

---

# Multilingual Pre-training with Universal Dependency Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The pre-trained language model (PrLM) demonstrates domination in downstream  
2 natural language processing tasks, in which multilingual PrLM takes advantage of  
3 language universality to alleviate the issue of limited resources for low-resource  
4 languages. Despite its successes, the performance of multilingual PrLM is still un-  
5 satisfactory, when multilingual PrLMs only focus on plain text and ignore obvious  
6 universal linguistic structure clues. Existing PrLMs have shown that monolingual  
7 linguistic structure knowledge may bring about better performance. Thus we pro-  
8 pose a novel multilingual PrLM that supports both explicit universal dependency  
9 parsing and implicit language modeling. Syntax in terms of universal dependency  
10 parse serves as not only pre-training objective but also learned representation in  
11 our model, which brings unprecedented PrLM interpretability and convenience  
12 in downstream task use. Our model outperforms two popular multilingual PrLM,  
13 multilingual-BERT and XLM-R, on cross-lingual natural language understanding  
14 (NLU) benchmarks and linguistic structure parsing datasets, demonstrating the  
15 effectiveness and stronger cross-lingual modeling capabilities of our approach.

## 16 1 Introduction

17 The pre-trained language model (PrLM) such as BERT [1] and many kinds of its variants [2, 3, 4] have  
18 proved their effectiveness in many downstream natural language processing (NLP) tasks including  
19 semantic textual similarity [5], question answering [6], sentiment classification [7], linguistic structure  
20 [4, 8] and so on. Most of these PrLM are aimed at languages that with a large amount of available  
21 linguistic resources and are widely used, such as English. However, it is not realistic to train an  
22 individual PrLM for all languages, especially for those low-resource languages. As a result, several  
23 multilingual PrLMs which take advantage of language universality have been published. Multilingual-  
24 BERT (m-BERT) [1] has shown good cross-lingual performance on several NLP tasks. But some  
25 work [9] has shown that, m-BERT does not have equally high-quality representation for all of the  
26 languages it includes and for some languages, m-BERT performs worse than non-BERT. Improvement  
27 can be made by explicit cross-language signals including bitext (XLM) [10] and word translation  
28 pairs from a dictionary [11]. This suggests that the effectiveness of multilingual PrLM can be further  
29 improved by integrating more explicit universal linguistic characteristics.

30 Existing PrLMs [12, 13] have tried to incorporate monolingual linguistic structure knowledge to  
31 improve the performance across multiple linguistics tasks by Multi-Task Learning (MTL) [14].  
32 However, the combination of universal linguistic structure knowledge has not been explored in the  
33 multilingual area. Learning universal knowledge across languages is more complex than learning  
34 monolingual knowledge, so a better integrating method than MTL needs to be explored.

35 Syntactic dependency parsing disclosing syntactic relations between words in a sentence, has been  
36 found to be extremely useful for many NLP tasks [15, 16, 17]. The syntactic dependency parsing

37 is also limited by low-resource languages. To meet the huge demand for training syntactic parser  
38 among various languages, the project of universal dependencies (UD) Treebanks was launched [18]  
39 which provides a uniform syntactic parsing structure for different languages. Therefore, UD offers an  
40 excellent linguistic structure for the multilingual PrLM as a universal structure characteristic.

41 In this paper, we propose a multilingual PrLM that supports both explicit universal dependency  
42 parsing and implicit language modeling. Unlike using MTL in monolingual works, we embed a  
43 parsing scorer in our PrLM, and directly optimizes this scorer and the encoders below it with UD  
44 pre-training; meantime, we propose and structural encoder to encode the predicted structure given by  
45 the parsing scorer and integrated it into the final representation for other pre-training or downstream  
46 training process. Our approach can be smoothly applied to a variety of multilingual PrLM such as  
47 m-BERT [1] and XLM-R [19].

48 To verify the cross-lingual modeling capabilities of our model, we carry on the experiments on  
49 both cross-lingual NLU benchmarks: XNLI and XQuAD, and linguistic structure parsing datasets:  
50 UD<sup>1</sup> v2.7, SPMRL'14 [20], English Penn Treebank (PTB) 3.0 [21] and the Chinese Penn Treebank  
51 (CTB) 5.1 [22]. Our empirical results show that universal structure knowledge learnt and integrated  
52 can indeed help the multilingual PrLM obtain better universal linguistic word representations and  
53 outperform m-BERT and XLM-R baselines in all the above tasks.

## 54 **2 Related Work**

55 Previous works have tried to improve mono-lingual PrLM by introducing linguistic knowledge  
56 through parallel multi-task learning in pre-training process. These tasks can be roughly divided into  
57 the following three categories: word-aware pre-training tasks, structure-aware pre-training tasks  
58 and semantic-aware pre-training tasks. Word-aware tasks have word structure objectives such as  
59 word-level ordering [23] and word definition objectives [24] such as knowledge masking task [2].  
60 Structure-aware tasks include sentence structure objectives such as sentence-level ordering [23] and  
61 sentence distance task [2]. Semantic-aware tasks involve discourse relation [2], syntactic constituent  
62 and dependency parsing and span and dependency semantic role labeling [13]. Apart from this,  
63 linguistic knowledge imported by richer positional information [25] can also enhance the PrLM.  
64 However, these works all merely set the linguistic tasks as pre-training objectives or features, which  
65 does not fully conform with the nature of multilingual pre-training. In our multilingual PrLM, we  
66 use the universal structure knowledge not only for pre-training objective but also for improving the  
67 contextual representation. In addition, our model gets rid of the limitation of the external parser, no  
68 matter it is in the pre-training stage or the downstream application stage because our pre-training  
69 language model can be also conveniently adopted as a multilingual UD parser.

70 There are also many improvements to the multilingual PrLM such as explicit cross-lingual knowledge  
71 including bitext [10, 26] and word translation pairs from a dictionary [11]. Cross-lingual pre-training  
72 tasks including cross-lingual word recovery, cross-lingual paraphrase classification and cross-lingual  
73 masked language model can take advantage of bitext to learn the mappings among different languages  
74 from more perspectives [26]. Explicit cross-lingual knowledge induced from a bitext [27] has also  
75 been used to make all source languages share the same feature space. However, as far as we know,  
76 there are few works that directly integrate linguistic knowledge into the multilingual PrLM to make a  
77 further improvement.

## 78 **3 Universal Dependency as Language Modeling**

79 In this work, we chose UD parse as our universal structure knowledge. Our model includes five  
80 modules: token representation, multi-layer Transformers, universal structure learning (USL) layer,  
81 universal structure integration layer and pre-training objectives. Please refer to Appendix A.1 for the  
82 figure of our full model architecture.

### 83 **3.1 Token Representation**

84 Take BERT [1] as an example, in the token representation layer, given input sentence  $X$ , the sentence  
85 is concatenated with two special tokens “[CLS]” and “[SEP]”: [CLS],  $x_1, x_2, \dots, x_N$ , [SEP], and

<sup>1</sup><https://lindat.cz/repository/xmlui/handle/11234/1-3424>

86 [CLS] is also used as the dummy ROOT node in UD training process. The input  $X$  is mapped into  
 87 a sequence of input embedding vectors  $[e_1, e_2, \dots, e_{|X|}]$ , one for each token, which is a sum of the  
 88 corresponding word and positional embeddings.

89 Since the UD parse tree is annotated in word-level and the input sequence  $X$  in PrLM is based on  
 90 subword tokenization, to make the representations in the USL layer compatible with the whole model,  
 91 we propose a conversion strategy extending the parsing tree from the word level:  $\hat{Y}$  to the subword  
 92 level:  $Y$ , and use the subword-level parsing tree as the training objective for USL. Detailed strategy  
 93 description is shown in Appendix A.2.

### 94 3.2 Multi-layer Transformers

95 The multi-layer Transformers architecture in our model is adapted from Vaswani et al.[28], which  
 96 transforms the input embedding vectors into a sequence of contextualized representation vectors  
 97  $\mathbf{H} = [h_1, h_2, \dots, h_{|X|}]$  shared across different tasks. We use a Transformer architecture with  $L$  layers,  
 98  $A$  self-attention heads for each block and hidden size  $H$ :

$$\mathbf{H}^{(L)} = \text{Transformers}(\text{Emb}(X) + \text{PosEncoding}(X))$$

### 99 3.3 Universal Structure Learning

100 Our USL layer follows the state-of-the-art graph-based deep biaffine dependency parser [29]. We  
 101 replace the BiLSTM encoder with the multi-layer Transformers architecture and use the hidden state  
 102 of its last layer as the output of encoder  $\mathbf{H}^{(L)} = [h_1^{(L)}, h_2^{(L)}, \dots, h_{|X|}^{(L)}]$ .

