

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

Wanxiang Che, Yijia Liu, Yuxuan Wang, Bo Zheng, Ting Liu

Research Center for Social Computing and Information Retrieval

Harbin Institute of Technology, China

{car, yjliu, yxwang, bzheng, tliu}@ir.hit.edu.cn

Abstract

This paper describes our system (HIT-SCIR) submitted to the CoNLL 2018 shared task on Multilingual Parsing from Raw Text to Universal Dependencies. We base our submission on Stanford’s winning system for the CoNLL 2017 shared task and make two effective extensions: 1) incorporating deep contextualized word embeddings into both the part of speech tagger and parser; 2) ensembling parsers trained with different initialization. We also explore different ways of concatenating treebanks for further improvements. Experimental results on the development data show the effectiveness of our methods. In the final evaluation, our system was ranked first according to LAS (75.84%) and outperformed the other systems by a large margin.

1 Introduction

In this paper, we describe our system (HIT-SCIR) submitted to CoNLL 2018 shared task on Multilingual Parsing from Raw Text to Universal Dependencies (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 shown very competitive performance in both the shared task (Dozat et al., 2017) and previous parsing works (Ma and Hovy, 2017; Shi et al., 2017a; Liu et al., 2018b; Ma et al., 2018). A nature question that raises is how can we further improve their part of speech (POS) tagger and parser via simple yet effective technique. 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 representation (§3);
- Ensembling parsers trained with different initialization (§4).

Both these two extensions result in substantial improvements.

For some languages in the shared task, multiple treebanks of different domains are provided. Treebanks which are of the same language families are provided as well. Letting these treebanks help each other has been shown an effective way to improve parsing performance in both the cross-lingual-cross-domain parsing community and last year’s shared tasks (Ammar et al., 2016; Guo et al., 2015; 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)¹ to map *fasttext*² word embeddings (Bojanowski et al., 2016) of the source rich-resource language and target low-resource language into the same space. The transferred parser trained on the source language is used for the target low-resource language.

We conduct experiments on the development data to study the effects of ELMo, parser ensemble, and treebank concatenation. Experimental re-

¹https://github.com/Babylonpartners/fastText_multilingual

²<https://github.com/facebookresearch/fastText>

sults 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 significantly outperforms the others by a large margin.³

2 Deep Biaffine Parser

We based our system on the tagger and parser of Dozat et al. (2017). The core idea of the tagger and parser is using an 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

$$\mathbf{h}_i = \text{BiLSTM}(\mathbf{h}_0, (\mathbf{v}_1^{(word)}, \dots, \mathbf{v}_n^{(word)}))_i$$

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

$$\begin{aligned} \mathbf{h}_i^{(pos)} &= \text{MLP}^{(pos)}(\mathbf{h}_i) \\ \mathbf{s}_i^{(pos)} &= W \cdot \mathbf{h}_i^{(pos)} + \mathbf{b}^{(pos)} \\ y_i^{(pos)} &= \arg\max_j s_{i,j}^{(pos)} \end{aligned}$$

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

$$\begin{aligned} \mathbf{x}_i &= \mathbf{v}_i^{(word)} \oplus \mathbf{v}_i^{(tag)} \\ \mathbf{h}_i &= \text{BiLSTM}(\mathbf{h}_0, (\mathbf{x}_1, \dots, \mathbf{x}_n))_i \end{aligned}$$

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

$$\begin{aligned} \mathbf{s}_i^{(arc)} &= H^{(arc-head)} W^{(arc)} \mathbf{h}_i^{(arc-dep)} \\ &\quad + H^{(arc-head)} \mathbf{b}^{(arc)} \\ y^{(arc)} &= \arg\max_j s_{i,j}^{(arc)} \end{aligned}$$

where $\mathbf{h}_i^{(arc-dep)}$ is computed by feeding \mathbf{h}_i into an MLP and $H^{(arc-head)}$ is the stack of $\mathbf{h}_i^{(arc-head)}$ which is calculated in the same way with $\mathbf{h}_i^{(arc-dep)}$

but using another MLP. After getting the head $y^{(arc)}$ word, its relation with i -th word decided by calculating

