

Top-down Tree Long Short-Term Memory Networks

Xingxing Zhang, Liang Lu, Mirella Lapata

School of Informatics, University of Edinburgh

12th June, 2016

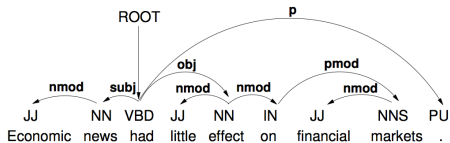
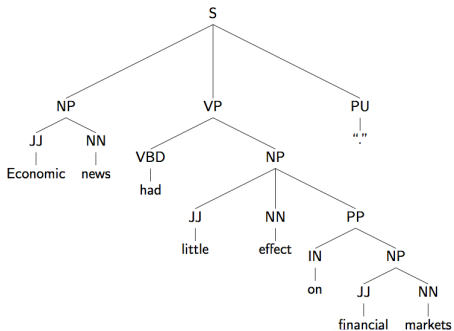
Sequential Language Models

$$P(S = w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{1:i-1}) \quad (1)$$

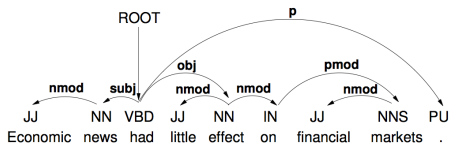
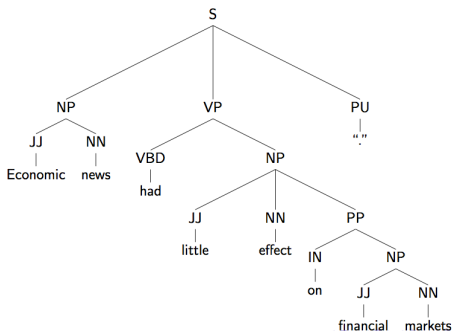
- State of the Art
 - based on Long Short Term Memory Network Language Model (Hochreiter and Schmidhuber, 1997; Sundermeyer et al., 2012)
 - **Billion word benchmark results** reported in Jozefowicz et al., (2016)

Models	PPL
KN5	67.6
LSTM	30.6
LSTM+CNN INPUTS	30.0

Will tree structures help LMs?



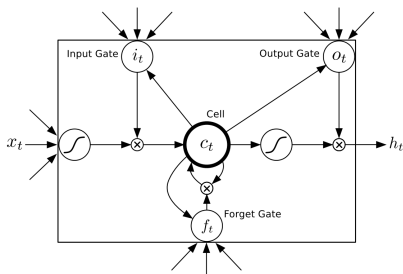
Will tree structures help LMs?



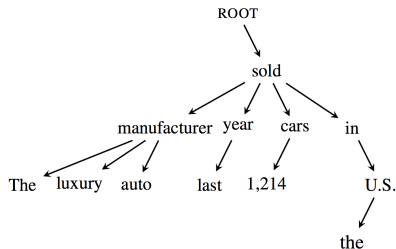
- Probably yes

- LMs based on Constituency Parsing (Chelba and Jelinek, 2000; Roark, 2001; Charniak, 2001)
- LMs based on Dependency Parsing (Shen et al., 2008; Zhang, 2009; Sennrich, 2015)

LSTMs + Dependency Trees = TreeLSTMs



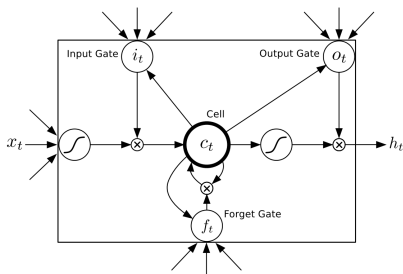
+



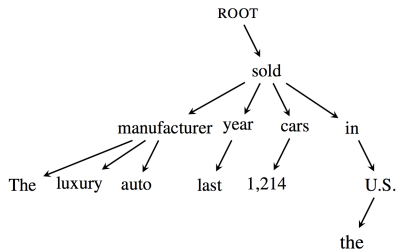
- Why?

- Sentence Length N v.s. Tree Height $\log(N)$

LSTMs + Dependency Trees = TreeLSTMs



+



- Why?
 - Sentence Length N v.s. Tree Height $\log(N)$
- How?
 - Top-down Generation
 - Breadth-first search
 - reminiscent of Eisner (1996)

Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.
ROOT

Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.

