

# Graph Neural Networks for Syntax Encoding in Cross-Lingual Semantic Role Labeling

Anonymous ACL submission

## Abstract

Recent models in cross-lingual semantic role labeling (SRL) barely consider the applicability of their network selection. They rely on LSTMs as their encoders, even though LSTMs do not transfer effectively to distant languages. We evaluate the effectiveness of different graph neural networks (GNNs) enriched with universal dependency trees, i.e., transformer-based, graph convolutional network-based, and graph attention network (GAT)-based models, and compare them with a BiLSTM-based model. We investigate which dependency-aware GNNs transfer best as an alternative encoder to LSTMs in cross-lingual SRL. We focus our study on a zero-shot setting by training the models in English and evaluating the models in 23 target languages in the Universal Proposition Bank. We consistently show that syntax from universal dependency trees is essential for cross-lingual SRL models to achieve better transferability. Dependency-aware self-attention with relative position representations (SAN-RPRs) transfer best across languages, especially in the long-range dependency distance. Furthermore, our proposed dependency-aware two-attention relational GATs perform better than SAN-RPRs in languages where most arguments lie in the 1 – 2 dependency distance.

## 1 Introduction

Semantic role labeling (SRL) is a task to assign semantic roles to words or phrases in a sentence concerning a specific predicate, as shown in Figure 1. SRL supports many natural language processing (NLP) tasks, e.g., information extraction (Christensen et al., 2010; Stanovsky and Dagan, 2016), abstractive summarization (Khan et al., 2015), and machine translation (Rapp, 2022). However, SRL resource availability is still low, hindering the performance of other NLP tasks in diverse languages. Cross-lingual SRL models try to solve the problem by training the model in resource-rich languages and applying the model to resource-poor languages.

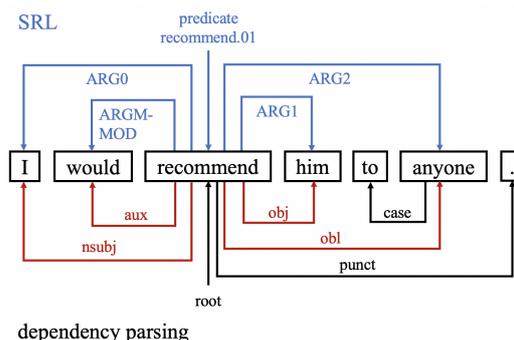


Figure 1: An example of an SRL task (top) and dependency parsing task (bottom) applied to a sentence taken from UPB. The red color indicates path intersections in both tasks.

Recent studies in cross-lingual SRL try to improve the model performance by separating language-universal and language-specific components (Fei et al., 2020; Conia et al., 2021) or minimizing dependence on external tools (Cai and Lapata, 2020). Nevertheless, these models barely consider the applicability of their network selection. They rely on LSTMs as their encoders, even though a study in cross-lingual dependency parsing (Ahmad et al., 2019) has stated that LSTMs do not transfer effectively to distant languages. LSTMs encode sentences sequentially, making them sensitive to word orders that vary across languages.

In this work, we investigate the effectiveness of various dependency-aware graph neural networks (GNNs) to find an alternative encoder for building cross-lingual SRL models. We encode universal dependency trees in dependency-aware GNNs, for the following 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 provides a general representation across languages. (3) Universal de-

069 dependency trees help cross-lingual models achieve  
 070 better transferability (Ahmad et al., 2021a; Ahmad  
 071 et al., 2021b; Zhang et al., 2021).

072 We conduct comprehensive experiments on vari-  
 073 ous networks as encoders, including transformer-  
 074 based, graph convolutional network (GCN)-  
 075 based, graph attention network (GAT)-based, and  
 076 BiLSTM-based encoders. We choose to investi-  
 077 gate transformer-based models because they have  
 078 been proven effective in performing cross-lingually  
 079 in dependency parsing (Ahmad et al., 2019) and  
 080 event argument role labeling (EARL) (Ahmad et al.,  
 081 2021b). Furthermore, we also investigate GCN-  
 082 based and GAT-based models because different  
 083 NLP tasks, e.g., monolingual SRL (Marcheggiani  
 084 and Titov, 2017), aspect-based sentiment analysis  
 085 (ABSA) (Wang et al., 2020; Jiang et al., 2021),  
 086 and relation prediction (Nathani et al., 2019), have  
 087 shown the effectiveness of exploiting the networks  
 088 to encode dependency trees in their models.

089 Following previous work (Fei et al., 2020), we  
 090 limit our exploration to argument detection and ar-  
 091 gument labeling in the dependency-based SRL. We  
 092 conduct experiments in a zero-shot setting to find  
 093 the most transferable network across languages.  
 094 We train and evaluate the models in 23 languages  
 095 provided by Universal Proposition Bank (UPB) v2.  
 096 We show that: (1) Universal dependency trees are  
 097 essential for cross-lingual SRL models to achieve  
 098 better transferability. (2) Transformer-based model  
 099 with dependency relation embedding (DRE) in the  
 100 node representation and relative position represen-  
 101 tation (RPR) in the edge representation, i.e., SAN-  
 102 RPRs, outperforms other models, especially as the  
 103 dependency distance increases. (3) Two-attention  
 104 relational GATs (TAGATs) with structural absolute  
 105 position embedding (SAPE) in the node represen-  
 106 tation, and also dependency relation representation  
 107 (DR) and RPR in the edge representation, perform  
 108 best in languages where most arguments lie in the  
 109 1 – 2 dependency distance.

## 110 2 Background

### 111 2.1 Universal Proposition Bank

112 Universal Proposition Bank (UPB) is a corpus  
 113 containing SRL annotations for diverse languages.  
 114 UPB v2 (Jindal et al., 2022) provides SRL an-  
 115 notations for 43 treebanks consisting of 23 lan-  
 116 guages, shown in Table 1. UPB is annotated semi-  
 117 automatically through filtered annotation projec-  
 118 tion and bootstrap training (Akbik et al., 2015).

Target Languages in UPB v2		
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 v2.

119 UPB v2 has significantly improved over UPB v1  
 120 regarding SRL annotation quality, language scope,  
 121 and availability of span-based SRL annotations  
 122 (Jindal et al., 2022). We use dependency-based  
 123 SRL annotations in UPB v2 that are annotated ac-  
 124 cording to UD v2.9 throughout our experiments.

### 125 2.2 Universal Dependencies

126 Universal Dependencies (UD) is a corpus contain-  
 127 ing consistent syntactic annotations for diverse lan-  
 128 guages, i.e., part-of-speech (POS) tags, morpho-  
 129 logical features, and dependency tree annotations.  
 130 UD v1 (Nivre et al., 2016) and UD v2 (Nivre et al.,  
 131 2020) have different annotation schemes<sup>1</sup> in terms  
 132 of word segmentation, POS tags, morphological  
 133 features, and syntactic relations. UD v1 and UD  
 134 v2 have 40<sup>2</sup> and 37<sup>3</sup> universal dependency rela-  
 135 tions, respectively. UD v2.9 (Zeman et al., 2021a)  
 136 contains dependency tree annotations for 217 tree-  
 137 banks of 122 languages.

### 138 2.3 Dependency-based Semantic Role 139 Labeling

140 Instead of labeling the whole argument span with a  
 141 semantic role, dependency-based SRL only labels  
 142 the argument head, i.e., the head of the argument  
 143 span according to the dependency tree. For ex-  
 144 ample, in Figure 1, the phrase “to anyone” is the  
 145 “ARG2” argument of the predicate “recommend”.  
 146 Based on the dependency tree at the bottom of the  
 147 figure, “anyone” is the head of the phrase “to any-  
 148 one”. Therefore, dependency-based SRL annotates  
 149 the edge that connects “recommend” to “anyone”  
 150 with the “ARG2” argument.

### 151 2.4 Related Work

152 Recent cross-lingual SRL models encode sentences  
 153 sequentially using BiLSTM-based models, even  
 154 though LSTMs do not transfer effectively across

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

<sup>2</sup><https://universaldependencies.org/docsv1/u/dep/>

<sup>3</sup><https://universaldependencies.org/u/dep/>

languages. Fei et al. (2020) introduce a parameter generation network to vanilla BiLSTMs that isolates language-specific parameters from the universal parameters in BiLSTMs. Meanwhile, Cai and Lapata (2020) exploit BiLSTM-based semantic role labeler and compressor to build a cross-lingual SRL model and utilize parallel sentences to help the semantic role compressor achieve generalization. Finally, Conia et al. (2021) introduce a BiLSTM-based universal sentence encoder and predicate-argument encoder, and also language-specific decoders to build a universal SRL model.

GNNs have been used to encode dependency trees in building models for different NLP tasks, e.g., monolingual SRL, ABSA, EARL, and relation prediction. In monolingual SRL, Marcheggiani and Titov (2017) employ syntactic GCNs (SGCNs) on top of BiLSTMs to incorporate dependency trees as graphs. In ABSA, Wang et al. (2020) and Jiang et al. (2021) apply relational GATs (R-GATs) and attention-based relational GCNs (ARGCNs), respectively, over modified dependency trees to establish direct connections between aspects and their corresponding words. In EARL, Ahmad et al. (2021b) modify vanilla Transformers to encode syntactic structures from dependency trees, i.e., graph attention transformer encoders (GATES). Finally, Nathani et al. (2019) propose KBGATs as a modification to original GATs (Veličković et al., 2018) to encode nodes and edge relations for relation prediction in knowledge graphs.

### 3 Model

#### 3.1 Architecture

We apply an encoder-decoder architecture to compare transformer-based, GCN-based, GAT-based, and BiLSTM-based cross-lingual SRL models, consisting of an input layer, an encoder, and a decoder.

##### 3.1.1 Input Layer

To produce the final word representation,  $h_i$ , for each word in a sentence, we optionally concatenate: (1) Predicate indicator embedding (PIE),  $p_i$ , represents whether a word is a predicate or not (Fei et al., 2020). (2) POS tag embedding (POSE),  $o_i$ , the POS tag of each word. (3) Absolute position embedding (APE),  $a_i$ , the position of each word in the sentence (Vaswani et al., 2017). (4) Structural absolute position embedding (SAPE),  $s_i$ , the dependency depth of each word relative to the root of the dependency tree (Wang et al., 2019b). (5)

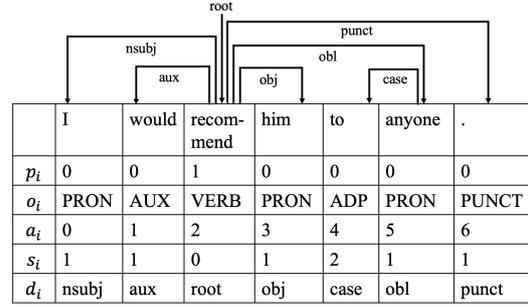


Figure 2: The illustration of predicate indicator, POS tag, absolute position, structural absolute position, and dependency relation of each word in a sentence.

Dependency relation embedding (DRE),  $d_i$ , the dependency relation of each word as a dependent. (6) Contextualized multilingual word embedding,  $c_i$ , obtained from the concatenation of the last four hidden layers (Conia et al., 2021) in multilingual BERT (mBERT) (Devlin et al., 2019). Figure 2 shows the example to obtain each embedding. As shown in Equation 1, we apply a dropout to the concatenation result, and then a linear transformation following GATES (Ahmad et al., 2021b).

$$h_i = W(\text{dropout}([p_i, o_i, a_i, s_i, d_i, c_i])) + b \quad (1)$$

#### 3.1.2 Encoder

We experiment with transformer-based, GCN-based, GAT-based, and BiLSTM-based encoders. For GCN-based and GAT-based encoders, we apply an activation function, dropout, and residual connection in consecutive order after each layer. As proposed by GATs (Veličković et al., 2018), we can either apply or drop the activation function at the final layer. For the BiLSTM-based encoder, we only apply the consecutive operations in the last layer. Through our comprehensive experiments, we find the best combination of node and edge representation in each encoder, as shown in Table 2.

**Transformer-Based Encoders** In transformer-based encoders, we experiment with different modifications to vanilla Transformers, i.e., self-attention with APE in the node representation (Trans) (Vaswani et al., 2017), GATES (Ahmad et al., 2021b), self-attention with RPR in the edge representation (SAN-RPRs) (Shaw et al., 2018), self-attention with SAPE in the node representation (SAN-SAPRs) (Wang et al., 2019b),

Model	Node Representation			Edge Representation		
	APE	SAPE	DRE	RPR	DR	SRPR
GATEs	x	✓	✓	x	x	x
Trans	✓	x	✓	x	x	x
SAN-RPRs	x	x	✓	✓	x	x
SAN-SAPRs	x	✓	x	x	x	x
SAN-SRPRs	x	x	x	x	x	✓
Trans-SRPRs	✓	x	✓	x	x	✓
SAPR-RPRs	x	✓	x	✓	x	x
Trans-SRPR-DRs	✓	x	✓	x	B	✓
SAPR-RPR-DRs	x	✓	x	x	B	x
SGCNs	x	x	✓	x	A	x
R-GCNs	x	✓	x	x	A	x
ARGCNs	x	x	✓	✓	A	x
GATs	x	✓	✓	x	x	x
S-HGNs	x	x	x	x	A	x
TAGATs	x	✓	x	✓	A	x
KBGATs	x	x	✓	✓	B	x
BiLSTMs	x	✓	✓	x	x	x

Table 2: Best combination of embeddings in node and representations in edge for each model in transformer-based, GCN-based, and GAT-based models.

and self-attention with structural relative position representation (SRPR) in the edge representation (SAN-SRPRs) (Wang et al., 2019b). We also combine the node representation with the edge representation, i.e., SAN-SAPRs with SAN-RPRs (SAPR-RPRs) and Transformers with SAN-SRPRs (Trans-SRPRs). Following Ahmad et al. (2019), we take the absolute value when calculating RPR and SRPR to make the models more robust to word order differences. Furthermore, we experiment with incorporating dependency relation representation (DR) into the edge representation. We explain how to construct DR in Section 3.2.

**GCN-Based Encoders** In GCN-based encoders, we experiment with different types of GCNs, i.e., syntactic GCNs (SGCNs) (Marcheggiani and Titov, 2017), relational GCNs (RGCNs) (Schlichtkrull et al., 2018), and attention-based relational GCNs (ARGCNs) (Jiang et al., 2021). Since SGCNs and RGCNs encode the edge representation by differentiating their weight matrices based on the edge relations, they can only incorporate A-DR (refer to Section 3.2).