$$\begin{aligned} \mathbf{s}_i^{(rel)} &= \mathbf{h}_{y^{(arc)}}^{T(rel-head)} \mathbf{U}^{(rel)} \mathbf{h}_i^{(rel-dep)} \\ &\quad + W^{(rel)} (\mathbf{h}_i^{(rel-dep)} \oplus \mathbf{h}_{y^{(arc)}}^{T(rel-head)}) \\ &\quad + \mathbf{b}^{(rel)}, \\ y^{(rel)} &= \arg\max_j s_{i,j}^{(rel)} \end{aligned}$$

where $\mathbf{h}^{(rel-head)}$ and $\mathbf{h}^{(rel-dep)}$ are calculated in the same way with $\mathbf{h}_i^{(arc-dep)}$ and $\mathbf{h}_i^{(arc-head)}$.

This decoding process can lead to cycles in the result. (Dozat et al., 2017) employed an iterative fixing methods on the cycles. We encourage the reader of this paper refer their paper for more training and decoding details.

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

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

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 γ as

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

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 $\text{BiLSTM}^{(LM)}$.

³<http://universaldependencies.org/conll18/results.html>

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 γ , which means

$$\mathbf{ELMo}_i = \sum_{j=0}^2 \mathbf{h}_{i,j}^{(LM)}.$$

In our preliminary experiments, using $\mathbf{h}_{i,0}^{(LM)}$ for \mathbf{ELMo}_i yields better performance on some treebanks. In our final submission, we decide using either $\sum_{j=0}^2 \mathbf{h}_{i,j}^{(LM)}$ or $\mathbf{h}_{i,0}^{(LM)}$ by their development performance.

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.

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 softmaxed output scores as

$$\mathbf{s}_i^{(rel)} = \frac{1}{3} \sum_{m=1}^3 \text{softmax}(\mathbf{s}_i^{(m,rel)}).$$

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*.⁴ 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.⁵ We show the experimental results of treebank concatenation in Section 9.3.

6 Low Resources Languages

In the shared task, 5 languages are presented with training set of less than 50 sentences. 4 languages do not even have any training data. It’s difficult to 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 embeddings using the same tokens shared by 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

⁴We opt out *cs*, *fi*, and *pl* because all the treebanks of these languages are relatively large which have more than 10K training sentences.

⁵We use *udpipe* for this part of experiments 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.

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

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

training data.⁶ For the treebank without any training data, we choose the source language according to their language family.⁷

Naija presents an exception for our method since it does not 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* instead.

7 Preprocessing

Besides improving the tagger and parser, we also consider the preprocessing as an important factor to the final performance try to improve it by using the state-of-the-art system for sentence segmentation, and developing our own word segmentor for languages whose tokenizations are non-trivial.

7.1 Sentence Segmentation

For some treebanks, sentence segmentation can be problematic since there is no explicitly sentence delimiters. [de Lhoneux et al. \(2017\)](#) and [Shao \(2017\)](#) presented a joint tokenization and sentence segmentation model⁸ 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.⁹ Using the Uppsala segmentor leads to a development improvement of 7.67 F-score in these treebanks over *udpipe* baseline and it was ranked first according to sentence segmentation in the final evaluation.

7.2 Tokenization for Chinese, Japanese, and Vietnamese

Tokenization is non-trivial for languages which do not have explicit word boundary markers, like

⁶We use *udpipe* for this test. When training the parser, the small set of target training data is also used.

⁷Thai does not have a treebank in the same family. We choose Chinese as source language because of geographical closeness and both these two languages are SVO in typology.

⁸<https://github.com/yanshao9798/segmenter/>, noted as Uppsala segmentor henceforth.

⁹We use Uppsala segmentor for *it_postwita*, *got_proiel*, *la_poroiel*, *cu_proiel*, *grc_proiel*, *sl_ssj*, *nl_lassysmall*, *fi_tdt*, *pt_bosque*, *da_dtd*, *id_gsd*, *el_gdt*, and *et_edt*.

Chinese, Japanese, and Vietnamese. We develop our own tokenizer for these three languages.¹⁰ Following [Che et al. \(2017\)](#) and [Zheng et al. \(2017\)](#), we model the tokenization as labeling the word boundary tag¹¹ on character and use features derived from large-scale unlabeled data to further improve the performance.¹² In addition to the pointwise mutual information (PMI), we also incorporate the character ELMo into our tokenizer. Embeddings of these features are concatenated along with a bigram character embeddings as input. These techniques lead to the best tokenization performance on all the related treebanks and the average improvement over *udpipe* baseline is 7.5 in tokenization F-score.¹³

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 and obtaining Thai words requires tokenization. To this remedy, we use the whitespace 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 by preserving the top 10% frequent words.

