

Bilingual Lexicon Induction for Low-Resource Languages using Graph Matching via Optimal Transport

Anonymous ACL submission

Abstract

Bilingual lexicons form a critical component of various NLP applications, including unsupervised and semisupervised machine translation and crosslingual information retrieval. In this work, we improve bilingual lexicon induction performance across 32 diverse language pairs with a graph-matching method based on optimal transport. The method is especially strong with very low amounts of supervision.

1 Introduction

Bilingual lexicon induction (BLI) from word embedding spaces is a popular task with a large body of existing literature (e.g. Artetxe et al., 2018; Conneau et al., 2018; Patra et al., 2019; Shi et al., 2021; Zhao et al., 2020b). The goal is to extract a dictionary of translation pairs given separate language-specific embedding spaces, which can then be used to bootstrap downstream tasks such as cross-lingual information retrieval and unsupervised/semi-supervised machine translation.

A great challenge across NLP is maintaining performance in low-resource scenarios. A common criticism of the BLI and low-resource MT literature is that while claims are made about diverse and under-resourced languages, research is often performed on down-sampled corpora of high-resource, highly-related languages on similar domains (Artetxe et al., 2020). Such corpora are not good proxies for true low-resource languages owing to data challenges such as dissimilar scripts, domain shift, noise, and lack of sufficient bitext (Marchisio et al., 2020). These differences can lead to dissimilarity between the embedding spaces (decreasing isometry), causing BLI to fail (Søgaard et al., 2018; Nakashole and Flaiger, 2018; Ormazabal et al., 2019; Glavaš et al., 2019; Vulić et al., 2019; Patra et al., 2019; Marchisio et al., 2020).

There are two axes by which a language dataset is considered “low-resource”. First, the language itself may be a low-resource language: one for which

little bitext and/or monolingual text exists. Even for high-resource languages, the long tail of words may have poorly trained word embeddings due to rarity in the dataset (Gong et al., 2018; Czarnowska et al., 2019). In the data-poor setting of true low-resource languages, a great majority of words have little representation in the corpus, resulting in poorly-trained embeddings for a large proportion of them. The second axis is low-supervision. Here, there are few ground-truth examples from which to learn. For BLI from word embedding spaces, low-supervision means there are few seeds from which to induce a relationship between spaces, regardless of the quality of the spaces themselves.

We bring a new algorithm for graph-matching based on optimal transport (OT) to the NLP and BLI literature. We evaluate using 32 diverse language pairs under varying amounts of supervision. The method works strikingly well across language pairs, especially in low-supervision contexts. As low-supervision on low-resource languages reflects the real-world use case for BLI, this is an encouraging development on realistic scenarios.

2 Background

The typical baseline approach for BLI from word embedding spaces assumes that spaces can be mapped via linear transformation. Such methods typically involve solutions to the Procrustes problem (see Gower et al. (2004) for a review). Alternatively, a graph-based view considers words as nodes in undirected weighted graphs, where edges are the distance between words. Methods taking this view do not assume a linear mapping of the spaces exists, allowing for more flexible matching.

BLI from word embedding spaces Assume separately-trained monolingual word embedding spaces: $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{Y} \in \mathbb{R}^{m \times d}$ where n/m are the source/target language vocabulary sizes and d is the embedding dimension. We build the matrices $\bar{\mathbf{X}}$

and \bar{Y} of seeds from X and Y , respectively, such that given s seed pairs $(x_1, y_1), (x_2, y_2), \dots, (x_s, y_s)$, the first row of \bar{X} is x_1 , the second row is x_2 , etc. We build \bar{Y} analogously for the y -component of each seed pair. The goal is to recover matches for the $X \setminus \bar{X}$ and/or $Y \setminus \bar{Y}$ non-seed words.

Procrustes Many BLI methods use solutions to the Procrustes problem (e.g. Artetxe et al., 2019b; Conneau et al., 2018; Patra et al., 2019). These compute the optimal transform W to map seeds:

$$\min_{W \in \mathbb{R}^{d \times d}} \|\bar{X}W - \bar{Y}\|_F^2 \quad (1)$$

Once solved for W , then XW and Y live in a shared space and translation pairs can be extracted via nearest-neighbor search. Constrained to the space of orthogonal matrices, Eq. 1 has a simple closed-form solution (Schönemann, 1966):

$$W = VU^T \quad U\Sigma V = \text{SVD}(\bar{Y}^T \bar{X})$$

Graph View Here, words are nodes in monolingual graphs $G_x, G_y \in \mathbb{R}^{n \times n}$, and cells in G_x, G_y are edge weights representing distance between words. We use cosine similarity, which is common in NLP. The objective function is Eq. 2, where Π is the set of *permutation matrices*.¹ Intuitively, PG_yP^T finds the optimal relabeling of G_y to align with G_x . This “minimizes edge-disagreements” between G_x and G_y . This graph-matching objective is NP-Hard. Eq. 3 is equivalent.

$$\min_{P \in \Pi} \|G_x - PG_yP^T\|_F^2 \quad (2)$$

$$\max_{P \in \Pi} \text{trace}(G_x^T PG_yP^T) \quad (3)$$

Ex. Take source words x_1, x_2 . We wish to recover valid translations y_{x_1}, y_{x_2} . If $\text{distance}(x_1, x_2) = \text{distance}(y_{x_1}, y_{x_2})$, a solution P can have an edge-disagreement of 0 here. We then extract y_{x_1}, y_{x_2} as translations of x_1, x_2 . In reality, though, it is unlikely that $\text{distance}(x_1, x_2) = \text{distance}(y_{x_1}, y_{x_2})$. Because Eq. 2 finds the ideal P to minimize edge disagreements over the entire graphs, we hope that nodes paired by P are valid translations. If G_x and G_y are isomorphic and there is a unique solution, then P correctly recovers all translations.

Graph-matching is an active research field and is computationally prohibitive on large graphs, but approximation algorithms exist. BLI involves matching large, non-isomorphic graphs—among the greatest challenges for graph-matching.

¹A permutation matrix represents a one-to-one mapping: There is a single 1 in each row and column, and 0 elsewhere.

2.1 FAQ Algorithm for Graph Matching

Vogelstein et al. (2015)’s Fast Approximate Quadratic Assignment Problem algorithm (FAQ) uses gradient ascent to approximate a solution to Eq. 2. Motivated by “connectonomics” in neuroscience (the study of brain graphs with biological [groups of] neurons as nodes and neuronal connections as edges), FAQ was designed to perform accurately and efficiently on large graphs.

