

# Natural Language Processing

Yue Zhang  
Westlake University



## Chapter 15

# Neural Structured Prediction



- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization



# Contents

- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization



# Structure Prediction

- Sequence labelling (Chapter 7, 8)

Sentence	POS tag sequence
<i>Jamie went to the shop yesterday .</i>	<i>NNP VBD TO DT NN AD .</i>
<i>What would you like to eat ?</i>	<i>WP MD PRP VB TO VB .</i>
<i>Tim is talking with Mary .</i>	<i>NNP VBZ VBG IN NNP .</i>
<i>I really appreciate it .</i>	<i>PRP RB VBP PRP .</i>
<i>John is a famous athlete .</i>	<i>NNP VBZ DT JJ NN .</i>

- Sequence Segmentation ( Chapter 9)

Word segmentation

其中国外企业  
中国外企业务

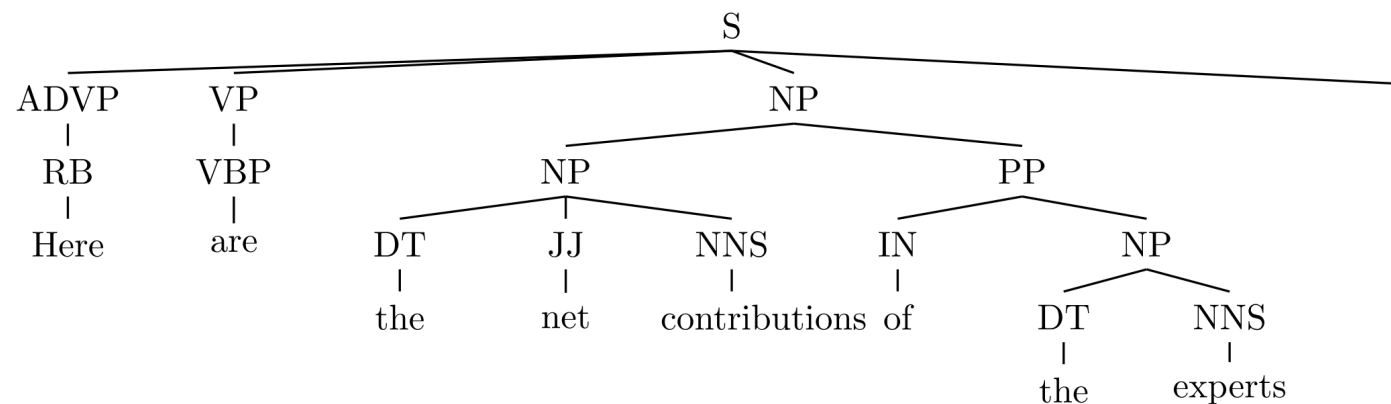
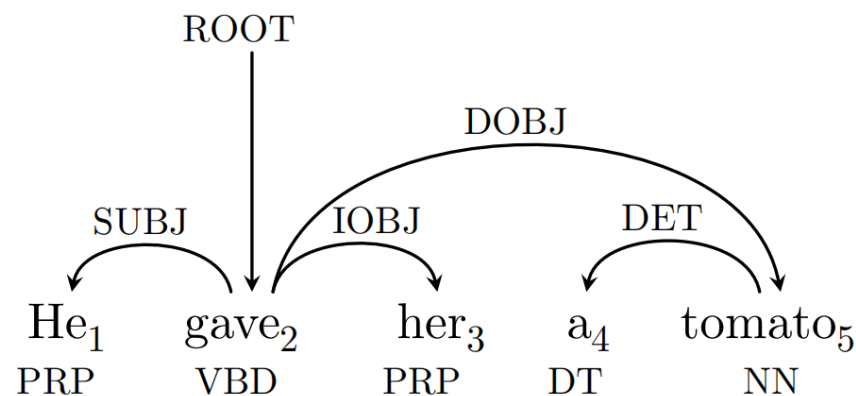
其中 国外 企业  
中国 外企 业务

Input	Output
<i>Michael Jordan is a Professor at University of Berkeley, located near Silicon Valley, USA.</i>	<i>[Michael Jordan]<sub>PER</sub> is a Professor at [University of Berkeley]<sub>ORG</sub>, located near [silicon valley]<sub>LOC</sub>, [USA]<sub>GPE</sub>.</i>
<i>Mary went to Chicago to meet her boyfriend John Smith.</i>	<i>[Mary]<sub>PER</sub> went to [Chicago]<sub>LOC</sub> to meet her boyfriend [John Smith]<sub>PER</sub>.</i>



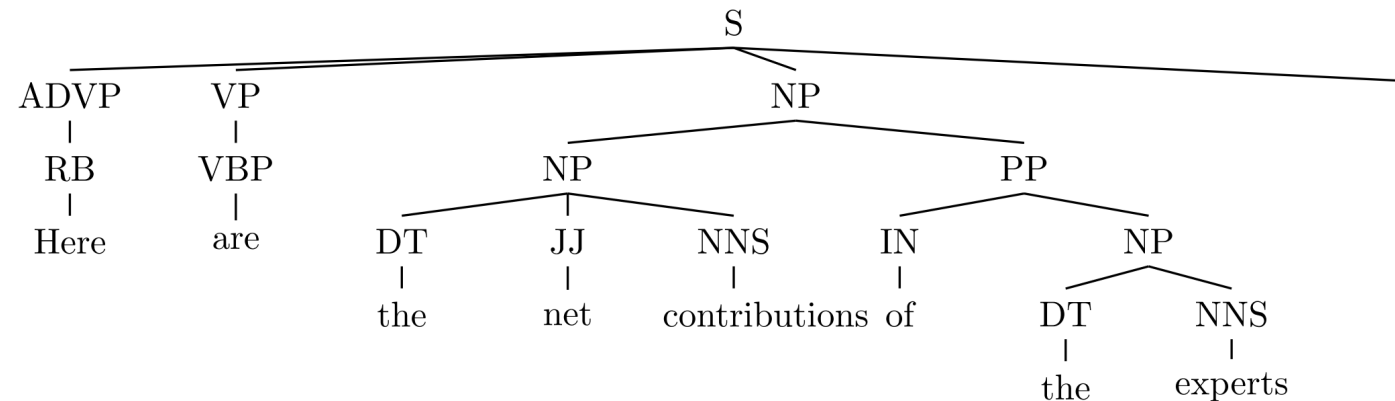
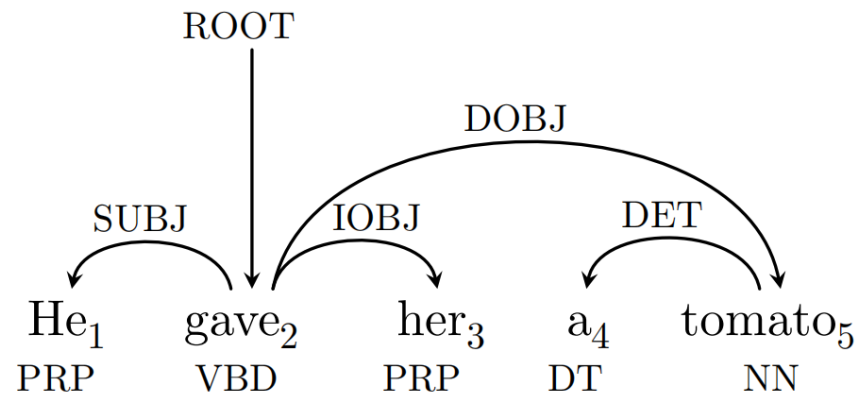
# Structure Prediction