7.4 Lemmatization and Morphology Tagging

We did not make an effort on lemmatization and morphology tagging, but only use the baseline model. This lags our performance in the MLAS and BLEX evaluation, in which we were ranked 6th and 2nd correspondingly. However, since our method, especially incorporating ELMo, is not limited to particular task, we expect it to improve both the lemmatization and morphology tagging and achieve better MLAS and BLEX scores.

8 Implementation Details

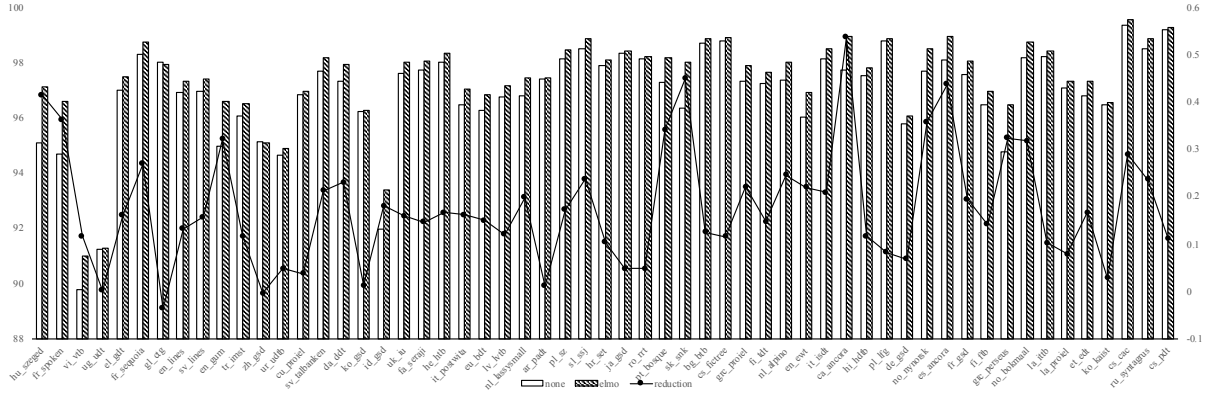
Pretrained Word Embeddings. We use the 100-dimensional pretrained word embeddings released by the shared task for the large languages. For the small languages and low-resource languages where cross-lingual transfer is required,

¹⁰note as SCIR segmentor henceforth.

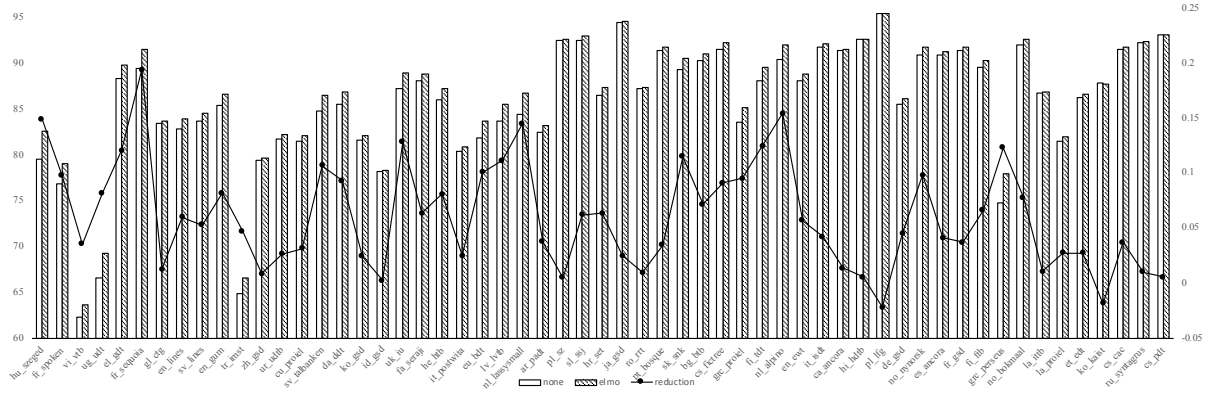
¹¹We use the BIES scheme.