FAQ relaxes the search space of Eq. 3 to allow any doubly-stochastic matrix (the set \mathcal{D}). Each cell in a doubly-stochastic matrix is a non-negative real number and each row/column sums to 1. The set \mathcal{D} thus contains Π but is much larger. Relaxing the search space makes it easier to optimize Eq. 3 via gradient ascent/descent.² FAQ solves the objective with the Frank-Wolfe method (Frank et al., 1956) then projects back to a permutation matrix.

Algorithm 1 is FAQ; $f(P) = \text{tr}(G_x^T PG_yP^T)$. Step 2 finds a permutation matrix approximation Q^{i} to P^{i} in the direction of the gradient. Finding such a P requires approximation when P is high-dimensional. Here, it is solved via the **Hungarian Algorithm** (Kuhn, 1955; Jonker and Volgenant, 1987), whose solution is a permutation matrix. Finally, P^n is projected back onto to the space of permutation matrices. Seeded Graph Matching (SGM; Fishkind et al., 2019) is a variant of FAQ allowing for supervision, and was recently shown to be effective for BLI by Marchisio et al. (2021).

Algorithm 1 FAQ Algorithm for Graph Matching

Let: $G_x, G_y \in \mathbb{R}^{n \times n}$, $P^{0} \in \mathcal{D}$ (dbl-stoch.)

while stopping criterion not met **do**

1. Calculate $\nabla f(P^{i})$:

$$\nabla f(P^{i}) = G_x P^{i} G_y^T + G_x^T P^{i} G_y$$

2. Q^{i} = permutation matrix approx. to $\nabla f(P^{i})$ via Hungarian Algorithm

3. Calculate step size:

$$\arg \max_{\alpha \in [0,1]} f(\alpha P^{i} + (1 - \alpha)Q^{i})$$

4. Update $P^{i+1} := \alpha P^{i} + (1 - \alpha)Q^{i}$

end while

return permutation matrix approx. to P^{n} via Hung. Alg.

Strengths/Weaknesses FAQ/SGM perform well solving the exact graph-matching problem: where graphs are isomorphic and a full matching exists. In reality, however, large graphs are rarely isomorphic. For BLI, languages have differing vocabulary size, synonyms/antonyms, and idiosyncratic

²“descent” for the Quadratic Assignment Problem, “ascent” for the Graph Matching Problem. The optimization objectives are equivalent: See Vogelstein et al. (2015) for a proof.

concepts; it is more natural to assume that an exact matching between word spaces does *not* exist, and that multiple matchings may be equally valid. This is an inexact graph-matching problem. FAQ generally performs poorly finding non-seeded inexact matchings (Saad-Eldin et al., 2021).

2.2 GOAT

Graph Matching via Optimal Transport (GOAT) (Saad-Eldin et al., 2021) is a new graph-matching method which uses advances in OT. Similar to SGM, GOAT amends FAQ and can use seeds. GOAT has been successful for the inexact graph-matching problem on non-isomorphic graphs: whereas FAQ rapidly fails on non-isomorphic graphs, GOAT maintains strong performance.

Optimal Transport OT is an optimization problem concerned with the most efficient way to transfer probability mass from distribution μ to distribution ν . Formally, discrete³ OT is the minimization of the inner product of a transportation “plan” matrix P with a cost matrix C , as in Eq. 4. $\langle \cdot, \cdot \rangle$ is the Frobenius inner product.

$$P^* = \arg \min_{P \in \mathcal{U}(r,c)} \langle P, C \rangle \quad (4)$$

P is an element of the “transportation polytope” $\mathcal{U}(r, c)$ —the set of matrices whose rows sum to r and columns sum to c . The Hungarian Algorithm approximately solves OT, but the search space is restricted to permutation matrices.

Sinkhorn: Lightspeed OT Cuturi (2013) introduce Sinkhorn distance, an approximation of OT distance that can be solved quickly and accurately by adding an entropy penalty h to Eq. 4. Adding h makes the objective easier and more efficient to compute, and encourages “intermediary” solutions similar to that seen in the **Intuition** subsection.

$$P^\lambda = \arg \min_{P \in \mathcal{U}(r,c)} \langle P, C \rangle - \frac{1}{\lambda} h(P) \quad (5)$$

As $\lambda \rightarrow \infty$, P^λ approaches the ideal transportation matrix P^* . Cuturi (2013) show that Eq. 5 can be computed using Sinkhorn’s algorithm (Sinkhorn, 1967). The interested reader can see details of the algorithm in Cuturi (2013); Peyre and Cuturi (2019). Unlike the Hungarian Algorithm, Sinkhorn has no restriction to a permutation matrix solution and can be solved over any $\mathcal{U}(r, c)$.

Sinkhorn in GOAT GOAT uses Cuturi (2013)’s algorithm to solve Eq. 5 over $\mathcal{U}(1, 1)$, the set

³As ours is, as we compute over matrices.

of doubly-stochastic matrices \mathcal{D} . They call this the “doubly stochastic OT problem”, and the algorithm that solves it “Lightspeed Optimal Transport” (LOT). Although Sinkhorn distance was created for efficiency, Saad-Eldin et al. (2021) find that using the matrix P^λ that minimizes Sinkhorn distance also improves matching performance on large and non-isometric graphs. Algorithm 2 is GOAT.

Algorithm 2 GOAT

Let: $\mathbf{G}_x, \mathbf{G}_y \in \mathbb{R}^{n \times n}$, $P^{\{0\}} \in \mathcal{D}$ (dbl-stoch.)

while stopping criterion not met **do**

1. Calculate $\nabla f(P^{\{i\}})$:

$$\nabla f(P^{\{i\}}) = G_x P^{\{i\}} G_y^T + G_x^T P^{\{i\}} G_y$$

2. $Q^{\{i\}}$ = dbl-stoch. approx. to $\nabla f(P^{\{i\}})$ via LOT.

3. Calculate step size:

$$\arg \max_{\alpha \in [0,1]} f(\alpha P^{\{i\}} + (1-\alpha)Q^{\{i\}})$$

4. Update $P^{\{i+1\}} := \alpha P^{\{i\}} + (1-\alpha)Q^{\{i\}}$

end while

return permutation matrix approx. to $P^{\{n\}}$ via Hung. Alg.

Intuition The critical difference between SGM/FAQ and GOAT is how each calculates step direction based on the gradient. Under the hood, each algorithm maximizes $\text{trace}(Q^T \nabla f(P^{\{i\}}))$ to compute $Q^{\{i\}}$ (the step direction) in Step 2 of their respective algorithms. See Saad-Eldin et al. (2021) or Fishkind et al. (2019) for a derivation. FAQ uses the Hungarian Algorithm and GOAT uses LOT.