- Sequence labelling (Chapter 7, 8)
- Sequence Segmentation (Chapter 9)
- Tree Structure Prediction (Chapter 10, 11)



# Structure Prediction

- Sequence labelling (Chapter 7, 8)
- Sequence Segmentation (Chapter 9)
- Tree Structure Prediction (Chapter 10, 11)



Graph-based models (Chapters 7 – 10) vs Transition-based models (Chapter 11)  
Local model vs Global model



# Local Graph-Based Models

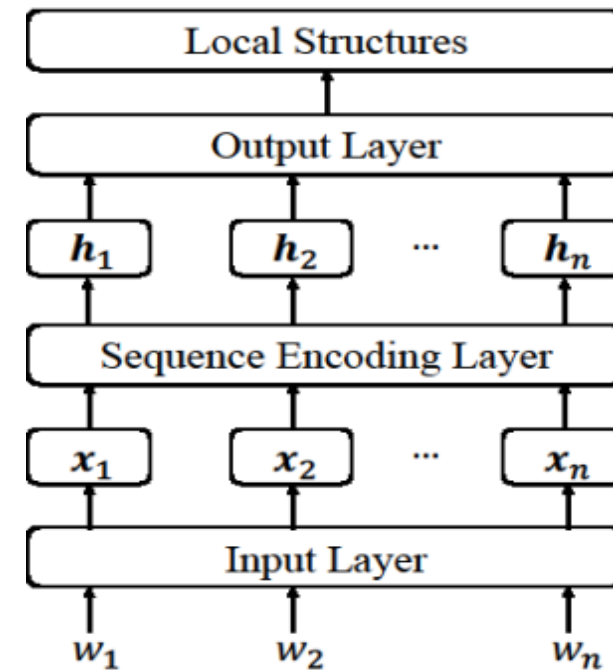
- Comparison with linear models
  - No complex discrete features
  - No dynamic program

- Overview

Input:  $W_{1:n} = w_1, w_2, \dots, w_n$

Embeddings:  $X_{1:n} = x_1, x_2, \dots, x_n$

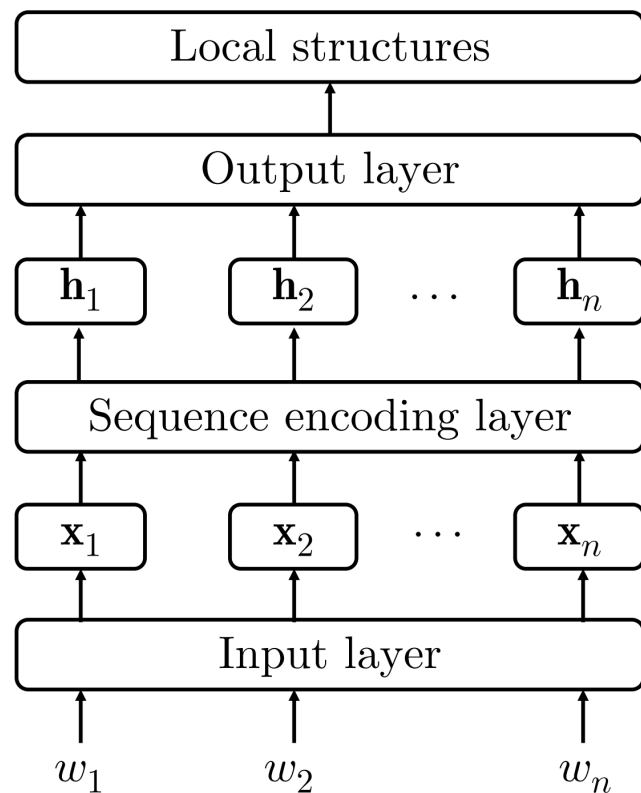
Hidden:  $H_{1:n} = h_1, h_2, \dots, h_n$



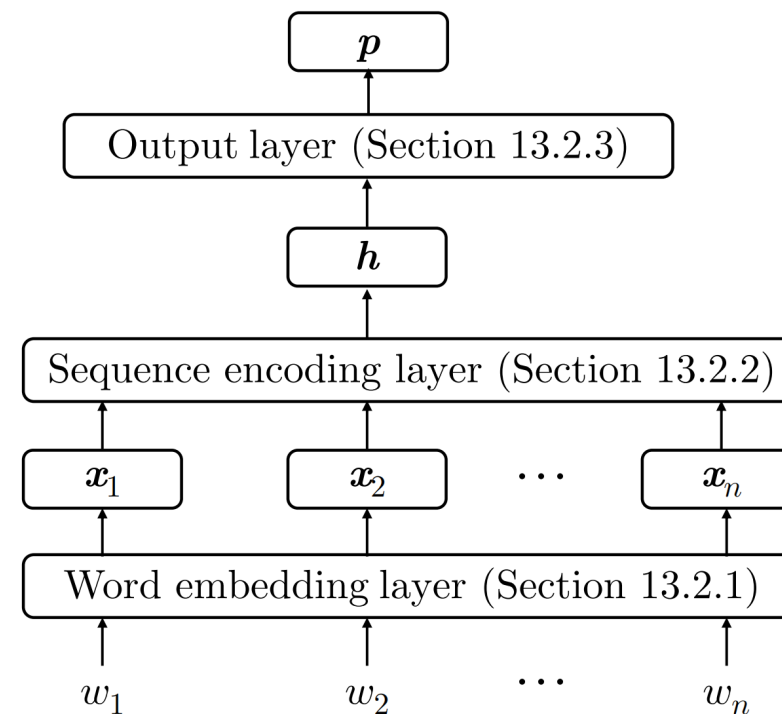


# Local Graph-Based Models

- Structure prediction *vs.* Classification



Structure prediction



Classification



- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization

# Local Graph-Based Models

- Sequence labelling

- Task:

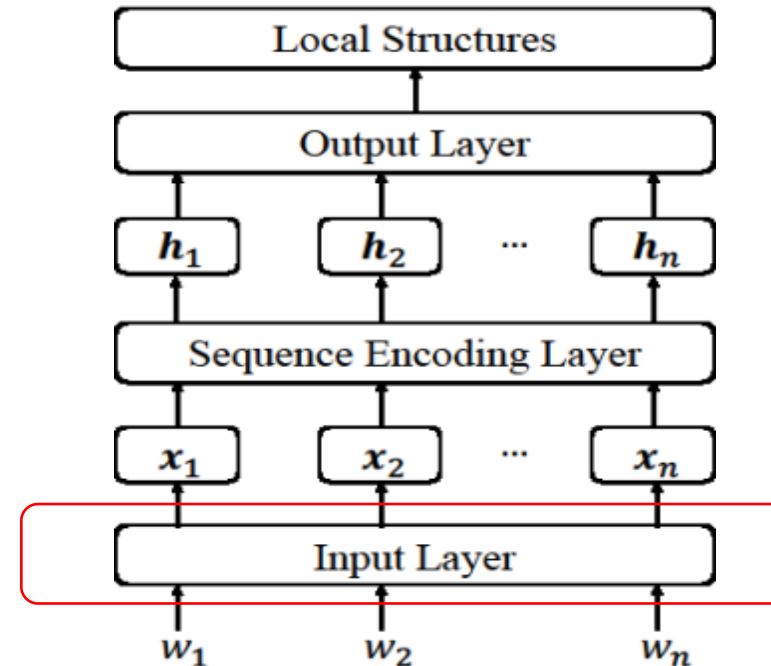
$$W_{1:n} \rightarrow T_{1:n} = t_1, t_2, \dots, t_n$$

- Input layer

$$x_i = \text{emb}(w_i)$$

- OOV words?

- add a special token <OOV> to represent OOV words.
    - during training, randomly flip infrequent words into <OOV>.



# Local Graph-Based Models

- Sequence labelling

- Task:

$$W_{1:n} \rightarrow T_{1:n} = t_1, t_2, \dots, t_n$$

- Input layer

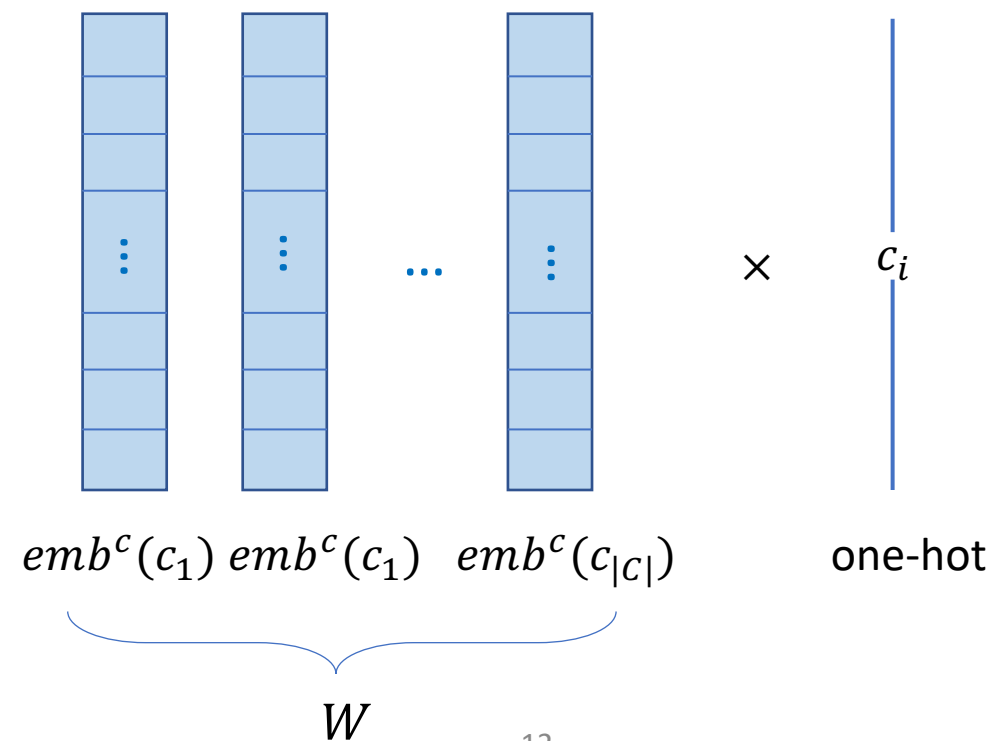
$$w_i = C_{1:|w_i|} = c_1^i, c_2^i, \dots, c_{|w_i|}^i$$

$$\mathbf{x}_{w_i}^c = [emb^c(c_1^i) \dots emb^c(c_{|w_i|}^i)]$$

$$chr(w_i) = Encoder(\mathbf{x}_{w_i}^c)$$

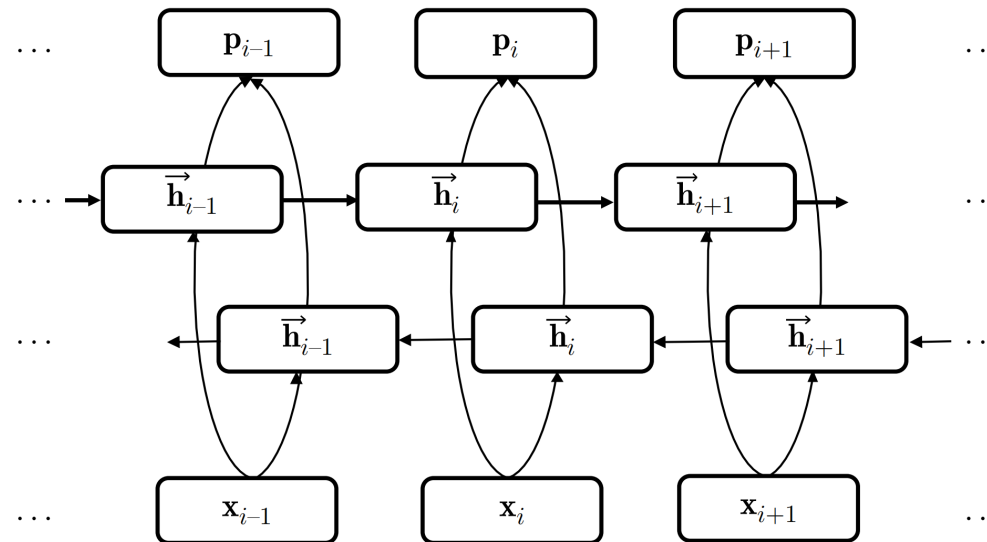
$$\mathbf{x}_i = chr(w_i) \oplus emb(w_i)$$