103 For both arc and label predictions, two separate MLPs are used to distinguish two kinds of low-  
 104 dimensional vectors as head and dependent representations respectively.

$$r_i^m = \text{ReLU}(\text{MLP}^m(h_i^{(L)})), m \in [\text{head}, \text{dep}], i = 1, 2, \dots, |X|$$

105 The scores of all possible head-dependent pairs for arc and all head-dependent-label triples for label  
 106 are computed via the Variable-class biaffine classifier [29]:

$$R_m = [r_1^m; r_2^m; \dots; r_{|X|}^m], m \in [\text{head}, \text{dep}]$$

$$S^k = \text{Softmax}(R_{\text{dep}}^T U_1 R_{\text{head}} + u_2^T R_{\text{head}} + u_3^T R_{\text{dep}} + b), k \in [\text{arc}, \text{label}].$$

107 For arc,  $U_1 \in \mathcal{R}^{H_{\text{dep}} \times H_{\text{head}}}$ ,  $u_2 \in \mathcal{R}^{H_{\text{head}}}$ ,  $u_3 \in \mathcal{R}^{H_{\text{dep}}}$ . For label,  $U_1 \in \mathcal{R}^{|D| \times H_{\text{dep}} \times H_{\text{head}}}$ ,  $u_2 \in$   
 108  $\mathcal{R}^{|D| \times H_{\text{head}}}$ ,  $u_3 \in \mathcal{R}^{|D| \times H_{\text{dep}}}$  where  $H_{\text{head}}$  is the dimension of the head representations,  $H_{\text{dep}}$  is the  
 109 dimension of the dependent representations and  $D$  is the label set. So that  $S^{\text{arc}} \in \mathcal{R}^{B \times |X|_{\text{dep}} \times |X|_{\text{head}}}$   
 110 and  $S^{\text{label}} \in \mathcal{R}^{B \times |D| \times |X|_{\text{dep}} \times |X|_{\text{head}}}$  where  $B$  denotes the batch size.

During training, we aim to optimize the following probability for UD parsing:

$$P_\theta(Y|X) = \prod_{i=1}^{|X|} P_\theta(y_i^{\text{label}} | x_i, y_i^{\text{arc}}) P_\theta(y_i^{\text{arc}} | x_i),$$

111 where  $\theta$  denotes the learnable parameters and  $y_i^{\text{arc}}, y_i^{\text{label}}$  denote the gold-standard head and depen-  
 112 dency relation for subword  $x_i$  in subword-level parsing tree  $Y$ . The training objective for UD parsing  
 113 is the cross-entropy, which minimizes the negative log-likelihood:

$$\mathcal{L}^{USL} = - \sum_{i=1}^{|X|} (\log P_\theta(y_i^{\text{arc}} | x_i) + \log P_\theta(y_i^{\text{label}} | x_i, y_i^{\text{arc}})).$$

114 For evaluation, we restore the subword-level score tensors:  $S^{\text{arc}}$  and  $S^{\text{label}}$  to word-level:  $\hat{S}^{\text{arc}}$  and  
 115  $\hat{S}^{\text{label}}$  by extracting the first subword of each word. Then, we judge whether the prediction result of  
 116  $\hat{S}^{\text{arc}}$  is a valid parsing tree. If so, we directly extract the corresponding prediction label from  $\hat{S}^{\text{label}}$ .  
 117 Otherwise, we use the max spanning tree (MST) algorithm to find the maximum spanning tree based  
 118 on  $\hat{S}^{\text{arc}}$ .

---

**Algorithm 1:** Training Process

---

**Input:** MLM training data  $\hat{X}_{MLM}$ , UD Treebanks  $(\hat{X}_{UD}, \hat{Y}_{UD})$ , Parameters:  $\theta = (\rho, \gamma, \omega, \phi)$ , Probability of training USL:  $p$ .

```
1  $D_{MLM}, X_{UD}, Y_{UD} = \mathbf{Token}(\hat{X}_{MLM} \cup \hat{X}_{UD}), \mathbf{Token}(\hat{X}_{UD}), \mathbf{Strategy}(\hat{Y}_{UD}, \mathbf{Token})$  ;
2 Initialize  $\theta_0$  randomly ;
3 for  $t \leftarrow 1$  to  $m$  do
4    $\theta_{t+1} \leftarrow \mathbf{Opt}(\theta_t, \mathcal{L}^{USL}(\rho_t + \gamma_t, X_{UD}^t, Y_{UD}^t) + \mathcal{L}^{MLM}(\rho_t + \phi_t, D_{MLM}^t))$  ;
5 for  $t \leftarrow m$  to  $m + n$  do
6   if  $\mathbf{random.uniform}(0, 1) < p$  then
7      $\theta_{t+1} \leftarrow \mathbf{Opt}(\theta_t, \mathcal{L}^{USL}(\rho_t + \gamma_t, X_{UD}^t, Y_{UD}^t) + \mathcal{L}^{MLM}(\rho_t + \gamma_t + \omega_t + \phi_t, D_{MLM}^t))$  ;
8   else
9      $\theta_{t+1} \leftarrow \mathbf{Opt}(\theta_t, \mathcal{L}^{MLM}(\rho_t + \gamma_t + \omega_t + \phi_t, D_{MLM}^t))$  ;
```

**Output:**  $\theta_n$ .

---

### 119 3.4 Universal Structure Integration

120 In order to better integrate linguistic structure knowledge into the output representation of our PrLM,  
121 we propose the universal structure integration layer, which combines  $S^{arc}$  and  $S^{label}$  obtained by the  
122 USL layer with  $\mathbf{H}^{(L)}$  as the final output representations.

123 We first combine  $S^{arc}$  and  $S^{label}$  into a full label scoring matrix  $S^{US}$  by dot product.  $S^{US} \in$   
124  $\mathcal{R}^{B \times |X|_{dep} \times |X|_{head} \times |D|}$ , in fact, stores information about the label probability of each head-  
125 dependent pair in the sentence. Then we use this label scoring matrix  $S^{US}$  as the attention score to  
126 obtain a dependency head and label specific representation by product summation operation to  $\mathbf{H}^{(L)}$ ,  
127 the result is  $\hat{\mathbf{H}}^{US} \in \mathcal{R}^{B \times |X| \times H \times |D|}$ .

$$S^{US} = S^{arc} \cdot S^{label}, \quad \hat{\mathbf{H}}^{US} := S_{bijk}^{US} \times \mathbf{H}_{bih}^{(L)} \rightarrow O_{bihk},$$

128 where  $[\cdot]_{bijk} \times [\cdot]_{bih} \rightarrow [\cdot]_{bihk}$  indicates the Einstein summation notation.

129 Then we employ a weight tensor  $W \in \mathcal{R}^{|D| \times H \times H}$  to aggregate and map the head and label specific  
130 representation to the final dependency tree-aware representation:  $\mathbf{H}^{US} \in \mathcal{R}^{B \times |X| \times H}$ .

$$\mathbf{H}^O := \hat{\mathbf{H}}_{bihk}^{US} \times W_{khm} \rightarrow O_{bihk}, \quad \mathbf{H}^{US} = \mathbf{GELU}(\mathbf{Linear}(\mathbf{H}^O)),$$

131 where  $[\cdot]_{bihk} \times [\cdot]_{khm} \rightarrow [\cdot]_{bihk}$  indicates the Einstein summation notation, the dimensions of  
132  $\mathbf{H}^O \in \mathcal{R}^{B \times |X| \times H \times |D|}$ , and will be flattened to shape  $\mathcal{R}^{B \times |X| \times (H \times |D|)}$  before input to **Linear**.

133 We also do a residual connection to avoid losing the information in  $\mathbf{H}^{(L)}$ , and we use an additional  
134 Transformer layer to get the final representation.

$$\mathbf{H}^{(L+1)} = \mathbf{Transformer}(\mathbf{H}^{US} + \mathbf{H}^{(L)})$$

### 135 3.5 Pre-training Objectives

136 We use Masked LM (MLM) as the only pre-training objective other than USL. In MLM, a random  
137 sample of the tokens in the input sequence is selected and replaced with the special token [MASK]. As  
138 described in BERT [1], 15% of the input tokens are uniformly selected for possible replacement. Of  
139 the selected tokens, 80% are replaced with [MASK], 10% are left unchanged, and 10% are replaced  
140 by a randomly selected vocabulary token. The MLM objective is a cross-entropy loss on predicting  
141 the masked tokens:

$$\mathcal{L}^{MLM} = - \sum_{x \in m(X)} \log P_\theta(x | X \setminus m(X)),$$

142 where  $m(X)$  and  $X \setminus m(X)$  denote the masked words from  $X$  and the rest of words respectively. The  
143 training loss of our model is calculated by adding the losses of USL and MLM.

$$\mathcal{L} = \mathcal{L}^{USL} + \mathcal{L}^{MLM}$$

### 144 3.6 Training Details

145 Algorithm 1 shows the training process of our model, in which **Token** denotes the tokenization,  
146 **Strategy** denotes the strategy creating subword-level parsing tree, **Opt** denotes the optimization  
147 strategy and  $\rho, \gamma, \omega, \phi$  denote the learnable parameters in Transformer encoder (and token representa-  
148 tion), USL layer, universal structure integration layer and MLM decoder respectively. We randomly  
149 initialize the model parameters  $\theta$ . In the first  $m$  epochs, USL and MLM are trained as two parallel  
150 tasks sharing the parameters in Transformer encoder. That is to say, we do not use universal structure  
151 integration layer to integrate universal structure knowledge for MLM because the parsing capability  
152 of the model is too weak in this phase. In the next  $n$  steps, we open the universal structure integration  
153 layer. At this time, the USL objective has converged well, so we reduce the frequency of training  
154 USL appropriately according to a certain probability  $p$ .