For $\nabla f(P^{\{i\}})$ below, there are *two* valid permutation matrices Q_1 and Q_2 that maximize the trace. When multiple solutions exist, the Hungarian Algorithm chooses one arbitrarily. Thus, updates of P in FAQ are constrained to be permutation matrices.

$$\nabla f(P^{\{i\}}) = \begin{pmatrix} 0 & 3 & 0 \\ 2 & 1 & 2 \\ 0 & 0 & 0 \end{pmatrix}$$

$$Q_1 = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}, \quad Q_2 = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$\text{trace}(Q_1^T \nabla f(P^{\{i\}})) = \text{trace}(Q_2^T \nabla f(P^{\{i\}})) = 5$$

Concerningly, Saad-Eldin et al. (2021) find that seed order influences the solution in a popular implementation of the Hungarian Algorithm. Since BLI is a high-dimensional many-to-many task, arbitrary choices could meaningfully affect the result.

GOAT, on the other hand, can step in the direction of a doubly-stochastic matrix. Saad-Eldin et al. (2021) prove that given multiple permutation matrices that equally approximate the gradient at $P^{\{i\}}$,

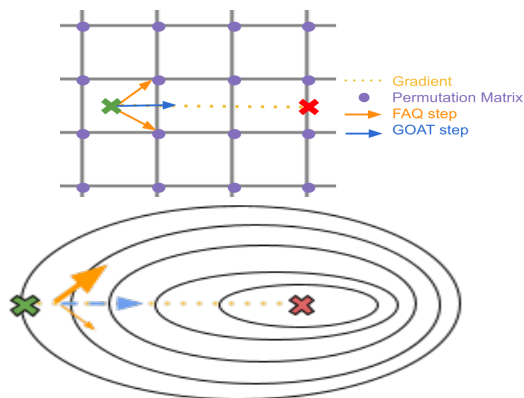


Figure 1: Optimization step of FAQ vs. GOAT. FAQ arbitrarily chooses the direction of a permutation matrix. GOAT averages perm. matrices to take a smoother path.

any convex linear combination is a doubly stochastic matrix that equally approximates the gradient:

$$P_\lambda = \sum_i^n \lambda_i P_i \quad \text{s.t. } \lambda_1 + \dots + \lambda_n = 1; \lambda_i \in [0, 1]$$

P_λ is a weighted combination of *many* valid solutions—obviating the need to arbitrarily select one for the gradient-based update. LOT’s output of a doubly-stochastic matrix in Step 2 is similar to finding a P_λ in that it needn’t discretize to a single permutation matrix. In this way, GOAT can be thought of as taking a step that incorporates many possible permutation solutions. For instance, GOAT may select $Q_{ds} = \frac{1}{2}Q_1 + \frac{1}{2}Q_2$, which also maximizes $\text{trace}(Q_{ds}^T \nabla f(P^{\{i\}}))$.

$$Q_{ds} = \begin{pmatrix} 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \end{pmatrix}$$

$$\text{trace}(Q_{ds}^T \nabla f(P^{\{i\}})) = 5$$

Thus whereas FAQ takes non-deterministic “choppy” update steps, GOAT optimizes smoothly and deterministically. Figure 1 is an illustration.

3 Experimental Setup

We run Procrustes, SGM, and GOAT on 32 language pairs. We also run system combination experiments similar to Marchisio et al. (2021). We evaluate with the standard precision@1 (P@1).

We induce lexicons using (1) the closed-form solution to the orthogonal Procrustes problem of Eq. 1, extracting nearest neighbors using CSLS (Conneau et al., 2018), (2) SGM, solving the seeded version of Eq. 2, and (3) GOAT. Word graphs are $\mathbf{G}_x = \mathbf{X}\mathbf{X}^T$, $\mathbf{G}_y = \mathbf{Y}\mathbf{Y}^T$.

System Combination We perform system combination experiments analogous to those of Marchi-

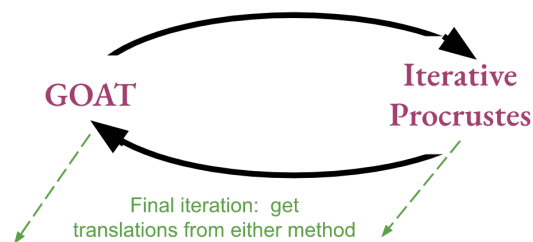


Figure 2: Combo system: Iter.Proc. & GOAT. (1) run GOAT in fwd/rev directions (2) intersect hypotheses, pass to IterProc, (3) run IterProc fwd/rev (4) intersect hypotheses, pass to Step (1). Repeat N cycles. (End) Get final translations from fwd GOAT or IterProc.

sio et al. (2021), incorporating GOAT. Figure 2 shows the system, which is made of two components: GOAT run in forward and reverse directions, and “Iterative Procrustes with Stochastic-Add” from Marchisio et al. (2021). This iterative version of Procrustes runs Procrustes in source→target and target→source directions and feeds H random hypotheses from the intersection of both directions into another run of Procrustes with the gold seeds. The process repeats for I iterations, adding H more random hypotheses each time until all are chosen. We set $H = 100$ and $I = 5$, as in the original work.

3.1 Data & Software

We use publicly-available fastText word embeddings (Bojanowski et al., 2017)⁴ which we normalize, mean-center, and renormalize (Artetxe et al., 2018; Zhang et al., 2019) and bilingual dictionaries from MUSE⁵ filtered to be one-to-one.⁶ For languages with 200,000+ embeddings, we use the first 200,000. Dictionary and embeddings space sizes are in Appendix Table A1. Each language pair has ~4100-4900 translation pairs post-filtering. We choose 0-4000 pairs in frequency order as seeds for experiments, leaving the rest as the test set.⁷ For SGM and GOAT, we use the publicly-available implementations from the GOAT repository⁸ with default hyperparameters (barycenter initialization). We set reg=500 for GOAT. For system combination experiments, we amend the code from Marchisio et al. (2021)⁹ to incorporate GOAT.

⁴<https://fasttext.cc/docs/en/pretrained-vectors.html>

⁵<https://github.com/facebookresearch/MUSE>

⁶For each source word, keep the first unused target word. Targets are in arbitrary order, so this is random sampling.

⁷Ex. En-De with 100 seeds has 4803 test items. With 1000 seeds, the test set contains 3903 items.

⁸<https://github.com/neurodata/goat>. Some exps. used SGM from Graspologic (github.com/microsoft/graspologic; Chung et al., 2019), but they are mathematically equal.

