

# Deep Learning and Lexical, Syntactic and Semantic Analysis

Wanxiang Che (HIT)

Yue Zhang (SUTD)

# Part 5: Beam-search Decoding

# A transition system

- Automata
  - State
    - Start state — an empty structure
    - End state — the output structure
    - Intermediate states — partially constructed structures
  - Actions
    - Change one state to another

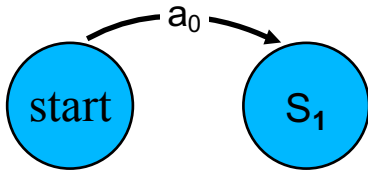
# A transition system

- Automata



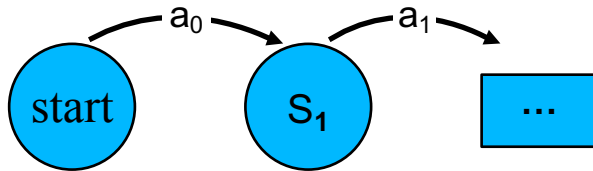
# A transition system

- Automata



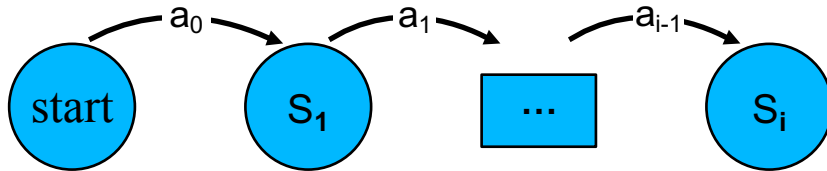
# A transition system

- Automata



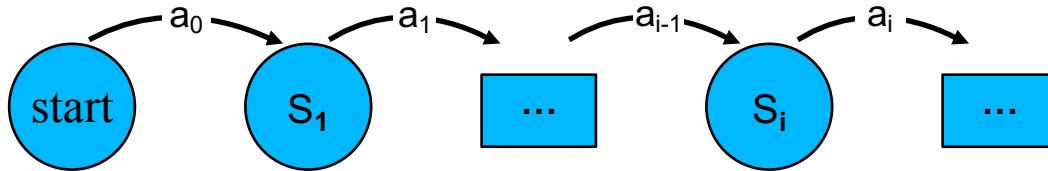
# A transition system

- Automata



# A transition system

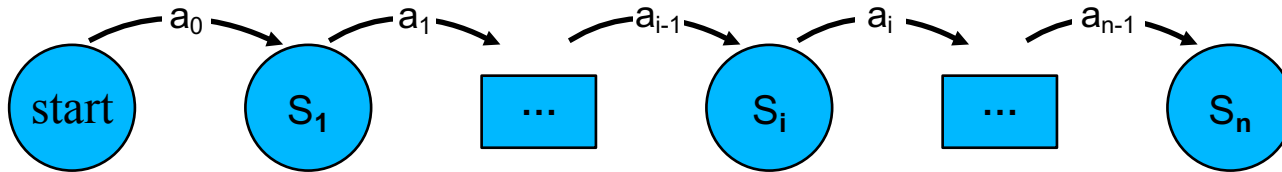
- Automata





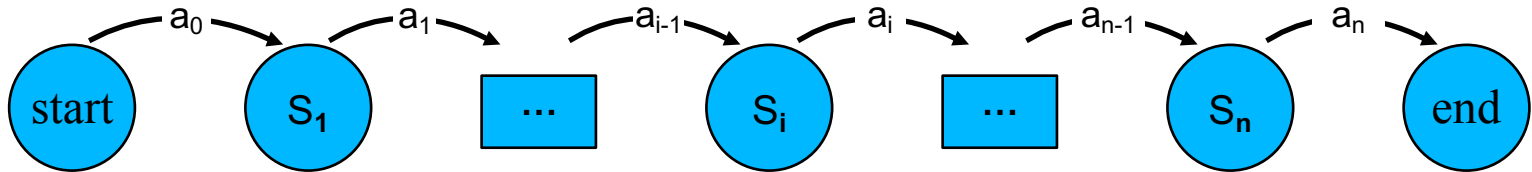
# A transition system

- Automata



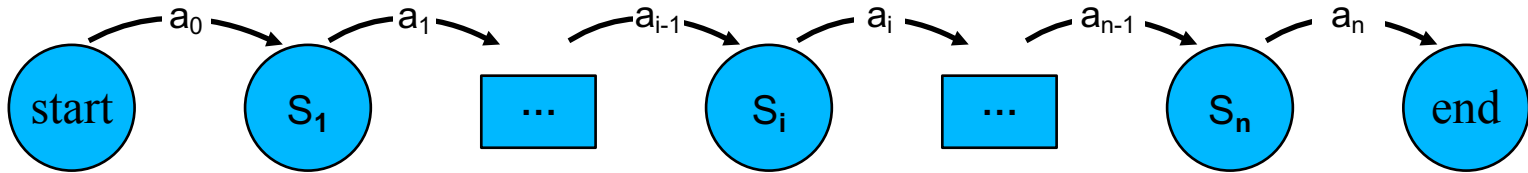
# A transition system

- Automata



# A transition system

- State
  - Corresponds to partial results during decoding
    - start state, end state,  $S_i$



- Actions
  - The operations that can be applied for state transition
  - Construct output incrementally
    - $a_i$

# A transition-based POS-tagging example

- POS tagging

I like reading books → I/PRON like/VERB reading/VERB books/NOUN

- Transition system

- State

- Partially labeled word-POS pairs
    - Unprocessed words

- Actions

- TAG(t)  $w_1/t_1 \cdots w_i/t_i \rightarrow w_1/t_1 \cdots w_i/t_i w_{i+1}/t$

# A transition-based POS-tagging example

- Start State



I like reading books

# A transition-based POS-tagging example

- TAG(PRON)

I/PRON

like reading books

# A transition-based POS-tagging example

- TAG(VERB)

I/PRON like/VERB

reading books

# A transition-based POS-tagging example

- TAG(VERB)

I/PRON like/VERB reading/VERB

books



# A transition-based POS-tagging example

- TAG (NOUN)

I/PRON like/VERB reading/VERB books/NOUN

# A transition-based POS-tagging example

- End State

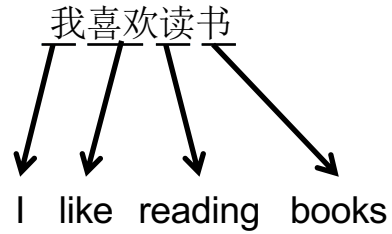
I/PRON like/VERB reading/VERB books/NOUN

# Word segmentation

- State
  - Partially segmented results
  - Unprocessed characters
- Two candidate actions
  - Separate    ## ## → ## ## #
  - Append      ## ## → ## ## #

# Word segmentation

- Initial State



# Word segmentation

- Separate

我

喜欢读书

# Word segmentation

- Separate

我 喜

欢读书

# Word segmentation

- Append

我 喜欢

读书

# Word segmentation

- Separate

我 喜欢 读

书



# Word segmentation

- Separate

我 喜欢 读 书

# Word segmentation

- End State

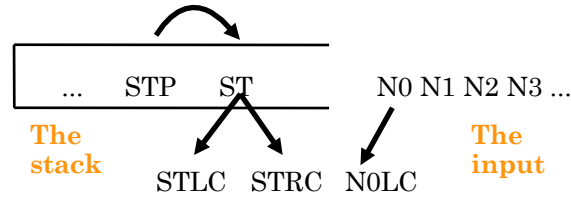
我 喜欢 读 书

# The arc-eager transition system

- State
  - A stack to hold partial structures
  - A queue of next incoming words
- Actions
  - SHIFT, REDUCE, ARC-LEFT, ARC-RIGHT

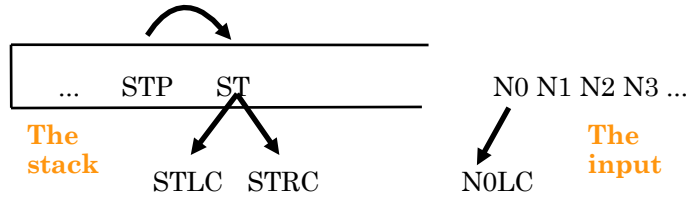
# The arc-eager transition system

- State



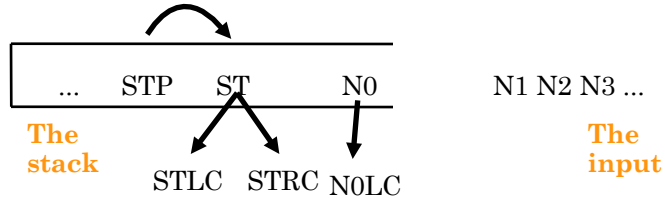
# The arc-eager transition system

- Actions
  - Shift



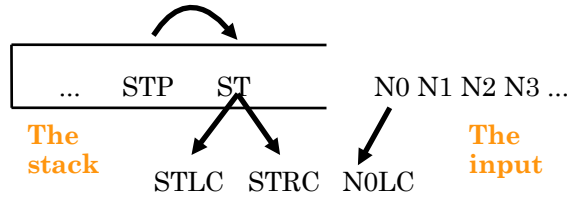
# The arc-eager transition system

- Actions
  - Shift
    - Pushes stack



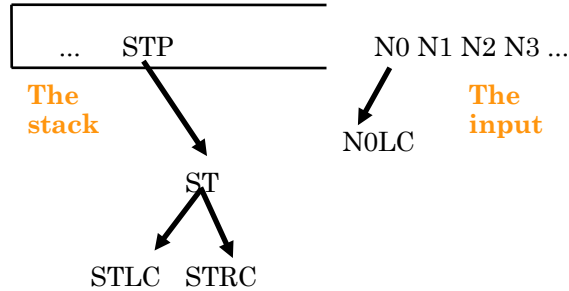
# The arc-eager transition system

- Actions
- Reduce



# The arc-eager transition system

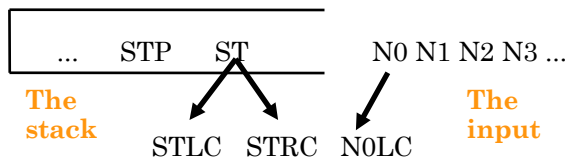
- Actions
  - Reduce
    - Pops stack





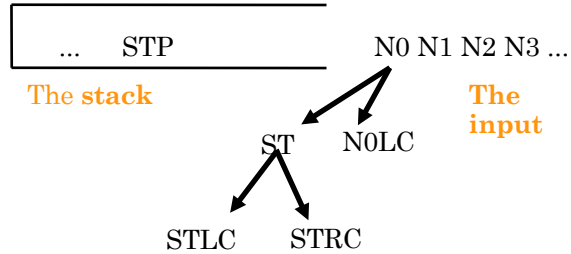
# The arc-eager transition system

- Actions
- Arc-Left



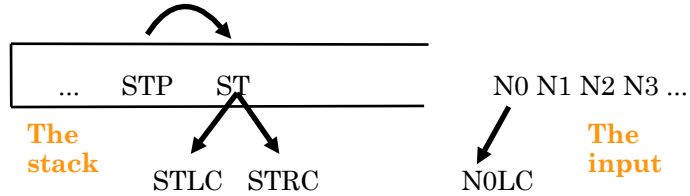
# The arc-eager transition system

- Actions
  - Arc-Left
    - Pops stack
    - Adds link



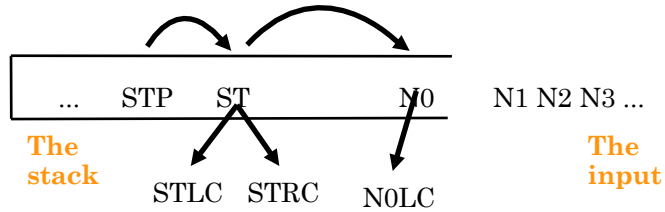
# The arc-eager transition system

- Actions
  - Arc-right



# The arc-eager transition system

- Actions
  - Arc-right
    - Pushes stack
    - Adds link



# The arc-eager transition system

- An example
  - S – Shift
  - R – Reduce
  - AL – ArcLeft
  - AR – ArcRight

He does it here

# The arc-eager transition system

- An example

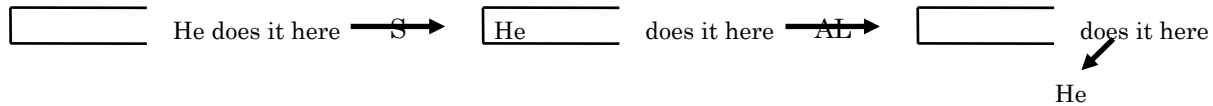
- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight

He does it here  $\xrightarrow{S}$   He  does it here

# The arc-eager transition system

- An example

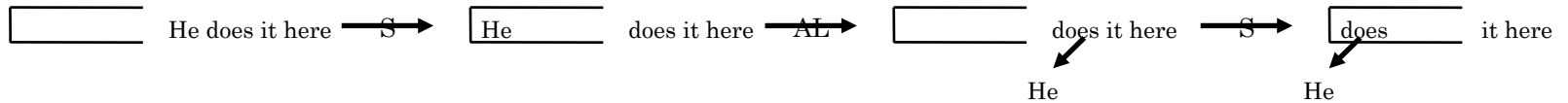
- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight



# The arc-eager transition system

- An example

- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight

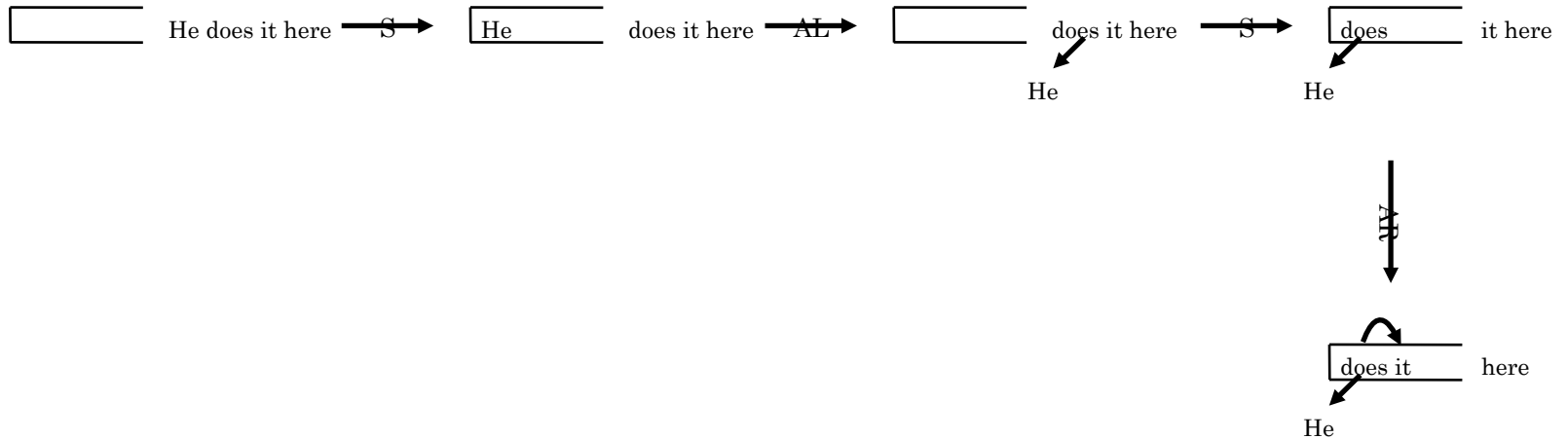




# The arc-eager transition system

- An example

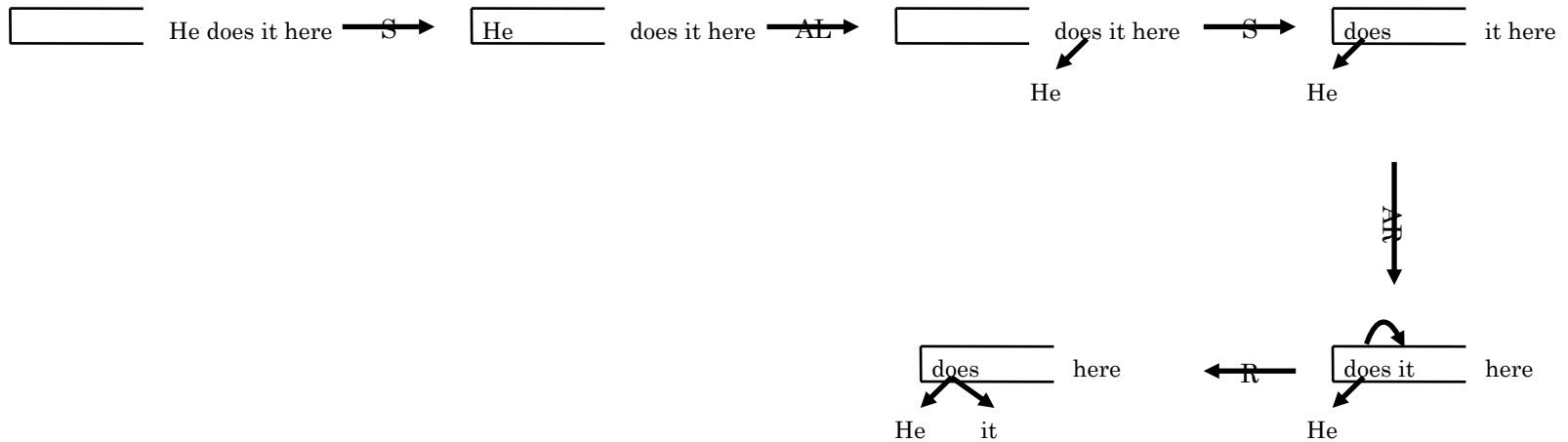
- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight



# The arc-eager transition system

- An example

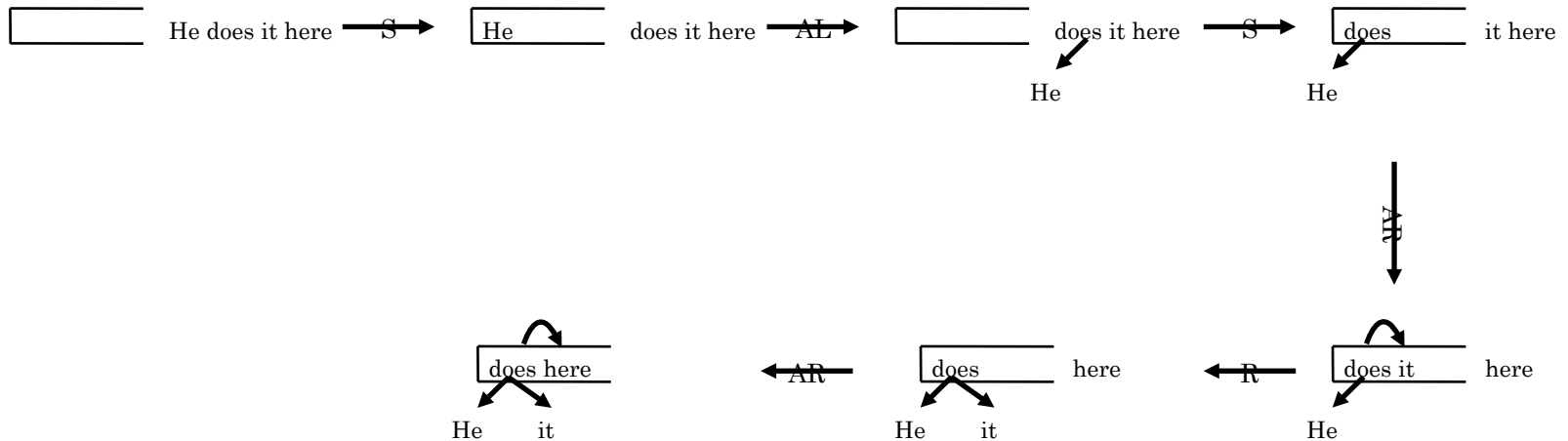
- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight



# The arc-eager transition system

- An example

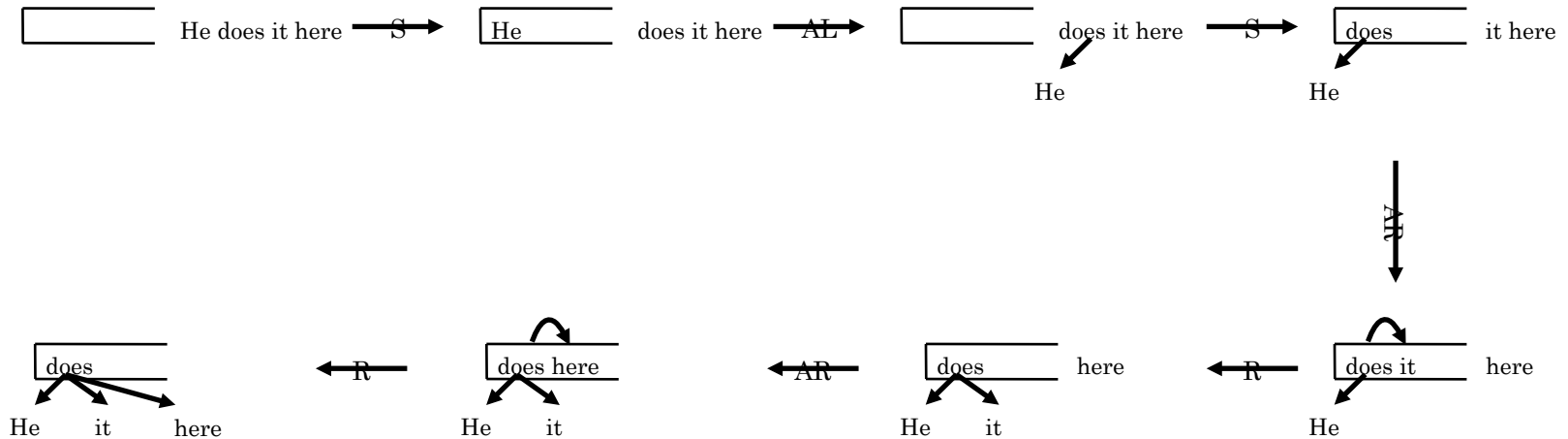
- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight



# The arc-eager transition system

- An example

- S – Shift
- R – Reduce
- AL – ArcLeft
- AR – ArcRight



# Other examples

- Language generation
- Translation
  - Word by word
  - Phrase by phrase
  - Syntax tree synthesis

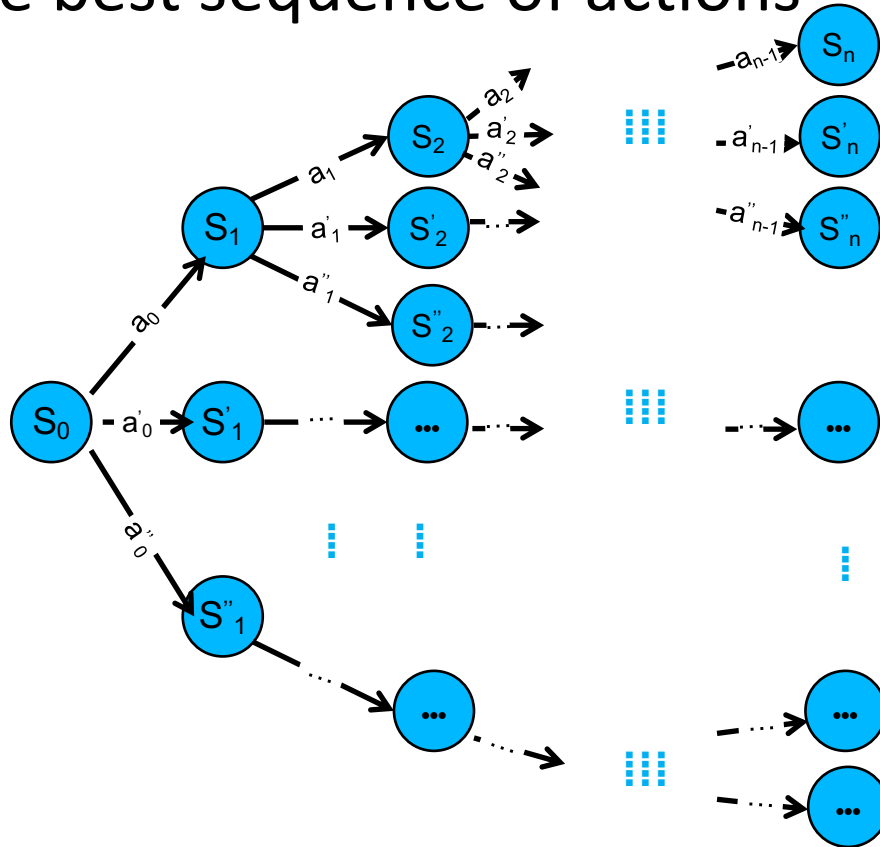
# Part 5.1: Beam-search Decoding

## ——learning to search

(Zhang and Clark,2011)

