# Graph Neural Networks for Multiparallel Word Alignment

## Anonymous ACL submission

## Abstract

After a period of decrease, interest in word alignments is increasing again for their usefulness in domains such as typological research, cross-lingual annotation projection and machine translation. Generally, alignment algorithms only use bitext and do not make use of the fact that many parallel corpora are multiparallel. We propose to use graph neural networks (GNNs) and community detection algorithms to exploit the graph structure of multiparallel word alignments. Our GNN approach (i) utilizes information about the meaning, position and language of the input words, (ii) incorporates information from multiple parallel sentences, (iii) can remove edges from the initial alignments, and (iv) provides a prediction model that can generalize beyond the sentences it is trained on. We show that community detection algorithms can provide valuable information for multiparallel word alignment. We show on three word alignment datasets and on a downstream task that our method outperforms previous work.
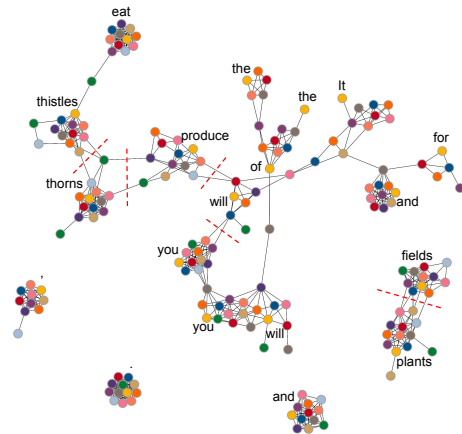
## 1 Introduction

Word alignments are crucial for statistical machine translation (Koehn et al., 2003) and useful for many other multilingual tasks such as neural machine translation (Alkhouli and Ney, 2017; Alkhouli et al., 2016), typological analysis (Lewis and Xia, 2008; Östling, 2015; Asgari and Schütze, 2017), annotation projection (Yarowsky and Ngai, 2001; Fossum and Abney, 2005; Wisniewski et al., 2014; Huck et al., 2019), bilingual lexicon induction (Lample et al., 2018; Artetxe et al., 2018; Shi et al., 2021) and creation of multilingual embeddings (Dufter et al., 2018). The rise of deep learning initially led to a temporary plateau, but interest in word alignments is now increasing, demonstrated by several recent publications (Jalili Sabet et al., 2020; Chen et al., 2020; Dou and Neubig, 2021; Marchisio et al., 2021; Wu and Dredze, 2020;



Figure 1: Alignment graph for the verse "It will produce thorns and thistles for you, and you will eat the plants of the field." in a 12-way multiparallel corpus. Colors represent languages. Each English (yellow) node is annotated with its word. Red dashed lines sever links that incorrectly connect distinct concepts. We exploit community detection algorithms to remove intercommunity and add intra-community links.

Imani et al., 2021)

Most prior work on word alignments uses bitext, but Imani et al. (2021) exploit the fact that many parallel corpora are multiparallel (i.e., they contain more than two parallel corpora). They introduce MPWA (MultiParallel Word Alignment), a framework that makes use of multiparallelism for better word alignments. They represent sets of bilingual word alignments as graphs and cast the word alignment task as an edge prediction problem. To exploit the graph structure, they apply two standard graph algorithms, Adamic-Adar (AdAd) and non-negative matrix factorization (NMF), and achieve improved results. However, these standard graph algorithms are applied to individual multiparallel sentences independently and therefore cannot accumulate knowledge from multiple sentences. Moreover, their edge predictions are solely based on the structure of the graph and do not take ad-

vantage of other beneficial signals such as a word's language, relative position and word meaning.

In this paper, we propose to use graph neural networks (GNNs) to exploit the graph structure of multiparallel word alignments and address the limitations of prior work. GNNs were proposed to extend the powerful current generation of neural network models to processing graph-structured data (Scarselli et al., 2009) and they have gained increasing popularity in many domains, such as social networks (Wu et al., 2020), natural science (Sanchez-Gonzalez et al., 2018), knowledge graphs (Hamaguchi et al., 2017), and recommender systems (He et al., 2020). In contrast to other graph algorithms, GNNs can incorporate heterogeneous sources of signal in the form of node and edge features.

Since the nodes in the graph are words that are translations of each other, we expect them to create densely connected regions or *communities*. Our analysis of the structure of the multiparallel alignment graph confirms this intuition; see Figure 1. We use the community detection algorithms GMC and LPC (see below) to find communities and show that pruning inter-community and adding intra-community edges is helpful. We use community information as node features for our GNN.

A limitation of (Imani et al., 2021) is that it only adds links and does not remove any. We address this by proposing a new method to infer alignments from the alignment probability matrix. Our method predicts new alignment links independently of initial edges. Therefore it is not limited to adding new edges to some initial bilingual alignments.

For our experiments, we follow the setup of Imani et al. (2021). We obtain bilingual alignments using the statistical word aligner Eflomal (Östling and Tiedemann, 2016). We train a GNN model on the resulting graph with a link prediction objective. We show improved results for three language pairs on word alignment (English-French, Finnish-Hebrew and Finnish-Greek). As a demonstration of the importance of high-quality alignments, we use our word alignments to project annotations from high-resource languages to low-resource languages. We improve the performance of a part-of-speech tagger for the Yoruba language by training it over a high-quality dataset, which is created using annotation projection.

**Contributions: i)** We propose graph neural networks that can incorporate a diverse set of features for word alignments in multiparallel corpora and show that GNNs establish a new state of the art in word alignment. **ii)** We show that community detection algorithms improve multiparallel word alignment. **iii)** We show that the improved alignments improve performance on a downstream task. **iv)** We propose a new method to infer alignments from the alignment probability matrix. **v)** We will make our code publicly available.

