# Towards Better UD Parsing: Deep Contextualized Word Embeddings, Ensemble, and Treebank Concatenation

## Abstract

This paper describes our system (HIT-SCIR) for CoNLL 2018 shared task on Universal Dependency parsing. We base our system on Stanford's winning system for the CoNLL 2017 shared task and made two effective extensions: ensembling parsers trained with different initialization and incorporating deep contextualized word embeddings into both the POS tagger and parser. We also explore different ways of concatenating treebanks and it leads further improvements. Experimental results on the development data show the effectiveness of our methods. In the evaluation on the test data, our system was ranked first according to LAS (75.84%) and outperformed the secondly ranked system by 2.56%.

## 1 Introduction

In this paper, we describe our system (HIT-SCIR) submitted to CoNLL 2018 shared task on Universal Dependency parsing (Zeman et al., 2018). We base our system on Stanford's winning system (Dozat et al., 2017, §2) for the CoNLL 2017 shared task (Zeman et al., 2017). Dozat and Manning (2016) and its extension (Dozat et al., 2017) have show very competitive performance in previous parsing works (Ma and Hovy, 2017; Shi et al., 2017a; Liu et al., 2018b; Ma et al., 2018). In this paper, we would like to know can we further improve their POS tagger and parser with cheap and universal solution.

In our system, we make two noteworthy extensions to their parser and tagger, which includes

- Incorporating the deep contextualized word embeddings (Peters et al., 2018, ELMo) into the word representaton (§3).

- Ensembling parsers trained with different initialization (§4).

In the shared task, multiple treebanks of different domains are provided for some languages. There are also treebanks which are of the same language families. Letting these treebanks help each other has been shown an effective way to improve parsing performance in last year's shared tasks (Che et al., 2017; Shi et al., 2017b; Björkelund et al., 2017). In our system, we apply the simple concatenation to the treebanks that are potentially helpful to each other and explore different ways of concatenation to improve the parser's performance (§5).

In dealing with the small languages and low-resource languages (§6), we adopt the word embedding transfer idea in the cross-lingual dependency parsing (Guo et al., 2015) and use the bilingual word vectors transformation technique (Smith et al., 2017) to map *fasttext* word embeddings (Bojanowski et al., 2016) of the source rich-resource language and target low-resource language into the same space and transfer the parser trained on the source language into the target language.

We conduct experiments on the development set to study the effects of ELMo, parser ensemble, and treebank concatenation. Experimental results show that these techniques substantially improve the parsing performance. Using these techniques, our system achieved an averaged LAS of 75.84 on the official test set and was ranked the first according to LAS (Zeman et al., 2018). This result outperforms the seondly ranked by 2.56%.

## 2 Dozat et al. (2017)

We based our system on the tagger and parser of Dozat et al. (2017). The core idea of the tagger and parser is using a LSTM network to produce

the vector representation for each word and then predict POS tags and dependency relations using the representation. For the tagger whose input is the word alone, this representation is calculated as

$$\mathbf{r}_i = \text{BiLSTM}(\mathbf{r}_0, (\mathbf{v}_1^{(word)}, ..., \mathbf{x}_n^{(word)}))_i$$
$$\mathbf{h}_i, \mathbf{c}_i = \text{split}(\mathbf{r}_i),$$

where $\mathbf{v}_i^{(word)}$ is the word embeddings. After getting $\mathbf{h}_i$, the scores of tags are calculated as

$$\mathbf{h}_i^{(pos)} = \text{MLP}^{(pos)}(\mathbf{h}_i)$$
$$\mathbf{s}_i^{(pos)} = W \cdot \mathbf{h}_i^{(pos)} + \mathbf{b}^{(pos)}$$

where each element in $\mathbf{s}_i^{(pos)}$ represents the possibility that $i$-th word is assigned with corresponding tag.

