

# Representation Alignment and Adversarial Networks for Cross-lingual Dependency Parsing

Anonymous ACL submission

## Abstract

Thanks to the strong representation capability of pre-trained language models, dependency parsing in rich-resource language has achieved remarkable improvements. However, the parsing accuracy drops sharply when the model is transferred to low-resource language due to distribution shifts. To alleviate this issue, we propose a representation alignment and adversarial model to filter out useful knowledge from rich-resource language and ignore useless ones. Our proposed model consists of two components, i.e., an alignment network in the input layer for selecting useful language-specific representation features and an adversarial network in the encoder layer for augmenting the language-invariant contextualized features. Experiments on the benchmark datasets show that our proposed model outperforms RoBARTa-enhanced strong baseline models by 1.37 LAS and 1.34 UAS. Detailed analysis shows that both alignment and adversarial networks are equally important in alleviating the distribution shifts problem and can benefit from each other. In addition, the comparative experiments demonstrate that both the alignment and adversarial networks can substantially facilitate extracting and utilizing relevant target language features, thereby increasing the adaptation capability of our proposed model.

## 1 Introduction

Dependency parsing, as an important fundamental task of natural language processing, aims to identify grammatical and syntax relationships between two words in the input sentence via a dependency tree. Figure 1 shows a dependency tree instance, where a dependency from the headword “voi (elephant)” to the modified word “thông minh (intelligent)” with the relation label “amod” means “thông minh (intelligent)” as an adjective modifies “voi (elephant)”. Dependency trees are widely applied

to various artificial intelligence tasks, such as machine translation (Zhang et al., 2019), grammatical error correction (Zhang et al., 2022), and information extraction (Tian et al., 2022).

In the past decades, pre-trained language model enhanced dependency parsers have achieved outstanding performances in rich-resource languages (Clark et al., 2018; Li et al., 2022; Nishida and Matsumoto, 2022; Mohammadshahi and Henderson, 2021; Yan et al., 2020). Most significantly, Dozat and Manning (2017) propose a BiAffine parser that leverages multi-layer BiLSTMs to encode input sentences and a BiAffine operation to compute scores, thus achieving better performance on various languages. Then, Li et al. (2019) develop a self-attentive BiAffine parser and further improve the model performance with ELMo and BERT representations. However, these model performances drop sharply in low-resource languages due to the lack of annotated data (Wang et al., 2020; Effland and Collins, 2023; Rotman and Reichart, 2019; Vania et al., 2019).

As shown in Figure 1, both sentences from Vietnamese and Chinese have a similar core grammatical structure “subject-predicate-object”, but they also have differences in the attributive positions where Vietnamese adopts “post-modifier” while Chinese is the opposite. Hence, how to construct the discrepancy and similarity between different languages becomes the key challenge for cross-lingual dependency parsing (Ahmad et al., 2019; Üstün et al., 2022; Ozaki et al., 2021; Liu et al., 2020; Xu and Koehn, 2021). A series of previous works have explored feature transfer to improve low-resource parsing. Most recently, Al Ghiffari et al. (2023) propose a hierarchical transfer learning (HTL) approach to exploit a source and an intermediate language to improve the parsing accuracy in low-resource languages. Similarly, Choudhary and O’riordan (2023) incorporate linguistic typology knowledge as an auxiliary task, further im-

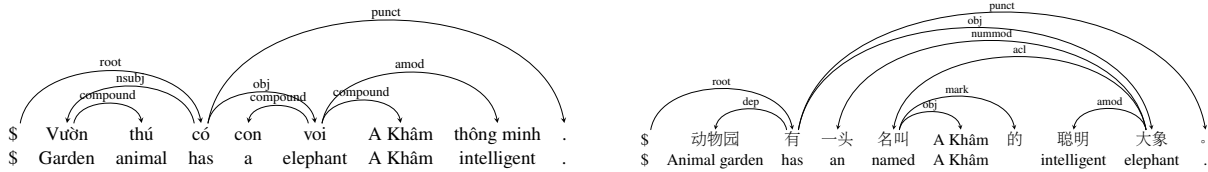


Figure 1: Examples of dependency tree from Universal Dependencies (UD) dataset, where the left sentence is from the low-resource Vietnamese treebanks (VTB) and the right one is from the rich-resource simplified Chinese treebanks (GSDSimp).

082 proving the low-resource dependency parsing per- 123  
 083 formances. Although transfer learning from rich- 124  
 084 resource to low-resource language has shown its 125  
 085 promising advantages, how to further emphasize 126  
 086 the helpful knowledge and filter out the harmful 127  
 087 ones automatically is still an important problem. 128

088 To address this issue, we propose a novel repre- 129  
 089 sentation alignment and adversarial networks for 130  
 090 cross-lingual dependency parsing. On the one hand, 131  
 091 we propose an alignment network on the input 132  
 092 layer to select useful language-specific word in- 133  
 093 formation. On the other hand, a language-aware 134  
 094 adversarial network is applied on the encoder layer 135  
 095 to excavate potential language-invariant knowl- 136  
 096 edge. Experiments on the benchmark dataset show 137  
 097 that our proposed model achieves notable per- 138  
 098 formance improvements, leading to new state-of- 139  
 099 the-art results. Detailed analysis shows that align- 140  
 100 ment and adversarial networks are complementary 141  
 101 and can benefit from each other. In-depth com- 142  
 102 parative experiments demonstrate that both align- 143  
 103 ment and adversarial networks are equally impor- 144  
 104 tant for filtering out effective knowledge from the 145  
 105 source language. In addition, our codes are released 146  
 106 at <https://github.com/noteljj/align> to facilitate 147  
 107 future research. 148

## 108 2 Related Work 149

109 **Cross-Lingual Dependency Parsing.** Cross- 150  
 110 lingual dependency parsing has emerged as a cru- 151  
 111 cial component of natural language processing, 152  
 112 with distinct methodologies contributing to its ad- 153  
 113 vancement. Among these, three primary categories 154  
 114 stand out: *transfer learning*, *multilingual model* 155  
 115 *adaptation*, and *subword representation alignment*. 156  
 116 Transfer learning techniques, epitomized by the 157  
 117 work of Chen et al. (2019), Liu et al. (2023b) and 158  
 118 Niu et al. (2022), leverage resources from rich- 159  
 119 resource languages to improve parsing accuracy 160  
 120 in low-resource languages, demonstrating the ver- 161  
 121 satility of transferring syntactic knowledge across 162  
 122 linguistic boundaries. In multilingual model adap- 163  
 164