⁹<https://github.com/kellymarchisio/euc-v-graph-bli>

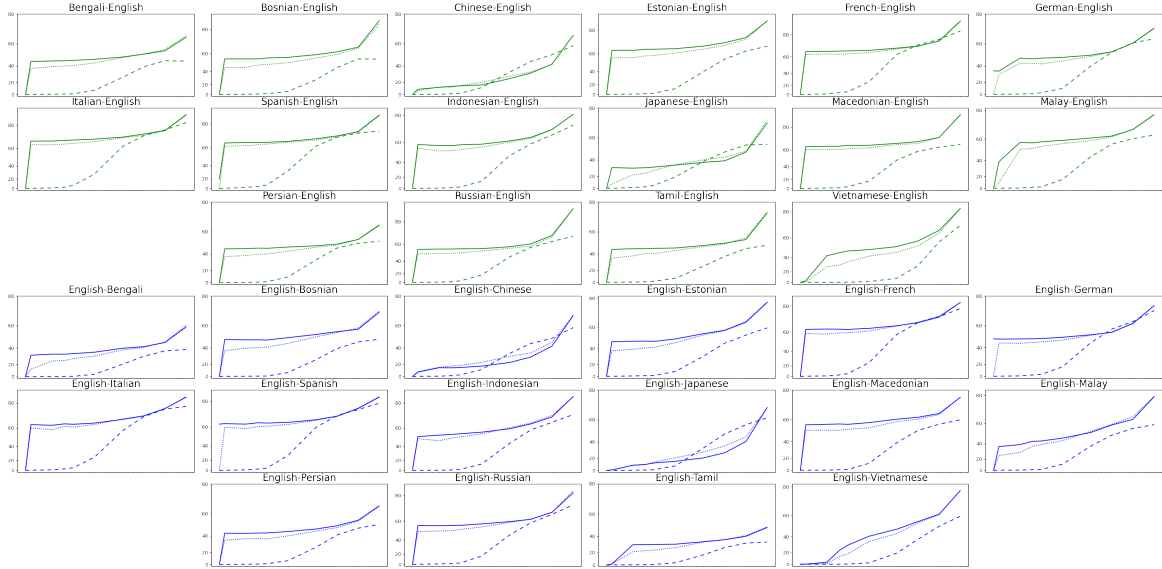


Figure 3: Procrustes (Proc) (dashed) vs. SGM (dotted) vs. GOAT (solid), visualized. X-axis: # of seeds (log scale). Y-axis: Precision@1 (\uparrow is better). Procrustes grows slowly, needing many seeds. GOAT is typically best.

4 Results

Results of Procrustes vs. SGM vs. GOAT are in Table 2, visualized in Figure 3.

Procrustes vs. SGM Marchisio et al. (2021) conclude that SGM strongly outperforms Procrustes for English \rightarrow German and Russian \rightarrow English with 100+ seeds. We find that the trend holds across diverse language pairs, with the effect even stronger with less supervision. SGM performs reasonably with only 50 seeds for nearly all language, and with only 25 seeds in many. Chinese \leftrightarrow English and Japanese \leftrightarrow English perform relatively worse, and highly-related languages perform best: French, Spanish, and Italian. German \leftrightarrow English performance is low relative to some less-related languages, which have surprisingly strong performance from SGM: Indonesian \leftrightarrow English and Macedonian \leftrightarrow English score $P@1 \approx 50$ -60, even with low supervision. Procrustes does not perform above ~ 10 for any language pair with ≤ 100 seeds, whereas SGM exceeds $P@1 = 10$ with only 25 seeds for 25 of 32 pairs.

SGM vs. GOAT GOAT improves considerably over SGM for nearly all language pairs, and the effect is particularly strong with very low amounts of seeds and less-related languages. GOAT improves upon SGM by +19.0, +8.5, and +7.9 on English \rightarrow Bengali with 25, 50, and 75 seeds, respectively. As the major use case of low-resource BLI and MT is dissimilar languages with low supervision, this is an encouraging result for real-world applications. It generally takes 200+ seeds

	EVS	GH		EVS	GH
bn	37.79	0.49	it	22.42	0.20
bs	35.93	0.41	ja	894.20	0.55
de	11.49	0.31	mk	151.02	0.19
es	9.91	0.21	ms	153.42	0.49
et	35.22	0.68	ru	14.19	0.46
fa	86.98	0.39	ta	56.66	0.26
fr	27.92	0.17	vi	256.28	0.42
id	188.98	0.39	zh	519.82	0.61

Table 1: Degree of isomorphism of embedding spaces in relation to English. EVS = Eigenvector Similarity. GH = Gromov-Hausdorff Distance. \downarrow : more isomorphic.

for SGM to achieve similar scores to GOAT with just 25 seeds. For some highly-related languages, GOAT performs well even with no seeds (unsupervised), where both SGM and Procrustes fail. GOAT scores 48.8 on English \rightarrow German, 34.5 on German \rightarrow English, 62.4 on English \rightarrow Spanish, and 19.6 on Spanish \rightarrow English with no supervision.

4.1 Isomorphism of Embedding Spaces

Eigenvector similarity (EVS; Søgaard et al., 2018) measures isomorphism of embedding spaces based on the difference of Laplacian eigenvalues. Gromov-Hausdorff distance (GH) measures distance based on nearest neighbors after an optimal orthogonal transformation (Patra et al., 2019). EVS and GH are symmetric, and lower means more isometric spaces. Refer to the original papers for mathematical descriptions. We compute the metrics over the word embedding using scripts from Vulić et al. (2020)¹⁰ and show results in Table 1. We observe a moderate correlation between EVS and GH (Spearman’s $\rho = 0.434$, Pearson’s $r = 0.44$).

¹⁰<https://github.com/cambridgeltl/iso-study/scripts>

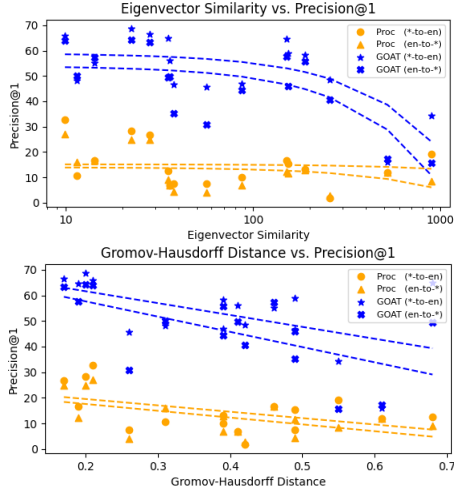


Figure 4: X-axis: Eigenvector Similarity (EVS) / Gromov-Hausdorff (GH) Distance of language compared to English. Y-axis: Precision@1 from Procrustes & GOAT with 200 seeds. \downarrow EVS/GH = \uparrow isomorphic.