For the parser whose inputs are the word and POS tag, such representation is calculated as

$$\mathbf{x}_i = \mathbf{v}_i^{(word)} \oplus \mathbf{v}_i^{(tag)}$$
$$\mathbf{r}_i = \text{BiLSTM}(\mathbf{r}_0, (\mathbf{x}_1, ..., \mathbf{x}_n))_i$$
$$\mathbf{h}_i, \mathbf{c}_i = \text{split}(\mathbf{r}_i).$$

And a pair of representations are fed into a biaffine classifier to predict if there is a dependency arc between these two words. The scores over all head words are calculated as

$$\mathbf{s}_i^{(rel)} = \mathbf{h}_{y^{(arc)}}^{T(rel-head)} \mathbf{U}^{(rel)} \mathbf{h}_i^{(rel-dep)}$$
$$+ W^{(rel)}(\mathbf{h}_i^{(rel-dep)} \oplus \mathbf{h}_{y^{(arc)}}^{T(rel-head)})$$
$$+ \mathbf{b}^{(rel)}$$

For both the biaffine tagger and parser, the word embedding $\mathbf{v}_i^{(word)}$ is obtained by summarizing a fine-tuned token embedding $\mathbf{w}_i$, a fixed word2vec embedding $\mathbf{p}_i$, and a LSTM-encoded character representation $\hat{\mathbf{v}}_i$ as

$$\mathbf{v}_i^{(word)} = \mathbf{w}_i + \mathbf{p}_i + \hat{\mathbf{v}}_i$$

## 3 Deep Contextualized Word Embeddings

Deep contextualized word embeddings (Peters et al., 2018, ELMo) has shown to be very effective on a range of syntactic and semantic tasks and it's straightforward to achieve by using a LSTM network to encode words in a sentence and training the LSTM network with language modeling objective on large-scale raw text. More specifically,

the $\mathbf{ELMo}_i$ is computed by first computing the hidden representation $\mathbf{h}_i^{(LM)}$ as

$$\mathbf{r}_i^{(LM)} = \text{BiLSTM}^{(LM)}(\mathbf{r}_0^{(LM)}, (\tilde{\mathbf{v}}_1, ..., \tilde{\mathbf{v}}_n))_i$$
$$\mathbf{h}_i^{(LM)}, \mathbf{c}_i^{(LM)} = \text{split}(\mathbf{r}_i^{(LM)}),$$

where $\tilde{\mathbf{v}}_i$ is the output of a CNN over characters, then attentively summarizing and scaling different layers of $\mathbf{h}_{i,j}^{(LM)}$ with $s_j$ and $\gamma$ as

$$\mathbf{ELMo}_i = \gamma \sum_{j=0}^{L} s_j \mathbf{h}_{i,j}^{(LM)},$$

where $L$ is the number of layers and $\mathbf{h}_{i,0}^{(LM)}$ is identical to $\tilde{\mathbf{v}}_i$. In our system, we follow Peters et al. (2018) and use a two-layer bidirectional LSTM as our $\text{BiLSTM}^{(LM)}$.

In this paper, we study the usage of ELMo for improving both the tagger and parser and make several simplifications. Different from Peters et al. (2018), we treat the output of ELMo as a fixed representation and do not tune its parameters during tagger and parser training. Thus, we cancel the layer-wise attention scores $s_j$ and the scaling factor $\gamma$. As Peters et al. (2018) pointed, the second layer of $\text{BiLSTM}^{(LM)}$ mainly captures the semantic information, and we care more about its syntactic part (i.e. the first layer of $\text{BiLSTM}^{(LM)}$). In our system, we only use $\mathbf{h}_{i,0}^{(LM)}$ and $\mathbf{h}_{i,1}^{(LM)}$ for $\mathbf{ELMo}_i$ which means

$$\mathbf{ELMo}_i = \mathbf{h}_{i,0}^{(LM)} + \mathbf{h}_{i,1}^{(LM)}.$$

After getting $\mathbf{ELMo}_i$, we project it to the same dimension as $\mathbf{v}_i^{(word)}$ and use it as an additional word embeddings. The calculation of $\mathbf{v}_i^{(word)}$ becomes

$$\mathbf{v}_i^{(word)} = \mathbf{w}_i + \mathbf{p}_i + \hat{\mathbf{v}}_i + W^{(ELMo)} \cdot \mathbf{ELMo}_i$$

for both the tagger and parser. We need to note that training the tagger and parser includes $W^{(ELMo)}$. To avoid overfitting, we impose a dropout function on projected vector $W^{(ELMo)} \cdot \mathbf{ELMo}_i$ during training.

We use the same hyperparameter settings as Peters et al. (2018) for $\text{BiLSTM}^{(LM)}$ and the character CNN. We train their parameters as training a bidirectional language model on a randomly sampled subset of raw text released by the shared task. More specifically, for each language, we randomly sample 20 million words. The training of ELMo on one language takes roughly 3 days on an NVIDIA P100 GPU.

