



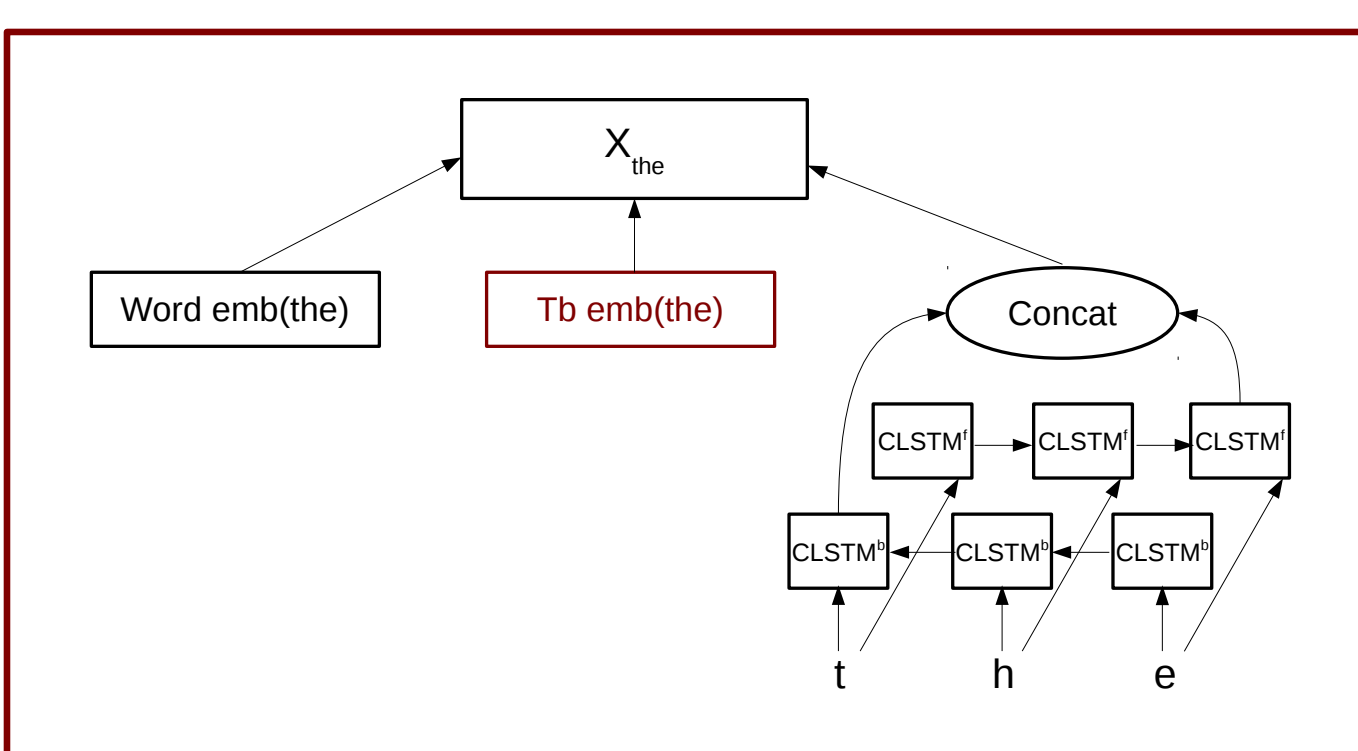
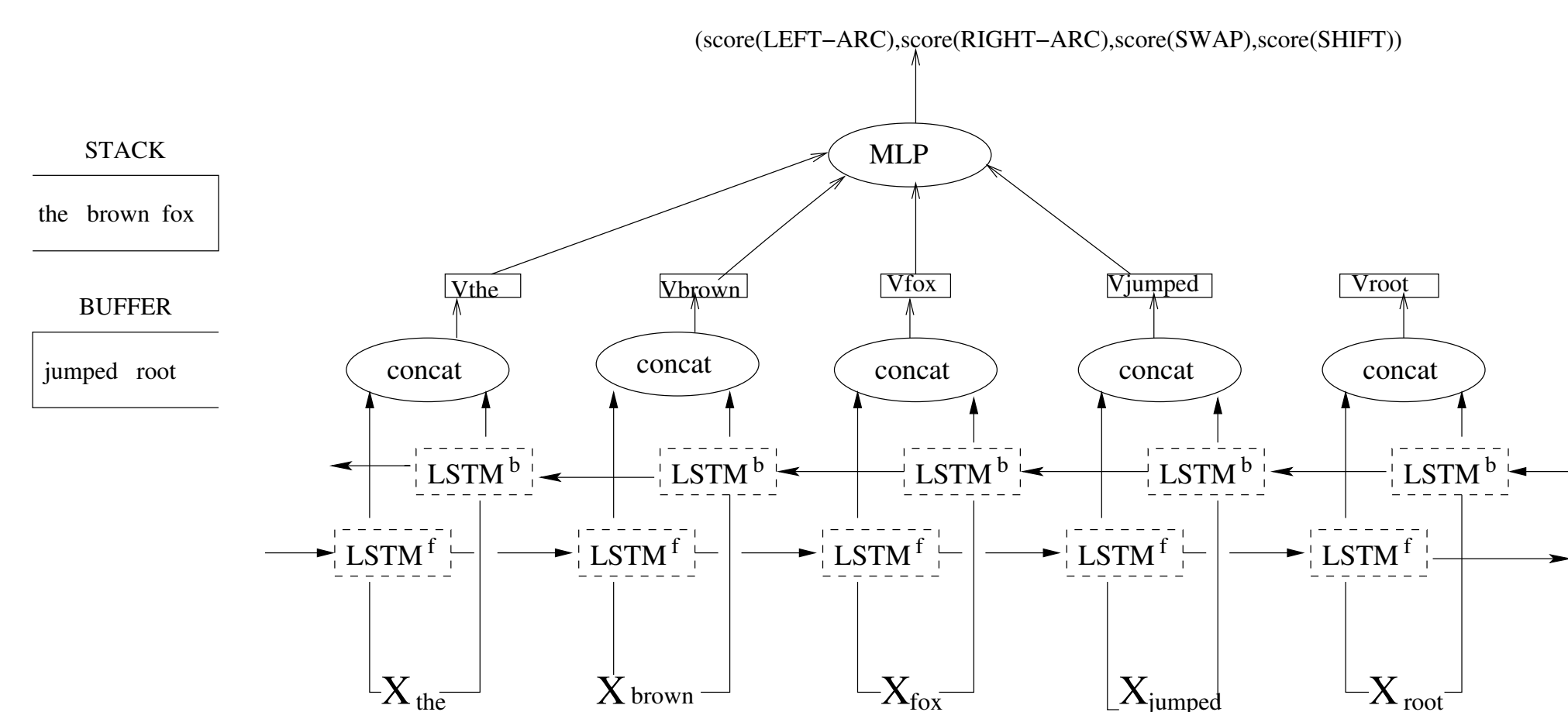
# Parser Training with Heterogeneous Treebanks