Figure 4 shows the relationship between relative isomorphism of each language vs. English, and performance of Procrustes/GOAT at 200 seeds. Trends indicate that higher isomorphism varies with higher precision from Procrustes and GOAT. GH shows a moderate to strong negative Pearson’s correlation with performance from Procrustes and GOAT: $r = -0.47$ and $r = -0.53$, respectively, for *-to-en and -0.55 and -0.61 for en-to-*. EVS correlates weakly negatively with performance from Procrustes (*-to-en: -0.06, en-to-*: -0.28) and strongly negatively with GOAT (*-to-en: -0.67, en-to-*: -0.75). As higher GH/EVS indicates less isomorphism, negative correlations imply that lower degrees of isomorphism correlate with lower scores from Procrustes/GOAT.

4.2 System Combination

System combination results are in Table 3. Similar to Marchisio et al. (2021)’s findings for their combined Procrustes/SGM system, we find (1) our combined Procrustes/GOAT system outperforms Procrustes and GOAT alone, (2) ending with the Iterative Procrustes is best for moderate amounts of seeds, (3) ending with GOAT is best for very low or very high number of seeds.

Whether we end with Iterative Procrustes vs. GOAT is critically important for the lowest seed sizes: -EndGOAT (-EG) usually fails with 25 seeds; all language pairs except German \leftrightarrow English and Russian \leftrightarrow English score $P@1 < 15.0$, and most score $P@1 < 2.0$. Simply switching the order of processing in the combination system, however, boosts performance dramatically: ex. from 0.6

for StartProc-EndGOAT to 61.5 for StartGOAT-EndProc for Bosnian \rightarrow English with 25 seeds.

There are some language pairs such as English \rightarrow Persian and Russian \leftrightarrow English where a previous experiment with no seeds had reasonable performance, but the combined system failed. It is worth investigating where this discrepancy arises.

5 Discussion

We have seen GOAT’s strength in low-resource scenarios and in non-isomorphic embedding spaces. As the major use case of low-resource BLI and MT is dissimilar languages with low supervision, GOAT’s strong performance is an encouraging result for real-world applications. Furthermore, GOAT outperforms SGM. As the graph-matching objective is NP-hard so all algorithms are approximate, GOAT does a better job by making a better calculation of step direction. Chinese \leftrightarrow English and Japanese \leftrightarrow English are outliers, which is worthy of future investigation. Notably, these languages have very poor isomorphism scores in relation to English.

Why might GOAT work? The goal for Procrustes is to find the ideal linear transformation $W_{ideal} \in \mathbb{R}^{d \times d}$ to map the spaces, where d is the word embedding dimension. Seeds in Procrustes solve Eq. 1 to find an approximation W to W_{ideal} . Accordingly, the seeds can be thought of as samples from which one deduces the optimal linear transformation. This is a supervised learning problem, so when there are few seeds/samples, it is difficult to estimate W_{ideal} . Furthermore, the entire space X is mapped by W to a shared space with Y meaning that *every* point in X is subject to a potentially inaccurate mapping W : the mapping extrapolates to the entire space. GOAT does not suffer this issue, and can induce non-linear relationships. Graph methods can be thought of as a semi-supervised learning problem: even words that don’t serve as seeds are incorporated in the matching process. The graph manifold provides additional information that can be exploited.

Secondly, the dimension of the relationship between words in GOAT is much lower than for Procrustes. For GOAT, the relationship is one-dimensional: distance. As words for the Procrustes method are embedded in d -dimensional Euclidean space, their relationships have a magnitude *and* a direction: they are $\{d + 1\}$ -dimensional. It is possible that the lower dimension in GOAT makes it ro-

	<i>Prev</i>	<i>-EP</i>	<i>-EG</i>	<i>Prev</i>	<i>-EP</i>	<i>-EG</i>	<i>Prev</i>	<i>-EP</i>	<i>-EG</i>	<i>Prev</i>	<i>-EP</i>	<i>-EG</i>	<i>Prev</i>	<i>-EP</i>	<i>-EG</i>	<i>Prev</i>	<i>-EP</i>	<i>-EG</i>
<i>Seeds</i>	<u>en-bn</u>			<u>bn-en</u>			<u>en-bs</u>			<u>bs-en</u>			<u>en-de</u>			<u>de-en</u>		
0	0.1	0.0	0.1	0.0	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	48.8	61.8	0.0	34.5	59.2	0.3
25	31.9	44.3	0.5	44.4	53.1	6.2	48.1	58.4	0.4	54.6	61.5	0.6	48.5	61.7	59.1	34.1	59.3	56.8
75	33.1	44.7	39.5	45.2	53.5	49.1	47.9	58.1	55.1	54.6	61.7	57.2	48.8	62.3	59.7	47.0	59.4	57.1
100	33.8	45.3	39.7	45.5	53.9	48.1	47.7	58.1	55.4	55.5	61.3	57.5	49.0	62.2	59.7	47.6	59.4	56.8
2000	45.9	48.7	49.3	55.2	56.6	56.3	58.1	59.9	60.5	65.9	66.6	69.1	63.1	66.3	69.4	60.7	65.1	67.9
4000	60.3	50.5	61.2	65.9	55.2	68.9	70.6	63.4	71.8	86.1	69.7	85.7	74.2	72.9	79.5	71.4	67.2	77.6
	<u>en-et</u>			<u>et-en</u>			<u>en-fa</u>			<u>fa-en</u>			<u>en-id</u>			<u>id-en</u>		
0	0.2	0.2	0.0	0.1	0.0	0.0	0.5	0.5	0.0	0.1	0.1	0.0	2.9	4.4	0.0	1.0	1.2	0.0
25	47.0	60.4	5.9	63.2	70.2	15.0	42.7	54.4	4.5	45.1	55.1	2.0	51.3	65.8	0.6	58.0	66.7	1.8
75	47.5	60.6	58.6	64.2	69.8	66.5	42.9	54.1	51.9	45.8	55.3	52.0	53.6	66.0	63.3	57.0	66.7	64.2
100	47.3	61.1	59.0	64.3	70.2	66.7	43.0	54.3	52.3	45.6	55.5	52.7	54.3	66.2	63.2	57.8	67.1	63.9
2000	64.0	66.6	67.4	74.2	74.1	75.0	54.7	58.0	58.4	53.8	57.5	56.9	70.7	72.1	74.2	70.1	72.5	72.9
4000	77.4	71.7	80.2	86.4	80.7	87.2	65.9	62.4	67.4	65.5	60.1	67.0	84.3	76.8	86.2	80.5	78.2	83.7
	<u>en-mk</u>			<u>mk-en</u>			<u>en-ms</u>			<u>ms-en</u>			<u>en-ru</u>			<u>ru-en</u>		
0	0.1	0.0	0.0	0.0	0.0	0.1	0.1	0.0	0.0	0.0	0.1	0.0	0.0	0.1	0.1	2.8	4.4	0.0
25	55.8	63.9	8.0	63.2	68.8	0.6	36.6	62.6	0.9	38.6	65.3	0.2	55.7	67.7	66.1	54.4	63.9	62.0
75	56.5	64.3	63.3	63.6	68.8	67.9	42.6	62.6	59.8	56.8	65.6	62.8	55.7	68.1	67.0	54.8	63.9	61.6
100	56.2	64.1	64.2	64.4	69.4	67.5	42.9	63.0	58.6	57.7	65.8	63.1	55.9	67.9	66.4	55.1	63.8	61.3
2000	64.6	66.8	67.9	71.7	71.0	73.1	65.0	67.0	68.5	67.4	69.4	69.7	67.5	72.6	74.2	68.5	69.3	72.5
4000	75.6	68.9	77.1	88.8	74.1	91.1	79.3	70.7	79.5	77.4	70.0	79.1	83.3	79.3	86.5	89.3	77.4	89.3
	<u>en-ta</u>			<u>ta-en</u>			<u>en-vi</u>			<u>vi-en</u>			<u>en-zh</u>			<u>zh-en</u>		
0	0.0	0.0	0.1	0.5	0.6	0.0	0.1	0.0	0.1	0.1	0.1	0.9	0.0	0.1	0.0	0.0	0.0	0.0
25	1.8	2.2	0.6	44.4	51.4	2.4	0.4	0.4	0.2	3.3	5.3	0.2	8.7	52.7	1.7	9.3	48.1	0.8
75	30.5	40.4	35.7	45.1	52.4	46.9	22.9	55.0	1.2	45.2	59.6	54.4	17.4	51.6	46.4	14.1	51.1	48.0
100	30.6	40.2	36.8	45.4	52.5	47.9	30.2	55.6	36.2	47.1	59.3	56.3	18.5	51.6	47.3	15.2	51.0	48.0
2000	40.6	42.7	44.2	55.1	55.3	56.2	61.1	67.3	65.8	66.3	73.5	71.5	49.3	58.0	57.7	41.8	57.6	56.8
4000	49.1	44.0	51.7	73.8	65.3	71.6	77.9	73.0	80.5	82.1	80.1	84.3	67.5	66.4	75.1	66.2	65.3	73.3