## 2   Related Work

**Bilingual Word Aligners.** Much work on bilingual word alignment is based on probabilistic models, mostly implementing variants of the IBM models of Brown et al. (1993): e.g., the aligners Giza++ (Och and Ney, 2003), fast-align (Dyer et al., 2013) and Eflomal (Östling and Tiedemann, 2016). They use statistical similarities between word distributions in sentence aligned parallel corpora to learn word alignment models. More recent work, including SimAlign (Jalili Sabet et al., 2020), Awesome-align (Dou and Neubig, 2021), Bidir+CL (Zenkel et al., 2020), and SHIFT-ATT/SHIFT-AET (Chen et al., 2020), uses pretrained neural language and machine translation models. Although neural models achieve superior performance compared to statistical aligners, they are only applicable for less than two hundred high-resource languages that are supported by multilingual language models like BERT (Devlin et al., 2019) and XLM-R (Conneau et al., 2020). This makes statistical models the only option for the majority of the world's languages. Due to its good performance, we use Eflomal as our initial bilingual aligner.

**Multiparallel Corpora.** Prior applications of using multiparallel corpora include reliable translations from small datasets (Cohn and Lapata, 2007), and phrase-based machine translation (PBMT) (Kumar et al., 2007). Multiparallel corpora are also used for language comparison (Mayer and Cysouw, 2012), typological studies (Östling, 2015; Asgari and Schütze, 2017) and PBMT (Nakov and Ng, 2012; Bertoldi et al., 2008; Dyer et al., 2013).

Despite the usefulness of multiparallel corpora, most past work on word alignment has focused on bilingual corpora. Östling (2014)[1] proposed a word alignment method specifically designed for multiparallel corpora. However, this method is outperformed by Eflomal (Östling and Tiedemann, 2016), a "biparallel" method from the same author. Re-

---

[1] `github.com/robertostling/eflomal`

cently, Imani et al. (2021) proposed MPWA (MultiParallel Word Alignment, see §3), which refines the graph structure of an initial multiparallel word alignment using standard graph algorithms. Our work improves on MPWA.

**Graph Neural Networks (GNNs)** have been used to address many problems that are inherently graph-like such as traffic networks, social networks, and physical and biological systems (Liu and Zhou, 2020). Their popularity increased rapidly after the efficient integration of powerful deep learning techniques, such as convolutional neural networks and attention mechanisms, into GNNs (Kipf and Welling, 2016, 2017; Velickovic et al., 2018). GNNs achieve impressive performance in many domains, including social networks (Wu et al., 2020), natural science (Sanchez-Gonzalez et al., 2018), knowledge graphs (Hamaguchi et al., 2017), and recommender systems (He et al., 2020). This has motivated the NLP community to apply GNNs to tasks such as sentence classification (Huang et al., 2020), named entity recognition (Luo and Zhao, 2020), question generation (Pan et al., 2020) and summarization (Fernandes et al., 2019). As far as we know, our work is the first to apply GNNs to word alignment.

## 3 Background

### 3.1 MPWA

MPWA (MultiParallel Word Alignment) aims to utilize the synergy between multiple language pairs to improve bilingual word alignments (Imani et al., 2021). The rationale is that some of the missing alignment edges between a source and a target language can be recovered by using their alignments with words in other languages.

The first step in MPWA is to create bilingual alignments for all language pairs in a multiparallel corpus using a bilingual word aligner. Then the bilingual alignments for a given multiparallel sentence are represented as a graph where words are nodes and initial word alignments are edges. Figure 1 gives an example: a bilingual alignment graph for a 12-way multiparallel corpus.

MPWA tries to infer missing alignment links based on the graph structure, casting the word alignment task as an edge prediction problem. Imani et al. (2021) use two traditional graph algorithms, Adamic-Adar and non-negative matrix factorization, for edge prediction. We replace them here with more powerful GNNs.

### 3.2 Community detection (CD)

The nodes in the alignment graph are words that are translations of each other. If the initial bilingual alignments are of good quality, we expect these translated words to form densely connected regions or *communities* in the graph; see Figure 1. We expect these communities to be disconnected, each corresponding to a distinct connected component. In other words, ideally, words representing a concept should be densely connected, but there should be no links between different concepts.

To examine to what extent this expectation is met, we count the components in the original (Eflomal-generated) graph. Table 1 shows that, for most sentences, the average number of components per sentence is less than three. But intuitively, the number of components (representing the concepts in the sentence) should be roughly equal to sentence length (or at least the number of content words). This indicates that there are many links that incorrectly connect different concepts. To detect such links, we use community detection (CD) algorithms.

CD algorithms find subnetworks of nodes that form tightly knit groups that are only loosely connected with a small number of links (Girvan and Newman, 2002). These algorithms try to maximise the modularity measure Newman and Girvan (2004). Modularity measures how beneficial a division of a community into two communities is, in the sense that there are many links within communities and only a few between them. Given a graph $G$ with $n$ nodes and $m$ edges, and $G$'s adjacency matrix $A \in \mathbb{R}^{n \times n}$. modularity is defined as:

$$mod = \frac{1}{2m} \sum_{ij} \left( A_{ij} - \gamma \frac{d_i d_j}{2m} \right) I(c_i, c_j) \quad (1)$$

Where $d_i$ is the degree of node $i$, and $I(c_i, c_j)$ is 1 if nodes $i$ and $j$ are in the same community and 0 otherwise.

We experiment with two CD algorithms:

- Greedy modularity communities (GMC). This method uses Clauset-Newman-Moore greedy modularity maximization (Clauset et al., 2004). GMC begins with each node in its own community and greedily joins the pair of communities that most increases modularity until no such pair exists.