Sara Stymne, Miryam de Lhoneux, Aaron Smith and Joakim Nivre

## Introduction

- **Problem:** How can we improve parsing when there are several, potentially heterogeneous treebanks for a language?
- Treebank diversity
  - Annotation scheme
  - Language variant
  - Spoken/written language
  - Genres and domains
  - Treebank size
  - Annotation quality and consistency
- **This work:**
  - Investigate previously proposed strategies
  - Introduce treebank embeddings

## System Architecture



Dims word emb	100
Dims char emb	12
Dims treebank emb	12
Word LSTM dims	125
Char LSTM dims	50
LSTM dropout	0.33
Word dropout ( $\alpha$ )	0.25
Epochs	30
Epochs fine-tuning	10

## Strategies

### Single

- One model per treebank
- + Simple
- Does not take advantage of all data
- Separate models for each treebank

### Concatenation

- One model per language, on concatenated data
- + Simple
- Does not take treebank differences into account
- + A single model per language

### Concatenation + fine-tuning

- Fine-tune a different model for each treebank, based on the concatenation (Che et al., 2017, Shi et al., 2017)
- Needs more training than previous models
- Separate models for each treebank
- + Takes treebank differences into account

### Treebank embeddings

- Train a single model per language, but use a treebank embedding to represent the treebank each word comes from. Similar to language embeddings (Ammar et al., 2016)
- + Simple
- + Takes treebank differences into account
- + A single model per language

### Other approaches (not in this paper)

- 1-hot treebank representation: similar to our approach, but with 1-hot representation rather than embedding (Lim & Poibeau, 2017).
- Adversarial learning: combine treebank specific models with a joint model where treebank identification is an adversarial task (Sato et al., 2017). Effective, especially on small treebanks, but more complicated than our model.

## Parsing Unseen Data

When parsing unseen data, we need to choose an existing treebank: **proxy treebank**

- Single: the treebank used to train a model
- Concatenation: N/A
- Concatenation + fine-tuning: the treebank used for fine-tuning
- Treebank embeddings: the treebank embedding to use in the model

## Experiments

- Universal dependencies version 2.1
- Standardized annotation scheme, but still many differences
- 9 languages:
  - at least 2 training treebanks
  - test set without training data (PUD)