123 tion, researchers like Pfeiffer et al. (2021). Wang 124  
 125 et al. (2020) and Dione (2021) have adapted mul- 126  
 127 tilingual BERT models to enhance parsing perfor- 127  
 128 mance across various languages, illustrating the 128  
 129 power of transformer-based methods in handling 129  
 130 diverse linguistic environments. Meanwhile, the 130  
 131 subword representation alignment approach, as ex- 131  
 132 plored by Schuster et al. (2019); Yaari et al. (2022), 132  
 133 focuses on the fine-grained alignment of word or 133  
 134 subword representations between languages, ad- 134  
 135 dressing the challenge of representing low-resource 135  
 136 languages in pre-trained models. Collectively, these 136  
 137 approaches underscore the dynamism and complex- 137  
 138 ity of cross-lingual dependency parsing, highlight- 138  
 139 ing both its progress and the ongoing challenges 139  
 140 of syntactic alignment and resource disparity. This 140  
 141 landscape sets the stage for our investigation into 141  
 142 the effective transfer of subword representations 142  
 143 from Chinese to Vietnamese, a venture that seeks 143  
 144 to mitigate the representation gap for low-resource 144  
 145 languages and contribute to the evolving narrative 145  
 146 of linguistic adaptability in computational models. 146  
 147

148 **Adversarial Learning.** Adversarial learning has 149  
 150 become increasingly central in NLP, notably for its 150  
 151 role in fortifying model robustness and counteract- 151  
 152 ing data biases (Lowd and Meek, 2005), Zalmout 152  
 153 and Habash (2019) and Chen et al. (2021) have 153  
 154 demonstrated the efficacy of adversarial examples 154  
 155 in bolstering the resilience of NLP models to lin- 155  
 156 guistic variations and malicious attacks. Extending 156  
 157 this, Lu et al. (2023) and Zou et al. (2021) have 157  
 158 successfully integrated adversarial learning into 158  
 159 domain adaptation, effectively reducing domain- 159  
 160 specific biases. A recent novel approach by Han 160  
 161 et al. (2021) and Zhang et al. (2018) involves using 161  
 162 adversarial training to mitigate biases in training. 162  
 163 Additionally, the advent of adversarial data aug- 163  
 164 mentation, as investigated by Tan et al. (2022), has 164  
 165 shown promise in diversifying training datasets, 165  
 166 further enhancing model robustness. Despite these 166  
 167 advancements, adversarial learning still confronts 167  
 168 challenges in balancing model stability and per- 168

165 formance, particularly when dealing with highly  
 166 complex and nuanced linguistic data, underscoring  
 167 the need for ongoing research and development in  
 168 this dynamic area of NLP.

169 **Feature Alignment and Transfer.** In the field  
 170 of feature alignment and transfer, existing research  
 171 can be categorized into *deep learning-based meth-*  
 172 *ods*, *instance-based methods*, and *model-based*  
 173 *methods*. Deep learning-based methods automati-  
 174 cally learn feature mapping relationships between  
 175 source and target domains through neural networks,  
 176 such as aligning feature distributions in the space  
 177 through adversarial training (Riemer et al., 2015),  
 178 (Kumar et al., 2023) and (Hazem et al., 2022).  
 179 Instance-based methods select and weight exam-  
 180 ples from the source domain to have a greater im-  
 181 pact in the target domain, like instance selection  
 182 based on conditional adversarial learning (Basu  
 183 Roy Chowdhury et al., 2019; Glavaš and Vulić,  
 184 2020). Model-based methods focus on how to use  
 185 the source domain’s model to assist learning in the  
 186 target domain, such as progressive neural networks  
 187 that learn to transfer knowledge across domains  
 188 (Chawla and Yang, 2020; Liu et al., 2023a). These  
 189 methods have their own advantages and can effec-  
 190 tively improve the performance of cross-domain  
 191 learning in different scenarios.

### 192 3 Our Approach

193 Considering not all rich-source language informa-  
 194 tion is equally important for cross-lingual depen-  
 195 dency parsing, we propose the alignment and adver-  
 196 sarial networks for effective representation selec-  
 197 tion. Concretely, we first leverage the multi-lingual  
 198 pre-trained language model XLM-RoBERTa to im-  
 199 prove the word representation capability of both  
 200 source and target languages. Then, a representa-  
 201 tion alignment network is applied on the input layer  
 202 to emphasize useful language-specific information  
 203 and ignore the harmful one. Next, we exploit an  
 204 adversarial network on the encoder layer to en-  
 205 hance language-invariant representations. Finally,  
 206 all selected representations are utilized to search  
 207 for the best dependency tree. Figure 2 illustrates the  
 208 framework of our proposed model, which is orga-  
 209 nized into three components, i.e., *Input layer based*  
 210 *on the alignment network*, *Encoder layer enhanced*  
 211 *with an adversarial network*, *MLP and BiAffine*  
 212 *layers*.

### 213 3.1 Input Layer Based on Representation 214 Alignment Network

215 Given an input sentence  $w_1, w_2, \dots, w_n$ , the input  
 216 layer maps them into dense vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ .  
 217 For the source language Chinese, we directly use  
 218 the normal embeddings as its input vectors. For  
 219 the target language Vietnamese, we exploit a repre-  
 220 sentation alignment network to select helpful Chi-  
 221 nese word information, further enhancing the Viet-  
 222 namese representation capability.

223 **Input vectors for Chinese.** As shown in Equa-  
 224 tion 1, each Chinese vector  $\mathbf{x}_i^{ch}$  is the concatenation  
 225 of its word representation and corresponding char-  
 226 acter representation  $\mathbf{word}_i^{char}$ , where word repre-  
 227 sentation is the addition of XLM-RoBERTa rep-  
 228 resentation  $\mathbf{rep}_i^{XLM-R}$  and a random initialization  
 229 word embedding  $\mathbf{emb}_i^{word}$ . The character repre-  
 230 sentation  $\mathbf{word}_i^{char}$  is generated by a BiLSTM net-  
 231 work, which first encodes the constituent characters  
 232 of each word  $w_i^{ch}$ , and then combines the hidden  
 233 vectors of two directions (Lample et al., 2016).

$$234 \mathbf{x}_i^{ch} = (\mathbf{rep}_i^{XLM-R} + \mathbf{emb}_i^{word}) \oplus \mathbf{word}_i^{char} \quad (1)$$

235 **Input vectors for Vietnamese.** Different from  
 236 Chinese input vectors, Vietnamese input vector  
 237  $\mathbf{x}_i^{vi}$  utilizes an additional aligned representation  
 238  $\mathbf{emb}_i^{vi-FT}$  to fuse more useful Chinese word infor-  
 239 mation, which is calculated in Equation 2,

$$240 \mathbf{x}_i^{vi} = (\mathbf{emb}_i^{vi-FT} + \mathbf{rep}_i^{XLM-R} + \mathbf{emb}_i^{word}) \oplus \mathbf{word}_i^{char} \quad (2)$$

241 where  $\mathbf{emb}_i^{vi-FT}$  is generated by our alignment  
 242 network and other representations are obtained sim-  
 243 ilarly to Chinese.