- Label propagation communities (LPC). This method finds communities in a graph using

| | FIN-HEB | | FIN-GRC | | ENG-FRA | |
|---|---|---|---|---|---|---|
| | #CC | $F_1$ | #CC | $F_1$ | #CC | $F_1$ |
| Eflomal intersection | 2.2 | 0.404 | 1.6 | 0.646 | 2.2 | 0.678 |
| GMC | 13.7 | 0.396 | 10.1 | 0.375 | 13.5 | 0.411 |
| LPC | 41.5 | 0.713 | 37.1 | 0.754 | 46.0 | 0.767 |
| Sentence length | 25.7 | | 23.2 | | 27.4 | |

Table 1: Effect of community detection algorithms on alignment prediction. #CC: average number of connected components. $F_1$: word alignment performance. LPC consistently increases the number of components and increases $F_1$.

label propagation (Cordasco and Gargano, 2010). It begins by giving a label to each node of the network. Then each node's label is updated by the most frequent label among its neighbors in each iteration. LPC's semi-synchronous algorithm, which at each step performs label propagation on a portion of nodes, quickly converges to a stable labeling.

After detecting communities, we link all nodes inside a community and remove all inter-community links. GMC (LPC) on average removes 3% (7%) of the edges. Table 1 reports the average number of graph components per sentence before and after runing GMC and LPC, as well as the corresponding $F_1$ for word alignment. We see that the number of communities found is lower for GMC than for LPC; therefore, LPC identifies more candidate links for deletion.[2] Comparing the number of communities detected with the average sentence length, GMC seems to have failed to detect enough communities to split different concepts properly. The $F_1$ scores confirm this observation and show that LPC performs well at detecting the communities we are looking for.

These results indicate that CD algorithms can provide valuable information. To exploit this in our GNN model, we add a node's community information as a GNN node feature of that node.

## 4 Methods

### 4.1 GNN in MPWA

GNNs can be used in transductive or inductive settings. In a transductive setting, nodes are represented as node IDs, and the final model can only be used for inference over the same graph that it is trained on. In an inductive setting, nodes are represented as feature vectors, and the final model has the advantage of being applicable to a different graph in inference. We use the inductive setting. GNNs can incorporate different sources of signal in the form of node and edge features. We only use node features. All are trained (or finetuned) during GNN training.

#### 4.1.1 Node Features

We use three main types of node features: (i) graph structural features, (ii) community-based features and (iii) word content features.

**Graph structural features.** We use *degree, closeness* (Freeman, 1978) *, betweenness* (Brandes, 2001) *, load* (Newman, 2001) and *harmonic centrality* (Boldi and Vigna, 2014) features as additional information about the graph structure. These features are continuous numbers, providing information about the position and connectivity of the nodes within the graph. We standardize (i.e., z-score) each of these features across all nodes in the graph, and train an embedding of size four for each feature.[3]

**Community-based features.** We use the community detection algorithms GMC and LPC (see §3.2) to identify communities in the graph. Then we take the community membership information of the nodes as one-hot vectors and learn an embedding of size 32 for each of the two algorithms.

**Word content features.** We train embeddings for *word position* (size 32) and *word language* (size 20). We learn 100-dimensional multilingual *word embeddings* using Levy et al. (2017)'s sentence-ID method on the 84 PBC languages selected by Imani et al. (2021). Word embeddings serve as initialization and are updated during GNN training.

After concatenating these features, each node is represented by a 236 dimensional vector that is then fed to the encoder.

#### 4.1.2 Model Architecture

Our model is inspired by the Graph Auto Encoder (GAE) model of Kipf and Welling (2016) for the link prediction task. The architecture consists of an encoder and a decoder.

The **encoder** is a graph attention network (GAT) (Veličković et al., 2018) with two GATConv layers

---

[2] LPC may detect more communities than average sentence length because of null words: words that have no translation in the other languages, giving rise to separate communities.

[3] Learning a size-four embedding instead of a single number gives the feature a weight similar to other features – which have a feature vector of about the same size.
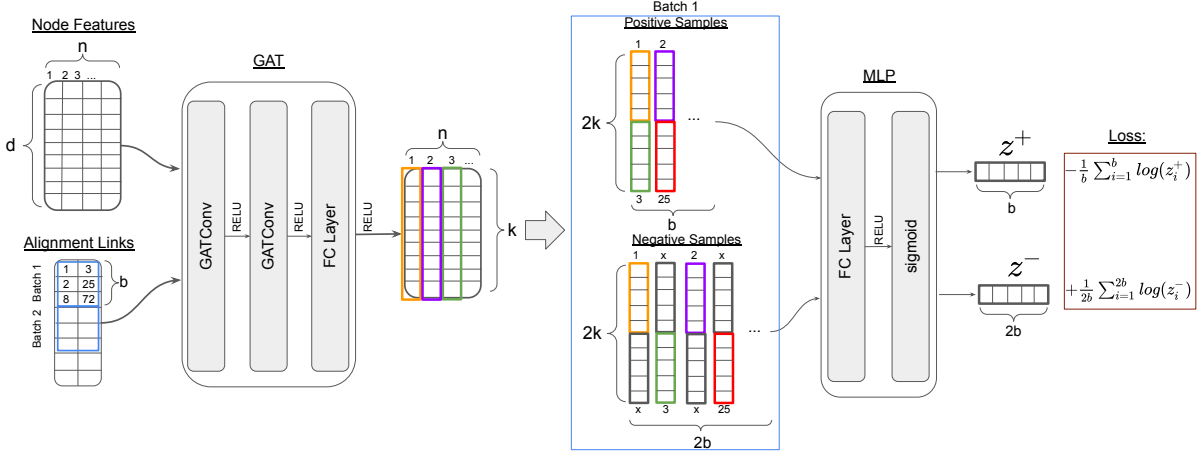
Figure 2: GNN training. At each training step, node features and alignment links of a multiparallel sentence are fed to a graph attention network (GAT) that creates hidden representations for all nodes. On the decoder side, at each step, one batch of alignment links and hidden node representations is used to create positive and negative samples, which are then processed and classified by a multi-layer perceptron (MLP). Parameters of GAT and MLP are updated after each training step. After all alignment links are processed in multiple training steps, training continues with the next multiparallel sentence in the training set. FC = fully connected.

