^{s_i} \in C_{s_i}$, and T_{jk} represents the transport plan between $v_j^{t_i}$ and $v_k^{s_i}$.

Based on the K pairs of subspaces, we calculate a bi-direction inter-cluster contrastive learning loss as follows:

$$\begin{aligned} \mathcal{L}_{s2t} &= -\frac{1}{K} \left\{ \log(e^{-\mathbf{D}_w(C_{s_i}, C_{t_i})/\tau}) \right. \\ &\quad \left. + \sum_{j \neq i} \log(1 - e^{-\mathbf{D}_w(C_{s_i}, C_{t_i})/\tau}) \right\} \\ \mathcal{L}_{t2s} &= -\frac{1}{K} \left\{ \log(e^{-\mathbf{D}_w(C_{t_i}, C_{s_i})/\tau}) \right. \\ &\quad \left. + \sum_{j \neq i} \log(1 - e^{-\mathbf{D}_w(C_{t_i}, C_{s_i})/\tau}) \right\} \quad (2) \end{aligned}$$

where τ is a temperature parameter. To be specific, the aforementioned process is applied to the sampled distribution of subspace, where the proportion of samples is determined by a preset threshold.

Finally, we obtain the final inter-cluster contrastive loss \mathcal{L}_{inter} as below, where λ is the trade-off set to be 0.5 between two directions:

$$\mathcal{L}_{inter} = \lambda * \mathcal{L}_{s2t} + (1 - \lambda) * \mathcal{L}_{t2s} \quad (3)$$

3.3.2 Intra-cluster Contrastive Learning

Given the subspace pair $\{C_{s_i}, C_{t_i}\}$, intra-cluster contrastive learning is to ensure word pair $(v_j^{s_i}, v_k^{t_i})$ are closer, which are translations to each other in C_{s_i} and C_{t_i} .

Based on the mapping matrices W_X and W_Y , we can initially construct a bilingual dictionary \mathbf{D} by retrieving the translation of each target word in the source language, where $\mathbf{D} = \{(v_1^{t_i}, v_1^{s_i}), (v_2^{t_i}, v_2^{s_i}), \dots, (v_n^{t_i}, v_n^{s_i})\}$ and n is the number of words in \mathbf{D} .

However, the quality of \mathbf{D} depends on the quality of mapping matrices. To alleviate the noise brought by the current solution, we selectively sample high-confidence word translation pairs from \mathbf{D} , where confidence is determined by the similarity gap between the selected translation and the second candidate translation with the source word.

Based on the sampled translation pairs, the intra-cluster contrastive learning loss can be defined as:

$$\mathcal{L}_{intra} = -\sum_{i=1}^s \log \frac{e^{\text{sim}(v_i^s, v_i^t)/\tau}}{\sum_{j=1}^s e^{\text{sim}(v_i^s, v_j^t)/\tau}} \quad (4)$$

Where s is the number of sampled translation pairs and τ is a temperature parameter. Ultimately, the loss of the whole contrastive refinement can be

defined as follows:

$$\mathcal{L} = \mathcal{L}_{inter} + \mathcal{L}_{intra} \quad (5)$$

3.4 Multiple Subspaces Dynamic Updating

A single round of subspace assignment may introduce noise from the initial solution, potentially causing CLWEs to fall into local optima. Therefore, we propose to dynamically adjust the subspace assignment of each word in target language during the process of updating W_X and W_Y .

To clarify, the assignment of multiple subspaces in source language $C_s = \{C_{s_1}, C_{s_2}, \dots, C_{s_K}\}$ is fixed once the clustering process is completed. For word v_i^t in target language, its translation from source language v_i^s is retrieved based on XW_X and YW_Y . The subspace index of v_i^s will be assigned to v_i^t . Upon updating W_X and W_Y , the subspace assignment of v_i^t will be adjusted accordingly to maintain consistency whenever its translation changes.

As we mentioned before, the whole refinement process will be operated subspace by subspace. For each subspace C_{t_i} in target language, the whole dynamic updating procedure stops until convergence is reached. Convergence can be determined by measuring the overlap of target words within C_{t_i} between the current and previous rounds. Besides, once a subspace has achieved convergence, its assignments are finalized, ensuring that the words within it remain unchanged in their respective subspaces. The whole methodology is summarised in Algorithm 1.