244 **Alignment network.** The key to our alignment  
 245 network is to enhance the Vietnamese word repre-  
 246 sentation capability by emphasizing useful Chinese  
 247 words and ignoring harmful ones. First, we con-  
 248 struct an alignment matrix based on a new high-  
 249 quality bilingual dictionary to map Vietnamese and  
 250 Chinese representations into a close space.

251 Since the bilingual dictionary significantly af-  
 252 fects the performance of our alignment matrix, we  
 253 adopt automatic generation and manual annotation  
 254 strategy to ensure the quality of the Vietnamese-  
 255 Chinese dictionary. Concretely, we first download  
 256 the dump data backup file from Wikipedia<sup>1</sup> and a

<sup>1</sup><https://en.wikipedia.org/>

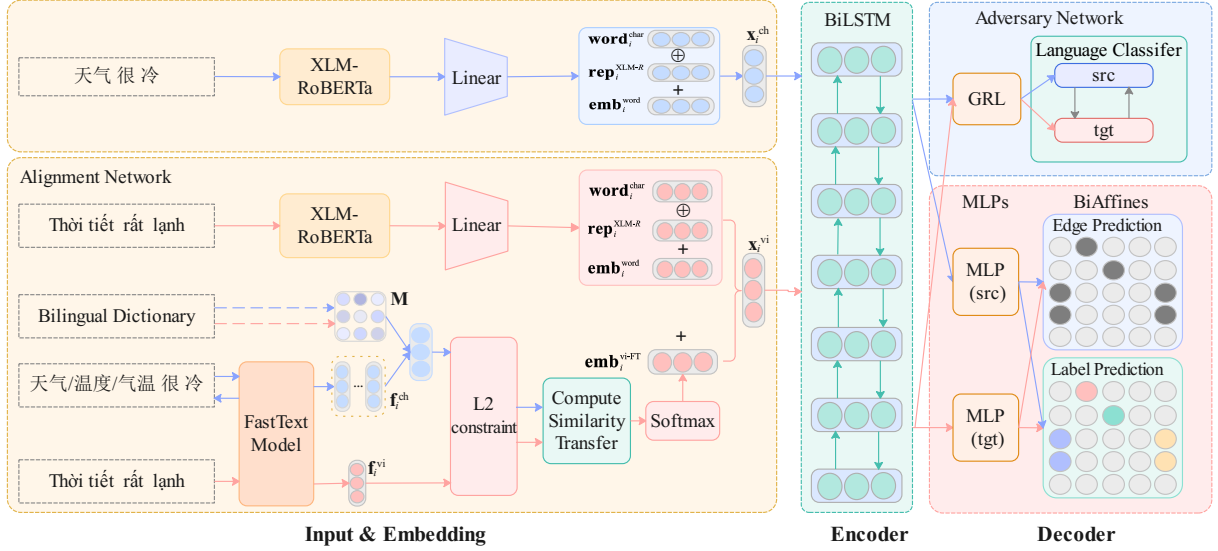


Figure 2: Framework of our proposed model.

simple bilingual dictionary<sup>2</sup>. Second, we use regular expressions to iteratively match and extract the Vietnamese-Chinese alignment titles and subheadings. Third, the alignment word pairs are used to augment the original bilingual dictionary. Finally, the automatic generation dictionary is manually proofread by Vietnamese speakers, thus obtaining a high-quality Vietnamese-Chinese dictionary that contains about 20,000 word pairs. based on the new dictionary, we use the pre-trained Fasttext models<sup>3</sup> to obtain Vietnamese matrix  $V \in \mathcal{R}^{n \times d_1}$  and Chinese matrix  $C \in \mathcal{R}^{n \times d_1}$  where  $n$  is the number of our dictionary and  $d_1$  denotes the dimension of Fasttext representations. Meanwhile, we exploit an orthogonal similarity transformation to obtain our alignment matrix  $M \in \mathcal{R}^{d_1 \times d_1}$  that can be regarded as a linear mapping between Vietnamese and Chinese based on the semantic similarity.

Given a Vietnamese sentence, we first utilize Fasttext models to obtain word segmentation sequences. Then, for each Vietnamese word, we select multiple corresponding Chinese words based on our dictionary. Next, all selected words are dotted with an alignment matrix  $M$ , and L2 constraint is applied on them to yield stable and aligned word representations  $\hat{\mathbf{f}}_i$ . The formula for this operation is as follows,

$$\hat{\mathbf{f}}_i = \frac{\mathbf{f}_i}{\sqrt{\sum_{i=1}^n \mathbf{f}_i^2 + \varepsilon}} \quad (3)$$

where  $\mathbf{f}_i$  represents the  $i$ -th word vector from the FastText model,  $\varepsilon$  is a very small positive number used to prevent division by zero. Considering each Vietnamese word may align with several Chinese words, we employ the cosine function to compute semantic similarity as alignment weights. The formulas are shown as follows,

$$S_{i,j}^{ch,vi} = \frac{(\hat{\mathbf{f}}_i^{ch})^T \hat{\mathbf{f}}_j^{vi}}{\|\hat{\mathbf{f}}_i^{ch}\| \|\hat{\mathbf{f}}_j^{vi}\|} \quad (4)$$

$$\mathbf{w}_{i,j}^{ch,vi} = \exp(S_{i,j}^{ch,vi} / \tau)$$

where  $S_{i,j}^{ch,vi}$  denotes the similarity score between the Chinese word  $i$  and the Vietnamese word  $j$ .  $\tau$  denotes the temperature coefficient.  $\mathbf{w}_{i,j}^{ch,vi}$  is the corresponding weight. Finally, We construct the final alignment Vietnamese representation  $\mathbf{emb}_i^{vi-FT}$  using constrained word vectors and alignment weights to emphasize useful words and ignore harmful ones. The formula is as follows,

$$\mathbf{emb}_i^{vi-FT} = \frac{\sum_{ch \in \mathcal{J}_{vi}} \mathbf{w}_{i,j}^{ch,vi} \cdot \hat{\mathbf{f}}_i^{ch}}{\sum_{ch \in \mathcal{J}_{vi}} \mathbf{w}_{i,j}^{ch,vi}} \quad (5)$$

where  $\mathcal{J}_{vi}$  represents a collection of Chinese words that exhibit the highest degree of similarity to a Vietnamese word.

<sup>2</sup><https://github.com/CPJKU/wechsel/tree/main/dicts/data>

<sup>3</sup><https://fasttext.cc/docs/en/crawl-vectors.html/>



### 3.2 Encoder Layer Enhanced with Adversarial Network

Different from the traditional BiLSTM encoder, we employ an adversarial network above the encoder to ensure it imply more potential language-invariant knowledge.