followed by a fully-connected layer. Layers are connected by RELU non-linearities. A GATConv layer computes its output $\mathbf{x}'_i$ for a node $i$ from its input $\mathbf{x}_i$ as

$$\mathbf{x}'_i = \alpha_{i,i}\mathbf{W}\mathbf{x}_i + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j}\mathbf{W}\mathbf{x}_j, \qquad (2)$$

where $\mathbf{W}$ is a weight matrix, $\mathcal{N}(i)$ is some *neighborhood* of node $i$ in the graph, and $\alpha_{i,j}$ is the attention coefficient indciating the importance of node $j$'s features to node $i$. $\alpha_{i,j}$ is computed as

$$\alpha_{i,j} = \frac{\exp\left(g\left(\mathbf{a}^\top[\mathbf{W}\mathbf{x}_i \,\|\, \mathbf{W}\mathbf{x}_j]\right)\right)}{\sum_{k \in \mathcal{N}(i) \cup \{i\}} \exp\left(g\left(\mathbf{a}^\top[\mathbf{W}\mathbf{x}_i \,\|\, \mathbf{W}\mathbf{x}_k]\right)\right)} \qquad (3)$$

where $\cdot^T$ is transposition and $\|$ is concatanation, $g$ is the LeakyReLU non-linearity, and $\mathbf{a}$ is a weight vector. Given the features for the nodes and their alignment edges, the encoder creates a contextualized hidden representation for each node.

Based on the hidden representations of two nodes, the **decoder** predicts whether a link connects them. The decoder architecture consists of a fully connected layer, a RELU non-linearity and a sigmoid layer.

**Training.** Figure 2 displays our GNN model and the training process. The outer loop iterates over the multiparallel sentences in the training set. The training set contains one graph for each sentence; the graph is constructed using the bilingual alignment edges between all language pairs.

Each graph is divided into multiple batches. Each batch contains a random subset of the graph's edges as positive samples. The negative samples are created as follows: Given a sentence $u_1u_2u_3 \ldots u_n$ in language $U$ and its translation $v_1v_2v_3 \ldots v_m$ in language $V$, for each alignment edge $u_i : v_j$ in the current batch, two negative edges $u_i : v'_j$ and $u'_i : v_j$ ($j' \neq j$, $i' \neq i$) are randomly sampled.

For each training batch, the encoder takes the batch's whole graph (i.e., node features for all graph nodes and all graph edges) as input and computes hidden representations for the nodes. On the decoder side, for each link of the batch, the hidden representations of the attached nodes are concatenated to create the decoder's input. The decoder's target is the link's class: 1 (resp. 0) for positive (resp. negative) links. We train with a binary classification objective:

$$\mathcal{L} = -\frac{1}{b}\sum_{i=1}^{b}\log(p_i^+) + \frac{1}{2b}\sum_{i=1}^{2b}\log(p_i^-) \qquad (4)$$

where $b$ is the batch size and $p_i^+$ and $p_i^-$ are the model predictions for the $i^{th}$ positive and negative samples within the batch. Parameters of the encoder and decoder as well as the node-embedding feature layer are updated after each training step.

### 4.1.3 Inducing Alignment Edges

When our trained GNN model is used to predict alignment edges between a source sentence

5

$\hat{x} = x_1, x_2, \ldots, x_m$ in language $X$ and a target sentence $\hat{y} = y_1, y_2, \ldots, y_l$ in language $Y$, it produces an alignment probability matrix $S^4$ of size $m \times l$ where $S_{ij}$ is the predicted alignment probability between words $x_i$ and $y_j$. Using these values directly to infer alignment edges is usually suboptimal; therefore, more sophisticated methods have been suggested (Ayan and Dorr, 2006; Liang et al., 2006). Here we propose a new approach: it combines Koehn et al. (2005)'s Grow-Diag-Final-And (GDFA) with Dou and Neubig (2021)'s probability thresholding. We modify the latter to account for the variable size of the probability matrix (i.e., length of source/target sentences). Our method is not limited to adding new edges to some initial bilingual alignments, a limitation of prior work. As we predict each edge independently, some initial links can be discarded from the final alignment.

We start by creating a set of *forward* (source-to-target) alignment edges and a set of *backward* (target-to-source) alignment edges. To this end, first, inspired by probability thresholding (Dou and Neubig, 2021), we apply softmax to $S$, and zero out probabilities below a threshold to get a source-to-target probability matrix $S^{XY}$:

$$S^{XY} = S * (\text{softmax}(S) > \frac{\alpha}{l}) \qquad (5)$$

Analogously, we compute the target-to-source probability matrix $S^{YX}$:

$$S^{YX} = S^\top * (\text{softmax}(S^\top) > \frac{\alpha}{m}) \qquad (6)$$

where $\alpha$ is a sensitivity hyperparameter, e.g., $\alpha = 1$ means that we pick edges with a probability higher than average. We experimentally set $\alpha = 2$. Next, from each row of $S^{XY}$ ($S^{YX}$), we pick the cell with the highest value (if any exists) and add this edge to the *forward* (*backward*) set.