## 4 Parser Ensemble

According to Reimers and Gurevych (2017), neural network training can be sensitive to initialization and Liu et al. (2018a) shows that ensemble neural network trained with different initialization leads to performance improvements. We follow their works and train three parsers with different initialization, then ensemble these parsers by averaging their output scores as $\mathbf{s}_i^{(rel)} = \frac{1}{3}\sum_{i=1}^{3}\mathbf{s}_i^{(i,rel)}$.

## 5 Treebank Concatenation

For 15 out of the 58 languages in the shared task, multiple treebanks from different domains are provided. There are also treebanks that comes from the same language families. Taking the advantages of the relation between treebanks has been shown a promising direction in both the research community (Ammar et al., 2016; Guo et al., 2015, 2016) and in the CoNLL 2017 shared task (Che et al., 2017; Björkelund et al., 2017; Shi et al., 2017b). In our system, we adopt the treebank concatenation technique as (Ammar et al., 2016) but limit that only a group of treebanks from the same language (*cross-domain concatenation*) or a pair of treebanks that are typologically or geographically correlated (*cross-lingual concatenation*) is concatenated.

In our system, we tried cross-domain concatenation on *nl*, *sv*, *ko*, *it*, *en*, *fr*, *gl*, *la*, *ru*, and *sl*. We also tried cross-lingual concatenation on *ug-tr*, *uk-ru*, *ga-en*, and *sme-fi* following (Che et al., 2017). However, due to the variance in vocabulary, grammatical genre, and even annotation, treebank concatenation does not guarantee to improve the model's performance. We decide the usage of concatenation by examining their development set performance. For some small treebanks which do not have development set, whether using treebank concatenation is decided through 5-fold cross validation.[1] We show the experimental results of treebank concatenation in Section 8.3.

## 6 Low Resources Languages

In the shared task, 5 languages are presented with training set of less than 50 sentences. 4 languages don't even have any training data. It's difficult to

| target | br | fo | hy | kk | bxr | kmr | hsb | th |
|--------|----|----|----|----|-----|-----|-----|-----|
| source | ga | no | et | tr | hi  | fa  | pl  | zh |

Table 1: Cross-lingual transfer settings for low-resource target languages.

train reasonable parser on these low-resource languages We deal with these treebanks by adopting the word embedding transfer idea of Guo et al. (2015). We transfer the word embeddings of the rich-resource language to the space of low-resource language using the bilingual word vectors transformation technique (Smith et al., 2017) and trained a parser using the source treebank with only pretrained word embeddings on the transformed space as $\mathbf{v}_i^{(word)} = \mathbf{p}_i$. The transformation matrix is automatically learned on the *fasttext* word embedding using the same tokens between two languages (like punctuations).

Table 1 shows our source languages for the target low-resource languages. For the treebank with a few training data, its source language is decided by testing the source parser's performance on the training data. For the treebank without any training data, we choose the source language according to their language family.

*Naija* presents an exception for our method since it doesn't have *fasttext* word embeddings and embedding transformation is infeasible. Since it's a dialect of English, we use the full pipeline of *en_ewt* for *pcm_nsc*.

## 7 Preprocessing

Besides improving the tagging and parsing performance, we also consider the preprocessing as an important factor to the final performance.

### 7.1 Sentence Segmentation

For some treebanks, sentence segmentation can be problematic since there is no explicitly sentence delimiters. de Lhoneux et al. (2017) presented a joint tokenization and sentence segmentation model[2] that outperformed the baseline model in last year's shared task (Zeman et al., 2017). We select a set of treebanks whose *udpipe* sentence segmentation F-scores are lower than 95 on the development set and use Uppsala segmentor instead.[3] Using the Uppsala segmentor leads to an

---

[1] We use *udpipe* (Straka et al., 2016) for the experiments of these part because we consider the effect of treebank concatenation as being irrelevant to the parser architecture and *udpipe* has the speed advantage in both training and testing.

[2] noted as Uppsala segmentor henceforth.