¹²For Vietnamese where whitespaces occur both inter- and intra-words, we treat the whitespace-separated token as a character.

¹³*ja_gsd*, *ja_modern*, *vi_vtb*, and *zh_gsd*.



(a) The effects of ELMo on POS tagging



(b) The effects of ELMo on parsing

Figure 1: The effects of ELMo. Treebanks are sorted by the number of training sentences from left to right.

we use the 300-dimensional *fasttext* word embeddings. *fro_srcmf* presents the only exceptions and we use the French embeddings instead. For all the embeddings, we filter the top 10% frequent words.

ELMo. We use the same hyperparameter settings as Peters et al. (2018) for BiLSTM^(LM) and the character CNN. We train their parameters as training a bidirectional language model on a set of 20 million words data randomly sampled from the raw text released by the shared task for each language. Similar to Peters et al. (2018), we use the *sample softmax* technique to make training on large vocabulary feasible (Jean et al., 2015). However, we use a window of words¹⁴ surrounding the target word as negative samples and it shows better performance in our preliminary experiments. The training of ELMo on one language takes roughly 3 days on an NVIDIA P100 GPU.

Biaffine Parser. We use the same hyperparameter settings as Dozat et al. (2017). When trained

with ELMo, we use a dropout of 33% on the projected vectors.

SCIR Segmentor. We use a 50-dimensional character bigram embeddings. For the character ELMo whose input is a character, the language model predict next character in the same way with the word ELMo. The final model is an ensemble of five single segmentors.

Uppsala Segmentor. We use the default settings for the Uppsala segmentor and the final model is an ensemble of three single segmentors.

9 Results

9.1 Effects of ELMo

We study the effect of ELMo on the large treebanks and report the results of single tagger and parser with and without ELMo. Figure 1a shows the tagger results on the development set and Figure 1b 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

¹⁴The number of words within the window is 8192.

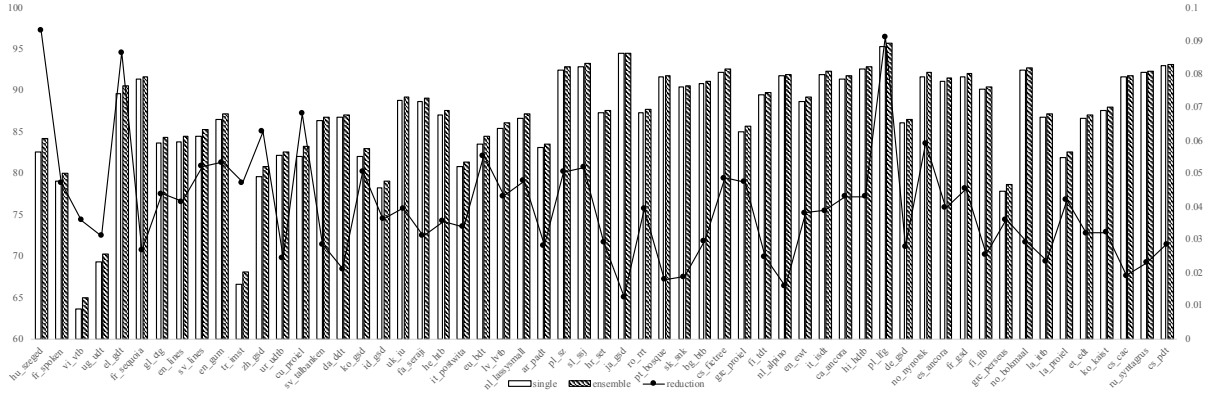


Figure 2: The effects of ensemble on parsing. Treebanks are sorted according to the number of training sentences from left to right.

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

<i>en</i>	<i>ewt</i>	<i>gum</i>	<i>lines</i>				<i>fr</i>	<i>gsd</i>	<i>sequoia</i>	<i>spoken</i>
# train	12.5	2.9	2.7				# train	14.6	2.2	1.2
single	88.75	86.52	83.86				single	91.64	91.44	79.06
concat.	88.74	85.65	85.30				concat.	91.44	90.51	81.99

Table 2: The development performance with cross-domain concatenation for languages which has multiple treebanks. *single* means training the parser on it own treebank without concatenation. *# train* shows the number of training sentences in the treebank measured in thousand.

reduction is 17.83%. 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. The line in Figure 1a and Figure 1b shows the error reduction from using ELMo on each treebank. However, no clear relation is revealed between the treebank size and the gains using ELMo.