We create the final set of alignment edges by applying the GDFA symmetrization method (Koehn et al., 2005) to *forward* and *backward* sets. The gist of GDFA is to use the intersection of *forward* and *backward* as initial alignment edges and add more edges from the union of *forward* and *backward* based on a number of heuristics. We call this method *TGDFA* (Thresholding GDFA).

In addition to the TGDFA alignments, we also experiment with combining them with the original bilingual GDFA alignments. We do so by adding

bilingual GDFA edges to the union of *forward* and *backward* before performing the GDFA heuristics. We refer to these alignments as *TGDFA+orig*.

## 4.2 Annotation Projection

Annotation projection automatically creates linguistically annotated corpora for low-resource languages. A model trained on data with "annotation-projected" labels can perform better than full unsupervision. Here, we focus on universal part-of-speech (UPOS) tagging (Petrov et al., 2012) for the low resource target language *Yoruba*; this language only has a small set of annotated sentences in Universal Dependencies (Nivre et al., 2020) and has poor POS results in unsupervised settings (Kondratyuk and Straka, 2019).

The quality of the target annotated corpus depends on the quality of the annotations in the source languages and the quality of the word alignments between sources and target. We use the Flair (Akbik et al., 2019) POS taggers for three high resource languages, English, German and French (Akbik et al., 2018), to annotate 30K verses whose Yoruba translations are available in PBC. We then transfer the POS tags from source to target using three different approaches: (i) We directly transfer annotations from English to the target. (ii) For each word in the target language, we get its alignments in the three source languages and predict the majority POS to annotate the target word. (iii) We repeat (ii) using alignments from our GNN (TGDFA) model instead of the original bilingual alignments. In all three approaches, we discard any target sentence from the POS tagger training data if more than 50% of its words are annotated with the "X" (other) tag.

We train a Flair SequenceTagger model on the target annotated data using mBERT embeddings (Devlin et al., 2019) and evaluate on Yoruba test from Universal Dependencies.[5]

# 5 Experimental Setup

## 5.1 Word Alignment Datasets

Following Imani et al. (2021), we use PBC as our multiparallel corpus. PBC contains 1758 editions of the Bible in 1334 languages, aligned at the verse level. A verse can contain more than one sentence, but we take it as one unit since sentence level alignments are not available.

For our main evaluation, we use the two word alignment gold datasets for PBC published by

---

[4]For inference, we feed all possible alignment links between source and target to the decoder.

[5]https://universaldependencies.org/

| Method | FIN-HEB | | | | FIN-GRC | | | | ENG-FRA | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Prec. | Rec. | $F_1$ | AER | Prec. | Rec. | $F_1$ | AER | Prec. | Rec. | $F_1$ | AER |
| Eflomal (intersection) | **0.818** | 0.269 | 0.405 | 0.595 | **0.897** | 0.506 | 0.647 | 0.353 | **0.971** | 0.521 | 0.678 | 0.261 |
| Eflomal (GDFA) | 0.508 | 0.448 | 0.476 | 0.524 | 0.733 | 0.671 | 0.701 | 0.300 | 0.856 | 0.710 | 0.776 | 0.221 |
| WAdAd (intersection) | 0.781 | 0.612 | 0.686 | 0.314 | 0.849 | 0.696 | 0.765 | 0.235 | 0.938 | 0.689 | 0.794 | 0.203 |
| NMF (intersection) | 0.780 | 0.576 | 0.663 | 0.337 | 0.864 | 0.669 | 0.754 | 0.248 | 0.948 | 0.624 | 0.753 | 0.245 |
| WAdAd (GDFA) | 0.546 | **0.693** | 0.611 | 0.389 | 0.707 | **0.783** | 0.743 | 0.257 | 0.831 | **0.796** | 0.813 | 0.186 |
| NMF (GDFA) | 0.548 | 0.646 | 0.593 | 0.407 | 0.720 | 0.759 | 0.739 | 0.261 | 0.844 | 0.767 | 0.804 | 0.195 |
| GNN (TGDFA) | 0.811 | 0.648 | **0.720** | **0.280** | 0.845 | 0.724 | **0.780** | **0.220** | 0.926 | 0.711 | 0.804 | 0.192 |
| GNN (TGDFA+orig) | 0.622 | 0.683 | 0.651 | 0.349 | 0.738 | 0.780 | 0.758 | 0.242 | 0.863 | 0.789 | **0.824** | **0.174** |

Table 2: Word alignment results on PBC for GNN and baselines. The best result in each column is in bold. GNN outperforms the baselines as well as the graph algorithms WAdAd and NMF on $F_1$ and AER.

Imani et al. (2021): Blinker (Melamed, 1998) and HELFI (Yli-Jyrä et al., 2020).

**The HELFI dataset** contains the Hebrew Bible, Greek New Testament and their translations into Finnish. The Finnish-Hebrew dataset has word level alignments for 22,291 verses and the Finnish-Greek dataset for 7,909. We use Imani et al. (2021)'s train/dev/test splits. **The Blinker dataset** provides word level alignments between English and French for 250 Bible verses.

The graph algorithms used by Imani et al. (2021) operate on each multiparallel sentence separately. In contrast, our approach allows for an inductive setting where a model is trained on a training set and then evaluated on a separate test set. This allows our model to learn from multiple training samples and use its accumulated knowledge on the new test samples with fast inference. We combine the verses in training sets of Finnish-Hebrew and Finnish-Greek for a combined train set size of size 24,159[6].