**GAT-Based Encoders** In GAT-based encoders, we experiment with different types of GATs, i.e., GATs (Veličković et al., 2018), simple heterogeneous GNNs (SHGNs) (Lv et al., 2021), relational GATs (RGATs) (Wang et al., 2020), and knowledge-based GATs (KBGATs) (Nathani et al., 2019). RGATs calculate the second attention weight,  $\beta$ , using position-wise feed-forward network (FFN) (Vaswani et al., 2017). However, we find that the original dot-product equation

used to calculate the attention weight proposed by Veličković et al. (2018) works better for this task. Therefore, we modify RGATs into TAGATs. To calculate the second attention weight,  $\beta$ , TAGATs slightly modify Equation 3 in the original GAT paper (Veličković et al., 2018) to incorporate the DR,  $r_{ij}$ , as shown in Equation 2, where  $\mathcal{N}_i$  is neighbor nodes of node  $i$ ,  $W_r$  is a weight matrix to linearly transform the DR,  $r_{ij}$ , and LR is a Leaky ReLU. We provide a detailed explanation of TAGATs in Appendix C.

$$\beta_{ij} = \text{softmax}_{j \in \mathcal{N}_i}(\text{LR}(a^T[W_r r_{ij}])) \quad (2)$$

### 3.1.3 Decoder

We apply a linear scorer as the decoder. For each word, we concatenate sentence representation,  $h_s$ , predicate node representation,  $h_p$ , and node representation,  $h_i$ .  $h_s$  is obtained by applying max-pooling over node representations in the sentence (Ahmad et al., 2021b). Meanwhile,  $h_p$  is taken from the node representation of the sentence’s predicate. After that, following GATEs (Ahmad et al., 2021b), we apply two feed-forward neural networks (FFNNs), each followed by a ReLU (RL), to produce the final node representation,  $h_f$ , with the dimension equal to the number of arguments,  $c$ , as shown in Equation 3.

$$h_f = \text{RL}(W_2(\text{RL}(W_1[h_s, h_p, h_i] + b_1)) + b_2) \quad (3)$$

Finally, we apply a softmax function to produce the probability for each argument,  $z$ , as shown in Equation 4. We train the model to minimize the cross-entropy loss.

$$P(z) = \text{softmax}_z(h_f), z \in [1, c] \quad (4)$$

## 3.2 Graph Construction

First, we explain how to construct dependency relation representation (DR) to be encoded in the edges of the graphs. DR encodes the dependency direction, i.e., self-connection, head-to-dependent, and dependent-to-head, and the dependency relation between a pair of nodes. There are two ways of generating DR to be encoded in the edge representation of the graphs, i.e., A-DR and B-DR. In A-DR, following Marcheggiani and Titov (2017), we assign two completely different representations for edges with the same dependency relation but have different dependency directions (Appendix B). In B-DR, we first generate the representations

separately for each edge’s dependency relation and direction with  $d_r$  and  $d_d$  dimensions, respectively. Then, we concatenate both representations to produce the DR with  $(d_r + d_d)$  dimension.

In transformer-based models, we encode a sentence as a **fully-connected graph**. Shaw et al. (2018) modify the vanilla Transformers to incorporate RPR among the edges. We use the same approach to encode the DR by adding the DR to the RPR. Edges in the fully-connected graph that do not have the corresponding edges in the dependency tree do not have dependency relations. Therefore, we label these edges as “norel” (short for no relation) when constructing DRs.

In GCN-based and GAT-based models, we encode a sentence by forming a **dependency graph** based on its dependency tree. We follow the method proposed by Marcheggiani and Titov (2017) (Appendix B). They convert a dependency tree to a graph by adding edges that flow in the opposite direction of the original dependency direction and edges that flow from nodes to themselves. Furthermore, we encode either A-DR or B-DR in the edge representation.

## 4 Experiments

### 4.1 Corpus

We conduct experiments using corpus from UPB v2<sup>4</sup> (Jindal et al., 2022). UPB v2 contains SRL annotations based on dependency tree annotations in UD v2.9 (Zeman et al., 2021a). Some treebanks in UPB v2 have enhanced dependency tree annotations that cause new tokens (i.e., enhanced tokens) to be added to the sentences. The enhanced tokens cause some SRL annotations in UPB v2 to be shifted when merged with UD v2.9, resulting in the shifted predicate or semantic role annotations. Therefore, we run some preprocessing steps to fix the shifted annotation problem (Appendix A.2). When running our experiments, we merge all treebanks that belong to the same language.

### 4.2 Settings

We focus on conducting experiments in a zero-shot setting. We train the model in English and evaluate the model in 23 target languages provided by UPB v2, as shown in Table 1. In each experiment, we choose the final model from the epoch whose model performs best in the English validation set. However, we take the average F1 scores from the

<sup>4</sup><https://github.com/UniversalPropositions>

validation sets of 23 languages when choosing the best hyperparameter setting. We run exhaustive experiments to find the best hyperparameter setting to obtain the most transferable cross-lingual SRL model (Appendix D.2).

We use predicted dependency trees and POS tags for model evaluation (Appendix D.1). To obtain the predicted dependency trees and POS tags, we use pre-trained models trained on UD 2.8<sup>5</sup> (Zeman et al., 2021b) provided by Stanza (Qi et al., 2020).

We train the models for 100 epochs with 32 batch size. We use SGD optimizer (Kiefer and Wolfowitz, 1952) with a 0.1 learning rate. Following GATEs (Ahmad et al., 2021b), we apply early stopping if there is no improvement after 20 consecutive epochs. If the validation performance decreases, we decrease the learning rate by 10%. The training process will be stopped if the learning rate falls below 0.00001 after the decrement. We freeze mBERT as our contextualized word embedding to observe each model’s effectiveness solely. Our preliminary experiments show that taking the average of subword embeddings from mBERT works best for our task. Therefore, we are going to apply this setting to our experiments. We report the average F1 scores from five runs with the standard deviation for model comparison.

### 4.3 Comparison Among Transformer-Based Models

Table 3 compares transformer-based models with the best hyperparameter setting. Trans-SRPR-DRs and SAPR-RPR-DRs perform worst among the transformer-based models. To obtain DR in the fully-connected graph, we label the edges that do not reflect the edges in the corresponding dependency tree as “norel”. However, there are a lot more edges labeled as “norel” than edges that are labeled as the actual dependency relations. The imbalanced proportion of dependency relations might cause the model to overfit the limited actual dependency relations.

Furthermore, the table shows the superiority of SAN-RPRs over Trans, indicating that encoding the position of each word relative to another word in the edge representation produces a more general model than encoding the absolute position of

<sup>5</sup>We train the models from scratch for Japanese-GSDLUW and French-Rhapsodie treebanks because the pre-trained models are unavailable, and for English-EWT treebank because the SRL annotations are annotated based on an older version of UD, i.e., UD v2.5.

Model	PR	EN	AVG
GATEs	12.2M	78.96 $\pm$ 0.31	52.57 $\pm$ 0.23
Trans	12.2M	76.16 $\pm$ 0.51	52.07 $\pm$ 0.21
SAN-RPRs	12.2M	78.26 $\pm$ 0.40	<b>52.73<math>\pm</math>0.40</b>
SAN-SAPRs	12.2M	75.93 $\pm$ 0.47	51.98 $\pm$ 0.10
SAN-SRPRs	12.1M	78.27 $\pm$ 0.50	51.51 $\pm$ 0.27
Trans-SRPRs	12.2M	79.03 $\pm$ 0.32	52.21 $\pm$ 0.30
SAPR-RPRs	12.2M	78.11 $\pm$ 0.42	52.64 $\pm$ 0.38
Trans-SRPR-DRs	12.2M	79.85 $\pm$ 0.21	50.60 $\pm$ 0.14
SAPR-RPR-DRs	12.2M	79.83 $\pm$ 0.19	50.69 $\pm$ 0.21
SGCNs	5.99M	79.94 $\pm$ 0.27	<b>52.52<math>\pm</math>0.38</b>
R-GCNs	3.23M	78.28 $\pm$ 0.29	51.48 $\pm$ 0.35
ARGCNs	3.52M	77.84 $\pm$ 0.44	52.13 $\pm$ 0.32
GATs	5.09M	79.81 $\pm$ 0.19	52.66 $\pm$ 0.14
S-HGNs	5.06M	78.84 $\pm$ 0.35	52.61 $\pm$ 0.26
TAGATs	6.31M	79.07 $\pm$ 0.19	<b>52.78<math>\pm</math>0.14</b>
KBGATs	7.73M	79.53 $\pm$ 0.31	52.31 $\pm$ 0.32

Table 3: F1 scores (%) of transformer-based, GCN-based, GAT-based models evaluated on UPB v2 test set with predicted parsers. AVG indicates the average F1 scores of a specific model evaluated in target languages. The bold score and underlined score indicate the highest and second-highest scores in each group.

each word in the sentence in the node representation. This finding aligns with the results in [Ahmad et al. \(2019\)](#). On the other hand, SAN-SAPRs outperform SAN-SRPRs, indicating that encoding the dependency distance of each word relative to the root of the dependency tree in the node representation produces a more general model than encoding each word’s dependency distance relative to another word in the edge representation. We further combine the information regarding the position of each word according to the sentence and the sentence’s dependency tree, i.e., SAPR-RPRs and Trans-SRPRs. Consistent with the previous results, SAPR-RPRs, consisting of features from SAN-SAPRs and SAN-RPRs, outperform Trans-SRPRs, consisting of features from Trans and SAN-SRPRs.

According to the average F1 score, SAN-RPRs and SAPR-RPRs are strong models with 52.73% and 52.64% average F1 scores, respectively, outperforming GATEs ([Ahmad et al., 2021b](#)), i.e., the model proposed for cross-lingual EARL. Furthermore, SAN-RPRs which employ DRE in the node representation, perform better than SAPR-RPRs. Nevertheless, SAPR-RPRs perform best among transformer-based models in more languages than SAN-RPRs (Appendix E.1.1). We take both models for comparison with other best models in Section 4.6.

#### 4.4 Comparison Among GCN-Based Models

Table 3 compares GCN-based models with the best hyperparameter setting. SGCNs significantly outperform the other GCN-based models, i.e., RGCNs and ARGCNs, with a 52.52% average F1 score. RGCNs perform the worst among GCN-based models because RGCNs are the only network not applying the attention mechanism. ARGCNs apply a self-attention mechanism, while SGCNs realize the attention mechanism in the form of a gating mechanism. Self-attention or gating mechanism measures how much attention each node should pay to other nodes when the network updates each node’s representation. Those mechanisms help emphasize the edges in the dependency graph that intersect with the predicate-argument paths.

#### 4.5 Comparison Among GAT-Based Models

Table 3 compares GAT-based models with the best hyperparameter setting. TAGATs perform better than SHGNs indicating that the GAT-based model learns better when we separate the attention weight calculation based on node representations and edge representations. Moreover, GATs also perform better than SHGNs indicating that encoding the SAPE and DRE in the node representation is a better way to encode dependency features than combining node representations with edge representations when calculating the attention weight. According to the average F1 score, TAGATs and GATs are strong models with 52.78% and 52.66% average F1 scores, respectively. Therefore, we compare both models with other best models in Section 4.6.

On the other hand, according to the average F1 score, KBGATs perform worst among GAT-based models. The significant difference between KBGATs and the others is that KBGATs update each node representation with (1) neighbor node representations and (2) surrounding edges’ representations that contain information about DR and RPR. Meanwhile, the other models update each node representation only with (1). We conjecture that the approach of KBGATs might overpopulate each node in every update with information too specific to a particular language the network learns from.

#### 4.6 Comparison Among Best Models

We compare the best models from each group, i.e., SAN-RPRs, SAPR-RPRs, SGCNs, GATs, and TAGATs, with the BiLSTM-based model, i.e., BiLSTMs, in Table 4. We calculate each model’s supe-