Table 3: P@1 of Combination Exps. -EP starts with GOAT, ends with IterProc. -EG: IterProc, ends with GOAT. *Prev* is previous best of prior experiments. Some seed sizes omitted for brevity (see Appendix).

bust to noise, explaining why GOAT outperforms Procrustes in low-resource settings. This hypothesis should be investigated in follow-up studies.

6 Related Work

BLI Recent years have seen a proliferation of the BLI literature (e.g. Ruder et al., 2018; Aldarmaki et al., 2018; Joulin et al., 2018; Doval et al., 2018; Artetxe et al., 2019a; Huang et al., 2019; Patra et al., 2019; Zhang et al., 2020; Biesialska and Ruiz Costa-Jussà, 2020). Many use Procrustes-based solutions, which assume that embedding spaces are roughly isomorphic. Wang et al. (2021) argue that the mapping can only be piece-wise linear, and induce multiple mappings. Ganesan et al. (2021) learn an “invertible neural network” as a non-linear mapping of spaces, and Cao and Zhao (2018) align spaces using point set registration. Many approaches address only high-resource languages. The tendency to evaluate on similar languages with high-quality data from similar domains hinders advancement in the field (Artetxe et al., 2020).

BLI with OT Most similar to ours are BLI approaches which incorporate OT formulations using the Sinkhorn and/or Hungarian algorithms (e.g. Alvarez-Melis and Jaakkola, 2018; Alaux et al., 2018). Grave et al. (2019) optimize “Procrustes in

Wasserstein Distance”, iteratively updating a linear transformation and permutation matrix using Frank-Wolfe on samples from embedding spaces \mathbf{X} and \mathbf{Y} . Zhao et al. (2020b) and Zhang et al. (2017) also use an iterative procedure. Ramírez et al. (2020) combine Procrustes and their Iterative Hungarian algorithms. Xu et al. (2018) use Sinkhorn distance in the loss function, and (Zhang et al., 2017) use Sinkhorn to minimize distance between spaces. Haghighi et al. (2008) use the Hungarian Algorithm for BLI from text. Lian et al. and Alaux et al. (2018) align all languages to a common space for multilingual BLI. The latter use Sinkhorn to approximate a permutation matrix in their formulation. Zhao et al. (2020a) incorporate OT for semi-supervised BLI.

7 Conclusion

We perform bilingual lexicon induction from word embedding spaces of 32 diverse language pairs, utilizing the newly-developed GOAT algorithm for graph-matching. Performance is strong across all pairs, especially on dissimilar languages with low-supervision. As the major use case of low-resource BLI and MT is dissimilar languages with low supervision, the strong performance of GOAT is an encouraging result for real-world applications.

494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550

References

Jean Alaux, Edouard Grave, Marco Cuturi, and Armand Joulin. 2018. Unsupervised hyper-alignment for multilingual word embeddings. In *International Conference on Learning Representations*.

Hanan Aldarmaki, Mahesh Mohan, and Mona Diab. 2018. Unsupervised word mapping using structural similarities in monolingual embeddings. *Transactions of the Association for Computational Linguistics*, 6:185–196.

David Alvarez-Melis and Tommi Jaakkola. 2018. Gromov-Wasserstein alignment of word embedding spaces. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1881–1890, Brussels, Belgium. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 789–798, Melbourne, Australia. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2019a. Bilingual lexicon induction through unsupervised machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5002–5007, Florence, Italy. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2019b. An effective approach to unsupervised machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 194–203, Florence, Italy. Association for Computational Linguistics.

Mikel Artetxe, Sebastian Ruder, Dani Yogatama, Gorka Labaka, and Eneko Agirre. 2020. A call for more rigor in unsupervised cross-lingual learning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7375–7388, Online. Association for Computational Linguistics.

Magdalena Marta Biesialska and Marta Ruiz Costa-Jussà. 2020. Refinement of unsupervised cross-lingual word embeddings. In *ECAI 2020, 24th European Conference on Artificial Intelligence: 29 August–8 September 2020, Santiago de Compostela, Spain: including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020): proceedings*, pages 1–4. Ios Press.

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.

Hailong Cao and Tiejun Zhao. 2018. Point set registration for unsupervised bilingual lexicon induction. In *IJCAI*, pages 3991–3997.

Jaewon Chung, Benjamin D Pedigo, Eric W Bridgeford, Bijan K Varjavand, Hayden S Helm, and Joshua T Vogelstein. 2019. Graspy: Graph statistics in python. *Journal of Machine Learning Research*, 20(158):1–7.

Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. Word translation without parallel data. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

Marco Cuturi. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26:2292–2300.