### 5.2 Initial Word Alignments

We use the Eflomal statistical word aligner to obtain bilingual alignments. We do not consider SimAlign (Jalili Sabet et al., 2020) since it is shown to perform poorly for languages whose representations in the multilingual pretrained language model are of low quality, which includes the target languages in the HELFI dataset (Imani et al., 2021). We evaluate on the same subset of 84 languages as Imani et al. (2021). To train Eflomal for a target language pair, we use all available translations; e.g., for a language pair with two and four different versions of the Bible, Eflomal is trained on all eight translation pairs.

We use Eflomal asymmetrical alignments post-processed with the intersection heuristic to get high precision bilingual alignments as input to the GNN.

### 5.3 Training Details

We use NetworkX[7] for graph structural and community-based feature extraction; PyTorch Geometric[8] to construct and train the GNN, and Gensim[9] to train sentence-ID embeddings. The model's hidden layer size is 512 for both GATConv and Linear layers. We train for one epoch on the train set – a small portion of the train set is enough to learn good embeddings (see §6.1.1). For training, we use a batch size of 400 and learning rate of .001 with AdamW (Loshchilov and Hutter, 2017).

## 6 Experiments and Results

### 6.1 Multiparallel corpus results

Table 2 shows results on Blinker and HELFI for our GNNs and the baselines: bilingual alignments and the traditional graph algorithms WAdAd and NMF from (Imani et al., 2021). Our GNNs provide a better trade-off between precision and recall, most likely thanks to their ability to remove edges, and achieve the best $F_1$ and AER on all three datasets, outperforming WAdAd and NMF.

GNN (TGDFA) achieves the best results on HELFI (FIN-HEB, FIN-GRC) while GNN (TGDFA+orig) is best on Blinker (ENG-FRA). As argued in (Imani et al., 2021), this is mostly due to the different ways these two datasets were annotated. Most HELFI alignments are one-to-one, while many Blinker alignments are many-to-many: phrase-level alignments where every word in a source phrase is aligned with every word in a target phrase. This suggests that one can choose between

---

[6]Note that we don't use any gold alignments for training the GNN. We use these sets only to ensure that our training sentences are different than test sentences
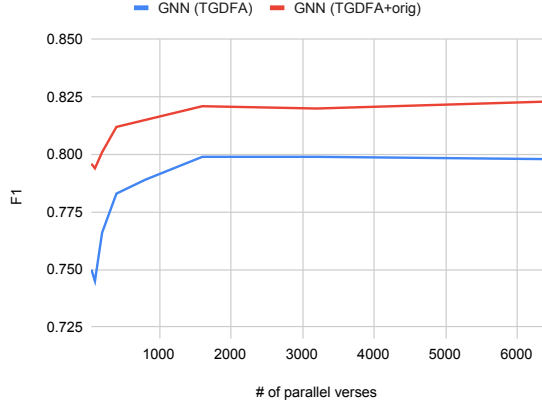
[7]https://networkx.org/
[8]pytorch-geometric.readthedocs.io
[9]https://radimrehurek.com/gensim/

Figure 3: $F_1$ of GNN (TGDFA) and GNN (TGDFA+orig) on Blinker as a function of train size

| | FIN-HEB | FIN-GRC | ENG-FRA |
|---|---|---|---|
| GNN (TGDFA) | 0.720 | 0.780 | 0.804 |
| ¬ language | -0.323 | -0.280 | -0.370 |
| ¬ position | -0.068 | -0.045 | -0.066 |
| ¬ centrality | -0.636 | -0.730 | -0.772 |
| ¬ community | -0.204 | -0.238 | -0.253 |
| ¬ word-embedding | -0.139 | -0.103 | -0.129 |
| GNN (TGDFA+orig) | 0.651 | 0.758 | 0.824 |
| ¬ language | -0.238 | -0.077 | -0.162 |
| ¬ position | -0.088 | +0.029 | -0.032 |
| ¬ centrality | -0.556 | -0.530 | -0.617 |
| ¬ community | -0.156 | -0.039 | -0.083 |
| ¬ word-embedding | -0.135 | +0.002 | -0.058 |

Table 3: $F_1$ for GNNs and $\Delta F_1$ for five ablations

| Model | Yoruba YTB |
|---|---|
| Unsupervised (Kondratyuk and Straka, 2019) | 50.86 |
| Eflomal Inter - eng | 43.45 |
| Eflomal GDFA - eng | 55.13 |
| Eflomal Inter - majority | 54.13 |
| Eflomal GDFA - majority | 60.27 |
| GNN (TGDFA) - majority | **65.74** |
| GNN (TGDFA+orig) - majority | 64.55 |

Table 4: POS tagging with annotation projection for Yoruba. Apart from "Unsupervised", all lines show a sequence tagger trained on pseudo-labels induced by word alignments. GNN-based pseudo-labels outperform prior work by 5%.

GNN (TGDFA) and GNN (TGDFA+orig) based on the characteristics of the desired alignments.

### 6.1.1 Effect of Training Set Size

To investigate the effect of training set size, we train the GNN on subsets of our training data with increasing sizes. Figure 3 shows results. Performance improves fast until around 2,000 verses; then it stays mostly constant. Indeed, using more than 6,400 samples does not change the performance at all. Therefore, in the other experiments we use 6,400 randomly sampled verses from the training set to train GNNs.

### 6.2 Ablation Experiments

To examine the importance of node features, we ablate language, position, centrality, community and word embedding features. Table 3 shows that removal of graph structural features drastically reduces performance. Community features and language information are also important. Removal of word position information and word embeddings – which store semantic information about words – has the least effect. Based on these results, it can be argued that the lexical information contained in the initial alignments and in the community features provides a very strong signal regarding words relatedness. The novel information that is crucial is about the overall graph structure which goes beyond local word associations, which are captured by word position and word embeddings.