## 155 4 Experiments

### 156 4.1 Setup

157 **Pre-training Data** Similar to that of m-BERT, we chose the top 104 languages with the largest  
158 Wikipedias and apply exponentially smoothed weighting to these languages to balance the Wikipedia  
159 size of each language for training MLM in the resulted UD-BERT according to our proposed approach.  
160 We do not use the NSP objective for UD-BERT. For UD-XLM-R<sub>base</sub> and UD-XLM-R<sub>large</sub>, we use the  
161 same training set from CommonCrawl Corpus as in Conneau et al.(2019) [19]. We also used MLM as  
162 the only objective other than USL for UD-XLM-R<sub>base</sub> and UD-XLM-R<sub>large</sub>. For structure learning, we  
163 concatenate all the training TreeBanks covering 60 languages in Universal Dependencies Treebanks  
164 (v2.2) [30] as the training set. In addition, we add the sentences in the training TreeBanks of UD  
165 to the training set of MLM letting language modeling directly help our model learn the structure  
166 knowledge.

167 **Pre-training Settings** Our UD-BERT and UD-XLM-R<sub>base</sub> use a Transformer architecture with  
168  $L = 12, H = 768$  and  $A = 12$  with a vocabulary of 110k and 250k respectively. Our UD-XLM-R<sub>large</sub>  
169 uses a large Transformer architecture with  $L = 24, H = 1024$  and  $A = 16$  with a 250k vocabulary.  
170 We use WordPiece [31] tokenization of UD-BERT and SentencePiece [32] tokenization for UD-  
171 XLM-R<sub>base</sub> and UD-XLM-R<sub>large</sub>. We randomly initialize the model parameters rather than using  
172 the pre-trained parameters of m-BERT or XLM-R. We train our models with the Adam optimizer  
173 [33] using the parameters:  $\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1e - 6$  and  $L_2$  weight decay of 0.01, a linear  
174 warmup [28], GELU activation [34] and a dropout rate of 0.1. Models are trained for  $m = 600,000$   
175 and  $n = 600,000$  epochs in each phase respectively, with  $B = 6$ , and the probability of training  
176 USL in the second phase is  $p = 0.8$ . The max sequence length of MLM is 384 and the max sequence  
177 length for UD parsing is 256. In the USL layer, we set  $H_{head} = 128$  and  $H_{dep} = 64$ .

178 **XNLI: Cross-lingual Natural Language Inference** takes two sentences as input and determines  
179 whether one entails the other, contradicts it or neither. XNLI is defined on 15 languages. Each  
180 language contains a development set with 2,490 sentence pairs and a test set with 5,010 sentence  
181 pairs. Only English has training data, which is a crowd-sourced collection of 433k sentence pairs  
182 from MultiNLI [35]. The performance is evaluated by classification accuracy.

183 **XQuAD: Cross-lingual Question Answering Dataset** [36] is a benchmark dataset for evaluating  
184 cross-lingual question answering performance. The dataset consists of a subset of 240 paragraphs  
185 and 1,190 question-answer pairs from the development set of SQuAD v1.1 [6] together with their  
186 professional translations into 10 languages. The performance is evaluated by F1 and exact match  
187 (EM) scores.

188 **Universal Linguistic Structure Parsing:** For universal dependency parsing, we evaluate our  
189 model on 22 languages in Universal Dependencies Treebanks (v2.7) [30] whose detail information  
190 is shown in Appendix 8. We use the graph-based deep biaffine dependency parsing model [29] as  
191 our dependency parser. For monolingual evaluation, we train a model for each language on their  
192 training set using word, character and POS tag embeddings of dimension 100 and representation  
193 from PrLM of dimension 300. For cross-lingual evaluation, we train a single model on English

194 training set using POS tag and representation from PrLM. Unlabeled Attachment Scores (UAS) and  
 195 Labeled Attachment Scores (LAS) are adopted as the evaluation metrics. For universal constituent  
 196 parsing, we use the constituent-based datasets in SPMRL Shared Task 2014 which focus on parsing  
 197 nine morphologically rich languages, from different language families, in both a constituent-based  
 198 and dependency-based format. We also evaluate our model on PTB and CTB. We use the CRF  
 199 constituency parsing model [37] as our constituent parser with word and character embeddings of  
 200 dimension 100 and representation from PrLM of dimension 300 for monolingual evaluation. For  
 201 cross-lingual evaluation, we train a single model on PTB using only the representation from PrLM.

## 202 4.2 Results and Analysis

203 To evaluate the multilingual performance and cross-lingual transfer effect of the PrLM that learns the  
 204 universal linguistic structure and integrates the universal linguistic structure into the representation ex-  
 205 plicitly, we conducted experiments on the two typical tasks: universal natural language understanding  
 206 and the universal linguistic structure parsing.

207 **Universal Natural Language Understanding** In Table 1, we show the cross-lingual transfer re-  
 208 sults (*Cross-Transfer*) of the baselines and our proposed model on the cross-lingual text classification  
 209 benchmark - XNLI. Meanwhile, we also list the multilingual performance (*Train-Trans-FT*, *Test-*  
 210 *Trans-Eval*, and *All-FT*) as a reference. First, compare the results of the source language - English, our  
 211 UD-BERT, UD-XLM-R<sub>base</sub>, and UD-XLM-R<sub>large</sub> outperform the corresponding m-BERT, XLM-R<sub>base</sub>,  
 212 and XLM-R<sub>large</sub> baselines, demonstrating the universal linguistic structure as pre-training objective  
 213 and explicitly syntactic structure integration improve the model pre-training and final representations.

Table 1: Results on cross-lingual text classification task. We report the accuracy on each of the 15 XNLI languages and the average accuracy. Results with † are from [26].

Model	en	fr	es	de	el	bg	ru	tr	ar	vi	th	zh	hi	sw	ur	Avg
<i>Train-Trans-FT: Fine-tune multilingual model on each training set translated from English</i>																
XLM [10]	82.9	77.6	77.9	77.9	77.1	75.7	75.5	72.6	71.2	75.8	73.1	76.2	70.4	66.5	62.4	74.2
<i>Test-Trans-Eval: Translate test sets to English and use English-only model for evaluation</i>																
BERT-en	88.8	81.4	82.3	80.1	80.3	80.9	76.2	76.0	75.4	72.0	71.9	75.6	70.0	65.8	65.8	76.2
RoBERTa	91.3	82.9	84.3	81.2	81.7	83.1	78.3	76.8	76.6	74.2	74.1	77.5	70.9	66.7	66.8	77.8
<i>All-FT: Fine-tune multilingual model on all training sets</i>																
XLM [10]	84.5	80.1	81.3	79.3	78.6	79.4	77.5	75.2	75.6	78.3	75.7	78.3	72.1	69.2	67.7	76.9
XLM [10]†	85.0	80.8	81.3	80.3	79.1	80.9	78.3	75.6	77.6	78.5	76.0	79.5	72.9	72.8	68.5	77.8
Unicoder [26]	85.6	81.1	82.3	80.9	79.5	81.4	79.7	76.8	78.2	77.9	77.1	80.5	73.4	73.8	69.6	78.5
XLM-R <sub>base</sub>	85.4	81.4	82.2	80.3	80.4	81.3	79.7	78.6	77.3	79.7	77.9	80.2	76.1	73.1	73.0	79.1
XLM-R <sub>large</sub>	89.1	85.1	86.6	85.7	85.3	85.9	83.5	83.2	83.1	83.7	81.5	83.7	81.6	78.0	78.1	83.6
<i>Cross-Transfer: Fine-tune multilingual model on English training set</i>																
XLM [10]	85.0	78.7	78.9	77.8	76.6	77.4	75.3	72.5	73.1	76.1	73.2	76.5	69.6	68.4	67.3	75.1
Unicoder [26]	85.1	79.0	79.4	77.8	77.2	77.2	76.3	72.8	73.5	76.4	73.6	76.2	69.4	69.7	66.7	75.4
m-BERT [1]	82.1	73.8	74.3	71.1	66.4	68.9	69.0	61.6	64.9	69.5	55.8	69.3	60.0	50.4	58.0	66.3
<b>UD-BERT</b>	<b>82.7</b>	<b>74.9</b>	<b>75.2</b>	<b>72.0</b>	<b>67.4</b>	<b>69.2</b>	<b>70.3</b>	<b>62.7</b>	<b>65.8</b>	<b>70.3</b>	<b>59.6</b>	<b>69.7</b>	<b>61.4</b>	<b>51.2</b>	<b>58.7</b>	<b>67.4</b>
XLM-R <sub>base</sub>	85.8	79.7	80.7	78.7	77.5	79.6	78.1	74.2	73.8	76.5	74.6	76.7	72.4	66.5	68.3	76.2
<b>UD-XLM-R<sub>base</sub></b>	<b>86.5</b>	<b>80.3</b>	<b>81.6</b>	<b>79.8</b>	<b>78.4</b>	<b>80.0</b>	<b>78.9</b>	<b>75.1</b>	<b>74.4</b>	<b>77.1</b>	<b>75.2</b>	<b>77.3</b>	<b>73.0</b>	<b>67.0</b>	<b>68.8</b>	<b>76.9</b>
XLM-R <sub>large</sub>	89.1	84.1	85.1	83.9	82.9	84.0	81.2	79.6	79.8	80.8	78.1	80.2	76.9	73.9	73.8	80.9
<b>UD-XLM-R<sub>large</sub></b>	<b>89.4</b>	<b>84.8</b>	<b>85.6</b>	<b>84.5</b>	<b>83.6</b>	<b>84.7</b>	<b>81.2</b>	<b>80.0</b>	<b>81.0</b>	<b>81.9</b>	<b>78.6</b>	<b>80.7</b>	<b>76.8</b>	<b>74.4</b>	<b>74.3</b>	<b>81.4</b>