**BiLSTM encoder.** Following Dozat and Manning (2017), we also adopt a three-layer BiLSTM network as the encoder to generate original contextualized vectors. Since BiLSTM is able to encode the words in a sentence from two directions, each word can obtain contextualized information  $\mathbf{h}_i$ .

$$\mathbf{h}_i = \text{BiLSTM}(\mathbf{x}_i, \theta_{\text{BiLSTM}}) \quad (6)$$

where  $\theta_{\text{BiLSTM}}$  is the BiLSTM parameters.

**Adversarial network.** The adversarial network mainly contains three components, i.e., the shared BiLSTM encoder, the Gradient Reversal Layer (GRL), and a language classifier. First, Sentence from Chinese or Vietnamese are fed into the shared BiLSTM layer to obtain contextualized word representations  $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n$ . Then, they pass the GRL which inverts the gradient during backpropagation, thus fostering BiLSTM to learn more shared features between Vietnamese and Chinese. The forward and backward propagation equations for GRL are as follows,

$$\begin{aligned} \text{GRL}_\gamma(\mathbf{h}_i) &= \mathbf{h}_i \\ \frac{d\text{GRL}_\gamma(\mathbf{h}_i)}{d(\mathbf{h}_i)} &= -\gamma \mathbf{I} \end{aligned} \quad (7)$$

where  $\gamma$  is a hyperparameter to balance the impact of adversarial learning and dependency parsing on the shared BiLSTM. Then, we use a multilayer perceptron (MLP) to compute the language distribution scores and a softmax function to obtain the language distribution probabilities. The formula is as follows,

$$\mathbf{re}_i = \text{softmax}(\text{MLP}(\mathbf{h}_i)) \quad (8)$$

Finally, we employ a standard cross-entropy loss to optimize all parameters of the adversarial network,

$$\mathcal{L}^{adv} = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m (\tilde{\mathbf{r}}_{i,j}) \log((\mathbf{re}_{i,j})) \quad (9)$$

where  $m$  is the number of languages,  $n$  is the word number of input sentence, and  $\tilde{\mathbf{r}}_{i,j}$  represents the gold-standard language distribution vector, where only one element is 1 corresponding to the language index where the sentence comes from.

### 3.3 MLP and BiAffine Layer

The MLP layer employs the enhanced contextualized vector  $\mathbf{h}_i$  as its input and reduce the dimension of  $\mathbf{h}_i$ , extracting its head representation  $\mathbf{r}_i^h$  and modifier representation  $\mathbf{r}_i^d$  for each word  $w_i$ .

$$\begin{aligned} \mathbf{r}_i^h &= \text{MLP}_h(\mathbf{h}_i) \\ \mathbf{r}_i^d &= \text{MLP}_d(\mathbf{h}_i) \end{aligned} \quad (10)$$

where  $\text{MLP}_h(*)$  and  $\text{MLP}_d(*)$  have a single hidden layer with the ReLU activation function. Then, a BiAffine computes  $\text{score}(i \leftarrow j)$  between the current word  $w_i$  and the other word  $w_j$ . Simultaneously,  $\text{score}(i \leftarrow^l j)$  is calculated by another separated BiAffine layer as equation 11

$$\begin{aligned} \text{score}(i \leftarrow j) &= \begin{bmatrix} \mathbf{r}_i^d \\ 1 \end{bmatrix}^T \mathbf{U}_1 \mathbf{r}_j^h \\ \text{score}(i \leftarrow^l j) &= \mathbf{r}_j^h \mathbf{U}_2 \mathbf{r}_i^d + (\mathbf{r}_j^h \oplus \mathbf{r}_i^d) \mathbf{U}_3 + b \end{aligned} \quad (11)$$

where  $\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3$ , and  $b$  are parameters.  $l$  denotes the relation label. After obtaining the scores of dependency arcs and dependency labels, we use the typical Maximum Spanning Tree (MST) algorithm to find the highest-score tree as our final parsing result. Finally, for each position  $i$ , if the gold-standard head of word  $w_i$  is word  $w_j$  and its corresponding gold relation label is  $l$ , the parsing loss is computed as follows,

$$\begin{aligned} \mathcal{L}^{\text{par}} &= -\log \frac{e^{\text{score}(i \leftarrow j)}}{\sum_{0 \leq k \leq n, k \neq i} e^{\text{score}(i \leftarrow k)}} \\ &\quad - \log \frac{e^{\text{score}(i \leftarrow^l j)}}{\sum_{l' \in \mathcal{L}} e^{\text{score}(i \leftarrow^{l'} j)}} \end{aligned} \quad (12)$$

where  $\text{score}(i \leftarrow k)$  denotes the dependency arc score from head  $w_i$  to modifier  $w_k$ .  $\mathcal{L}$  refers to the collection of all dependency labels  $l'$ .

### 3.4 Cyclic Cross-lingual Training

In this work, we propose a cyclic training strategy to mitigate data imbalance between source and target languages, as outlined in Algorithm 1. Considering the data scale of the source language is much larger than the target one, we divide the first  $n_1$  mini-batches of the source language as  $s^f$  and the last as  $s^l$  where  $n_1$  is the mini-batch number of the target language. During training, we take turns

---

**Algorithm 1:** Cyclic Training Procedure

---

**Input:** Source language data  $S$ , target language data  $T$   
**Hyper-parameters:** Loss weight  $\alpha$ , training iterations  $k$   
1: Initialize  $iter = 0$   
2: **Repeat**  
3: Sample mini-batch  $x$  alternately from  $S$  or  $T$   
4: **if**  $x \in S^f$ :  
5: Update parameter by minimizing  $\mathcal{L}^{par} + \alpha\mathcal{L}^{adv}$   
6: **elif**  $x \in S^l$ :  
7: Update parameter by minimizing  $\mathcal{L}^{par}$   
8: **else**  $x \in T$ :  
9: Compute  $\mathbf{emb}_i^{vi-FT} = \text{alignment}(\theta_s)$   
11: Update parameters by minimizing  $\mathcal{L}^{par} + \alpha\mathcal{L}^{adv}$   
12:  $iter + = 1$   
13: **until**  $iter = k$  or convergence

---

Table 1: Cyclic Cross-lingual Training Procedure.

384 to sample mini-batch  $x$  of source and target lan-  
385 guages. If  $x$  comes from the first part of the source  
386 language  $S^f$ , we update parsing and adversarial  
387 parameters by minimizing parsing and adversarial  
388 losses. While  $x$  belongs to  $S^l$ , we only update the  
389 parser parameters  $\theta_1$  by minimizing the parsing  
390 loss. If  $x$  comes from the target language  $T$ , we  
391 compute an alignment representation  $\mathbf{emb}_i^{vi-FT}$   
392 via an alignment network. and update all parame-  
393 ters by minimizing parsing and adversarial losses.  
394 Finally, we iteratively train all the data until it con-  
395 verges or stops prematurely.

| Dataset | Train | Dev | Test |
|---------|-------|-----|------|
| GSDSimp | 3,997 | 500 | 500  |
| VTB     | 1,400 | 800 | 800  |

Table 2: Dataset statistics in sentence number.

