

# Relational Graph Attention Networks for Syntax Encoding in Zero-shot Cross-lingual Semantic Role Labeling

Anonymous ACL submission

## Abstract

Recent models in cross-lingual semantic role labeling (SRL) rely heavily on BiLSTMs, a derivation of RNNs, as their main encoders. However, a previous study in dependency parsing has shown that RNN-based cross-lingual models are ineffective in distant languages. Therefore, we propose graph neural networks (GNNs) built on dependency trees to replace BiLSTMs' role as the encoder for cross-lingual models. We hypothesize that encoding sentences based on their dependency trees helps cross-lingual SRL models achieve better generalization. Through a simple encoder-decoder architecture, we compare various GNNs, i.e., gated graph convolutional networks (GGCNs), graph attention networks (GATs), two-attention relational GATs (2ATT-GATs), and modified self-attention networks from Transformer (SATs). We focus on a zero-shot setting and evaluate the models in 23 languages available in Universal Proposition Bank. The evaluation shows that 2ATT-GATs outperform other GNNs. Moreover, comparisons against BiLSTM-based models show that 2ATT-GATs are more effective for building cross-lingual SRL models, especially in languages with different word orders.

## 1 Introduction

Semantic role labeling (SRL) is a task to determine the predicates of a sentence and argument roles for each corresponding predicate, as shown in Figure 1. SRL supports many natural language processing (NLP) tasks, e.g., information extraction (Christensen et al., 2010), abstractive summarization (Khan et al., 2015), and machine translation (Rapp, 2022). However, SRL resource availability is low, hindering the performance of other NLP tasks in diverse languages. Cross-lingual SRL models try to solve this problem by training the models in resource-rich languages and transferring the models to resource-poor languages.

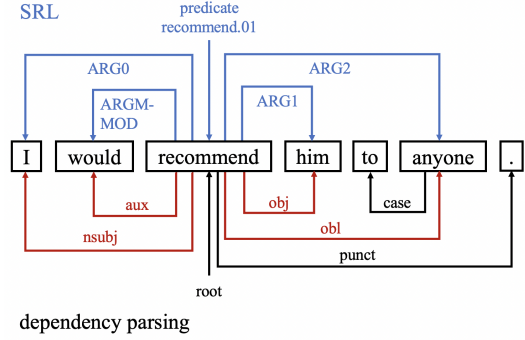


Figure 1: An example of the SRL task (top) and the dependency parsing task (bottom) applied to a sentence taken from UPB. The red color indicates path intersections in both tasks.

A study in cross-lingual dependency parsing (Ahmad et al., 2019) proves that recurrent neural network (RNN)-based cross-lingual models are sensitive to word orders, making them unable to transfer effectively to distant languages. Despite the developments in cross-lingual SRL, recent models still heavily rely on BiLSTMs, the derivation of RNNs, as their encoders, e.g., Fei et al. (2020), Cai and Lapata (2020), and Conia et al. (2021).

We propose to apply GNNs over universal dependency trees provided by Universal Dependencies (UD) (Nivre et al., 2016b, 2020) as the encoder for cross-lingual models. We hypothesize that encoding sentences based on their dependency trees makes cross-lingual SRL models generalize better. We provide two main reasons: (1) Many predicate-argument paths and argument roles in SRL intersect with dependency paths and dependency relations in dependency parsing (Marcheggiani and Titov, 2017), as shown in Figure 1. (2) Universal dependency tree, representing a sentence's grammatical structure in a language-universal scheme, is a more generalized representation across languages than word sequences.

We experiment on various GNNs as an encoder to extract language-universal information from

dependency trees, including GGCNs (Marcheggiani and Titov, 2017), graph attention networks (GATs) (Veličković et al., 2017), modified relational GATs (Wang et al., 2020) (2ATT-GATs), and modified self-attention networks from Transformer (Shaw et al., 2018) (SATs). We apply GGCNs and 2ATT-GATs since they have been proven useful to encode dependency trees in monolingual SRL (Marcheggiani and Titov, 2017) and sentiment analysis (Wang et al., 2020), respectively. We compare 2ATT-GATs with GATs that treat the dependency tree as unlabeled. In addition, we employ SATs as it is taken from self-attention in the popular Transformer (Vaswani et al., 2017). Furthermore, we compare the best GNN-based model with BiLSTMs-based models to show the effectiveness of the GNN-based model in cross-lingual SRL.

SRL consists of four steps, i.e., predicate detection, predicate sense disambiguation, argument detection, and argument labeling. Following previous work in cross-lingual SRL (Fei et al., 2020), we focus on argument detection and argument labeling in the dependency-based SRL. We conduct experiments in a zero-shot setting to find the most transferable networks across languages.

We train and evaluate the models in seven and 23 languages provided by Universal Proposition Bank (UPB) v1.0 and UPB v2.0, respectively. Throughout the paper, we show that in cross-lingual SRL:

1. 2ATT-GATs outperform other GNN-based SRL models, indicating that 2ATT-GATs transfer more effectively across languages than other GNNs.
2. 2ATT-GATs perform better than BiLSTM-based SRL models, even when built on inaccurate dependency trees, especially in target languages with different word orders than the source language.

## 2 Background

### 2.1 Universal Proposition Bank

Universal Proposition Bank (UPB) is a dataset containing SRL annotations across languages. UPB v1.0 (Akbik et al., 2015, 2016) provides SRL annotations for nine treebanks and eight languages, including English and other seven languages shown at the left side of Table 1. UPB v2.0 (Jindal et al., 2022) provides SRL annotations for 43 treebanks and 23 languages, shown in Table 1. UPB is annotated semi-automatically through filtered annota-

v1.0 and v2.0	v2.0	
Chinese (ZH)	Czech (CS)	Dutch (NL)
Finnish (FI)	Greek (EL)	Polish (PL)
Italian (IT)	Korean (KO)	Telugu (TE)
Spanish (ES)	Romanian (RO)	Indonesian (ID)
French (FR)	Hindi (HI)	Japanese (JA)
German (DE)	Marathi (MR)	Russian (RU)
Portuguese (PT)	Tamil (TA)	Ukrainian (UK)
	Hungarian (HU)	Vietnamese (VI)

Table 1: List of target languages available in UPB v1.0 and UPB v2.0.

tion projection and bootstrap training (Akbik et al., 2015). UPB v2.0 has significantly improved over UPB v1.0 regarding SRL annotation quality, language scope, and availability of span-based SRL annotations (Jindal et al., 2022).

### 2.2 Universal Dependencies

Universal Dependencies is a dataset that contains consistent syntactic annotations across languages, i.e., part-of-speech (POS) tags, morphological features, and syntactic dependencies. UD v1 (Nivre et al., 2016b) and UD v2 (Nivre et al., 2020) have different annotation schemes<sup>1</sup> in word segmentation, pos tags, morphological features, and syntactic relations. UD v1 and UD v2 have 40 and 37 universal dependencies relations, respectively. UPB v1.0 and UPB v2.0 that we use throughout the experiments are annotated based on UD v1.4 (Nivre et al., 2016a) and UD v2.9 (Zeman et al., 2021), respectively.

### 2.3 Dependency-based Semantic Role Labeling

Dependency-based SRL labels the argument heads for each predicate in a sentence based on its dependency tree. For example, in Figure 1, the phrase "to anyone" is the argument "ARG2" of the predicate "recommend". Based on the dependency tree, "anyone" is the head of the phrase "to anyone". Therefore, dependency-based SRL annotates the edge that connects "recommend" and "anyone" with "ARG2" label.

### 2.4 Related Works

Existing models in cross-lingual SRL rely on BiLSTMs as their main encoders despite the findings in Ahmad et al. (2019), which prove that RNN-based models perform ineffectively in distant languages for dependency parsing task. Fei

<sup>1</sup><https://universaldependencies.org/v2/summary.html>

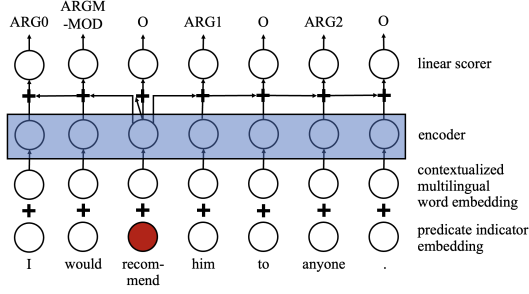


Figure 2: The architecture applied to a sentence with "recommend" as the predicate.

et al. (2020) propose parameter generation network (PGN)-BiLSTMs to build a cross-lingual SRL model. Cai and Lapata (2020) propose BiLSTM-based models as semantic role labeler and compressor in their architecture. Conia et al. (2021) propose BiLSTM-based universal sentence encoder and BiLSTM-based universal predicate-argument encoder to encode predicate-related and predicate-argument information.

GNNs have been used to encode dependency trees in monolingual SRL and aspect-based sentiment analysis (ABSA). In monolingual SRL, Marcheggiani and Titov (2017) employ GGCNs on top of BiLSTMs to incorporate dependency trees as graphs. In ABSA, Wang et al. (2020) and Jiang et al. (2021) employ relational GATs (R-GATs) and attention-based relational GCNs (ARGCNs), respectively. They apply GNNs on top of modified dependency trees to establish direct connections between aspects and their corresponding words.

### 3 Model

We apply a common encoder-decoder architecture for comparing various GNN-based and BiLSTM-based cross-lingual SRL models. The architecture consists of an input layer (i.e., predicate indicator embedding and contextualized multilingual word embedding), an encoder, and a decoder (i.e., linear scorer), as shown in Figure 2.

#### 3.1 Input Layer

For each word in a sentence, we concatenate predicate indicator embedding,  $\vec{p}_i$ , and contextualized multilingual word embedding,  $\vec{c}_i$ , to produce the final word representation,  $\vec{h}_i$ , as shown in Equation 1.

$$\vec{h}_i = [\vec{p}_i || \vec{c}_i] \quad (1)$$

Predicate indicator embedding is an embedding

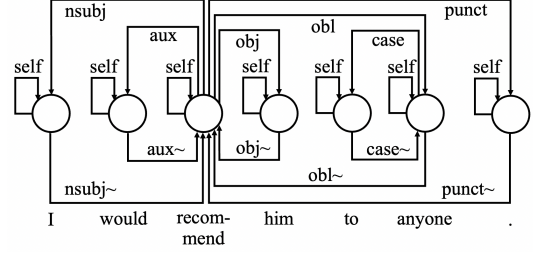


Figure 3: Dependency graph of a sentence converted from its dependency tree.

that represents whether a word is a predicate or not (Fei et al., 2020). We compare contextualized multilingual word embedding from two language models, i.e., multilingual BERT (mBERT) (Devlin et al., 2019) and XLM-RoBERTa (XLM-R) (Conneau et al., 2020). Both mBERT and XLM-R produce embedding at a subword level according to WordPiece tokenization (Wu et al., 2016). We take the left-most subword as the representation for the corresponding word following Wang et al. (2019). To generate the word embedding,  $\vec{c}_i$ , we adopt the method proposed by existing cross-lingual SRL. Conia et al. (2021) concatenate the last four hidden layers obtained from corresponding pre-trained language models. Given the result of concatenation,  $\vec{v}_i$ , they apply a feed-forward neural network (FNN) and Swish activation function (Ramachandran et al., 2018), as shown in Equation 2. They also apply a dropout after the activation function.

$$\vec{c}_i = \text{Swish}(\mathbf{W}\vec{v}_i + \mathbf{b}) \quad (2)$$

#### 3.2 Encoder

We experiment with various GNN-based and BiLSTM-based encoders. GNNs that we apply are graph attention networks (GATs), two-attention relational GATs (2ATT-GATs), gated graph convolutional networks (GGCNs), and modified self-attention networks from Transformer (SATs).

In GNN-based models, we encode a sentence by forming a dependency graph based on its dependency tree, which consists of dependency arcs and dependency relations. We follow the method proposed by Marcheggiani and Titov (2017). They convert a dependency tree to a graph by adding edges that flow in the opposite direction of the original dependency arcs and edges that flow from nodes to themselves. Figure 3 displays the dependency graph of the dependency tree at the bottom of Figure 1.