Paula Czarrowska, Sebastian Ruder, Édouard Grave, Ryan Cotterell, and Ann Copestake. 2019. Don’t forget the long tail! a comprehensive analysis of morphological generalization in bilingual lexicon induction. 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 974–983.

Yerai Doval, Jose Camacho-Collados, Luis Espinosa-Anke, and Steven Schockaert. 2018. Improving cross-lingual word embeddings by meeting in the middle. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 294–304, Brussels, Belgium. Association for Computational Linguistics.

Donniell E Fishkind, Sancar Adali, Heather G Patsolic, Lingyao Meng, Digvijay Singh, Vince Lyzinski, and Carey E Priebe. 2019. Seeded graph matching. *Pattern recognition*, 87:203–215.

Marguerite Frank, Philip Wolfe, et al. 1956. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110.

Ashwinkumar Ganesan, Francis Ferraro, and Tim Oates. 2021. Learning a reversible embedding mapping using bi-directional manifold alignment. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3132–3139, Online. Association for Computational Linguistics.

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

Chengyue Gong, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2018. Frage: Frequency-agnostic word representation. *arXiv preprint arXiv:1809.06858*.

606	John C Gower, Garnt B Dijksterhuis, et al. 2004. <i>Procrustes problems</i> , volume 30. Oxford University Press on Demand.	660
607		661
608		662
609	Edouard Grave, Armand Joulin, and Quentin Berthet. 2019. Unsupervised alignment of embeddings with wasserstein procrustes. In <i>The 22nd International Conference on Artificial Intelligence and Statistics</i> , pages 1880–1890. PMLR.	663
610		664
611		665
612		666
613		
614	Aria Haghighi, Percy Liang, Taylor Berg-Kirkpatrick, and Dan Klein. 2008. Learning bilingual lexicons from monolingual corpora . In <i>Proceedings of ACL-08: HLT</i> , pages 771–779, Columbus, Ohio. Association for Computational Linguistics.	667
615		668
616		669
617		670
618		671
619	Jiaji Huang, Qiang Qiu, and Kenneth Church. 2019. Hubless nearest neighbor search for bilingual lexicon induction . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4072–4080, Florence, Italy. Association for Computational Linguistics.	672
620		673
621		674
622		675
623		676
624		
625	Roy Jonker and Anton Volgenant. 1987. A shortest augmenting path algorithm for dense and sparse linear assignment problems. <i>Computing</i> , 38(4):325–340.	677
626		678
627		679
628		680
629	Armand Joulin, Piotr Bojanowski, Tomas Mikolov, Hervé Jégou, and Edouard Grave. 2018. Loss in translation: Learning bilingual word mapping with a retrieval criterion . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2979–2984, Brussels, Belgium. Association for Computational Linguistics.	681
630		682
631		683
632		684
633		685
634		686
635	Harold W Kuhn. 1955. The hungarian method for the assignment problem. <i>Naval research logistics quarterly</i> , 2(1-2):83–97.	687
636		688
637		689
638	Xin Lian, Kshitij Jain, Jakub Truszkowski, Pascal Poupart, and Yaoliang Yu. Unsupervised multilingual alignment using wasserstein barycenter.	690
639		691
640		692
641	Kelly Marchisio, Kevin Duh, and Philipp Koehn. 2020. When does unsupervised machine translation work? In <i>Proceedings of the Fifth Conference on Machine Translation</i> , pages 571–583, Online. Association for Computational Linguistics.	693
642		694
643		695
644		696
645		697
646	Kelly Marchisio, Youngser Park, Ali Saad-Eldin, Anton Alyakin, Kevin Duh, Carey Priebe, and Philipp Koehn. 2021. An analysis of Euclidean vs. graph-based framing for bilingual lexicon induction from word embedding spaces . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 738–749, Punta Cana, Dominican Republic. Association for Computational Linguistics.	698
647		699
648		700
649		701
650		702
651		703
652		704
653		705
654	Ndapa Nakashole and Raphael Flauger. 2018. Characterizing departures from linearity in word translation . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 221–227, Melbourne, Australia. Association for Computational Linguistics.	706
655		707
656		708
657		709
658		710
659		711
	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.	712
		713
		714
		715
	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.	
	Gabriel Peyre and Marco Cuturi. 2019. Computational optimal transport. <i>Foundations and Trends in Machine Learning</i> , 11(5-6):355–607.	
	Guillem Ramírez, Rumen Dangovski, Preslav Nakov, and Marin Soljačić. 2020. On a novel application of wasserstein-procrustes for unsupervised cross-lingual learning. <i>arXiv preprint arXiv:2007.09456</i> .	
	Sebastian Ruder, Ryan Cotterell, Yova Kementchedjhiya, and Anders Søgaard. 2018. A discriminative latent-variable model for bilingual lexicon induction . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 458–468, Brussels, Belgium. Association for Computational Linguistics.	
	Ali Saad-Eldin, Benjamin D Pedigo, Carey E Priebe, and Joshua T Vogelstein. 2021. Graph matching via optimal transport. <i>arXiv preprint arXiv:2111.05366</i> .	
	Peter H Schönemann. 1966. A generalized solution of the orthogonal procrustes problem. <i>Psychometrika</i> , 31(1):1–10.	
	Haoyue Shi, Luke Zettlemoyer, and Sida I. Wang. 2021. Bilingual lexicon induction via unsupervised bitext construction and word alignment . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 813–826, Online. Association for Computational Linguistics.	
	Richard Sinkhorn. 1967. Diagonal equivalence to matrices with prescribed row and column sums. <i>The American Mathematical Monthly</i> , 74(4):402–405.	
	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.	
	Joshua T Vogelstein, John M Conroy, Vince Lyzinski, Louis J Podrazik, Steven G Kratzer, Eric T Harley, Donniell E Fishkind, R Jacob Vogelstein, and Carey E Priebe. 2015. Fast approximate quadratic	