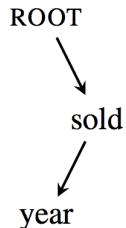
ROOT



sold

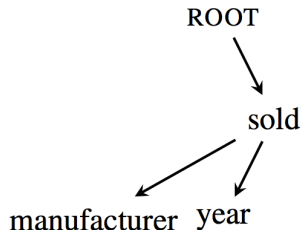
Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.



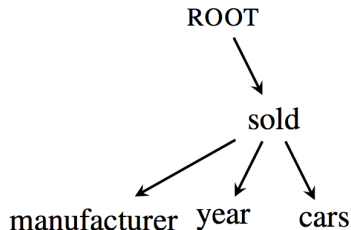
Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.



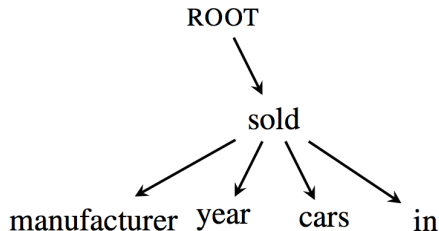
Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.



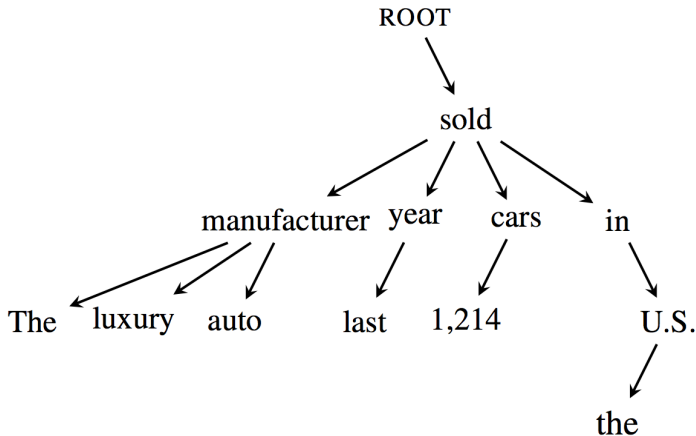
Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.



Generation Process (Unlabeled Trees)

The luxury auto manufacturer last year sold 1,214 cars in the U.S.



Tree LSTM

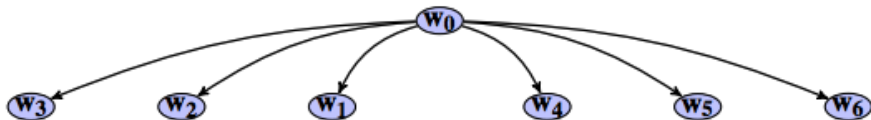
$$P(S) = \prod_{i=1}^n P(w_i | w_{1:i-1}) \quad (2)$$



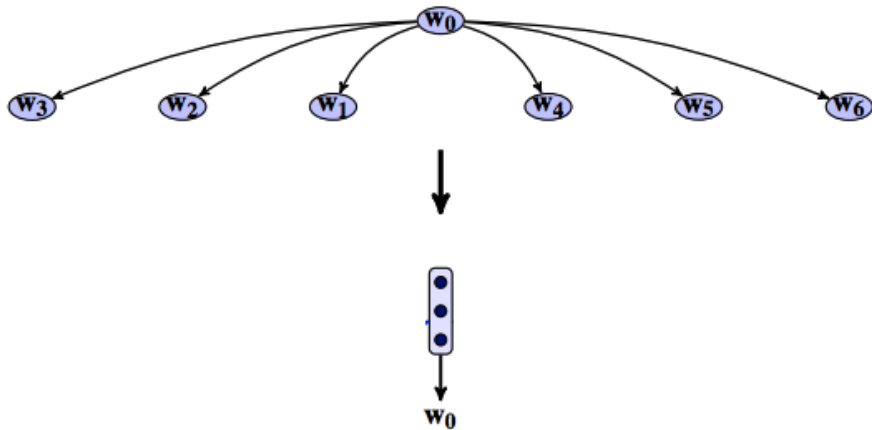
$$P(S|T) = \prod_{w \in \text{BFS}(T) \setminus \text{ROOT}} P(w | \mathcal{D}(w)) \quad (3)$$

- $\mathcal{D}(w)$ is the *Dependency Path* of w .
- $\mathcal{D}(w)$ is a generated sub-tree.
- Works on **projective** and **unlabeled** dependency trees.