## 396 4 Experiments

### 397 4.1 Settings

398 **Datasets.** To compare with previous work fairly,  
399 we use the shared multi-language Universal De-  
400 pendencies (UD) 2.12 treebank as our benchmark  
401 datasets<sup>4</sup>. Concretely, we choose Chinese as our  
402 source language and Vietnamese as our target lan-  
403 guage. The detailed illustrations of our datasets are  
404 shown in Table 2.

405 **Evaluation.** Following Hajic et al. (2009), we  
406 employ the Labeled Attachment Score (LAS) and  
407 Unlabeled Attachment Score (UAS) as our evalua-  
408 tion indicators. Each model is trained for at most  
409 1, 000 iterations, and the performance is evaluated  
410 on the dev data after each iteration for model selec-

<sup>4</sup><https://universaldependencies.org/>

411 tion. We stop the training if the peak performance  
412 does not increase in 100 consecutive iterations.

413 **Hyper-parameter choices.** We mostly maintain  
414 the hyper-parameter settings of Li et al. (2019),  
415 such as MLP and BiAffine dimensions, dropout  
416 ratios, and so on. The adversary loss weight  $\alpha$ ,  
417 neighbor, and temperature, which are set as 1, 10,  
418 and 0.1 respectively. The character embeddings are  
419 initialized randomly with a dimension of 100.

420 **Baseline.** To validate the advantages and effec-  
421 tiveness of our proposed model, we choose the  
422 following approaches as our strong baselines.

- 423 • **Pre-training method.** BiAffine parser is first  
424 proposed by Dozat and Manning (2017), and  
425 then is widely used on various dependency  
426 parsing tasks. Different from the original Bi-  
427 Affine parser, we first exploit the Vietnamese  
428 pre-trained language model XLM-RoBERTa-  
429 base<sup>5</sup> to enhance the parsing performance.  
430 Then, we pre-train the enhanced BiAffine  
431 parser exclusively on the Vietnamese Univer-  
432 sal Dependencies (UD) dataset, which is used  
433 as our strong baseline model.
- 434 • **Fine-tuning method.** Shi et al. (2022) pro-  
435 pose to fine-tune the basic model twice and  
436 achieve selective differential privacy for large  
437 language models. In this work, we also uti-  
438 lize the idea of fine-tuning method to improve  
439 the adaptation capability of the enhanced Bi-  
440 Affine parser in Vietnamese. We first use the  
441 Chinese dataset for initial training, and then  
442 fine-tune the pre-trained model with the Viet-  
443 namese dataset, thus transferring the syntactic  
444 knowledge contained in the Chinese treebank  
445 to Vietnamese.
- 446 • **Adversarial learning method.** Li et al.  
447 (2021) apply the adversarial network on the  
448 BiAffine parser, thus achieving impressive re-  
449 sults on cross-domain dependency parsing. In  
450 this work, we attempt to apply an adversarial  
451 network on BiAffine parser to capture more  
452 similarities between Chinese and Vietnamese.

### 453 4.2 Main Results

454 Table 3 displays the final results on our test data and  
455 gives a detailed comparison with previous works.  
456 First, we find that our model outperforms the “Ad-  
457 versary” model, demonstrating that our alignment

<sup>5</sup><https://huggingface.co/xlm-roberta-base>

| Model                     | LAS          | UAS          |
|---------------------------|--------------|--------------|
| Results of previous works |              |              |
| UDPipe (2019)             | 62.56        | 70.38        |
| UDify(2019)               | 66.00        | 74.11        |
| UDPipe2.0+WCBF(2019)      | 65.41        | 72.94        |
| TOWER (2021a)             | 63.50        | 72.40        |
| Pre-training              | 67.61        | 75.47        |
| Fine-tuning               | 68.09        | 75.93        |
| Adversary                 | 68.47        | 76.39        |
| Our model                 | <b>68.98</b> | <b>76.81</b> |

Table 3: Main results on the Vietnamese UD test dataset. network can emphasize useful language-specific representations from the source language and ignore the harmful ones, thus further improving the cross-lingual dependency parsing accuracy. Second, compared with the “Fine-tuning” model, the “Adversary” model achieves better performance, revealing that an adversarial network can extract potential language-invariant knowledge to construct the in-depth relationship between source and target languages. Finally, we can see that our proposed model outperforms all strong baselines, indicating that our proposed representation alignment and adversarial networks are extremely useful for cross-language dependency parsing.

We also compare with previous works in the top block. [Kondratyuk and Straka \(2019\)](#) first propose the UDpipe model, which integrates a tokenizer, morphological analyzer, POS tagger, lemmatizer, and dependency parser into a single model for comprehensive natural language processing. Then, they propose a UDify framework based on a multilingual BERT self-attention model with tagging and parser joint training, which fine-tunes a multilingual pre-trained model with 104 languages to improve parsing accuracy. [Straka et al. \(2019\)](#) enhance the UDpipe model by incorporating various embeddings, including BERT and Flair. Lastly, [Glavaš and Vulić \(2021b\)](#) propose a TOWER model, which uses hierarchical language clustering to improve the low-resource dependency parsing performance. Compared with these works, we find that our model can achieve the best performance with only a single target language, highlighting the efficiency and powerful parsing capabilities of our proposed model.

### 4.3 Ablation Study

Results of ablation studies are shown in Table 4. First, we find that removing either the adversarial

network or the representation alignment network can lead to a decrement in parsing performance. This outcome suggests that each module plays a crucial role in mitigating the potential conflicts arising from direct language transfer. Second, removing adversarial and alignment modules simultaneously leads to a significant decline in dependency parsing accuracy, revealing that the two modules are complementarity and benefit from each other. Most notably, the performance deteriorates to its lowest when the source language is excluded altogether, affirming that the source language encompasses valuable information beneficial for the target language. This observation not only emphasizes the importance of preserving source language features but also reinforces the necessity of their strategic filtration.

| Model              | LAS          | UAS          |
|--------------------|--------------|--------------|
| Our model          | <b>68.98</b> | <b>76.81</b> |
| w/o Adv            | 68.71        | 76.53        |
| w/o Ali            | 68.47        | 76.39        |
| w/o Adv & Ali      | 68.09        | 75.93        |
| w/o Adv & Ali & Ch | 67.61        | 75.47        |

Table 4: Ablation study on reducing the component of our model on test data, where “w/o Adv”, “w/o Ali”, and “w/o Ch” mean removing the adversarial network, representation alignment network or the Chinese UD training dataset.