[3] We use Uppsala segmentor for *it_postwita*, *got_proiel*, *la_poroiel*, *cu_proiel*, *grc_proiel*, *sl_ssj*, *nl_lassysmall*, *fi_tdt*, *pt_bosque*, *da_ddt*, *id_gsd*, *el_gdt*, and *et_edt*.
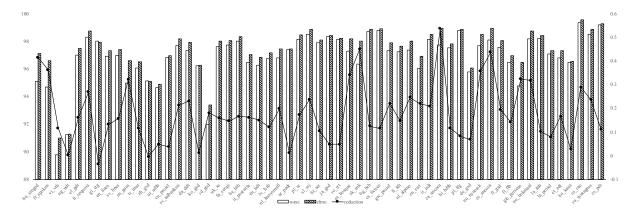
Figure 1: The effects of ELMo on POS tagging. Treebanks are sorted according to the number of training sentences from left to right.
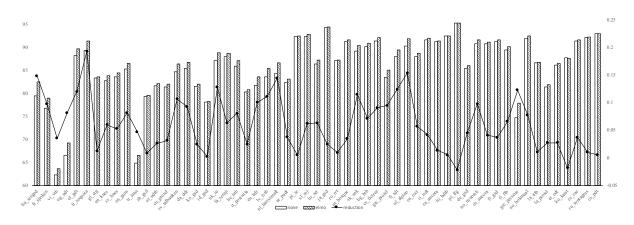


Figure 2: The effects of ELMo on parsing. Treebanks are sorted according to the number of training sentences from left to right.

improvement of 7.67 F-score in these treebanks over *udpipe* baseline.

## 7.2 Tokenization for Chinese, Japanese, and Vietnamese

Tokenization is non-trivial for languages which do not have explicit word boundary markers, like Chinese, Japanese, and Vietnamese. We develop our own tokenizer for these three languages. Following Che et al. (2017) and Zheng et al. (2017), we model the tokenization as labeling the word boundary tag on character and use features derived from large-scale unlabeled data to further improve the performance.[4] In addition to the pointwise mutual information (PMI), we also incorporate the character ELMO into our tokenizer. These techniques lead to the best tokenization performance on all the related treebanks and the average improvement over *udpipe* baseline is 7.5 in tokeniza-

tion F-score.[5]

## 7.3 Preprocessing for Thai

Thai language presents unique challenge in the preprocessing. Our survey on the Thai Wikipedia indicates that there is no explicit sentence delimiter for Thai language and obtaining Thai words requires tokenization. To this remedy, we use space as sentence delimiter and use the lexicon-based word segmentation – forward maximum matching algorithm for Thai tokenization. Our lexicon is derived from the *fasttext* word embeddings.

## 8 Results

### 8.1 Effects of ELMo

We study the effect of ELMo on the large treebanks and report the results of singe tagger and parser with and without ELMo. Figure 1 shows

---

[4]For Vietnamese where whitespaces occur both inter- and intra-words, we treat the whitespace-separated token as a character.

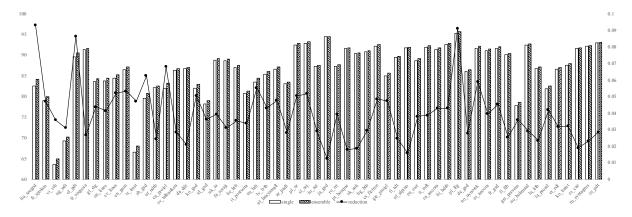[5]*ja_gsd*, *ja_modern*, *vi_vtb*, and *zh_gsd*.

Figure 3: The effects of ensemble on parsing. Treebanks are sorted according to the number of training sentences from left to right.

the tagger results on the development set and Figure 2 shows the parser results. Using ELMo in the tagger leads to a macro-averaged improvement of 0.56% in UPOS and the macro-averaged error reduction is 17.83%. and using ELMo in the parser leads to a macro-averaged improvement of 0.84% in LAS and the macro-averaged error reduction is 7.88%.

ELMo improves the tagging performance almost on every treebank, except for *zh_gsd* and *gl_ctg*. Similar trends are witnessed in the parsing experiments with *ko_kaist* and *pl_lfg* being the only treebanks that ELMo slightly worsens the performance.

We also study the relative improvements against the size of the treebank. However, no clear relation is revealed between the treebank size and the gains using ELMo.

## 8.2 Effects of Ensemble

We also test the effect of ensemble and show the results in Figure 3. Parser ensemble leads to an averaged improvement of 0.55% in LAS and the averaged error reduction is 4.0%. These results indicate that ensemble is a effective way of improving the parsing performance.

## 8.3 Effects of Treebank Concatenation

As mentioned in Section 5, we study the effects of both the *cross-domain concatenation* and *cross-lingual concatenation*.