Algorithm 1: Dynamic Multiple Subspaces Alignment for Unsupervised BLI

Input: Monolingual word embedding spaces X, Y

Output: $\{W_{x_i}^*\}_{i=1}^K, \{W_{y_i}^*\}_{i=1}^K$

- 1 $\{C_{s_i}\}_{i=1}^K \leftarrow$ Apply Clustering on X ;
 - 2 $W_X, W_Y \leftarrow$ Initial Alignment ;
 - 3 $\{C_{t_i}\}_{i=1}^K \leftarrow$ Calculate XW_X, YW_Y ;
 - 4 **for** $i \leq K$ **do**
 - 5 Initialize W_{x_i}, W_{y_i} with W_X, W_Y ;
 - 6 **while not convergence do**
 - 7 $W_{x_i}, W_{y_i} \leftarrow$ Optimise loss
 - $C_{s_i} \leftarrow$ Keep C_{s_i} fixed
 - $C_{t_i} \leftarrow$ Update C_{t_i} with W_{x_i}, W_{y_i}
 - 8 **return** $\{W_{x_i}^*\}_{i=1}^K, \{W_{y_i}^*\}_{i=1}^K$;
-

4 Experiment Setup

We evaluate our framework in both supervised and unsupervised BLI tasks on 12 language pairs, which contain 6 rich-resource language pairs: Spanish (ES), German (DE), Russian (RU), Arabic (AR), Japanese (JA) and Chinese (ZH), all cross-lingual to English (EN) and six low-resource language pairs: Finnish (FI), Hindi (HI), Turkish (TR), Indonesian (ID), Bulgarian (BG) and Catalan (CA), all cross-lingual to English (EN).

4.1 Dataset

We use fastText vectors trained on full Wikipedias for each language (Bojanowski et al., 2016) as monolingual word embeddings. We use the widely used MUSE bilingual lexicon (Conneau et al., 2017), released by Facebook, as ground truth lexicon. MUSE provides 110 bilingual lexicons and each lexicon contains the 6,500 most frequently used words in each language, split in a test set of 1,500 words and a training set of 5,000.

4.2 Baselines

Baselines are divided into supervised and unsupervised two lines as described below. We run the released code of each baseline in our experiments.

Supervised BLI

MUSE: Conneau et al. (2017) learned an orthogonal map by minimizing the Euclidean distance between the supervised translation pairs.

VecMap: Artetxe et al. (2018) used a multi-step framework consisting of several steps: whitening, orthogonal mapping, re-weighting, de-whitening, and dimensionality reduction.

BLISS: Patra et al. (2019) proposed a semi-supervised approach with a weak orthogonality constraint in the form of a back-translation loss.

CL-BLI: Li et al. (2023) proposed a robust and effective two-stage contrastive learning framework to combine static and contextual embeddings.

Unsupervised BLI

MUSE: Unsupervised MUSE (Conneau et al., 2017) used adversarial training and iterative Procrustes refinement.

VecMap: Unsupervised VecMap (Artetxe et al., 2018) used intra-linguistic word similarity information to induce initial solution.

Ad. : Mohiuddin and Joty (2019) proposed an adversarial auto-encoder framework, where adversarial mapping was done at the latent embedding space.

