

Deep Learning and Lexical, Syntactic and Semantic Analysis

Wanxiang Che (HIT)

Yue Zhang (SUTD)

Part 4: Dynamic Programming Decoding

Part 4.1: Dynamic Programming

Decoding for Tagging

Word-Level Log-Likelihood

Approach	POS (PWA)	Chunking (F1)	NER (F1)	SRL (F1)
Benchmark Systems	97.24	94.29	89.31	77.92
NN+WLL	96.31	89.13	79.53	55.40

- **WLL: Word-Level Log-Likelihood**
 - Each word in a sentence is considered independently

Sentence-Level Log-Likelihood

- Considering dependencies between tags in a sentence
- Conditional likelihood by **normalizing** all possible paths (CRF)
- Sentence score for one tag path

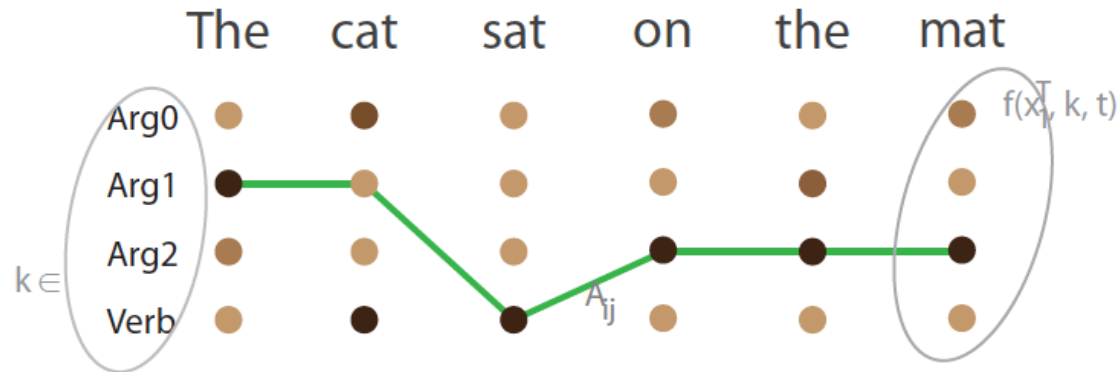
$$\log p([y]_1^T \mid [x]_1^T, \tilde{\theta}) = s([x]_1^T, [y]_1^T, \tilde{\theta}) - \underset{\forall [j]_1^T}{\text{logadd}} s([x]_1^T, [j]_1^T, \tilde{\theta})$$

– where $A_{[i][j]}$ is a transition score for jumping from tag i to j

$$s([x]_1^T, [i]_1^T, \tilde{\theta}) = \sum_{t=1}^T \left(A_{[i]_{t-1}[i]_t} + f([x]_1^T, [i]_t, t, \theta) \right)$$

Sentence-Level Log-Likelihood

- Decoding: finding the max scored path
 - Viterbi algorithm



Results

Approach	POS (PWA)	Chunking (F1)	NER (F1)	SRL (F1)
Benchmark Systems	97.24	94.29	89.31	77.92
NN+WLL	96.31	89.13	79.53	55.40
NN+SLL	96.37	90.33	81.47	70.99

- SLL helps, but fair performance for POS

Improvements

- Supervised word embeddings

Approach	POS (PWA)	CHUNK (F1)	NER (F1)	SRL (F1)
Benchmark Systems	97.24	94.29	89.31	77.92
NN+WLL	96.31	89.13	79.53	55.40
NN+SLL	96.37	90.33	81.47	70.99
NN+WLL+LM1	97.05	91.91	85.68	58.18
NN+SLL+LM1	97.10	93.65	87.58	73.84
NN+WLL+LM2	97.14	92.04	86.96	58.34
NN+SLL+LM2	97.20	93.63	88.67	74.15

- More (embedding) features

Approach	POS (PWA)	CHUNK (F1)	NER (F1)	SRL
Benchmark Systems	97.24	94.29	89.31	77.92
NN+SLL+LM2	97.20	93.63	88.67	74.15
NN+SLL+LM2+Suffix2	97.29	–	–	–
NN+SLL+LM2+Gazetteer	–	–	89.59	–
NN+SLL+LM2+POS	–	94.32	88.67	–
NN+SLL+LM2+CHUNK	–	–	–	74.72

Speed

System	RAM (Mb)	Time (s)
Toutanova, 2003	1100	1065
Shen, 2007	2200	833
SENNA	32	4