**Cross-Domain Concatenation.** For the treebanks which have development set, the development performances are shown in Table 2. Numbers of sentences in the training set are also shown

in this table. The general trend is that for the treebank with small training set, cross-domain concatenation achieves better performance. While for those with large training set, concatenation doesn't improve the performance or even worsen the results.

For the small treebanks which do not have development set, the 5 fold cross validation results are shown in Table 3 in which concatenation improves most of the treebanks except for *gl_treegal*.

**Cross-Lingual Concatenation.** The experimental results of cross-lingual concatenation are shown in Table 4. Unfortunately, concatenating treebanks from different languages only achieves improved performance on *uk_iu*. This results also indicate that in cross lingual parsing, sophisticated methods like word embeddings transfer (Guo et al., 2015) and treebank transfer (Guo et al., 2016) are still necessary.

## 8.4 Effects of Better Preprocessing

We also study how preprocessing contributes to the final parsing performance. The experimental results on the development set are shown in Table 5. From this table, the performance of word segmentation is almost linearly correlated with the final performance. Similar trends on sentence segmentation performance are witnessed but *el_gdt* and *pt_bosque* presents some exceptions where better preprocess leads drop in the final parsing performance.

## 9 Conclusion

Our system submitted to the CoNLL 2018 shared task made several improvements on last year's

| nl | apino | lassysmall | | sv | lines | talbanken | | ko | gsd | kaist | | it | isdt | postwita |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # train | 12.2 | 5.8 | | # train | 2.7 | 4.3 | | # train | 4.4 | 23.0 | | # train | 13.1 | 5.4 |
| apino | 91.87 | | | lines | 84.64 | | | gsd | 82.05 | | | isdt | **92.01** | |
| lassysmall | | 86.82 | | talbanken | | 86.39 | | kaist | | **87.83** | | postwita | | 80.79 |
| concat. | **92.08** | **89.34** | | concat. | **85.76** | **86.77** | | concat. | **83.73** | 87.61 | | concat. | 91.80 | **82.54** |

| en | ewt | gum | lines | | fr | gsd | sequoia | spoken |
|---|---|---|---|---|---|---|---|---|
| # train | 12.5 | 2.9 | 2.7 | | # train | 14.6 | 2.2 | 1.2 |
| ewt | **88.75** | | | | gsd | **91.64** | | |
| gum | | **86.52** | | | sequoia | | **91.44** | |
| lines | | | 83.86 | | spoken | | | 79.06 |
| concat. | 88.74 | 85.65 | **85.30** | | concat. | 91.44 | 90.51 | **81.99** |

Table 2: The developement performance with cross-domain concatenation for languages which has multiple treebanks. *# train* shows the number of training sentences in the treebank. We opt out *cs*, *fi*, and *pl* because all the treebanks of these languages are relatively large (# train > 10K).

| gl | treegal | | la | perseus | | no | nynorsklia | | ru | taiga | | sl | sst |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # train | 0.6 | | # train | 1.3 | | # train | 0.3 | | # train | 0.9 | | # train | 2.1 |
| treegal | **66.71** | | perseus | 44.05 | | nynorsklia | 51.05 | | taiga | 54.70 | | sst | 55.15 |
| +ctg | 56.73 | | +proiel | **50.78** | | +nynorsk | **58.49** | | +syntagrus | **60.75** | | +ssj | **59.52** |

Table 3: The 5-fold cross validation results for the cross-domain concatenation for treebank which doesn't have development set.

winning system from Dozat et al. (2017), including incorporating deep contextualized word embeddings, parser ensemble, and treebank concatenation. Experimental results on the development set show the effectiveness of our methods. Using these techniques, our system achieved an averaged LAS of 75.84% and obtained the first place in LAS in the final evaluation.

## 10 Credits

There are a few references we would like to give proper credit, especially to data providers: the core Universal Dependencies paper from LREC 2016 (Nivre et al., 2016), the UD version 2.2 datasets (Nivre et al., 2018), the baseline *udpipe* model released by Straka et al. (2016), the deep contextualized word embeddings code released by Peters et al. (2018), the biaffine tagger and parser released by Dozat et al. (2017), the joint sentence segmentor and tokenizer released by de Lhoneux et al. (2017), and the evaluation platform TIRA (Potthast et al., 2014).