### 3.2.1 Graph Attention Networks

Given a graph that consists of nodes,  $i$ , and edges,  $e$ , GATs (Veličković et al., 2017) update each node representation,  $\vec{h}_i$ , according to its neighbor node representations, using multi-head attention mechanism that employs  $K$  heads. GATs utilize an attention weight,  $\alpha$ , to measure the contribution of neighbor node representations when updating the corresponding node representation,  $h_i$ . The attention weight,  $\alpha_{ij}$ , for the edge that connects node  $i$  and node  $j$ , is calculated by taking the dot-product between a weight vector,  $\vec{\alpha}$ , with the concatenation of linearly transformed  $\vec{h}_i$  and  $\vec{h}_j$ . The result of the dot-product is passed to a LeakyReLU activation function, LR, and softmax function, SM, as shown in Equation 3.

$$\alpha_{ij}^k = \text{SM}_j(\text{LR}(\vec{\alpha}^k \cdot [\mathbf{W}^k \vec{h}_i || \mathbf{W}^k \vec{h}_j])) \quad (3)$$

### 3.2.2 Two-attention Relational Graph Attention Networks

Since GATs treat a graph as unlabeled (Veličković et al., 2017), Wang et al. (2020) modify GATs to use two attention weights. The first attention weight contains node representations, while the second includes dependency-type representations taken from dependency types that represent dependency arcs and relations in the dependency trees. Instead of using different equations for both attentions as in Wang et al. (2020), we find that applying the FNNs (Veličković et al., 2017) for calculating both attentions works best for our task. Therefore, we explain the modification in this section.

Equation 3 shows how we calculate the first attention weight,  $\alpha$ . To calculate the second attention weight,  $\beta$ , we slightly modify Equation 3 to encode dependency-type representation,  $\vec{r}_{ij}$ , as shown in Equation 4.

$$\beta_{ij}^k = \text{SM}_j(\text{LR}(\vec{\alpha}^k \cdot \mathbf{W}^k \vec{r}_{ij})) \quad (4)$$

We obtain node representations from attention weight  $\alpha$  and  $\beta$ , as shown in Equation 5 and Equation 6.  $\mathcal{N}_i$  indicates the set of neighbor nodes of node  $i$ .  $\vec{h}_{\alpha,i}^l$  and  $\vec{h}_{\beta,i}^l$  are the node representations in layer  $l$  obtained from attention weight  $\alpha$  and  $\beta$ , respectively.

$$\vec{h}_{\alpha,i}^l = \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j^{l-1} \quad (5)$$

$$\vec{h}_{\beta,i}^l = \sum_{j \in \mathcal{N}_i} \beta_{ij}^k \mathbf{W}^k \vec{h}_j^{l-1} \quad (6)$$

Finally, we calculate the node representation in layer  $l$ ,  $\vec{h}_i^l$ , by applying an FNN to the concatenation of node representation  $\vec{h}_{\alpha,i}^l$  and  $\vec{h}_{\beta,i}^l$ . We concatenate the node representation from each attention head,  $k$  (Equation 7), except in the final layer where we take the average of node representations (Equation 8).  $L$  indicates the number of layers, while  $\mathbf{W}$  and  $\mathbf{b}$  indicate weight matrix and bias in the FNN. We optionally apply an activation function,  $\sigma$ .

$$\vec{h}_i^l = \sigma(\frac{1}{K} \sum_{k=1}^K (\mathbf{W}^k [\vec{h}_{\alpha,i}^l || \vec{h}_{\beta,i}^l] + \mathbf{b}^k)), l < L \quad (7)$$

$$\vec{h}_i^l = \sigma(\frac{1}{K} \sum_{k=1}^K (\mathbf{W}^k [\vec{h}_{\alpha,i}^l || \vec{h}_{\beta,i}^l] + \mathbf{b}^k)), l = L \quad (8)$$

### 3.2.3 Gated Graph Convolutional Networks

Marcheggiani and Titov (2017) propose GGCNs to encode syntactic features from dependency trees in monolingual SRL models. Equation 9 shows how to calculate node representation,  $\vec{h}_i^l$ , in layer  $l$ . GGCNs separate the weight matrices,  $\mathbf{W}_{dir}$ , according to the direction of dependency arcs, i.e., original direction, opposite direction, and self-direction. Meanwhile,  $\mathbf{b}_{rel}$  represents the dependency relation. Unlike regular GCNs (Kipf and Welling, 2017), GGCNs employ a scalar gate,  $g_{ij}$ , to measure the importance of neighbor node representations when updating the corresponding node representation.

$$\vec{h}_i^l = \text{ReLU}(\sum_{j \in \mathcal{N}_i} g_{ij} (\mathbf{W}_{dir_{ij}} \vec{h}_j^{l-1} + \mathbf{b}_{rel_{ij}})) \quad (9)$$

### 3.2.4 Modified Self-attention Networks

Self-attention with relative position representations (Shaw et al., 2018) is an extension of the self-attention in Transformer (Vaswani et al., 2017) to include edges that represent relative positions between words. The self-attention can be categorized as GNNs since they receive input as a graph consisting of nodes and edges. Referring to the paper (Shaw et al., 2018), we modify edge representations from representing relative positions to representing dependency types. Instead of a fully connected graph as input, SATs take a graph formed over the sentence's dependency tree, as shown in Figure 3.

### 3.3 Decoder

We apply a linear scorer as the decoder. We concatenate node representation for each word,  $\vec{h}_i$ , with the predicate node representation,  $\vec{h}_{p,i}$ . Predicate node representation is taken from the node representation,  $\vec{h}_i$ , of the sentence’s predicate. After that, we apply an FNN to produce the final node representation with an embedding size equal to the number of arguments,  $n$ . We then apply a softmax function, SM, to produce the probability for each label,  $z$ , as shown in Equation 10. We train the model to minimize the cross-entropy loss.

$$P(z) = \text{SM}(\mathbf{W}[\vec{h}_i || \vec{h}_{p,i}] + \mathbf{b})_z, z \in [1, n] \quad (10)$$

## 4 Experiment Results

### 4.1 Datasets

We conduct experiments using datasets from UPB v1.0 and UPB v2.0. The dataset distribution can be found in Appendix A.1. In UPB v1.0, UP\_English-EWT is annotated based on UD v2, while the other languages are annotated based on UD v1.4. Since UP\_English-EWT in UPB v1.0 is annotated based on UD v2, we use its SRL annotations to construct a cross-lingual SRL model for target languages in UPB v2.0. To construct a cross-lingual SRL model for UPB v1.0, we convert syntactic annotations of UP\_English-EWT (UPB v1.0) to UD v1.4 using the script available at Zhang et al. (2021)<sup>2</sup>. We merge annotations from English Web Treebank (EWT) (Bies et al., 2012), PropBank v3 (Kingsbury and Palmer, 2002; Palmer et al., 2005; Gildea and Palmer, 2002), and UD v1.4 (Nivre et al., 2016a).

### 4.2 Settings

We focus on conducting experiments in a zero-shot setting to examine the model’s transferability across languages. We train the model in English and evaluate the model in seven languages and 23 languages from UPB v1.0 and UPB v2.0, respectively, as shown in Table 1. Furthermore, we compare the performance of the cross-lingual model against the monolingual model to illustrate the generalization achieved by each model.

We use gold and predicted dependency trees for model evaluation. We train dependency parsers from scratch using Stanza (Qi et al., 2020) to produce predicted dependency trees for languages in

UPB v1.0. For languages in UPB v2.0, we use pre-trained models provided by Stanza<sup>3</sup>. Appendix B.1 shows each dependency parser’s unlabeled attachment score (UAS) and labeled attachment score (LAS).

Following Conia et al. (2021), we train the model for 30 epochs by increasing the learning rate linearly for 1 epoch and decreasing linearly for 15 epochs. We use AdamW (Loshchilov and Hutter, 2019) as the optimizer. We run the experiments five times and report the average F1 scores and standard errors. We choose the best epoch from each experiment based on the F1 score of the English validation set. Appendix B.2 explains the implementation details and hyperparameters we use throughout the experiments. We choose most hyperparameter values based on previous studies.

### 4.3 Comparison Among GNN-based SRL Models

We compare the performance of various GNN-based SRL models, i.e., GGCNs, SATs, GATs, and 2ATT-GATs, by substituting the encoder part in the architecture (Figure 2) with these GNNs. We stack three layers of GNNs and freeze mBERT as our contextualized word embedding to solely observe the influence of GNNs.

The left side of Table 2 shows the F1 score of each model evaluated using the test sets in UPB v1.0 and UPB v2.0 with predicted dependency trees. The results with gold dependency trees are available in Appendix C. In the monolingual setting, as indicated by EN, SATs achieve the highest F1 score for both datasets. Meanwhile, in the cross-lingual setting, as shown in AVG, GGCNs perform the best in UPB v1.0, while 2ATT-GATs perform the best in UPB v2.0.

Although GGCNs outperform 2ATT-GATs in UPB v1.0, if we extract the F1 scores of languages available in UPB v1.0 from UPB v2.0, i.e., FI, IT, ES, FR, DE, PT, and ZH, 2ATT-GATs still outperform GGCNs in all these languages. As we know, UPB v1.0 and UPB v2.0 are annotated semi-automatically (Akbik et al., 2015). UPB v2.0 has significantly improved over UPB v1.0 regarding SRL annotations quality (Jindal et al., 2022). Therefore, we believe the evaluation against UPB v2.0 is more crucial.

<sup>3</sup>Except for UP\_Japanese-GSDLUW and UP\_French-Rhapsodie, where we train from scratch since the models are unavailable.