# Search

- Find the best sequence of actions

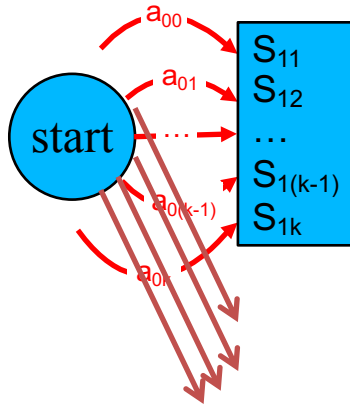


# Beam-search decoding

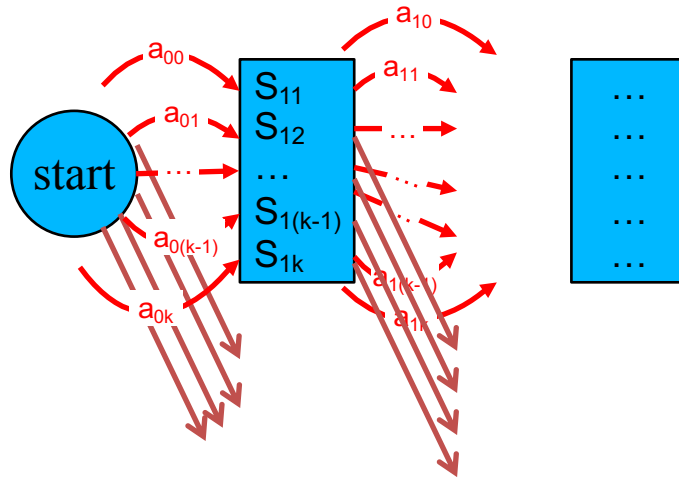




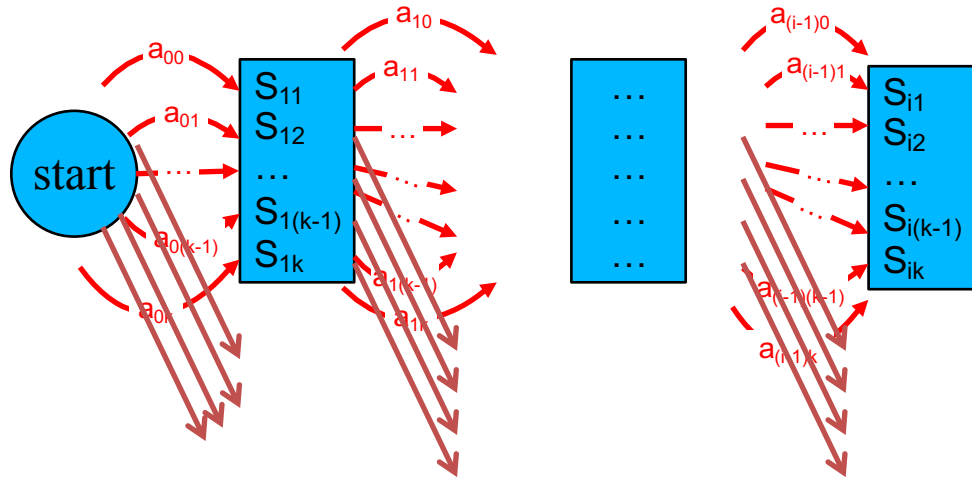
# Beam-search decoding



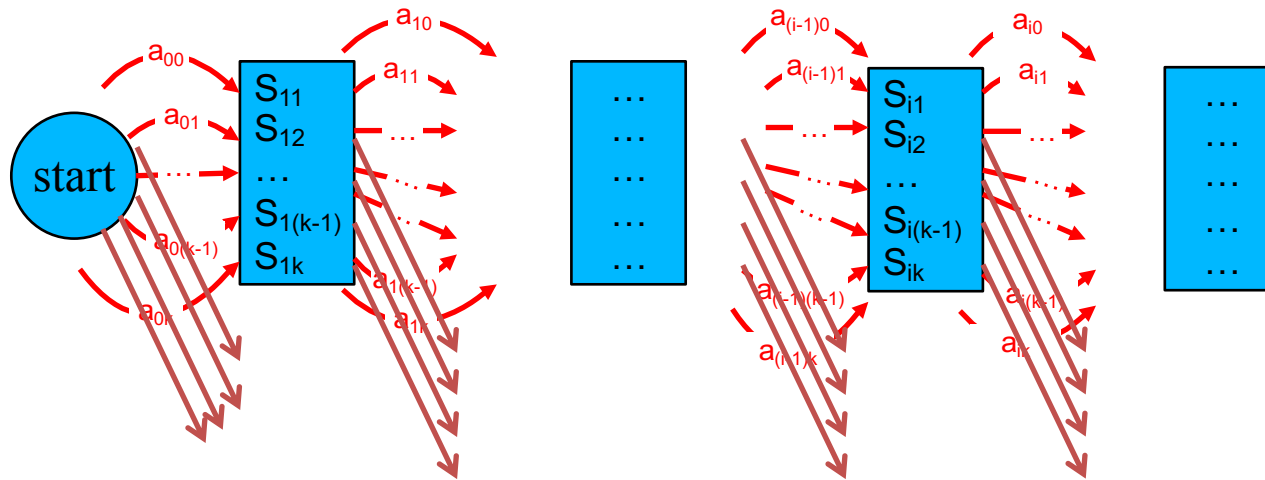
# Beam-search decoding



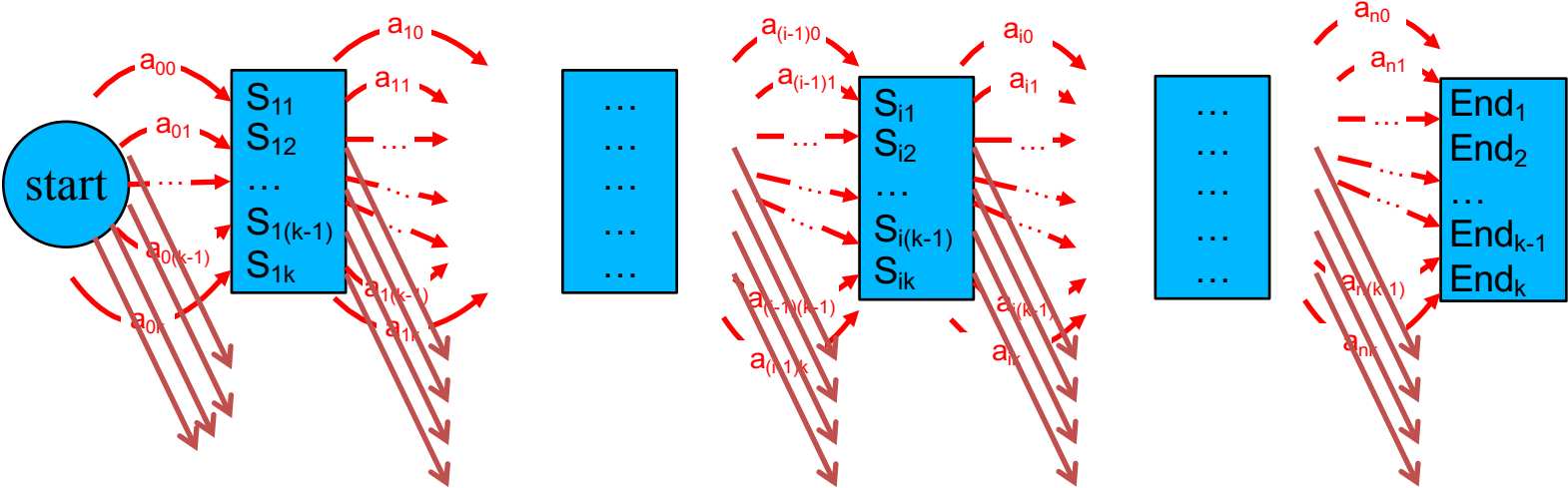
# Beam-search decoding



# Beam-search decoding



# Beam-search decoding



# Beam-search decoding

**function** BEAM-SEARCH(*problem, agenda, candidates, B*)

*candidates* ← {STARTITEM(*problem*)}

*agenda* ← CLEAR(*agenda*)

**loop do**

**for each** *candidate* **in** *candidates*

*agenda* ← INSERT(EXPAND(*candidate, problem*), *agenda*)

*best* ← TOP(*agenda*)

**if** GOALTEST(*problem, best*)

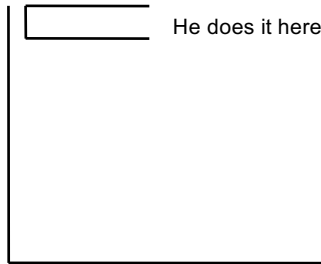
**then return** *best*

*candidates* ← TOP-B(*agenda, B*)

*agenda* ← CLEAR(*agenda*)