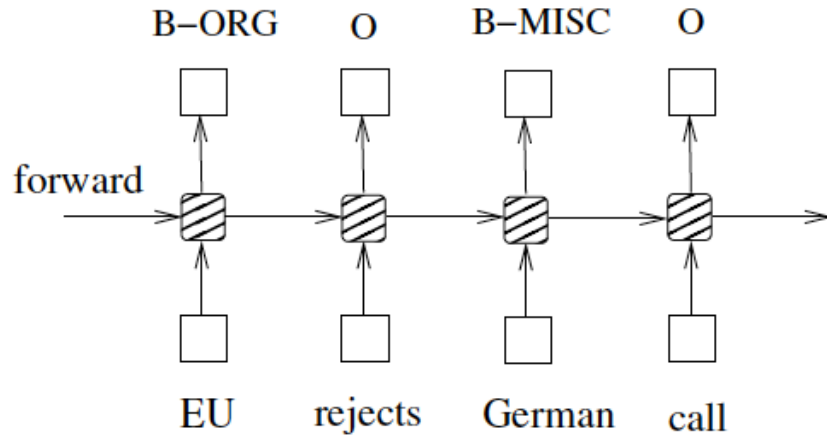
(a) POS

System	RAM (Mb)	Time (s)
Koomen, 2005	3400	6253
SENNA	124	52

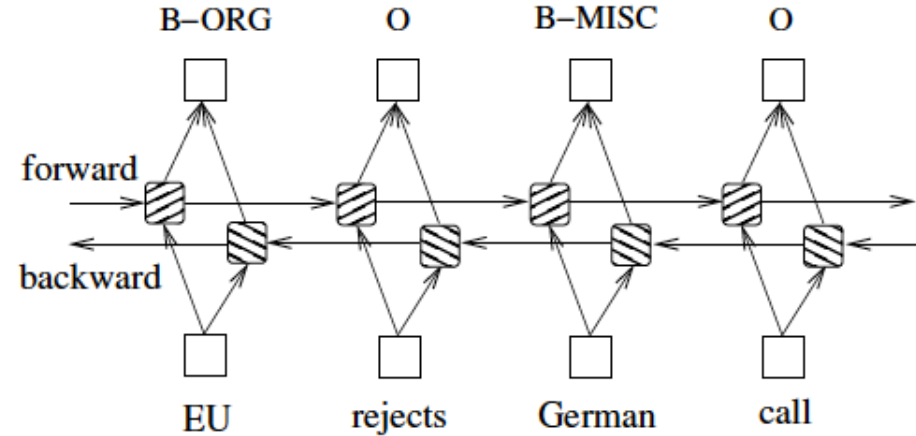
(b) SRL

RNNs for Tagging

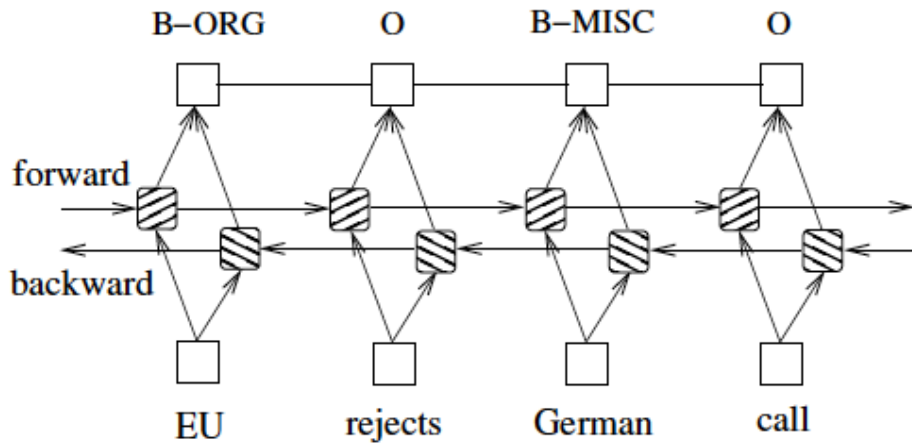
- LSTM



- Bi-LSTM



Bi-LSTM-CRF



Algorithm 1 Bidirectional LSTM CRF model training procedure

```
1: for each epoch do
2:   for each batch do
3:     1) bidirectional LSTM-CRF model forward pass:
4:       forward pass for forward state LSTM
5:       forward pass for backward state LSTM
6:     2) CRF layer forward and backward pass
7:     3) bidirectional LSTM-CRF model backward pass:
8:       backward pass for forward state LSTM
9:       backward pass for backward state LSTM
10:    4) update parameters
11:   end for
12: end for
```

Results

		POS	CoNLL2000	CoNLL2003
Random	Conv-CRF (Collobert et al., 2011)	96.37	90.33	81.47
	LSTM	97.10	92.88	79.82
	BI-LSTM	97.30	93.64	81.11
	CRF	97.30	93.69	83.02
	LSTM-CRF	97.45	93.80	84.10
	BI-LSTM-CRF	97.43	94.13	84.26
Senna	Conv-CRF (Collobert et al., 2011)	97.29	94.32	88.67 (89.59)
	LSTM	97.29	92.99	83.74
	BI-LSTM	97.40	93.92	85.17
	CRF	97.45	93.83	86.13
	LSTM-CRF	97.54	94.27	88.36
	BI-LSTM-CRF	97.55	94.46	88.83 (90.10)

BI-LSTM-CRF for SRL

- End-to-end tagging model
 - 8 layer bi-directional LSTM
 - No parsing features
- Features
 - Argument
 - Predicate
 - Predicate-context
 - Region-mark
- Achieving new SOTA

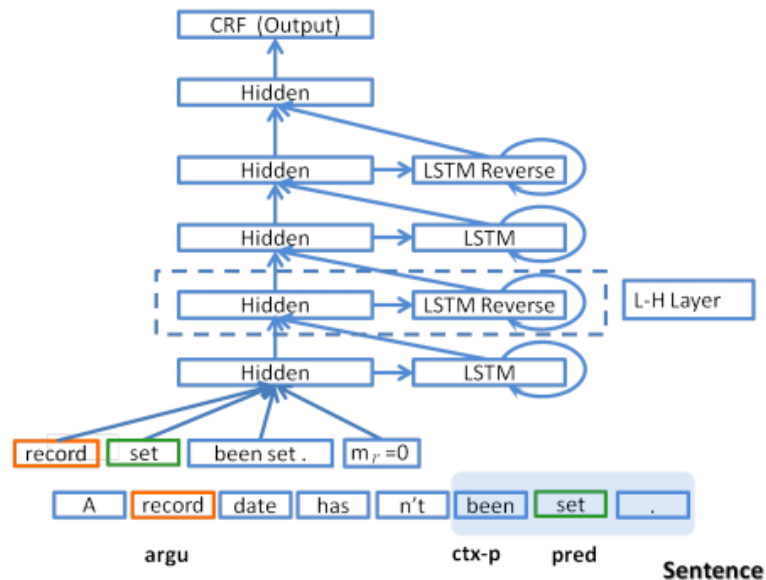
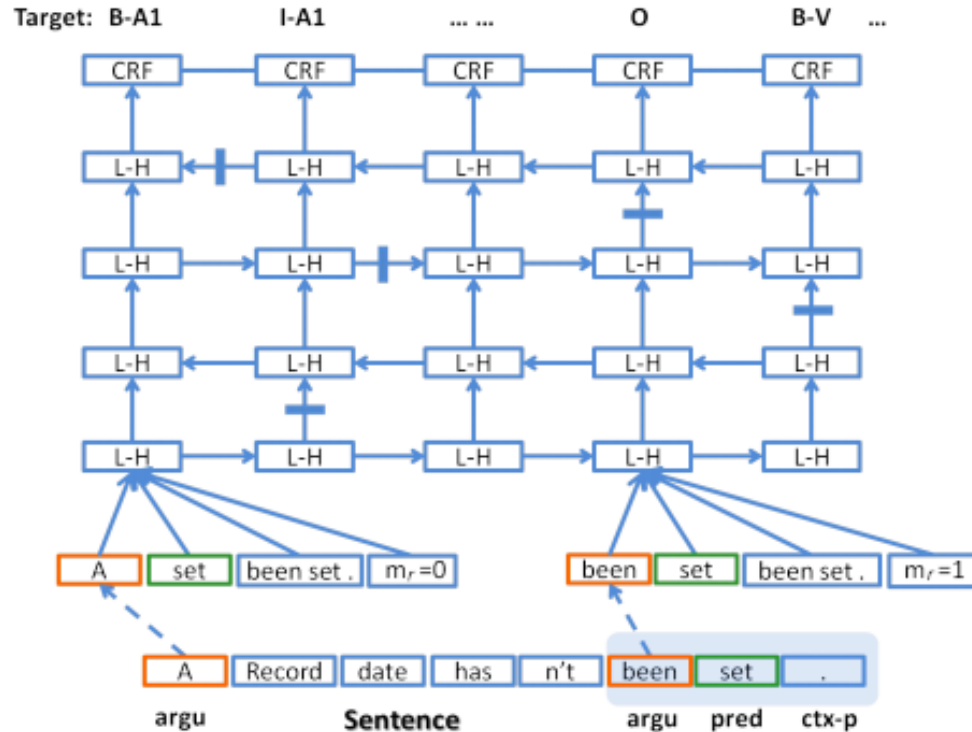
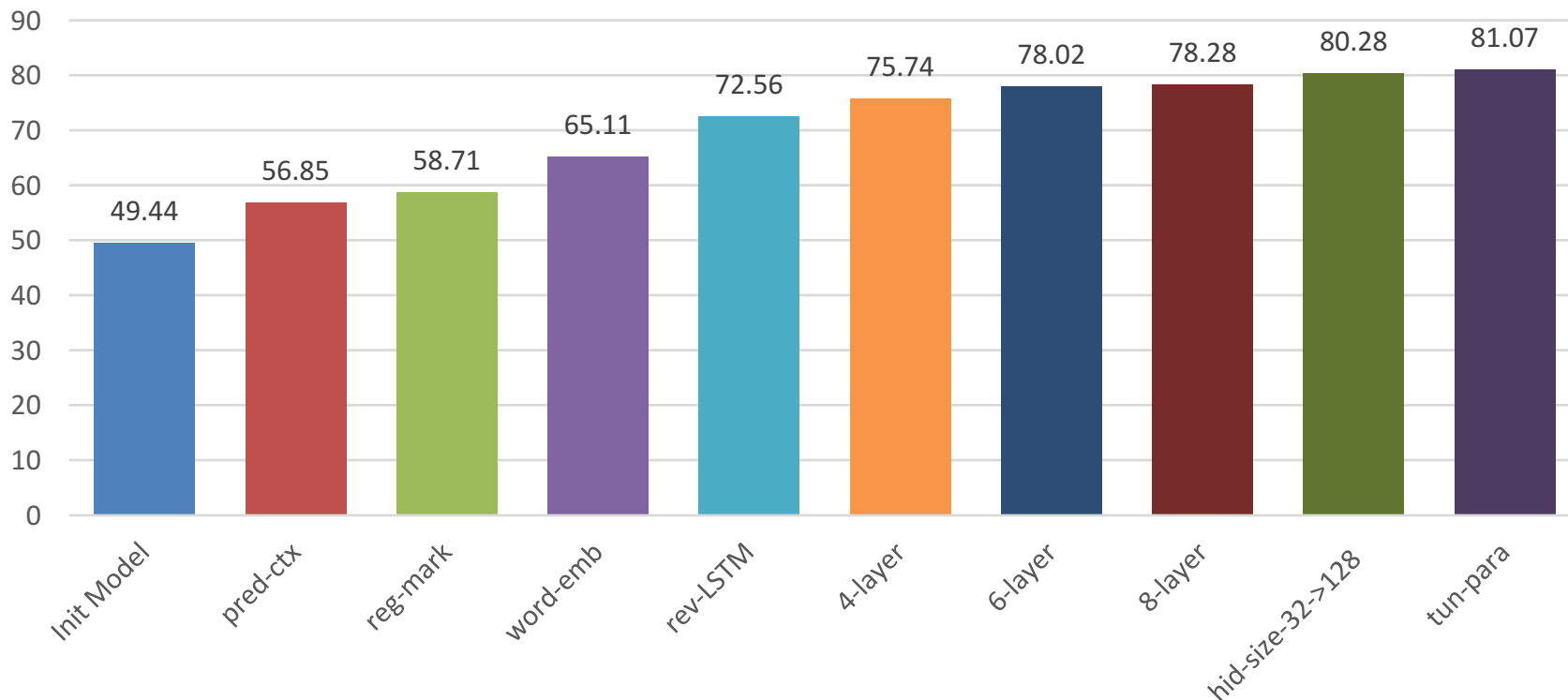