<sup>2</sup><https://github.com/zsformlp/zmsp>

Lang	GGCNs	SATs	GATs	2ATT-GATs	BiLSTM+ 2ATT-GATs	BiLSTM	3BiLSTMs
UPB v1.0							
EN	80.50 $\pm$ 0.09	<b>80.89<math>\pm</math>0.04</b>	79.77 $\pm$ 0.05	80.04 $\pm$ 0.05	82.19 $\pm$ 0.10	79.72 $\pm$ 0.05	80.23 $\pm$ 0.08
FI	<b>54.10<math>\pm</math>0.09</b>	53.95 $\pm$ 0.17	48.95 $\pm$ 0.32	53.17 $\pm$ 0.10	52.83 $\pm$ 0.15	39.85 $\pm$ 0.19	40.18 $\pm$ 0.22
IT	<b>57.57<math>\pm</math>0.13</b>	57.33 $\pm$ 0.14	54.91 $\pm$ 0.26	<u>57.12<math>\pm</math>0.08</u>	56.45 $\pm$ 0.10	47.20 $\pm$ 0.19	46.74 $\pm$ 0.30
ES	54.32 $\pm$ 0.11	<b>54.88<math>\pm</math>0.17</b>	50.44 $\pm$ 0.15	<u>54.69<math>\pm</math>0.12</u>	52.41 $\pm$ 0.16	41.77 $\pm$ 0.15	41.52 $\pm$ 0.05
FR	<b>46.92<math>\pm</math>0.12</b>	46.18 $\pm$ 0.12	44.60 $\pm$ 0.28	46.28 $\pm$ 0.20	45.58 $\pm$ 0.23	40.64 $\pm$ 0.32	40.70 $\pm$ 0.31
DE	58.75 $\pm$ 0.21	58.65 $\pm$ 0.05	57.06 $\pm$ 0.10	<b>58.80<math>\pm</math>0.07</b>	56.18 $\pm$ 0.13	41.45 $\pm$ 0.24	41.99 $\pm$ 0.12
PT	<b>53.60<math>\pm</math>0.04</b>	53.12 $\pm$ 0.13	52.50 $\pm$ 0.13	<u>53.33<math>\pm</math>0.07</u>	52.86 $\pm$ 0.16	44.38 $\pm$ 0.24	44.10 $\pm$ 0.12
ZH	37.83 $\pm$ 0.14	37.54 $\pm$ 0.26	37.07 $\pm$ 0.11	<b>38.42<math>\pm</math>0.16</b>	40.01 $\pm$ 0.24	32.87 $\pm$ 0.37	32.53 $\pm$ 0.23
AVG	<b>51.87<math>\pm</math>0.05</b>	51.67 $\pm$ 0.09	49.36 $\pm$ 0.11	51.69 $\pm$ 0.07	50.90 $\pm$ 0.07	41.17 $\pm$ 0.13	41.11 $\pm$ 0.11
UPB v2.0							
EN	80.20 $\pm$ 0.08	<b>80.71<math>\pm</math>0.06</b>	79.63 $\pm$ 0.08	79.65 $\pm$ 0.02	81.99 $\pm$ 0.10	79.51 $\pm$ 0.08	80.08 $\pm$ 0.08
CS	56.60 $\pm$ 0.20	56.86 $\pm$ 0.12	55.93 $\pm$ 0.06	<b>57.79<math>\pm</math>0.11</b>	55.19 $\pm$ 0.23	47.13 $\pm$ 0.06	47.50 $\pm$ 0.13
EL	58.39 $\pm$ 0.30	58.56 $\pm$ 0.52	59.68 $\pm$ 0.27	<u>60.69<math>\pm</math>0.36</u>	59.48 $\pm$ 0.21	56.05 $\pm$ 0.16	55.53 $\pm$ 0.27
KO	39.31 $\pm$ 0.24	38.26 $\pm$ 0.77	41.25 $\pm$ 0.25	<b>42.20<math>\pm</math>0.45</b>	35.01 $\pm$ 0.54	27.75 $\pm$ 0.29	27.84 $\pm$ 0.99
RO	52.43 $\pm$ 0.07	52.65 $\pm$ 0.19	<b>53.51<math>\pm</math>0.20</b>	53.24 $\pm$ 0.17	54.41 $\pm$ 0.09	48.71 $\pm$ 0.30	47.83 $\pm$ 0.31
HI	47.26 $\pm$ 0.07	47.93 $\pm$ 0.08	44.80 $\pm$ 0.22	<b>48.40<math>\pm</math>0.14</b>	42.69 $\pm$ 0.44	27.49 $\pm$ 0.84	27.30 $\pm$ 0.44
MR	35.43 $\pm$ 0.59	37.86 $\pm$ 0.86	<b>39.53<math>\pm</math>1.70</b>	38.94 $\pm$ 1.01	36.17 $\pm$ 0.96	27.22 $\pm$ 1.89	26.31 $\pm$ 1.60
TA	31.64 $\pm$ 0.27	33.74 $\pm$ 0.54	<b>34.58<math>\pm</math>0.66</b>	34.25 $\pm$ 0.40	30.70 $\pm$ 0.75	23.64 $\pm$ 0.84	22.20 $\pm$ 0.88
HU	50.08 $\pm$ 0.27	50.90 $\pm$ 0.26	48.27 $\pm$ 0.14	<b>51.00<math>\pm</math>0.12</b>	46.65 $\pm$ 0.21	39.05 $\pm$ 0.29	39.52 $\pm$ 0.39
PL	57.73 $\pm$ 0.12	57.54 $\pm$ 0.08	57.57 $\pm$ 0.14	<b>59.85<math>\pm</math>0.16</b>	57.01 $\pm$ 0.33	51.65 $\pm$ 0.11	51.19 $\pm$ 0.33
TE	46.27 $\pm$ 0.64	45.09 $\pm$ 0.55	42.66 $\pm$ 0.72	<b>47.01<math>\pm</math>0.82</b>	40.12 $\pm$ 0.99	35.97 $\pm$ 1.30	32.84 $\pm$ 0.55
NL	62.94 $\pm$ 0.20	62.81 $\pm$ 0.21	63.00 $\pm$ 0.19	<b>63.39<math>\pm</math>0.15</b>	58.39 $\pm$ 0.28	51.95 $\pm$ 0.36	52.42 $\pm$ 0.10
ID	53.35 $\pm$ 0.34	53.43 $\pm$ 0.23	<b>58.82<math>\pm</math>0.25</b>	54.24 $\pm$ 0.10	<u>56.80<math>\pm</math>0.32</u>	56.55 $\pm$ 0.34	55.00 $\pm$ 0.35
JA	36.07 $\pm$ 0.15	36.12 $\pm$ 0.39	36.30 $\pm$ 0.51	<b>36.82<math>\pm</math>0.11</b>	33.39 $\pm$ 1.08	23.76 $\pm$ 1.07	21.40 $\pm$ 1.07
RU	57.63 $\pm$ 0.22	58.69 $\pm$ 0.35	<b>59.43<math>\pm</math>0.27</b>	58.88 $\pm$ 0.23	<u>58.95<math>\pm</math>0.29</u>	55.49 $\pm$ 0.24	55.48 $\pm$ 0.23
UK	57.39 $\pm$ 0.04	58.45 $\pm$ 0.16	57.18 $\pm$ 0.26	<b>58.62<math>\pm</math>0.12</b>	57.59 $\pm$ 0.27	52.80 $\pm$ 0.39	52.98 $\pm$ 0.17
ZH	42.88 $\pm$ 0.33	44.13 $\pm$ 0.31	<b>45.55<math>\pm</math>0.33</b>	44.77 $\pm$ 0.16	47.73 $\pm$ 0.19	<u>48.54<math>\pm</math>0.46</u>	47.95 $\pm$ 0.34
VI	25.78 $\pm$ 0.05	26.11 $\pm$ 0.14	<b>27.48<math>\pm</math>0.03</b>	26.51 $\pm$ 0.14	27.32 $\pm$ 0.08	<u>29.39<math>\pm</math>0.29</u>	27.34 $\pm$ 0.36
FI	53.54 $\pm$ 0.11	54.09 $\pm$ 0.13	52.92 $\pm$ 0.13	<b>54.46<math>\pm</math>0.04</b>	54.37 $\pm$ 0.17	51.54 $\pm$ 0.13	51.83 $\pm$ 0.23
IT	57.31 $\pm$ 0.10	57.81 $\pm$ 0.15	58.00 $\pm$ 0.12	<b>58.52<math>\pm</math>0.07</b>	58.14 $\pm$ 0.19	55.23 $\pm$ 0.16	54.97 $\pm$ 0.12
ES	54.10 $\pm$ 0.06	54.85 $\pm$ 0.10	<b>55.82<math>\pm</math>0.06</b>	<u>55.68<math>\pm</math>0.09</u>	55.02 $\pm$ 0.21	51.02 $\pm$ 0.13	50.51 $\pm$ 0.06
FR	61.55 $\pm$ 0.14	61.82 $\pm$ 0.10	62.05 $\pm$ 0.19	<b>62.43<math>\pm</math>0.04</b>	60.60 $\pm$ 0.25	59.81 $\pm$ 0.15	59.88 $\pm$ 0.26
DE	58.34 $\pm$ 0.06	58.80 $\pm$ 0.14	59.34 $\pm$ 0.08	<b>59.39<math>\pm</math>0.07</b>	53.92 $\pm$ 0.31	41.40 $\pm$ 0.38	43.10 $\pm$ 0.23
PT	65.55 $\pm$ 0.08	65.82 $\pm$ 0.10	<b>66.37<math>\pm</math>0.07</b>	<u>66.09<math>\pm</math>0.07</u>	65.85 $\pm$ 0.11	64.05 $\pm$ 0.17	63.71 $\pm$ 0.10
AVG	50.50 $\pm$ 0.04	50.97 $\pm$ 0.13	51.31 $\pm$ 0.08	<b>51.88<math>\pm</math>0.09</b>	49.81 $\pm$ 0.16	44.62 $\pm$ 0.11	44.11 $\pm$ 0.23

Table 2: F1 scores (%) in UPB v1.0 and UPB v2.0 test sets with predicted dependency trees (frozen mBERT). The bold score in each language indicates the highest F1 score among GNN-based models, i.e., GGCNs, SATs, GATs, and 2ATT-GATs. The underlined score in each language indicates the highest F1 score among 2ATT-GATs, BiLSTM+2ATT-GATs, BiLSTM, and 3BiLSTMs. AVG indicates the average F1 score of each model for all languages except English.

In UPB v2.0, 2ATT-GATs outperform other GNNs, indicating that 2ATT-GATs can generalize better across languages. However, GATs sometimes perform better in languages that have relatively low LAS, i.e., MR, TA, ID, VI, and RU<sup>4</sup>. We conjecture that GATs are more robust to inaccurate dependency relation labels because GATs treat dependency trees as unlabeled. Clearer evidence can be shown in ID, where the F1 scores decrease significantly when we replace gold with predicted dependency trees except for GATs<sup>5</sup> be-

<sup>4</sup>UAS/LAS of each language: 79.85/70.63 (MR), 80.89/72.30 (TA), 87.31/77.33 (ID), 77.58/74.16 (VI), 84.42/81.41 (RU Taiga), and 90.44/87.2 (RU GSD).

<sup>5</sup>F1 scores of ID using gold trees: 57.18 $\pm$ 0.26 (GGCNs), 57.05 $\pm$ 0.20 (SATs), 59.05 $\pm$ 0.21 (GATs), and 58.06 $\pm$ 0.13 (2ATT-GATs).

cause the UAS of predicted dependency trees in ID is relatively high compared to LAS, i.e., 87.31 and 77.33.

#### 4.4 Comparison of 2ATT-GATs Against BiLSTM-based SRL Models

We compare 2ATT-GATs with the widely used network in cross-lingual SRL, i.e., BiLSTMs. We build the first BiLSTM-based SRL model by stacking one BiLSTM layer with two layers of 2ATT-GATs (BiLSTM+2ATT-GATs). We also compare 2ATT-GATs with syntax-agnostic models, i.e., a model with one layer of BiLSTM as the encoder (BiLSTM) and a model with three layers of BiLSTMs as the encoder (3BiLSTMs).

The right side of Table 2 shows the F1 score of each model. The results with gold dependency

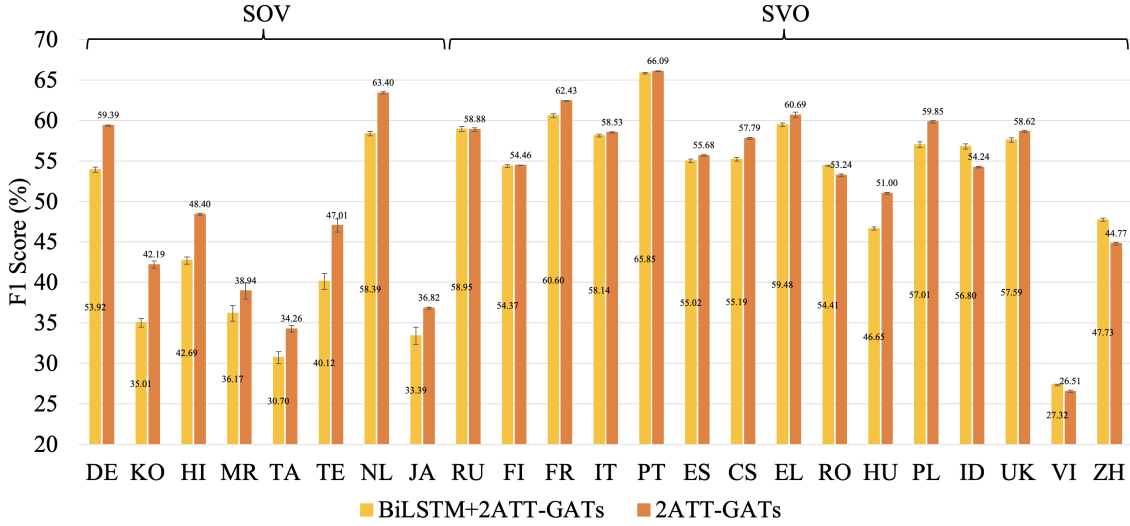


Figure 4: F1 scores (%) in UPB v2.0 test set using predicted dependency trees. Each language has two bars, i.e., the left bar indicates the F1 score for BiLSTM+2ATT-GATs, and the right bar indicates the F1 score for 2ATT-GATs.

trees are available in Appendix C. From EN, we can see that BiLSTM+2ATT-GATs perform best in the monolingual setting. BiLSTM+2ATT-GATs also perform better than syntax-agnostic models, i.e., BiLSTM and 3BiLSTMs, in the cross-lingual setting. This indicates that the help of syntax improves performance in both settings even though the syntax provided is not accurate.

Among the syntax-aware models, 2ATT-GATs perform the best in the cross-lingual setting. The context for each word learned through BiLSTMs, i.e., what words precede and follow each word in a sentence, might be too specific to the language it is trained with. By encoding the sentence over its dependency tree, we let information flow based on the sentence’s grammatical structure, which is more universal across languages.

Figure 4 compares the F1 scores of BiLSTM+2ATT-GATs and 2ATT-GATs for each language in UPB v2.0. We know that the models are trained in English that has a subject-verb-object (SVO) word order. If we look deeper, for certain languages with subject-object-verb (SOV) word order, i.e., DE, KO, HI, MR, TA, TE, NL, and JA, 2ATT-GATs show significant improvements over BiLSTM+2ATT-GATs. This proves that 2ATT-GATs have better transferability to languages with diverse word orders than the BiLSTM-based models.

BiLSTM+2ATT-GATs sometimes perform better than 2ATT-GATs in languages with the same word order as English, i.e., SVO word order, including RU, RO, ID, VI, and ZH. We confirm that in

RU, the predicted dependency trees decrease the F1 score a little more significantly for 2ATT-GATs<sup>6</sup>. For the other languages, the quality of the dependency trees (especially the dependency relations), which is relatively low in the original dataset, might be why 2ATT-GATs perform worst.