$$emb^c(c) = W \cdot \text{one-hot}(c)$$

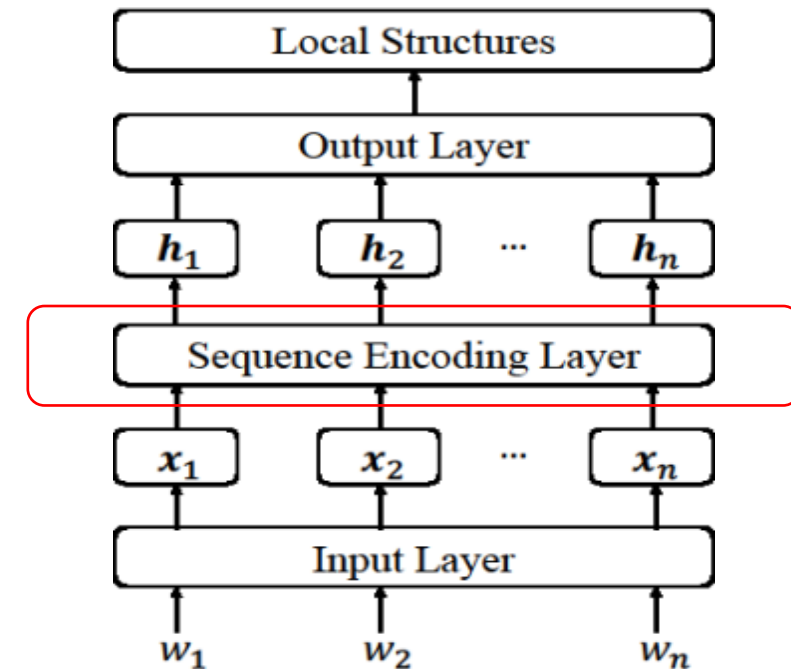


# Local Graph-Based Models

- Sequence labelling
  - Sequence representation layer



$$\mathbf{H}_{1:n} = BiLSTM(\mathbf{X}_{1:n})$$



- Can stack multi layers.

# Local Graph-Based Models

- Sequence labelling
  - Each  $\mathbf{x}_i$  is classified using  $\mathbf{h}_i$
  - There are  $|L|$  labels.
- Output layer

$$\mathbf{o}_i = \mathbf{W}\mathbf{h}_i + \mathbf{b} \quad \mathbf{p}_i = \text{softmax}(\mathbf{o}_i)$$

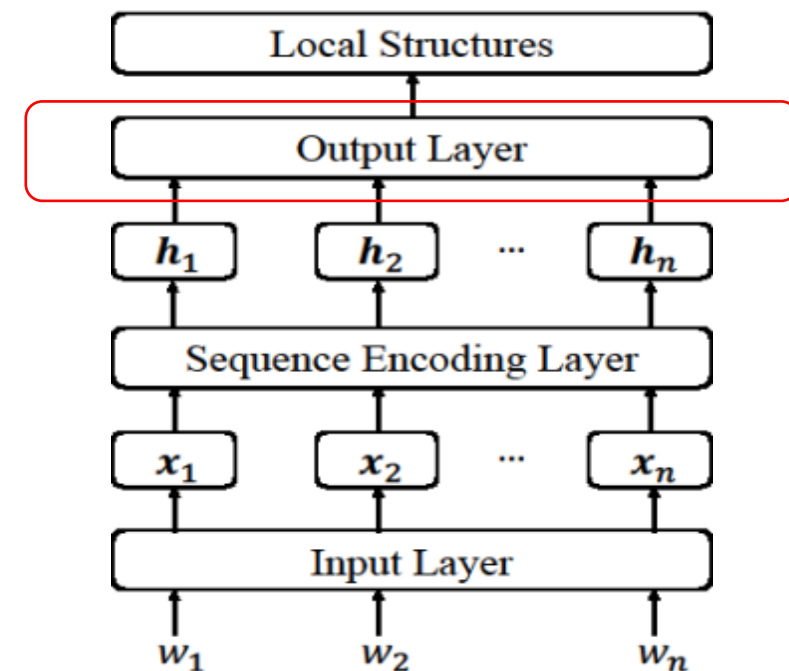
$\mathbf{o}_i$  and  $\mathbf{p}_i$  have  $|L|$  dimensions