| Method | Precision@1 | | | | | | Precision@5 | | | | | | Avg. |
|----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | FI-* | HI-* | TR-* | ID-* | BG-* | CA-* | FI-* | HI-* | TR-* | ID-* | BG-* | CA-* | |
| Supervised | | | | | | | | | | | | | |
| MUSE | 46.50 | 25.65 | 39.82 | 35.56 | 39.28 | 46.19 | 66.07 | 39.17 | 57.56 | 50.92 | 56.62 | 60.52 | 46.99 |
| BLISS | 49.94 | 28.17 | 41.45 | 38.49 | 42.21 | 47.26 | 68.97 | 42.43 | 59.39 | 54.05 | 59.51 | 61.94 | 49.48 |
| VecMap | <u>58.12</u> | <u>34.07</u> | <u>49.37</u> | <u>44.72</u> | <u>49.13</u> | <u>54.35</u> | 75.43 | <u>48.40</u> | 66.24 | 59.52 | <u>64.62</u> | 66.84 | <u>55.90</u> |
| CL-BLI | 57.78 | 32.62 | 48.52 | 43.43 | 47.34 | 53.89 | <u>75.97</u> | 47.02 | 59.93 | 58.63 | 64.20 | <u>67.09</u> | 54.70 |
| DM-BLI | 60.29 | 35.57 | 53.09 | 48.24 | 50.80 | 56.47 | 77.08 | 49.24 | 69.11 | 62.09 | 66.16 | 68.57 | 58.06 |
| Unsupervised | | | | | | | | | | | | | |
| MUSE | 0.05 | 0.00 | 36.82 | 36.35 | 38.31 | 46.07 | 0.05 | 0.05 | 54.76 | 51.65 | 55.05 | 60.51 | 31.64 |
| VecMap | 54.71 | 28.19 | 48.92 | 45.65 | 45.69 | <u>53.52</u> | <u>71.72</u> | 41.54 | <u>65.25</u> | <u>59.76</u> | 61.24 | <u>65.63</u> | 53.49 |
| Ad. | 0.45 | 0.01 | 46.69 | 0.09 | 0.03 | 53.06 | 1.47 | 0.03 | 63.08 | 0.31 | 0.11 | 65.55 | 19.24 |
| BLOOM _{7B} | 23.43 | 28.30 | 30.82 | 45.45 | 16.75 | 43.89 | 25.75 | 28.54 | 34.08 | 49.77 | 16.94 | 48.01 | 32.64 |
| Llama _{13B} | 40.98 | 30.68 | 44.90 | 48.63 | <u>56.86</u> | 48.83 | 41.64 | 30.69 | 45.24 | 48.95 | 57.16 | 49.19 | 45.31 |
| GPT-3.5 | 60.37 | 56.11 | 54.49 | <u>48.37</u> | 67.51 | 45.15 | 64.33 | 57.40 | 55.99 | 49.35 | 69.53 | 45.78 | 56.19 |
| DM-BLI | <u>57.48</u> | <u>30.80</u> | <u>51.98</u> | 48.81 | 47.63 | 56.15 | 74.10 | <u>43.75</u> | 67.95 | 62.46 | <u>63.36</u> | 67.61 | <u>56.00</u> |

Table 1. Precision@1 and Precision@5 for the BLI task on six low-resource language pairs, where * represents EN(English). The best score is shown in **bold**, and the suboptimal score is shown in underlined.

BLOOM_{7B} (Workshop et al., 2022): It is a decoder-only Transformer language model that supports 46 natural languages. 7B parameters version was used in our experiment.

Llama_{13B} (Touvron et al., 2023): It is a decoder-only LLM which supports 20 languages. 13B parameters version was used in our experiment.

GPT-3.5 (Brown et al., 2020): It is a decoder-only LLM with 175B parameters, supported by 38 languages. GPT-3.5-turbo was used in our experiment.

4.3 Implementation details

We choose the most 7,500 frequent vocabularies of each language. The normalization procedure for pre-trained embedding contains three steps: length normalizes the embeddings, then mean centers each dimension, and then length normalizes them again.

For multiple subspaces discovery, the number of subspaces is set to be 9 and we will discuss the impact of this setting later. For inter-cluster contrastive learning, only words with weight above 0.45 are sampled to represent the subspace distribution. For intra-cluster contrastive learning, we only sample the top 20% of word translation pairs sorted descending by confidence.

Following the previous research (Patra et al., 2019), the prompt template for **Llama**_{13B} is defined as: "Translate from L^x to L_y : $w^x \Rightarrow$ "; the prompt template for **GPT-3.5** is defined as: "Translate the L^x word w^x into L_y ". Both of them are provided as the best template for each of them in

Li et al. (2023).

The evaluation for BLI is done by comparing the bilingual lexicon constructed by each model with the benchmark lexicon MUSE (Conneau et al., 2017) and reporting precision Precision@N for $N = 1, 5$. Precision@N accounts for accuracy for which the correct translation of the source words is in the N -th nearest neighbors based on CSLS (Conneau et al., 2017).

5 Result and Discussion

5.1 Results in Low-resource Languages

Table 1 summarizes the results of the supervised and unsupervised BLI tasks in low-resource language pairs. In both tasks, our proposed method shows significant improvements, particularly in Precision@5, with an average of 2.16 points higher than the strongest baseline VecMap in the supervised task. In the unsupervised task, our method performs nearly as well as the strong baseline GPT-3.5.

In the supervised task, DM-BLI outperforms all the baseline methods on all language pairs, demonstrating the robustness and effectiveness of our framework on low-resource language pairs. In the unsupervised task, DM-BLI outperforms all the baseline methods on four out of six language pairs and archives suboptimal scores in the remaining pairs at Precision@5. It demonstrates that our method is competitive even compared with GPT-3.5, which has 175B parameters and supports 38 languages. The unsatisfied performance