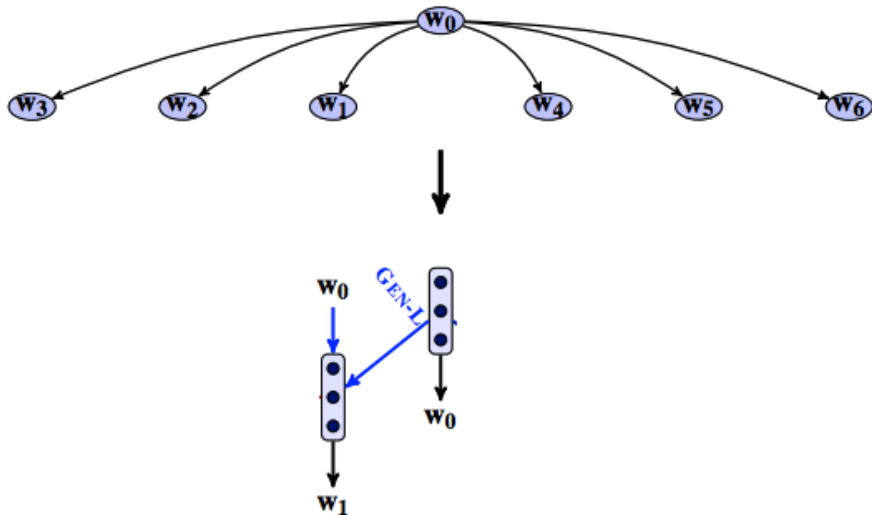
Tree LSTM



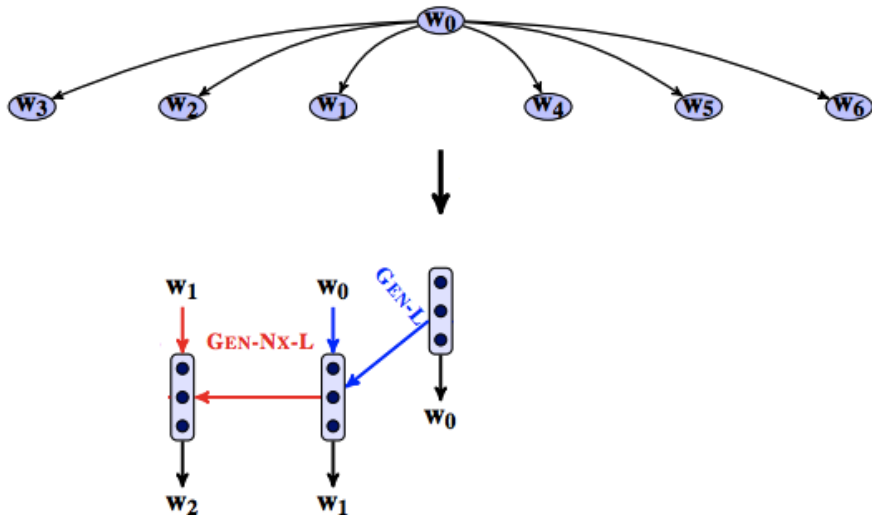
Tree LSTM



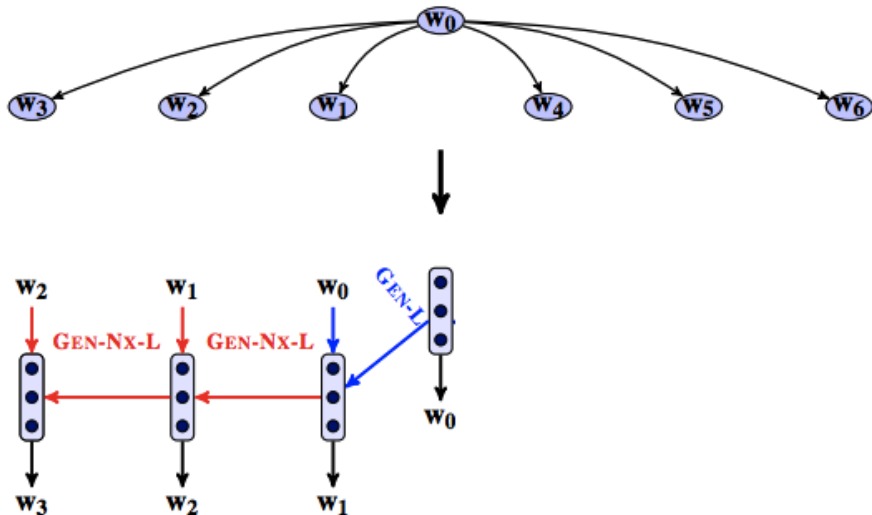
Tree LSTM



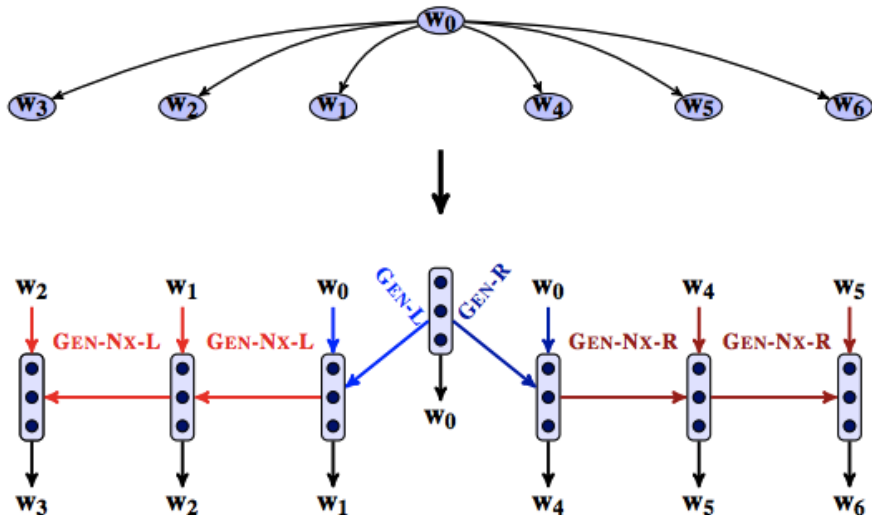
Tree LSTM



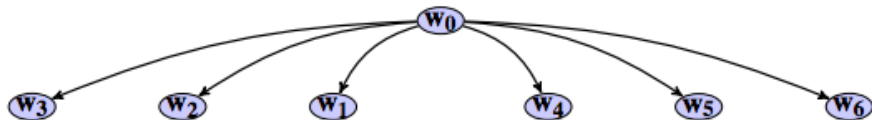
Tree LSTM



Tree LSTM