### 6.2.1 Annotation Projection

Table 4 presents accuracies for POS tagging in Yoruba. Unsupervised baseline performance is 50.86%. Supervised training using pseudo-labels

mostly outperforms the unsupervised baseline. Projecting the majority POS labels to Yoruba improves over projecting English labels. Using the GNN model to project labels works best and outperforms Eflomal-GDFA-majority (the unsupervised baseline) by 5% (15%) absolute improvement.

## 7 Conclusion and Future Work

We introduced graph neural networks and community detection algorithms for multiparallel word alignment. By incorporating signals from diverse sources as node features, including community features, our GNN model outperformed the baselines and prior work, establishing new state-of-the-art results on three PBC gold standard datasets. We also showed that our GNN model improves downstream task performance in low-resource languages through annotation projection.

We have only used node features to provide signals to GNNs. In the future, other signals can be added in the form of edge features to further boost the performance.

# References

Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019. FLAIR: An easy-to-use framework for state-of-the-art NLP. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59, Minneapolis, Minnesota. Association for Computational Linguistics.

Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Tamer Alkhouli, Gabriel Bretschner, Jan-Thorsten Peter, Mohammed Hethnawi, Andreas Guta, and Hermann Ney. 2016. Alignment-based neural machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 1, Research Papers*, pages 54–65, Berlin, Germany. Association for Computational Linguistics.

Tamer Alkhouli and Hermann Ney. 2017. Biasing attention-based recurrent neural networks using external alignment information. In *Proceedings of the Second Conference on Machine Translation*, pages 108–117, Copenhagen, Denmark. 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.

Ehsaneddin Asgari and Hinrich Schütze. 2017. Past, present, future: A computational investigation of the typology of tense in 1000 languages. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 113–124, Copenhagen, Denmark. Association for Computational Linguistics.

Necip Fazil Ayan and Bonnie J. Dorr. 2006. A maximum entropy approach to combining word alignments. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 96–103, New York City, USA. Association for Computational Linguistics.

Nicola Bertoldi, Madalina Barbaiani, Marcello Federico, and Roldano Cattoni. 2008. Phrase-based statistical machine translation with pivot languages. In *International Workshop on Spoken Language Translation (IWSLT) 2008*.

Paolo Boldi and Sebastiano Vigna. 2014. Axioms for centrality. *Internet Mathematics*, 10(3-4):222–262.

Ulrik Brandes. 2001. A faster algorithm for betweenness centrality. *Journal of mathematical sociology*, 25(2):163–177.

Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.

Yun Chen, Yang Liu, Guanhua Chen, Xin Jiang, and Qun Liu. 2020. Accurate word alignment induction from neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 566–576, Online. Association for Computational Linguistics.

Aaron Clauset, Mark EJ Newman, and Cristopher Moore. 2004. Finding community structure in very large networks. *Physical review E*, 70(6):066111.

Trevor Cohn and Mirella Lapata. 2007. Machine translation by triangulation: Making effective use of multi-parallel corpora. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 728–735, Prague, Czech Republic. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Gennaro Cordasco and Luisa Gargano. 2010. Community detection via semi-synchronous label propagation algorithms. In *2010 IEEE international workshop on: business applications of social network analysis (BASNA)*, pages 1–8. IEEE.

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

Zi-Yi Dou and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2112–2128, Online. Association for Computational Linguistics.

Philipp Dufter, Mengjie Zhao, Martin Schmitt, Alexander Fraser, and Hinrich Schütze. 2018. Embedding learning through multilingual concept induction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1:*

9

*Long Papers)*, pages 1520–1530, Melbourne, Australia. Association for Computational Linguistics.

Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. A simple, fast, and effective reparameterization of IBM model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.

Patrick Fernandes, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Structured neural summarization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Victoria Fossum and Steven Abney. 2005. Automatically inducing a part-of-speech tagger by projecting from multiple source languages across aligned corpora. In *Second International Joint Conference on Natural Language Processing: Full Papers*.

Linton C Freeman. 1978. Centrality in social networks conceptual clarification. *Social networks*, 1(3):215–239.

Michelle Girvan and Mark EJ Newman. 2002. Community structure in social and biological networks. *Proceedings of the national academy of sciences*, 99(12):7821–7826.

Takuo Hamaguchi, Hidekazu Oiwa, Masashi Shimbo, and Yuji Matsumoto. 2017. Knowledge transfer for out-of-knowledge-base entities : A graph neural network approach. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 1802–1808.

Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-Dong Zhang, and Meng Wang. 2020. *LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation*, page 639–648. Association for Computing Machinery, New York, NY, USA.

Lianzhe Huang, Xin Sun, Sujian Li, Linhao Zhang, and Houfeng Wang. 2020. Syntax-aware graph attention network for aspect-level sentiment classification. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 799–810, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Matthias Huck, Diana Dutka, and Alexander Fraser. 2019. Cross-lingual annotation projection is effective for neural part-of-speech tagging. In *Proceedings of the Sixth Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 223–233, Ann Arbor, Michigan. Association for Computational Linguistics.

Ayyoob Imani, Masoud Jalili Sabet, Lütfi Kerem Şenel, Philipp Dufter, François Yvon, and Hinrich Schütze. 2021. Graph algorithms for multiparallel word alignment.

Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. 2020. SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1627–1643, Online. Association for Computational Linguistics.

Thomas N Kipf and Max Welling. 2016. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*.

Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*.