214 Second, our UD-BERT and UD-XLM-R performed better in most cases of the 14 transferring target  
 215 languages. UD-BERT has an average increase of 1.1 when compared to the baseline, UD-XLM-R<sub>base</sub>  
 216 has an average increase of 0.7, and UD-XLM-R<sub>large</sub> has an average increase of 0.5. This improvement  
 217 highlights the fact that the cross-lingual transferring ability of our multilingual PrLM has improved as  
 218 a result of the employment of universal linguistic structures. Furthermore, the cross-lingual transfer  
 219 effect of our UD-XLM-R<sub>base</sub> outperforms BERT with a similar model structure and parameters, who  
 220 translated the test set for evaluation. The results of our UD-XLM-R<sub>large</sub> achieved better results than  
 221 all the methods which leveraging a monolingual language model on the translations, which shows  
 222 that cross-lingual transfer is a more promising mode.

223 Text classification is a relatively simple and intuitive NLU task. To further verify our conclusions, we  
 224 conducted experiments on a more complex task - cross-lingual Machine Reading Comprehension  
 225 (MRC). The results on XQuAD dataset are shown in Table 2. Similarly, we first compare the  
 226 MRC results on the source language English. The performance of UD-BERT and UD-XLM-R has

Table 2: F<sub>1</sub> / EM scores on XQuAD with English as the source language for each target language.

<i>Cross-Transfer</i>		en	ar	de	el	es	hi	ru	th	tr	vi	zh	ro	Avg
m-BERT	F <sub>1</sub>	83.5	61.5	70.6	62.6	75.5	59.2	71.3	42.7	55.4	69.5	58.0	72.7	65.2
	EM	72.2	45.1	54.0	44.9	56.9	46.0	53.3	33.5	40.1	49.6	48.3	59.9	50.3
UD-BERT	F <sub>1</sub>	<b>83.9</b>	<b>62.8</b>	<b>72.3</b>	<b>62.9</b>	<b>75.7</b>	59.0	<b>71.6</b>	<b>48.6</b>	<b>55.5</b>	<b>69.9</b>	<b>58.7</b>	<b>73.7</b>	<b>66.2</b>
	EM	<b>72.5</b>	<b>46.3</b>	<b>56.0</b>	<b>45.8</b>	<b>57.3</b>	<b>46.3</b>	<b>54.0</b>	<b>38.5</b>	<b>40.2</b>	<b>49.7</b>	<b>48.6</b>	<b>60.3</b>	<b>51.3</b>
XLM-R <sub>base</sub>	F <sub>1</sub>	83.6	66.8	74.4	73.0	76.4	68.2	74.3	66.5	68.3	73.7	51.3	77.8	71.2
	EM	72.1	49.1	60.1	55.7	58.3	51.7	58.1	56.7	52.8	53.8	42.0	62.8	56.1
UD-XLM-R <sub>base</sub>	F <sub>1</sub>	<b>84.0</b>	<b>69.1</b>	<b>74.9</b>	<b>73.5</b>	<b>77.0</b>	<b>68.4</b>	<b>74.3</b>	<b>66.9</b>	<b>69.3</b>	<b>74.1</b>	<b>51.8</b>	<b>78.0</b>	<b>71.8</b>
	EM	<b>72.5</b>	<b>51.2</b>	<b>60.5</b>	<b>56.0</b>	<b>58.3</b>	<b>51.9</b>	<b>58.5</b>	<b>57.1</b>	<b>52.9</b>	<b>54.2</b>	<b>42.6</b>	<b>63.3</b>	<b>56.6</b>
XLM-R <sub>large</sub>	F <sub>1</sub>	86.5	68.6	80.4	79.8	82.0	76.7	80.1	74.2	75.9	79.1	59.3	83.6	77.2
	EM	75.7	49.0	63.4	61.7	63.9	59.7	64.3	62.8	59.3	59.0	50.0	69.7	61.5
UD-XLM-R <sub>large</sub>	F <sub>1</sub>	<b>86.8</b>	<b>75.2</b>	<b>80.9</b>	<b>80.0</b>	<b>82.3</b>	<b>77.1</b>	<b>80.3</b>	73.8	<b>76.3</b>	<b>79.5</b>	<b>59.4</b>	<b>83.9</b>	<b>78.0</b>
	EM	<b>75.9</b>	<b>58.2</b>	<b>63.8</b>	61.7	<b>64.0</b>	<b>59.8</b>	<b>64.5</b>	62.2	<b>60.5</b>	<b>59.8</b>	49.9	<b>69.9</b>	<b>62.5</b>

Table 3: The monolingual UD parsing results (UAS/LAS) on 22 UD Treebanks.

<i>All-FT</i>	m-BERT		Li et al. [38]		UD-BERT		XLM-R <sub>base</sub>		UD-XLM-R <sub>base</sub>		XLM-R <sub>large</sub>		UD-XLM-R <sub>large</sub>	
	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS
bg	94.75	90.88	95.43	91.27	<b>96.12</b>	<b>93.09</b>	96.42	93.56	<b>96.57</b>	<b>93.80</b>	96.53	93.72	<b>96.69</b>	<b>94.04</b>
ca	95.36	93.25	95.38	93.57	<b>95.61</b>	<b>94.27</b>	95.53	94.26	<b>95.78</b>	<b>94.61</b>	95.75	94.58	<b>95.94</b>	<b>94.72</b>
cs	94.38	91.62	95.08	91.52	<b>95.62</b>	<b>93.17</b>	95.60	93.30	<b>95.82</b>	<b>93.43</b>	95.87	93.69	<b>95.99</b>	<b>93.87</b>
nl	94.74	92.72	95.32	92.82	<b>95.38</b>	<b>93.61</b>	95.36	93.44	<b>95.79</b>	<b>93.82</b>	95.63	93.78	<b>96.13</b>	<b>94.50</b>
en	92.52	91.29	<b>93.17</b>	90.79	93.01	<b>91.43</b>	93.60	91.83	<b>94.15</b>	<b>92.55</b>	93.47	91.73	<b>94.19</b>	<b>92.54</b>
et	90.88	88.95	—	—	<b>91.65</b>	<b>89.73</b>	92.53	90.78	<b>92.78</b>	<b>91.02</b>	93.16	91.50	<b>93.28</b>	<b>91.68</b>
fi	92.98	90.65	—	—	<b>94.07</b>	<b>91.89</b>	94.99	93.14	<b>95.15</b>	<b>93.44</b>	95.17	93.46	<b>95.66</b>	<b>94.01</b>
fr	94.12	90.75	93.95	91.42	<b>96.08</b>	<b>94.24</b>	96.10	94.34	<b>96.48</b>	<b>94.69</b>	96.01	94.15	<b>96.53</b>	<b>94.75</b>
de	90.77	86.83	<b>91.33</b>	87.27	91.08	<b>87.35</b>	91.30	87.42	<b>91.58</b>	<b>87.63</b>	91.39	87.50	<b>91.72</b>	<b>87.94</b>
he	92.32	89.95	—	—	<b>93.45</b>	<b>91.06</b>	93.50	91.41	<b>93.87</b>	<b>91.64</b>	93.67	91.48	<b>93.99</b>	<b>91.79</b>
hi	96.54	94.22	—	—	<b>96.71</b>	<b>94.41</b>	96.73	94.51	<b>97.03</b>	<b>94.99</b>	96.93	94.76	<b>97.17</b>	<b>95.09</b>
id	<b>88.29</b>	<b>83.97</b>	—	—	87.96	83.72	88.25	84.06	<b>88.51</b>	<b>84.25</b>	88.38	<b>84.27</b>	<b>88.48</b>	84.19
it	95.63	94.01	95.73	93.52	<b>96.32</b>	<b>95.16</b>	96.15	94.93	<b>96.70</b>	<b>95.33</b>	96.26	95.01	<b>96.82</b>	<b>95.60</b>
ko	90.73	88.27	—	—	<b>90.99</b>	<b>88.67</b>	91.33	89.25	<b>91.42</b>	<b>89.30</b>	92.15	89.79	<b>92.16</b>	<b>89.98</b>
la	85.69	81.84	—	—	<b>87.01</b>	<b>83.45</b>	86.64	82.99	<b>89.86</b>	<b>86.91</b>	86.97	83.29	<b>90.44</b>	<b>87.33</b>
lv	91.55	89.06	—	—	<b>92.06</b>	<b>89.72</b>	93.53	91.25	<b>93.88</b>	<b>91.50</b>	94.23	91.90	<b>94.56</b>	<b>92.52</b>
no	95.62	93.88	95.87	94.15	<b>96.23</b>	<b>95.17</b>	96.57	95.58	<b>96.71</b>	<b>95.63</b>	96.70	95.64	<b>96.74</b>	<b>95.75</b>
pl	98.15	96.54	—	—	<b>98.54</b>	<b>97.14</b>	98.51	97.22	<b>98.80</b>	<b>97.77</b>	98.47	97.08	<b>98.79</b>	<b>97.74</b>
ro	92.70	86.39	92.72	86.16	<b>94.12</b>	<b>89.39</b>	94.25	89.67	<b>94.60</b>	<b>90.11</b>	94.39	89.85	<b>94.73</b>	<b>90.32</b>
ru	95.26	94.00	<b>95.88</b>	94.26	95.74	<b>94.58</b>	96.25	95.26	<b>96.54</b>	<b>95.59</b>	96.38	95.45	<b>96.62</b>	<b>95.72</b>
sk	94.93	91.40	—	—	<b>96.62</b>	<b>93.95</b>	95.96	93.21	<b>97.49</b>	<b>95.53</b>	95.56	92.74	<b>97.26</b>	<b>95.34</b>
es	94.69	92.89	94.83	92.36	<b>94.92</b>	<b>93.23</b>	94.98	93.44	<b>95.38</b>	<b>93.80</b>	95.33	93.80	<b>95.55</b>	<b>94.09</b>
Avg	93.30	90.61	—	—	<b>94.06</b>	<b>91.75</b>	94.28	92.04	<b>94.77</b>	<b>92.61</b>	94.47	92.24	<b>94.97</b>	<b>92.89</b>