| Method | Precision@1 | | | | | | Precision@5 | | | | | | Avg. |
|----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | ES-* | DE-* | RU-* | AR-* | JA-* | ZH-* | ES-* | DE-* | RU-* | AR-* | JA-* | ZH-* | |
| Supervised | | | | | | | | | | | | | |
| MUSE | 67.80 | 63.14 | 53.23 | 44.33 | 0.14 | 8.29 | 78.13 | 75.86 | 70.19 | 61.16 | 0.41 | 18.87 | 45.13 |
| BLISS | 68.46 | 63.49 | 54.88 | 45.70 | 0.01 | 6.43 | 78.86 | 76.69 | 71.28 | 62.47 | 0.04 | 14.00 | 45.19 |
| VecMap | 71.70 | 66.46 | 59.58 | 51.54 | <u>37.14</u> | 42.50 | 80.43 | 78.22 | 74.69 | 67.00 | 53.65 | 62.23 | 62.10 |
| CL-BLI | 73.02 | 69.00 | <u>61.31</u> | 53.14 | 35.07 | <u>42.44</u> | 81.71 | 80.28 | 77.10 | 68.95 | <u>50.68</u> | <u>62.26</u> | <u>62.91</u> |
| DM-BLI | <u>72.87</u> | <u>68.28</u> | 61.61 | <u>52.33</u> | 41.03 | 44.83 | 81.16 | <u>79.35</u> | <u>76.35</u> | <u>67.80</u> | 56.94 | 64.13 | 63.89 |
| Unsupervised | | | | | | | | | | | | | |
| MUSE | 67.89 | 63.27 | 50.49 | 0.03 | 0.09 | 0.01 | 78.37 | 75.87 | 67.10 | 0.08 | 0.37 | 0.04 | 33.63 |
| VecMap | 72.00 | 67.17 | 56.42 | 47.43 | 26.62 | 33.39 | 79.91 | 77.77 | 71.45 | 63.53 | 40.62 | 51.86 | 57.35 |
| Ad. | 71.93 | 66.63 | 55.50 | 0.00 | 0.00 | 0.00 | 79.99 | 77.59 | 70.56 | 0.00 | 0.01 | 0.01 | 35.19 |
| BLOOM _{7B} | 52.50 | 38.34 | 26.06 | 32.67 | 21.34 | 34.35 | 56.19 | 41.49 | 26.27 | 32.80 | 21.38 | 34.53 | 34.83 |
| Llama _{13B} | 60.58 | 57.80 | <u>64.44</u> | 22.13 | <u>38.56</u> | 32.28 | 61.09 | 58.51 | 65.10 | 22.14 | <u>38.57</u> | 32.29 | 46.12 |
| GPT-3.5 | 68.17 | 63.07 | 74.15 | 65.94 | 71.80 | 65.12 | 70.72 | 66.08 | 76.84 | 69.88 | 74.95 | 68.69 | 69.62 |
| DM-BLI | 72.94 | 68.67 | 58.91 | <u>48.58</u> | 32.42 | <u>37.34</u> | 80.65 | 78.92 | <u>73.45</u> | <u>64.70</u> | <u>47.98</u> | <u>56.45</u> | <u>60.08</u> |

Table 2. Precision@1 and Precision@5 for the BLI task on six rich-resource language pairs, where * represents EN(English). The best score for is shown in **bold**, and the suboptimal score is shown in underlined.

of BLOOM_{7B} and Llama_{13B} also suggests that the generalization of LLMs to low-resource languages remains an open challenge.

5.2 Results in Rich-resource Languages

Table 2 summarizes the main results of the supervised and the unsupervised BLI tasks on rich-resource language pairs.

In supervised tasks, our proposed method achieves significant improvements, with average nearly 1 point higher than the strongest baseline CL-BLI. We achieve the optimal or sub-optimal performance on all the language pairs. Notably, our method achieves a 6.26% improvement over CL-BLI on distant language pairs Japanese to English, demonstrating advantages of multiple subspace alignment on distant language pairs.

In unsupervised tasks, DM-BLI achieves the sub-optimal result on rich-resource language pairs. While it outperforms the previous mapping-based SOTA method VecMap but underperforms GPT-3.5. The outstanding performance of GPT-3.5 verifies the potential of the latest generation of LLMs for developing bilingual lexicons with sufficient training and a large amount of parameters. However, BLOOM_{7B} and Llama_{13B} are still far lagging behind the traditional mapping-based method even on rich-resource language pairs, which verifies that it is difficult to extract lexical information from large language models (Liu et al., 2021).