Philipp Koehn, Amittai Axelrod, Alexandra Birch, Chris Callison-Burch, Miles Osborne, and David Talbot. 2005. Edinburgh system description for the 2005 iwslt speech translation evaluation. *International Workshop on Spoken Language Translation*.

Philipp Koehn, Franz J Och, and Daniel Marcu. 2003. Statistical phrase-based translation. Technical report, UNIVERSITY OF SOUTHERN CALIFORNIA MARINA DEL REY INFORMATION SCIENCES INST.

Dan Kondratyuk and Milan Straka. 2019. 75 languages, 1 model: Parsing Universal Dependencies universally. 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 2779–2795, Hong Kong, China. Association for Computational Linguistics.

Shankar Kumar, Franz J. Och, and Wolfgang Macherey. 2007. Improving word alignment with bridge languages. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 42–50, Prague, Czech Republic. Association for Computational Linguistics.

Guillaume Lample, Alexis Conneau, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. Word translation without parallel data. In *International Conference on Learning Representations*.

Omer Levy, Anders Søgaard, and Yoav Goldberg. 2017. A strong baseline for learning cross-lingual word embeddings from sentence alignments. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 765–774, Valencia, Spain. Association for Computational Linguistics.

William D. Lewis and Fei Xia. 2008. Automatically identifying computationally relevant typological features. In *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-II*.

10

Percy Liang, Ben Taskar, and Dan Klein. 2006. Alignment by agreement. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 104–111, New York City, USA. Association for Computational Linguistics.

Zhiyuan Liu and Jie Zhou. 2020. Introduction to graph neural networks. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 14(2):1–127.

Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Ying Luo and Hai Zhao. 2020. Bipartite flat-graph network for nested named entity recognition. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6408–6418, Online. Association for Computational Linguistics.

Kelly Marchisio, Conghao Xiong, and Philipp Koehn. 2021. Embedding-enhanced giza++: Improving alignment in low-and high-resource scenarios using embedding space geometry. *arXiv preprint arXiv:2104.08721*.

Thomas Mayer and Michael Cysouw. 2012. Language comparison through sparse multilingual word alignment. In *Proceedings of the EACL 2012 Joint Workshop of LINGVIS & UNCLH*, pages 54–62, Avignon, France. Association for Computational Linguistics.

I. Dan Melamed. 1998. Manual annotation of translational equivalence: The blinker project. *CoRR*, cmp-lg/9805005.

Preslav Nakov and Hwee Tou Ng. 2012. Improving statistical machine translation for a resource-poor language using related resource-rich languages. *Journal of Artificial Intelligence Research*, 44:179–222.

Mark EJ Newman. 2001. Scientific collaboration networks. ii. shortest paths, weighted networks, and centrality. *Physical review E*, 64(1):016132.

Mark EJ Newman and Michelle Girvan. 2004. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113.

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. Universal Dependencies v2: An evergrowing multilingual treebank collection. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France. European Language Resources Association.

Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

Robert Östling. 2014. Bayesian word alignment for massively parallel texts. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, volume 2: Short Papers*, pages 123–127, Gothenburg, Sweden. Association for Computational Linguistics.

Robert Östling. 2015. Word order typology through multilingual word alignment. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 205–211, Beijing, China. Association for Computational Linguistics.

Robert Östling and Jörg Tiedemann. 2016. Efficient word alignment with Markov Chain Monte Carlo. *The Prague Bulletin of Mathematical Linguistics*, 106(1).

Liangming Pan, Yuxi Xie, Yansong Feng, Tat-Seng Chua, and Min-Yen Kan. 2020. Semantic graphs for generating deep questions. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1463–1475, Online. Association for Computational Linguistics.

Slav Petrov, Dipanjan Das, and Ryan McDonald. 2012. A universal part-of-speech tagset. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2089–2096, Istanbul, Turkey. European Language Resources Association (ELRA).

Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin Riedmiller, Raia Hadsell, and Peter Battaglia. 2018. Graph networks as learnable physics engines for inference and control. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4470–4479. PMLR.

Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2009. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80.

Haoyue Shi, Luke Zettlemoyer, and Sida I Wang. 2021. Bilingual lexicon induction via unsupervised bitext construction and word alignment. *arXiv preprint arXiv:2101.00148*.

Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In *International Conference on Learning Representations*.

11

Guillaume Wisniewski, Nicolas Pécheux, Souhir Gahbiche-Braham, and François Yvon. 2014. Cross-lingual part-of-speech tagging through ambiguous learning. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1779–1785, Doha, Qatar. Association for Computational Linguistics.

Shijie Wu and Mark Dredze. 2020. Do explicit alignments robustly improve multilingual encoders? *arXiv preprint arXiv:2010.02537*.

Yongji Wu, Defu Lian, Yiheng Xu, Le Wu, and Enhong Chen. 2020. Graph convolutional networks with markov random field reasoning for social spammer detection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01):1054–1061.

David Yarowsky and Grace Ngai. 2001. Inducing multilingual POS taggers and NP bracketers via robust projection across aligned corpora. In *Second Meeting of the North American Chapter of the Association for Computational Linguistics*.

Anssi Yli-Jyrä, Josi Purhonen, Matti Liljeqvist, Arto Antturi, Pekka Nieminen, Kari M. Räntilä, and Valtter Luoto. 2020. HELFI: a Hebrew-Greek-Finnish parallel Bible corpus with cross-lingual morpheme alignment. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4229–4236, Marseille, France. European Language Resources Association.

Thomas Zenkel, Joern Wuebker, and John DeNero. 2020. End-to-end neural word alignment outperforms GIZA++. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1605–1617, Online. Association for Computational Linguistics.