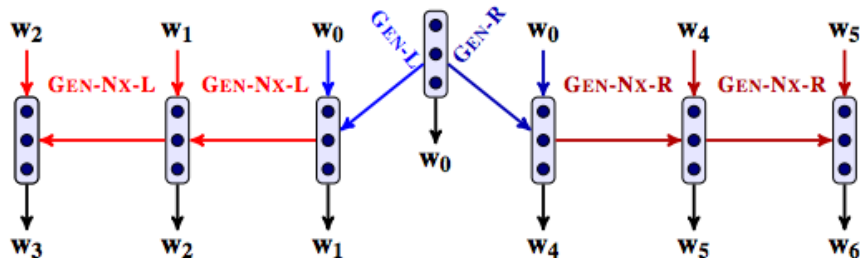


Tree LSTM

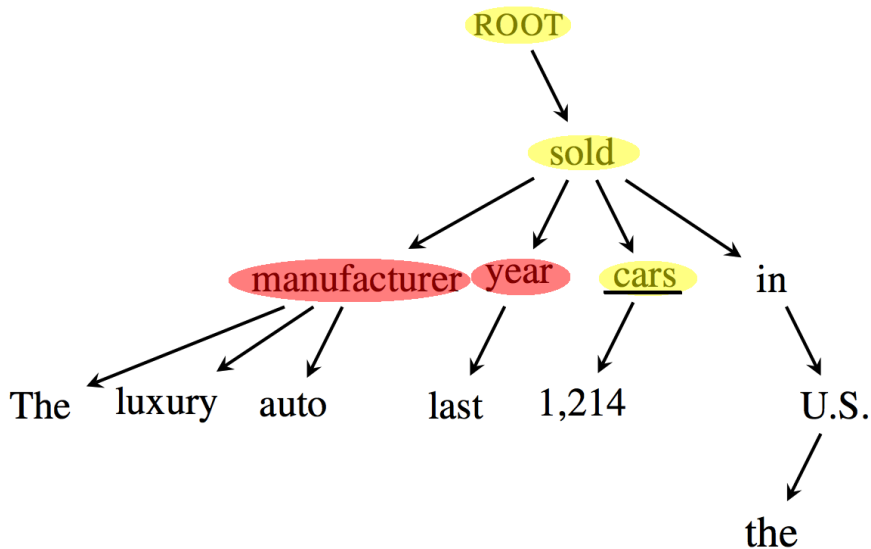


Generated by four LSTMs

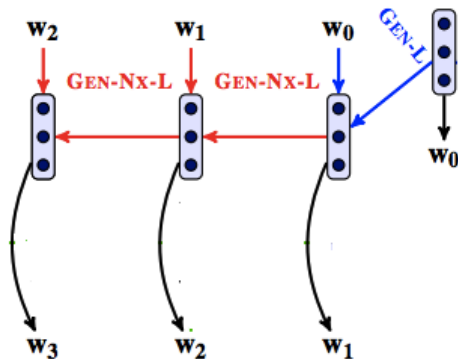
with tied W_e and tied W_{ho}



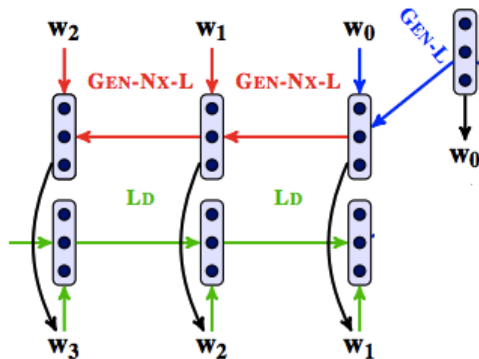
One Limitation of Tree LSTM



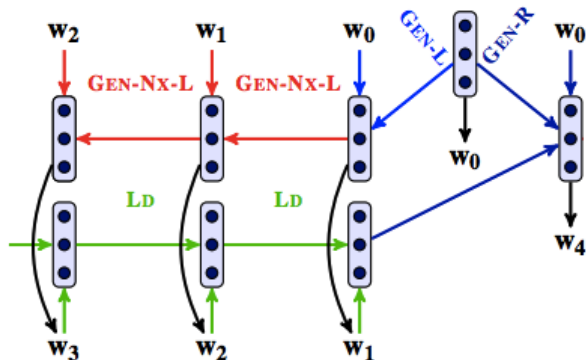