Figure 2: DB-LSTM network. Shadow part denote the predicate context within length 1.

Temporal Expanded

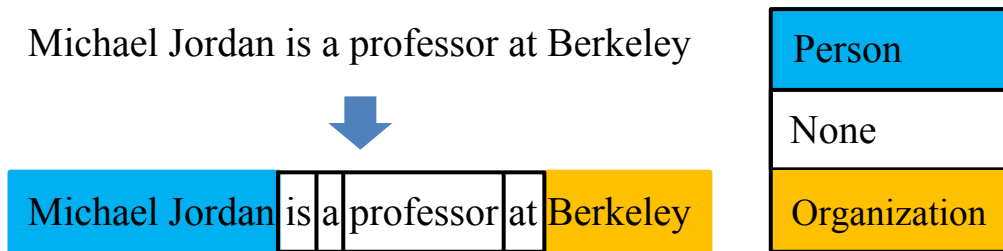


Results



Segmentation Models

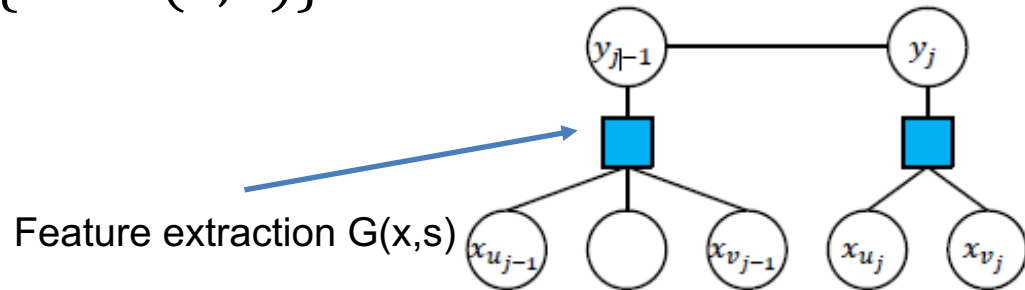
- Tagging models cannot extract segment information
 - E.g. the length of a segment
- Some tagging problems can be naturally modeled into segmentation task
 - E.g. word segmentation, named entity recognition



浦东开发与建设 → 浦东 / 开发 / 与 / 建设
Pudong development and construction

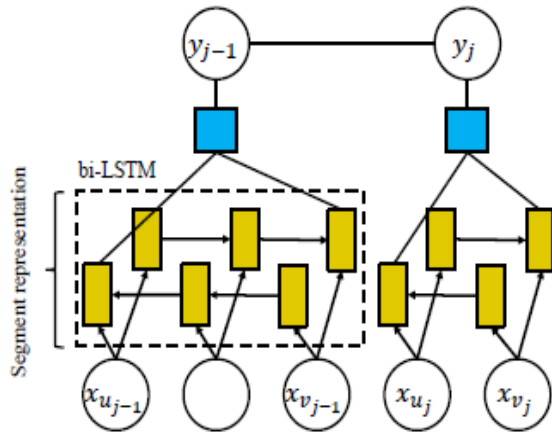
Semi-CRF

- A solution
 - Semi-Markov CRF [Sarawagi and Cohen, 2004]
 - Modeling segments directly
 - $p(\mathbf{s}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\{W \cdot G(\mathbf{x}, \mathbf{s})\}$

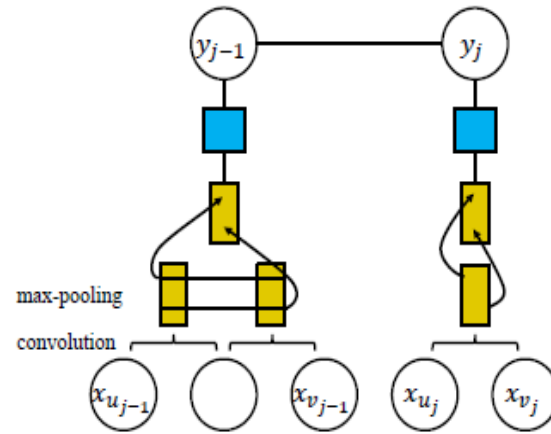


Can we represent segments with vectors?

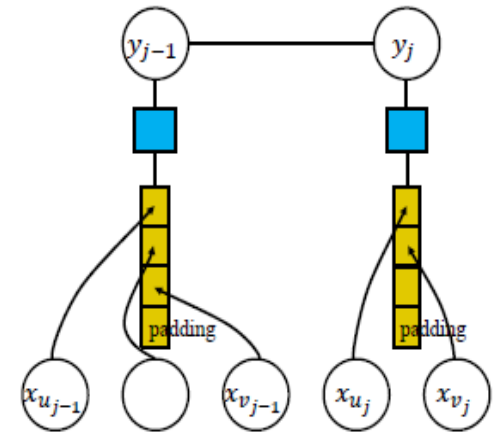
Compositional Segment Representation



(b) SRNN



(c) SCNN



(d) SCONCATE

Decoding Algorithm

Input: a sequence $X = (x_0, \dots, x_{n-1})$ of n units, the maximum length of the segment L

Output: the highest scored segmentation $S = (s_0, \dots, s_{m-1})$, where $s = (u, v, y)$ is a segment and u represents the starting position, v represents the ending position, and an optional tag y associate with the segment.