	$d \geq 3$	SAN-RPRs	SAPR-RPRs	SGCNs	GATs	TAGATs	BiLSTMs
EN	2.68	78.26 $\pm$ 0.40	78.11 $\pm$ 0.42	79.94 $\pm$ 0.27	79.81 $\pm$ 0.19	79.07 $\pm$ 0.19	76.86 $\pm$ 0.34
AVG	-	52.73 $\pm$ 0.40	52.64 $\pm$ 0.38	52.52 $\pm$ 0.38	52.66 $\pm$ 0.14	<b>52.78<math>\pm</math>0.14</b>	51.85 $\pm$ 0.09
TA	17.18	37.96 $\pm$ 1.75	<b>39.57<math>\pm</math>1.18</b>	34.32 $\pm$ 1.12	35.08 $\pm$ 0.58	35.68 $\pm$ 1.28	34.19 $\pm$ 1.22
HI	8.46	47.51 $\pm$ 0.62	45.04 $\pm$ 0.35	48.24 $\pm$ 0.61	<b>48.25<math>\pm</math>0.33</b>	47.65 $\pm$ 0.33	46.63 $\pm$ 0.38
ZH	8.41	50.37 $\pm$ 1.17	<b>50.96<math>\pm</math>0.88</b>	45.77 $\pm$ 0.67	46.12 $\pm$ 0.39	46.80 $\pm$ 0.64	47.56 $\pm$ 0.89
JA	8.11	37.69 $\pm$ 0.99	34.78 $\pm$ 1.29	37.43 $\pm$ 0.28	37.99 $\pm$ 0.52	<b>39.30<math>\pm</math>0.73</b>	37.40 $\pm$ 0.61
VI	7.89	28.69 $\pm$ 0.79	<b>29.10<math>\pm</math>0.45</b>	27.95 $\pm$ 0.56	28.06 $\pm$ 0.59	28.31 $\pm$ 0.55	28.18 $\pm$ 0.88
KO	6.88	42.61 $\pm$ 1.84	<b>45.24<math>\pm</math>1.23</b>	42.92 $\pm$ 0.64	43.22 $\pm$ 0.56	44.57 $\pm$ 0.24	41.77 $\pm$ 1.61
ID	5.52	58.78 $\pm$ 1.09	<b>59.97<math>\pm</math>0.53</b>	58.54 $\pm$ 0.82	58.33 $\pm$ 0.69	59.11 $\pm$ 0.87	56.11 $\pm$ 0.84
HU	5.38	49.76 $\pm$ 0.35	49.08 $\pm$ 0.34	50.76 $\pm$ 0.41	<b>51.10<math>\pm</math>0.51</b>	50.90 $\pm$ 0.37	50.64 $\pm$ 0.39
RO	5.32	54.23 $\pm$ 0.67	<b>54.46<math>\pm</math>0.52</b>	53.57 $\pm$ 0.47	54.12 $\pm$ 0.49	53.60 $\pm$ 0.45	53.26 $\pm$ 0.34
FR	4.55	<b>62.19<math>\pm</math>0.41</b>	62.11 $\pm$ 0.47	60.93 $\pm$ 0.38	61.64 $\pm$ 0.44	61.13 $\pm$ 0.22	61.22 $\pm$ 0.27
MR	4.08	<b>41.06<math>\pm</math>2.89</b>	40.36 $\pm$ 2.20	40.97 $\pm$ 3.40	38.06 $\pm$ 0.13	39.26 $\pm$ 1.20	37.18 $\pm$ 2.28
UK	4.06	58.92 $\pm$ 0.26	59.36 $\pm$ 0.72	59.66 $\pm$ 0.76	59.49 $\pm$ 0.56	<b>59.72<math>\pm</math>0.31</b>	58.96 $\pm$ 0.07
PT	3.75	66.05 $\pm$ 0.21	<b>66.49<math>\pm</math>0.33</b>	65.62 $\pm$ 0.43	65.99 $\pm$ 0.32	65.61 $\pm$ 0.15	64.40 $\pm$ 0.33
IT	3.73	<b>58.11<math>\pm</math>0.33</b>	57.80 $\pm$ 0.39	57.43 $\pm$ 0.42	58.00 $\pm$ 0.42	57.34 $\pm$ 0.34	58.02 $\pm$ 0.27
ES	3.67	63.71 $\pm$ 0.33	63.62 $\pm$ 0.25	63.87 $\pm$ 0.61	<b>64.29<math>\pm</math>0.36</b>	63.91 $\pm$ 0.27	62.48 $\pm$ 0.29
CS	3.66	56.87 $\pm$ 0.27	55.80 $\pm$ 0.51	57.95 $\pm$ 0.52	<b>58.02<math>\pm</math>0.21</b>	57.62 $\pm$ 0.28	56.59 $\pm$ 0.36
EL	3.59	60.59 $\pm$ 0.23	60.23 $\pm$ 0.40	60.56 $\pm$ 0.69	60.74 $\pm$ 0.48	<b>60.86<math>\pm</math>0.34</b>	59.76 $\pm$ 0.45
FI	3.35	<b>55.58<math>\pm</math>0.42</b>	55.29 $\pm$ 0.54	54.87 $\pm$ 0.40	54.62 $\pm$ 0.32	54.88 $\pm$ 0.20	54.62 $\pm$ 0.20
RU	3.07	60.18 $\pm$ 0.44	<b>61.13<math>\pm</math>0.50</b>	59.98 $\pm$ 0.34	60.14 $\pm$ 0.16	60.30 $\pm$ 0.22	59.73 $\pm$ 0.25
NL	3.05	62.84 $\pm$ 0.37	62.22 $\pm$ 0.58	62.94 $\pm$ 0.21	<b>63.53<math>\pm</math>0.64</b>	62.97 $\pm$ 0.37	62.47 $\pm$ 0.36
TE	2.49	44.66 $\pm$ 2.00	43.88 $\pm$ 1.57	46.08 $\pm$ 1.07	<b>46.96<math>\pm</math>1.49</b>	<b>46.96<math>\pm</math>1.82</b>	45.95 $\pm$ 0.51
DE	2.46	56.86 $\pm$ 0.32	56.98 $\pm$ 1.13	<b>58.61<math>\pm</math>0.30</b>	58.52 $\pm$ 0.26	58.52 $\pm$ 0.18	57.72 $\pm$ 0.23
PL	1.71	57.67 $\pm$ 0.36	57.28 $\pm$ 0.46	<b>59.08<math>\pm</math>0.31</b>	59.00 $\pm$ 0.44	58.92 $\pm$ 0.52	57.73 $\pm$ 0.23
SC	-	13	18	9	16	15	1
PR	-	12.2M	12.2M	5.99M	5.09M	6.31M	9.03M

Table 4: F1 scores (%) of best models evaluated on UPB v2 test set with predicted parsers. The bold score and underlined score indicate the highest and second-highest scores. AVG indicates the average F1 scores of a specific model evaluated in target languages. PR and SC are the number of parameters and the superiority score of each model.  $d \geq 3$  column indicates the proportion of gold arguments (%) that fall in  $\geq 3$  dependency distance. We use predicted dependency trees to measure the dependency distance.

riority score (SC) based on the model performance in target languages. We allocate 2 points if the model achieves the highest F1 score or 1 point if the model achieves the second-highest F1 score for a specific language.

TAGATs have the best average F1 score among the models with 52.78% average F1 score. However, SAPR-RPRs and GATs perform best among the models in slightly more languages, indicated by the higher SCs. Despite having the second-best average F1 score of 52.73%, SAN-RPRs have a lower SC than SAPR-RPRs, GATs, and TAGATs. Overall, transformer-based and GAT-based models outperform SGCNs, indicating that the self-attention mechanism is better at emphasizing the essential dependency paths than the gating mechanism. Finally, BiLSTMs perform the worst as the only network that encodes sentences sequentially.

According to hyperparameter search, stacking two layers of GAT-based encoders performs best. However, this approach has a drawback, as the information can only travel as far as two hops from the origin node. On the other hand, stacking three

layers of transformer-based encoders does not affect the traveling distance of the information. In transformer-based models, we construct a fully-connected graph of a sentence allowing the information to travel from one node to every other node despite how many layers are stacked together.

The second column in Table 4 shows the percentage of arguments in  $\geq 3$  dependency distance, i.e., the number of hops from the predicate node to a specific node according to the sentence’s dependency tree. Some languages like TA, HI, ZH, JA, and VI have a relatively high number of arguments ( $> 7\%$ ) in  $\geq 3$  dependency distance. Transformer-based models perform significantly better in TA ( $> 3\%$ ) and ZH ( $> 4\%$ ), slightly better in VI ( $< 1\%$ ), and slightly worse in HI ( $< 1\%$ ) and JA ( $< 2\%$ ). This evidence proves that transformer-based models are generally better in long-range dependency distance than GAT-based models. On the other hand, some languages have a relatively low number of arguments ( $< 3\%$ ) in  $\geq 3$  dependency distance, i.e., EN, TE, DE, and PL. For these languages, where most of the arguments lie in the

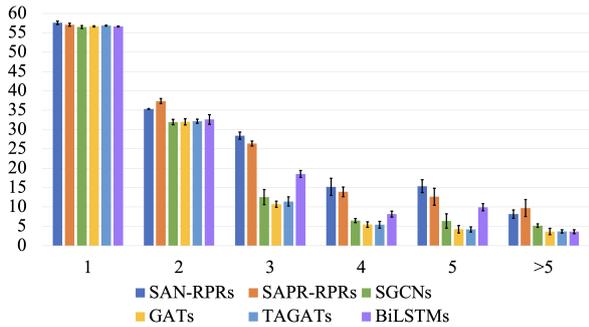


Figure 3: Average F1 scores (%) of best models grouped by the dependency distance evaluated on UPB v2 test set with predicted parsers.

Model	SAN-RPRs	SAPR-RPRs
base	$52.73 \pm 0.40$	$52.64 \pm 0.38$
w/o POSE	$52.73 \pm 0.32$	$51.91 \pm 0.46$
w/o PIE	$51.54 \pm 0.36$	$51.14 \pm 0.43$
w/o SAPE	-	$51.65 \pm 0.28$
w/o DRE	$51.65 \pm 0.28$	-

Table 5: Average F1 scores (%) of SAN-RPRs and SAPR-RPRs with certain embedding removed from the node representation evaluated on UPB v2 test set with predicted parsers.

1 – 2 dependency distance, GAT-based models perform better than transformer-based models.

Figure 3 shows the average F1 scores of each model grouped by the dependency distance. Although BiLSTMs perform worst among the models, BiLSTMs are better in long-range dependency distance than GCN-based and GAT-based models. The figure also shows the superiority of transformer-based models as the dependency distance increases, even when compared to BiLSTMs.

#### 4.7 Ablation Study

We conduct ablation studies for the best transformer-based models, i.e., SAN-RPRs and SAPR-RPRs, as shown in Table 5. We experiment with removing DRE, POSE, or PIE from the node representation in SAN-RPRs. Removing either DRE or PIE from the node representation reduces the performance of SAN-RPRs. However, SAN-RPRs without POSE perform better in most languages (Appendix E.1.2).

Furthermore, we also experiment with removing SAPE, POSE, or PIE from the node representation in SAPR-RPRs. Removing SAPE, POSE, or PIE from the node representation reduces the performance of SAPR-RPRs. Unlike SAN-RPRs, SAPR-RPRs do not have DRE in their node representation. We conjecture that the combination of POSE and

SAPE in SAPR-RPRs is necessary to replace the role of DRE.

In Table 4, according to the SC, we can see that SAPR-RPRs perform best in more languages than SAN-RPRs, even though the average F1 score of SAN-RPRs is better than SAPR-RPRs. However, as discussed, SAN-RPRs without POSE perform better in most languages than SAN-RPRs with POSE, which we use for comparison in Table 4. Therefore, we re-compare SAN-RPRs without POSE with the other best models in Table 4. After removing POSE from SAN-RPRs, we find that the model performs best in more languages than SAPR-RPRs (Appendix E.2.1).

## 5 Conclusions and Future Work

Through comprehensive experiments, we consistently show that incorporating syntax from dependency trees can improve the transferability of cross-lingual SRL models across languages. Overall, we show that the transformer-based model, i.e., SAN-RPRs that encode DRE without POSE in the node representation and RPR in the edge representation, stacked in three layers, performs the best among all models, especially as the dependency distance increases. However, TAGATs that encode SAPE in the node representation, and also DR and RPR in the edge representation, stacked in two layers, perform better than SAN-RPRs in languages where most of the arguments lie in the 1 – 2 dependency distance.

In the future, we can extend our model to incorporate language-specific components and modify the objective function to maximize the learning of universal features without ignoring the specific features that appear in each language. This can be useful if we want to extend the model training to a few-shot setting where we include a certain proportion of target sentences in the training set.

### 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.

## Ethical 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 to help solve the problems in this diverse society.

## References

- Wasi Ahmad, Haoran Li, Kai-Wei Chang, and Yashar Mehdad. 2021a. [Syntax-augmented multilingual BERT for cross-lingual transfer](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4538–4554, Online. Association for Computational Linguistics.
- 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.
- Wasi Uddin Ahmad, Nanyun Peng, and Kai-Wei Chang. 2021b. [Gate: Graph attention transformer encoder for cross-lingual relation and event extraction](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):12462–12470.
- 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.
- 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.

- 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.
- 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.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, page 1243–1252. JMLR.org.
- 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 summarization based on semantic role labelling](#). *Applied Soft Computing*, 30:737–747.
- Jack Kiefer and Jacob Wolfowitz. 1952. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*, pages 462–466.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations*,