Left Dependent Tree LSTM



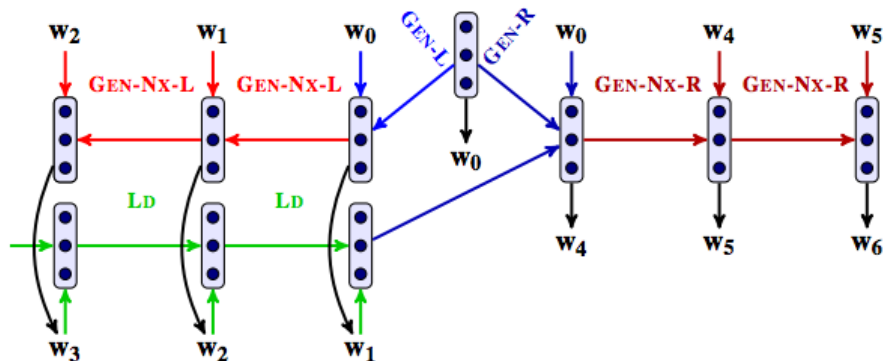
Left Dependent Tree LSTM



Left Dependent Tree LSTM



Left Dependent Tree LSTM

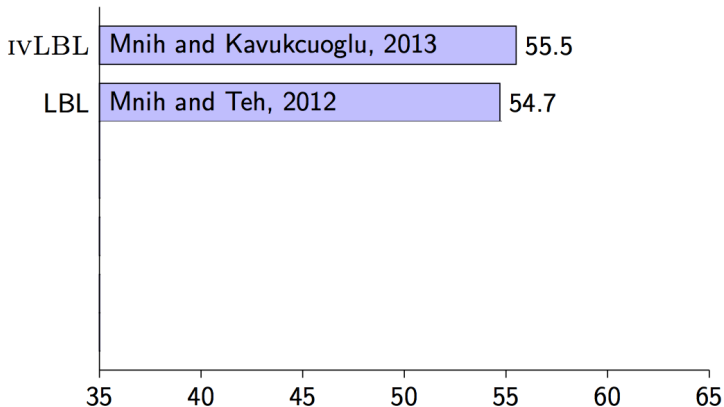


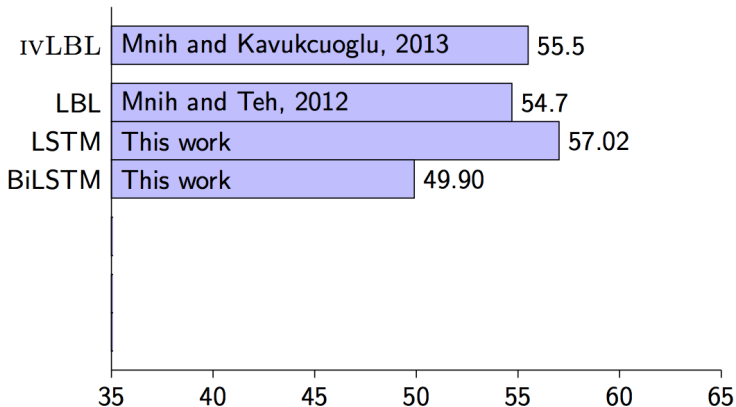
Experiments

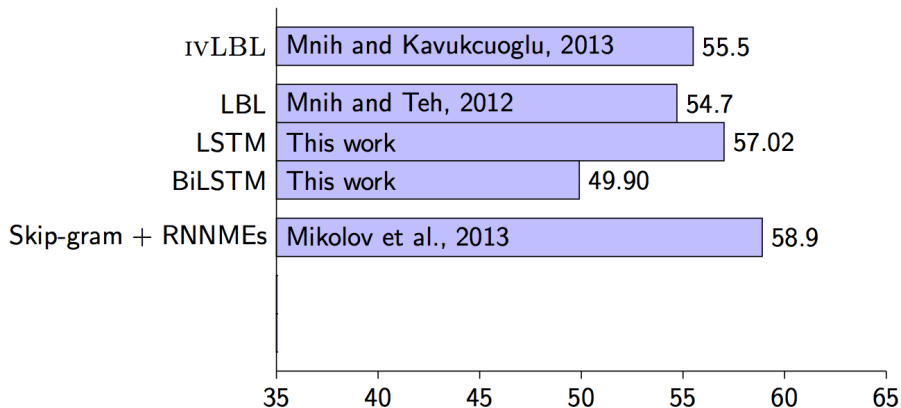
MSR Sentence Completion Challenge