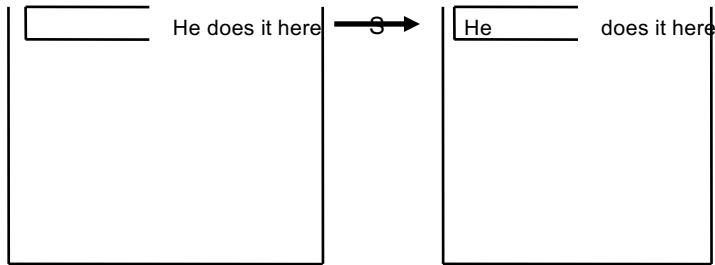
# Beam-search decoding

- Our parser
- Decoding



# Beam-search decoding

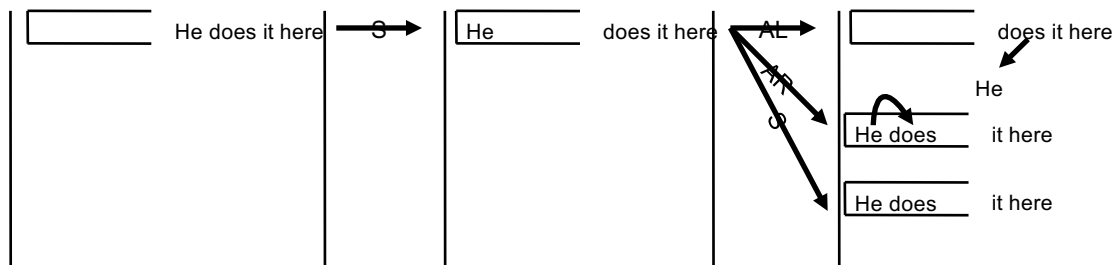
- Our parser
- Decoding





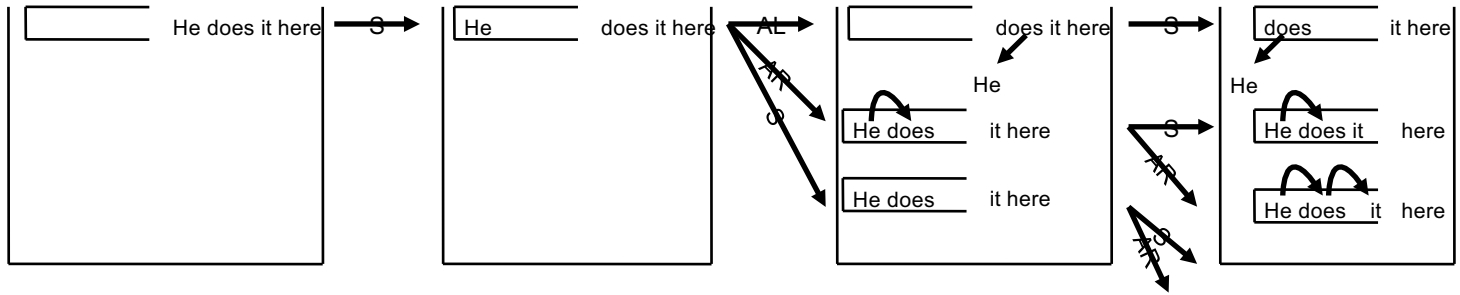
# Beam-search decoding

- Our parser
- Decoding



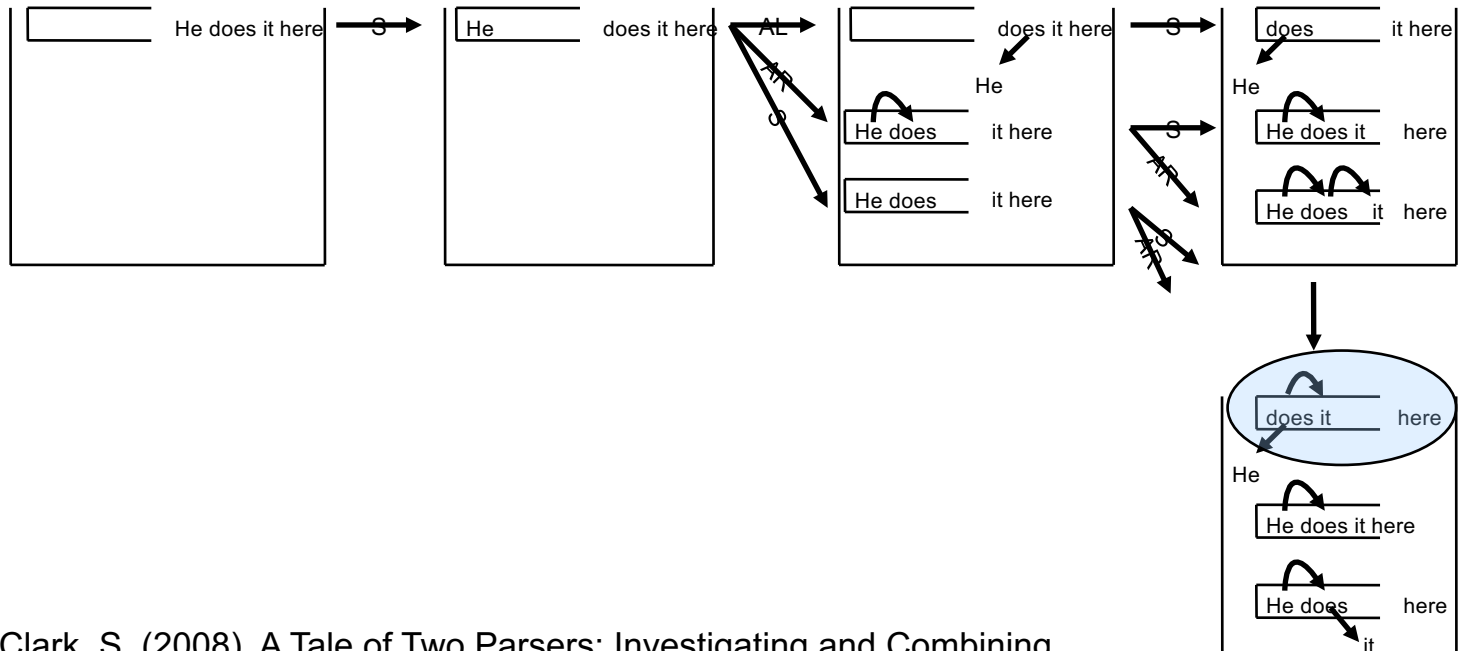
# Beam-search decoding

- Our parser
- Decoding



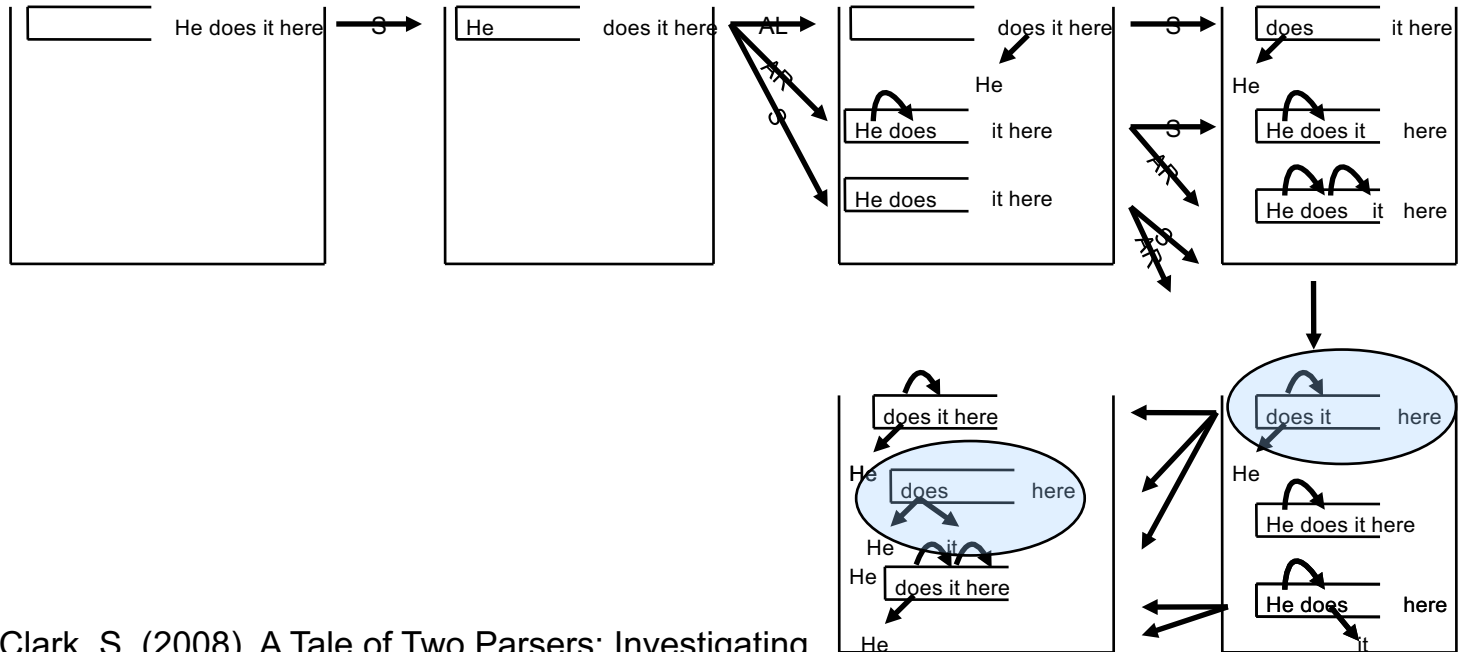
# Beam-search decoding

- Our parser
- Decoding



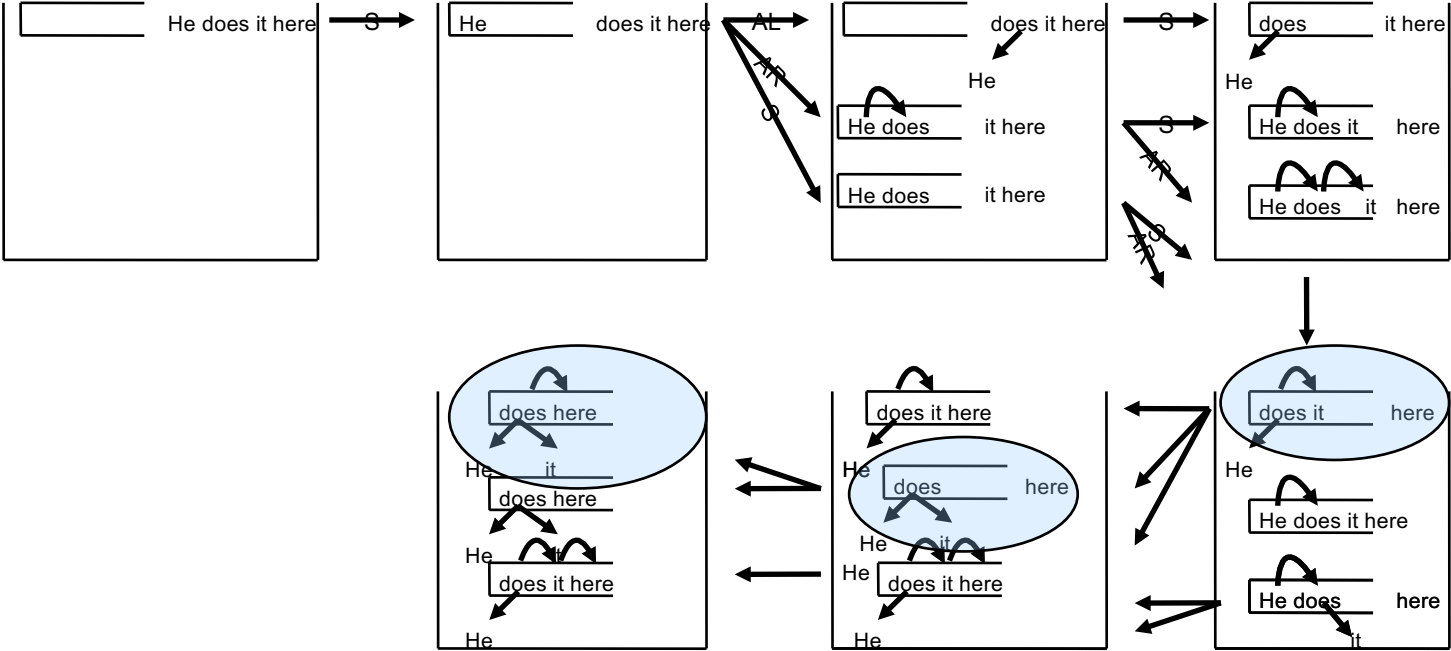
# Beam-search decoding

- Our parser
- Decoding



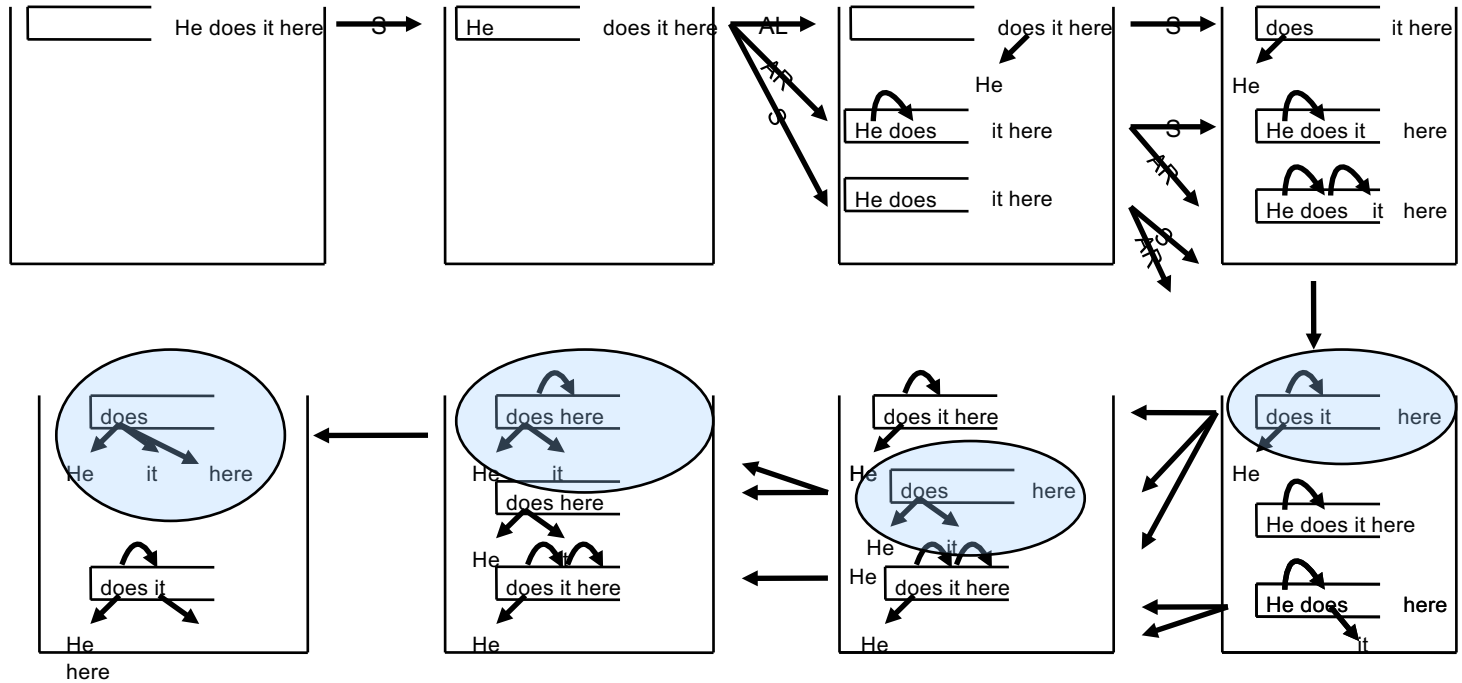
# Beam-search decoding

- Our parser
- Decoding



# Beam-search decoding

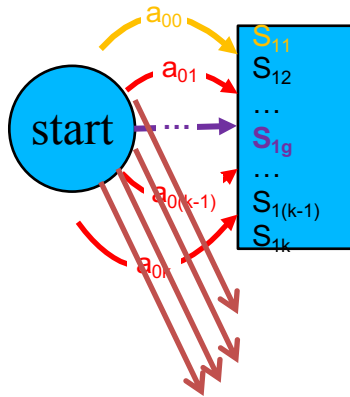
- Our parser
- Decoding



# Online learning

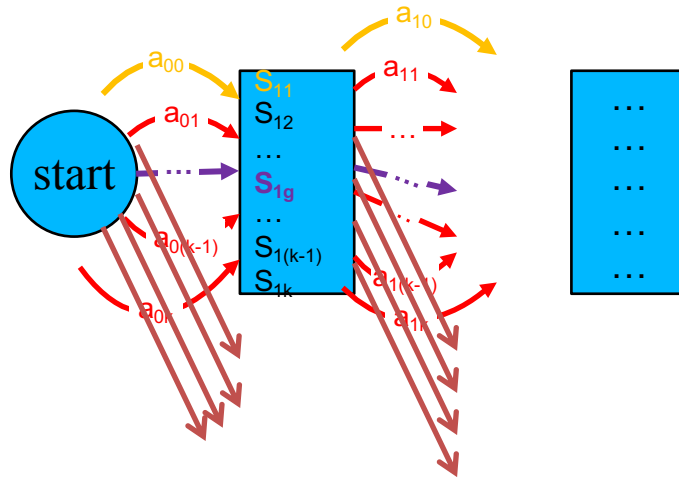


# Online learning

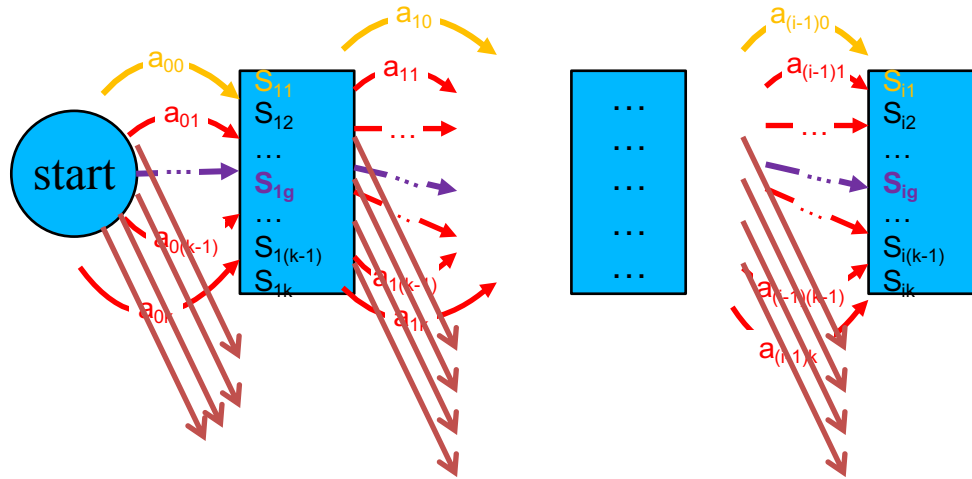




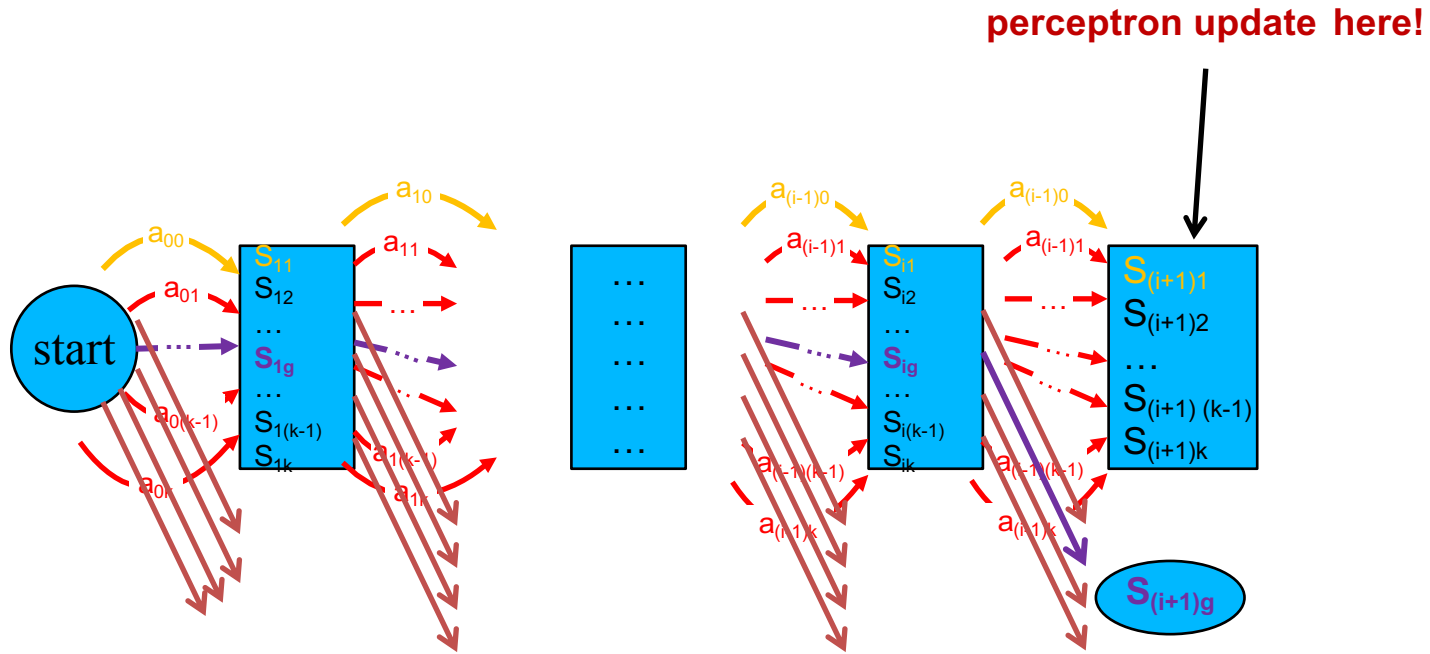
# Online learning



# Online learning



# Online learning



# Online learning

**Inputs:** training examples  $(x_i, y_i = \{S_0^i S_1^i \cdots S_m^i\})$  is a state sequence  $_1^N$

**Initialization:** set  $\vec{w} = 0$

**Algorithm:**

**for**  $r = 1 \cdots P, i = 1 \cdots N$  **do**

$candidates \leftarrow \{S_0^i\}$

$agenda \leftarrow CLEAR(agenda)$

**for**  $k = 1 \cdots m, m$  corresponds to a specific training example. **do**

**for** each *candidate* in *candidates* **do**

$agenda \leftarrow INSERT(EXPAND(candidate), agenda)$

$candidates \leftarrow TOP - B(agenda, B)$

$best \leftarrow TOP(agenda)$

**if**  $S_k^i$  is not in *candidates* or ( $best \neq S_m^i$  and  $k$  equals  $m$ ) **then**

$\vec{w} = \vec{w} + \Phi(S_k^i) - \Phi(best)$

**end if**

**end for**

**end for**

**end for**

**Output:**  $\vec{w}$

# The main strengths

- Fast
- Arbitrary nonlocal features
- Learning fixes search

# State-of-the-art results

- Chinese

- Word segmentation

- Yue Zhang and Stephen Clark. Chinese Segmentation Using a Word-Based Perceptron Algorithm. In proceedings of ACL 2007. Prague, Czech Republic. June.