#### 4.5 Obtaining the Best Model

We conduct thorough experiments with the number of layers and contextualized multilingual word embedding that we use. Using frozen mBERT as contextualized multilingual word embedding, we compare the result of stacking one layer, two layers, and three layers of 2ATT-GATs as the encoder. We find that two layers of 2ATT-GATs and three layers of 2ATT-GATs give the best F1 score in UPB v1.0 and UPB v2.0, respectively.

We further experiment with two and three layers of 2ATT-GATs using mBERT and XLM-R as the contextualized multilingual word embeddings and fine-tune them. We provide the evaluation results in Appendix C. In UPB v1.0, a combination of frozen mBERT and two layers of 2ATT-GATs gives the best F1 score. Meanwhile, in UPB v2.0, fine-tuned XLM-R with two layers of 2ATT-GATs gives the best F1 score.

We further analyze why two layers of 2ATT-GATs perform better than three layers of 2ATT-GATs. In Figure 5, we group the F1 score difference between two 2ATT-GATs and three 2ATT-GATs for each dependency range, i.e., the number

<sup>6</sup>F1 scores of RU using gold trees:  $59.93 \pm 0.18$  (2ATT-GATs) and  $59.59 \pm 0.36$  (BiLSTM+2ATT-GATs).

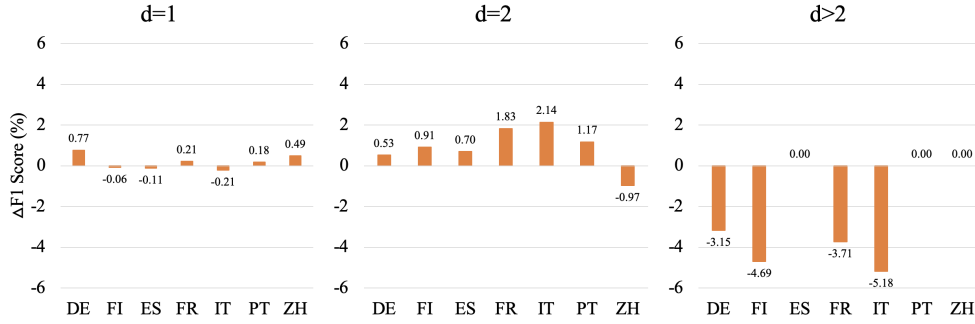


Figure 5: F1 score differences (%) between 2layers+mBERT and 3layers+mBERT in UPB v1.0 dev set (predicted).

of edges lies between the predicate and its argument. The evaluation shows that 2ATT-GATs perform better than three 2ATT-GATs because they perform better in  $d = 1$  and  $d = 2$ , and the dependency range ( $d$ ) in languages mostly lies in 1-2.

#### 4.6 Comparison of the Best Model Against Existing Works

We compare the best model on UPB v1.0 with existing works, i.e., Fei et al. (2020) and Zhang et al. (2021), as shown in Table 3. Fei et al. (2020) translate and project the SRL annotations from the source language to the target language to be included as part of the training set. Zhang et al. (2021) employ multi-task learning consisting of dependency parsing task and SRL task and include the dependency parsing task of the source and target language as part of the training process. Since both previous works have not been evaluated using UPB v2.0, we compare the F1 score in UPB v1.0. The evaluation results of our best model in UPB v2.0 can be seen in Appendix C.

It can be seen in Table 3 that 2ATT-GATs underperform previous models. Since our goal is to find the most transferable networks to build the cross-lingual model, unlike the other models, we do not include any knowledge of sentences from the target language in the training process. Our work can complement these previous works and provide important insights regarding the architecture design for building future cross-lingual SRL models.

## 5 Conclusions and Future Works

Through a simple encoder-decoder architecture, we show that GNNs are better than BiLSTMs for building cross-lingual SRL models, especially in distant languages. Encoding sentences based on their dependency trees helps create a more generalized cross-lingual SRL model rather than using word

Lang	2layers+mBERT (gold)	2layers+mBERT (predicted)	Fei	Zhang
FI	57.47 $\pm$ 0.24	53.82 $\pm$ 0.21	54.5	<b>59.9</b>
FR	48.80 $\pm$ 0.28	46.27 $\pm$ 0.25	<b>64.8</b>	56.6
DE	62.50 $\pm$ 0.24	58.78 $\pm$ 0.13	<b>65.0</b>	60.2
IT	<b>60.96<math>\pm</math>0.10</b>	57.73 $\pm$ 0.13	58.7	60.6
PT	57.21 $\pm$ 0.04	53.41 $\pm$ 0.12	56.0	<b>59.5</b>
ES	57.74 $\pm$ 0.10	54.47 $\pm$ 0.07	<b>62.5</b>	57.3
ZH	43.03 $\pm$ 0.10	38.87 $\pm$ 0.11	-	-

Table 3: F1 scores (%) in UPB v1.0 test set.

sequences. Furthermore, through empirical experiments comparing four types of GNNs using 23 languages available in UPB, we conclude that two-attention relational GATs are the most effective GNNs.

In the future, we can extend our model to incorporate language-specific information to distinguish characteristics between languages. This can be useful for training the model in a few-shot setting where we include the target sentences in the training set. Furthermore, in the experiments, we compare two-attention relational GATs with modified self-attention networks from Transformer. We can further analyze the effect of incorporating two-attention relational GATs to replace self-attention in Transformer in the cross-lingual SRL domain.

## Limitations

The limitation of this work is that we focus on argument detection and argument labeling in cross-lingual SRL, assuming that the sentences' gold predicates are easy to obtain. Furthermore, we focus on conducting experiments in a zero-shot setting. The availability of target sentences in the training set might affect the models' behavior, which should be investigated further.

## Ethics Statement

We believe there is no ethical issue raised in this work. SRL is a low-level task to support other advanced NLP applications. Therefore, increasing the coverage of SRL models in various languages is beneficial for developing NLP tools that can help solve the problems in this diverse society.

## References

- Wasi Ahmad, Zhisong Zhang, Xuezhe Ma, Eduard Hovy, Kai-Wei Chang, and Nanyun Peng. 2019. [On difficulties of cross-lingual transfer with order differences: A case study on dependency parsing](#). 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 2440–2452, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alan Akbik, Laura Chiticariu, Marina Danilevsky, Yunyao Li, Shivakumar Vaithyanathan, and Huaiyu Zhu. 2015. [Generating high quality proposition Banks for multilingual semantic role labeling](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 397–407, Beijing, China. Association for Computational Linguistics.
- Alan Akbik, Vishwajeet Kumar, and Yunyao Li. 2016. [Towards semi-automatic generation of proposition Banks for low-resource languages](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 993–998, Austin, Texas. Association for Computational Linguistics.
- Ann Bies, Justin Mott, Colin Warner, and Seth Kulick. 2012. [English web treebank](#).
- Rui Cai and Mirella Lapata. 2020. [Alignment-free cross-lingual semantic role labeling](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3883–3894, Online. Association for Computational Linguistics.
- Janara Christensen, Mausam, Stephen Soderland, and Oren Etzioni. 2010. [Semantic role labeling for open information extraction](#). In *Proceedings of the NAACL HLT 2010 First International Workshop on Formalisms and Methodology for Learning by Reading*, pages 52–60, Los Angeles, California. Association for Computational Linguistics.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. 2016. [Fast and accurate deep network learning by exponential linear units \(elus\)](#). In *International Conference on Learning Representations*.

- Simone Conia, Andrea Bacciu, and Roberto Navigli. 2021. [Unifying cross-lingual semantic role labeling with heterogeneous linguistic resources](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 338–351, Online. 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.
- 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.
- Hao Fei, Meishan Zhang, and Donghong Ji. 2020. [Cross-lingual semantic role labeling with high-quality translated training corpus](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7014–7026, Online. Association for Computational Linguistics.
- Daniel Gildea and Martha Palmer. 2002. [The necessity of parsing for predicate argument recognition](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 239–246, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Aleksa Gordić. 2020. [pytorch-gat](#). <https://github.com/gordicaleksa/pytorch-GAT>.
- Junfeng Jiang, An Wang, and Akiko Aizawa. 2021. [Attention-based relational graph convolutional network for target-oriented opinion words extraction](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1986–1997, Online. Association for Computational Linguistics.
- Ishan Jindal, Alexandre Rademaker, Michał Ulewicz, Ha Linh, Huyen Nguyen, Khoi-Nguyen Tran, Huaiyu Zhu, and Yunyao Li. 2022. [Universal Proposition Bank 2.0](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1700–1711, Marseille, France. European Language Resources Association.
- Atif Khan, Naomie Salim, and Yogan Jaya Kumar. 2015. [A framework for multi-document abstractive summa-](#)