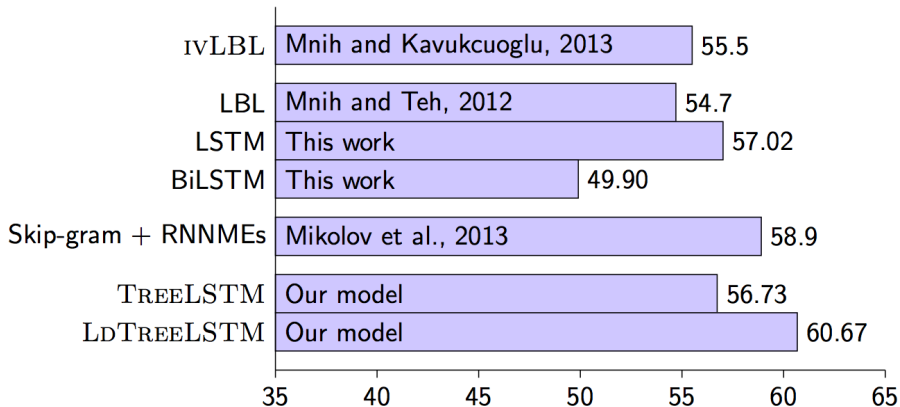
1) I have seen it on him , and could _____ to it.
a) write b) migrate c) climb d) swear e) contribute

- Training set: 49 million words (around 2 million sentences)
- development set: 4000 sentences
- test set: 1040 completion questions.