# State-of-the-art results

- Chinese

- Joint segmentation and POS-tagging

- Yue Zhang and Stephen Clark. Joint Word Segmentation and POS Tagging Using a Single Perceptron. In proceedings of ACL 2008. Ohio, USA. June.

- Yue Zhang and Stephen Clark. A Fast Decoder for Joint Word Segmentation and POS-tagging Using a Single Discriminative Model. In proceedings of EMNLP 2010. Massachusetts, USA. October.

# State-of-the-art results

- Chinese

- Joint segmentation, POS-tagging and chunking

- Chen Lyu, Yue Zhang and Donghong Ji. Joint Word Segmentation, POS-Tagging and Syntactic Chunking. In Proceedings of the AAIL 2016, Phoenix, Arizona, USA, February



# State-of-the-art results

- Chinese

- Joint segmentation, POS-tagging and dependency parsing

- Meishan Zhang, Yue Zhang, Wanxiang Che and Ting Liu. Character-Level Chinese Dependency Parsing. In Proceedings of ACL 2014. Baltimore, USA, June.

# State-of-the-art results

- Chinese

- Joint segmentation, POS-tagging and constituent parsing

- Meishan Zhang, Yue Zhang, Wanxiang Che and Ting Liu. Chinese Parsing Exploiting Characters. In proceedings of ACL 2013. Sophia, Bulgaria. August.

# State-of-the-art results

- Chinese

- Joint segmentation, POS-tagging and normalization

- Tao Qian, Yue Zhang, Meishan Zhang and Donghong Ji. A Transition-based Model for Joint Segmentation, POS-tagging and Normalization. In proceedings of EMNLP 2015, Lisboa, Portugal, September.

# State-of-the-art results

- All Languages

- Constituent parsing

- Yue Zhang and Stephen Clark. Transition-Based Parsing of the Chinese Treebank Using a Global Discriminative Model. In proceedings of IWPT 2009. Paris, France. October.

- Muhua Zhu, Yue Zhang, Wenliang Chen, Min Zhang and Jingbo Zhu. Fast and Accurate Shift-Reduce Constituent Parsing. In proceedings of ACL 2013. Sophia, Bulgaria. August.

# State-of-the-art results

- All Languages

- Dependency parsing

- Yue Zhang and Stephen Clark. Joint Word Segmentation and POS Tagging Using a Single Perceptron. In proceedings of ACL 2008. Ohio, USA. June.

- Yue Zhang and Joakim Nivre. Transition-Based Dependency Parsing with Rich Non-Local Features. In proceedings of ACL 2011, short papers. Portland, USA. June.

- Yue Zhang and Joakim Nivre. Analyzing the Effect of Global Learning and Beam-Search for Transition-Based Dependency Parsing. In proceedings of COLING 2012, posters. Mumbai, India. December.

- Ji Ma, Yue Zhang and Jingbo Zhu. Punctuation Processing for Projective Dependency Parsing. In Proceedings of ACL 2014. Baltimore, USA, June.

# State-of-the-art results

- All Languages

- CCG parsing

- Yue Zhang and Stephen Clark. Shift-Reduce CCG Parsing. In proceedings of ACL 2011. Portland, USA. June.

- Wenduan Xu, Stephen Clark and Yue Zhang. Shift-Reduce CCG Parsing with a Dependency Model. In Proceedings of ACL 2014. Baltimore, USA, June.

# State-of-the-art results

- All Languages

- Natural language synthesis

- Yijia Liu, Yue Zhang, Wanxiang Che and Bing Qin. Transition-Based Syntactic Linearization. In Proceedings of NAACL 2015, Denver, Colorado, USA, May.

- Jiangming Liu and Yue Zhang. An Empirical Comparison Between N-gram and Syntactic Language Models for Word Ordering. In proceedings of EMNLP 2015, Lisboa, Portugal, September.

- Ratish Puduppully, Yue Zhang and Manish Shrivastava. Transition-Based Syntactic Linearization with Lookahead Features. In Proceedings of the NAACL 2016, San Diego, USA, June.

# State-of-the-art results

- All Languages

- Joint morphological generation and text linearization

- Linfeng Song, Yue Zhang, Kai Song and Qun Liu. Joint Morphological Generation and Syntactic Linearization. In Proceedings of AAI 2014. Quebec City, Canada, July.



# State-of-the-art results

- All Languages

- Joint entity and relation extraction

- Fei Li, Yue Zhang, Meishan Zhang and Donghong Ji. Joint Models for Extracting Adverse Drug Events from Biomedical Text. In Proceedings of IJCAI 2016. New York City, USA, July.

# Part 5.2: A Neural Network Version

# Neural Network Model

- Use NN to substitute perceptron
- Why?
  - Better non-linear power
  - Unsupervised word embeddings
  - Automatic feature combination
  - Shown useful in greedy models

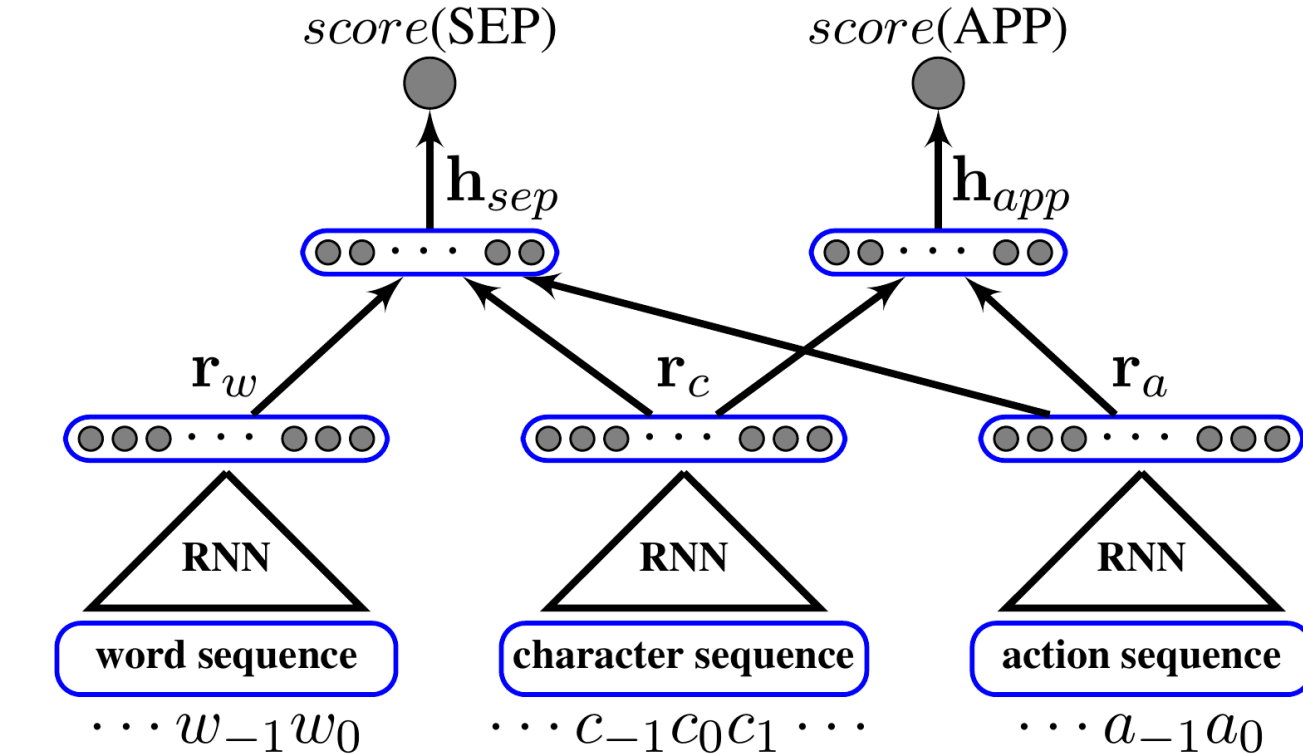
# Word segmentation

step	action	buffer( $\dots w_{-1}w_0$ )	queue( $c_0c_1\dots$ )
0	-	$\phi$	中 国 ...
1	<i>SEP</i>	中	国 外 ...
2	<i>APP</i>	中国	外 企 ...
3	<i>SEP</i>	中国 外	企 业 ...
4	<i>APP</i>	中国 外企	业 务 ...
5	<i>SEP</i>	中国 外企 业	务 发 ...
6	<i>APP</i>	中国 外企 业务	发 展 ...
7	<i>SEP</i>	... 业务 发	展 迅 速
8	<i>APP</i>	... 业务 发展	迅 速
9	<i>SEP</i>	... 发展 迅	速
10	<i>APP</i>	... 发展 迅速	$\phi$

# Word segmentation

<b>Feature templates</b>	<b>Action</b>
$c_{-1}c_0$	<i>APP, SEP</i>
$w_{-1}, w_{-1}w_{-2}, w_{-1}c_0, w_{-2}len(w_{-1})$ $start(w_{-1})c_0, end(w_{-1})c_0$ $start(w_{-1})end(w_{-1}), end(w_{-2})end(w_{-1})$ $w_{-2}len(w_{-1}), len(w_{-2})w_{-1}$ $w_{-1}, \text{ where } len(w_{-1}) = 1$	<i>SEP</i>

# Word segmentation



# Word segmentation

Models	P	R	F
word-based models			
discrete	95.29	95.26	95.28
neural	95.34	94.69	95.01
combined	<b>96.11</b>	<b>95.79</b>	<b>95.95</b>
character-based models			
discrete	95.38	95.12	95.25
neural	94.59	94.92	94.76
combined	95.63	95.60	95.61
other models			
Zhang et al. (2014)	N/A	N/A	95.71
Wang et al. (2011)	95.83	95.75	95.79
Zhang and Clark (2011)	95.46	94.78	95.13

Main results on CTB60 test dataset

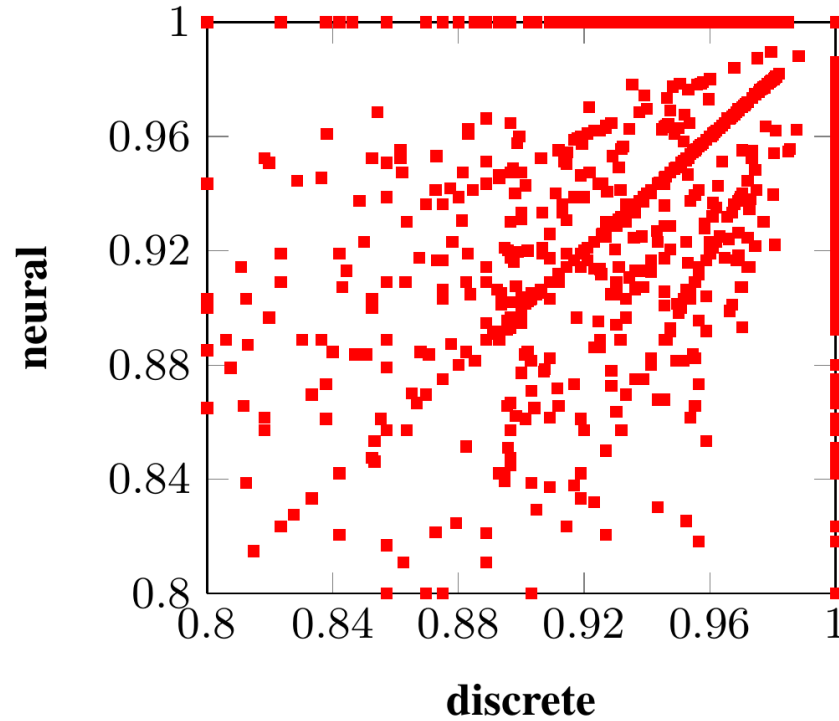
# Word segmentation

Models	PKU	MSR
our word-based models		
discrete	95.1	97.3
neural	95.1	97.0
combined	95.7	<b>97.7</b>
character-based models		
discrete	94.9	96.8
neural	94.4	97.2
combined	95.4	97.2
other models		
Cai and Zhao (2016)	95.5	96.5
Ma and Hinrichs (2015)	95.1	96.6
Pei et al. (2014)	95.2	97.2
Zhang et al. (2013a)	<b>96.1</b>	97.5
Sun et al. (2012)	95.4	97.4
Zhang and Clark (2011)	95.1	97.1
Sun (2010)	95.2	96.9
Sun et al. (2009)	95.2	97.3