9.2 Effects of Ensemble

We also test the effect of ensemble and show the results in Figure 2. 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 an effective way to improve the parsing performance. The relationship

between gains using ensemble and treebank size is also studied in this figure and the trend is that small treebank benefit more from ensemble. We address this to the fact the ensemble improve the model’s generalization ability in which the parser trained on small treebank is weak due to overfitting.

9.3 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 does not improve the performance or even worsen the results.

For the small treebanks which do not have development set, the 5 fold cross validation results

<i>gl</i>	<i>treegal</i>		<i>la</i>	<i>perseus</i>		<i>no</i>	<i>nynorsklika</i>		<i>ru</i>	<i>taiga</i>		<i>sl</i>	<i>sst</i>
# train	0.6		# train	1.3		# train	0.3		# train	0.9		# train	2.1
<i>treegal</i>	66.71		<i>perseus</i>	44.05		<i>nynorsklika</i>	51.05		<i>taiga</i>	54.70		<i>sst</i>	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 of treebank which does not have development set.

	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_syntagrus	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.

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

Table 5: The effects of improved preprocessing on the parsing performance. The first block shows the effects of sentence segmentation improvement. Δ -sent. means the sentence segmentation F-score difference between Uppsala segmentor and *udpipe*. The second block shows the effects of word segmentation improvement. Δ -word means the word segmentation in F-score difference between SCIR segmentor and *udpipe*.

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.

9.4 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.

9.5 Parsing Strategies and Test Evaluation

Using the development set and cross validation, we choose the best model and data combination and the choices are shown in Table 6 along with the test evaluation. From this table, we can see that our system gains more improvements when both ELMo and parser ensemble are used. For some treebanks, concatenation also contributes to the improvements and the Japanese, Vietnamese, and Chinese parsing clearly benefit from better word segmentation,

We also report the time and memory consumption. A full run over the 82 test sets on the TIRA virtual machine (Potthast et al., 2014) takes about 40 hours and consumes about 4G RAM memory.

10 Conclusion

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

11 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

This work was supported by the National Key Basic Research Program of China via grant 2014CB340503 and the National Natural Science Foundation of China (NSFC) via grant 61300113 and 61632011.

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.
- 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*.
- Sébastien Jean, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. 2015. *On using very large target vocabulary for neural machine translation*. In *Proc. of ACL*. <http://www.aclweb.org/anthology/P15-1001>.
- 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/SUPPLYTHENPERMANENTIDHERE!>
- 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*.
- Yan Shao. 2017. [Cross-lingual word segmentation and morpheme segmentation as sequence labelling](#). *CoRR* abs/1709.03756. <http://arxiv.org/abs/1709.03756>.
- 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 Luotolahti, Sampo Pyysalo, Slav Petrov, Martin Potthast, 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 Leung, Marie-Catherine de Marneffe, Manuela Sanguinetti, Maria Simi, Hiroshi Kanayama, Valeria de Paiva, Kira Droganova, Héctor Martínez Alonso, Çağrı Çöltekin, Umut Sulubacak, Hans Uszkoreit, 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 Strnadova, Esha Banerjee, Ruli Manurung, Antonio Stella, Atsuko Shimada, Sookyoung Kwak, Gustavo Mendonça, Tatiana Lando, Rattima Nitisoroj, 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*.