831	Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia Ackermann, Noëmi Aepli, Hamid Aghaei, Željko Agić, Amir Ahmadi, Lars Ahrenberg, Chika Kennedy Ajede, Gabrielë Aleksandravičiūtė, Ika Alfina, Lene Antonsen, Katya Aplonova, Angelina Aquino, Carolina Aragon, Maria Jesus Aranzabe, Bilge Nas Arican, Hórunn Arnardóttir, Gashaw Arutie, Jessica Naraiswari Arwidarasti, Masayuki Asahara, Deniz Baran Aslan, Luma Ateyah, Furkan Atmaca, Mohammed Attia, Aitziber Atutxa, Liesbeth Augustinus, Elena Badmaeva, Keerthana Balasubramani, Miguel Ballesteros, Esha Banerjee, Sebastian Bank, Verginica Barbu Mititelu, Starkadur Barkarson, Rodolfo Basile, Victoria Basmov, Colin Batchelor, John Bauer, Seyyit Talha Bedir, Kepa Bengoetxea, Gözde Berk, Yevgeni Berzak, Irshad Ahmad Bhat, Riyaz Ahmad Bhat, Erica Biagetti, Eckhard Bick, Agnė Bielinskienė, Kristín Bjarnadóttir, Rogier Blokland, Victoria Bobicev, Loïc Boizou, Emanuel Borges Völker, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd, Anouck Braggaar, Kristina Brokaitė, Aljoscha Burchardt, Marie Candito, Bernard Caron, Gauthier Caron, Lauren Cassidy, Tatiana Cavalcanti, Gülşen Cebiroğlu Eryiğit, Flavio Massimiliano Cecchini, Giuseppe G. A. Celano, Slavomír Čéplö, Neslihan Cesur, Savas Cetin, Özlem Çetinoğlu, Fabricio Chalub, Shweta Chauhan, Ethan Chi, Taishi Chika, Yongseok Cho, Jinho Choi, Jayeol Chun, Juyeon Chung, Alessandra T. Cignarella, Silvie Cinková, Aurélie Collomb, Çağrı Çoltekin, Miriam Connor, Marine Courtin, Mihaela Cristescu, Philemon Daniel, Elizabeth Davidson, Marie-Catherine de Marneffe, Valeria de Paiva, Mehmet Oguz Derin, Elvis de Souza, Arantza Diaz de Ilaraza, Carly Dickerson, Arawinda Dinakaramani, Elisa Di Nuovo, Bamba Dione, Peter Dirix, Kaja Dobrovoljc, Timothy Dozat, Kira Droganova, Puneet Dwivedi, Hanne Eckhoff, Sandra Eiche, Marhaba Eli, Ali Elkahky, Binyam Ephrem, Olga Erina, Tomaz Erjavec, Aline Etienne, Wograine Evelyn, Sidney Facundes, Richárd Farkas, Jannatul Ferdousi, Marília Fernanda, Hector Fernandez Alcalde, Jennifer Foster, Cláudia Freitas, Kazunori Fujita, Katarína Gajdošová, Daniel Galbraith, Marcos Garcia, Moa Gärdenfors, Sebastian Garza, Fabrício Ferraz Gerardi, Kim Gerdes, Filip Ginter, Gustavo Godoy, Iakes Goenaga, Koldo Gojenola, Memduh Gökırmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Bernadeta Griciūtė, Matias Grioni, Loïc Grobol, Normunds Grūzītis, Bruno Guillaume, Céline Guillot-Barbance, Tunga Güngör, Nizar Habash, Hinrik Hafsteinsson, Jan Hajič, Jan Hajič jr., Mika Hämäläinen, Linh Hà Mỹ, Na-Rae Han, Muhammad Yudistira Hanifmuti, Sam Hardwick, Kim Harris, Dag Haug, Johannes Heinicke, Oliver Hellwig, Felix Hennig, Barbora Hladká, Jaroslava Hlaváčová, Florinel Hociung, Petter Hohle, Eva Huber, Jena Hwang, Takumi Ikeda, Anton Karl Ingason, Radu Ion, Elena Irimia, Olájídé Ishola, Kaoru Ito, Siratun Jannat, Tomáš Jelínek, Apoorva Jha, Anders Johannsen, Hildur Jónsdóttir, Fredrik Jørgensen, Markus Juutinen, Sarveswaran K, Hüner	Kaşıkara, Andre Kaasen, Nadezhda Kabaeva, Sylvain Kahane, Hiroshi Kanayama, Jenna Kanerva, Neslihan Kara, Boris Katz, Tolga Kayadelen, Jessica Kenney, Václava Kettnerová, Jesse Kirchner, Elena Klementieva, Elena Klyachko, Arne Köhn, Abdullatif Köksal, Kamil Kopacewicz, Timo Korkiakangas, Mehmet Köse, Natalia Kotsyba, Jolanta Kovalevskaitė, Simon Krek, Parameswari Krishnamurthy, Sandra Kübler, Oğuzhan Kuyrukçu, Asli Kuzgun, Sookyung Kwak, Veronika Laippala, Lucia Lam, Lorenzo Lambertino, Tatiana Lando, Septina Dian Larasati, Alexei Lavrentiev, John Lee, Phûông Lê Hồng, Alessandro Lenci, Saran Lertpradit, Herman Leung, Maria Levina, Cheuk Ying Li, Josie Li, Keying Li, Yuan Li, KyungTae Lim, Bruna Lima Padovani, Krister Lindén, Nikola Ljubešić, Olga Loginova, Stefano Lusito, Andry Luthfi, Mikko Luukko, Olga Lyashevskaya, Teresa Lynn, Vivien Macketanz, Menel Mahamdi, Jean Maillard, Aibek Makazhanov, Michael Mandl, Christopher Manning, Ruli Manurung, Büşra Marşan, Cătălina Mărănduc, David Mareček, Katrin Marheinecke, Héctor Martínez Alonso, Lorena Martín-Rodríguez, André Martins, Jan Mašek, Hiroshi Matsuda, Yuji Matsumoto, Alessandro Mazzei, Ryan McDonald, Sarah McGuinness, Gustavo Mendonça, Tatiana Merzhevich, Niko Miekka, Karina Mischenkova, Margarita Misirpashayeva, Anna Missilä, Cătălin Mititelu, Maria Mitrofan, Yusuke Miyao, AmirHosseini Mojiri Foroushani, Judit Molnár, Amirsaeid Moloodi, Simonetta Montemagni, Amir More, Laura Moreno Romero, Giovanni Moretti, Keiko Sophie Mori, Shinsuke Mori, Tomohiko Morioka, Shigeki Moro, Bjartur Mortensen, Bohdan Moskalevskiy, Kadri Muischnek, Robert Munro, Yugo Murawaki, Kaili Müürisep, Pinkey Nainwani, Mariam Nakhlé, Juan Ignacio Navarro Horňiáček, Anna Nedoluzhko, Gunta Nešpore-Bērzkalne, Manuela Nevaci, Lũông Nguyễn Thị, Huyên Nguyễn Thị Minh, Yoshihiro Nikaido, Vitaly Nikolaev, Rattima Nitisaroj, Alireza Nourian, Hanna Nurmi, Stina Ojala, Atul Kr. Ojha, Adédayò Olúòkun, Mai Omura, Emeka Onwuegbuzia, Petya Osenova, Robert Östling, Lilja Övrelid, Şaziye Betül Özateş, Merve Özçelik, Arzucan Özgür, Balkız Öztürk Başaran, Hyunji Hayley Park, Niko Partanen, Elena Pascual, Marco Passarotti, Agnieszka Patejuk, Guilherme Paulino-Passos, Angelika Peljak-Łapińska, Siyao Peng, Cenel-Augusto Perez, Natalia Perkova, Guy Perrier, Slav Petrov, Daria Petrova, Jason Phelan, Jussi Piitulainen, Tommi A Pirinen, Emily Pitler, Barbara Plank, Thierry Poibeau, Larisa Ponomareva, Martin Popel, Lauma Pretkalniņa, Sophie Prévost, Prokopis Prokopidis, Adam Przepiórkowski, Tiina Puolakainen, Sampo Pyysalo, Peng Qi, Andriela Rääbis, Alexandre Rademaker, Mizanur Rahoman, Taraka Rama, Loganathan Ramasamy, Carlos Ramisch, Fam Rashel, Mohammad Sadegh Rasooli, Vinit Ravishankar, Livy Real, Petru Rebeja, Siva Reddy, Mathilde Regnault, Georg Rehm, Ivan Riabov, Michael Rießler, Erika Rimkutė, Larissa Rinaldi, Laura Rituma, Putri Rizqiyah, Luisa Rocha, Eiríkur Rögnvaldsson, Mykhailo Romanenko, Rudolf	894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956
-----	--	--	---

957	Rosa, Valentin Roşca, Davide Rovati, Olga Rudina, Jack Rueter, Kristján Rúnarsson, Shoal Sadde,	1019
958	Pegah Safari, Benoît Sagot, Aleksí Sahala, Shadi Saleh, Alessio Salomoni, Tanja Samardžić, Stephanie Samson, Manuela Sanguinetti, Ezgi Sanıyar, Dage Särg, Baiba Saulīte, Yanin Sawanakunanon, Shefali Saxena, Kevin Scannell, Salvatore Scarlata, Nathan Schneider, Sebastian Schuster, Lane Schwartz, Djamé Seddah, Wolfgang Seeker, Mojgan Seraji, Syeda Shahzadi, Mo Shen, Atsuko Shimada, Hiroyuki Shirasu, Yana Shishkina, Muh Shohibussirri, Dmitry Sichinava, Janine Siewert, Einar Freyr Sigurdsson, Aline Silveira, Natalia Silveira, Maria Simi, Radu Simionescu, Katalin Simkó, Mária Šimková, Kiril Simov, Maria Skachedubova, Aaron Smith, Isabela Soares-Bastos, Shafi Sourov, Carolyn Spadine, Rachele Sprugnoli, Steinhórfur Steingrímsson, Antonio Stella, Milan Straka, Emmett Strickland, Jana Strnadová, Alane Suhr, Yogi Lesmana Sulestio, Umut Sulubacak, Shingo Suzuki, Zsolt Szántó, Chihiro Taguchi, Dima Taji, Yuta Takahashi, Fabio Tamburini, Mary Ann C. Tan, Takaaki Tanaka, Dipta Tanaya, Samson Tella, Isabelle Tellier, Marinella Testori, Guillaume Thomas, Liisi Torga, Marsida Toska, Trond Trosterud, Anna Trukhina, Reut Tsarfaty, Utku Türk, Francis Tyers, Sumire Uematsu, Roman Untilov, Zdeňka Urešová, Larraitz Uria, Hans Uszkoreit, Andrius Utka, Sowmya Vajjala, Rob van der Goot, Martine Vanhove, Daniel van Niekerk, Gertjan van Noord, Viktor Varga, Eric Villemonte de la Clergerie, Veronika Vincze, Natalia Vlasova, Aya Wakasa, Joel C. Wallenberg, Lars Wallin, Abigail Walsh, Jing Xian Wang, Jonathan North Washington, Maximilian Wendt, Paul Widmer, Sri Hartati Wijono, Seyi Williams, Mats Wirén, Christian Wittern, Tsegay Woldemariam, Tak-sum Wong, Alina Wróblewska, Mary Yako, Kayo Yamashita, Naoki Yamazaki, Chunxiao Yan, Koichi Yasuoka, Marat M. Yavrumyan, Arife Betül Yenice, Olcay Taner Yıldız, Zhuoran Yu, Arlisa Yuliawati, Zdeněk Žabokrtský, Shorouq Zahra, Amir Zeldes, He Zhou, Hanzhi Zhu, Anna Zhuravleva, and Rayan Ziane. 2021a. <a href="#">Universal dependencies 2.9</a> . LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.	1020
959		1021
960		1022
961		1023
962		1024
963		1025
964		1026
965		1027
966		1028
967		1029
968		1030
969		1031
970		1032
971		1033
972		1034
973		1035
974		1036
975		1037
976		1038
977		1039
978		1040
979		1041
980		1042
981		1043
982		1044
983		1045
984		1046
985		1047
986		1048
987		1049
988		1050
989		1051
990		1052
991		1053
992		1054
993		1055
994		1056
995		1057
996		1058
997		1059
998		1060
999		1061
1000		1062
1001		1063
1002		1064
1003	Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia Ackermann, Noëmi Aepli, Hamid Aghaei, Željko Agić, Amir Ahmadi, Lars Ahrenberg, Chika Kennedy Ajede, Gabrielė Aleksandravičiūtė, Ika Alfina, Lene Antonsen, Katya Aplonova, Angelina Aquino, Carolina Aragon, Maria Jesus Aranzabe, Bilge Nas Arıcan, Hórunn Arnardóttir, Gashaw Arutie, Jessica Naraiswari Arwidarasti, Masayuki Asahara, Deniz Baran Aslan, Luma Ateyah, Furkan Atmaca, Mohammed Attia, Aitziber Atutxa, Liesbeth Augustinus, Elena Badmaeva, Keerthana Balasubramani, Miguel Ballesteros, Esha Banerjee, Sebastian Bank, Verginica Barbu Mititelu, Starkadur Barkarson, Victoria Basmov, Colin Batchelor, John Bauer, Seyyit Talha Bedir, Kepa Bengoetxea, Gözde Berk, Yevgeni Berzak, Irshad Ahmad Bhat, Riyaz Ah-	1065
1004		1066
1005		1067
1006		1068
1007		1069
1008		1070
1009		1071
1010		1072
1011		1073
1012		1074
1013		1075
1014		1076
1015		1077
1016		1078
1017		1079
1018		1080
	mad Bhat, Erica Biagetti, Eckhard Bick, Agnė Bielinskienė, Kristín Bjarnadóttir, Rogier Blokland, Victoria Bobicev, Loïc Boizou, Emanuel Borges Völker, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd, Anouck Braggaa, Kristina Brokaitė, Aljoscha Burchardt, Marie Candito, Bernard Caron, Gauthier Caron, Lauren Cassidy, Tatiana Cavalcanti, Gülşen Cebiroğlu Eryiğit, Flavio Massimiliano Cecchini, Giuseppe G. A. Celano, Slavomír Čéplö, Neslihan Cesur, Savas Cetin, Özlem Çetinoğlu, Fabricio Chalub, Shweta Chauhan, Ethan Chi, Taishi Chika, Yongseok Cho, Jinho Choi, Jayeol Chun, Alessandra T. Cignarella, Silvie Cinková, Aurélie Collomb, Çağrı Çöltekin, Miriam Connor, Marine Courtin, Mihaela Cristescu, Philemon. Daniel, Elizabeth Davidson, Marie-Catherine de Marneffe, Valeria de Paiva, Mehmet Oguz Derin, Elvis de Souza, Arantza Diaz de Ilarraza, Carly Dickerson, Arawinda Dinakaramani, Elisa Di Nuovo, Bamba Dione, Peter Dirix, Kaja Dobrovoljc, Timothy Dozat, Kira Droganova, Puneet Dwivedi, Hanne Eckhoff, Sandra Eiche, Marhaba Eli, Ali Elkahky, Binyam Ephrem, Olga Erina, Tomaz Erjavec, Aline Etienne, Wograine Evelyn, Sidney Facundes, Richárd Farkas, Marília Fernanda, Hector Fernandez Alcalde, Jennifer Foster, Cláudia Freitas, Kazunori Fujita, Katarína Gajdošová, Daniel Galbraith, Marcos Garcia, Moa Gärdenfors, Sebastian Garza, Fabrício Ferraz Gerardi, Kim Gerdes, Filip Ginter, Gustavo Godoy, Iakes Goenaga, Koldo Gojenola, Memduh Gökırmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Bernadeta Griciūtė, Matias Gtrioni, Loïc Grobol, Normunds Grūzītis, Bruno Guillaume, Céline Guillot-Barbance, Tunga Güngör, Nizar Habash, Hinrik Hafsteinsson, Jan Hajič, Jan Hajič jr., Mika Hämäläinen, Linh Hà Mỹ, Na-Rae Han, Muhammad Yudistira Hanifmuti, Sam Hardwick, Kim Harris, Dag Haug, Johannes Heinicke, Oliver Hellwig, Felix Hennig, Barbora Hladká, Jaroslava Hlaváčová, Florinel Hociung, Petter Hohle, Eva Huber, Jena Hwang, Takumi Ikeda, Anton Karl Ingason, Radu Ion, Elena Irimia, Olájídé Ishola, Kaoru Ito, Tomáš Jelínek, Apoorva Jha, Anders Johannsen, Hildur Jónsdóttir, Fredrik Jørgensen, Markus Juutinen, Sarveswaran K, Hüner Kaşıkara, Andre Kaasen, Nadezhda Kabaeva, Sylvain Kahane, Hiroshi Kanayama, Jenna Kanerva, Neslihan Kara, Boris Katz, Tolga Kayadelen, Jessica Kenney, Václava Kettnerová, Jesse Kirchner, Elena Klementieva, Arne Köhn, Abdullatif Köksal, Kamil Kopacewicz, Timo Korhikangas, Natalia Kotsyba, Jolanta Kovalevskaitė, Simon Krek, Parameswari Krishnamurthy, Oğuzhan Kuyrukçu, Aslı Kuzgun, Sookyoung Kwak, Veronika Laippala, Lucia Lam, Lorenzo Lambertino, Tatiana Lando, Septina Dian Larasati, Alexei Lavrentiev, John Lee, Phùng Lê Hồng, Alessandro Lenci, Saran Lertpradit, Herman Leung, Maria Levina, Cheuk Ying Li, Josie Li, Keying Li, Yuan Li, KyungTae Lim, Bruna Lima Padovani, Krister Lindén, Nikola Ljubešić, Olga Loginova, Andry Luthfi, Mikko Luukko, Olga Lyashevskaya, Teresa Lynn, Vivien Macketanz, Aibek Makazhanov, Michael Mandl, Christopher	1081