## Acknowledgments

## References

Waleed Ammar, George Mulcaire, Miguel Ballesteros, Chris Dyer, and Noah Smith. 2016. Many languages, one parser. *TACL* 4.

Anders Björkelund, Agnieszka Falenska, Xiang Yu, and Jonas Kuhn. 2017. Ims at the conll 2017 ud shared task: Crfs and perceptrons meet neural networks. In *Proc. of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. http://www.aclweb.org/anthology/K17-3004.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching word vectors with subword information. *CoRR* abs/1607.04606. http://arxiv.org/abs/1607.04606.

Wanxiang Che, Jiang Guo, Yuxuan Wang, Bo Zheng, Huaipeng Zhao, Yang Liu, Dechuan Teng, and Ting Liu. 2017. The hit-scir system for end-to-end parsing of universal dependencies. In *Proc. of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. http://www.aclweb.org/anthology/K17-3005.

Miryam de Lhoneux, Yan Shao, Ali Basirat, Eliyahu Kiperwasser, Sara Stymne, Yoav Goldberg, and Joakim Nivre. 2017. From raw text to universal dependencies - look, no tags! In *Proc. of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. http://www.aclweb.org/anthology/K17-3022.

Timothy Dozat and Christopher D. Manning. 2016. Deep biaffine attention for neural dependency parsing. *CoRR* abs/1611.01734.

| | ug_udt | | | uk_iu | | | ga_idt | | | sme_giella | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ug_udt | **69.27** | | uk_iu | 88.84 | | ga_idt | **62.84** | | sme_giella | **66.33** |
| +tr_imst | 19.27 | | +ru_syntagus | **90.74** | | +en_ewt | 51.00 | | +fi_ftb | 59.86 |

Table 4: Cross-lingual concatenation results. The results for *ug_udt* and *uk_iu* are obtained on the development set. The results for *ga_idt* and *sme_giella* are obtained with *udpipe* by 5-fold cross validation.

| | $\Delta$-sent | *udpipe* | *improved* |
|---|---|---|---|
| fi_tdt | +0.69 | 88.13 | **88.67** |
| et_edt | +1.22 | 86.33 | **86.36** |
| nl_lassysmall | +1.39 | 88.08 | **88.60** |
| da_ddt | +1.56 | 86.21 | **86.51** |
| el_gdt | +1.57 | **90.08** | 89.96 |
| cu_proiel | +1.72 | 72.79 | **74.04** |
| pt_bosque | +1.83 | **90.73** | 90.20 |
| id_gsd | +2.46 | 74.14 | **78.83** |
| la_proiel | +4.82 | 73.21 | **74.22** |
| got_proiel | +5.36 | 67.55 | **68.40** |
| grc_proiel | +5.86 | 79.67 | **80.72** |
| sl_ssj | +18.81 | 88.43 | **92.27** |
| it_postwita | +30.40 | 74.91 | **79.26** |
| | $\Delta$-word | *udpipe* | *improved* |
| ja_gsd | +4.07 | 80.53 | **85.23** |
| zh_gsd | +7.16 | 66.16 | **75.78** |
| vi_vtb | +9.02 | 48.58 | **57.53** |

Table 5: The effect of improved preprocessing.

Timothy Dozat, Peng Qi, and Christopher D. Manning. 2017. Stanford's graph-based neural dependency parser at the conll 2017 shared task. In *Proc. of CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*.

Jiang Guo, Wanxiang Che, Haifeng Wang, and Ting Liu. 2016. A universal framework for inductive transfer parsing across multi-typed treebanks. In *Proc. of Coling*. http://www.aclweb.org/anthology/C16-1002.

Jiang Guo, Wanxiang Che, David Yarowsky, Haifeng Wang, and Ting Liu. 2015. Cross-lingual dependency parsing based on distributed representations. In *Proc. of ACL*.

Yijia Liu, Wanxiang Che, Huaipeng Zhao, Bing Qin, and Ting Liu. 2018a. Distilling knowledge for search-based structured prediction. *CoRR* abs/1805.11224. http://arxiv.org/abs/1805.11224.

Yijia Liu, Yi Zhu, Wanxiang Che, Bing Qin, Nathan Schneider, and Noah A. Smith. 2018b. Parsing tweets into universal dependencies. In *Proc. of NAACL*. http://aclweb.org/anthology/N18-1088.

Xuezhe Ma and Eduard Hovy. 2017. Neural probabilistic model for non-projective mst parsing. In *Proc. of IJCNLP*. http://www.aclweb.org/anthology/I17-1007.