Defining $V(i, y)$ which represents the best sub-segmentation that ends with x_i (not included) and $V(i, y)$ can be calculated as:

$$V(i, y) = \begin{cases} \max_{y', d=1 \dots L} V(i-d, y') + \text{score}(i-d, i, y), & \text{if } i > 0 \\ 0, & \text{if } i = 0 \\ -\infty, & \text{if } i < 0 \end{cases}$$

for $i \leftarrow 1 \dots n$

for $y \in \mathcal{Y}$:

for $d \leftarrow 1 \dots L$

 if $i - d = 0$:

$V(i, y) \leftarrow \text{score}(i - d, i, y)$

 else:

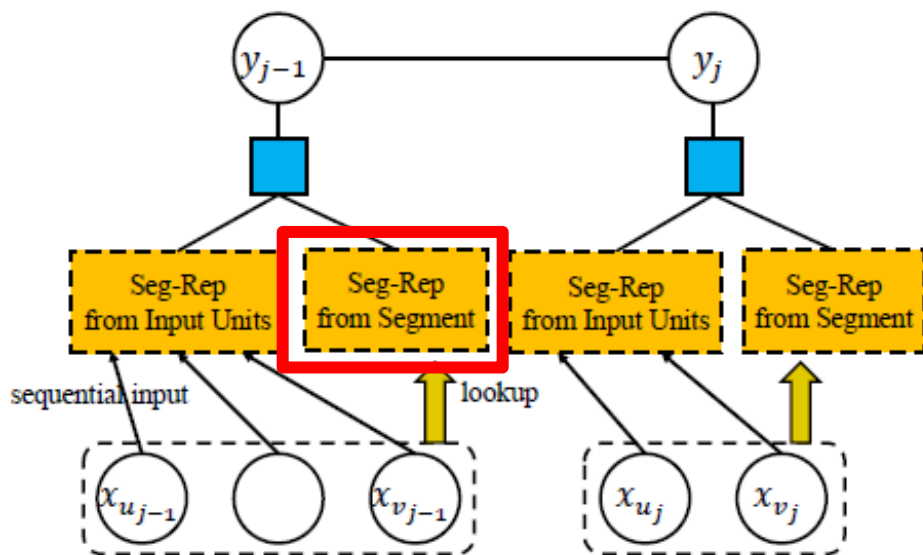
$best_{i-d} \leftarrow \max_{y'} V(i - d, y')$

$V(i, y) \leftarrow \max(V(i, y), best_{i-d} + \text{score}(i - d, i, y))$

Results

		NER CoNLL03		CTB6		CWS PKU		MSR		spd
		dev	test	dev	test	dev	test	dev	test	
<i>baseline</i>	NN-LABELER	93.03	88.62	93.70	93.06	93.57	92.99	93.22	93.79	3.30
	NN-CRF	93.06	89.08	94.33	93.65	94.09	93.28	93.81	94.17	2.72
	SPARSE-CRF	88.87	83.43	95.68	95.08	95.85	95.06	96.09	96.54	
<i>neural semi-CRF</i>	SRNN	92.97	88.63	94.56	94.06	94.86	93.91	94.38	95.21	0.62
	SCONCATE	92.96	89.07	94.34	93.96	94.41	93.57	94.05	94.53	1.08
	SCNN	91.53	87.68	87.82	87.51	79.64	80.75	85.04	85.79	1.46

Segment-level Representation



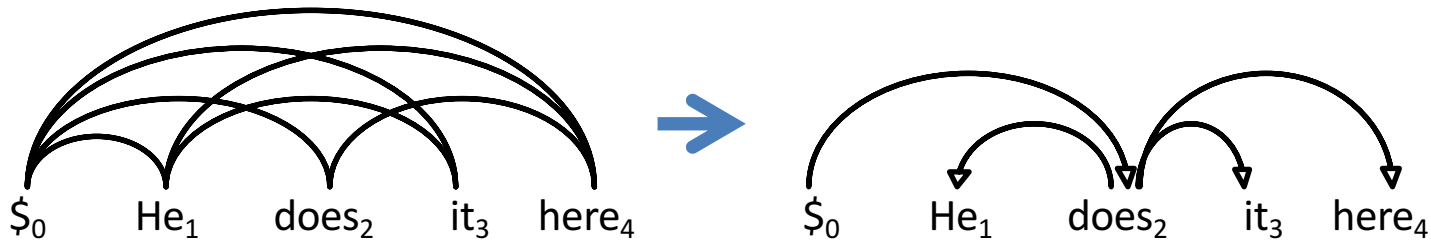
<i>model</i>	CoNLL03	CTB6	PKU	MSR
NN-LABELER	88.62	93.06	92.99	93.79
NN-CRF	89.08	93.65	93.28	94.17
SPARSE-CRF	83.43	95.08	95.06	96.54
SRNN	88.63	94.06	93.91	95.21
+SEMB-HETERO	89.59	95.48	95.60	97.39
	+0.96	+1.42	+1.69	+2.18
SCONCATE	89.07	93.96	93.57	94.53
+SEMB-HETERO	89.77	95.42	95.67	97.58
	+0.70	+1.43	+2.10	+3.05

Part 4.2: Dynamic Programming

Decoding for Parsing

Graph-based Dependency Parsing

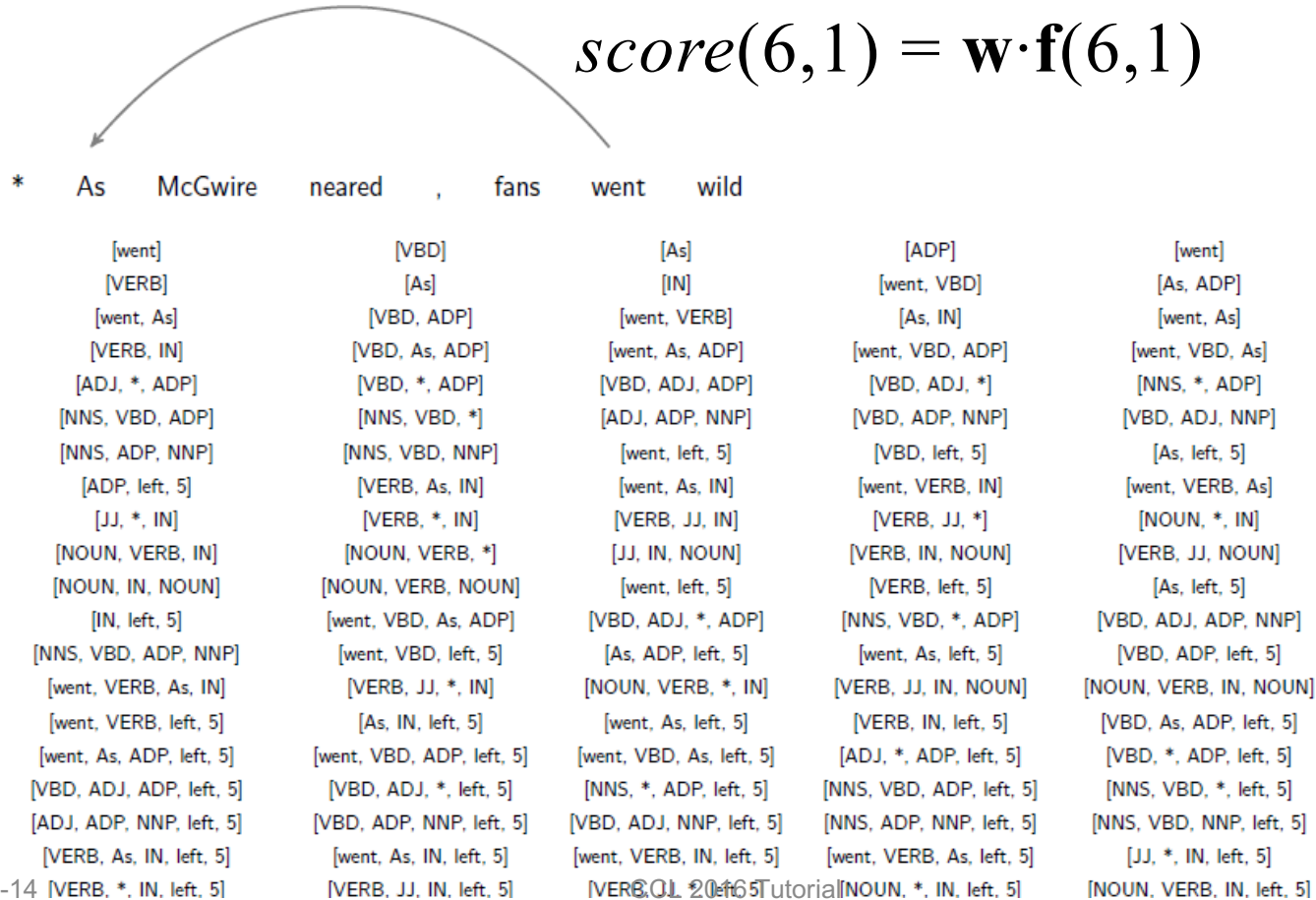
- Find the highest scoring tree from a complete graph
- Dynamic Programming Decoding
 - E.g. Eisner Algorithm