667	rization based on semantic role labelling. <i>Applied</i>	
668	<i>Soft Computing</i> , 30:737–747.	
669	Paul Kingsbury and Martha Palmer. 2002. <a href="#">From Tree-</a>	
670	<a href="#">Bank to PropBank</a> . In <i>Proceedings of the Third In-</i>	
671	<i>ternational Conference on Language Resources and</i>	
672	<i>Evaluation (LREC'02)</i> , Las Palmas, Canary Islands	
673	- Spain. European Language Resources Association	
674	(ELRA).	
675	Thomas N. Kipf and Max Welling. 2017. <a href="#">Semi-</a>	
676	<a href="#">supervised classification with graph convolutional</a>	
677	<a href="#">networks</a> . In <i>International Conference on Learning</i>	
678	<i>Representations</i> .	
679	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled</a>	
680	<a href="#">weight decay regularization</a> . In <i>International Confer-</i>	
681	<i>ence on Learning Representations</i> .	
682	Diego Marcheggiani and Ivan Titov. 2017. <a href="#">Encoding</a>	
683	<a href="#">sentences with graph convolutional networks for se-</a>	
684	<a href="#">mantic role labeling</a> . In <i>Proceedings of the 2017</i>	
685	<i>Conference on Empirical Methods in Natural Lan-</i>	
686	<i>guage Processing</i> , pages 1506–1515, Copenhagen,	
687	Denmark. Association for Computational Linguis-	
688	tics.	
689	Joakim Nivre, Željko Agić, Lars Ahrenberg, Maria Je-	
690	sus Aranzabe, Masayuki Asahara, Aitziber Atutxa,	
691	Miguel Ballesteros, John Bauer, Kepa Bengoetxea,	
692	Yevgeni Berzak, Riyaz Ahmad Bhat, Eckhard	
693	Bick, Carl Birstell, Cristina Bosco, Gosse Bouma,	
694	Sam Bowman, Gülşen Cebiroğlu Eryiğit, Giuseppe	
695	G. A. Celano, Fabricio Chalub, Çağrı Çöl-	
696	tekin, Miriam Connor, Elizabeth Davidson, Marie-	
697	Catherine de Marneffe, Arantza Diaz de Ilarraza,	
698	Kaja Dobrovoljc, Timothy Dozat, Kira Drostanova,	
699	Puneet Dwivedi, Marhaba Eli, Tomaž Erjavec,	
700	Richárd Farkas, Jennifer Foster, Claudia Freitas,	
701	Katarína Gajdošová, Daniel Galbraith, Marcos Gar-	
702	cia, Moa Gärdenfors, Sebastian Garza, Filip Ginter,	
703	Iakes Goenaga, Koldo Gojenola, Memduh Gökır-	
704	mak, Yoav Goldberg, Xavier Gómez Guinovart,	
705	Berta Gonzáles Saavedra, Matias Grioni, Normunds	
706	Grūzītis, Bruno Guillaume, Jan Hajič, Linh Hà Mỹ,	
707	Dag Haug, Barbora Hladká, Radu Ion, Elena Ir-	
708	imnia, Anders Johannsen, Fredrik Jørgensen, Hüner	
709	Kaşıkar, Hiroshi Kanayama, Jenna Kanerva, Boris	
710	Katz, Jessica Kenney, Natalia Kotsyba, Simon Krek,	
711	Veronika Laippala, Lucia Lam, Phuong Lê Hồng,	
712	Alessandro Lenci, Nikola Ljubešić, Olga Lya-	
713	shevskaya, Teresa Lynn, Aibek Makazhanov, Christo-	
714	pher Manning, Cătălina Măranduc, David Mareček,	
715	Héctor Martínez Alonso, André Martins, Jan Mašek,	
716	Yuji Matsumoto, Ryan McDonald, Anna Missilä,	
717	Verginica Mititelu, Yusuke Miyao, Simonetta Monte-	
718	magni, Keiko Sophie Mori, Shunsuke Mori, Bohdan	
719	Moskalevskyi, Kadri Muischnek, Nina Mustafina,	
720	Kaili Müürisep, Luong Nguyễn Thị, Huyền Nguyễn	
721	Thị Minh, Vitaly Nikolaev, Hanna Nurmi, Petya	
722	Osenova, Robert Östling, Lilja Øvrelid, Valeria Paiva,	
723	Elena Pascual, Marco Passarotti, Cenel-Augusto	
724	Perez, Slav Petrov, Jussi Piitulainen, Barbara Plank,	
	Martin Popel, Lauma Pretkalniņa, Prokopis Proko-	725
	pidis, Tiina Puolakainen, Sampo Pyysalo, Alexan-	726
	dre Rademaker, Loganathan Ramasamy, Livy Real,	727
	Laura Rituma, Rudolf Rosa, Shadi Saleh, Baiba	728
	Saulite, Sebastian Schuster, Wolfgang Seeker, Moj-	729
	gan Seraji, Lena Shakurova, Mo Shen, Natalia Sil-	730
	veira, Maria Simi, Radu Simionescu, Katalin Simkó,	731
	Mária Šimková, Kiril Simov, Aaron Smith, Carolyn	732
	Spadine, Alane Suhr, Umut Sulubacak, Zsolt Szántó,	733
	Takaaki Tanaka, Reut Tsarfaty, Francis Tyers, Sumire	734
	Uematsu, Larraitz Uria, Gertjan van Noord, Vik-	735
	tor Varga, Veronika Vincze, Lars Wallin, Jing Xian	736
	Wang, Jonathan North Washington, Mats Wirén,	737
	Zdeněk Žabokrtský, Amir Zeldes, Daniel Zeman,	738
	and Hanzhi Zhu. 2016a. <a href="#">Universal dependencies</a>	739
	1.4. LINDAT/CLARIAH-CZ digital library at the	740
	Institute of Formal and Applied Linguistics (ÚFAL),	741
	Faculty of Mathematics and Physics, Charles Univer-	742
	sity.	743
	Joakim Nivre, Marie-Catherine de Marneffe, Filip Gin-	744
	ter, Yoav Goldberg, Jan Hajič, Christopher D. Man-	745
	ning, Ryan McDonald, Slav Petrov, Sampo Pyysalo,	746
	Natalia Silveira, Reut Tsarfaty, and Daniel Zeman.	747
	2016b. <a href="#">Universal Dependencies v1: A multilingual</a>	748
	<a href="#">treebank collection</a> . In <i>Proceedings of the Tenth In-</i>	749
	<i>ternational Conference on Language Resources and</i>	750
	<i>Evaluation (LREC'16)</i> , pages 1659–1666, Portorož,	751
	Slovenia. European Language Resources Association	752
	(ELRA).	753
	Joakim Nivre, Marie-Catherine de Marneffe, Filip Gin-	754
	ter, Jan Hajič, Christopher D. Manning, Sampo	755
	Pyysalo, Sebastian Schuster, Francis Tyers, and	756
	Daniel Zeman. 2020. <a href="#">Universal Dependencies v2:</a>	757
	<a href="#">An evergrowing multilingual treebank collection</a> . In	758
	<i>Proceedings of the Twelfth Language Resources and</i>	759
	<i>Evaluation Conference</i> , pages 4034–4043, Marseille,	760
	France. European Language Resources Association.	761
	Martha Palmer, Daniel Gildea, and Paul Kingsbury.	762
	2005. <a href="#">The Proposition Bank: An annotated corpus of</a>	763
	<a href="#">semantic roles</a> . <i>Computational Linguistics</i> , 31(1):71–	764
	106.	765
	Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and	766
	Christopher D. Manning. 2020. <a href="#">Stanza: A python</a>	767
	<a href="#">natural language processing toolkit for many human</a>	768
	<a href="#">languages</a> . In <i>Proceedings of the 58th Annual Meet-</i>	769
	<i>ing of the Association for Computational Linguistics:</i>	770
	<i>System Demonstrations</i> , pages 101–108, Online. As-	771
	sociation for Computational Linguistics.	772
	Prajit Ramachandran, Barret Zoph, and Quoc V. Le.	773
	2018. <a href="#">Searching for activation functions</a> .	774
	Reinhard Rapp. 2022. <a href="#">Using semantic role labeling to</a>	775
	<a href="#">improve neural machine translation</a> . In <i>Proceedings</i>	776
	<i>of the Thirteenth Language Resources and Evalua-</i>	777
	<i>tion Conference</i> , pages 3079–3083, Marseille, France.	778
	European Language Resources Association.	779
	Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018.	780
	<a href="#">Self-attention with relative position representations</a> .	781
	In <i>Proceedings of the 2018 Conference of the North</i>	782

783	<i>American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 464–468, New Orleans, Louisiana. Association for Computational Linguistics.	
784		
785		
786		
787		
788	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	
789	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	
790	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	
791	<a href="#">you need</a> . In <i>Advances in Neural Information Pro-</i>	
792	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	
793	Petar Veličković, Guillem Cucurull, Arantxa Casanova,	
794	Adriana Romero, Pietro Liò, and Yoshua Bengio.	
795	2017. <a href="#">Graph attention networks</a> .	
796	Kai Wang, Weizhou Shen, Yunyi Yang, Xiaojun Quan,	
797	and Rui Wang. 2020. <a href="#">Relational graph attention net-</a>	
798	<a href="#">work for aspect-based sentiment analysis</a> . In <i>Pro-</i>	
799	<i>ceedings of the 58th Annual Meeting of the Asso-</i>	
800	<i>ciation for Computational Linguistics</i> , pages 3229–	
801	3238, Online. Association for Computational Lin-	
802	guistics.	
803	Yuxuan Wang, Wanxiang Che, Jiang Guo, Yijia Liu, and	
804	Ting Liu. 2019. <a href="#">Cross-lingual BERT transformation</a>	
805	<a href="#">for zero-shot dependency parsing</a> . In <i>Proceedings of</i>	
806	<i>the 2019 Conference on Empirical Methods in Natu-</i>	
807	<i>ral Language Processing and the 9th International</i>	
808	<i>Joint Conference on Natural Language Processing</i>	
809	<i>(EMNLP-IJCNLP)</i> , pages 5721–5727, Hong Kong,	
810	China. Association for Computational Linguistics.	
811	Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le,	
812	Mohammad Norouzi, Wolfgang Macherey, Maxim	
813	Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff	
814	Klingner, Apurva Shah, Melvin Johnson, Xiaobing	
815	Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato,	
816	Taku Kudo, Hideto Kazawa, Keith Stevens, George	
817	Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason	
818	Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals,	
819	Greg Corrado, Macduff Hughes, and Jeffrey Dean.	
820	2016. <a href="#">Google’s neural machine translation system:</a>	
821	<a href="#">Bridging the gap between human and machine trans-</a>	
822	<a href="#">lation</a> .	
823	Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia	
824	Ackermann, Noëmi Aepli, Hamid Aghaei, Željko	
825	Agić, Amir Ahmadi, Lars Ahrenberg, Chika Kennedy	
826	Ajede, Gabrielè Aleksandravičiūtė, Ika Alfina, Lene	
827	Antonsen, Katya Aplonova, Angelina Aquino, Car-	
828	olina Aragon, Maria Jesus Aranzabe, Bilge Nas	
829	Arıcan, Hórunn Arnardóttir, Gashaw Arutie, Jes-	
830	sica Naraiswari Arwidarasti, Masayuki Asahara,	
831	Deniz Baran Aslan, Luma Ateyah, Furkan Atmaca,	
832	Mohammed Attia, Aitziber Atutxa, Liesbeth Au-	
833	gustinus, Elena Badmaeva, Keerthana Balasubra-	
834	mani, Miguel Ballesteros, Esha Banerjee, Sebastian	
835	Bank, Verginica Barbu Mititelu, Starkaður Barkar-	
836	son, Rodolfo Basile, Victoria Basmov, Colin Batch-	
837	elor, John Bauer, Seyyit Talha Bedir, Kepa Ben-	
838	goetxea, Gözde Berk, Yevgeni Berzak, Irshad Ah-	
839	mad Bhat, Riyaz Ahmad Bhat, Erica Biagetti, Eck-	
840	hard Bick, Agnė Bielinskienė, Kristín Bjarnadóttir,	
841	Rogier Blokland, Victoria Bobicev, Loïc Boizou,	
	Emanuel Borges Völker, Carl Börstell, Cristina	842
	Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd,	843
	Anouck Braggaar, Kristina Brokaitė, Aljoscha Bur-	844
	chardt, Marie Candito, Bernard Caron, Gauthier	845
	Caron, Lauren Cassidy, Tatiana Cavalcanti, Gülşen	846
	Cebiroğlu Eryiğit, Flavio Massimiliano Cecchini,	847
	Giuseppe G. A. Celano, Slavomír Čéplö, Nesli-	848
	han Cesur, Savas Cetin, Özlem Çetinoğlu, Fabri-	849
	cio Chalub, Shweta Chauhan, Ethan Chi, Taishi	850
	Chika, Yongseok Cho, Jinho Choi, Jayeol Chun,	851
	Juyeon Chung, Alessandra T. Cignarella, Silvie	852
	Cinková, Aurélie Collomb, Çağrı Çöltekin, Miriam	853
	Connor, Marine Courtin, Mihaela Cristescu, Phile-	854
	mon Daniel, Elizabeth Davidson, Marie-Catherine	855
	de Marneffe, Valeria de Paiva, Mehmet Oguz De-	856
	rin, Elvis de Souza, Arantza Diaz de Ilaraza,	857
	Carly Dickerson, Arawinda Dinakaramani, Elisa	858
	Di Nuovo, Bamba Dione, Peter Dirix, Kaja Do-	859
	brovoljc, Timothy Dozat, Kira Droganova, Puneet	860
	Dwivedi, Hanne Eckhoff, Sandra Eiche, Marhaba	861
	Eli, Ali Elkahky, Binyam Ephrem, Olga Erina,	862
	Tomaž Erjavec, Aline Etienne, Wograinne Evelyn,	863
	Sidney Facundes, Richárd Farkas, Jannatul Fer-	864
	daousi, Marília Fernanda, Hector Fernandez Alcalde,	865
	Jennifer Foster, Cláudia Freitas, Kazunori Fujita,	866
	Katarína Gajdošová, Daniel Galbraith, Marcos Gar-	867
	cia, Moa Gärdenfors, Sebastian Garza, Fabrício Fer-	868
	raz Gerardi, Kim Gerdes, Filip Ginter, Gustavo	869
	Godoy, Iakes Goenaga, Koldo Gojenola, Memduh	870
	Gökırmak, Yoav Goldberg, Xavier Gómez Guino-	871
	vart, Berta González Saavedra, Bernadeta Griciūtė,	872
	Matias Grioni, Loïc Grobol, Normunds Grūzītis,	873
	Bruno Guillaume, Céline Guillot-Barbance, Tunga	874
	Güngör, Nizar Habash, Hinrik Hafsteinsson, Jan Ha-	875
	jič, Jan Hajič jr., Mika Hämäläinen, Linh Hà Mỹ,	876
	Na-Rae Han, Muhammad Yudistira Hanifmuti, Sam	877
	Hardwick, Kim Harris, Dag Haug, Johannes Hei-	878
	neck, Oliver Hellwig, Felix Hennig, Barbora Hladká,	879
	Jaroslava Hlaváčová, Florinel Hociung, Petter Hohle,	880
	Eva Huber, Jena Hwang, Takumi Ikeda, Anton Karl	881
	Ingason, Radu Ion, Elena Irimia, Olájidé Ishola,	882
	Kaoru Ito, Siratun Jannat, Tomáš Jelínek, Apoorva	883
	Jha, Anders Johannsen, Hildur Jónsdóttir, Fredrik	884
	Jørgensen, Markus Juutinen, Sarveswaran K, Hüner	885
	Kaşıkar, Andre Kaasen, Nadezhda Kabaeva, Syl-	886
	vain Kahane, Hiroshi Kanayama, Jenna Kanerva,	887
	Neslihan Kara, Boris Katz, Tolga Kayadelen, Jes-	888
	sica Kenney, Václava Kettnerová, Jesse Kirchner,	889
	Elena Klementieva, Elena Klyachko, Arne Köhn,	890
	Abdullatif Köksal, Kamil Kopacewicz, Timo Korki-	891
	akangas, Mehmet Köse, Natalia Kotsyba, Jolanta	892
	Kovalevskaitė, Simon Krek, Parameswari Krishna-	893
	murthy, Sandra Kübler, Oğuzhan Kuyrukçu, Asli	894
	Kuzgun, Sookyoung Kwak, Veronika Laippala,	895
	Lucia Lam, Lorenzo Lambertino, Tatiana Lando,	896
	Septina Dian Larasati, Alexei Lavrentiev, John Lee,	897
	Phuong Lê Hồng, Alessandro Lenci, Saran Lertpra-	898
	dit, Herman Leung, Maria Levina, Cheuk Ying Li,	899
	Josie Li, Keying Li, Yuan Li, KyungTae Lim, Bruna	900
	Lima Padovani, Krister Lindén, Nikola Ljubešić,	901
	Olga Loginova, Stefano Lusito, Andry Luthfi, Mikko	902
	Luukko, Olga Lyashevskaya, Teresa Lynn, Vivien	903
	Macketanz, Menel Mahamdi, Jean Maillard, Aibek	904