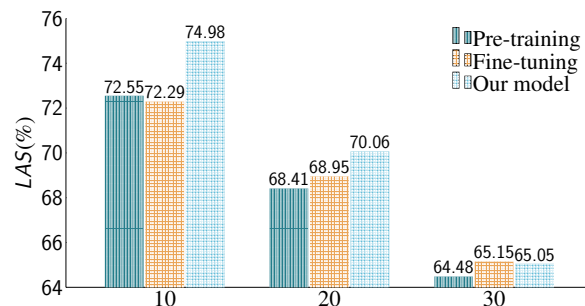


Figure 3: LAS regarding diverse sentence lengths.

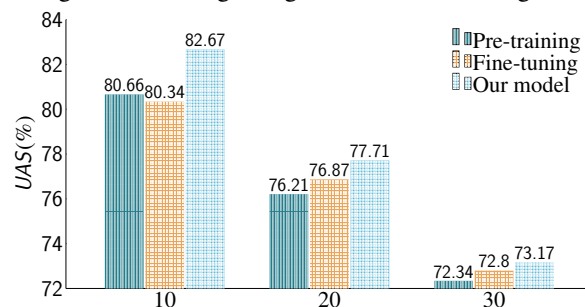


Figure 4: UAS regarding diverse sentence lengths.

#### 4.4 Error Analysis

**Sentence length.** Figure 3 and Figure 4 present the LAS and UAS scores regarding diverse sentence lengths. First, it is clear that all models perform better with shorter sentences. For sentences under 10 words, the LAS and UAS scores hover around 73 and 82, respectively. However, there is a noticeable drop of over 9 points in scores for sentences approximately 30 words in length, indicating that the parsing difficulty is sharply improved with the increase in sentence length. Then, we can see that the “Pre-training” model records the lowest scores across all length categories. Notably, incorporating the Chinese corpus enhances its performance across most lengths, except for the 10-word category. The reason may be that pronounced structural disparities between short Chinese and Vietnamese sentences. Finally, our model significantly mitigates the performance decline observed with the “Fine-tuning” model, achieving substantial improvements across all sentence lengths.

| DEP       | Precision (%) |              |              |
|-----------|---------------|--------------|--------------|
|           | Pre-training  | Fine-tuning  | Our          |
| amod      | 67.45         | 63.78        | <b>67.97</b> |
| cc        | 87.34         | 86.74        | <b>88.64</b> |
| ccomp     | 54.33         | 54.64        | <b>56.45</b> |
| compound  | 73.03         | 73.47        | <b>74.75</b> |
| conj      | 63.69         | 64.50        | <b>66.60</b> |
| cop       | 81.35         | 81.94        | <b>82.05</b> |
| discourse | 44.12         | <b>53.57</b> | 52.78        |
| mark      | 73.00         | 73.33        | <b>73.58</b> |
| nmod      | 70.84         | 71.99        | <b>73.12</b> |
| nsubj     | 83.42         | 83.47        | <b>83.85</b> |
| obj       | 79.86         | 81.17        | <b>81.67</b> |
| root      | 79.64         | 79.71        | <b>80.14</b> |

Table 5: Precisions of dependency labels on different models.

**Dependency labels.** Table 5 presents the precisions of main dependency labels on different models. These models include the Chinese training dataset to analyze inter-language connections. First, the “Pre-training” model registers the lowest scores across all dependency labels. Then, the “Fine-tuning” model achieves better performance on most dependency labels. The reason may be that the dependency trees in the target language contain abundant language-specific syntax information. Finally, our proposed model consistently obtains the highest scores on almost all labels, further proving

the effectiveness of our proposed model.

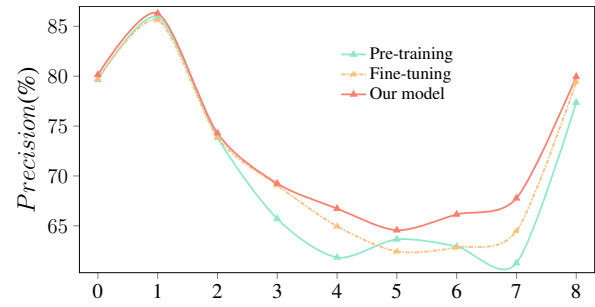


Figure 5: Precision of diverse models regarding different binned head absolute distances with punctuation.

**Absolute distance.** Figure 5 shows the effects of absolute distances from the head word to the modifier word on dev data. First, the “Pre-training” model achieves the lowest performance at most absolute distances, revealing that not all knowledge of source language is equally important to improve cross-lingual dependency performance. Second, compared with the “Pre-training” model, the “Fine-tuning” model achieves better performance at distances above 6, demonstrating that target language data can facilitate our model to capture the long dependency relationship. Finally, our model substantially enhances performances on all absolute distances, highlighting the importance of filtering source language information.

## 5 Conclusion

We propose a feature selection approach to emphasize useful representative features and ignore the useless ones, thus improving the performance of cross-lingual dependency parsing. Our model not only exploits a representation alignment network that selectively filters advantageous source language representations at the input layer but also utilizes an adversarial network to strengthen context-invariant features within the encoding layer. Experiments on a benchmark dataset illustrate that our proposed model significantly outperforms several strong baseline models. Detailed comparative experiments show that both the alignment and adversarial networks can substantially facilitate extracting and utilizing relevant target language features, thereby increasing the adaptation capability of our model. Furthermore, in-depth analysis reveals that our model achieves notable improvements in parsing long-distance dependencies and exhibits robustness capabilities, confirming its comprehensive applicative value in cross-lingual settings.



## 584 Limitations

585 Our proposed representation alignment and adver-  
586 sarial networks require a high-quality bilingual dic-  
587 tionary to facilitate language associations through  
588 matrix alignment. Hence, when there exists a bilin-  
589 gual dictionary, our method can be easily adapted  
590 to other cross-lingual dependency parsing tasks.  
591 Meanwhile, our constructed Vietnamese-Chinese  
592 bilingual dictionary will be released to facilitate  
593 future researches.

## 594 References

595 Wasi Uddin Ahmad, Zhisong Zhang, Xueze Ma, Kai-  
596 Wei Chang, and Nanyun Peng. 2019. [Cross-lingual  
597 dependency parsing with unlabeled auxiliary lan-  
598 guages](#). In *Proceedings of CoNLL*, pages 372–382.

599 Fadli Aulawi Al Ghiffari, Ika Alfina, and Kurniawati  
600 Azizah. 2023. [Cross-lingual transfer learning for  
601 javanese dependency parsing](#). In *Proceedings of  
602 IJCNLP-AACL*, pages 1–9.