$$Y^* = \arg \max_{Y \in \Phi(X)} score(X, Y)$$

How to Score an Arc?

$$\text{score}(6,1) = \mathbf{w} \cdot \mathbf{f}(6,1)$$



NN for Graph-based Parsing

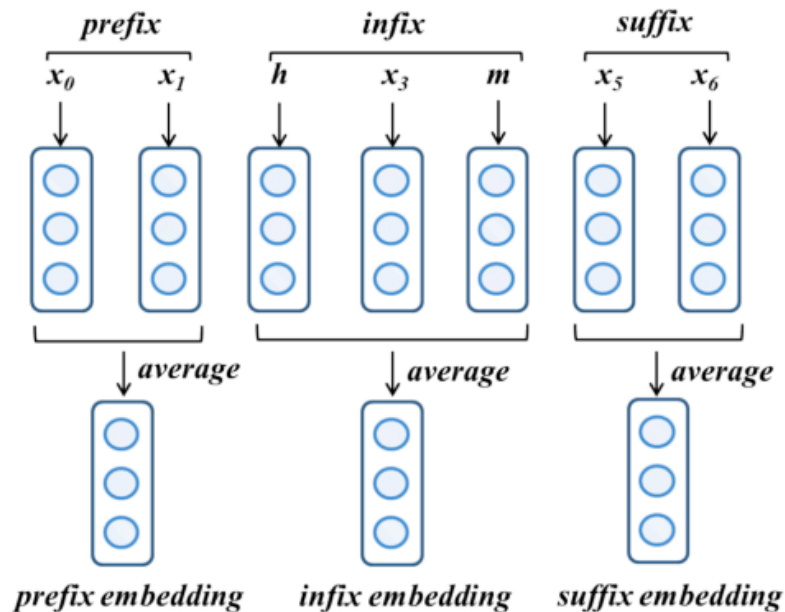
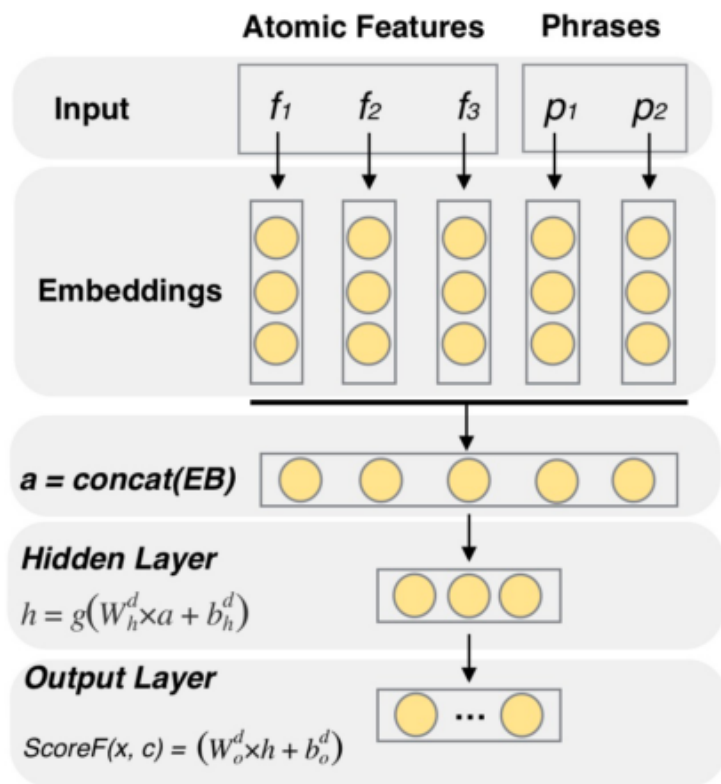


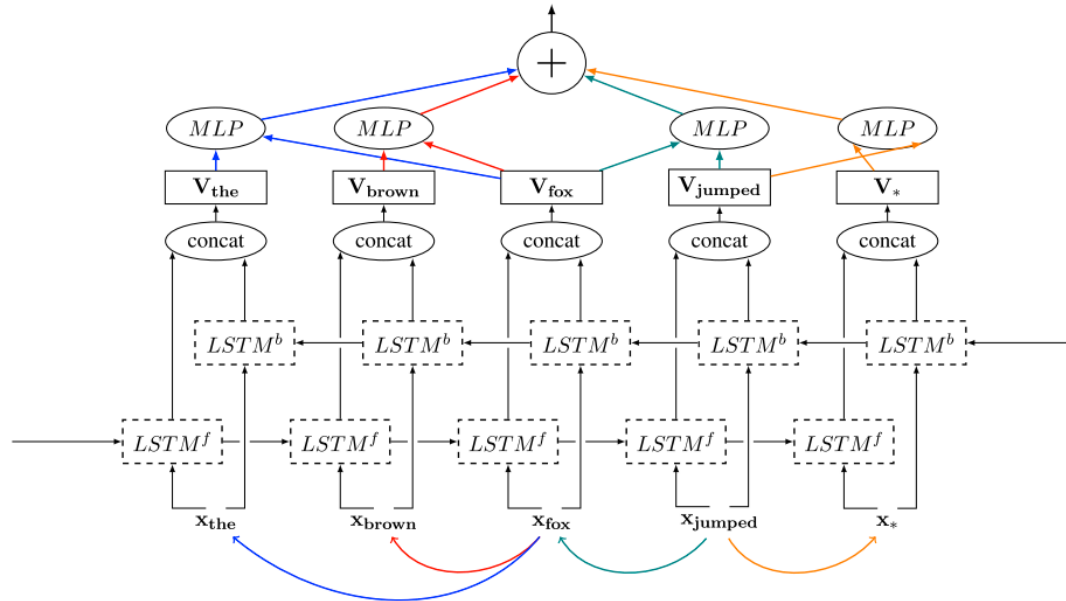
Figure 3: Illustration for phrase embeddings. h , m and x_0 to x_6 are words in the sentence.