1082	Manning, Ruli Manurung, Büşra Marşan, Cătălina	1145
1083	Mărânduc, David Mareček, Katrin Marheinecke,	1146
1084	Héctor Martínez Alonso, André Martins, Jan Mašek,	1147
1085	Hiroshi Matsuda, Yuji Matsumoto, Alessandro	1148
1086	Mazzei, Ryan McDonald, Sarah McGuinness, Gus-	1149
1087	tavo Mendonça, Niko Miekka, Karina Mischenkova,	1150
1088	Margarita Misirpashayeva, Anna Missilä, Cătălin	1151
1089	Mititelu, Maria Mitrofan, Yusuke Miyao, AmirHos-	1152
1090	sein Mojiri Foroushani, Judit Molnár, Amirsaeid	1153
1091	Moloodi, Simonetta Montemagni, Amir More, Laura	1154
1092	Moreno Romero, Giovanni Moretti, Keiko Sophie	1155
1093	Mori, Shinsuke Mori, Tomohiko Morioka, Shigeki	1156
1094	Moro, Bjartur Mortensen, Bohdan Moskalevskyi,	1157
1095	Kadri Muischnek, Robert Munro, Yugo Murawaki,	1158
1096	Kaili Müürisep, Pinkey Nainwani, Mariam Nakhlé,	1159
1097	Juan Ignacio Navarro Horñiacek, Anna Nedoluzhko,	1160
1098	Gunta Nešpore-Bērzkalne, Manuela Nevaci, Lùòng	1161
1099	Nguyên Thị, Huyên Nguyên Thị Minh, Yoshihiro	1162
1100	Nikaido, Vitaly Nikolaev, Rattima Nitisaroj, Alireza	1163
1101	Nourian, Hanna Nurmi, Stina Ojala, Atul Kr. Ojha,	1164
1102	Adédayó Olúòkun, Mai Omura, Emeka Onwueg-	1165
1103	buzia, Petya Osenova, Robert Östling, Lilja Øvre-	1166
1104	lid, Şaziye Betül Özateş, Merve Özçelik, Arzu-	
1105	can Özgür, Balkız Öztürk Başaran, Hyunji Hay-	
1106	ley Park, Niko Partanen, Elena Pascual, Marco	
1107	Passarotti, Agnieszka Patejuk, Guilherme Paulino-	
1108	Passos, Angelika Peljak-Łapińska, Siyao Peng,	
1109	Cenel-Augusto Perez, Natalia Perkova, Guy Per-	
1110	rier, Slav Petrov, Daria Petrova, Jason Phelan, Jussi	
1111	Piitulainen, Tommi A Pirinen, Emily Pitler, Bar-	
1112	bara Plank, Thierry Poibeau, Larisa Ponomareva,	
1113	Martin Popel, Lauma Pretkalniņa, Sophie Prévost,	
1114	Prokopis Prokopidis, Adam Przepiórkowski, Tiina	
1115	Puolakainen, Sampo Pyysalo, Peng Qi, Andriela	
1116	Rääbis, Alexandre Rademaker, Taraka Rama, Lo-	
1117	ganathan Ramasamy, Carlos Ramisch, Fam Rashel,	
1118	Mohammad Sadegh Rasooli, Vinit Ravishankar, Livy	
1119	Real, Petru Rebeja, Siva Reddy, Georg Rehm, Ivan	
1120	Riabov, Michael Rießler, Erika Rimkutė, Larissa Ri-	
1121	naldi, Laura Rituma, Luisa Rocha, Eiríkur Rögnvalds-	
1122	son, Mykhailo Romanenko, Rudolf Rosa, Valentin	
1123	Roşca, Davide Rovati, Olga Rudina, Jack Rueter,	
1124	Kristján Rúnarsson, Shoal Sadde, Pegah Safari,	
1125	Benoît Sagot, Aleksí Sahala, Shadi Saleh, Alessio	
1126	Salomoni, Tanja Samardžić, Stephanie Samson,	
1127	Manuela Sanguinetti, Ezgi Sanyar, Dage Särg, Baiba	
1128	Saulīte, Yanin Sawanakunanon, Shefali Saxena,	
1129	Kevin Scannell, Salvatore Scarlata, Nathan Schnei-	
1130	der, Sebastian Schuster, Lane Schwartz, Djamé Sed-	
1131	dah, Wolfgang Seeker, Mojgan Seraji, Mo Shen,	
1132	Atsuko Shimada, Hiroyuki Shirasu, Yana Shishk-	
1133	ina, Muh Shohibussirri, Dmitry Sichinava, Janine	
1134	Siewert, Einar Freyr Sigurd'sson, Aline Silveira,	
1135	Natalia Silveira, Marija Simi, Radu Simionescu,	
1136	Katalin Simkó, Mária Šimková, Kiril Simov, Maria	
1137	Skachedubova, Aaron Smith, Isabela Soares-Bastos,	
1138	Carolyn Spadine, Rachele Sprugnoli, Steinhórf Ste-	
1139	ingrímsson, Antonio Stella, Milan Straka, Emmett	
1140	Strickland, Jana Strnadová, Alane Suhr, Yogi Les-	
1141	mana Sulestio, Umut Sulubacak, Shingo Suzuki,	
1142	Zsolt Szántó, Dima Taji, Yuta Takahashi, Fabio Tam-	
1143	burini, Mary Ann C. Tan, Takaaki Tanaka, Sam-	
1144	son Tella, Isabelle Tellier, Marinella Testori, Guil-	
	laume Thomas, Liisi Torga, Marsida Toska, Trond	1145
	Trosterud, Anna Trukhina, Reut Tsarfaty, Utku Türk,	1146
	Francis Tyers, Sumire Uematsu, Roman Untilov,	1147
	Zdeňka Uřešová, Larraitz Uriá, Hans Uszkoreit, An-	1148
	drius Utká, Sowmya Vajjala, Rob van der Goot,	1149
	Martine Vanhove, Daniel van Niekerk, Gertjan van	1150
	Noord, Viktor Varga, Eric Villemonte de la Clerg-	1151
	erie, Veronika Vincze, Natalia Vlasova, Aya Wakasa,	1152
	Joel C. Wallenberg, Lars Wallin, Abigail Walsh,	1153
	Jing Xian Wang, Jonathan North Washington, Max-	1154
	imilan Wendt, Paul Widmer, Seyi Williams, Mats	1155
	Wirén, Christian Wittern, Tsegay Woldemariam, Tak-	1156
	sum Wong, Alina Wróblewska, Mary Yako, Kayo	1157
	Yamashita, Naoki Yamazaki, Chunxiao Yan, Koichi	1158
	Yasuoka, Marat M. Yavrumyan, Arife Betül Yenice,	1159
	Olçay Taner Yıldız, Zhuoran Yu, Zdeněk Žabokrt-	1160
	ský, Shorouq Zahra, Amir Zeldes, Hanzhi Zhu, Anna	1161
	Zhuravleva, and Rayan Ziane. 2021b. <a href="#">Universal de-</a>	1162
	<a href="#">pendencies 2.8</a> . LINDAT/CLARIAH-CZ digital li-	1163
	brary at the Institute of Formal and Applied Linguis-	1164
	tics (ÚFAL), Faculty of Mathematics and Physics,	1165
	Charles University.	1166
	Zhisong Zhang, Emma Strubell, and Eduard Hovy. 2021.	1167
	<a href="#">On the benefit of syntactic supervision for cross-</a>	1168
	<a href="#">lingual transfer in semantic role labeling</a> . In <i>Proceed-</i>	1169
	<i>ings of the 2021 Conference on Empirical Methods</i>	1170
	<i>in Natural Language Processing</i> , pages 6229–6246,	1171
	Online and Punta Cana, Dominican Republic. Asso-	1172
	ciation for Computational Linguistics.	1173

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

Table 6: Number of sentences available in each language in UPB v2.

## A Artifacts

### A.1 Corpus Distribution

Table 6 shows the corpus distribution in UPB v2. 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 Corpus Preprocessing

Some treebanks have enhanced dependency tree annotations that cause new tokens to be added to the sentences. These tokens are called enhanced tokens. The enhanced tokens cause some SRL annotations in UPB v2 to be shifted when merged with UD v2.9, resulting in the wrong predicate or semantic role annotations. For example, look at the example of wrong predicate annotation taken from the dev set in Finnish-TDT (UPB v2). Token 10.1 is an enhanced token.

```
# sent_id = w063.9
# text = Osasto.....
1 _ _ _
2 _ _ _
3 be.01 A1:2|AM-LOC:7
A1:1-2|AM-LOC:6-7
4 _ _ _
5 _ _ _
6 _ _ _
7 _ _ _
```

```
8 _ _ _
9 _ _ _
10 _ _ _
10.1 _ _ _
11 be.01 A1:9 A1:9-10
12 _ _ _
13 _ _ _
```

The corresponding annotation in UD v2.9 is as follows. Note that we only present the tokenized words, lemmas, and POS tags here since we only present the UD annotation to highlight the shifted annotation problem.

```
# sent_id = w063.9
# text = Osasto N7 sijaitsee
samassa korttelissa
Naistenklinikan rakennuksessa
ja osasto LV37 Kätilöopiston
sairaalassa.
1 Osasto osasto NOUN
2 N7 N7 SYM
3 sijaitsee sijaita VERB
4 samassa sama PRON
5 korttelissa kortteli NOUN
6 Naistenklinikan nais#klinikka
NOUN
7 rakennuksessa rakennus NOUN
8 ja ja CONJ
9 osasto osasto NOUN
10 LV37 LV37 SYM
10.1 sijaitsee sijaita VERB
11 Kätilöopiston kätilö#opisto
NOUN
12 sairaalassa sairaala NOUN
13 . . PUNCT
```

If we compare the two annotations from UPB v2 and UD v2.9, we can see that the first predicate annotated on token 3, i.e., “sijaitsee”, is correct. However, the second predicate annotated on the token 11, i.e., “Kätilöopiston”, is wrong. The correct second predicate is token 10.1, i.e., “sijaitsee”. The annotation is somehow shifted because of the enhanced token added, token 10.1. Therefore, we fix the annotation in UPB v2 to be as follows.

```
# sent_id = w063.9
# text = Osasto.....
1 _ _ _
2 _ _ _
3 be.01 A1:2|AM-LOC:7
A1:1-2|AM-LOC:6-7
4 _ _ _
5 _ _ _
```

1252 6 \_ \_ \_  
 1253 7 \_ \_ \_  
 1254 8 \_ \_ \_  
 1255 9 \_ \_ \_  
 1256 10 \_ \_ \_  
 1257 10.1 be.01 A1:9 A1:9-10  
 1258 11 \_ \_ \_  
 1259 12 \_ \_ \_  
 1260 13 \_ \_ \_

1261 In some cases, not only the predicate annotations  
 1262 indicated in the third column but also the semantic  
 1263 role annotations indicated in the fourth column  
 1264 are shifted. We run a script to fix the annota-  
 1265 tion problems in all treebanks with enhanced  
 1266 dependency tree annotations. The treebanks  
 1267 with enhanced dependency tree annotations are  
 1268 Czech-CAC, Czech-FicTree, Czech-PDT,  
 1269 Dutch-Alpino, Dutch-LassySmall,  
 1270 Finnish-TDT, Italian-ISDT,  
 1271 Spanish-AnCora, and Ukrainian-IU.

1272 After we fix the shifted predicate and se-  
 1273 mantic role annotations, we notice that some  
 1274 predicates and their semantic roles are anno-  
 1275 tated in the enhanced tokens that do not appear  
 1276 in the original sentence. The treebanks that  
 1277 contain this phenomenon are Dutch-Alpino,  
 1278 Dutch-LassySmall, Finnish-TDT,  
 1279 Ukrainian-IU, and Spanish-AnCora. We  
 1280 cannot accommodate these annotations since we  
 1281 build the models based on the sentence’s original  
 1282 tokens. Therefore, we omit the predicates and  
 1283 their corresponding semantic roles annotated on  
 1284 enhanced tokens in our experiments.

### 1285 A.3 License

1286 Complete UPB v2 contains annotations from UD  
 1287 v2.9. Table 7 shows the license for each treebank in  
 1288 UD v2.9. Despite licenses inherited from UD v2.9,  
 1289 UPB v2 also has a CDLA-Sharing-1.0 license.

1290 We refer to publicly available codes to build the  
 1291 corpus and models for experiments. We provide the  
 1292 list of GitHub repositories with their corresponding  
 1293 licenses, as follows.

- 1294 1. [UniversalPropositions/tools](#): Apache-2.0
- 1295 2. [diegma/neural-dep-srl](#) (Marcheggiani and  
 1296 Titov, 2017): Apache-2.0
- 1297 3. [AnWang-AI/towe-eacl](#) (Jiang et al., 2021): No  
 1298 License
- 1299 4. [dmlc/dgl](#) (Wang et al., 2019a): Apache-2.0

Treebank	License
English-EWT	CC BY-SA 4.0
Chinese-GSD	CC BY-SA 4.0
Czech-CAC	CC BY-SA 4.0
Czech-CLTT	CC BY-SA 4.0
Czech-FicTree	CC BY-NC-SA 4.0
Czech-PDT	CC BY-NC-SA 3.0
Dutch-Alpino	CC BY-SA 4.0
Dutch-LassySmall	CC BY-SA 4.0
Finnish-FTB	CC BY 4.0
Finnish-TDT	CC BY-SA 4.0
French-GSD	CC BY-SA 4.0
French-Rhapsodie	CC BY-SA 4.0
French-Sequoia	LGPL-LR
German-GSD	CC BY-SA 4.0
German-HDT	CC BY-SA 4.0
Greek-GDT	CC BY-NC-SA 3.0
Hindi-HDTB	CC BY-NC-SA 4.0
Hungarian-Szeged	CC BY-NC-SA 3.0
Indonesian-GSD	CC BY-SA 4.0
Italian-ISDT	CC BY-NC-SA 3.0
Italian-ParTUT	CC BY-NC-SA 4.0
Italian-POSTWITA	CC BY-NC-SA 4.0
Italian-TWITTIRO	CC BY-SA 4.0
Italian-VIT	CC BY-NC-SA 3.0
Japanese-GSD	CC BY-SA 4.0
Japanese-GSDLUW	CC BY-SA 4.0
Korean-GSD	CC BY-SA 4.0
Korean-Kaist	CC BY-SA 4.0
Marathi-UFAL	CC BY-SA 4.0
Polish-LFG	GNU GPL 3.0
Polish-PDB	CC BY-NC-SA 4.0
Portuguese-Bosque	CC BY-SA 4.0
Portuguese-GSD	CC BY-SA 4.0
Romanian-Nonstandard	CC BY-SA 4.0
Romanian-RRT	CC BY-SA 4.0
Romanian-SiMoNEro	CC BY-SA 4.0
Russian-GSD	CC BY-SA 4.0
Russian-Taiga	CC BY-SA 4.0
Spanish-AnCora	CC BY 4.0
Spanish-GSD	CC BY-SA 4.0
Tamil-TTB	CC BY-NC-SA 3.0
Telugu-MTG	CC BY-SA 4.0
Ukrainian-IU	CC BY-NC-SA 4.0
Vietnamese-VTB	CC BY-SA 4.0

Table 7: License for each treebank in UD v2.9.