603 Somnath Basu Roy Chowdhury, Annervaz M, and  
604 Ambedkar Dukkipati. 2019. [Instance-based induc-  
605 tive deep transfer learning by cross-dataset querying  
606 with locality sensitive hashing](#). In *Proceedings of  
607 DeepLo*, pages 183–191.

608 Kunal Chawla and Diyi Yang. 2020. [Semi-supervised  
609 formality style transfer using language model dis-  
610 criminator and mutual information maximization](#). In  
611 *Findings of EMNLP*, pages 2340–2354.

612 Guandan Chen, Kai Fan, Kaibo Zhang, Boxing Chen,  
613 and Zhongqiang Huang. 2021. [Manifold adversarial  
614 augmentation for neural machine translation](#). In  
615 *Findings of ACL-IJCNLP*, pages 3184–3189.

616 Xilun Chen, Ahmed Hassan Awadallah, Hany Hassan,  
617 Wei Wang, and Claire Cardie. 2019. [Multi-source  
618 cross-lingual model transfer: Learning what to share](#).  
619 In *Proceedings of ACL*, pages 3098–3112.

620 Chinmay Choudhary and Colm O’riordan. 2023. [Multi-  
621 lingual end-to-end dependency parsing with linguis-  
622 tic typology knowledge](#). In *Proceedings of the 5th  
623 Workshop on Research in Computational Linguistic  
624 Typology and Multilingual NLP*, pages 12–21.

625 Kevin Clark, Minh-Thang Luong, Christopher D. Man-  
626 ning, and Quoc Le. 2018. [Semi-supervised sequence  
627 modeling with cross-view training](#). In *Proceedings  
628 of EMNLP*, pages 1914–1925.

629 Cheikh M. Bamba Dione. 2021. [Multilingual depen-  
630 dency parsing for low-resource African languages:  
631 Case studies on Bambara, Wolof, and Yoruba](#). In  
632 *Proceedings of IWPT*, pages 84–92.

633 Timothy Dozat and Christopher D. Manning. 2017.  
634 [Deep biaffine attention for neural dependency pars-  
635 ing](#). In *Proceedings of ICLR*.

Thomas Effland and Michael Collins. 2023. [Improv-  
ing low-resource cross-lingual parsing with expected  
statistic regularization](#). *TACL*, pages 122–138. 636  
637  
638

Goran Glavaš and Ivan Vulić. 2020. [Non-linear  
instance-based cross-lingual mapping for non-  
isomorphic embedding spaces](#). In *Proceedings of  
ACL*, pages 7548–7555. 639  
640  
641  
642

Goran Glavaš and Ivan Vulić. 2021a. [Climbing the  
tower of treebanks: Improving low-resource depen-  
dency parsing via hierarchical source selection](#). In  
*Findings of ACL-IJCNLP*, pages 4878–4888. 643  
644  
645  
646

Goran Glavaš and Ivan Vulić. 2021b. [Climbing the  
tower of treebanks: Improving low-resource depen-  
dency parsing via hierarchical source selection](#). In  
*Findings of ACL-IJCNLP*, pages 4878–4888. 647  
648  
649  
650

Jan Hajic, Massimiliano Ciaramita, Richard Johansson,  
Daisuke Kawahara, M Antònia Martí, Lluís Màrquez,  
Adam Meyers, Joakim Nivre, Sebastian Padó, Jan  
Štěpánek, et al. 2009. [The conll-2009 shared task:  
Syntactic and semantic dependencies in multiple lan-  
guages](#). In *Proceedings of CoNLL*, pages 1–18. 651  
652  
653  
654  
655  
656

Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021.  
[Diverse adversaries for mitigating bias in training](#). In  
*Proceedings of EACL*, pages 2760–2765. 657  
658  
659

Amir Hazem, Mérieme Bouhandi, Florian Boudin, and  
Béatrice Daille. 2022. [Cross-lingual and cross-  
domain transfer learning for automatic term extrac-  
tion from low resource data](#). In *Proceedings of  
LREC*, pages 648–662. 660  
661  
662  
663  
664

Dan Kondratyuk and Milan Straka. 2019. [75 languages,  
1 model: Parsing Universal Dependencies univer-  
sally](#). In *Proceedings of EMNLP-IJCNLP*, pages  
2779–2795. 665  
666  
667  
668

Shanu Kumar, Soujanya Abbaraju, Sandipan Dandapat,  
Sunayana Sitaram, and Monojit Choudhury. 2023.  
[DiTTO: A feature representation imitation approach  
for improving cross-lingual transfer](#). In *Proceedings  
of EACL*, pages 385–406. 669  
670  
671  
672  
673

Guillaume Lample, Miguel Ballesteros, Sandeep Sub-  
ramanian, Kazuya Kawakami, and Chris Dyer. 2016.  
[Neural architectures for named entity recognition](#). In  
*Proceedings of NAACL-HLT*, pages 260–270. 674  
675  
676  
677

Ying Li, Shuaike Li, and Min Zhang. 2022. [Semi-  
supervised domain adaptation for dependency pars-  
ing with dynamic matching network](#). In *Proceedings  
of ACL*, pages 1035–1045. 678  
679  
680  
681

Ying Li, Zhenghua Li, Min Zhang, Rui Wang, Sheng Li,  
and Luo Si. 2019. [Self-attentive biaffine dependency  
parsing](#). In *Proceedings of IJCAI*, pages 5067–5073. 682  
683  
684

Ying Li, Meishan Zhang, Zhenghua Li, Min Zhang,  
Zhefeng Wang, Baoxing Huai, and Nicholas Jing  
Yuan. 2021. [APGN: Adversarial and parameter gen-  
eration networks for multi-source cross-domain de-  
pendency parsing](#). In *Findings of EMNLP*, pages  
1724–1733. 685  
686  
687  
688  
689  
690