$\mathbf{p}_i[j]$  denotes  $p(t_i = l_j | w_{1:n})$

- Training

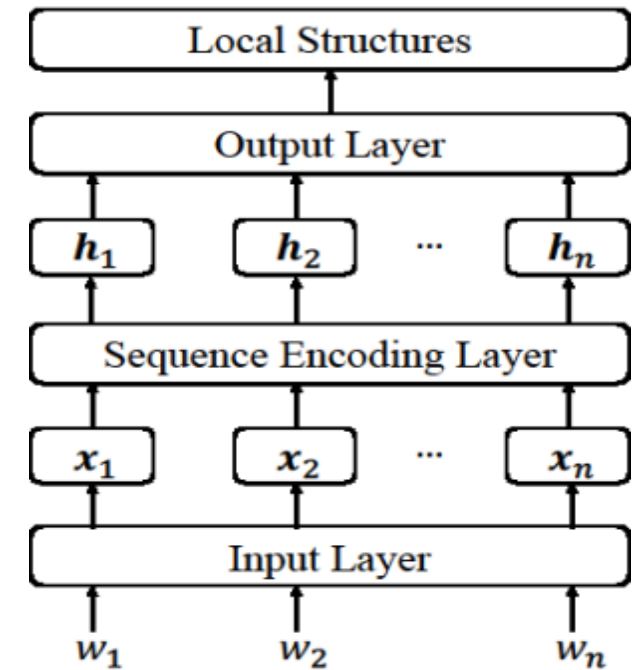
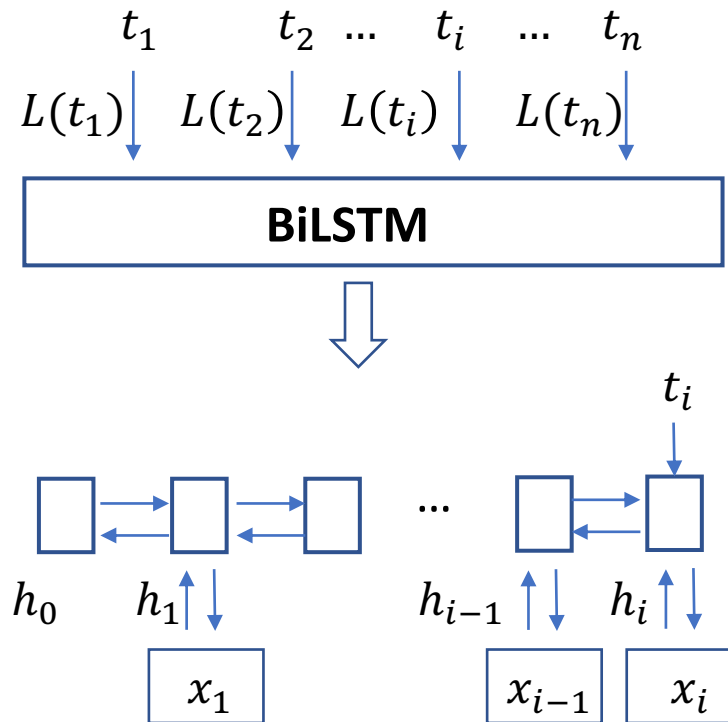
$$D = \{(W_i, T_i)\}_{i=1}^N$$

$$L = - \sum_i^N \sum_{j=1}^{|W_i|} \log(\mathbf{p}_j^i[t_j^i])$$



# Local Graph-Based Models

- There is a loss at every  $t_i$  ( $i \in [1, \dots, n]$ ).
- Gradient Propagation for RNN-based models



- Gradients from each label are accumulated.

- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - **15.1.2 Dependency Parsing**
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization