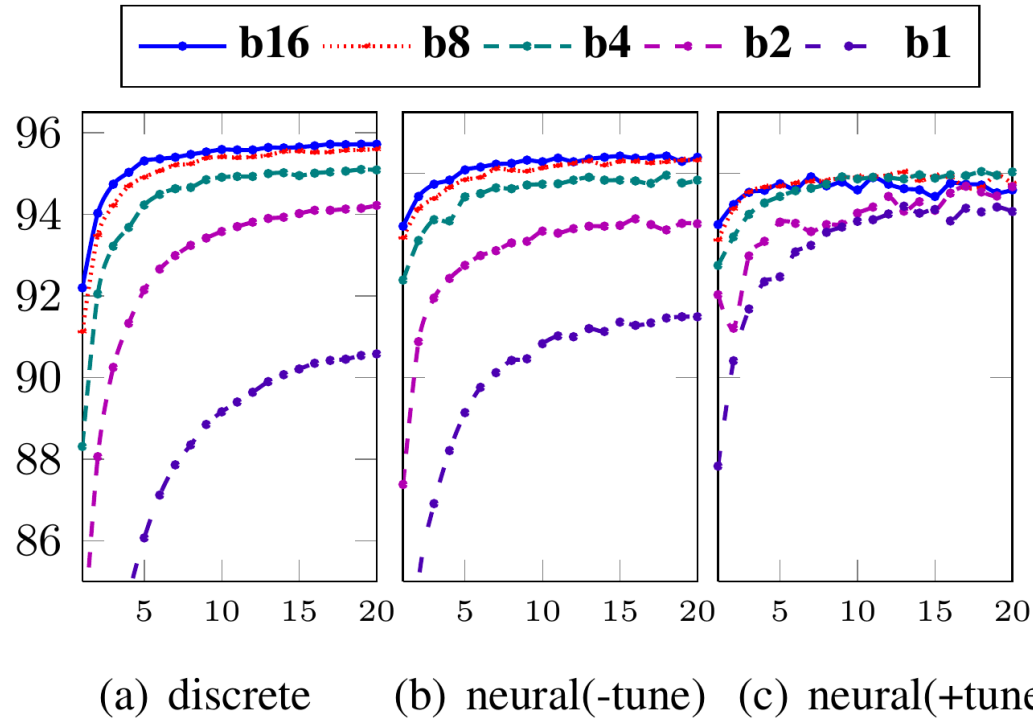
Main results on PKU and MSR test dataset



# Word segmentation



# Word segmentation



(a) discrete (b) neural(-tune) (c) neural(+tune)

# Word segmentation

- Cai and Zhao (2016) presents a similar idea

# Dependency Parsing

- Zhang & Nivre (2011)

$$y = \arg \max_{y' \in \text{GEN}(x)} \text{score}(y')$$

$$\text{score}(y) = \sum_{a \in y} \theta \cdot \Phi(a)$$

# Dependency Parsing

- Chen and Manning (2014)

$$h = (W_1x + b_1)^3$$

$$p = \text{softmax}(o)$$

$$o = W_2h$$

# Dependency Parsing

- What does not work

$$s(y) = \sum_{a \in y} \log p_a$$

$$L(\theta) = \max(0, \delta - s(y_g) + s(y_p)) + \frac{\lambda}{2} \|\theta\|^2$$

# Dependency Parsing

- Sentence-level log likelihood

$$p(y_i | x, \theta) = \frac{e^{f(x, \theta)_i}}{\sum_{y_j \in \text{GEN}(x)} e^{f(x, \theta)_j}}$$

$$f(x, \theta)_i = \sum_{a_k \in y_i} o(x, y_i, k, a_k)$$

# Dependency Parsing

- Contrastive Estimation

$$\begin{aligned} L(\theta) &= - \sum_{(x_i, y_i) \in (X, Y)} \log p(y_i | x_i, \theta) \\ &= - \sum_{(x_i, y_i) \in (X, Y)} \log \frac{e^{f(x_i, \theta)_i}}{Z(x_i, \theta)} \\ &= \sum_{(x_i, y_i) \in (X, Y)} \log Z(x_i, \theta) - f(x_i, \theta)_i \end{aligned}$$

$$Z(x, \theta) = \sum_{y_j \in \text{GEN}(x)} e^{f(x, \theta)_j}$$



# Dependency Parsing

- Contrastive Estimation

$$\begin{aligned}L'(\theta) &= - \sum_{(x_i, y_i) \in (X, Y)} \log p'(y_i | x_i, \theta) \\ &= - \sum_{(x_i, y_i) \in (X, Y)} \log \frac{e^{f(x_i, \theta)_i}}{Z'(x_i, \theta)} \\ &= \sum_{(x_i, y_i) \in (X, Y)} \log Z'(x_i, \theta) - f(x_i, \theta)_i \\ Z'(x, \theta) &= \sum_{y_j \in \text{BEAM}(x)} e^{f(x, \theta)_j}\end{aligned}$$

# Dependency Parsing

- Results

Description	UAS	
Baseline	91.63	
	structured	greedy
beam = 1	74.90	91.63
beam = 4	84.64	91.92
beam = 16	91.53	91.90
beam = 64	93.12	91.84
beam = 100	93.23	91.81

# Dependency Parsing

- Results

Description	UAS
greedy neural parser	91.47
ranking model	89.08
beam contrastive learning	93.28

# Dependency Parsing

## •Results

System	UAS	LAS	Speed	
baseline greedy parser	91.47	90.43	0.001	
Huang and Sagae (2010)	92.10		0.04	
Zhang and Nivre (2011)	92.90	91.80	0.03	
Choi and McCallum (2013)	92.96	91.93	0.009	
Ma et al. (2014)	93.06			
Bohnet and Nivre (2012) <sup>†‡</sup>	93.67	92.68	0.4	
Suzuki et al. (2009) <sup>†</sup>	93.79			
Koo et al. (2008) <sup>†</sup>	93.16			
Chen et al. (2014) <sup>†</sup>	93.77			
beam size				
training	decoding			
100	100	<b>93.28</b>	<b>92.35</b>	0.07
100	64	93.20	92.27	0.04
100	16	92.40	91.95	0.01

# Google

- Andor et al. follows this method
  - Offers theorem
  - Tries more tasks
  - Get better results

# Google

- Dependency parsing

Method	WSJ		Union-News		Union-Web		Union-QTB	
	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS
Martins et al. (2013)*	92.89	90.55	93.10	91.13	88.23	85.04	94.21	91.54
Zhang and McDonald (2014)*	93.22	91.02	93.32	91.48	88.65	85.59	93.37	90.69
Weiss et al. (2015)	93.99	92.05	93.91	92.25	89.29	86.44	94.17	92.06
Alberti et al. (2015)	94.23	92.36	94.10	92.55	89.55	86.85	94.74	93.04
Our Local (B=1)	92.95	91.02	93.11	91.46	88.42	85.58	92.49	90.38
Our Local (B=32)	93.59	91.70	93.65	92.03	88.96	86.17	93.22	91.17
Our Global (B=32)	<b>94.61</b>	<b>92.79</b>	<b>94.44</b>	<b>92.93</b>	<b>90.17</b>	<b>87.54</b>	<b>95.40</b>	<b>93.64</b>
Parsey McParseface (B=8)	-	-	94.15	92.51	89.08	86.29	94.77	93.17

# Google

- Dependency parsing

Method	Catalan		Chinese		Czech		English		German		Japanese		Spanish	
	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS
Best Shared Task Result	-	87.86	-	79.17	-	80.38	-	89.88	-	87.48	-	92.57	-	87.64
Ballesteros et al. (2015)	90.22	86.42	80.64	76.52	79.87	73.62	90.56	88.01	88.83	86.10	93.47	92.55	90.38	86.59
Zhang and McDonald (2014)	91.41	87.91	82.87	78.57	86.62	80.59	92.69	90.01	89.88	87.38	92.82	91.87	90.82	87.34
Lei et al. (2014)	91.33	87.22	81.67	76.71	88.76	81.77	92.75	90.00	90.81	87.81	<b>94.04</b>	91.84	91.16	87.38
Bohnet and Nivre (2012)	92.44	89.60	82.52	78.51	88.82	83.73	92.87	90.60	<b>91.37</b>	<b>89.38</b>	93.67	92.63	92.24	89.60
Alberti et al. (2015)	92.31	89.17	83.57	79.90	88.45	83.57	92.70	90.56	90.58	88.20	93.99	<b>93.10</b>	92.26	89.33
Our Local (B=1)	91.24	88.21	81.29	77.29	85.78	80.63	91.44	89.29	89.12	86.95	93.71	92.85	91.01	88.14
Our Local (B=16)	91.91	88.93	82.22	78.26	86.25	81.28	92.16	90.05	89.53	87.4	93.61	92.74	91.64	88.88
Our Global (B=16)	<b>92.67</b>	<b>89.83</b>	<b>84.72</b>	<b>80.85</b>	<b>88.94</b>	<b>84.56</b>	<b>93.22</b>	<b>91.23</b>	90.91	89.15	93.65	92.84	<b>92.62</b>	<b>89.95</b>

# Google

- POS-tagging

Method	En	En-Union			CoNLL '09							Avg
	WSJ	News	Web	QTB	Ca	Ch	Cz	En	Ge	Ja	Sp	-
Linear CRF	97.17	97.60	94.58	96.04	98.81	94.45	98.90	97.50	97.14	97.90	98.79	97.17
Ling et al. (2015)	<b>97.78</b>	97.44	94.03	96.18	98.77	94.38	99.00	97.60	<b>97.84</b>	97.06	98.71	97.16
Our Local (B=1)	97.44	97.66	94.46	96.59	98.91	94.56	98.96	97.36	97.35	98.02	98.88	97.29
Our Local (B=8)	97.45	97.69	94.46	96.64	98.88	94.56	98.96	97.40	97.35	98.02	98.89	97.30
Our Global (B=8)	97.44	<b>97.77</b>	<b>94.80</b>	<b>96.86</b>	<b>99.03</b>	<b>94.72</b>	<b>99.02</b>	<b>97.65</b>	97.52	<b>98.37</b>	<b>98.97</b>	<b>97.47</b>
Parsey McParseface	-	97.52	94.24	96.45	-	-	-	-	-	-	-	-



# Google

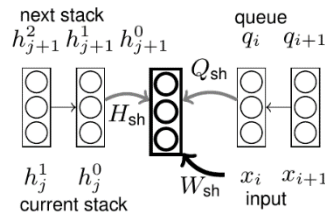
- Compression

Method	Generated corpus		Human eval	
	A	F1	read	info
Filippova et al. (2015)	<b>35.36</b>	<b>82.83</b>	4.66	4.03
Automatic	-	-	4.31	3.77
Our Local (B=1)	30.51	78.72	4.58	4.03
Our Local (B=8)	31.19	75.69	-	-
Our Global (B=8)	35.16	81.41	<b>4.67</b>	<b>4.07</b>

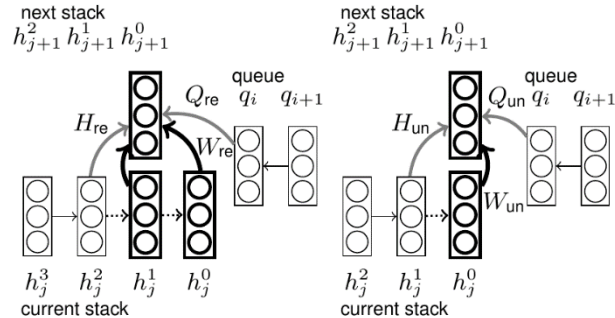
# Part 5.3: Similar methods by others

# Other methods ( I )

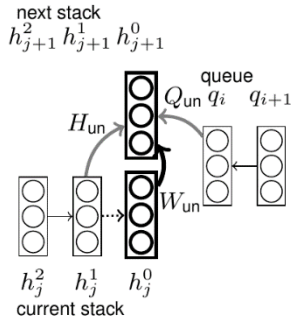
- Constituent parsing



(a) shift- $X$  action



(b) reduce- $X$  action



(c) unary- $X$  action

# Other methods ( I )

- Update at max-violation

$$j^* = \arg \min_j \left\{ \rho_{\theta}(y_0^j) - \max_{\mathbf{d} \in B_j} \rho_{\theta}(\mathbf{d}) \right\}$$

- Using expected loss from all violations

$$L(\mathbf{w}, \mathbf{y}; \mathbf{B}, \theta) = \max \left\{ 0, 1 - \rho_{\theta}(y_0^{j^*}) + \mathbb{E}_{\tilde{B}_{j^*}}[\rho_{\theta}] \right\}$$

$$\tilde{B}_{j^*} = \left\{ \mathbf{d} \in B_{j^*} \mid \rho_{\theta}(\mathbf{d}) > \rho_{\theta}(y_0^{j^*}) \right\}$$

$$p_{\theta}(\mathbf{d}) = \frac{\exp(\rho_{\theta}(\mathbf{d}))}{\sum_{\mathbf{d}' \in \tilde{B}_{j^*}} \exp(\rho_{\theta}(\mathbf{d}'))}$$

$$\mathbb{E}_{\tilde{B}_{j^*}}[\rho_{\theta}] = \sum_{\mathbf{d} \in \tilde{B}_{j^*}} p_{\theta}(\mathbf{d}) \rho_{\theta}(\mathbf{d}).$$