5.3 Influence of Translation Direction

In this subsection, we examine how the translation direction affects BLI results in unsupervised setup. The language pairs we choose as examples are Japanese (JA), Chinese (ZH), Finish (FI), Indonesian (ID) from and to English (EN), as shown in Table 3.

| Method | EN-JA | | EN-ZH | | EN-FI | | EN-TR | |
|---------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | → | ← | → | ← | → | ← | → | ← |
| MUSE | 0.01 | 0.37 | 0.01 | 0.04 | 0.06 | 0.05 | 30.73 | 54.76 |
| VecMap | 35.63 | 40.62 | 32.62 | 56.45 | 43.08 | <u>71.72</u> | 40.10 | <u>65.29</u> |
| GPT-3.5 | 57.06 | 74.98 | 42.56 | 68.69 | 58.97 | 64.33 | 52.63 | 55.99 |
| DM-BLI | <u>39.43</u> | <u>47.98</u> | <u>34.69</u> | <u>56.45</u> | <u>44.30</u> | 74.10 | <u>41.90</u> | 67.95 |

Table 3. Precision@5 for the bi-direction unsupervised BLI task on four language pairs. The best score is shown in **bold**, the suboptimal score is shown in underlined.

From Table 3, we observe the performance differences in the two directions of the language pair. Specifically, the results from English to other languages significantly lag behind those from other languages to English. A part of the reason is that there are more unique English words than non-English words in the evaluation set (Xu et al., 2018). It also proves that LLMs exhibit unbalanced capacities across languages, performing better at translating into English than translating into non-English (Zhu et al., 2023b).

5.4 Influence of the Number of Subspaces

In this section, we discuss the impact of the number of subspaces on performance of DM-BLI, taking distant language pair JA2EN as an example.

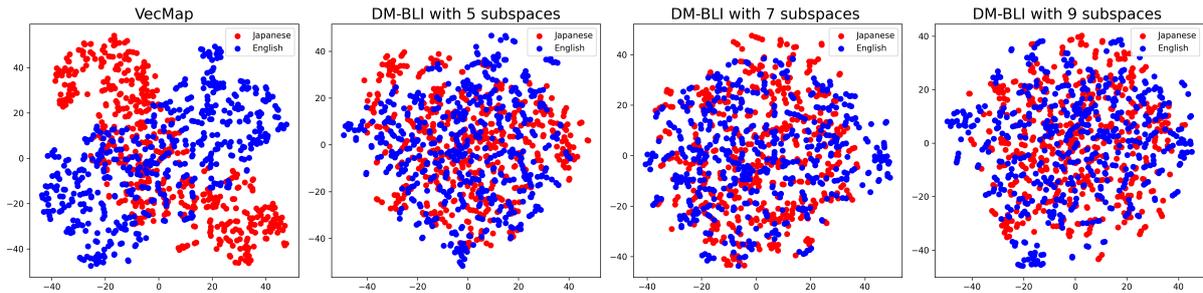


Figure 3. t-SNE visualization of sampled CLWEs derived from VecMap and DM-BLI, where visualization of CLWE derived from DM-BLI is based on different numbers of multiple subspaces.

As shown in Figure 3, compared with VecMap who only use a global mapping, our method lets word with same meaning from different languages get much closer in a shared CLWEs space via multiple subspace-level alignments.

Notably, from Figure 3, we can find that even using different numbers of subspaces, DM-BLI still achieved nearly the same results, which shows that it is not sensitive to the number of subspaces and further proves the robustness of our method.

5.5 Effect of Multiple Subspaces Alignment

Notice that our method focuses on leveraging multiple subspace alignments to achieve better performance for BLI. In this subsection, we discuss the advantages of multiple subspaces alignment from our method DM-BLI, taking low-resource language pair CA2EN as an example.

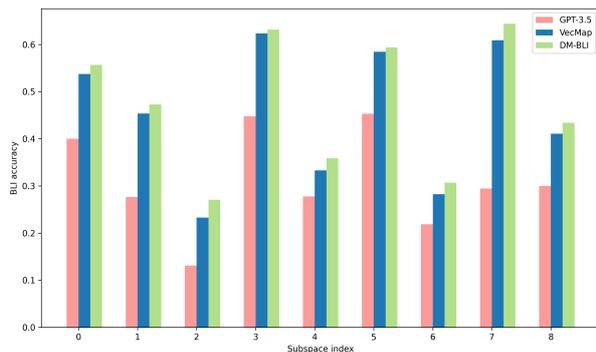


Figure 4. Precision@1 for unsupervised BLI from Catalan to English in different English subspaces.

As shown in Figure 4, on low-resource language pair like CA2EN, we can find that BLI accuracies for all subspaces based on DM-BLI are higher than the strongest mapping-based baseline VecMap. Notably, we also find that unbalanced alignments occur in a generative way via GPT-3.5 as well. Furthermore, LLM’s capability on BLI is still far lagging behind mapping-based approach.