Xuezhe Ma, Zecong Hu, Jingzhou Liu, Nanyun Peng, Graham Neubig, and Eduard H. Hovy. 2018. Stack-pointer networks for dependency parsing. *CoRR* abs/1805.01087. http://arxiv.org/abs/1805.01087.

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, and Daniel Zeman. 2016. Universal Dependencies v1: A multilingual treebank collection. In *Proc. of LREC-2016*.

Joakim Nivre et al. 2018. Universal Dependencies 2.2. LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics, Charles University, Prague, http://hdl.handle.net/11234/SUPPLYTHENEWPERMANENTIDHERE!

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*. http://aclweb.org/anthology/N18-1202.

Martin Potthast, Tim Gollub, Francisco Rangel, Paolo Rosso, Efstathios Stamatatos, and Benno Stein. 2014. Improving the reproducibility of PAN's shared tasks: Plagiarism detection, author identification, and author profiling. In Evangelos Kanoulas, Mihai Lupu, Paul Clough, Mark Sanderson, Mark Hall, Allan Hanbury, and Elaine Toms, editors, *Information Access Evaluation meets Multilinguality, Multimodality, and Visualization. 5th International Conference of the CLEF Initiative (CLEF 14)*. Springer, Berlin Heidelberg New York, pages 268–299. https://doi.org/10.1007/978-3-319-11382-1_22.

Nils Reimers and Iryna Gurevych. 2017. Reporting score distributions makes a difference: Performance study of lstm-networks for sequence tagging. In *Proc. of EMNLP*.

Tianze Shi, Liang Huang, and Lillian Lee. 2017a. Fast(er) exact decoding and global training for transition-based dependency parsing via a minimal feature set. In *Proc. of EMNLP*. https://www.aclweb.org/anthology/D17-1002.

Tianze Shi, Felix G. Wu, Xilun Chen, and Yao Cheng. 2017b. Combining global models for parsing universal dependencies. In *Proc. of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. http://www.aclweb.org/anthology/K17-3003.

Samuel L. Smith, David H. P. Turban, Steven Hamblin, and Nils Y. Hammerla. 2017. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *CoRR* abs/1702.03859. http://arxiv.org/abs/1702.03859.

Milan Straka, Jan Hajič, and Jana Straková. 2016. UD-Pipe: trainable pipeline for processing CoNLL-U files performing tokenization, morphological analysis, POS tagging and parsing. In *Proc. of LREC-2016*.

Daniel Zeman, Filip Ginter, Jan Hajič, Joakim Nivre, Martin Popel, and Milan Straka. 2018. CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies. In *Proc. of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*.

Daniel Zeman, Martin Popel, Milan Straka, Jan Hajič, Joakim Nivre, Filip Ginter, Juhani Luoto-lahti, Sampo Pyysalo, Slav Petrov, Martin Pot-thast, Francis Tyers, Elena Badmaeva, Memduh Gökırmak, Anna Nedoluzhko, Silvie Cinková, Jan Hajič jr., Jaroslava Hlaváčová, Václava Kettnerová, Zdeňka Urešová, Jenna Kanerva, Stina Ojala, Anna Missilä, Christopher Manning, Sebastian Schuster, Siva Reddy, Dima Taji, Nizar Habash, Herman Le-ung, Marie-Catherine de Marneffe, Manuela Sanguinetti, Maria Simi, Hiroshi Kanayama, Valeria de Paiva, Kira Droganova, Héctor Martínez Alonso, Çağrı Çöltekin, Umut Sulubacak, Hans Uszkor-eit, Vivien Macketanz, Aljoscha Burchardt, Kim Harris, Katrin Marheinecke, Georg Rehm, Tolga Kayadelen, Mohammed Attia, Ali Elkahky, Zhuoran Yu, Emily Pitler, Saran Lertpradit, Michael Mandl, Jesse Kirchner, Hector Fernandez Alcalde, Jana Str-nadova, Esha Banerjee, Ruli Manurung, Antonio Stella, Atsuko Shimada, Sookyoung Kwak, Gustavo Mendonça, Tatiana Lando, Rattima Nitisaroj, and Josie Li. 2017. CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies.

Bo Zheng, Wanxiang Che, Jiang Guo, and Ting Liu. 2017. Enhancing lstm-based word segmentation using unlabeled data. In *Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data*.