- 716 programming for graph matching. *PLOS one*,
717 10(4):e0121002.
- 718 Ivan Vulić, Goran Glavaš, Roi Reichart, and Anna Ko-
719 rhonen. 2019. [Do we really need fully unsuper-](#)
720 [vised cross-lingual embeddings?](#) In *Proceedings of*
721 *the 2019 Conference on Empirical Methods in Natu-*
722 *ral Language Processing and the 9th International*
723 *Joint Conference on Natural Language Processing*
724 *(EMNLP-IJCNLP)*, pages 4407–4418, Hong Kong,
725 China. Association for Computational Linguistics.
- 726 Ivan Vulić, Sebastian Ruder, and Anders Søgaard. 2020.
727 [Are all good word vector spaces isomorphic?](#) In
728 *Proceedings of the 2020 Conference on Empirical*
729 *Methods in Natural Language Processing (EMNLP)*,
730 pages 3178–3192, Online. Association for Computa-
731 tional Linguistics.
- 732 Haozhou Wang, James Henderson, and Paola Merlo.
733 2021. [Multi-adversarial learning for cross-lingual](#)
734 [word embeddings.](#) In *Proceedings of the 2021 Con-*
735 *ference of the North American Chapter of the Asso-*
736 *ciation for Computational Linguistics: Human Lan-*
737 *guage Technologies*, pages 463–472, Online. Associ-
738 ation for Computational Linguistics.
- 739 Ruochen Xu, Yiming Yang, Naoki Otani, and Yuexin
740 Wu. 2018. [Unsupervised cross-lingual transfer of](#)
741 [word embedding spaces.](#) In *Proceedings of the 2018*
742 *Conference on Empirical Methods in Natural Lan-*
743 *guage Processing*, pages 2465–2474, Brussels, Bel-
744 gium. Association for Computational Linguistics.
- 745 Meng Zhang, Yang Liu, Huanbo Luan, and Maosong
746 Sun. 2017. [Earth mover’s distance minimization for](#)
747 [unsupervised bilingual lexicon induction.](#) In *Pro-*
748 *ceedings of the 2017 Conference on Empirical Meth-*
749 *ods in Natural Language Processing*, pages 1934–
750 1945, Copenhagen, Denmark. Association for Com-
751 putational Linguistics.
- 752 Mozhi Zhang, Yoshinari Fujinuma, Michael J. Paul, and
753 Jordan Boyd-Graber. 2020. [Why overfitting isn’t](#)
754 [always bad: Retrofitting cross-lingual word embed-](#)
755 [dings to dictionaries.](#) In *Proceedings of the 58th An-*
756 *nuual Meeting of the Association for Computational*
757 *Linguistics*, pages 2214–2220, Online. Association
758 for Computational Linguistics.
- 759 Mozhi Zhang, Keyulu Xu, Ken-ichi Kawarabayashi, Ste-
760 fanie Jegelka, and Jordan Boyd-Graber. 2019. [Are](#)
761 [girls neko or shōjo? cross-lingual alignment of non-](#)
762 [isomorphic embeddings with iterative normalization.](#)
763 In *Proceedings of the 57th Annual Meeting of the As-*
764 *sociation for Computational Linguistics*, pages 3180–
765 3189, Florence, Italy. Association for Computational
766 Linguistics.
- 767 Xu Zhao, Zihao Wang, Hao Wu, and Yong Zhang.
768 2020a. [Semi-supervised bilingual lexicon induction](#)
769 [with two-way interaction.](#) In *Proceedings of the 2020*
770 *Conference on Empirical Methods in Natural Lan-*
771 *guage Processing (EMNLP)*, pages 2973–2984, On-
772 line. Association for Computational Linguistics.
- Xu Zhao, Zihao Wang, Yong Zhang, and Hao Wu. 773
2020b. [A relaxed matching procedure for unsuper-](#) 774
[vised BLI.](#) In *Proceedings of the 58th Annual Meet-* 775
ing of the Association for Computational Linguistics, 776
pages 3036–3041, Online. Association for Computa- 777
tional Linguistics. 778

Appendix

	-to-en		en-to-		# Embs
	Full	1-1	Full	1-1	
bn	7588	4299	8467	4556	145350
bs	6164	4294	8153	4795	166505
de	10866	4451	14677	4903	200000
en	n/a	n/a	n/a	n/a	200000
es	8667	4445	11977	4866	200000
et	6509	4352	8261	4738	200000
fa	8510	4582	8869	4595	200000
fr	8270	4548	10872	4827	200000
id	9677	4563	9407	4573	200000
it	7364	4478	9657	4815	200000
ja	6819	4112	7135	4351	200000
mk	7197	4259	10075	4820	176947
ms	8140	4650	7394	4454	155629
ru	7452	4084	10887	4812	200000
ta	6850	4225	8091	4744	200000
vi	7251	4775	6353	4507	200000
zh	8891	4450	8728	4381	200000

Table A1: Size of train/test sets before (Full) & after making one-to-one (1-1), with # of embeddings used.

Seeds	en-de		ru-en	
	Rand.	Bary.	Rand.	Bary.
100	45.7	45.9	49.6	50.4
200	47.4	47.5	52.5	52.5
500	52.3	51.9	55.4	55.6
1000	54.6	54.9	58.3	58.1
2000	61.5	61.4	67.1	67.1
4000	74.2	74.2	89.3	89.3

Table A2: SGM with barycenter vs. randomized initialization for languages used in Marchisio et al. (2021). The difference is negligible.

Iterative Results of Iterative Procrustes (IterProc), Iterative SGM (IterSGM), and Iterative GOAT (IterGOAT) are in Table A3. We run the Iterative Procrustes and Iterative SGM procedures of Marchisio et al. (2021) with stochastic-add. Here, Procrustes [or SGM] is run in source \leftrightarrow target directions, hypotheses are intersected, and H random hypotheses are added to the gold seeds and fed into subsequent runs of Procrustes [SGM]. The next iteration adds $2H$ hypotheses, repeating until all hypotheses are chosen. We set $H = 100$ and create an analogous iterative algorithm for GOAT, which we call Iterative GOAT.

IterSGM/GOAT perform similarly across conditions, with a few exceptions where either performs very strongly with no supervision: IterGOAT scores 49.2, 45.2, 34.4, 58.2, and 55.9 for En-De, En-Fa, De-En, Id-En, and Ru-En, respectively, and IterSGM scores 57.3 for En-Ru. On Chinese \leftrightarrow English, IterGOAT underperforms IterSGM, similar to GOAT’s underperformance of SGM in the single run.

Similar to Marchisio et al. (2021), we find that IterProc compensates for an initial poor first run and outperforms IterSGM with a moderate amount of seeds (100+). Extending to the very lowest seeds sizes (0-75), however, IterSGM/IterGOAT are superior. With 25 seeds, IterProc fails for all language pairs except En \leftrightarrow De and En \leftrightarrow Ru, scoring $P@1 < 5$. IterSGM and IterGOAT, however, perform reasonably well for most language pairs with 25 seeds, suggesting that the graph-based framing is the better approach for low-seed levels. At the highest supervision level (2000+ seeds), IterSGM/IterGOAT again tends to be superior.

Differences are minor btwn using GOAT or SGM in the iterative or system combination experiments. This results suggests that mixing Procrustes and graph-based framings is helpful for BLI, regardless of which algorithm one picks. It is interesting to contemplate what other problems might benefit from examination from multiple mathematical framings in one solution, as each may have complementary benefits.