## Results

Language	Treebank	Size	Same treebank test set				PUD (unseen) test set			
			SINGLE	CONCAT	C+FT	TB-EMB	SINGLE	CONCAT	C+FT	TB-EMB
Czech	PDT	68495	86.7	87.5 <sup>+</sup>	<b>88.3<sup>*</sup></b>	87.2 <sup>+</sup>	<b>81.7</b>		81.6	81.2
	CAC	23478	86.0	87.8 <sup>+</sup>	88.1 <sup>+</sup>	<b>88.5<sup>+</sup></b>	75.0	<b>81.7</b>	81.3	81.1
	FicTree	10160	84.3	89.3 <sup>+</sup>	<b>89.5<sup>+</sup></b>	89.2 <sup>+</sup>	66.1		79.8	80.3
	CLTT	860	72.5	86.2 <sup>+</sup>	<b>86.9<sup>+</sup></b>	86.0 <sup>+</sup>	42.1		80.8	80.9
English	EWT	12543	82.2	82.1	82.5	<b>83.0</b>	80.7		<b>81.7<sup>*</sup></b>	<b>81.9<sup>*</sup></b>
	LinES	2738	72.1	76.7 <sup>+</sup>	<b>77.3<sup>+</sup></b>	<b>77.3<sup>+</sup></b>	62.6	80.0	75.9	74.5
	ParTUT	1781	80.5	83.5 <sup>+</sup>	85.4 <sup>+</sup>	<b>85.7<sup>+</sup></b>	68.0		78.1	76.9
Finnish	FTB	14981	76.4 <sup>×</sup>	74.4	80.1 <sup>*</sup>	<b>80.6<sup>*</sup></b>	46.7		54.6	53.1
	TDT	12217	78.1 <sup>×</sup>	70.6	<b>80.6<sup>*</sup></b>	80.3 <sup>*</sup>	78.6 <sup>×</sup>	73.0	<b>81.3<sup>*</sup></b>	<b>80.9<sup>*</sup></b>
French	FTB	14759	83.2	83.2	83.9 <sup>*</sup>	<b>84.1<sup>*</sup></b>	72.0		76.7	74.1
	GSD	14554	84.5	84.1	85.3	<b>85.6<sup>×</sup></b>	79.1	79.4	<b>80.2<sup>*</sup></b>	<b>80.3<sup>*</sup></b>
	Sequoia	2231	84.0	86.0 <sup>+</sup>	<b>89.8<sup>*</sup></b>	89.1 <sup>*</sup>	69.5		78.1	77.6
Italian	ParTUT	803	79.8	80.5	89.1 <sup>*</sup>	<b>90.3<sup>*</sup></b>	63.4		78.8	77.5
	ISDT	12838	87.7	<b>87.9</b>	87.7	87.6	85.4		85.7	86.0
Portuguese	PoSTWITA	2808	71.4	76.7 <sup>+</sup>	76.8 <sup>+</sup>	<b>77.0<sup>+</sup></b>	68.5	86.0	85.7	85.3
	ParTUT	1781	83.4	89.2 <sup>+</sup>	<b>89.3<sup>+</sup></b>	88.8 <sup>+</sup>	77.4		85.8 <sup>+</sup>	<b>86.1<sup>+</sup></b>
	GSD	9664	88.3	87.3	89.0 <sup>*</sup>	<b>89.1<sup>*</sup></b>	74.0		75.2	74.9
Russian	Bosque	8331	84.7	84.2	86.2 <sup>×</sup>	<b>86.3<sup>*</sup></b>	75.2	76.8 <sup>+</sup>	77.5 <sup>+</sup>	<b>77.6<sup>+</sup></b>
	SynTagRus	48814	90.2 <sup>×</sup>	89.4	<b>90.4<sup>×</sup></b>	<b>90.4<sup>×</sup></b>	66.0		66.3	66.4
Spanish	GSD	3850	74.7 <sup>×</sup>	73.4	79.8 <sup>*</sup>	<b>80.8<sup>*</sup></b>	70.1 <sup>×</sup>	68.7	<b>77.6<sup>*</sup></b>	<b>78.0<sup>*</sup></b>
	AnCora	14305	87.2 <sup>×</sup>	86.2	87.5 <sup>×</sup>	<b>87.6<sup>×</sup></b>	75.2	79.9	77.7	76.4
Swedish	GSD	14187	84.7	83.0	85.8 <sup>×</sup>	<b>86.2<sup>*</sup></b>	79.8		80.8 <sup>+</sup>	<b>80.9<sup>*</sup></b>
	Talbanken	4303	79.6	79.1	80.2	<b>80.6<sup>×</sup></b>	70.3	72.0 <sup>+</sup>	<b>73.2<sup>*</sup></b>	<b>73.6<sup>*</sup></b>
Average	LinES	2738	74.3	76.8	<b>77.3<sup>+</sup></b>	<b>77.1<sup>+</sup></b>	64.0		70.0	69.0
	Average		81.4	82.7 <sup>+</sup>	<b>84.9<sup>*</sup></b>	<b>84.9<sup>*</sup></b>	77.9	77.5	80.0 <sup>*</sup>	<b>80.1<sup>*</sup></b>

<sup>+</sup> significantly better than SINGLE    <sup>×</sup> significantly better than CONCAT    <sup>\*</sup> significantly better than SINGLE+CONCAT

## Conclusion

- Combining treebanks is beneficial, especially for small treebanks
- Treebank embeddings successful
  - At least on par with other methods
  - Simple model
  - Works for many different scenarios
- Choice of proxy treebank very important

## References

Waleed Ammar, George Mulcaire, Miguel Ballesteros, Chris Dyer, and Noah Smith. 2016. Many languages, one parser. *TACL*, 4:431444.

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. *CoNLL 2017*.

KyungTae Lim and Thierry Poibeau. 2017. A system for multilingual dependency parsing based on bidirectional LSTM feature representations. *CoNLL 2017*.

Motoki Sato, Hitoshi Manabe, Hiroshi Noji, and Yuji Matsumoto. 2017. Adversarial training for cross-domain universal dependency parsing. *CoNLL 2017*.

Tianze Shi, Felix G. Wu, Xilun Chen, and Yao Cheng. 2017. Combining global models for parsing universal dependencies. *CoNLL 2017*.