5. [gordicaleksa/pytorch-GAT](#) (Gordić, 2020): 1300  
MIT 1301
6. [deepakn97/relationPrediction](#) (Nathani et al., 1302  
2019): No License 1303
7. [thudm/hgb](#) (Lv et al., 2021): No License 1304
8. [shenwzh3/RGAT-ABSA](#) (Wang et al., 2020): 1305  
MIT 1306
9. [wasiahmad/GATE](#) (Ahmad et al., 2021b): 1307  
MIT 1308

We access all the resources we mentioned above 1309  
 solely for academic research. We make sure that 1310  
 we obey the intended usage of each artifact. 1311

Treebank	Dev			Test		
	F1 <sub>POS</sub>	UAS	LAS	F1 <sub>POS</sub>	UAS	LAS
English-EWT	96.79	92.46	90.86	96.80	91.42	89.82
Chinese-GSD	95.35	85.11	83.19	95.52	87.06	85.13
Czech-CAC	99.26	92.97	91.62	98.70	93.43	91.68
Czech-CLTT	99.44	89.13	86.98	98.98	88.32	86.09
Czech-FicTree	98.43	94.68	93.11	98.34	94.61	92.76
Czech-PDT	98.77	93.74	92.24	98.63	93.50	91.87
Dutch-Alpino	98.36	94.53	92.24	97.33	92.87	90.42
Dutch-LassySmall	97.03	90.77	87.62	96.31	92.12	89.11
Finnish-FTB	96.90	93.77	92.29	96.87	94.03	92.41
Finnish-TDT	98.08	91.97	90.41	97.78	92.24	90.74
French-GSD	98.45	95.66	94.45	98.20	93.47	91.87
French-Rhapsodie	98.12	87.75	83.25	97.64	86.42	81.88
French-Sequoia	99.03	93.54	92.23	99.12	93.10	91.70
German-GSD	96.19	91.78	88.61	95.37	89.65	85.62
German-HDT	98.08	95.18	93.64	98.30	95.30	93.72
Greek-GDT	97.74	91.77	90.43	97.71	92.93	91.19
Hindi-HDTB	97.89	96.62	94.49	97.93	96.68	94.43
Hungarian-Szeged	96.66	87.64	84.10	96.06	86.72	83.25
Indonesian-GSD	94.64	86.49	76.25	94.73	87.31	77.33
Italian-ISDT	98.54	94.41	92.84	98.62	94.37	93.16
Italian-ParTUT	97.86	92.76	90.52	98.54	93.10	91.40
Italian-PostWITA	97.35	87.21	83.20	96.96	88.33	84.41
Italian-TWITTIRO	96.79	87.25	81.64	96.20	84.85	79.77
Italian-VIT	98.12	90.63	88.82	98.16	91.54	89.05
Japanese-GSD	98.34	96.09	95.47	98.10	95.11	94.21
Japanese-GSDLUW	98.54	96.12	95.82	98.58	95.35	95.12
Korean-GSD	95.79	88.22	85.41	96.27	89.65	87.07
Korean-Kaist	96.19	91.35	90.39	95.58	90.41	89.45
Marathi-UFAL	89.32	74.55	64.32	90.53	79.85	70.63
Polish-LFG	98.94	97.56	96.73	99.05	97.80	96.92
Polish-PDB	98.75	94.17	92.69	98.74	94.58	93.16
Portuguese-Bosque	97.92	94.25	92.51	98.10	94.85	93.54
Portuguese-GSD	98.36	94.44	93.34	98.28	94.21	93.23
Romanian-Nonstandard	96.77	93.18	90.04	96.40	91.43	87.75
Romanian-RRT	98.06	91.96	88.60	97.92	91.93	88.45
Romanian-SiMoNERo	98.19	93.38	91.21	98.23	93.78	91.86
Russian-GSD	98.38	90.55	87.80	98.09	90.44	87.21
Russian-Taiga	95.80	83.94	79.32	97.06	84.42	81.41
Spanish-AnCora	98.99	93.83	92.16	98.96	93.82	92.00
Spanish-GSD	97.13	91.91	89.79	97.26	91.93	89.58
Tamil-TTB	87.17	81.24	73.48	86.93	80.89	72.30
Telugu-MTG	94.41	92.90	86.25	94.45	93.07	85.58
Ukrainian-IU	98.08	91.14	89.34	97.67	90.10	88.24
Vietnamese-VTB	92.84	78.92	74.99	92.81	77.58	74.16

Table 8: F1<sub>POS</sub>, UAS, and LAS of each treebank’s POS tagger and dependency parser in UPB v2. F1<sub>POS</sub> indicates the F1 score of the POS tagger.

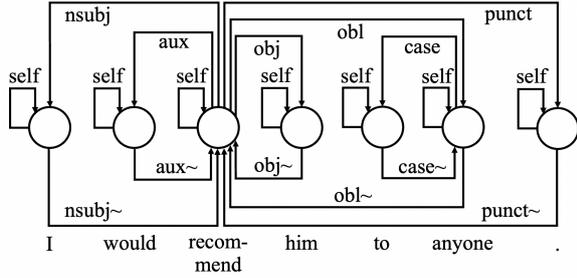


Figure 4: Dependency graph of a sentence converted from its dependency tree.

## B Dependency Tree to Dependency Graph

Marcheggiani and Titov (2017) convert a dependency tree to a dependency graph in the following steps:

1. Add the edges from the dependency tree to the graph. These edges flow from the heads to the dependents. Label each edge with the corresponding dependency relation from the dependency tree.
2. Add new edges that flow in the opposite directions of the original dependency directions. These new edges flow from the dependents to the heads.
3. Assign a unique dependency relation, derived from the original dependency relation, to each new edge added in step 2. For example, if the original dependency relation is “nsubj”, the edge that flows in the opposite direction is labeled as “nsubj~”.
4. Add new edges that flow from each node to itself (self-connection) and label them as “self”.

Figure 4 displays the dependency graph derived from the dependency tree at the bottom of Figure 1.

## C Two-Attention Relational Graph Attention Networks

Similar to RGATs (Wang et al., 2020), two-attention relational GATs (TAGATs) also apply two types of attention weights to measure the influence of neighbor nodes when updating the corresponding node representation. However, instead of using position-wise FFN to calculate the second attention weight as in RGATs (Wang et al., 2020), two-attention relational GATs (TAGATs)

apply dot-product equation proposed by Veličković et al. (2018) to calculate both attention weights. We explain the modification that we made to RGATs in this section.

TAGATs calculate the first attention weight,  $\alpha$ , using Equation 3 in the original GAT paper (Veličković et al., 2018). To calculate the second attention weight,  $\beta$ , TAGATs slightly modify the equation to incorporate the DR,  $r_{ij}$ , as shown in Equation 5, where  $k$  is the current attention head,  $l$  is the current layer,  $\mathcal{N}_i$  is the neighbor nodes of node  $i$ ,  $W_r$  is a weight matrix to linearly transform the DR,  $r_{ij}$ , and LR is a Leaky ReLU.

$$\beta_{ij}^{l,k} = \text{softmax}_{j \in \mathcal{N}_i} (\text{LR}(a^{l,kT} [W_r^{l,k} r_{ij}^l])) \quad (5)$$

Furthermore, TAGATs obtain node representation from attention weight  $\alpha$  using Equation 5 and Equation 6 in the original GAT paper (Veličković et al., 2018). To obtain node representation from attention weight  $\beta$ , TAGATs employ Equation 6 and Equation 7. TAGATs concatenate node representations from  $K$  heads in the intermediate layers, as shown in Equation 6. Meanwhile, in the final layer, TAGATs take the average of node representations from  $K$  heads, as shown in Equation 7, where  $L$  is the number of layers.

$$h_{i,\beta}^{l+1} = \sigma(\|_{k=1}^K \sum_{j \in \mathcal{N}_i} \beta_{ij}^{l,k} W^{l,k} h_j^l), \quad l < L \quad (6)$$

$$h_{i,\beta}^{l+1} = \sigma(\frac{1}{K} \sum_{k=1}^K \sum_{j \in \mathcal{N}_i} \beta_{ij}^{l,k} W^{l,k} h_j^l), \quad l = L \quad (7)$$

Finally, TAGATs calculate the node representation for layer  $(l + 1)$ ,  $h_i^{l+1}$ , by applying a linear transformation to the concatenation of node representation  $h_{i,\alpha}^{l+1}$  and node representation  $h_{i,\beta}^{l+1}$  and optionally apply an activation function,  $\sigma$ , on top of the linear transformation, as shown in Equation 8.

$$x_i^{l+1} = [h_{i,\alpha}^{l+1}, h_{i,\beta}^{l+1}] \quad (8)$$

$$h_i^{l+1} = \sigma(W_{l+1} x_i^{l+1} + b_{l+1})$$

## D Experiments

### D.1 Dependency Parsers and POS Taggers

Table 8 shows the POS tagger and dependency parser evaluation results on each treebank in UPB v2. For Japanese-GSDLUW,

Hyperparameter	Value
num_epochs	100
batch_size	32
optimizer	SGD
learning_rate	0.1
num_early_stop	20
num_decay_epoch	5
lr_decay	0.9
min_lr	0.00001
pos_dim	30
pred_ind_dim	30
emb_dropout	0.5
hid_dim	512
num_heads	8
d_k	64
d_v	64
d_ff	2048

Table 9: Basic hyperparameters applied in the experiments.

French-Rhapsodie, and English-EWT treebanks, we train the POS taggers and dependency parsers from scratch using Stanza (Qi et al., 2020) with a 0.0005 learning rate, 70,000 max steps, and 10,000 max steps before stopping. We measure the performance of POS taggers with the F1 score. Meanwhile, we measure the performance of dependency parsers with the unlabeled attachment score (UAS) and the labeled attachment score (LAS).

## D.2 Hyperparameter Search

We take the representative models from transformer-based models (i.e., SAN-RPRs), GCN-based models (i.e., ARGCNs and SGCNs), GAT-based models (i.e., GATs), and BiLSTM-based model (i.e., BiLSTMs), to experiment with the optimizers. We experiment with SGD (Kiefer and Wolfowitz, 1952), Adam (Kingma and Ba, 2015), and AdamW (Loshchilov and Hutter, 2019) as the optimizer. We also try different learning rates for each optimizer, i.e., 0.1, 0.01, and 0.001. The SGD optimizer with a 0.1 learning rate works the best in all the representative models. Therefore, we apply this setting to the rest of our experiments. Table 9 summarizes the fixed hyperparameters we use throughout our experiments. Algorithm 1 shows the logic for model training. Below, we explain each hyperparameter:

1. `num_epochs`: Number of epochs for model training.
2. `batch_size`: Batch size for model training.

3. `optimizer`: Type of optimizer for model training.
4. `learning_rate`: Initial learning rate for model training.
5. `num_early_stop`: Stop the training if there is no improvement after a certain number of consecutive epochs.
6. `num_decay_epoch`: The upper limit of epoch before we start decaying the learning rate.
7. `lr_decay`: The ratio to decay the learning rate.
8. `min_lr`: Minimum learning rate allowed. We stop the model training if the learning rate falls below this threshold.
9. `pos_dim`: Dimension of POSE,  $o$  (Equation 1). If the dimension is 0, we will not concatenate  $o$  in the input layer.
10. `pred_ind_dim`: Dimension of PIE,  $p$  (Equation 1). If the dimension is 0, we will not concatenate  $p$  in the input layer.
11. `emb_dropout`: Dropout applied in the input layer (Equation 1).
12. `hid_dim`: The dimension of the node representation that the encoder accepts.
13. `num_heads`: Number of heads applied in the multi-head self-attention mechanism present in transformer-based models, GAT-based models, and ARGCNs.
14. `d_k`: The dimension of keys applied in transformer-based models.
15. `d_v`: The dimension of values applied in transformer-based models.
16. `d_ff`: The output dimension of the first linear transformation in transformer-based model’s position-wise FFN.

The following sections will explain the hyperparameter search and hyperparameter values that work best in each model. Table 11, Table 12, Table 13, and Table 10 describe the hyperparameter search for transformer-based models, GCN-based models, GAT-based models, and BiLSTM-based

---

**Algorithm 1** Pseudocode of the model training.

---

```
Require: num_early_stop, num_decay_epoch, min_lr, lr_decay, num_epochs, learning_rate
best_f1 ← 0
no_improvement ← 0
for curr_epoch ← 1, num_epochs do
  curr_f1 ← train(learning_rate)
  if curr_f1 > best_f1 then
    best_f1 ← curr_f1
    no_improvement ← 0
  else
    no_improvement ← no_improvement + 1
    if no_improvement ≥ num_early_stop then
      break
    end if
    if curr_epoch > num_decay_epoch then
      learning_rate ← lr_decay * learning_rate
      if learning_rate < min_lr then
        break
      end if
    end if
  end if
end for
```

---

1462 models, respectively. Due to the number of hyper-  
1463 parameters, we divide the hyperparameter search  
1464 into groups indicated by the leftmost column, i.e.,  
1465 the column with a "No" header. We will search  
1466 for the best combination between the hyperparamete-  
1467 rs in the same group. For example, in Table 11,  
1468 num\_enc\_layers and enc\_dropout belong  
1469 to group 1, which means we experiment with dif-  
1470 ferent dropouts, i.e., 0.1, 0.2, 0.3, 0.4, and 0.5, for  
1471 each number of layers, i.e., 1, 2, 3, and 4. Below,  
1472 we explain each hyperparameter involved in the  
1473 hyperparameter search.

- 1474 1. num\_enc\_layers: Number of layers  
1475 stacked together.
- 1476 2. enc\_dropout: Dropout applied in the mod-  
1477 els, including dropout applied in the encoder  
1478 (Section 3.1.2).
- 1479 3. lstm\_num\_layers: Number of BiLSTM  
1480 layers stacked together.
- 1481 4. lstm\_dropout\_net: Dropout applied in  
1482 BiLSTMs, including dropout applied in the  
1483 encoder (Section 3.1.2).
- 1484 5. gnn\_activation: Activation function  
1485 used in GCN-based and GAT-based models,  
1486 applied in the encoder (Section 3.1.2).
- 1487 6. gnn\_activation\_at\_final\_layer:  
1488 Whether to apply the activation function in  
1489 the last layer, especially if we stack more than  
1490 one layer.