Results

	Models	Dev		Test		Speed (sent/s)
		UAS	LAS	UAS	LAS	
First-order	MSTParser-1-order	92.01	90.77	91.60	90.39	20
	1-order-atomic-rand	92.00	90.71	91.62	90.41	55
	1-order-atomic	92.19	90.94	92.14	90.92	55
	1-order-phrase-rand	92.47	91.19	92.25	91.05	26
	1-order-phrase	92.82	91.48	92.59	91.37	26
Second-order	MSTParser-2-order	92.70	91.48	92.30	91.06	14
	2-order-phrase-rand	93.39	92.10	92.99	91.79	10
	2-order-phrase	93.57	92.29	93.29	92.13	10
Third-order	(Koo and Collins, 2010)	93.49	N/A	93.04	N/A	N/A

BI-LSTM for Graph-based Parsing-I

- Each dependency arc in a sentence is scored using MLP that is fed the BI-LSTM encoding of the words at the arc's end points



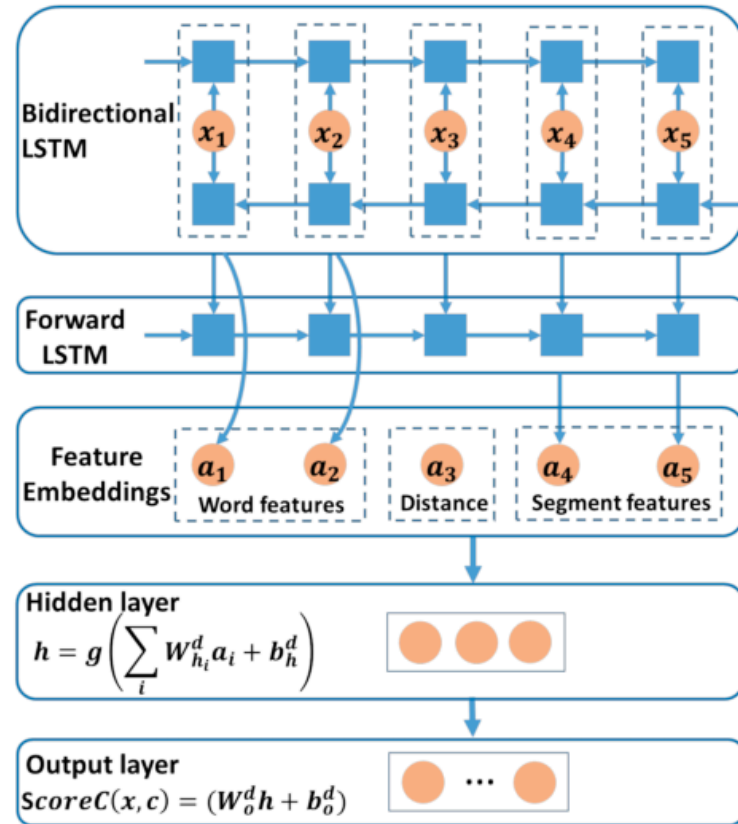
Results

System	Method	Representation	Emb	PTB-YM	PTB-SD		CTB	
				UAS	UAS	LAS	UAS	LAS
This work	graph, 1st order	2 BiLSTM vectors	–	–	93.1	91.0	86.6	85.1
This work	transition (greedy, dyn-oracle)	4 BiLSTM vectors	–	–	93.1	91.0	86.2	85.0
This work	transition (greedy, dyn-oracle)	11 BiLSTM vectors	–	–	93.2	91.2	86.5	84.9
ZhangNivre11	transition (beam)	large feature set (sparse)	–	92.9	–	–	86.0	84.4
Martins13 (TurboParser)	graph, 3rd order+	large feature set (sparse)	–	92.8	93.1	–	–	–
Pei15	graph, 2nd order	large feature set (dense)	–	93.0	–	–	–	–
Dyer15	transition (greedy)	Stack-LSTM + composition	–	–	92.4	90.0	85.7	84.1
Ballesteros16	transition (greedy, dyn-oracle)	Stack-LSTM + composition	–	–	92.7	90.6	86.1	84.5
This work	graph, 1st order	2 BiLSTM vectors	YES	–	93.0	90.9	86.5	84.9
This work	transition (greedy, dyn-oracle)	4 BiLSTM vectors	YES	–	93.6	91.5	87.4	85.9
This work	transition (greedy, dyn-oracle)	11 BiLSTM vectors	YES	–	93.9	91.9	87.6	86.1
Weiss15	transition (greedy)	large feature set (dense)	YES	–	93.2	91.2	–	–
Weiss15	transition (beam)	large feature set (dense)	YES	–	94.0	92.0	–	–
Pei15	graph, 2nd order	large feature set (dense)	YES	93.3	–	–	–	–
Dyer15	transition (greedy)	Stack-LSTM + composition	YES	–	93.1	90.9	87.1	85.5
Ballesteros16	transition (greedy, dyn-oracle)	Stack-LSTM + composition	YES	–	93.6	91.4	87.6	86.2
LeZuidema14	reranking /blend	inside-outside recursive net	YES	93.1	93.8	91.5	–	–
Zhu15	reranking /blend	recursive conv-net	YES	93.8	–	–	85.7	–

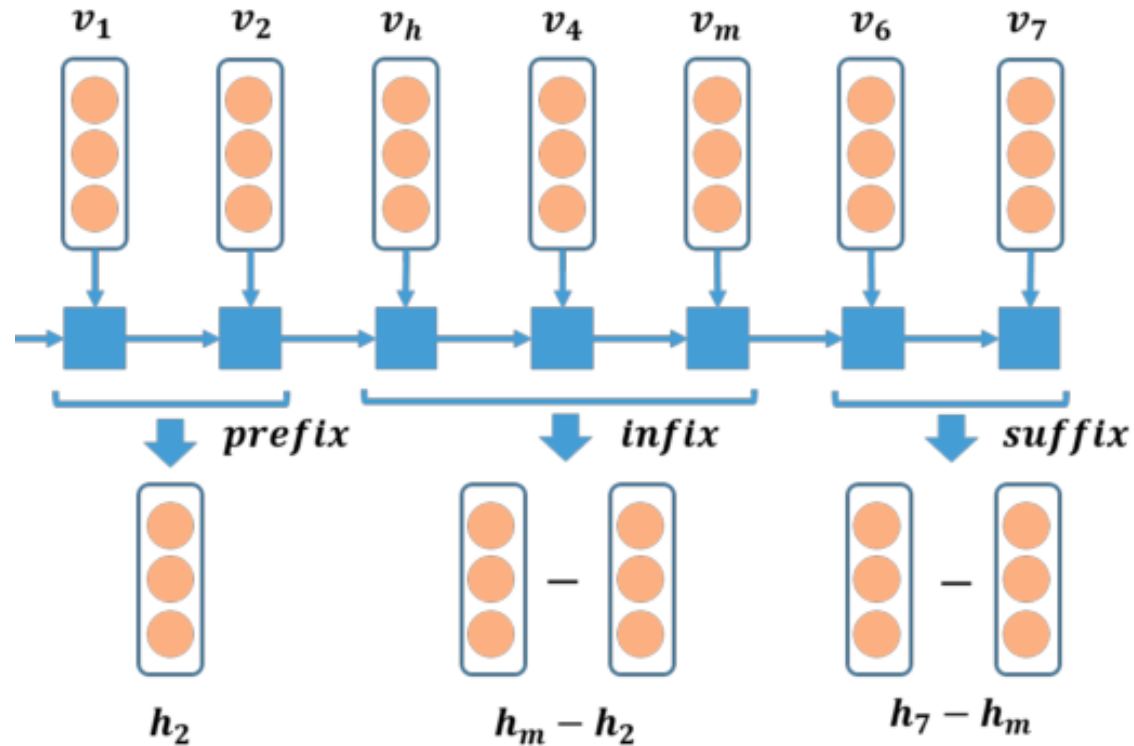
Kiperwasser, E., & Goldberg, Y. (2016). Simple and Accurate Dependency Parsing Using Bidirectional LSTM Feature Representations. *TACL*.