<i>ltcode</i>	sent+tokenize	tagger	parser	LAS	Other LAS	rank	diff.
af.afribooms	udpipe: self	biaffine (none): self	biaffine (none)*3: self	85.47	85.45	1	0.02
ar.padt	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	73.63	77.06	2	-3.43
bg.btb	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	91.22	90.41	1	0.81
br.keb	udpipe: self	biaffine.trans: self+ga.idt	biaffine.trans*3: self+ga.idt	8.54	38.64	21	-30.1
bxr.bdt	udpipe: self	biaffine.trans: self+hi.hdtb	biaffine.trans*3: self+hi.hdtb	15.44	19.53	6	-4.09
ca.ancora	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	91.61	90.82	1	0.79
cs.cac	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	91.61	91.00	1	0.61
cs.fictree	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	92.02	91.83	1	0.19
cs.pdt	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	91.68	90.57	1	1.11
cs.pud	udpipe: cs.pdt	biaffine (h_0): cs.pdt	biaffine ($h_{0,1,2}$)*3: cs.pdt	86.13	85.35	1	0.78
cu.proiel	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	74.29	75.73	3	-1.44
da.ddt	uppsala: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	86.28	84.88	1	1.40
de.gsd	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	80.36	79.03	1	1.33
el.gdt	uppsala: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	89.65	89.59	1	0.06
en.ewt	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	84.57	84.02	1	0.55
en.gum	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	84.42	85.05	2	-0.63
en.lines	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self+en.ewt+en.gum	81.97	81.44	1	0.53
en.pud	udpipe: en.ewt	biaffine (h_0): en.ewt	biaffine ($h_{0,1,2}$)*3: en.ewt	87.73	87.89	2	-0.16
es.ancora	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	90.93	90.47	1	0.46
et.edt	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	85.35	84.15	1	1.20
eu.bdt	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self	84.22	83.13	1	1.09
fa.seraji	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	88.11	86.18	1	1.93
fi.ftb	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self	88.53	87.86	1	0.67
fi.pud	udpipe: fi.tdt	biaffine (h_0): fi.tdt	biaffine (h_0)*3: fi.tdt	90.23	89.37	1	0.86
fi.tdt	uppsala: self	biaffine (h_0): self	biaffine (h_0)*3: self	88.73	87.64	1	1.09
fo.oft	udpipe: no.bokmaal	biaffine.trans: no.bokmaal	biaffine.trans*3: self	44.05	49.43	1	-5.38
fr.gsd	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	86.89	86.46	1	0.43
fr.sequoia	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	89.65	89.89	1	-0.24
fr.proiel	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self+fr.gsd+fr.sequoia	75.78	74.31	1	1.47
fro.srcmf	udpipe: self	biaffine (none): self	biaffine (none)*3: self	87.07	87.12	2	-0.05
ga.idt	udpipe: self	biaffine (none): self	biaffine (none)*3: self	68.57	70.88	5	-2.31
gl.ctg	udpipe: self	biaffine (none): self	biaffine (none)*3: self	82.35	82.76	2	-0.41
gl.treel	udpipe: self	biaffine (none): self	biaffine (none)*3: self	72.88	74.25	4	-1.37
got.proiel	uppsala: self	biaffine (none): self	biaffine (none)*3: self	69.26	69.55	3	-0.29
grc.perseus	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	79.39	74.29	1	5.10
grc.proiel	uppsala: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	79.25	76.76	1	2.49
he.hdtb	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	67.05	76.09	3	-9.04
hi.hdtb	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self	92.41	91.75	1	0.66
hr.set	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	87.36	86.76	1	0.60
hsb.ufal	udpipe: self	biaffine.trans: self+pl.lfg	biaffine.trans*3: self+pl.lfg	37.68	46.42	4	-8.74
hu.szeged	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	82.66	79.47	1	3.19
hy.armtdp	udpipe: self	biaffine.trans: self+et.edt	biaffine.trans*3: self+et.edt	33.90	37.01	3	-3.11
id.gsd	uppsala: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	80.05	79.13	1	0.92
it.isdt	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	92.00	91.47	1	0.53
it.postwita	uppsala: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self+it.isdt	79.39	78.62	1	0.77
ja.gsd	udpipe+scir: self	biaffine (h_0): self	biaffine (h_0)*3: self	83.11	79.97	1	3.14