905	Makazhanov, Michael Mandl, Christopher Manning,	968
906	Ruli Manurung, Büşra Marşan, Cătălina Mărănduc,	969
907	David Mareček, Katrin Marheinecke, Héctor	970
908	Martínez Alonso, Lorena Martín-Rodríguez, An-	971
909	dré Martins, Jan Mašek, Hiroshi Matsuda, Yuji	972
910	Matsumoto, Alessandro Mazzei, Ryan McDonald,	973
911	Sarah McGuinness, Gustavo Mendonça, Tatiana	974
912	Merzhevich, Niko Miekka, Karina Mischenkova,	975
913	Margarita Misirpashayeva, Anna Missilä, Cătălin	976
914	Mititelu, Maria Mitrofan, Yusuke Miyao, AmirHos-	977
915	sein Mojiri Froushani, Judit Molnár, Amirsaeid	978
916	Moloodi, Simonetta Montemagni, Amir More, Laura	979
917	Moreno Romero, Giovanni Moretti, Keiko Sophie	980
918	Mori, Shinsuke Mori, Tomohiko Morioka, Shigeki	981
919	Moro, Bjartur Mortensen, Bohdan Moskalevskyi,	982
920	Kadri Muischnek, Robert Munro, Yugo Murawaki,	983
921	Kaili Müürisep, Pinkey Nainwani, Mariam Nakhlé,	984
922	Juan Ignacio Navarro Horňáček, Anna Nedoluzhko,	985
923	Gunta Nešpore-Běrzkalne, Manuela Nevaci, Luong	986
924	Nguyễn Thị, Huyền Nguyễn Thị Minh, Yoshihiro	987
925	Nikaido, Vitaly Nikolaev, Rattima Nitisaroj, Alireza	988
926	Nourian, Hanna Nurmi, Stina Ojala, Atul Kr. Ojha,	989
927	Adédáyò Olúòkun, Mai Omura, Emeka Onwueg-	990
928	buzia, Petya Osenova, Robert Östling, Lilja Øvre-	991
929	lid, Şaziye Betül Özateş, Merve Özçelik, Arzu-	992
930	can Özgür, Balkız Öztürk Başaran, Hyunji Hay-	993
931	ley Park, Niko Partanen, Elena Pascual, Marco	994
932	Passarotti, Agnieszka Patejuk, Guilherme Paulino-	
933	Passos, Angelika Peljak-Łapińska, Siyao Peng,	995
934	Cenel-Augusto Perez, Natalia Perkova, Guy Per-	996
935	rier, Slav Petrov, Daria Petrova, Jason Phelan, Jussi	997
936	Piitulainen, Tommi A Pirinen, Emily Pitler, Bar-	998
937	bara Plank, Thierry Poibeau, Larisa Ponomareva,	999
938	Martin Popel, Lauma Pretkalnina, Sophie Prévost,	1000
939	Prokopis Prokopidis, Adam Przepiórkowski, Ti-	1001
940	ina Puolakainen, Sampo Pyysalo, Peng Qi, An-	
941	driela Rääbis, Alexandre Rademaker, Mizanur Ra-	
942	homan, Taraka Rama, Loganathan Ramasamy, Car-	
943	los Ramisch, Fam Rashel, Mohammad Sadegh Ra-	
944	sooli, Vinit Ravishankar, Livy Real, Petru Rebeja,	
945	Siva Reddy, Mathilde Regnault, Georg Rehm, Ivan	
946	Riabov, Michael Rießler, Erika Rimkutė, Larissa Ri-	
947	naldi, Laura Rituma, Putri Rizqiyah, Luisa Rocha,	
948	Eirikur Rögnvaldsson, Mykhailo Romanenko, Rudolf	
949	Rosa, Valentin Roşca, Davide Rovati, Olga Rud-	
950	ina, Jack Rueter, Kristján Rúnarsson, Shoval Sadde,	
951	Pegah Safari, Benoît Sagot, Aleks Sahala, Shadi	
952	Saleh, Alessio Salomoni, Tanja Samardžić, Stephanie	
953	Samson, Manuela Sanguinetti, Ezgi Sanıyar, Dage	
954	Särg, Baiba Saulīte, Yanin Sawanakunanon, Shefali	
955	Saxena, Kevin Scannell, Salvatore Scarlata, Nathan	
956	Schneider, Sebastian Schuster, Lane Schwartz,	
957	Djamé Seddah, Wolfgang Seeker, Mojgan Seraji,	
958	Syeda Shahzadi, Mo Shen, Atsuko Shimada, Hi-	
959	royuki Shirasu, Yana Shishkina, Muh Shohibussirri,	
960	Dmitry Sichinava, Janine Siewert, Einar Freyr Sig-	
961	urðsson, Aline Silveira, Natalia Silveira, Maria Simi,	
962	Radu Simionescu, Katalin Simkó, Mária Šimková,	
963	Kiril Simov, Maria Skachodubova, Aaron Smith, Is-	
964	abela Soares-Bastos, Shafī Sourov, Carolyn Spadine,	
965	Rachele Sprugnoli, Steinhórfur Steingrímsson, Antonio	
966	Stella, Milan Straka, Emmett Strickland, Jana Str-	
967	nádová, Alane Suhr, Yogi Lesmana Sulestio, Umut	
	Sulubacak, Shingo Suzuki, Zsolt Szántó, Chihiro	968
	Taguchi, Dima Taji, Yuta Takahashi, Fabio Tam-	969
	burini, Mary Ann C. Tan, Takaaki Tanaka, Dipta	970
	Tanaya, Samson Tella, Isabelle Tellier, Marinella	971
	Testori, Guillaume Thomas, Liisi Torga, Marsida	972
	Toska, Trond Trosterud, Anna Trukhina, Reut Tsar-	973
	faty, Utku Türk, Francis Tyers, Sumire Uematsu, Ro-	974
	man Untilov, Zdeňka Uřešová, Larraitz Uria, Hans	975
	Uszkoreit, Andrius Utka, Sowmya Vajjala, Rob	976
	van der Goot, Martine Vanhove, Daniel van Niekerk,	977
	Gertjan van Noord, Viktor Varga, Eric Villemonte	978
	de la Clergerie, Veronika Vincze, Natalia Vlasova,	979
	Aya Wakasa, Joel C. Wallenberg, Lars Wallin, Abi-	980
	gail Walsh, Jing Xian Wang, Jonathan North Wash-	981
	ington, Maximilian Wendt, Paul Widmer, Sri Hartati	982
	Wijono, Seyi Williams, Mats Wirén, Christian Wit-	983
	tern, Tsegay Woldemariam, Tak-sum Wong, Alina	984
	Wróblewska, Mary Yako, Kayo Yamashita, Naoki	985
	Yamazaki, Chunxiao Yan, Koichi Yasuoka, Marat M.	986
	Yavrumyan, Arife Betül Yenice, Olcay Taner Yıldız,	987
	Zhuoran Yu, Arlisa Yuliawati, Zdeněk Žabokrtský,	988
	Shorouq Zahra, Amir Zeldes, He Zhou, Hanzhi Zhu,	989
	Anna Zhuravleva, and Rayan Ziane. 2021. <a href="#">Universal</a>	990
	<a href="#">dependencies 2.9</a> . LINDAT/CLARIAH-CZ digital	991
	library at the Institute of Formal and Applied Linguis-	992
	tics (ÚFAL), Faculty of Mathematics and Physics,	993
	Charles University.	994
	Zhisong Zhang, Emma Strubell, and Eduard Hovy. 2021.	995
	<a href="#">On the benefit of syntactic supervision for cross-</a>	996
	<a href="#">lingual transfer in semantic role labeling</a> . In <i>Proceed-</i>	997
	<i>ings of the 2021 Conference on Empirical Methods</i>	998
	<i>in Natural Language Processing</i> , pages 6229–6246,	999
	Online and Punta Cana, Dominican Republic. Asso-	1000
	ciation for Computational Linguistics.	1001

Lang	Train	Dev	Test
<b>UPB v1.0</b>			
English (EN)	12,542	1,974	2,060
Finnish (FI)	12,217	716	648
Italian (IT)	12,837	489	489
Spanish (ES)	28,492	3,206	1,995
French (FR)	14,553	1,596	298
German (DE)	14,118	799	977
Portuguese (PT)	7,494	938	936
Chinese (ZH)	3,997	500	500
<b>UPB v2.0</b>			
English (EN)	12,542	1,974	2,062
Czech (CS)	102,993	11,311	12,203
Greek (EL)	1,662	403	456
Korean (KO)	27,410	3,016	3,276
Romanian (RO)	35,911	2,247	2,272
Hindi (HI)	13,304	1,659	1,684
Marathi (MR)	373	46	47
Tamil (TA)	400	80	120
Hungarian (HU)	910	441	449
Polish (PL)	31,496	3,960	3,942
Telugu (TE)	1,051	131	146
Dutch (NL)	18,078	1,394	1,472
Indonesian (ID)	4,482	559	557
Japanese (JA)	14,100	1,014	1,086
Russian (RU)	19,894	1,525	1,482
Ukrainian (UK)	5,496	672	892
Chinese (ZH)	3,997	500	500
Vietnamese (VI)	1,400	800	800
Finnish (FI)	27,198	3,239	3,422
Italian (IT)	29,685	2,277	2,518
Spanish (ES)	28,474	3,054	2,147
French (FR)	17,968	2,970	1,712
German (DE)	166,849	19,233	19,436
Portuguese (PT)	16,633	2,376	2,367

Table 4: Number of sentences available for each language in UPB v1.0 and UPB v2.0.

Treebank	License
<b>UD v1.4</b>	
UD_English	CC BY-SA 4.0
UD_Chinese	CC BY-NC-SA 4.0
UD_Finnish	CC BY-SA 4.0
UD_Spanish	CC BY-NC-SA 3.0 US
UD_Spanish-AnCora	GNU GPL 3.0
UD_French	CC BY-NC-SA 3.0 US
UD_German	CC BY-NC-SA 3.0 US
UD_Italian	CC BY-NC-SA 3.0
UD_Portuguese-Bosque	CC BY-SA 4.0
<b>UD v2.9</b>	
UD_English-EWT	CC BY-SA 4.0
UD_Chinese-GSD	CC BY-SA 4.0
UD_Czech-CAC	CC BY-SA 4.0
UD_Czech-CLTT	CC BY-SA 4.0
UD_Czech-FicTree	CC BY-NC-SA 4.0
UD_Czech-PDT	CC BY-NC-SA 3.0
UD_Dutch-Alpino	CC BY-SA 4.0
UD_Dutch-LassySmall	CC BY-SA 4.0
UD_Finnish-FTB	CC BY 4.0
UD_Finnish-TDT	CC BY-SA 4.0
UD_French-GSD	CC BY-SA 4.0
UD_French-Rhapsodie	CC BY-SA 4.0
UD_French-Sequoia	LGPL-LR
UD_German-GSD	CC BY-SA 4.0
UD_German-HDT	CC BY-SA 4.0
UD_Greek-GDT	CC BY-NC-SA 3.0
UD_Hindi-HDTB	CC BY-NC-SA 4.0
UD_Hungarian-Szeged	CC BY-NC-SA 3.0
UD_Indonesian-GSD	CC BY-SA 4.0
UD_Italian-ISDT	CC BY-NC-SA 3.0
UD_Italian-ParTUT	CC BY-NC-SA 4.0
UD_Italian-PoSTWITA	CC BY-NC-SA 4.0
UD_Italian-TWITTIRO	CC BY-SA 4.0
UD_Italian-VIT	CC BY-NC-SA 3.0
UD_Japanese-GSD	CC BY-SA 4.0
UD_Japanese-GSDLUW	CC BY-SA 4.0
UD_Korean-GSD	CC BY-SA 4.0
UD_Korean-Kaist	CC BY-SA 4.0
UD_Marathi-UFAL	CC BY-SA 4.0
UD_Polish-LFG	GNU GPL 3.0
UD_Polish-PDB	CC BY-NC-SA 4.0
UD_Portuguese-Bosque	CC BY-SA 4.0
UD_Portuguese-GSD	CC BY-SA 4.0
UD_Romanian-Nonstandard	CC BY-SA 4.0
UD_Romanian-RRT	CC BY-SA 4.0
UD_Romanian-SiMoNERo	CC BY-SA 4.0
UD_Russian-GSD	CC BY-SA 4.0
UD_Russian-Taiga	CC BY-SA 4.0
UD_Spanish-AnCora	CC BY 4.0
UD_Spanish-GSD	CC BY-SA 4.0
UD_Tamil-TTB	CC BY-NC-SA 3.0
UD_Telugu-MTG	CC BY-SA 4.0
UD_Ukrainian-IU	CC BY-NC-SA 4.0
UD_Vietnamese-VTB	CC BY-SA 4.0

Table 5: License for each treebank in UD v1.4 and UD2.9 that we use in the experiments.