|     |  |   |  |
|-----|--|---|--|
| 691 | Jiduan Liu, Jiahao Liu, Qifan Wang, Jingang Wang, Xunliang Cai, Dongyan Zhao, Ran Wang, and Rui Yan. 2023a. Retrieval-based knowledge transfer: An effective approach for extreme large language model compression. In <i>Findings of EMNLP</i> , pages 8643–8657. | Tal Schuster, Ori Ram, Regina Barzilay, and Amir Globerson. 2019. Cross-lingual alignment of contextual word embeddings, with applications to zero-shot dependency parsing. In <i>Proceedings of NAACL</i> , pages 1599–1613.   | 744<br>745<br>746<br>747<br>748        |
| 697 | Lu Liu, Yi Zhou, Jianhan Xu, Xiaoqing Zheng, Kai-Wei Chang, and Xuan-Jing Huang. 2020. Cross-lingual dependency parsing by pos-guided word reordering. In <i>Findings of EMNLP</i> , pages 2938–2948.  | Weiyang Shi, Ryan Shea, Si Chen, Chiyuan Zhang, Ruoxi Jia, and Zhou Yu. 2022. Just fine-tune twice: Selective differential privacy for large language models. In <i>Proceedings of EMNLP</i> , pages 6327–6340.   | 749<br>750<br>751<br>752               |
| 701 | Yihong Liu, Haotian Ye, Leonie Weissweiler, Renhao Pei, and Hinrich Schuetze. 2023b. Crosslingual transfer learning for low-resource languages based on multilingual colexification graphs. In <i>Findings of EMNLP</i> , pages 8376–8401.                         | Milan Straka, Jana Straková, and Jan Hajic. 2019. Evaluating contextualized embeddings on 54 languages in POS tagging, lemmatization and dependency parsing. <i>CoRR</i> , abs/1908.07448.  | 753<br>754<br>755<br>756               |
| 706 | Daniel Lowd and Christopher Meek. 2005. Adversarial learning. In <i>Proceedings of ACM SIGKDD</i> , pages 641–647.   | Weiting Tan, Shuoyang Ding, Huda Khayrallah, and Philipp Koehn. 2022. Doubly-trained adversarial data augmentation for neural machine translation. In <i>Proceedings of AMTA</i> , pages 157–174.   | 757<br>758<br>759<br>760               |
| 709 | Menglong Lu, Zhen Huang, Yunxiang Zhao, Zhiliang Tian, Yang Liu, and Dongsheng Li. 2023. DaMSTF: Domain adversarial learning enhanced meta self-training for domain adaptation. In <i>Proceedings of ACL</i> , pages 1650–1668.                                    | Yuanhe Tian, Yan Song, and Fei Xia. 2022. Improving relation extraction through syntax-induced pre-training with dependency masking. In <i>Findings of ACL</i> , pages 1875–1886, Dublin, Ireland.  | 761<br>762<br>763<br>764               |
| 714 | Alireza Mohammadshahi and James Henderson. 2021. Recursive non-autoregressive graph-to-graph transformer for dependency parsing with iterative refinement. <i>TACL</i> , 9:120–138.  | Ahmet Üstün, Arianna Bisazza, Gosse Bouma, and Gertjan van Noord. 2022. Uadapter: Typology-based language adapters for multilingual dependency parsing and sequence labeling. <i>Computational Linguistics</i> , 48(3):555–592.   | 765<br>766<br>767<br>768<br>769        |
| 718 | Noriki Nishida and Yuji Matsumoto. 2022. Out-of-domain discourse dependency parsing via bootstrapping: An empirical analysis on its effectiveness and limitation. <i>TACL</i> , 10:127–144.  | Clara Vania, Yova Kementchedjheva, Anders Søgaard, and Adam Lopez. 2019. A systematic comparison of methods for low-resource dependency parsing on genuinely low-resource languages. In <i>Proceedings of EMNLP-IJCNLP</i> , pages 1105–1116.   | 770<br>771<br>772<br>773<br>774        |
| 722 | Tong Niu, Kazuma Hashimoto, Yingbo Zhou, and Caiming Xiong. 2022. OneAligner: Zero-shot cross-lingual transfer with one rich-resource language pair for low-resource sentence retrieval. In <i>Findings of ACL</i> , pages 2869–2882.                              | Zihan Wang, Karthikeyan K, Stephen Mayhew, and Dan Roth. 2020. Extending multilingual BERT to low-resource languages. In <i>Findings of EMNLP</i> , pages 2649–2656.  | 775<br>776<br>777<br>778               |
| 727 | Hiroaki Ozaki, Gaku Morio, Terufumi Morishita, and Toshinori Miyoshi. 2021. Project-then-transfer: Effective two-stage cross-lingual transfer for semantic dependency parsing. In <i>Proceedings of ACL</i> , pages 2586–2594.                                     | Haoran Xu and Philipp Koehn. 2021. Zero-shot cross-lingual dependency parsing through contextual embedding transformation. In <i>Proceedings of Adapt-NLP</i> .   | 779<br>780<br>781<br>782               |
| 732 | Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2021. UNKs everywhere: Adapting multilingual language models to new scripts. In <i>Proceedings of EMNLP</i> , pages 10186–10203.  | Adam Yaari, Jan DeWitt, Henry Hu, Bennett Stankovits, Sue Felshin, Yevgeni Berzak, Helena Aparicio, Boris Katz, Ignacio Cases, and Andrei Barbu. 2022. The aligned multimodal movie treebank: An audio, video, dependency-parse treebank. In <i>Proceedings of EMNLP</i> , pages 9531–9539. | 783<br>784<br>785<br>786<br>787<br>788 |
| 736 | Matthew Riemer, Sophia Krasikov, and Harini Srinivasan. 2015. A deep learning and knowledge transfer based architecture for social media user characteristic determination. In <i>Proceedings of SocialNLP</i> , pages 39–47.                                      | Hang Yan, Xipeng Qiu, and Xuanjing Huang. 2020. A graph-based model for joint Chinese word segmentation and dependency parsing. <i>TACL</i> , 8:78–92.  | 789<br>790<br>791                      |
| 741 | Guy Rotman and Roi Reichart. 2019. Deep contextualized self-training for low resource dependency parsing. <i>TACL</i> , 7:695–713.   | Nasser Zalmout and Nizar Habash. 2019. Adversarial multitask learning for joint multi-feature and multi-dialect morphological modeling. In <i>Proceedings of ACL</i> , pages 1775–1786.   | 792<br>793<br>794<br>795               |

- 796 Brian Hu Zhang, Blake Lemoine, and Margaret  
797 Mitchell. 2018. [Mitigating unwanted biases with](#)  
798 [adversarial learning](#). In *Proceedings of AAAI/ACM*  
799 *AIES*, pages 335–340.
- 800 Meishan Zhang, Zhenghua Li, Guohong Fu, and Min  
801 Zhang. 2019. [Syntax-enhanced neural machine](#)  
802 [translation with syntax-aware word representations](#).  
803 In *Proceedings of NAACL-HLT*, pages 1151–1161.
- 804 Yue Zhang, Bo Zhang, Zhenghua Li, Zuyi Bao, Chen Li,  
805 and Min Zhang. 2022. [SynGEC: Syntax-enhanced](#)  
806 [grammatical error correction with a tailored GEC-](#)  
807 [oriented parser](#). In *Proceedings of EMNLP*, pages  
808 2518–2531.
- 809 Han Zou, Jianfei Yang, and Xiaojian Wu. 2021. [Unsu-](#)  
810 [pervised energy-based adversarial domain adaptation](#)  
811 [for cross-domain text classification](#). In *Findings of*  
812 *ACL-IJCNLP*, pages 1208–1218.