# Other methods ( *I* )

parser	test
Collins (Collins, 1997)	87.8
Berkeley (Petrov and Klein, 2007)	90.1
SSN (Henderson, 2004)	90.1
ZPar (Zhu et al., 2013)	90.4
CVG (Socher et al., 2013)	90.4
Charniak-R (Charniak and Johnson, 2005)	<b>91.0</b>
This work: TNCP	90.7

## Other methods ( *I* )

parser	test
ZPar (Zhu et al., 2013)	83.2
Berkeley (Petrov and Klein, 2007)	83.3
Joint (Wang and Xue, 2014)	<b>84.9</b>
This work: TNCP	84.3

## Other methods ( II )

- CCG Parsing
- expected F1 training

$$\begin{aligned} J(\theta) &= -\text{x}F1(\theta) \\ &= - \sum_{y_i \in \Lambda(x_n)} p(y_i|\theta) F1(\Delta_{y_i}, \Delta_{x_n}^G) \end{aligned}$$

$$p(y_i|\theta) = \frac{\exp\{\rho(y_i)\}}{\sum_{y \in \Lambda(x_n)} \exp\{\rho(y)\}}$$

## Other methods ( II )

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta} &= - \sum_{y_i \in \Lambda(x_n)} \sum_{y_{ij} \in y_i} \frac{\partial J(\theta)}{\partial s_\theta(y_{ij})} \frac{\partial s_\theta(y_{ij})}{\partial \theta} \\ &= - \sum_{y_i \in \Lambda(x_n)} \sum_{y_{ij} \in y_i} \delta_{y_{ij}} \frac{\partial s_\theta(y_{ij})}{\partial \theta},\end{aligned}$$

$$\begin{aligned}\delta_{y_{ij}} &= - \frac{\partial xF1(\theta)}{\partial s_\theta(y_{ij})} \\ &= - \frac{\partial(G(\theta)/Z(\theta))}{\partial s_\theta(y_{ij})} \\ &= \frac{G(\theta)Z'(\theta) - G'(\theta)Z(\theta)}{Z^2(\theta)} \\ &= \frac{\exp\{\rho(y_i)\}}{Z(\theta)} (xF1(\theta) - F1(\Delta_{y_i}, \Delta_{x_n}^G)) \frac{1}{s_\theta(y_{ij})} \\ &= p(y_i|\theta) (xF1(\theta) - F1(\Delta_{y_i}, \Delta_{x_n}^G)) \frac{1}{s_\theta(y_{ij})},\end{aligned}$$

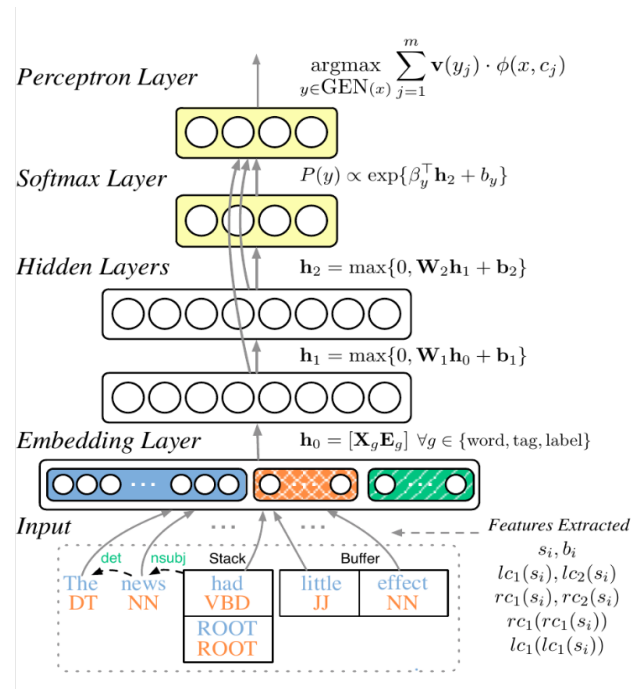


# Other methods ( II )

Model	Section 00				Section 23				Speed
	LP	LR	LF	CAT	LP	LR	LF	CAT	
C&C (normal)	85.18	82.53	83.83	92.39	85.45	83.97	84.70	92.83	97.90
C&C (hybrid)	86.07	82.77	84.39	92.57	86.24	84.17	85.19	93.00	95.25
Zhang and Clark (2011) ( $b = 16$ )	87.15	82.95	85.00	92.77	87.43	83.61	85.48	93.12	-
Zhang and Clark (2011)* ( $b = 16$ )	86.76	83.15	84.92	92.64	87.04	84.14	85.56	92.95	49.54
Xu et al. (2014) ( $b = 128$ )	86.29	<b>84.09</b>	85.18	92.75	87.03	<b>85.08</b>	86.04	93.10	12.85
RNN-greedy ( $b = 1$ )	88.12	81.38	84.61	93.42	88.53	81.65	84.95	93.57	337.45
RNN-greedy ( $b = 6$ )	87.96	82.27	85.02	93.47	88.54	82.77	85.56	93.68	96.04
RNN-xF1 ( $b = 8$ )	<b>88.20</b>	83.40	<b>85.73</b>	<b>93.56</b>	<b>88.74</b>	84.22	<b>86.42</b>	<b>93.87</b>	67.65

# Other methods (III)

- Dependency parsing



## Other methods (*III*)

- Using Chen and Manning features for perceptron training
- Back-propagation pre-training

$$L(\Theta) = - \sum_j \log P(y_j | c_j, \Theta) + \lambda \sum_i \|\mathbf{W}_i\|_2^2$$

- Structured perceptron training

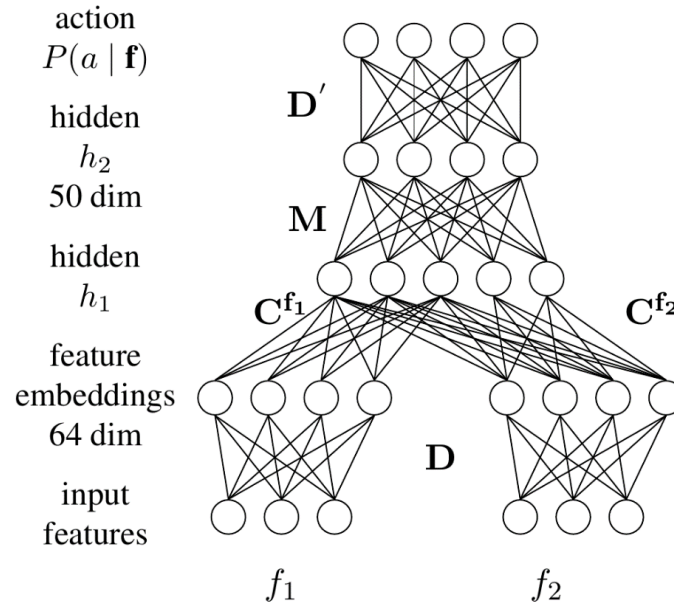
$$(h_1, h_2, P(y))$$

# Other methods ( *III* )

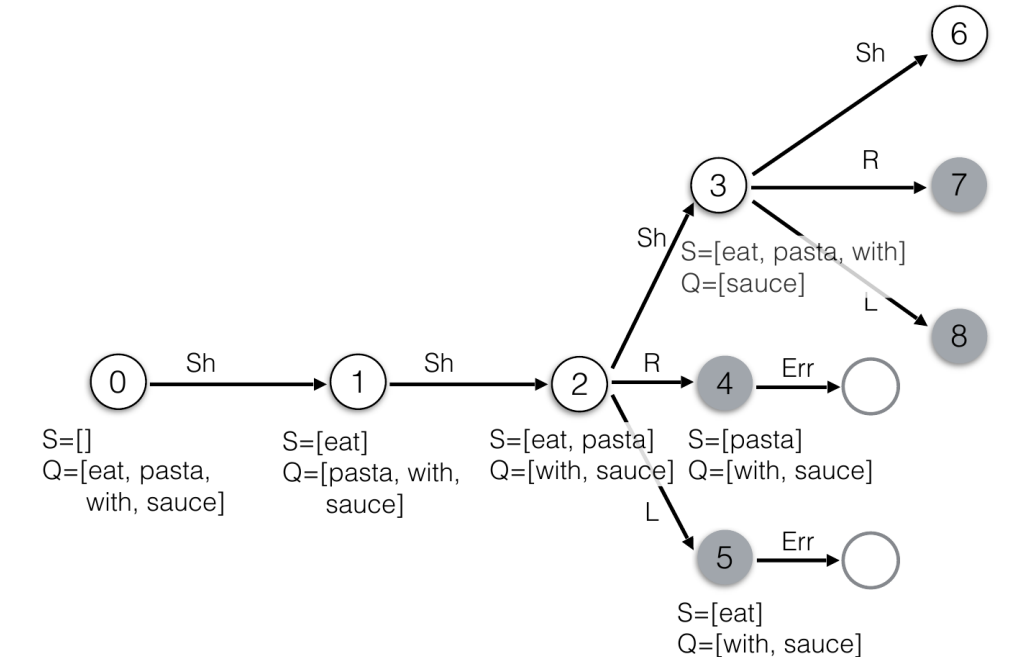
Method	UAS	LAS	Beam
<i>Graph-based</i>			
Bohnet (2010)	92.88	90.71	n/a
Martins et al. (2013)	92.89	90.55	n/a
Zhang and McDonald (2014)	93.22	91.02	n/a
<i>Transition-based</i>			
*Zhang and Nivre (2011)	93.00	90.95	32
Bohnet and Kuhn (2012)	93.27	91.19	40
Chen and Manning (2014)	91.80	89.60	1
S-LSTM (Dyer et al., 2015)	93.20	90.90	1
Our Greedy	93.19	91.18	1
Our Perceptron	<b>93.99</b>	<b>92.05</b>	8
<i>Tri-training</i>			
*Zhang and Nivre (2011)	92.92	90.88	32
Our Greedy	93.46	91.49	1
Our Perceptron	<b>94.26</b>	<b>92.41</b>	8

# Other methods (IV)

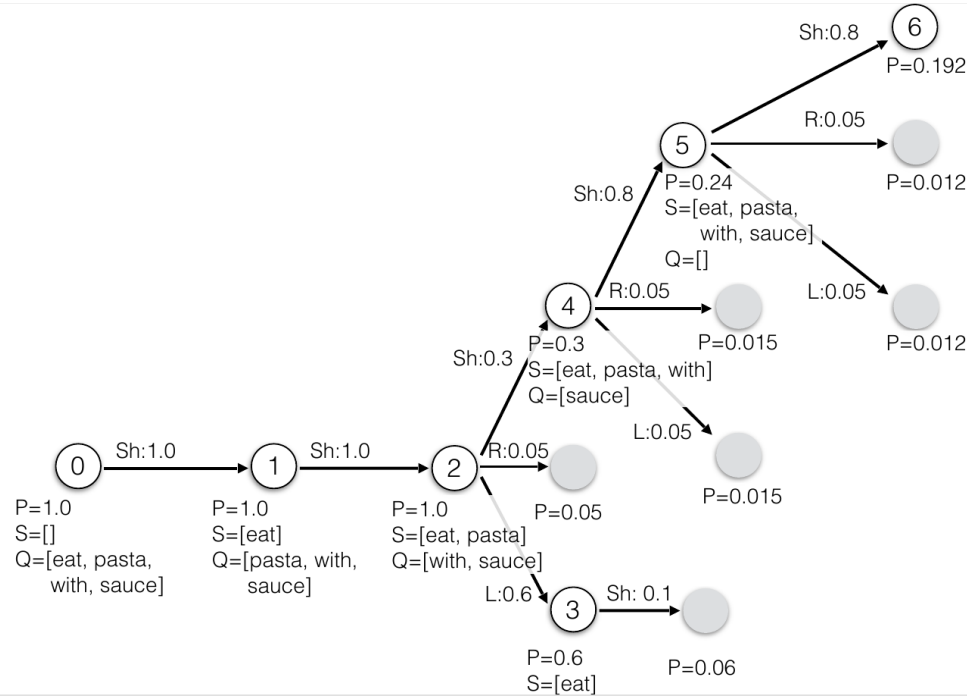
- Dependency parsing



# Other methods (IV)



# Other methods (IV)



## Other methods (*IV*)

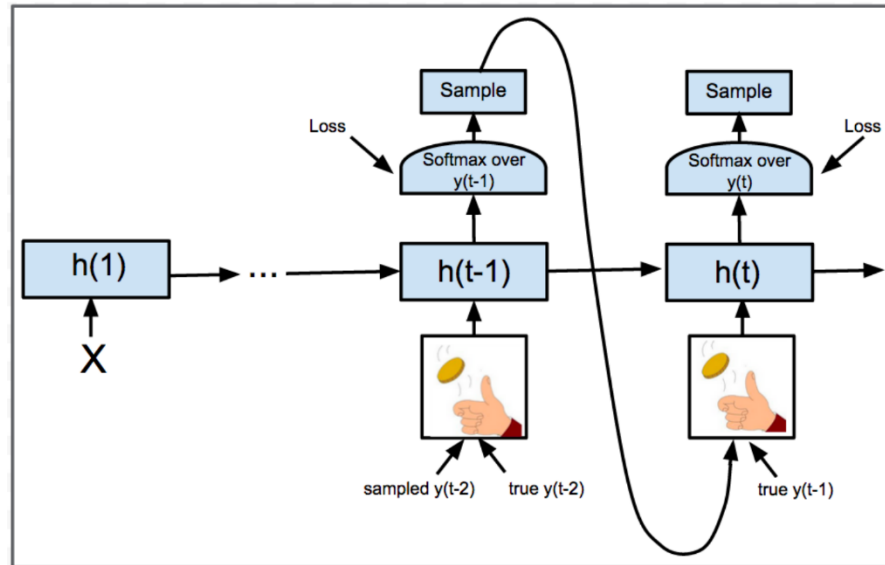
System	wsj23-S	wsj23-YM
<b>ErrSt-25-rand</b>	<b>92.17</b>	<b>92.16</b>
<b>ErrSt-25-pre*</b>	<b>93.61</b>	<b>93.21</b>
Chen & Manning*	91.8	–
Huang & Sagae	–	92.1
Zhang & Nivre	93.5	92.9
Weiss et al.*	93.99	–
Zhang & McDonald	93.71	93.57
Martins et al.	92.82	93.07
Koo et al. (dep2c)*	–	93.16