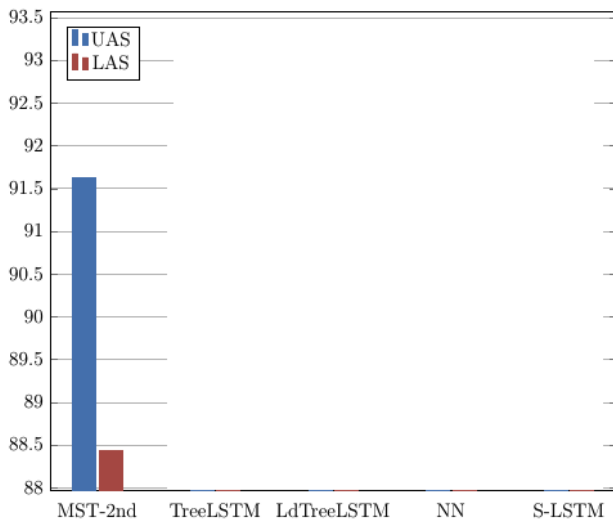




Dependency Parsing Reranking

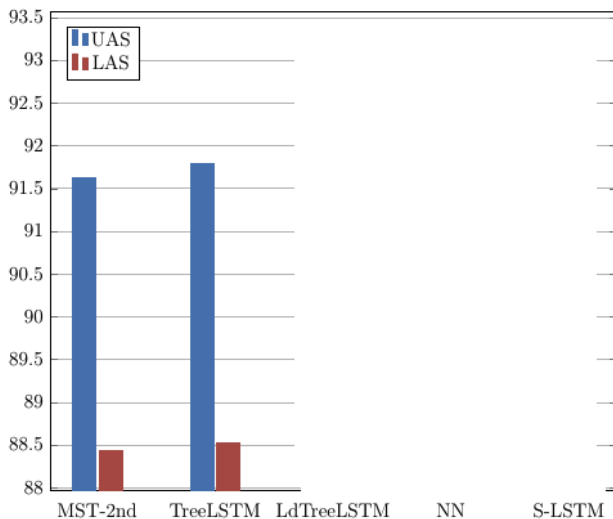
- Rerank 2nd Order MSTParser (McDonald and Pereira, 2006)
- We train TreeLSTM and LdTreeLSTM as language models.
- We only use words as input features; POS tags, dependency labels or composition features are not used.

Dependency Parsing Reranking



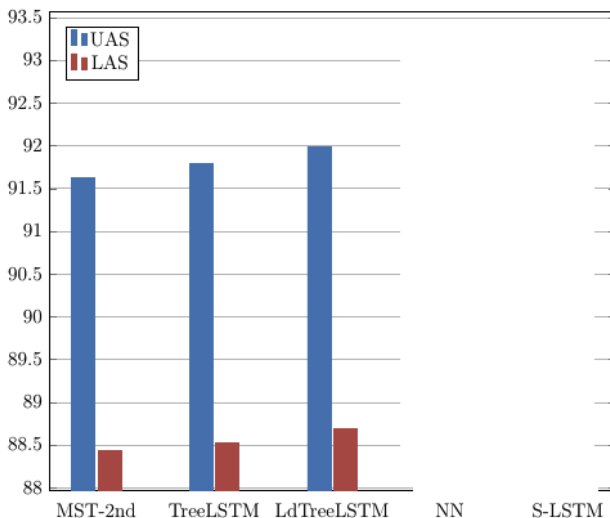
NN: Chen & Manning, 2014; S-LSTM: Dyer et al., 2015