7. lstm\_activation: Activation function  
used in BiLSTMs (Section 3.1.2). 1491 1492
8. lstm\_hidden\_size: Hidden size of BiL-  
STMs. Since the network is bidirectional, the  
dimension of the final hidden representation  
is  $2 \times \text{lstm\_hidden\_size}$ . 1493 1494 1495 1496
9. deprel\_dim: Dimension of DRE,  $d$  (Equa-  
tion 1). If the dimension is 0, we will not  
concatenate  $d$  in the input layer. 1497 1498 1499
10. abs\_position\_dim: Dimension of APE,  
 $a$ , and SAPE,  $s$  (Equation 1). 1500 1501
11. use\_dep\_abs\_position: Boolean  
value which indicates whether to concatenate  
SAPE,  $s$ , in the input layer (Equation 1). 1502 1503 1504
12. use\_word\_abs\_position: Boolean  
value which indicates whether to concatenate  
APE,  $a$ , in the input layer (Equation 1). 1505 1506 1507
13. att\_dim: Dimension of a trainable vector  
 $a$  in ARGCNs (Equation 7 in Jiang et al.  
(2021)). 1508 1509 1510
14. base\_size: Base size,  $B$ , in RGCNs  
(Equation 3 in Schlichtkrull et al. (2018)). 1511 1512
15. rel\_pos\_dim: Dimension of RPR in GCN-  
based and GAT-based models. We do not use  
this parameter for transformer-based models  
because the RPR in transformer-based models  
must have the same dimension as  $d_k$  and  $d_v$   
(Vaswani et al., 2017). 1513 1514 1515 1516 1517 1518

- 1519 16. num\_embed\_graph\_heads: Number of  
 1520 heads where we modify  $M$  matrix according  
 1521 to distance matrix,  $D$ , in GATEs (Equation 3  
 1522 in Ahmad et al. (2021b)). We apply a zero  
 1523 matrix for  $M$  for the other heads to connect  
 1524 all the nodes in the graph.
- 1525 17. max\_tree\_dists: The  $\delta$  parameter  
 1526 applied to each head in GATEs (Equation  
 1527 3 in Ahmad et al. (2021b)). The  
 1528 length of this parameter must equal  
 1529 num\_embed\_graph\_heads.
- 1530 18. max\_relative\_positions: The maxi-  
 1531 mum absolute value for relative position ( $k$   
 1532 in Shaw et al. (2018)) or structural relative  
 1533 position ( $r$  in Section Wang et al. (2019b)).
- 1534 19. use\_dep\_rel\_pos: The boolean value in-  
 1535 dicates whether to incorporate SRPR in the  
 1536 edge representation.
- 1537 20. use\_word\_rel\_pos: The boolean value  
 1538 indicates whether to incorporate RPR in the  
 1539 edge representation.
- 1540 21. deprel\_edge\_dim: Dimension of repre-  
 1541 sentation for dependency relation,  $d_r$ , when  
 1542 we use B-DR.
- 1543 22. deparc\_edge\_dim: Dimension of repre-  
 1544 sentation for direction,  $d_d$ , when we use B-  
 1545 DR.
- 1546 23. deprel\_ext\_edge\_dim: Dimension of  
 1547 dependency relation representation when we  
 1548 use A-DR.

1549 There are two ways of generating APE, SAPE,  
 1550 RPR, and SRPR, i.e., using learned positional em-  
 1551 bedding (Gehring et al., 2017) and using sine and  
 1552 cosine functions (Vaswani et al., 2017). We con-  
 1553 duct preliminary experiments and find that sine and  
 1554 cosine functions work better than learned positional  
 1555 embedding. Therefore, we generate the represen-  
 1556 tation for each type of position in the experiments  
 1557 using sine and cosine functions.

### 1558 D.2.1 Transformer-Based Models

1559 Table 11 shows the hyperparameter search in  
 1560 transformer-based models.

### 1561 D.2.2 GCN-Based Models

1562 Table 12 shows the hyperparameter search in GCN-  
 1563 based models.

No	Hyperparameter	Value
1	lstm_num_layers lstm_dropout_net	1, 2, 3, <b>4</b> 0.1, 0.2, <b>0.3</b> , 0.4, 0.5
2	lstm_activation	<b>ReLU</b> , Leaky ReLU, ELU
3	lstm_hidden_size	<b>256</b> , 512
4	deprel_dim abs_position_dim use_dep_abs_position	0, <b>30</b> 0, <b>30</b> T

Table 10: Hyperparameter search in BiLSTM-based models. The search is divided into groups shown in the "No" header. We search for the best combination of hyperparameters in the same group. The bold values indicate the results of the hyperparameter search.

### D.2.3 GAT-Based Models

Table 13 shows the hyperparameter search in GAT-based models.

### D.2.4 BiLSTM-Based Models

Table 10 shows the hyperparameter search in BiLSTM-based models.

## D.3 Computational Resource

We use Tesla P100 to train the models. Training time for GCN-based and GAT-based models takes around 5 hours. Meanwhile, training time for BiLSTM-based and transformer-based models takes around 10 hours. Hyperparameter search in GCN-based models costs around 775 GPU hours. Hyperparameter search in GAT-based models costs around 910 GPU hours. Hyperparameter search in transformer-based models costs around 1,720 GPU hours. Hyperparameter search in BiLSTM-based models costs around 290 GPU hours. After searching for the best hyperparameter setting for each model, we run the training five times for each model, spending around 1,075 GPU hours. Therefore, in total, we spend around 4,770 hours.

## E Supporting Results

### E.1 Transformer-Based Models

#### E.1.1 Comparison

Table 14 shows the detailed comparison of transformer-based models in each language. We calculate the superiority score (SC) of each model based on the model performance in target languages. We allocate 2 points if the model achieves the highest F1 score or 1 point if the model achieves the second-highest F1 score for a specific language.

No	Hyperparameter	Value
<b>General</b>		
1	num_enc_layers enc_dropout	1, 2, <b>3</b> , 4 0.1, <b>0.2</b> , 0.3, 0.4, 0.5
<b>Transformers</b>		
1	deprel_dim	0, <b>30</b>
<b>GATEs</b>		
1	num_embed_graph_heads max_tree_dists	4 <1, 1, 2, 2>, <2, 2, 4, 4>, < <b>4, 4, 8, 8</b> >, <1, 2, 4, 8>
2	deprel_dim abs_position_dim <use_dep_abs_position, use_word_abs_position>	0, <b>30</b> 0, <b>30</b> < <b>T, F</b> >, <F, T>
<b>SAN-RPRs</b>		
1	max_relative_positions	1, 2, 4, 8, <b>16</b>
2	deprel_dim	0, <b>30</b>
<b>SAN-SAPRs</b>		
1	deprel_dim	<b>0</b> , 30
<b>SAN-SRPRs</b>		
1	max_relative_positions	1, 2, 4, 8, 16
2	deprel_dim	<b>0</b> , 30
<b>Trans-SRPRs</b>		
1	deprel_dim	0, <b>30</b>
<b>Trans-SRPR-DRs</b>		
1	deprel_dim abs_position_dim use_word_abs_position use_dep_rel_pos <deprel_edge_dim, deparc_edge_dim, deprel_ext_edge_dim>	0, <b>30</b> 0, <b>30</b> T T, F <32, 32, 0>, <48, 16, 0>, <56, 8, 0>, < <b>60, 4, 0</b> >, <62, 2, 0>, <63, 1, 0> <0, 0, 64>
<b>SAPR-RPRs</b>		
1	deprel_dim	<b>0</b> , 30
<b>SAPR-RPR-DRs</b>		
1	deprel_dim abs_position_dim use_dep_abs_position use_word_rel_pos <deprel_edge_dim, deparc_edge_dim, deprel_ext_edge_dim>	<b>0</b> , 30 0, <b>30</b> T T, F <32, 32, 0>, <48, 16, 0>, <56, 8, 0>, < <b>60, 4, 0</b> >, <62, 2, 0>, <63, 1, 0> <0, 0, 64>

Table 11: Hyperparameter search in transformer-based models. The search is divided into groups shown in the "No" header. We search for the best combination of hyperparameters in the same group. The bold values indicate the results of the hyperparameter search.

No	Hyperparameter	Value
<b>SGCNs</b>		
1	num_enc_layers enc_dropout	1, 2, <b>3</b> , 4 0.1, <b>0.2</b> , 0.3, 0.4, 0.5
2	gnn_activation gnn_activation_at_final_layer	<b>ReLU</b> , Leaky ReLU, ELU <b>T</b> , <b>F</b>
3	deprel_dim abs_position_dim <use_dep_abs_position, use_word_abs_position>	<b>0</b> , <b>30</b> <b>0</b> , <b>30</b> < <b>T</b> , <b>F</b> >, < <b>F</b> , <b>T</b> >
<b>RGCNs</b>		
1	num_enc_layers enc_dropout	<b>1</b> , 2, 3, 4 0.1, 0.2, 0.3, 0.4, <b>0.5</b>
2	gnn_activation gnn_activation_at_final_layer	<b>ReLU</b> , Leaky ReLU, ELU <b>T</b> , <b>F</b>
3	base_size	<b>1</b> , <b>2</b> , 4, 8, 16, 32, 80
4	deprel_dim abs_position_dim <use_dep_abs_position, use_word_abs_position>	<b>0</b> , 30 <b>0</b> , <b>30</b> < <b>T</b> , <b>F</b> >, < <b>F</b> , <b>T</b> >
<b>ARGCNs</b>		
1	num_enc_layers enc_dropout	<b>1</b> , 2, 3, 4 0.1, 0.2, 0.3, 0.4, <b>0.5</b>
2	gnn_activation gnn_activation_at_final_layer	<b>ReLU</b> , Leaky ReLU, ELU <b>T</b> , <b>F</b>
3	deprel_ext_edge_dim att_dim rel_pos_dim use_word_rel_pos	1, 2, 4, <b>8</b> , 16 1, 2, 4, 8, <b>16</b> <b>64</b> , 128 <b>T</b>
4	deprel_dim abs_position_dim use_dep_abs_position	<b>0</b> , <b>30</b> <b>0</b> , 30 <b>T</b>

Table 12: Hyperparameter search in GCN-based models. The search is divided into groups shown in the "No" header. We search for the best combination of hyperparameters in the same group. The bold values indicate the results of the hyperparameter search.

### 1596 E.1.2 Ablation Study of SAN-RPRs

1597 We experiment with removing either DRE,  $d$ ,  
1598 POSE,  $o$ , or PIE,  $p$ , from the node representation  
1599 in SAN-RPRs. Table 15 shows the results of the  
1600 ablation study. Removing either DRE or PIE from  
1601 the node representation reduces the performance of  
1602 SAN-RPRs. However, SAN-RPRs without POSE  
1603 perform better in most languages.

## 1604 E.2 Best Models

### 1605 E.2.1 Comparison with SAN-RPRs w/o POSE

1606 In Table 4, according to the superiority score, we  
1607 can see that SAPR-RPRs perform best in more lan-  
1608 guages than SAN-RPRs, even though the average  
1609 F1 score of SAN-RPRs is better than SAPR-RPRs.  
1610 However, as discussed before, SAN-RPRs without  
1611 POSE perform better in most languages than SAN-  
1612 RPRs with POSE, which we use for comparison  
1613 in Table 4. Therefore, in Table 16, we re-compare  
1614 SAN-RPRs without POSE with the other best mod-  
1615 els. After removing POSE from SAN-RPRs, we  
1616 find that the model performs best in more languages

than SAPR-RPRs with 18 and 16 superiority scores,  
respectively.

### E.2.2 Fine-Tuned Models

We fine-tune the contextualized word embedding  
in SAN-RPRs and SAPR-RPRs, i.e., multilingual  
BERT (mBERT). Table 17 compares the average  
F1 scores of models with frozen and fine-tuned  
mBERT. Overall, fine-tuning increases the perfor-  
mance of both models. However, in the fine-tuned  
models, the variability of the average F1 score in  
each run increases, indicated by the higher stan-  
dard deviation in fine-tuned mBERT. This is ex-  
pected as when we fine-tune the mBERT, many  
parameters from mBERT are involved in the train-  
ing process, increasing the randomness variable  
in model training. The behavior of the models is  
similar before and after the fine-tuning. The fine-  
tuned SAN-RPRs perform better than SAPR-RPRs  
with 54.01% and 53.82% average F1 scores, re-  
spectively.

No	Hyperparameter	Value
<b>General</b>		
1	gnn_activation gnn_activation_at_final_layer	ReLU, <b>Leaky ReLU</b> , ELU T, F
<b>GATs</b>		
1	num_enc_layers enc_dropout	1, <b>2</b> , 3, 4 <b>0.1</b> , 0.2, 0.3, 0.4, 0.5
2	deprel_dim abs_position_dim <use_dep_abs_position, use_word_abs_position>	0, <b>30</b> 0, <b>30</b> <T, F>, <F, T>
<b>SHGNs</b>		
1	num_enc_layers enc_dropout	1, <b>2</b> , 3, 4 0.1, 0.2, <b>0.3</b> , 0.4, 0.5
2	deprel_ext_edge_dim rel_pos_dim use_word_rel_pos	<b>16</b> , 32, 64, 128 <b>16</b> , 32, 64, 128 T
3	deprel_dim abs_position_dim use_word_rel_pos use_dep_abs_position	<b>0</b> , 30 <b>0</b> , 30 T, F T
4	<deprel_ext_edge_dim, deparc_ext_edge_dim, deprel_ext_edge_dim>	<8, 8, 0>, <12, 4, 0>, <14, 2, 0>, <15, 1, 0>, < <b>0, 0, 16</b> >
<b>TAGATs</b>		
1	num_enc_layers enc_dropout	1, <b>2</b> , 3, 4 0.1, 0.2, <b>0.3</b> , 0.4, 0.5
2	deprel_ext_edge_dim rel_pos_dim use_word_rel_pos	<b>16</b> , 32, 64, 128 16, 32, 64, <b>128</b> T
3	deprel_dim abs_position_dim use_word_rel_pos use_dep_abs_position	<b>0</b> , 30 <b>0</b> , 30 T, F T
4	<deprel_ext_edge_dim, deparc_ext_edge_dim, deprel_ext_edge_dim>	<8, 8, 0>, <12, 4, 0>, <14, 2, 0>, <15, 1, 0>, < <b>0, 0, 16</b> >
<b>KBGATs</b>		
1	num_enc_layers enc_dropout	1, <b>2</b> , 3, 4 0.1, <b>0.2</b> , 0.3, 0.4, 0.5
2	deprel_ext_edge_dim rel_pos_dim use_word_rel_pos	16, <b>32</b> , 64, 128 16, <b>32</b> , 64, 128 T
3	deprel_dim abs_position_dim use_word_rel_pos use_dep_abs_position	<b>0</b> , 30 <b>0</b> , 30 T, F T
4	<deprel_ext_edge_dim, deparc_ext_edge_dim, deprel_ext_edge_dim>	<16, 16, 0>, < <b>24, 8, 0</b> >, <28, 4, 0>, <30, 2, 0>, <31, 1, 0>, <0, 0, 32>

Table 13: Hyperparameter search in GAT-based models. The search is divided into groups shown in the "No" header. We search for the best combination of hyperparameters in the same group. The bold values indicate the results of the hyperparameter search.