ja.modern	udpipe+scir: ja.gsd	biaffine (h_0): ja.gsd	biaffine (h_0)*3: ja.gsd	26.58	28.33	4	-1.75
kk.ktb	udpipe: self	biaffine.trans: self+tr.imst	biaffine.trans*3: self+tr.imst	23.92	31.93	10	-8.01
kmr.mg	udpipe: self	biaffine.trans: self+fa.seraji	biaffine.trans*3: self+fa.seraji	26.26	30.41	5	-4.15
ko.gsd	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	85.14	84.31	1	0.83
ko.kaist	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	86.91	86.84	1	0.07
la.ittb	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	87.08	86.54	1	0.54
la.perseus	udpipe: self	biaffine (h_0): self+la.proiel	biaffine ($h_{0,1,2}$)*3: self+la.proiel	72.63	68.07	1	4.56
la.proiel	uppsala: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	73.61	71.76	1	1.85
lv.lvtb	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	83.97	81.85	1	2.12
nl.alpino	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self+nl.lassysmall	89.56	87.49	1	2.07
nl.lassysmall	uppsala: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self+nl.alpino	86.84	84.27	1	2.57
no.bokmaal	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	91.23	90.37	1	0.86
no.nynorsk	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	90.99	89.46	1	1.53
no.nynorskliia	udpipe: self	biaffine (h_0): self+no.nynorsk	biaffine ($h_{0,1,2}$)*3: self+no.nynorsk	70.34	68.71	1	1.63
pcm.nsc	udpipe: en.ewt	biaffine (h_0): en.ewt	biaffine ($h_{0,1,2}$)*3: en.ewt	24.48	30.07	2	-5.59
pl.lfg	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	94.86	94.62	1	0.24
pl.sz	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	92.23	91.59	1	0.64
pt.bosque	uppsala: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	87.61	87.81	3	-0.20
ro.rtt	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	86.87	86.33	1	0.54
ru.syntagrus	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self	92.48	91.72	1	0.76
ru.taiga	udpipe: self	biaffine ($h_{0,1,2}$): self+ru.syntagrus	biaffine ($h_{0,1,2}$)*3: self+ru.syntagrus	71.81	74.24	3	-2.43
sk.snk	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine ($h_{0,1,2}$)*3: self	88.85	87.59	1	1.26
sl.ssj	uppsala: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	91.47	91.26	1	0.21
sl.sst	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self+sl.ssj	61.39	58.12	1	3.27
sme.giella	udpipe: self	biaffine (none): self	biaffine ($h_{0,1,2}$)*3: self	69.06	69.87	3	-0.81
sr.set	udpipe: self	biaffine (none): self	biaffine ($h_{0,1,2}$)*3: self	88.33	88.66	3	-0.33
sv.lines	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self+sv.talbanken	84.08	81.97	1	2.11
sv.pud	udpipe: sv.lines	biaffine (h_0): sv.lines	biaffine (h_0)*3: sv.lines+sv.talbanken	80.35	79.71	1	0.64
sv.talbanken	udpipe: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self+sv.lines	88.63	86.45	1	2.18
th.pud	thai	biaffine.trans: zh.gsd	biaffine.trans*3: zh.gsd	0.64	13.70	14	-13.06
tr.imst	udpipe: self	biaffine (h_0): self	biaffine ($h_{0,1,2}$)*3: self	66.44	64.79	1	1.65
ug.udt	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	67.05	65.23	1	1.82
uk.iu	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self+ru.syntagrus	88.43	85.16	1	3.27
ur.udtb	udpipe: self	biaffine (h_0): self	biaffine (h_0)*3: self	83.39	82.15	1	1.24
vi.vtb	udpipe+scir: self	biaffine ($h_{0,1,2}$): self	biaffine (h_0)*3: self	55.22	47.41	1	7.81
zh.gsd	udpipe+scir: self	biaffine (none): self	biaffine ($h_{0,1,2}$)*3: self	76.77	71.04	1	5.73

Table 6: The strategies used in the final submission. The *toolkit* and *model* are separated by colon. *uppsala* denotes the Uppsala segmentor; *scir* denotes our segmentor for Japanese, Vietnamese, and Chinese; *biaffine* denotes the biaffine tagger and parser; *biaffine.trans* denotes our transfer parser for low-resource languages. h_0 or $h_{0,1,2}$ in the brackets denotes the ELMo used to train the model where h_0 means using $\mathbf{h}_{i,0}^{(LM)}$ for ELMo and $h_{0,1,2}$ means using $\sum_{j=0}^2 \mathbf{h}_{i,j}^{(LM)}$. *self* notes that the model is trained with the treebank itself. If the model field is not filled with *self*, the model is trained through treebank concatenation. We also show the test LAS and the difference against the other best system.