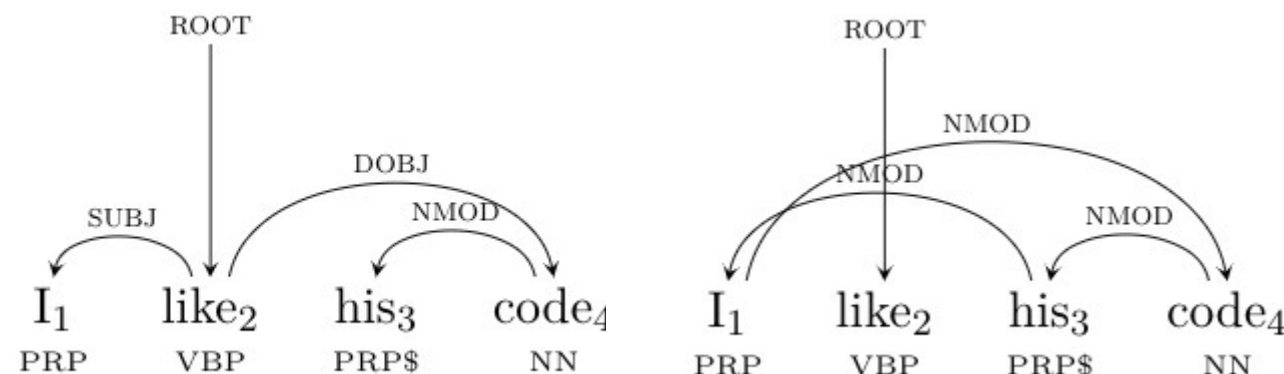
# Local Graph-Based Models

- Dependency parsing
  - Input

$$S = (w_1, t_1), (w_2, t_2), \dots, (w_n, t_n)$$

- Output

$$o = \{(i, h_i, l_i)\}_{i=1}^n$$



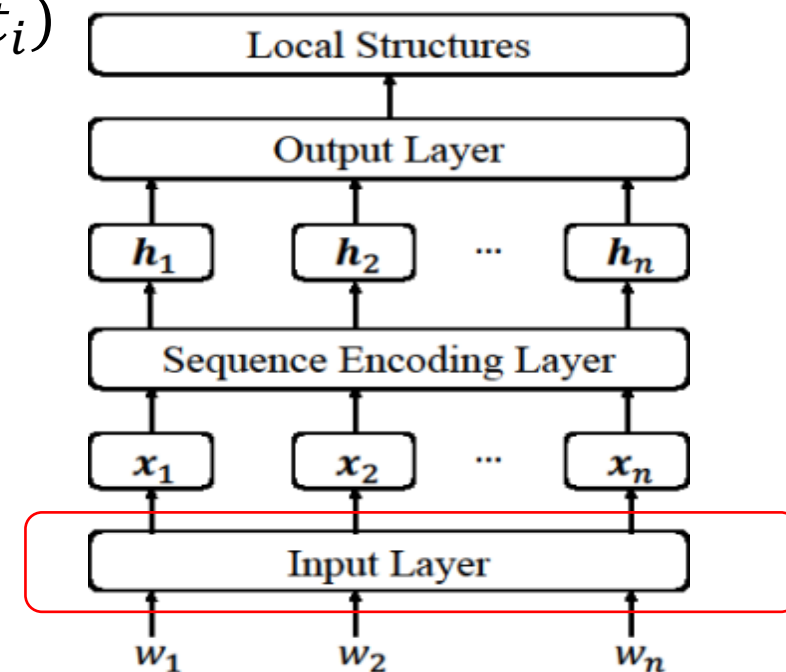
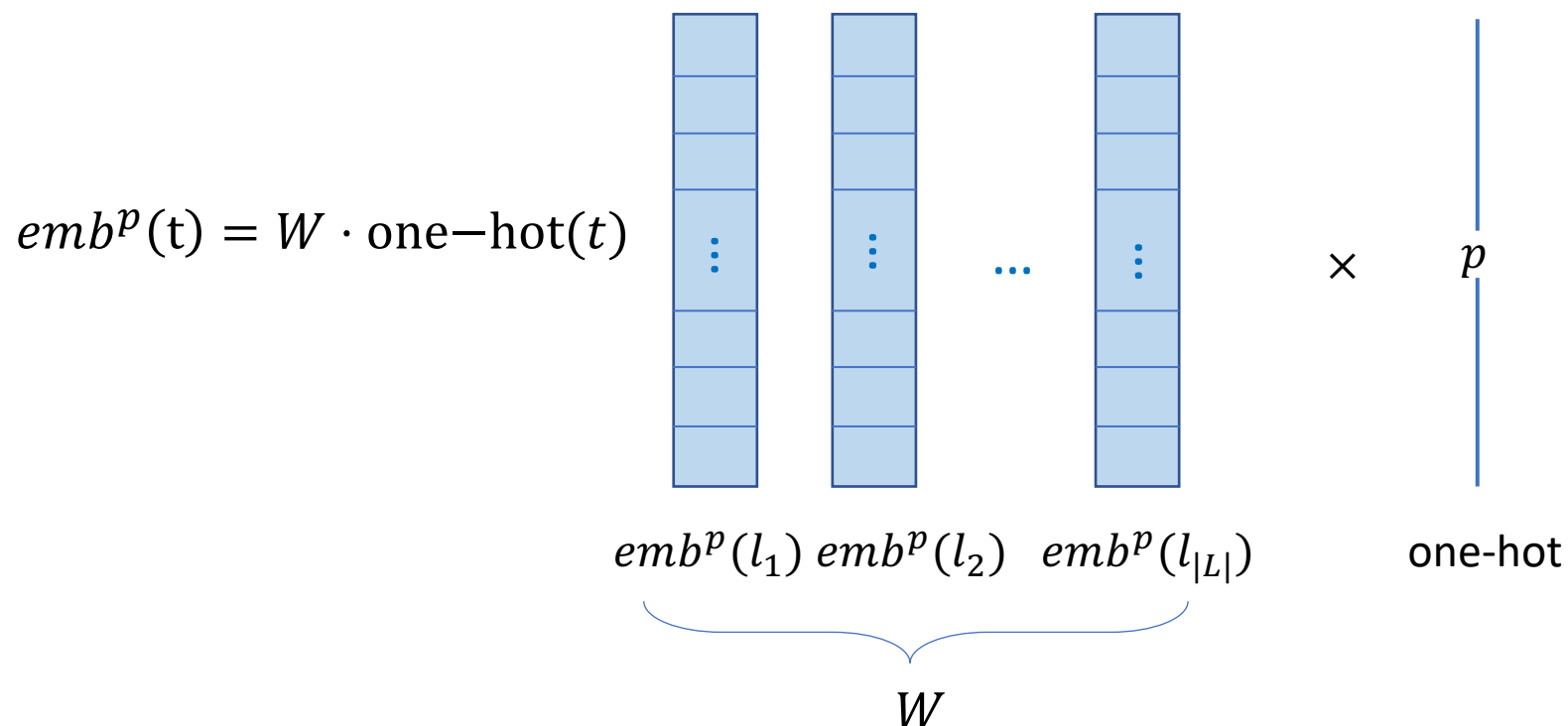
- $w_0 = ROOT$  pseudo root

# Local Graph-Based Models

- Dependency parsing
  - Input layer

$$\mathbf{x}_i = \text{emb}(w_i) \oplus \text{chr}(w_i) \oplus \text{emb}^p(t_i)$$

- POS embedding

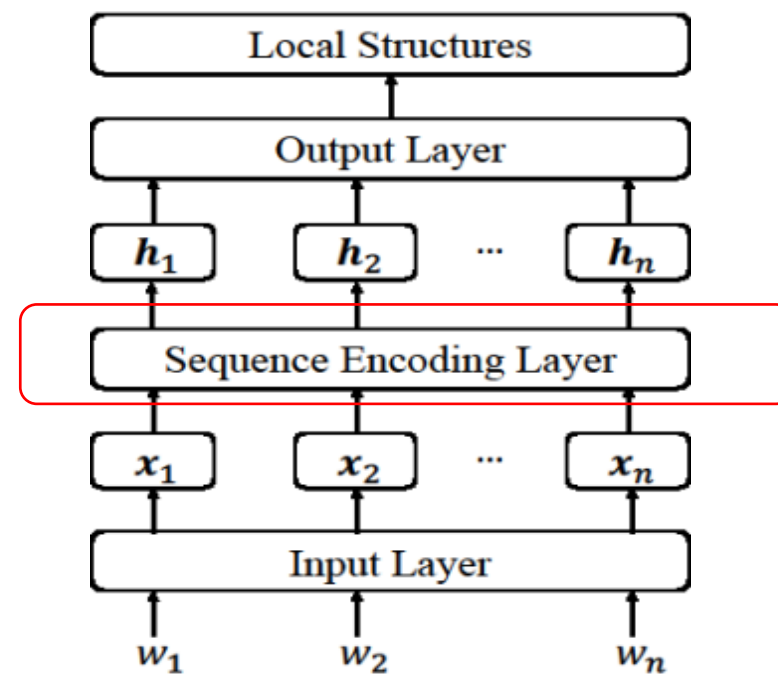




# Local Graph-Based Models

- Dependency parsing
  - Sequence encoding layer

$$\mathbf{H}_{1:n} = BiLSTM(\mathbf{X}_{1:n})$$



# Local Graph-Based Models

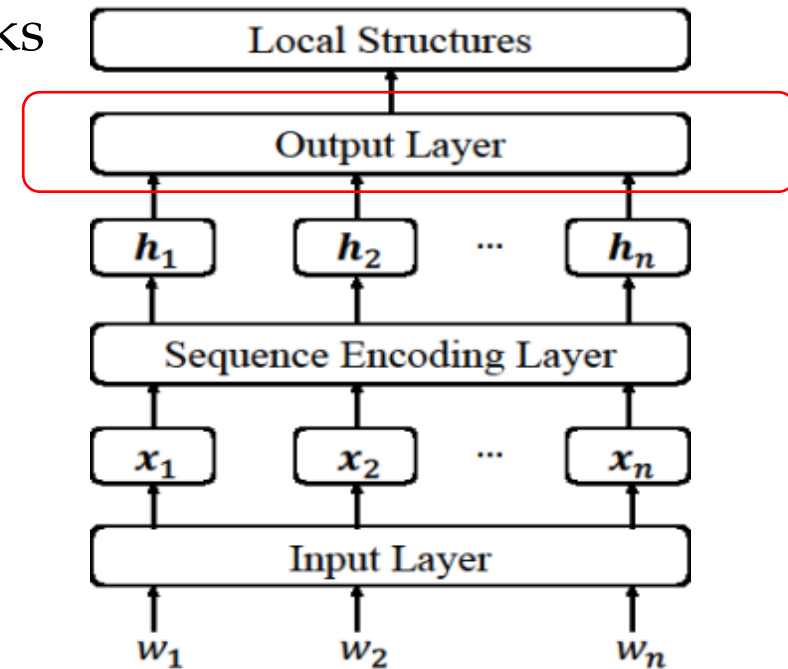
- Dependency parsing
  - Output layer
    - Structure assignment (each word  $i$  looks for a head  $j$ , including  $\langle \text{ROOT} \rangle = 0$ )