227 increased relative to the baseline, which verifies the conclusion that our universal linguistic structure  
228 improves the NLU ability of multilingual language model. For the 11 transferring target languages,  
229 we observed a similar improvement trend on UD-BERT and UD-XLM-R as in the XNLI task, and  
230 the improvement was even greater. The average improvement of UD-BERT, UD-XLM-R base and  
231 UD-XLM-R large reached 1.0, 0.6, 0.8 F<sub>1</sub> scores respectively. Among them, Arabic has the largest  
232 improvement, with 1.3, 2.3, and 6.6 F<sub>1</sub> scores respectively. All the results here reveal that universal  
233 syntactic structure information embedded is effective for cross-lingual MRC task.

234 **Universal Linguistic Structure Parsing** Since the universal dependency parsing structure and  
235 the dependency parse encoding structure are built into our multilingual PrLM, to demonstrate that  
236 our model learned how to extract linguistic structure and to encode the structure into the output  
237 representations, we performed the universal linguistic structure parsing on UD and multilingual  
238 constituent parsing tasks. It is worth noting that, to shield the influence of the increase in the number  
239 of parameters caused by the addition of the PrLM to the downstream task model, we kept all the  
240 PrLM parameters frozen in the universal NLU evaluation.

241 In Table 3, we evaluated the effect of using the multilingual PrLM to enhance the parsing model on  
242 22 languages of UD dataset respectively, in order to analyze how many features useful (i.e., syntactic-  
243 aware features) for parsing are provided by the output representation of the various multilingual  
244 PrLM. We also listed previous state-of-the-art (SOTA) results from [38] on UD for comparison.

Table 4: The cross-lingual UAS/LAS results on 22 languages of UD Treebanks.

<i>Cross-Transfer</i>	en	bg	ca	cs	nl	et	fi	fr
m-BERT	92.52 / 91.29	83.80 / 72.77	79.60 / 69.00	73.95 / 61.23	77.76 / 69.02	73.11 / 50.90	75.71 / 54.96	83.18 / 70.76
<b>UD-BERT</b>	93.01 / 91.43	89.64 / 79.02	87.18 / 76.06	86.07 / 70.38	86.79 / 78.74	82.38 / 58.53	83.86 / 60.97	89.07 / 75.86
XLM-R <sub>large</sub>	93.10 / 91.32	88.67 / 77.91	85.00 / 73.71	77.80 / 65.12	82.10 / 72.79	78.66 / 57.67	80.97 / 60.63	88.27 / 74.71
<b>UD-XLM-R<sub>large</sub></b>	94.19 / 92.54	92.66 / 82.15	90.76 / 79.35	88.98 / 75.67	90.94 / 82.33	87.57 / 65.90	90.09 / 68.95	93.19 / 78.93
	en	de	he	hi	id	it	ko	la
m-BERT	92.52 / 91.29	76.54 / 66.00	73.48 / 47.26	43.77 / 29.83	57.06 / 47.82	87.41 / 80.93	36.97 / 23.44	52.39 / 36.42
<b>UD-BERT</b>	93.01 / 91.43	84.54 / 74.88	87.33 / 55.98	66.99 / 42.82	76.18 / 62.37	92.74 / 86.59	53.89 / 36.94	71.66 / 51.91
XLM-R <sub>large</sub>	93.10 / 91.32	81.65 / 71.34	73.86 / 48.66	48.08 / 31.66	60.57 / 51.32	91.02 / 84.31	40.64 / 25.50	60.95 / 42.44
<b>UD-XLM-R<sub>large</sub></b>	94.19 / 92.54	88.05 / 79.34	86.42 / 56.69	65.57 / 47.13	77.65 / 64.52	95.39 / 88.83	52.98 / 37.52	76.26 / 57.14
	en	lv	no	pl	ro	ru	sk	es
m-BERT	92.52 / 91.29	76.51 / 54.52	87.28 / 78.67	88.76 / 76.43	75.94 / 62.83	72.73 / 62.40	79.14 / 67.86	80.08 / 70.64
<b>UD-BERT</b>	93.01 / 91.43	83.07 / 60.94	90.41 / 82.22	92.97 / 81.26	87.50 / 72.59	85.80 / 72.58	88.69 / 74.83	85.93 / 76.15
XLM-R <sub>large</sub>	93.10 / 91.32	83.94 / 62.52	90.50 / 82.56	90.68 / 77.90	81.75 / 69.01	76.64 / 66.58	81.63 / 70.60	85.00 / 74.36
<b>UD-XLM-R<sub>large</sub></b>	94.19 / 92.54	91.21 / 69.64	94.14 / 86.37	95.48 / 83.91	91.28 / 77.80	91.29 / 80.64	92.02 / 80.56	90.37 / 79.78

245 Comparing UD-BERT and m-BERT, our UD-BERT obtained generally better results, except for  
 246 Indonesian, with a 0.76/1.14 UAS/LAS average improvement. Compared with the previous SOTA  
 247 work [38], which also used m-BERT but with a more complex model architecture, although our  
 248 baseline is weaker after enhanced with UD-BERT, we surpassed [38] in most languages. This shows  
 249 that integrating more effective features is also a good option for improving the parsing performance  
 250 in addition to model designs. When we further compare UD-BERT to stronger XLM-R<sub>base</sub> and  
 251 XLM-R<sub>large</sub> baselines, we found that XLM-R<sub>base</sub> achieved similar results as UD-BERT. Since XLM-R  
 252 uses more data and longer pre-training time, from this point of view, UD-BERT has played a role  
 253 in reducing the data and time required for pre-training due to the addition of syntactic supervision  
 254 information. Further comparing UD-XLM-R with the baseline XLM-R, we found that our method is  
 255 still helpful on such strong baselines, which improves the parsing performance in each language. In  
 256 addition, the results on some languages on XLM-R<sub>base</sub> surpass XLM-R<sub>large</sub>.

257 Table 4 shows the zero-shot cross-lingual transfer experiments<sup>2</sup> of our PrLM on the 22 languages  
 258 of UD Treebanks. The findings from comparison reveal that our UD-BERT and UD-XLM-R have  
 259 significantly improved the transfer effect of this parser in which the average improvement of UD-  
 260 BERT and UD-XLM-R is 9.93/8.09, 8.86/8.32 UAS/LAS, respectively. This is due in large part to  
 261 the usage of UD annotations in the language model pre-training, but it also demonstrates that our  
 262 PrLM well learned UD parsing and encoded the structural information of the parse into the final  
 263 representations.