BI-LSTM for Graph-based Parsing-II

- Besides the word vectors, they used sentence segment (phrase) embeddings



Learning Segment Embeddings



Results

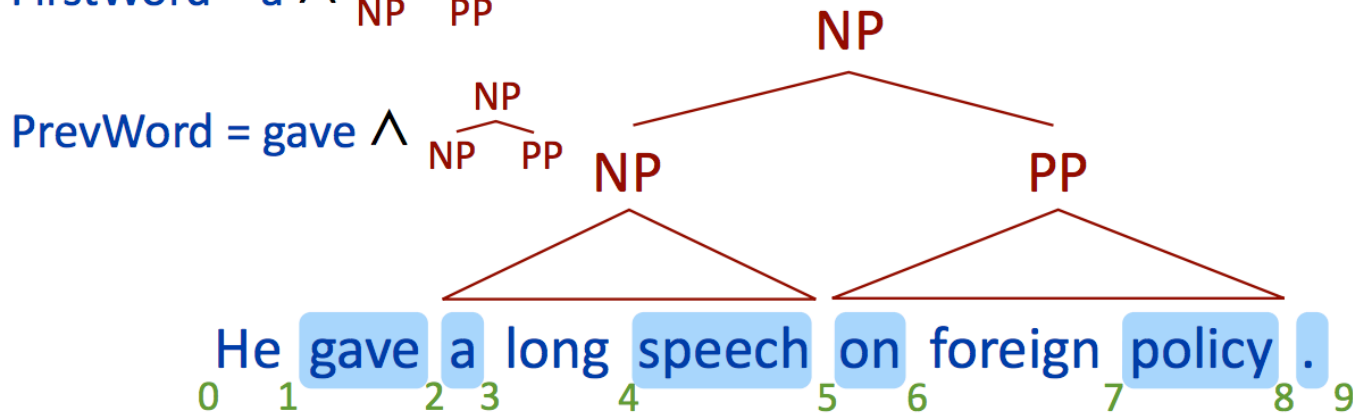
	Models	UAS	LAS	Speed(sent/s)
First-order	MSTParser	91.60	90.39	20
	1st-order atomic (Pei et al., 2015)	92.14	90.92	55
	1st-order phrase (Pei et al., 2015)	92.59	91.37	26
	Our basic model	93.09	92.03	61
	Our basic model + segment	93.51	92.45	26
Second-order	MSTParser	92.30	91.06	14
	2nd-order phrase (Pei et al., 2015)	93.29	92.13	10
Third-order	(Koo and Collins, 2010)	93.04	N/A	N/A
Fourth-order	(Ma and Zhao, 2012)	93.4	N/A	N/A
Unlimited-order	(Zhang and McDonald, 2012)	93.06	91.86	N/A
	(Zhang et al., 2013)	93.50	92.41	N/A
	(Zhang and McDonald, 2014)	93.57	92.48	N/A

Neural CRF for Constituency Parsing

- CRF Parsing with CKY decoding

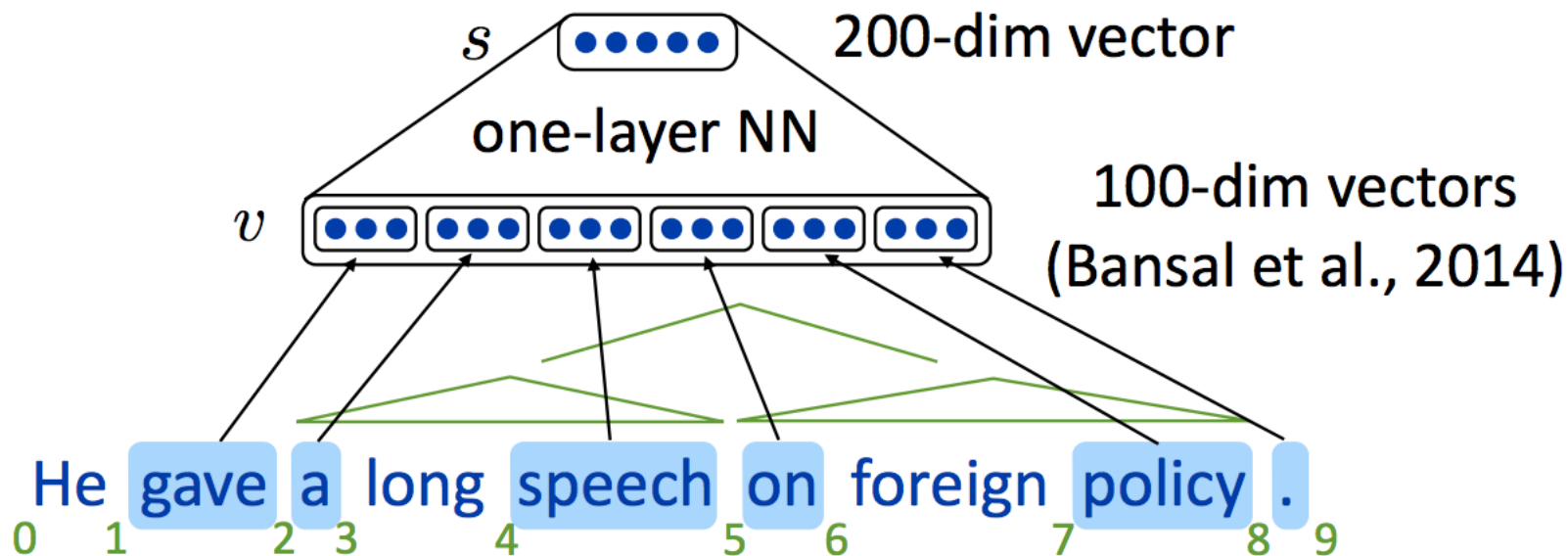
$$P(T|x) \propto \prod_{r \in T} \exp(\text{score}(r)) \quad \text{score} \left(\begin{array}{c} \text{NP} \\ \text{NP} \quad \text{PP} \\ 2 \quad 5 \quad 8 \end{array} \right) = w^\top f \left(\begin{array}{c} \text{NP} \\ \text{NP} \quad \text{PP} \\ 2 \quad 5 \quad 8 \end{array} \right)$$

FirstWord = a \wedge $\begin{array}{c} \text{NP} \\ \text{NP} \quad \text{PP} \end{array}$