In order to show effect of DM-BLI more intu-

itively, we sample 2 subspaces for visualization. As shown in Figure 5, via multiple subspaces alignment, translation pairs within the subspace stay closer together than applying a global mapping.

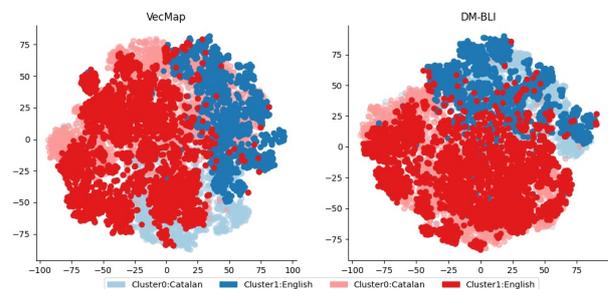


Figure 5. t-SNE visualization of two sampled subspaces in CLWE space derived from VecMap and DM-BLI on CA2EN. Within the subspace, dots denoted by the same color but different transparency are translation pairs.

6 Conclusion

In this paper, we propose a Dynamic Multiple subspaces alignment framework for unsupervised BLI, called DM-BLI. Our method utilizes multiple subspaces alignment instead of a single mapping alignment to achieve more accurate alignment on the subspace level. The experiments show that our method can significantly improve the bilingual word induction performance compared with strong baselines even including GPT-3.5, especially for distant and low-resource language pairs. At the same time, the unsatisfied performances of BLOOM_{7B} and Llama_{13B} on all language pairs also suggest that it is difficult to extract lexical information from large language models and the generalization of LLMs to low-resource languages remains an open challenge. In the future, we will consider combining our method with multilingual LLMs to take advantage of these two paradigms.

545 **Limitations**

546 First, due to our limited computing resources, we
547 did not conduct a comprehensive evaluation of the
548 BLI capabilities of multilingual LLMs. For open-
549 source LLMs, LLMs exceeding 13B parameters
550 were not evaluated in the experiment. For close-
551 source LLMs, experiments were mainly conducted
552 on GPT-3.5-turbo which is not the latest and best.

553 Second, public BLI datasets are not enough to
554 support a comprehensive evaluation. In the evalua-
555 tion standard dictionary, the proportion of ground-
556 truth translations in different categories is uneven.
557 As also discussed in (Li et al., 2023), current evalua-
558 tion will not work for words that are not included
559 in the gold translations.

560 **Acknowledgements**

561 **References**

562 David Alvarez-Melis and T. Jaakkola. 2018. Gromov-
563 wasserstein alignment of word embedding spaces. In
564 *Conference on Empirical Methods in Natural Lan-
565 guage Processing*.

566 Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018.
567 *A robust self-learning method for fully unsupervised*
568 *cross-lingual mappings of word embeddings*. In *Pro-
569 ceedings of the 56th Annual Meeting of the Associa-
570 tion for Computational Linguistics (Volume 1: Long*
571 *Papers)*, pages 789–798, Melbourne, Australia. As-
572 sociation for Computational Linguistics.

573 Agirre E Artetxe M, Labaka G. 2017. Learning bilin-
574 gual word embeddings with (almost) no bilingual
575 data. *Proceedings of the 55th Annual Meeting of the*
576 *Association for Computational Linguistics (Volume*
577 *1: Long Papers)*. Vancouver, Canada.

578 Piotr Bojanowski, Edouard Grave, Armand Joulin, and
579 Tomas Mikolov. 2016. Enriching word vectors with
580 subword information. *Transactions of the Associa-
581 tion for Computational Linguistics*, 5:135–146.

582 Eleftheria Briakou, Colin Cherry, and George Foster.
583 2023. Searching for needles in a haystack: On the
584 role of incidental bilingualism in palm’s translation
585 capability. *arXiv preprint arXiv:2305.10266*.

586 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
587 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
588 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
589 Askell, et al. 2020. Language models are few-shot
590 learners. *Advances in neural information processing*
591 *systems*, 33:1877–1901.

592 Alexis Conneau, Guillaume Lample, Marc’Aurelio Ran-
593 zato, Ludovic Denoyer, and Hervé Jégou. 2017.
594 Word translation without parallel data. *arXiv preprint*
595 *arXiv:1710.04087*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
Kristina Toutanova. 2019. *BERT: Pre-training of*
deep bidirectional transformers for language under-
standing. In *Proceedings of the 2019 Conference of*
the North American Chapter of the Association for
Computational Linguistics: Human Language Tech-
nologies, Volume 1 (Long and Short Papers), pages
4171–4186, Minneapolis, Minnesota. Association for
Computational Linguistics.