## A Artifacts

### A.1 Dataset Distribution

Table 4 shows the dataset distribution in Universal Proposition Bank (UPB). Since we run our experiments in a zero-shot setting, we only use the dev set and test set for languages other than English.

### A.2 Licenses

Complete UPB v1.0 and UPB v2.0 contain annotations from Universal Dependencies (UD) v1.4 and UD v2.9. Therefore, we provide the license for each UD treebank in Table 5. Despite UD’s inherited license, UPB v1.0 and UPB v2.0 also have a CDLA-Sharing-1.0 license. We also retrieve annotations from PropBank v3, English Web Treebank, and UD v1.4 to form English SRL annotations based on UD v1.4. Therefore, we provide the license for PropBank v3, i.e., CC BY-SA 4.0, and English Web Treebank, i.e., LDC User Agreement for Non-Members.

We refer to publicly available codes to help pre-

process the data and build the model. We provide the list of repositories with their corresponding licenses, i.e., Zhang et al. (2021)<sup>7</sup> (GPL-3.0), <https://github.com/UniversalPropositions/tools> (Apache-2.0), Gordić (2020)<sup>8</sup> (MIT), and Marcheggiani and Titov (2017)<sup>9</sup> (Apache-2.0).

We access all the resources we mentioned above solely for academic research. We make sure that we obey the intended use for each artifact.

## B Training

### B.1 Dependency Parsers

We train the dependency parsers using Stanza with a 0.0005 learning rate, 70,000 max steps, and 10,000 max steps before stopping. Table 6 shows

<sup>7</sup><https://github.com/zzsforlpl/zmsp>

<sup>8</sup><https://github.com/gordicaleksa/pytorch-GAT>

<sup>9</sup><https://github.com/diegma/neural-dep-srl>

Treebank	Dev		Test	
	UAS	LAS	UAS	LAS
<b>UPB v1.0</b>				
UP_English-EWT	92.18	90.27	91.39	89.48
UP_Chinese	87.28	85.11	88.31	86.13
UP_Finnish	92.15	90.82	90.28	89.04
UP_Spanish	91.68	89.75	91.15	88.95
UP_Spanish-AnCora	94.06	92.81	93.80	92.39
UP_French	92.50	91.08	89.89	87.61
UP_German	92.70	91.13	89.45	87.31
UP_Italian	94.12	92.75	93.70	92.36
UP_Portuguese-Bosque	93.20	92.02	92.45	91.05
<b>UPB v2.0</b>				
UP_English-EWT	92.46	90.86	91.42	89.82
UP_Chinese-GSD	85.11	83.19	87.06	85.13
UP_Czech-CAC	92.97	91.62	93.43	91.68
UP_Czech-CLTT	89.13	86.98	88.32	86.09
UP_Czech-FicTree	94.68	93.11	94.61	92.76
UP_Czech-PDT	93.74	92.24	93.50	91.87
UP_Dutch-Alpino	94.53	92.24	92.87	90.42
UP_Dutch-LassySmall	90.77	87.62	92.12	89.11
UP_Finnish-FTB	93.77	92.29	94.03	92.41
UP_Finnish-TDT	91.97	90.41	92.24	90.74
UP_French-GSD	95.66	94.45	93.47	91.87
UP_French-Rhapsodie	87.75	83.25	86.42	81.88
UP_French-Sequoia	93.54	92.23	93.10	91.70
UP_German-GSD	91.78	88.61	89.65	85.62
UP_German-HDT	95.18	93.64	95.30	93.72
UP_Greek-GDT	91.77	90.43	92.93	91.19
UP_Hindi-HDTB	96.62	94.49	96.68	94.43
UP_Hungarian-Szeged	87.64	84.10	86.72	83.25
UP_Indonesian-GSD	86.49	76.25	87.31	77.33
UP_Italian-ISDT	94.41	92.84	94.37	93.16
UP_Italian-ParTUT	92.76	90.52	93.10	91.40
UP_Italian-PoSTWITA	87.21	83.20	88.33	84.41
UP_Italian-TWITTRO	87.25	81.64	84.85	79.77
UP_Italian-VIT	90.63	88.82	91.54	89.05
UP_Japanese-GSD	96.09	95.47	95.11	94.21
UP_Japanese-GSDLUW	96.12	95.82	95.35	95.12
UP_Korean-GSD	88.22	85.41	89.65	87.07
UP_Korean-Kaist	91.35	90.39	90.41	89.45
UP_Marathi-UFAL	74.55	64.32	79.85	70.63
UP_Polish-LFG	97.56	96.73	97.80	96.92
UP_Polish-PDB	94.17	92.69	94.58	93.16
UP_Portuguese-Bosque	94.25	92.51	94.85	93.54
UP_Portuguese-GSD	94.44	93.34	94.21	93.23
UP_Romanian-Nonstandard	93.18	90.04	91.43	87.75
UP_Romanian-RRT	91.96	88.60	91.93	88.45
UP_Romanian-SiMoNERo	93.38	91.21	93.78	91.86
UP_Russian-GSD	90.55	87.80	90.44	87.21
UP_Russian-Taiga	83.94	79.32	84.42	81.41
UP_Spanish-AnCora	93.83	92.16	93.82	92.00
UP_Spanish-GSD	91.91	89.79	91.93	89.58
UP_Tamil-TTB	81.24	73.48	80.89	72.30
UP_Telugu-MTG	92.90	86.25	93.07	85.58
UP_Ukrainian-IU	91.14	89.34	90.10	88.24
UP_Vietnamese-VTB	78.92	74.99	77.58	74.16

Table 6: UAS and LAS of each treebank’s dependency parser.

UAS and LAS for each treebank in UPB v1.0 and UPB v2.0.

## B.2 SRL Models

We build the code for training and evaluation using PyTorch library<sup>10</sup>. We use the transformers library provided by Hugging Face<sup>11</sup> to produce contextualized multilingual word embedding from mBERT (i.e., bert-base-multilingual-cased) and XLM-R (i.e., xlm-roberta-base).

We refer to Gordić (2020) to implement GATs. Following the implementation, we use ELU (Clevert et al., 2016) as the activation function in each layer,  $\sigma$ , except in the final layer where we use *Swish*, the same activation function that we use in the input layer. We tried to use the same activation

<sup>10</sup><https://pytorch.org/>

<sup>11</sup><https://huggingface.co/>

Hyperparameter	Value
learning rate	$10^{-3}$
minimum learning rate	$10^{-5}$
weight decay	$10^{-4}$
word embedding learning rate	$10^{-5}$
word embedding minimum learning rate	$10^{-6}$
word embedding weight decay	$10^{-4}$
batch size	32
epochs	30
warmup epochs	1
cooldown epochs	15
number of attention heads, $K$	8
predicate indicator embedding size, $\vec{p}_i^z$	64
multilingual word embedding size, $\vec{c}_i^z$	512
edge embedding size, $\vec{r}_{ij}^z$	64
LSTM hidden size	512
GATs, SATs, and 2ATT-GATs hidden size	64
GATs, SATs, and 2ATT-GATs output size	same as input (576)
GGCNs hidden size	same as input (576)
GGCNs output size	same as input (576)
dropout in input layer	0.2
dropout before decoder	0.3
GATs best node and edge dropout	0.2
SATs best node and edge dropout	0.2
2ATT-GATs best node and edge dropout	0.3
GGCNs best dropout	0.5

Table 7: Hyperparameter values for cross-lingual and monolingual SRL models.

function for every layer, but this setting works best in this task. We use the same skeleton as GATs for implementing 2ATT-GATs and SATs, where node and edge dropouts are placed at the beginning of every layer.

In SATs, we experiment with two settings for the edge representations, i.e.,  $a_{ij}$ . First, we used the same edge representation for each attention head,  $k$ . Then, we used different edge representations for each attention head. The results show that using the same edge representation in each attention head works best for this task. Therefore, we use this setting when comparing SATs with other GNNs.

For the implementation of GGCNs, we refer to the original implementation, i.e., Marcheggiani and Titov (2017)<sup>12</sup>, and its reimplement in PyTorch<sup>13</sup>. Following the reimplement in PyTorch, we use LeakyReLU as the activation function for GGCNs in each layer, except in the final layer where we use *Swish*. We place the dropout at the end of every layer.

<sup>12</sup><https://github.com/diegma/neural-dep-srl>

<sup>13</sup><https://github.com/kdrivas/Graph-convolutional>

Lang	Gold	Predicted
CS	61.44 $\pm$ 0.08	61.42 $\pm$ 0.07
EL	68.80 $\pm$ 0.09	67.90 $\pm$ 0.16
KO	46.76 $\pm$ 0.19	47.06 $\pm$ 0.20
RO	55.55 $\pm$ 0.05	55.36 $\pm$ 0.04
HI	48.73 $\pm$ 0.15	48.97 $\pm$ 0.16
MR	41.78 $\pm$ 0.41	42.19 $\pm$ 0.80
TA	40.73 $\pm$ 0.34	37.79 $\pm$ 0.52
HU	54.86 $\pm$ 0.07	53.94 $\pm$ 0.10
PL	62.85 $\pm$ 0.10	63.01 $\pm$ 0.11
TE	56.51 $\pm$ 0.47	56.19 $\pm$ 0.46
NL	65.36 $\pm$ 0.14	64.69 $\pm$ 0.11
ID	63.83 $\pm$ 0.13	63.22 $\pm$ 0.07
JA	37.82 $\pm$ 0.75	38.22 $\pm$ 0.80
RU	62.40 $\pm$ 0.19	62.02 $\pm$ 0.16
UK	60.85 $\pm$ 0.17	60.63 $\pm$ 0.15
ZH	42.33 $\pm$ 0.72	41.68 $\pm$ 0.75
VI	28.80 $\pm$ 0.14	29.30 $\pm$ 0.15
FI	58.28 $\pm$ 0.10	57.95 $\pm$ 0.11
IT	60.82 $\pm$ 0.10	60.53 $\pm$ 0.07
ES	57.97 $\pm$ 0.09	58.29 $\pm$ 0.07
FR	64.38 $\pm$ 0.10	63.65 $\pm$ 0.11
DE	61.35 $\pm$ 0.11	60.86 $\pm$ 0.09
PT	67.75 $\pm$ 0.07	67.31 $\pm$ 0.07

Table 8: F1 scores (%) in UPB v2.0 test set with gold dependency trees (left) and predicted dependency trees (right) (2layers+XLM-R+fine-tuned).

We provide the hyperparameter values we use to train the SRL models in Table 7. For GNN-based models, we search for the best dropout in 0.1 – 0.5 with a 0.1 increment.

We use Tesla P100 to train the models. Training time for GNN-based models ranges from 3-4 hours. Training time for BiLSTM+2ATT-GATs and BiLSTM ranges from 4-5 hours. Meanwhile, for 3BiLSTMs, it takes around 10 hours to train a model. We run the experiment 5 times and produce two models for each setting, i.e., a model trained on UPB v1.0 and a model trained on UPB v2.0. The estimation of total GPU hours is 2,000.

## C Supporting Results

Table 9 shows the F1 scores for GNN-based and BiLSTM-based cross-lingual and monolingual SRL models using gold dependency trees.