264 Though we verified the linguistic feature extraction and encoding capabilities of our PrLM in the  
 265 UD parsing, since the evaluation task is the same as in the pre-training, it is hard to illustrate this  
 266 effect on the cross-task linguistic parsing task, we thus conduct further exploration on multilingual  
 267 constituent parsing benchmarks. Table 5 shows the enhanced ability of the language model for  
 268 monolingual constituent parsing. The comparison demonstrates that with the help of additional  
 269 universal dependency features in the representations, the constituent parsing performance of UD-  
 270 BERT and UD-XLM-R has been greatly increased, and the average improvement has reached  
 271 0.83, 0.57, and 0.74 respectively. This reflects that the universal syntactic features contained in the  
 272 representations are very helpful for downstream tasks. In addition, with our syntax-aware multilingual  
 273 PrLM, the integration of syntactic tree is no longer necessary to change the downstream task model  
 274 for adding extra syntactic encoders, which will greatly reduce the cost of syntactic tree information  
 275 application.

276 The multilingual boosting ability of our PrLM is shown in Table 5. Similarly, we evaluated the PrLM’s  
 277 cross-lingual transfer ability on constituent parsing in Table 6, and because there are no available  
 278 universal constituent label annotations, we only report unlabeled F1. Again, our UD-BERT and UD-  
 279 XLM-R achieve an improvement in all transferring target languages, which shows that the explicit  
 280 universal features contained in our PrLM can be used to help downstream cross-lingual transfer  
 281 tasks. Since the universal dependency structure can be derived directly by additional inference in our  
 282 PrLM’s intermediate output, our good cross-lingual transfer ability is interpretable. Furthermore, our

<sup>2</sup>In the strict sense, this cannot be totally categorized as zero-shot because the language model part of the full parser has been pre-trained by UD annotations.

Table 5: Labeled F<sub>1</sub> scores on PTB, CTB, and SPMRL test sets.

<i>All-FT</i>	en	ar	eu	fr	de	he	hu	ko	pl	sv	zh	Avg
m-BERT	94.87	88.60	91.06	84.98	90.20	83.35	92.33	89.36	96.36	83.05	90.83	89.54
<b>UD-BERT</b>	<b>95.09</b>	<b>89.97</b>	<b>92.23</b>	<b>85.57</b>	<b>91.06</b>	<b>84.36</b>	<b>93.46</b>	<b>89.72</b>	<b>96.59</b>	<b>85.06</b>	<b>90.99</b>	<b>90.37</b>
XLM-R <sub>base</sub>	95.50	89.67	91.68	85.63	91.54	85.20	93.92	90.82	96.72	86.00	91.84	90.77
<b>UD-XLM-R<sub>base</sub></b>	<b>95.83</b>	<b>90.22</b>	<b>92.67</b>	<b>86.10</b>	<b>91.98</b>	<b>86.03</b>	<b>94.61</b>	<b>91.32</b>	<b>96.95</b>	<b>86.71</b>	<b>92.30</b>	<b>91.34</b>
XLM-R <sub>large</sub>	95.74	90.41	92.54	86.47	92.10	86.11	94.56	91.77	97.15	86.80	92.24	91.44
<b>UD-XLM-R<sub>large</sub></b>	<b>96.15</b>	<b>90.77</b>	<b>94.13</b>	<b>86.96</b>	<b>92.65</b>	<b>86.97</b>	<b>95.41</b>	<b>92.16</b>	<b>97.33</b>	<b>89.06</b>	<b>92.43</b>	<b>92.18</b>

Table 6: Unlabeled F<sub>1</sub> scores on PTB, CTB, and SPMRL test sets.

<i>Cross-Transfer</i>	en	ar	eu	fr	de	he	hu	ko	pl	sv	zh	Avg
m-BERT	95.54	24.22	32.61	59.40	44.36	46.01	57.86	33.99	37.47	71.69	56.87	50.91
<b>UD-BERT</b>	<b>95.93</b>	<b>27.09</b>	<b>35.50</b>	<b>60.87</b>	<b>45.24</b>	<b>46.88</b>	<b>58.99</b>	<b>35.67</b>	<b>39.23</b>	<b>72.20</b>	<b>57.74</b>	<b>52.30</b>
XLM-R <sub>base</sub>	96.19	28.17	34.34	59.70	44.98	47.05	60.26	38.95	38.06	73.68	55.36	52.43
<b>UD-XLM-R<sub>base</sub></b>	<b>96.60</b>	<b>30.46</b>	<b>37.75</b>	<b>61.36</b>	<b>45.52</b>	<b>47.42</b>	<b>62.00</b>	<b>42.94</b>	<b>39.83</b>	<b>73.90</b>	<b>59.97</b>	<b>54.34</b>
XLM-R <sub>large</sub>	96.44	20.56	34.23	59.81	45.04	47.01	59.37	36.72	37.87	73.43	54.34	51.35
<b>UD-XLM-R<sub>large</sub></b>	<b>96.73</b>	<b>23.44</b>	<b>37.63</b>	<b>61.18</b>	<b>45.71</b>	<b>47.30</b>	<b>60.46</b>	<b>41.26</b>	<b>38.42</b>	<b>73.81</b>	<b>59.58</b>	<b>53.23</b>

283 studies found that the transfer effect of XLM-R<sub>base</sub> is superior to that of XLM-R<sub>large</sub>, implying that  
 284 larger model parameters may not always result in higher cross-lingual transfer capabilities.

## 285 5 Ablation Study

286 In our multilingual PrLM, we propose both novel model structure and novel pre-training strategy,  
 287 whose effectiveness is verified by the ablation study in Table 7. From the results, if universal  
 288 dependency parsing is only used as an additional pre-training process (*w/o Universal Structure  
 289 Integration*), there is only a slight performance improvement compared to the baseline, which shows  
 290 that just as a pre-training approach is not enough to fuse the universal syntax implicitly. When we only  
 291 employ the model structure without UD pre-training (*w/o UD Pre-training*), this additional parser  
 292 structure, even though it is not trained by the UD supervision, can still improve the performance,  
 293 and the improvement is greater than that of *w/o Universal Structure Integration*. Combining these  
 294 two results with that of our full UD-BERT, we can conclude that the supervised learning of the UD  
 295 structure and its explicit integration into the representation are indispensable options for performance  
 296 improvement.

Table 7: Ablation studies on language **ar** and **eu** of SPMRL test sets.

Method	ar			eu		
	LP	LR	LF <sub>1</sub>	LP	LR	LF <sub>1</sub>
m-BERT	89.79	87.43	88.60	90.70	91.42	91.06
<b>UD-BERT</b>	90.97	88.86	89.97	91.79	92.67	92.23
<i>w/o Universal Structure Integration</i>	88.25	89.45	88.85	90.94	91.65	91.29
<i>w/o UD Pre-training (Structure Only)</i>	89.32	88.82	89.07	91.22	91.79	91.50

## 297 6 Conclusion

298 In this work, we propose a multilingual PrLM that supports both explicit universal structure learning  
 299 and implicit language modeling. We chose the universal dependency parses as our universal structure  
 300 knowledge and evaluate the cross-lingual modeling capability of our model on two cross-lingual  
 301 NLU tasks and four syntactic parsing datasets. Our model outperforms m-BERT and XLM-R in all  
 302 tasks and achieves state-of-the-art results on syntactic parsing. Unlike other works that use the syntax  
 303 as an objective or feature, our structure learning is used not only for a pre-training objective but also  
 304 for improving the representation, which makes our model both a PrLM and a universal dependency  
 305 parser. This will greatly change the way that downstream NLP tasks use syntax, because we have  
 306 explicitly integrated the syntactic knowledge into the representation of the PrLM.