Tobias Eder, Viktor Hangya, and Alexander Fraser.
2021. *Anchor-based bilingual word embeddings for*
low-resource languages. In *Proceedings of the 59th*
Annual Meeting of the Association for Computational
Linguistics and the 11th International Joint Confer-
ence on Natural Language Processing (Volume 2:
Short Papers), pages 227–232, Online. Association
for Computational Linguistics.

Dyer C Faruqui M. 2014. Improving vector space word
representations using multilingual correlation. *Pro-*
ceedings of the 14th Conference of the European
Chapter of the Association for Computational Lin-
guistics. Gothenburg, Sweden: Association for Com-
putational Linguistics.

Marjan Ghazvininejad, Hila Gonen, and Luke Zettle-
moyer. 2023. Dictionary-based phrase-level prompt-
ing of large language models for machine translation.
arXiv preprint arXiv:2302.07856.

Goran Glavaš, Robert Litschko, Sebastian Ruder, and
Ivan Vulić. 2019. *How to (properly) evaluate cross-*
lingual word embeddings: On strong baselines, compar-
ative analyses, and some misconceptions. In *Pro-*
ceedings of the 57th Annual Meeting of the Associa-
tion for Computational Linguistics, pages 710–721,
Florence, Italy. Association for Computational Lin-
guistics.

Goran Glavaš and Ivan Vulić. 2020. *Non-linear*
instance-based cross-lingual mapping for non-
isomorphic embedding spaces. In *Proceedings of*
the 58th Annual Meeting of the Association for Com-
putational Linguistics, pages 7548–7555, Online. As-
sociation for Computational Linguistics.

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, Online. Association for Computational
Linguistics.

Yuehui Han, Le Hui, Haobo Jiang, Jianjun Qian, and
Jin Xie. 2022. Generative subgraph contrast for self-
supervised graph representation learning. In *Euro-*
pean Conference on Computer Vision, pages 91–107.
Springer.

Guillaume Lample and Alexis Conneau. 2019. Cross-
lingual language model pretraining. *arXiv preprint*
arXiv:1901.07291.