# Part 5.4: Beam-search Decoding for Sequence to Sequence Models

# Sequence to sequence ( $I$ )

- Scheduled Sampling



Beam Search Inference

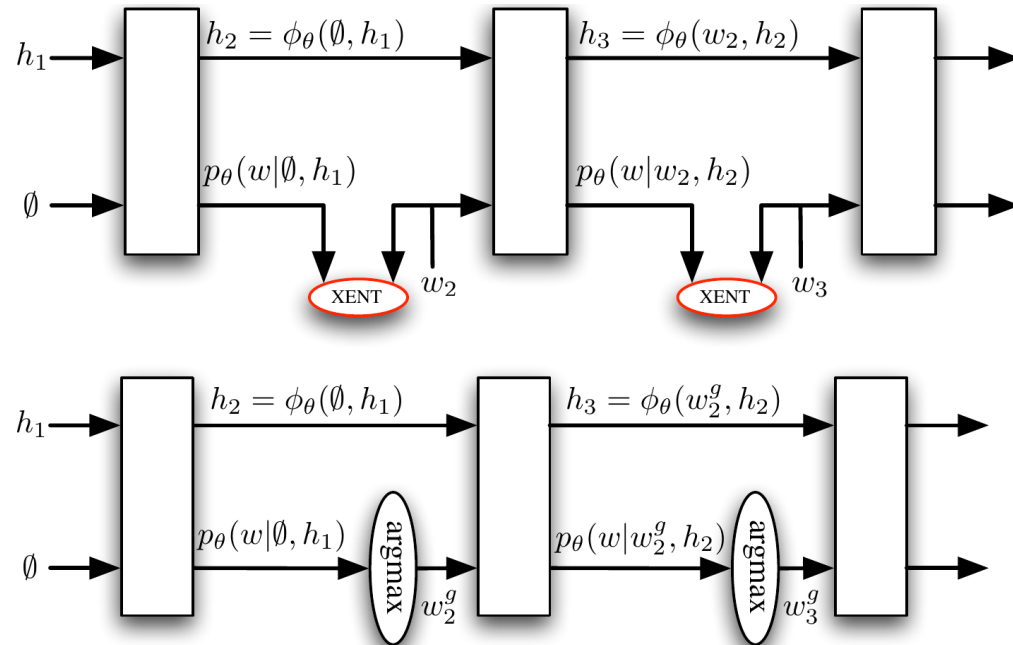
# Sequence to sequence ( $I$ )

- Scheduled Sampling

Approach	F1
Baseline LSTM	86.54
Baseline LSTM with Dropout	87.0
Always Sampling	-
Scheduled Sampling	<b>88.08</b>
Scheduled Sampling with Dropout	<b>88.68</b>

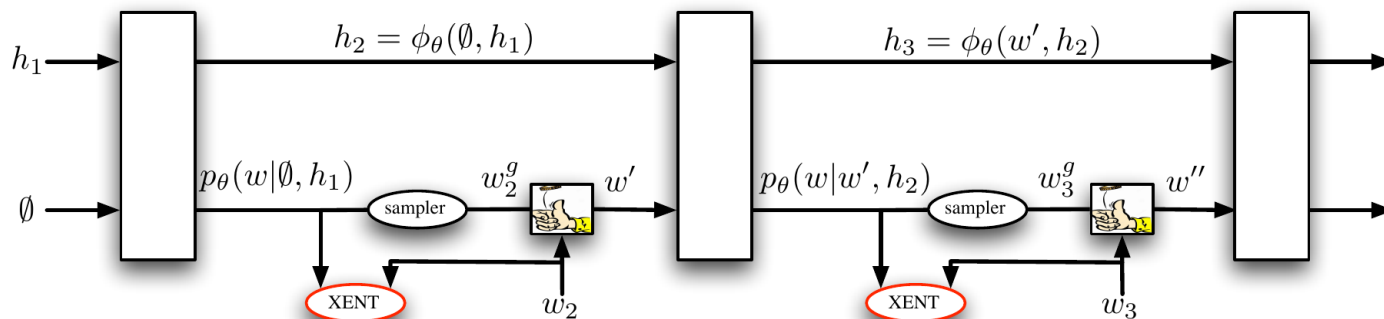
# Sequence to sequence ( II )

- Sequence-level training



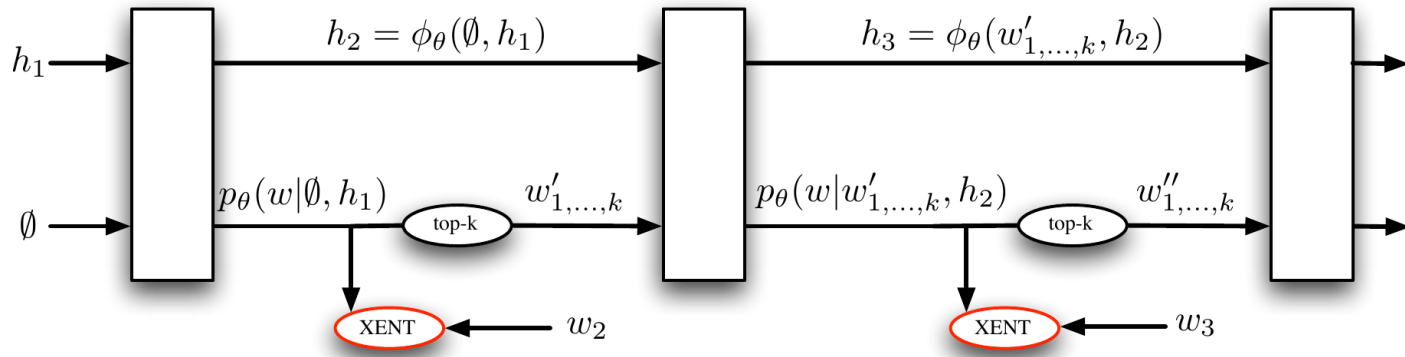
# Sequence to sequence ( II )

- Sequence-level training



# Sequence to sequence ( II )

- Sequence-level training



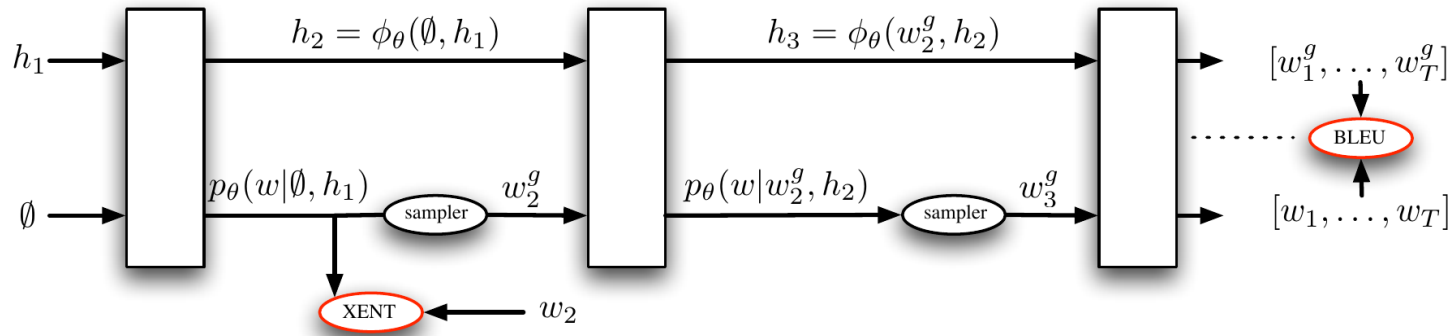
# Sequence to sequence ( II )

- Reinforce

$$L_{\theta} = - \sum_{w_1^g, \dots, w_T^g} p_{\theta}(w_1^g, \dots, w_T^g) r(w_1^g, \dots, w_T^g) = -\mathbb{E}_{[w_1^g, \dots, w_T^g] \sim p_{\theta}} r(w_1^g, \dots, w_T^g)$$

# Sequence to sequence ( II )

- Mixer





# Sequence to sequence ( II )

**Data:** a set of sequences with their corresponding context.

**Result:** RNN optimized for generation.

Initialize RNN at random and set  $N^{\text{XENT}}$ ,  $N^{\text{XE+R}}$  and  $\Delta$ ;

**for**  $s = T, 1, -\Delta$  **do**

**if**  $s == T$  **then**

        train RNN for  $N^{\text{XENT}}$  epochs using XENT only;

**else**

        train RNN for  $N^{\text{XE+R}}$  epochs. Use XENT loss in the first  $s$  steps, and REINFORCE (sampling from the model) in the remaining  $T - s$  steps;

**end**

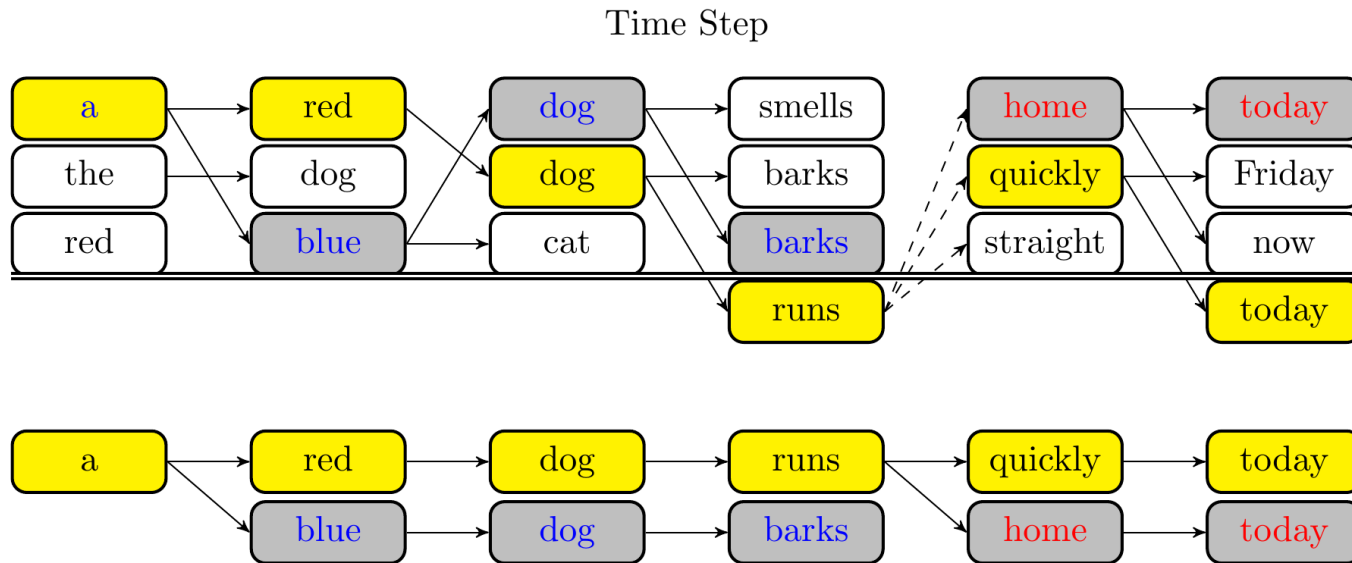
**end**

## Sequence to sequence ( II )

<i>TASK</i>	XENT	DAD	E2E	MIXER
<i>summarization</i>	13.01	12.18	12.78	<b>16.22</b>
<i>translation</i>	17.74	20.12	17.77	<b>20.73</b>
<i>image captioning</i>	27.8	28.16	26.42	<b>29.16</b>

# Sequence to sequence (III)

- Learning for Search



# Sequence to sequence (III)

$$\mathcal{L}(f) = \sum_{t=1}^T \Delta(\hat{y}_{1:t}^{(K)}) \left[ 1 - f(y_t, \mathbf{h}_{t-1}) + f(\hat{y}_t^{(K)}, \hat{\mathbf{h}}_{t-1}^{(K)}) \right]$$

# Sequence to sequence (*III*)

- Need greedy pre-training

# Sequence to sequence (*III*)

- Curriculum beam increase

## Sequence to sequence (*III*)

	Word Ordering (BLEU)		
	$K_{te} = 1$	$K_{te} = 5$	$K_{te} = 10$
seq2seq	25.2	29.8	31.0
BSO	28.0	33.2	34.3
ConBSO	<b>28.6</b>	<b>34.3</b>	<b>34.5</b>
LSTM-LM	15.4	-	26.8

# Sequence to sequence (III)

---

	Dependency Parsing (UAS/LAS)		
	$K_{te} = 1$	$K_{te} = 5$	$K_{te} = 10$
seq2seq	<b>87.33/82.26</b>	88.53/84.16	88.66/84.33
BSO	86.91/82.11	91.00/ <b>87.18</b>	91.17/ <b>87.41</b>
ConBSO	85.11/79.32	<b>91.25</b> /86.92	<b>91.57</b> /87.26
Andor	93.17/91.18	-	-

---



# Sequence to sequence (III)

	Machine Translation (BLEU)		
	$K_{te} = 1$	$K_{te} = 5$	$K_{te} = 10$
seq2seq	22.53	24.03	23.87
BSO, SB- $\Delta$	<b>23.83</b>	<b>26.36</b>	<b>25.48</b>
XENT	17.74	$\leq 20.5$	$\leq 20.5$
DAD	20.12	$\leq 22.5$	$\leq 23.0$
MIXER	20.73	-	$\leq 22.0$