307 **References**

- 308 [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional  
309 transformers for language understanding,” in *Proceedings of the 2019 Conference of the*  
310 *North American Chapter of the Association for Computational Linguistics: Human Language*  
311 *Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for  
312 Computational Linguistics, Jun. 2019, pp. 4171–4186.
- 313 [2] Y. Sun, S. Wang, Y. Li, S. Feng, H. Tian, H. Wu, and H. Wang, “Ernie 2.0: A continual  
314 pre-training framework for language understanding,” in *Proceedings of the AAAI Conference on*  
315 *Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 8968–8975.
- 316 [3] L. Dong, N. Yang, W. Wang, F. Wei, X. Liu, Y. Wang, J. Gao, M. Zhou, and  
317 H.-W. Hon, “Unified language model pre-training for natural language understanding  
318 and generation,” in *Advances in Neural Information Processing Systems*, H. Wallach,  
319 H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran  
320 Associates, Inc., 2019. [Online]. Available: [https://proceedings.neurips.cc/paper/2019/file/](https://proceedings.neurips.cc/paper/2019/file/c20bb2d9a50d5ac1f713f8b34d9aac5a-Paper.pdf)  
321 [c20bb2d9a50d5ac1f713f8b34d9aac5a-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/c20bb2d9a50d5ac1f713f8b34d9aac5a-Paper.pdf)
- 322 [4] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer,  
323 and V. Stoyanov, “Roberta: A robustly optimized BERT pretraining approach,” *CoRR*, vol.  
324 [abs/1907.11692](https://arxiv.org/abs/1907.11692), 2019. [Online]. Available: <http://arxiv.org/abs/1907.11692>
- 325 [5] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia, “SemEval-2017 task 1:  
326 Semantic textual similarity multilingual and crosslingual focused evaluation,” in *Proceedings*  
327 *of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*. Vancouver,  
328 Canada: Association for Computational Linguistics, Aug. 2017, pp. 1–14. [Online]. Available:  
329 <https://www.aclweb.org/anthology/S17-2001>
- 330 [6] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “SQuAD: 100,000+ questions for machine  
331 comprehension of text,” in *Proceedings of the 2016 Conference on Empirical Methods in*  
332 *Natural Language Processing*. Austin, Texas: Association for Computational Linguistics,  
333 Nov. 2016, pp. 2383–2392. [Online]. Available: <https://www.aclweb.org/anthology/D16-1264>
- 334 [7] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, “Recursive  
335 deep models for semantic compositionality over a sentiment treebank,” in *Proceedings of the*  
336 *2013 conference on empirical methods in natural language processing*, 2013, pp. 1631–1642.
- 337 [8] J. Zhou, Z. Li, and H. Zhao, “Parsing all: Syntax and semantics, dependencies and spans,” in  
338 *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association  
339 for Computational Linguistics, Nov. 2020, pp. 4438–4449.
- 340 [9] S. Wu and M. Dredze, “Are all languages created equal in multilingual BERT?” in *Proceedings of*  
341 *the 5th Workshop on Representation Learning for NLP*. Online: Association for Computational  
342 Linguistics, Jul. 2020, pp. 120–130.
- 343 [10] A. Conneau and G. Lample, “Cross-lingual language model pretraining,” in *Advances in*  
344 *Neural Information Processing Systems 32: Annual Conference on Neural Information*  
345 *Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M.  
346 Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds.,  
347 2019, pp. 7057–7067. [Online]. Available: [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/c04c19c2c2474dbf5f7ac4372c5b9af1-Abstract.html)  
348 [c04c19c2c2474dbf5f7ac4372c5b9af1-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/c04c19c2c2474dbf5f7ac4372c5b9af1-Abstract.html)
- 349 [11] A. Conneau, S. Wu, H. Li, L. Zettlemoyer, and V. Stoyanov, “Emerging cross-lingual structure  
350 in pretrained language models,” in *Proceedings of the 58th Annual Meeting of the Association*  
351 *for Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020,  
352 pp. 6022–6034. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.536>
- 353 [12] W. Wang, B. Bi, M. Yan, C. Wu, J. Xia, Z. Bao, L. Peng, and L. Si, “Structbert: Incorporating  
354 language structures into pre-training for deep language understanding,” in *8th International*  
355 *Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.  
356 OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=BJgQ4lSFPH>
- 357 [13] J. Zhou, Z. Zhang, H. Zhao, and S. Zhang, “Limit-bert : Linguistics informed multi-task bert,”  
358 in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020.
- 359 [14] R. Caruana, “Multitask learning: A knowledge-based source of inductive bias,” in *Proceedings*  
360 *of the Tenth International Conference on Machine Learning*. Morgan Kaufmann, 1993, pp.  
361 41–48.

- 362 [15] S. R. Bowman, J. Gauthier, A. Rastogi, R. Gupta, C. D. Manning, and C. Potts, “A fast unified  
363 model for parsing and sentence understanding,” in *Proceedings of the 54th Annual Meeting of  
364 the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany:  
365 Association for Computational Linguistics, Aug. 2016, pp. 1466–1477. [Online]. Available:  
366 <https://www.aclweb.org/anthology/P16-1139>
- 367 [16] G. Angeli, M. J. J. Premkumar, and C. D. Manning, “Leveraging linguistic structure for open  
368 domain information extraction,” in *Proceedings of the 53rd Annual Meeting of the Association  
369 for Computational Linguistics and the 7th International Joint Conference on Natural Language  
370 Processing (Volume 1: Long Papers)*, 2015, pp. 344–354.
- 371 [17] K. Toutanova, X. V. Lin, W.-t. Yih, H. Poon, and C. Quirk, “Compositional learning of em-  
372 beddings for relation paths in knowledge base and text,” in *Proceedings of the 54th Annual  
373 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016, pp.  
374 1434–1444.
- 375 [18] R. McDonald, J. Nivre, Y. Quirnbach-Brundage, Y. Goldberg, D. Das, K. Ganchev, K. Hall,  
376 S. Petrov, H. Zhang, O. Täckström, C. Bedini, N. Bertomeu Castelló, and J. Lee, “Universal  
377 Dependency annotation for multilingual parsing,” in *Proceedings of the 51st Annual Meeting  
378 of the Association for Computational Linguistics (Volume 2: Short Papers)*. Sofia, Bulgaria:  
379 Association for Computational Linguistics, Aug. 2013, pp. 92–97.
- 380 [19] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave,  
381 M. Ott, L. Zettlemoyer, and V. Stoyanov, “Unsupervised cross-lingual representation learning  
382 at scale,” in *Proceedings of the 58th Annual Meeting of the Association for Computational  
383 Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp. 8440–8451.  
384 [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.747>
- 385 [20] D. Seddah, S. Kübler, and R. Tsarfaty, “Introducing the SPMRL 2014 shared task on parsing  
386 morphologically-rich languages,” in *Proceedings of the First Joint Workshop on Statistical  
387 Parsing of Morphologically Rich Languages and Syntactic Analysis of Non-Canonical  
388 Languages*. Dublin, Ireland: Dublin City University, Aug. 2014, pp. 103–109. [Online].  
389 Available: <https://www.aclweb.org/anthology/W14-6111>
- 390 [21] M. P. Marcus, B. Santorini, and M. A. Marcinkiewicz, “Building a large annotated corpus of  
391 English: The Penn Treebank,” *Computational Linguistics*, vol. 19, no. 2, pp. 313–330, 1993.  
392 [Online]. Available: <https://www.aclweb.org/anthology/J93-2004>
- 393 [22] N. Xue, F.-D. Chiou, and M. Palmer, “Building a large-scale annotated Chinese corpus,” in  
394 *COLING 2002: The 19th International Conference on Computational Linguistics*, 2002.
- 395 [23] W. Wang, B. Bi, M. Yan, C. Wu, J. Xia, Z. Bao, L. Peng, and L. Si, “Structbert:  
396 Incorporating language structures into pre-training for deep language understanding,”  
397 in *International Conference on Learning Representations*, 2020. [Online]. Available:  
398 <https://openreview.net/forum?id=BJgQ4ISFPH>
- 399 [24] X. Ren, X. Sun, H. Wang, and Q. Liu, “Towards semantics-enhanced pre-training: Can lexicon  
400 definitions help learning sentence meanings?” 2021.
- 401 [25] H. Bai, P. Shi, J. Lin, Y. Xie, L. Tan, K. Xiong, W. Gao, and M. Li, “Segatron: Segment-aware  
402 transformer for language modeling and understanding,” 2020.
- 403 [26] H. Huang, Y. Liang, N. Duan, M. Gong, L. Shou, D. Jiang, and M. Zhou, “Unicoder: A  
404 universal language encoder by pre-training with multiple cross-lingual tasks,” in *Proceedings  
405 of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th  
406 International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong  
407 Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 2485–2494. [Online].  
408 Available: <https://www.aclweb.org/anthology/D19-1252>
- 409 [27] B. Ji, Z. Zhang, X. Duan, M. Zhang, B. Chen, and W. Luo, “Cross-lingual pre-training based  
410 transfer for zero-shot neural machine translation,” in *Proceedings of the AAAI Conference on  
411 Artificial Intelligence*, vol. 34, no. 01, 2020, pp. 115–122.
- 412 [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and  
413 I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing  
414 Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,  
415 and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017. [Online]. Available: <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>  
416