| | | | |
|-----|--|---|-----|
| 652 | Yaoyiran Li, Anna Korhonen, and Ivan Vulić. 2023. | Ndapa Nakashole. 2018. NORMA: Neighborhood sensitive maps for multilingual word embeddings . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 512–522, Brussels, Belgium. Association for Computational Linguistics. | 709 |
| 653 | On bilingual lexicon induction with large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9577–9599, Singapore. Association for Computational Linguistics. | | 710 |
| 654 | | | 711 |
| 655 | | | 712 |
| 656 | | | 713 |
| 657 | | | 714 |
| 658 | Yaoyiran Li, Fangyu Liu, Nigel Collier, Anna Korhonen, and Ivan Vulić. 2022. Improving word translation via two-stage contrastive learning . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4353–4374, Dublin, Ireland. Association for Computational Linguistics. | Aitor Ormazabal, Mikel Artetxe, Gorka Labaka, Aitor Soroa, and Eneko Agirre. 2019. Analyzing the limitations of cross-lingual word embedding mappings . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4990–4995, Florence, Italy. Association for Computational Linguistics. | 715 |
| 659 | | | 716 |
| 660 | | | 717 |
| 661 | | | 718 |
| 662 | | | 719 |
| 663 | | | 720 |
| 664 | | | 721 |
| 665 | Fangyu Liu, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2021. Fast, effective, and self-supervised: Transforming masked language models into universal lexical and sentence encoders . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 1442–1459, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. | Barun Patra, Joel Ruben Antony Moniz, Sarthak Garg, Matthew R. Gormley, and Graham Neubig. 2019. Bilingual lexicon induction with semi-supervision in non-isometric embedding spaces . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 184–193, Florence, Italy. Association for Computational Linguistics. | 722 |
| 666 | | | 723 |
| 667 | | | 724 |
| 668 | | | 725 |
| 669 | | | 726 |
| 670 | | | 727 |
| 671 | | | 728 |
| 672 | | | |
| 673 | James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations. In <i>Proceedings of the fifth Berkeley symposium on mathematical statistics and probability</i> , volume 1, pages 281–297. Oakland, CA, USA. | Shuo Ren, Shujie Liu, Ming Zhou, and Shuai Ma. 2020. A graph-based coarse-to-fine method for unsupervised bilingual lexicon induction . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3476–3485, Online. Association for Computational Linguistics. | 729 |
| 674 | | | 730 |
| 675 | | | 731 |
| 676 | | | 732 |
| 677 | | | 733 |
| 678 | Kelly Marchisio, Ali Saad-Eldin, Kevin Duh, Carey Priebe, and Philipp Koehn. 2022. Bilingual lexicon induction for low-resource languages using graph matching via optimal transport . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2545–2561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | Sebastian Ruder, Ivan Vulić, and Anders Søgaard. 2019. A survey of cross-lingual word embedding models. <i>Journal of Artificial Intelligence Research</i> , 65:569–631. | 735 |
| 679 | | | 736 |
| 680 | | | 737 |
| 681 | | | 738 |
| 682 | | | |
| 683 | | | 739 |
| 684 | | | 740 |
| 685 | | | 741 |
| 686 | Tomas Mikolov, Quoc V Le, and Ilya Sutskever. 2013a. Exploiting similarities among languages for machine translation. <i>arXiv preprint arXiv:1309.4168</i> . | Saqib Sarfraz, Vivek Sharma, and Rainer Stiefelhagen. 2019. Efficient parameter-free clustering using first neighbor relations. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 8934–8943. | 742 |
| 687 | | | 743 |
| 688 | | | |
| 689 | Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. <i>Advances in neural information processing systems</i> , 26. | Anders Søgaard, Sebastian Ruder, and Ivan Vulić. 2018. On the limitations of unsupervised bilingual dictionary induction . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 778–788, Melbourne, Australia. Association for Computational Linguistics. | 744 |
| 690 | | | 745 |
| 691 | | | 746 |
| 692 | | | 747 |
| 693 | | | 748 |
| 694 | Tasnim Mohiuddin, M Saiful Bari, and Shafiq Joty. 2020. LNMMap: Departures from isomorphic assumption in bilingual lexicon induction through non-linear mapping in latent space . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2712–2723, Online. Association for Computational Linguistics. | Hagai Taitelbaum, Gal Chechik, and Jacob Goldberger. 2019. Multilingual word translation using auxiliary languages. In <i>Proc. 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1330–1335, Hong Kong, China. | 751 |
| 695 | | | 752 |
| 696 | | | 753 |
| 697 | | | 754 |
| 698 | | | 755 |
| 699 | | | 756 |
| 700 | | | 757 |
| 701 | Tasnim Mohiuddin and Shafiq Joty. 2019. Revisiting adversarial autoencoder for unsupervised word translation with cycle consistency and improved training. In <i>Proc. 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3857–3867, Minneapolis, Minnesota. | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> . | 758 |
| 702 | | | 759 |
| 703 | | | 760 |
| 704 | | | 761 |
| 705 | | | 762 |
| 706 | | | 763 |
| 707 | | | |
| 708 | | | |

764 Ivan Vulić, Goran Glavaš, Roi Reichart, and Anna Korhonen. 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 4407–4418, Hong Kong, China. Association for Computational Linguistics.

772 Ivan Vulić, Edoardo Maria Ponti, Robert Litschko, Goran Glavaš, and Anna Korhonen. 2020. [Probing pretrained language models for lexical semantics.](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7222–7240, Online. Association for Computational Linguistics.

779 Moens M F Vulić I. 2015. Bilingual word embeddings from non-parallel document-aligned data applied to bi-lingual lexicon induction. *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Beijing, China.

784 Haozhou Wang, James Henderson, and Paola Merlo. 2020. Multi-adversarial learning for cross-lingual word embeddings. *arXiv preprint arXiv:2010.08432*.

787 BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Lucic, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.

793 Ruochen Xu, Yiming Yang, Naoki Otani, and Yuexin Wu. 2018. [Unsupervised cross-lingual transfer of word embedding spaces.](#) In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2465–2474, Brussels, Belgium. Association for Computational Linguistics.

799 Meng Zhang, Yang Liu, Huanbo Luan, and Maosong Sun. 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, Vancouver, Canada. Association for Computational Linguistics.

806 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023a. Multilingual machine translation with large language models: Empirical results and analysis. *arXiv preprint arXiv:2304.04675*.

811 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023b. Multilingual machine translation with large language models: Empirical results and analysis. *arXiv preprint arXiv:2304.04675*.

816 Cer D et al Zou W Y, Socher R. 2013. Bilingual word embeddings for phrase-based machine translation. *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Seattle,

Washington, USA: Association for Computational Linguistics.

820
821