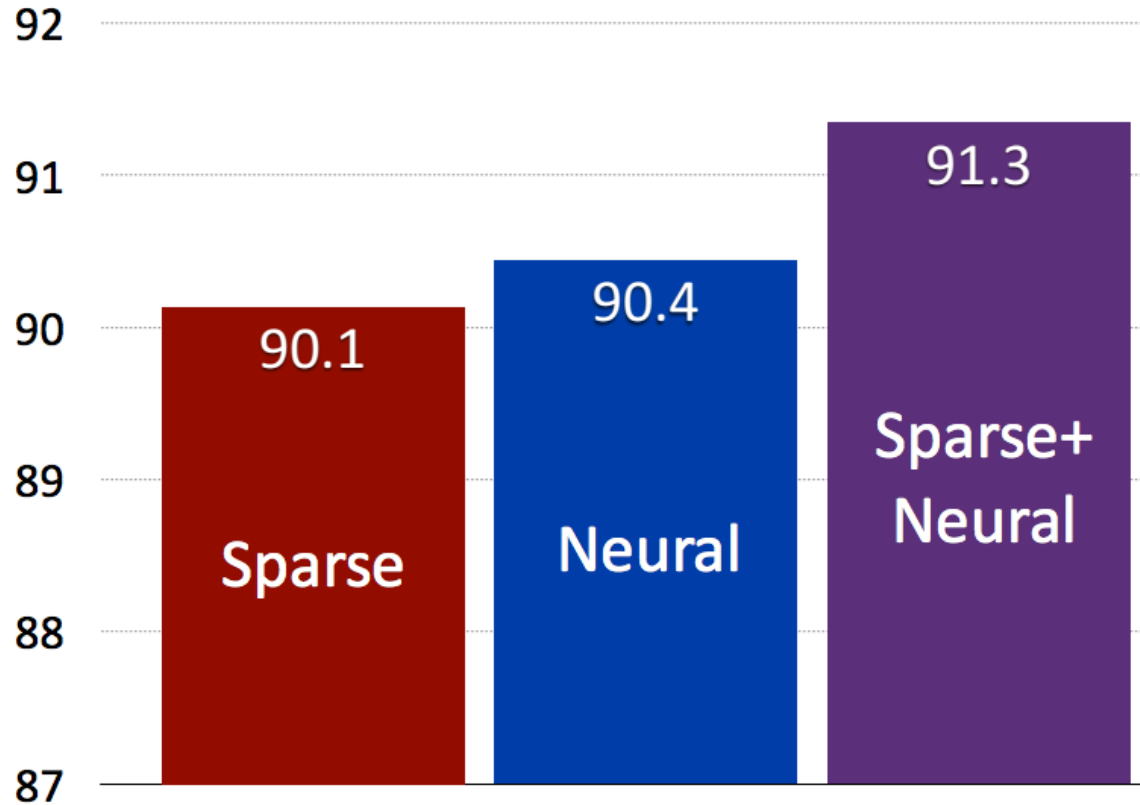


Neural CRF for Constituency Parsing

- Neural CRF Parsing

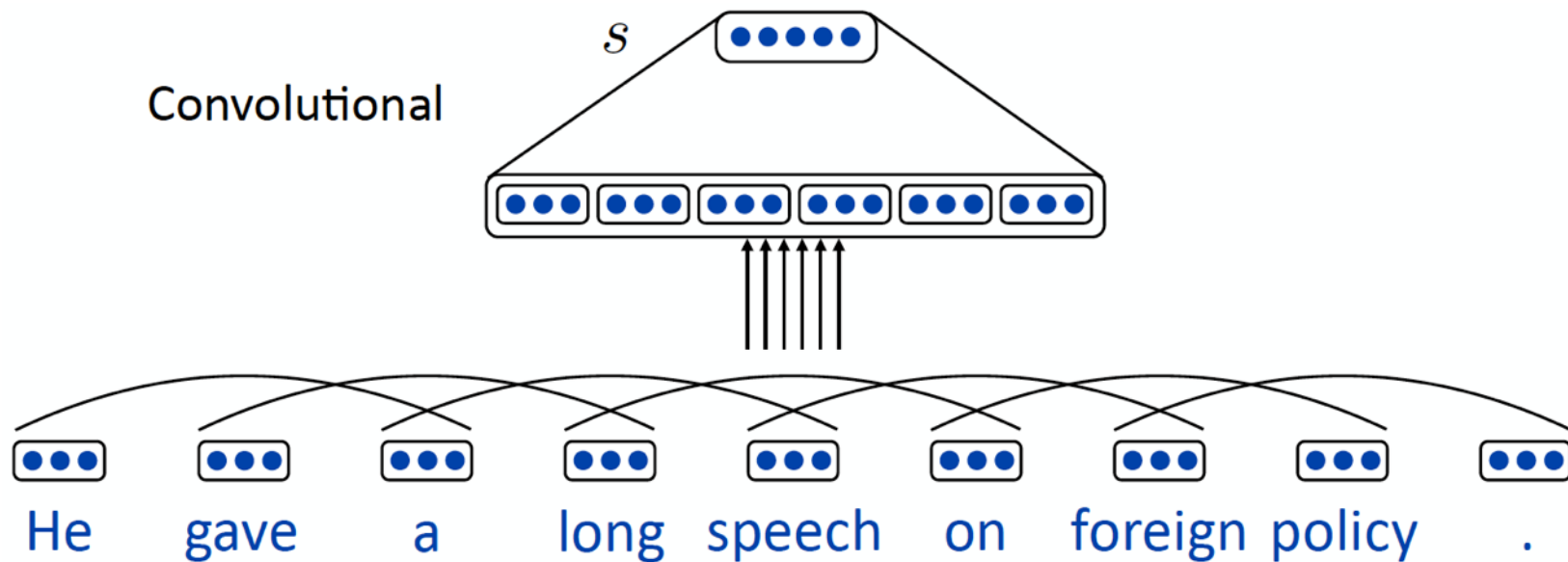


Results



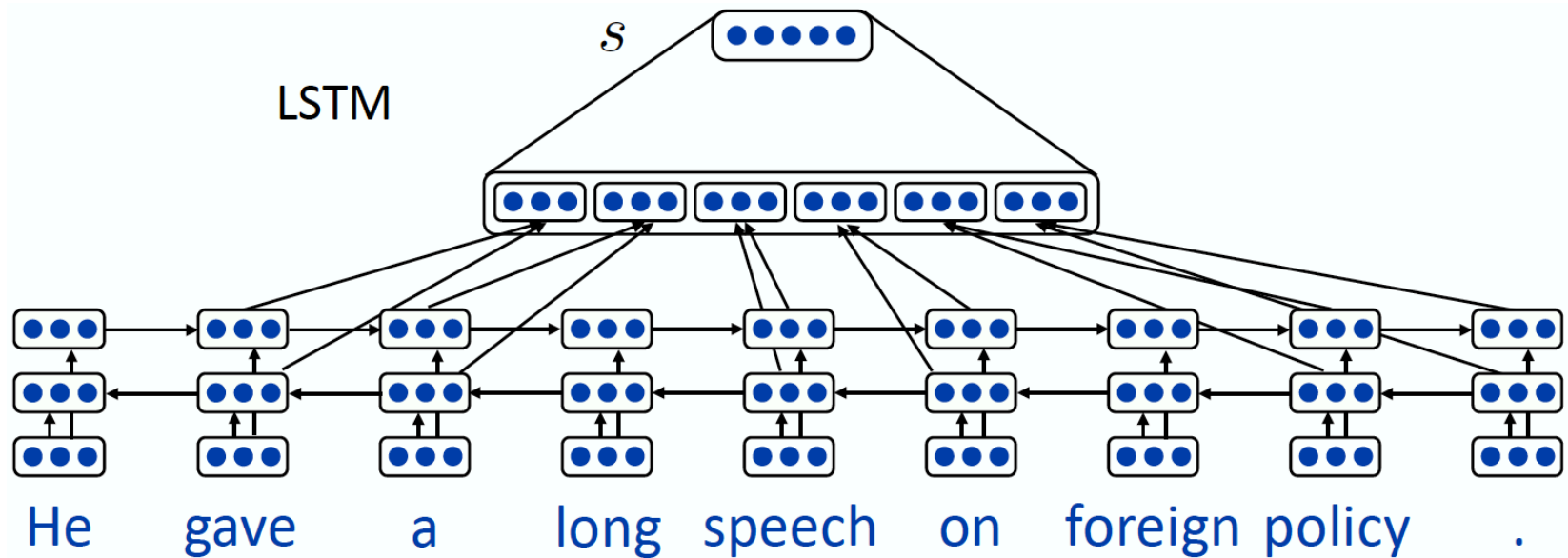
Neural CRF for Constituency Parsing

- More neural networks



Neural CRF for Constituency Parsing

- More neural networks



Conclusion

- Neural nets can provide continuous features in discrete structured models
- Inference and learning are almost unchanged from the purely discrete model