- 417 [29] T. Dozat and C. D. Manning, “Deep biaffine attention for neural dependency parsing,” *ICLR*,  
418 vol. abs/1611.01734, 2017.
- 419 [30] J. Nivre, M.-C. de Marneffe, F. Ginter, J. Hajič, C. D. Manning, S. Pyysalo, S. Schuster,  
420 F. Tyers, and D. Zeman, “Universal Dependencies v2: An evergrowing multilingual treebank  
421 collection,” in *Proceedings of the 12th Language Resources and Evaluation Conference*.  
422 Marseille, France: European Language Resources Association, May 2020, pp. 4034–4043.  
423 [Online]. Available: <https://www.aclweb.org/anthology/2020.lrec-1.497>
- 424 [31] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao,  
425 K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo,  
426 H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick,  
427 O. Vinyals, G. Corrado, M. Hughes, and J. Dean, “Google’s neural machine translation system:  
428 Bridging the gap between human and machine translation,” *CoRR*, vol. abs/1609.08144, 2016.  
429 [Online]. Available: <http://arxiv.org/abs/1609.08144>
- 430 [32] T. Kudo and J. Richardson, “Sentencepiece: A simple and language independent subword  
431 tokenizer and detokenizer for neural text processing,” in *Proceedings of the 2018*  
432 *Conference on Empirical Methods in Natural Language Processing, EMNLP 2018: System*  
433 *Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, E. Blanco and W. Lu,  
434 Eds. Association for Computational Linguistics, 2018, pp. 66–71. [Online]. Available:  
435 <https://doi.org/10.18653/v1/d18-2012>
- 436 [33] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International*  
437 *Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015,*  
438 *Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available:  
439 <http://arxiv.org/abs/1412.6980>
- 440 [34] D. Hendrycks and K. Gimpel, “Bridging nonlinearities and stochastic regularizers with  
441 gaussian error linear units,” *CoRR*, vol. abs/1606.08415, 2016. [Online]. Available:  
442 <http://arxiv.org/abs/1606.08415>
- 443 [35] A. Williams, N. Nangia, and S. Bowman, “A broad-coverage challenge corpus for sentence  
444 understanding through inference,” in *Proceedings of the 2018 Conference of the North American*  
445 *Chapter of the Association for Computational Linguistics: Human Language Technologies,*  
446 *Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics,  
447 Jun. 2018, pp. 1112–1122. [Online]. Available: <https://www.aclweb.org/anthology/N18-1101>
- 448 [36] M. Artetxe, S. Ruder, and D. Yogatama, “On the cross-lingual transferability of monolingual  
449 representations,” in *Proceedings of the 58th Annual Meeting of the Association for*  
450 *Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp.  
451 4623–4637. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.421>
- 452 [37] Y. Zhang, H. Zhou, and Z. Li, “Fast and accurate neural crf constituency parsing,” in  
453 *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-*  
454 *20*, C. Bessiere, Ed. International Joint Conferences on Artificial Intelligence Organization, 7  
455 2020, pp. 4046–4053, main track. [Online]. Available: <https://doi.org/10.24963/ijcai.2020/560>
- 456 [38] Z. Li, H. Zhao, and K. Parnow, “Global greedy dependency parsing,” in *The*  
457 *Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second*  
458 *Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI*  
459 *Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York,*  
460 *NY, USA, February 7-12, 2020*. AAAI Press, 2020, pp. 8319–8326. [Online]. Available:  
461 <https://aaai.org/ojs/index.php/AAAI/article/view/6348>

## 462 Checklist

- 463 1. For all authors...
- 464 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
465 contributions and scope? [Yes]
- 466 (b) Did you describe the limitations of your work? [Yes] As we mentioned in Introduction,  
467 the multilingual PrLMs are hard to have equally high-quality representation for all of  
468 the languages, although our model has improved from the baseline in most languages,  
469 there are still gaps in the performance of our model in different languages.

- 470 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We  
 471 consider our work does not have a potential negative societal impact beyond existed  
 472 works.
- 473 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 474 them? [Yes]
- 475 2. If you are including theoretical results...
- 476 (a) Did you state the full set of assumptions of all theoretical results? [Yes]  
 477 (b) Did you include complete proofs of all theoretical results? [N/A] There is not new  
 478 theories proposed in this paper.
- 479 3. If you ran experiments...
- 480 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 481 mental results (either in the supplemental material or as a URL)? [Yes] Please refer to  
 482 the supplemental material.
- 483 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 484 were chosen)? [Yes] Please refer to the Section 4.1
- 485 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 486 ments multiple times)? [No]
- 487 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 488 of GPUs, internal cluster, or cloud provider)? [No]
- 489 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 490 (a) If your work uses existing assets, did you cite the creators? [Yes]  
 491 (b) Did you mention the license of the assets? [Yes]  
 492 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
 493 (d) Did you discuss whether and how consent was obtained from people whose data you're  
 494 using/curating? [No] We obtain copyrighted data with LDC licence.  
 495 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 496 information or offensive content? [No]
- 497 5. If you used crowdsourcing or conducted research with human subjects...
- 498 (a) Did you include the full text of instructions given to participants and screenshots, if  
 499 applicable? [No]  
 500 (b) Did you describe any potential participant risks, with links to Institutional Review  
 501 Board (IRB) approvals, if applicable? [No]  
 502 (c) Did you include the estimated hourly wage paid to participants and the total amount  
 503 spent on participant compensation? [No]

## 504 A Appendix

### 505 A.1 Model Architecture

### 506 A.2 Extending Strategy for Parsing Treebanks

507 Denote a sentence as  $\hat{X} = [w_1, w_2, \dots, w_N]$  and its corresponding subword-level sentence as  $X =$   
 508  $[[CLS], w_{11}, \dots, w_{1s_1}, w_{21}, \dots, w_{2s_2}, \dots, w_{N1}, \dots, w_{Ns_N}, [SEP]]$ . In the word-level parsing tree, for  
 509  $w_i \in \hat{X}$ , suppose the head of  $w_i$  is  $w_{h_i}$  and the dependency label is  $l_i$ . We design three strategies to  
 510 extend the parsing tree from word level to subword level. All of them regard [CLS] as the ROOT  
 511 ( $w_{01}$ ), and the header of [SEP] as [CLS], and the label as the padding label "\_".

- 512 • Set the head of all the subwords of  $w_i$  as  $w_{h_{i1}}$ , and the label as  $l_i$ .
- 513 • Set the head of  $w_{i1}$  as  $w_{h_{i1}}$  and the label as  $l_i$ , and set the head of  $w_{i2}, \dots, w_{is_i}$  as  $w_{i1}$  and  
 514 label as a special label "APP".
- 515 • Set the head of  $w_{i1}$  as  $w_{h_{i1}}$  and the label as  $l_i$ , and set the head of  $w_{ij}$  as  $w_{i(j-1)}$  for  
 516  $i = 2, \dots, s_i$  and label as a special label "APP". As shown in Figure 2.

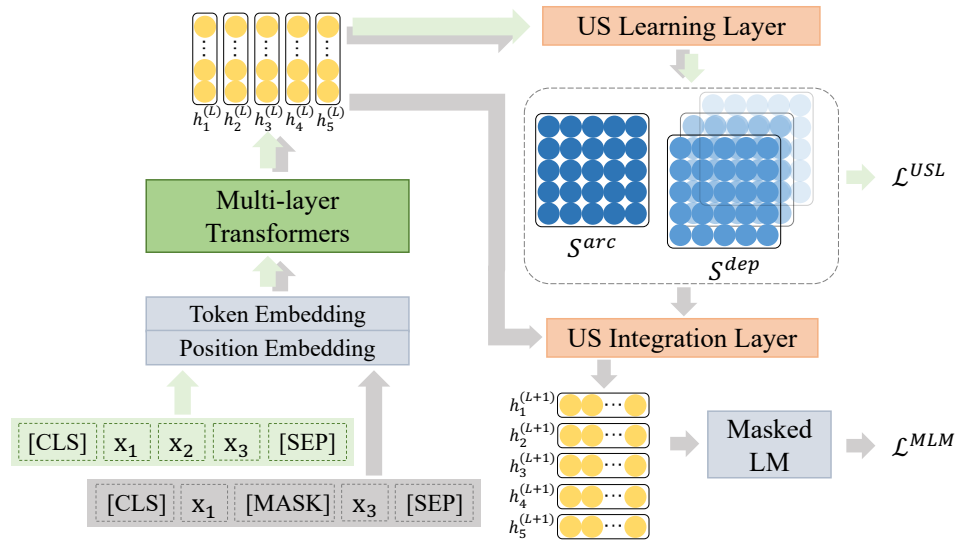


Figure 1: The model architecture for UD-BERT.

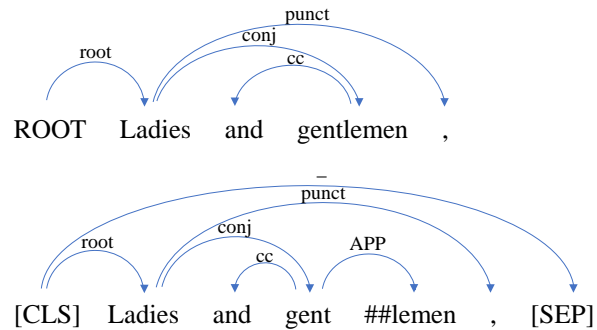


Figure 2: An example of extending a word-level parsing tree to a subword-level.

Table 8: Details of the selected languages in UD.

Language	Treebank	Sents
Bulgarian (bg)	BTB	8,907
Catalan (ca)	AnCora	13,123
Czech (cs)	PDT	102,993
Dutch (nl)	Alpino	18,058
English (en)	EWT	12,543
Estonian (et)	EDT	20,827
Finnish (fi)	TDT	12,217
French (fr)	GSD	14,554
German (de)	GSD	13,814
Hebrew (he)	HTB	5,241
Hindi (hi)	HDTB	13,304
Indonesian (id)	GSD	4,477
Italian (it)	ISDT	13,121
Korean (ko)	GSD	27,410
Latin (la)	PROIEL	15,906
Latvian (lv)	LVTB	5,424
Norwegian (no)	Bokmaal	29870
Polish (pl)	LFG	19,874
Romanian (ro)	RRT	8,043
Russian (ru)	SynTagRus	48,814
Slovak (sk)	SNK	8,483
Spanish (es)	AnCora	28,492