Dependency Parsing Reranking



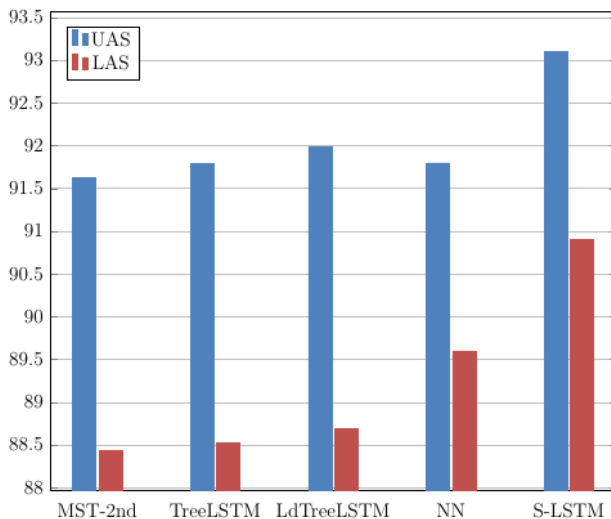
NN: Chen & Manning, 2014; S-LSTM: Dyer et al., 2015

Dependency Parsing Reranking



NN: Chen & Manning, 2014; S-LSTM: Dyer et al., 2015

Dependency Parsing Reranking

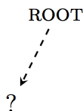


NN: Chen & Manning, 2014; S-LSTM: Dyer et al., 2015

Tree Generation

Four binary classifiers:

- Add Left? No!



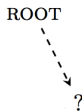
Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation

Four binary classifiers:

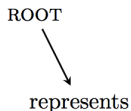
- Add Right? Yes!



Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation



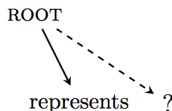
Four binary classifiers:

- Add Right? Yes!

Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation



Four binary classifiers:

- Add Next Right? No!

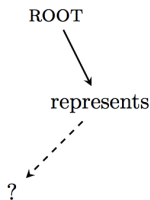
Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation

Four binary classifiers:

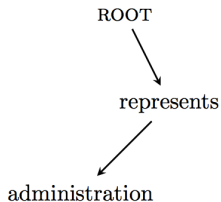
- Add Left? Yes!



Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation



Four binary classifiers:

- Add Left? Yes!

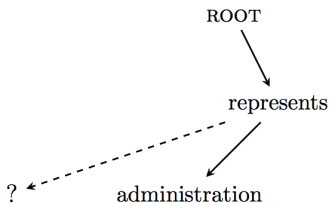
Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation

Four binary classifiers:

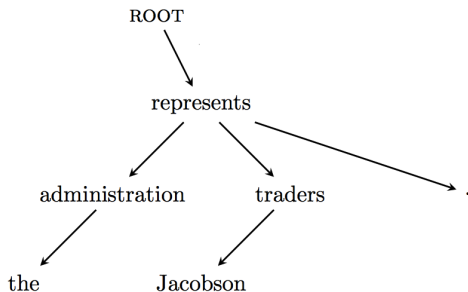
- Add Next Left? No!



Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation



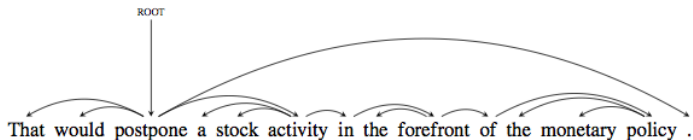
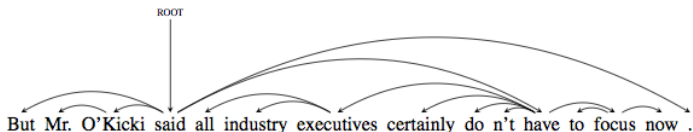
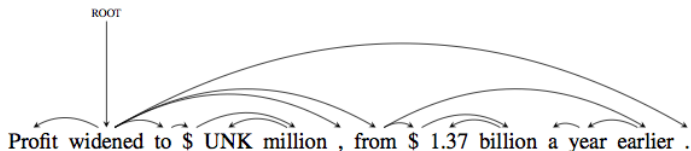
Four binary classifiers:

- Add Left?
- Add Right?
- Add Next Left?
- Add Next Right?

Features: hidden states and word embeddings

Classifiers	Accuracies
Add-Left	94.3
Add-Right	92.6
Add-Nx-Left	93.4
Add-Nx-Right	96.0

Tree Generation



Conclusions

- Syntax can help language modeling.
- Predicting tree structures with Neural Networks is possible.
- Next Steps:
 - Sequence to Tree Models
 - Tree to Tree Models
- code available:
<https://github.com/XingxingZhang/td-treelstm>

Thanks & Questions?