$$s_{i,j} = \mathbf{h}_i^T \mathbf{U} \mathbf{h}_j + \mathbf{v}^T([\mathbf{h}_i; \mathbf{h}_j])$$

$$\mathbf{o}_i^{\text{arc}} = \langle s_{i,1}, s_{i,2}, \dots, s_{i,n} \rangle$$

$$\mathbf{p}_i^{\text{arc}} = \text{softmax}(\mathbf{o}_i^{\text{arc}})$$

$$h_i = \underset{h}{\operatorname{argmax}} \mathbf{p}_i^{\text{arc}}[h]$$



- $\mathbf{o}, \mathbf{p}$  has  $n + 1$  dimensions, and  $\mathbf{p}_i^{\text{arc}}[j] = P(w_j \text{ is head of } w_i)$

# Local Graph-Based Models

- Dependency parsing

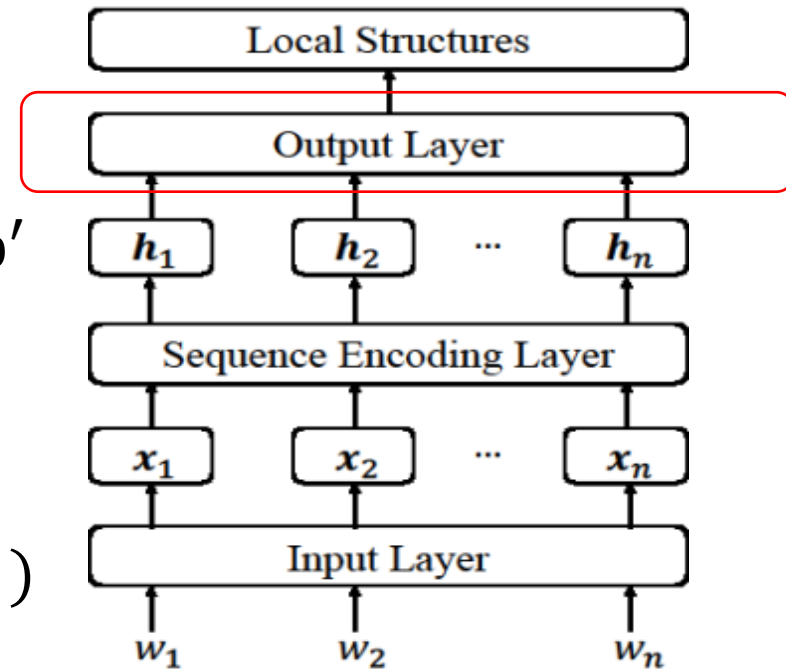
- Output layer

- Label assignment

$$\mathbf{o}_i^{label} = \mathbf{h}_i^T \mathbf{U}' \mathbf{h}_{h_i} + \mathbf{V}'([\mathbf{h}_i; \mathbf{h}_{h_i}]) + \mathbf{b}'$$

$$\mathbf{p}_i^{label} = \text{softmax}(\mathbf{o}_i^{label})$$

- $\mathbf{p}_i^{label}[j] = p(\text{arc } i \leftarrow h_i \text{ has label } l_j)$
- Results in invalid tree?
- Use a dynamic program.



# Local Graph-Based Models

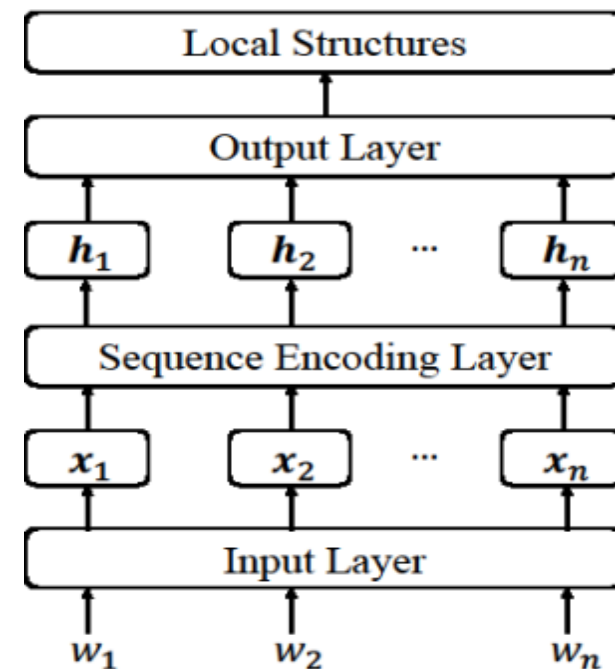
- Dependency parsing

- Training

$$D = \{(S_i, T_i)\}_{i=1}^N \quad T_i = \{(j, h_j^i, l_j^i)\}_{j=1}^{|S_i|}$$

$$L = - \sum_i^N \sum_{j=1}^{|W_i|} \left( \log \left( (\mathbf{p}_j^i)^{arc} [h_j^i] \right) + \log \left( (\mathbf{p}_j^i)^{label} [l_j^i] \right) \right)$$

- Loss from every arc accumulates via  $\mathbf{h}_i, \mathbf{h}_j$
    - Loss from every arc label accumulates via  $\mathbf{h}_i, \mathbf{h}_j$