Table 10 shows the average F1 scores of target languages in dev sets using different numbers of layers, i.e., two layers and three layers, and different contextualized multilingual word embeddings, i.e., mBERT and XLM-R.

Table 8 shows the F1 scores of the best model (2layers+XLM-R+fine-tuned) for each language in the UPB v2.0 test set.

Lang	GGCNs	SATs	GATs	2ATT-GATs	BiLSTM+ 2ATT-GATs	BiLSTM	3BiLSTMs
UPB v1.0							
EN	83.82 $\pm$ 0.06	<b>84.19<math>\pm</math>0.05</b>	82.77 $\pm$ 0.04	83.21 $\pm$ 0.07	<u>85.35<math>\pm</math>0.14</u>	79.72 $\pm$ 0.05	80.23 $\pm$ 0.08
FI	<b>58.02<math>\pm</math>0.09</b>	57.77 $\pm$ 0.19	52.07 $\pm$ 0.28	56.82 $\pm$ 0.08	56.18 $\pm$ 0.24	39.85 $\pm$ 0.19	40.18 $\pm$ 0.22
IT	<b>61.09<math>\pm</math>0.12</b>	60.77 $\pm$ 0.23	57.18 $\pm$ 0.29	<u>60.52<math>\pm</math>0.09</u>	59.00 $\pm$ 0.10	47.20 $\pm$ 0.19	46.74 $\pm$ 0.30
ES	57.39 $\pm$ 0.12	<b>58.12<math>\pm</math>0.17</b>	52.70 $\pm$ 0.18	<u>57.97<math>\pm</math>0.16</u>	55.02 $\pm$ 0.20	41.77 $\pm$ 0.15	41.52 $\pm$ 0.05
FR	<b>49.42<math>\pm</math>0.14</b>	48.99 $\pm$ 0.20	46.04 $\pm$ 0.19	<u>48.68<math>\pm</math>0.17</u>	47.87 $\pm$ 0.15	40.64 $\pm$ 0.32	40.70 $\pm$ 0.31
DE	<b>62.78<math>\pm</math>0.27</b>	62.38 $\pm$ 0.05	60.31 $\pm$ 0.15	<u>62.63<math>\pm</math>0.09</u>	59.55 $\pm$ 0.22	41.45 $\pm$ 0.24	41.99 $\pm$ 0.12
PT	<b>57.60<math>\pm</math>0.05</b>	56.90 $\pm$ 0.15	55.26 $\pm$ 0.17	<u>57.18<math>\pm</math>0.08</u>	56.35 $\pm$ 0.24	44.38 $\pm$ 0.24	44.10 $\pm$ 0.12
ZH	42.12 $\pm$ 0.20	42.00 $\pm$ 0.20	39.63 $\pm$ 0.14	<b>42.60<math>\pm</math>0.16</b>	43.51 $\pm$ 0.14	32.87 $\pm$ 0.37	32.53 $\pm$ 0.23
AVG	<b>55.49<math>\pm</math>0.06</b>	55.28 $\pm$ 0.13	51.88 $\pm$ 0.13	55.20 $\pm$ 0.08	53.93 $\pm$ 0.14	41.17 $\pm$ 0.13	41.11 $\pm$ 0.11
UPB v2.0							
EN	83.93 $\pm$ 0.08	<b>84.39<math>\pm</math>0.06</b>	82.82 $\pm$ 0.09	83.26 $\pm$ 0.07	85.37 $\pm$ 0.10	79.51 $\pm$ 0.08	80.08 $\pm$ 0.08
CS	56.42 $\pm$ 0.19	56.75 $\pm$ 0.11	55.83 $\pm$ 0.07	<b>57.70<math>\pm</math>0.10</b>	55.16 $\pm$ 0.24	47.13 $\pm$ 0.06	47.50 $\pm$ 0.13
EL	59.57 $\pm$ 0.33	59.64 $\pm$ 0.57	60.84 $\pm$ 0.26	<b>61.80<math>\pm</math>0.35</b>	60.66 $\pm$ 0.21	56.05 $\pm$ 0.16	55.53 $\pm$ 0.27
KO	39.07 $\pm$ 0.22	38.20 $\pm$ 0.80	40.83 $\pm$ 0.26	<b>42.14<math>\pm</math>0.48</b>	34.94 $\pm$ 0.59	27.75 $\pm$ 0.29	27.84 $\pm$ 0.99
RO	52.64 $\pm$ 0.09	52.83 $\pm$ 0.19	<b>53.82<math>\pm</math>0.20</b>	53.37 $\pm$ 0.17	<u>54.47<math>\pm</math>0.14</u>	48.71 $\pm$ 0.30	47.83 $\pm$ 0.31
HI	46.77 $\pm$ 0.08	47.41 $\pm$ 0.08	44.63 $\pm$ 0.23	<b>47.93<math>\pm</math>0.12</b>	42.45 $\pm$ 0.46	27.49 $\pm$ 0.84	27.30 $\pm$ 0.44
MR	34.57 $\pm$ 0.99	37.08 $\pm$ 0.68	37.43 $\pm$ 1.66	<b>39.16<math>\pm</math>1.02</b>	36.73 $\pm$ 1.44	27.22 $\pm$ 1.89	26.31 $\pm$ 1.60
TA	34.83 $\pm$ 0.34	36.51 $\pm$ 0.36	<b>38.25<math>\pm</math>0.57</b>	<u>37.74<math>\pm</math>0.63</u>	33.76 $\pm$ 0.79	23.64 $\pm$ 0.84	22.20 $\pm$ 0.88
HU	52.33 $\pm$ 0.29	<b>53.14<math>\pm</math>0.17</b>	49.34 $\pm$ 0.14	<u>52.81<math>\pm</math>0.11</u>	48.44 $\pm$ 0.25	39.05 $\pm$ 0.29	39.52 $\pm$ 0.39
PL	57.56 $\pm$ 0.17	57.26 $\pm$ 0.08	57.50 $\pm$ 0.13	<b>59.75<math>\pm</math>0.15</b>	56.80 $\pm$ 0.34	51.65 $\pm$ 0.11	51.19 $\pm$ 0.33
TE	<b>44.50<math>\pm</math>0.72</b>	42.94 $\pm$ 0.62	40.91 $\pm$ 0.68	<u>44.44<math>\pm</math>0.51</u>	39.21 $\pm$ 1.22	35.97 $\pm$ 1.30	32.84 $\pm$ 0.55
NL	63.81 $\pm$ 0.26	63.95 $\pm$ 0.21	63.38 $\pm$ 0.14	<b>64.34<math>\pm</math>0.19</b>	59.15 $\pm$ 0.23	51.95 $\pm$ 0.36	52.42 $\pm$ 0.10
ID	57.18 $\pm$ 0.26	57.05 $\pm$ 0.20	<b>59.05<math>\pm</math>0.21</b>	58.06 $\pm$ 0.13	60.49 $\pm$ 0.23	56.55 $\pm$ 0.34	55.00 $\pm$ 0.35
JA	35.70 $\pm$ 0.13	35.87 $\pm$ 0.43	35.96 $\pm$ 0.54	<b>36.35<math>\pm</math>0.09</b>	33.15 $\pm$ 1.13	23.76 $\pm$ 1.07	21.40 $\pm$ 1.07
RU	58.84 $\pm$ 0.23	59.82 $\pm$ 0.25	59.59 $\pm$ 0.24	<b>59.93<math>\pm</math>0.18</b>	59.59 $\pm$ 0.36	55.49 $\pm$ 0.24	55.48 $\pm$ 0.23
UK	57.89 $\pm$ 0.08	58.88 $\pm$ 0.18	57.53 $\pm$ 0.21	<b>59.14<math>\pm</math>0.15</b>	58.18 $\pm$ 0.25	52.80 $\pm$ 0.39	52.98 $\pm$ 0.17
ZH	44.22 $\pm$ 0.27	45.51 $\pm$ 0.20	45.89 $\pm$ 0.31	<b>45.99<math>\pm</math>0.22</b>	<u>48.68<math>\pm</math>0.20</u>	48.54 $\pm$ 0.46	47.95 $\pm$ 0.34
VI	25.14 $\pm$ 0.04	25.32 $\pm$ 0.15	<b>27.95<math>\pm</math>0.04</b>	26.04 $\pm$ 0.07	27.07 $\pm$ 0.11	29.39 $\pm$ 0.29	27.34 $\pm$ 0.36
FI	53.89 $\pm$ 0.11	54.37 $\pm$ 0.14	53.12 $\pm$ 0.11	<b>54.84<math>\pm</math>0.03</b>	54.74 $\pm$ 0.19	51.54 $\pm$ 0.13	51.83 $\pm$ 0.23
IT	57.80 $\pm$ 0.05	58.29 $\pm$ 0.18	58.24 $\pm$ 0.08	<b>58.91<math>\pm</math>0.09</b>	58.73 $\pm$ 0.20	55.23 $\pm$ 0.16	54.97 $\pm$ 0.12
ES	53.59 $\pm$ 0.07	54.37 $\pm$ 0.08	<b>55.47<math>\pm</math>0.07</b>	<u>55.27<math>\pm</math>0.11</u>	54.75 $\pm$ 0.20	51.02 $\pm$ 0.13	50.51 $\pm$ 0.06
FR	63.03 $\pm$ 0.11	63.13 $\pm$ 0.13	62.67 $\pm$ 0.19	<b>63.79<math>\pm</math>0.06</b>	61.52 $\pm$ 0.18	59.81 $\pm$ 0.15	59.88 $\pm$ 0.26
DE	58.89 $\pm$ 0.05	59.31 $\pm$ 0.15	59.80 $\pm$ 0.06	<b>59.89<math>\pm</math>0.07</b>	54.44 $\pm$ 0.31	41.40 $\pm$ 0.38	43.10 $\pm$ 0.23
PT	66.14 $\pm$ 0.06	66.31 $\pm$ 0.13	<b>66.74<math>\pm</math>0.09</b>	<u>66.51<math>\pm</math>0.03</u>	66.30 $\pm$ 0.13	64.05 $\pm$ 0.17	63.71 $\pm$ 0.10
AVG	51.06 $\pm$ 0.07	51.48 $\pm$ 0.13	51.51 $\pm$ 0.07	<b>52.43<math>\pm</math>0.07</b>	50.41 $\pm$ 0.19	44.62 $\pm$ 0.11	44.11 $\pm$ 0.23

Table 9: F1 scores (%) on UPB v1.0 and UPB v2.0 test sets with gold dependency trees (frozen mBERT). The bold score in each language indicates the highest F1 score among GNN-based models, i.e., GGCNs, SATs, GATs, and 2ATT-GATs. The underlined score in each language indicates the highest F1 score among 2ATT-GATs, BiLSTM+2ATT-GATs, BiLSTM, and 3BiLSTMs. AVG indicates the average F1 score of each model for all languages except English.

	2layers+ mBERT	2layers+ mBERT+ fine-tuned	2layers+ XLM-R	2layers+ XLM-R+ fine-tuned	3layers+ mBERT	3layers+ mBERT+ fine-tuned	3layers+ XLM-R	3layers+ XLM-R+ fine-tuned
UPB v1.0								
AVG (gold)	<b>56.45<math>\pm</math>0.05</b>	55.00 $\pm$ 0.05	55.87 $\pm$ 0.15	54.33 $\pm$ 0.10	56.15 $\pm$ 0.11	55.05 $\pm$ 0.08	56.00 $\pm$ 0.13	54.45 $\pm$ 0.13
AVG (pred)	<b>53.34<math>\pm</math>0.05</b>	52.41 $\pm$ 0.033	52.93 $\pm$ 0.13	51.95 $\pm$ 0.11	52.99 $\pm$ 0.11	52.43 $\pm$ 0.08	52.96 $\pm$ 0.09	51.95 $\pm$ 0.11
UPB v2.0								
AVG (gold)	52.41 $\pm$ 0.10	54.63 $\pm$ 0.11	54.05 $\pm$ 0.082	<b>55.34<math>\pm</math>0.06</b>	52.67 $\pm$ 0.08	54.30 $\pm$ 0.09	53.71 $\pm$ 0.08	55.26 $\pm$ 0.08
AVG (pred)	51.64 $\pm$ 0.12	53.95 $\pm$ 0.10	53.25 $\pm$ 0.12	<b>54.93<math>\pm</math>0.09</b>	51.93 $\pm$ 0.07	53.62 $\pm$ 0.07	52.89 $\pm$ 0.09	54.72 $\pm$ 0.08

Table 10: Average F1 scores (%) of target languages in UPB v1.0 and UPB v2.0 dev sets using gold and predicted dependency trees.