	GATEs	Trans	SAN-RPRs	SAN-SAPRs	SAN-SRPRs	Trans-SRPRs	SAPR-RPRs	Trans-SRPR-DRs	SAPR-RPR-DRs
EN	78.96±0.31	76.16±0.51	78.26±0.40	75.93±0.47	78.27±0.50	79.03±0.32	78.11±0.42	79.85±0.21	79.83±0.19
AVG	52.57±0.23	52.07±0.21	<b>52.73±0.40</b>	51.98±0.10	51.51±0.27	52.21±0.30	52.64±0.38	50.60±0.14	50.69±0.21
TA	36.72±0.89	35.32±1.56	37.96±1.75	38.07±1.47	35.93±2.16	35.74±1.47	<b>39.57±1.18</b>	32.23±1.14	32.91±1.22
HI	<b>48.56±0.37</b>	47.40±0.18	47.51±0.62	45.07±0.52	45.70±0.63	47.80±0.57	45.04±0.35	48.04±0.45	47.57±0.65
ZH	48.09±0.64	48.86±0.31	50.37±1.17	50.06±0.15	46.96±0.54	46.51±0.53	<b>50.96±0.88</b>	43.63±1.09	43.12±0.87
JA	37.45±0.86	<b>38.38±0.37</b>	37.69±0.99	37.91±0.88	36.29±1.60	36.46±1.20	34.78±1.29	33.86±0.61	33.56±1.47
VI	28.01±0.91	28.25±0.57	28.69±0.79	<b>29.70±0.50</b>	27.70±0.61	27.34±0.59	29.10±0.45	25.72±0.22	25.37±0.23
KO	43.44±0.90	43.13±2.26	42.61±1.84	<b>46.09±1.02</b>	43.38±0.84	43.62±0.57	45.24±1.23	39.34±1.19	40.76±0.51
ID	57.88±0.47	55.36±0.42	58.78±1.09	57.95±0.56	59.79±0.27	58.41±0.35	<b>59.97±0.53</b>	53.41±0.69	52.89±0.25
HU	<b>50.69±0.25</b>	50.13±0.85	49.76±0.35	49.68±0.56	49.28±0.20	50.60±0.49	49.08±0.34	49.51±0.47	49.31±0.72
RO	53.86±0.57	52.90±0.39	54.23±0.67	52.11±0.23	52.44±0.33	53.75±0.55	<b>54.46±0.52</b>	52.50±0.31	52.50±0.49
FR	61.67±0.30	60.23±0.39	<b>62.19±0.41</b>	60.52±0.27	60.34±0.37	61.30±0.50	62.11±0.47	60.82±0.20	60.91±0.28
MR	37.49±1.74	<b>42.85±2.67</b>	41.06±2.89	41.22±1.65	37.37±1.27	36.58±1.68	40.36±2.20	36.75±1.67	37.40±1.56
UK	58.89±0.53	58.83±0.20	58.92±0.26	58.87±0.48	58.23±0.69	58.93±0.72	<b>59.36±0.72</b>	57.94±0.42	58.23±0.39
PT	65.49±0.24	64.05±0.26	66.05±0.21	64.76±0.48	65.16±0.41	65.34±0.15	<b>66.49±0.33</b>	64.61±0.24	64.84±0.17
IT	57.52±0.38	57.42±0.43	<b>58.11±0.33</b>	56.35±0.45	56.09±0.33	57.42±0.32	57.80±0.39	56.60±0.31	56.72±0.42
ES	63.36±0.24	62.04±0.42	<b>63.71±0.33</b>	61.53±0.31	62.70±0.20	63.47±0.27	63.62±0.25	61.78±0.14	61.68±0.22
CS	<b>57.69±0.39</b>	56.19±0.48	56.87±0.27	55.27±0.57	56.58±0.27	57.61±0.34	55.80±0.51	56.44±0.26	56.35±0.49
EL	60.37±0.59	59.03±0.34	<b>60.59±0.23</b>	57.30±0.69	58.98±0.82	59.31±0.30	60.23±0.40	57.89±0.33	58.33±0.33
FI	55.00±0.34	54.72±0.33	<b>55.58±0.42</b>	54.74±0.47	53.91±0.37	54.86±0.31	55.29±0.54	54.11±0.17	54.15±0.25
RU	59.36±0.24	59.31±0.47	60.18±0.44	<u>60.33±0.37</u>	59.09±0.30	59.29±0.40	<b>61.13±0.50</b>	58.60±0.39	58.56±0.25
NL	<b>63.73±0.53</b>	62.48±0.76	62.84±0.37	61.75±0.26	61.73±0.41	63.07±0.49	62.22±0.58	62.23±0.20	62.43±0.39
TE	<b>46.32±0.98</b>	44.94±1.38	44.66±2.00	41.64±1.81	41.88±1.17	46.10±1.10	43.88±1.57	42.48±1.28	43.09±1.26
DE	<b>58.88±0.40</b>	58.29±0.70	56.86±0.32	58.39±0.34	58.14±0.30	<u>58.84±0.30</u>	56.98±1.13	58.21±0.26	58.30±0.34
PL	<b>58.73±0.64</b>	57.58±0.25	57.67±0.36	56.25±0.35	57.14±0.40	58.60±0.29	57.28±0.46	57.18±0.62	56.85±0.64
SC	15	4	13	8	1	7	20	1	0
PR	12.2M	12.2M	12.2M	12.2M	12.1M	12.2M	12.2M	12.2M	12.2M

Table 14: F1 scores (%) of transformer-based models evaluated on UPB v2 test set with predicted parsers. The bold score and underlined score indicate the highest and second-highest scores. AVG indicates the average F1 scores of a specific model evaluated in target languages. PR and SC are the number of parameters and the superiority score of each model.

	base	w/o DRE	w/o POSE	w/o PIE
EN	78.26±0.40	77.56±0.41	78.32±0.32	77.66±0.38
AVG	<b>52.73±0.40</b>	51.65±0.28	<b>52.73±0.32</b>	51.54±0.36
TA	<b>37.96±1.75</b>	<u>37.59±1.95</u>	36.53±1.32	32.03±1.68
HI	47.51±0.62	42.22±0.85	<b>47.69±1.10</b>	43.89±0.81
ZH	50.37±1.17	<u>50.42±0.38</u>	<b>50.81±0.51</b>	49.83±0.54
JA	<b>37.69±0.99</b>	33.25±2.01	<u>36.07±1.20</u>	30.67±2.90
VI	28.69±0.79	29.18±1.05	<u>28.30±0.92</u>	<b>29.34±0.96</b>
KO	42.61±1.84	<b>42.73±1.52</b>	42.06±1.50	40.13±0.65
ID	<b>58.78±1.09</b>	58.41±0.89	57.77±0.83	57.06±0.78
HU	49.76±0.35	48.25±0.62	<b>50.02±0.42</b>	49.27±0.86
RO	<u>54.23±0.67</u>	54.11±0.37	<b>54.94±0.50</b>	54.19±0.76
FR	<u>62.19±0.41</u>	61.74±0.32	<b>62.36±0.59</b>	61.85±0.60
MR	41.06±2.89	39.45±1.76	<b>41.83±3.50</b>	41.15±2.62
UK	58.92±0.26	58.75±0.94	58.84±0.59	<b>59.64±0.15</b>
PT	66.05±0.21	<b>66.67±0.53</b>	<u>66.65±0.44</u>	65.46±0.44
IT	58.11±0.33	57.43±0.24	<b>58.65±0.49</b>	<u>58.34±0.32</u>
ES	<u>63.71±0.33</u>	63.18±0.33	<b>64.07±0.51</b>	63.63±0.34
CS	<u>56.87±0.27</u>	54.94±0.40	<b>57.20±0.14</b>	56.40±0.28
EL	<u>60.59±0.23</u>	59.40±0.23	<b>60.86±0.56</b>	59.87±1.11
FI	<b>55.58±0.42</b>	54.78±0.25	<u>55.40±0.16</u>	55.18±0.24
RU	<b>60.18±0.44</b>	60.05±0.67	<u>60.12±0.36</u>	60.00±0.28
NL	62.84±0.37	60.97±0.84	<b>63.21±0.38</b>	60.76±0.68
TE	<b>44.66±2.00</b>	42.41±2.63	43.46±2.02	<u>44.40±3.32</u>
DE	<u>56.86±0.32</u>	55.34±1.44	<b>57.70±0.43</b>	54.36±0.52
PL	57.67±0.36	56.68±0.63	<b>58.16±0.35</b>	<u>58.01±0.31</u>

Table 15: F1 scores (%) of SAN-RPRs with certain embedding removed from the node representation evaluated on UPB v2 test set with predicted parsers. The bold score and underlined score indicate the highest and second-highest scores. AVG indicates the average F1 scores of a specific model evaluated in target languages.

	SAN-RPRs (w/o POSE)	SAPR-RPRs	SGCNs	GATs	TAGATs	BiLSTMs
EN	78.32±0.32	78.11±0.42	79.94±0.27	79.81±0.19	79.07±0.19	76.86±0.34
AVG	52.73±0.32	52.64±0.38	52.52±0.38	52.66±0.14	<b>52.78±0.14</b>	51.85±0.09
TA	<u>36.53±1.32</u>	<b>39.57±1.18</b>	34.32±1.12	35.08±0.58	35.68±1.28	34.19±1.22
HI	47.69±1.10	45.04±0.35	48.24±0.61	<b>48.25±0.33</b>	47.65±0.33	46.63±0.38
ZH	50.81±0.51	<b>50.96±0.88</b>	45.77±0.67	46.12±0.39	46.80±0.64	47.56±0.89
JA	36.07±1.20	34.78±1.29	37.43±0.28	37.99±0.52	<b>39.30±0.73</b>	37.40±0.61
VI	28.30±0.92	<b>29.10±0.45</b>	27.95±0.56	28.06±0.59	28.31±0.55	28.18±0.88
KO	42.06±1.50	<b>45.24±1.23</b>	42.92±0.64	43.22±0.56	<u>44.57±0.24</u>	41.77±1.61
ID	57.77±0.83	<b>59.97±0.53</b>	58.54±0.82	58.33±0.69	<u>59.11±0.87</u>	56.11±0.84
HU	50.02±0.42	49.08±0.34	50.76±0.41	<b>51.10±0.51</b>	<u>50.90±0.37</u>	50.64±0.39
RO	<b>54.94±0.50</b>	<u>54.46±0.52</u>	53.57±0.47	54.12±0.49	53.60±0.45	53.26±0.34
FR	<b>62.36±0.59</b>	<u>62.11±0.47</u>	60.93±0.38	61.64±0.44	61.13±0.22	61.22±0.27
MR	<b>41.83±3.50</b>	40.36±2.20	40.97±3.40	38.06±0.13	39.26±1.20	37.18±2.28
UK	58.84±0.59	59.36±0.72	<u>59.66±0.76</u>	59.49±0.56	<b>59.72±0.31</b>	58.96±0.07
PT	<b>66.65±0.44</b>	<u>66.49±0.33</u>	65.62±0.43	65.99±0.32	65.61±0.15	64.40±0.33
IT	<b>58.65±0.49</b>	57.80±0.39	57.43±0.42	58.00±0.42	57.34±0.34	58.02±0.27
ES	<u>64.07±0.51</u>	63.62±0.25	63.87±0.61	<b>64.29±0.36</b>	63.91±0.27	62.48±0.29
CS	57.20±0.14	55.80±0.51	<u>57.95±0.52</u>	<b>58.02±0.21</b>	57.62±0.28	56.59±0.36
EL	<b>60.86±0.56</b>	60.23±0.40	<u>60.56±0.69</u>	60.74±0.48	<b>60.86±0.34</b>	59.76±0.45
FI	<b>55.40±0.16</b>	<u>55.29±0.54</u>	54.87±0.40	54.62±0.32	54.88±0.20	54.62±0.20
RU	60.12±0.36	<b>61.13±0.50</b>	59.98±0.34	60.14±0.16	<u>60.30±0.22</u>	59.73±0.25
NL	63.21±0.38	62.22±0.58	62.94±0.21	<b>63.53±0.64</b>	<u>62.97±0.37</u>	62.47±0.36
TE	43.46±2.02	43.88±1.57	<u>46.08±1.07</u>	<b>46.96±1.49</b>	<b>46.96±1.82</b>	45.95±0.51
DE	57.70±0.43	56.98±1.13	<b>58.61±0.30</b>	<u>58.52±0.26</u>	<u>58.52±0.18</u>	57.72±0.23
PL	58.16±0.35	57.28±0.46	<b>59.08±0.31</b>	<u>59.00±0.44</u>	58.92±0.52	57.73±0.23
SC	18	16	-	-	-	-

Table 16: F1 scores (%) of best models evaluated on UPB v2 test set with predicted parsers. The bold score and underlined score indicate the highest and second-highest scores. AVG indicates the average F1 scores of a specific model evaluated in target languages. SC indicates the superiority score of each model.

	Frozen mBERT		Fine-Tuned mBERT	
	SAN-RPRs	SAPR-RPRs	SAN-RPRs	SAPR-RPRs
EN	78.26±0.40	78.11±0.42	79.37±0.60	79.04±0.57
AVG	52.73±0.40	52.64±0.38	54.01±0.95	53.82±1.17

Table 17: F1 scores (%) of SAN-RPRs and SAPR-RPRs with frozen mBERT and fine-tuned mBERT. AVG indicates the average F1 scores of a specific model evaluated in target languages.