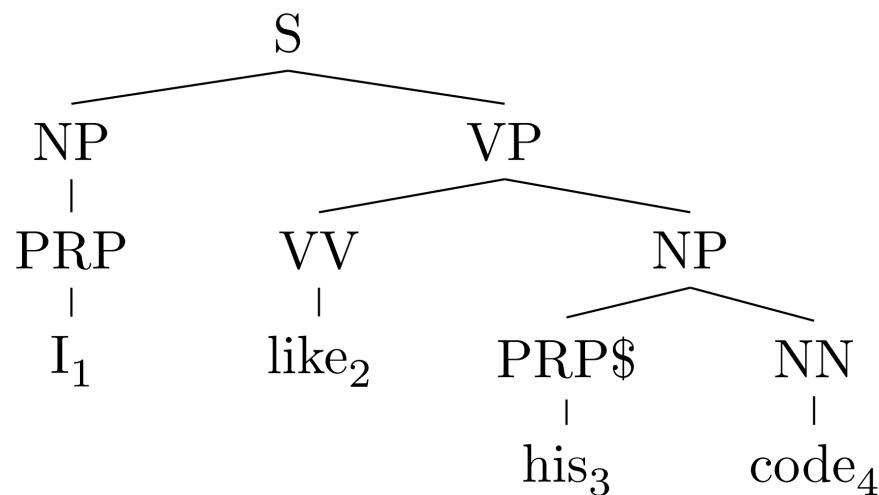


- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - **15.1.3 Constituent Parsing**
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization



# Local Graph-Based Models

- Constituent parsing
  - Task
    - Input  $S = (w_1, t_1), (w_2, t_2), \dots, (w_n, t_n)$
    - Output  $T = \{(b, e, c)\}, 1 \leq b \leq e \leq n, c = [b, e]$
  - Local model to classify each span.



# Local Graph-Based Models

- Constituent parsing

- Input layer

$$chr(w_i) = \tanh(\mathbf{W}_e^{char} \mathbf{ch}_i^l + \mathbf{W}_r^{char} \mathbf{ch}_i^r + \mathbf{b}^{char})$$

$$\mathbf{x}_i = emb(w_i) \oplus chr(w_i) \oplus emb^p(t_i)$$

- Sequence encoding layer (stacked bi-LSTM)

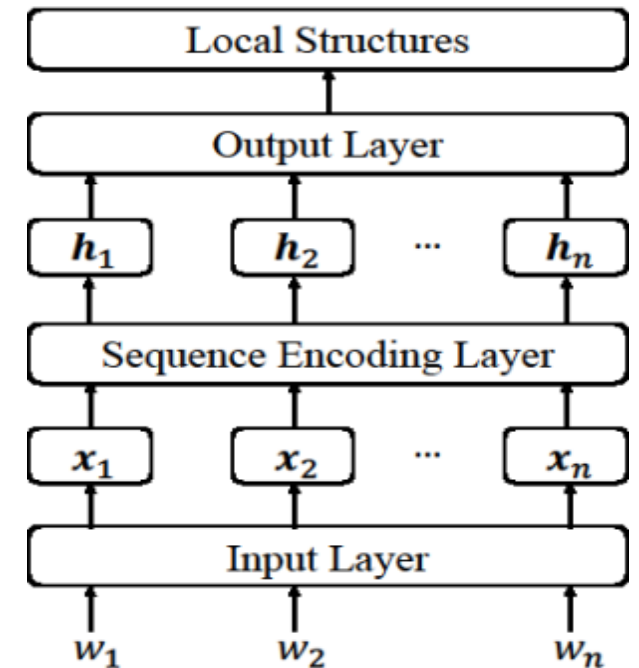
$$\mathbf{s}[b, e] = \mathbf{h}_b \oplus \mathbf{h}_e$$

- Output layer

$$\mathbf{h}[b, e] = \tanh(\mathbf{W}^h \mathbf{s}[b, e] + \mathbf{b}^h)$$

$$\mathbf{o}[b, e] = \mathbf{W}^o \mathbf{h}[b, e] + \mathbf{b}^o$$

$$\mathbf{p}[b, e] = softmax(\mathbf{o}[b, e])$$



- Constituent parsing
  - Decoding (find structure first)

$$P(y = 1 | S, b, e) = \sum_{c, c \neq \phi} P(c | S, b, e) = 1 - P(\phi | S, b, e)$$

$$P(y = 0 | S, b, e) = P(\phi | S, b, e)$$

- Rule (disregarding consistent labels)

$$W_{b:e} \rightarrow W_{b:b'-1} W_{b':e}$$

$$P(r | S, b, e) = P(y = 1 | S, b, b' - 1) P(y = 1 | S, b', e)$$

- CKY

# Local Graph-Based Models

- Constituent parsing
  - Training

$$D = \{(S_i, T_i)\}_{i=1}^N$$

$$S_i = (w_1^i, t_1^i), \dots, (w_{|S_i|}^i, t_{|S_i|}^i)$$

$$T_i = \left\{ \langle b_k^i, e_k^i, l_{b_k^i, e_k^i}^i \rangle \right\}_{k=1}^{|T_i|}$$

$$L = - \sum_{i=1}^N \sum_{1 \leq b \leq |S_i|, b \leq e \leq |S_i|} \log P(c|S, b, e)$$

$$c = l_{b,e} \text{ if } (b, e, c) \in T_i \text{ and } \phi \text{ otherwise.}$$

- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- **15.2 Local Transition-Based Models**
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization

# Local Transition-Based Models

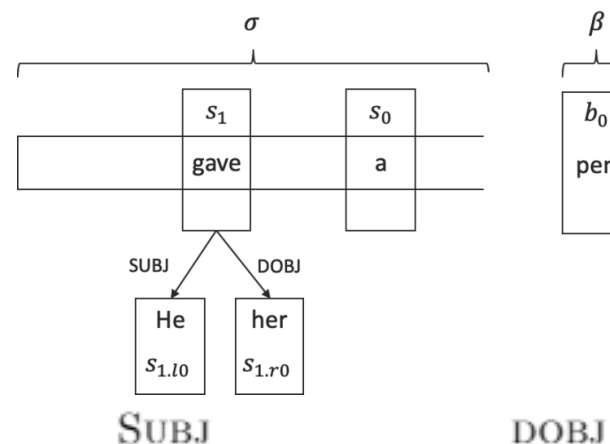
- Use state-transitions to model output (Chapter 11)
- State transitions
  - *State* represents a partially constructed output
  - *Transition actions* represent incremental steps for building structures

- Example

$$s = (\sigma, \beta, A)$$

$$\sigma = [\dots, s_1, s_0]$$

$$\beta = [b_0, b_1, \dots]$$



$$\sigma = [\text{gave}, \text{a}], \beta = [\text{pen}] \text{ and } A = \{(\text{He} \overset{\text{SUBJ}}{\curvearrowright} \text{gave}), (\text{gave} \overset{\text{DOBJ}}{\curvearrowright} \text{her})\}$$

# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: SHIFT



He gave her a pen

# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: SHIFT





# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: LEFT-ARC-SUBJ



# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: SHIFT

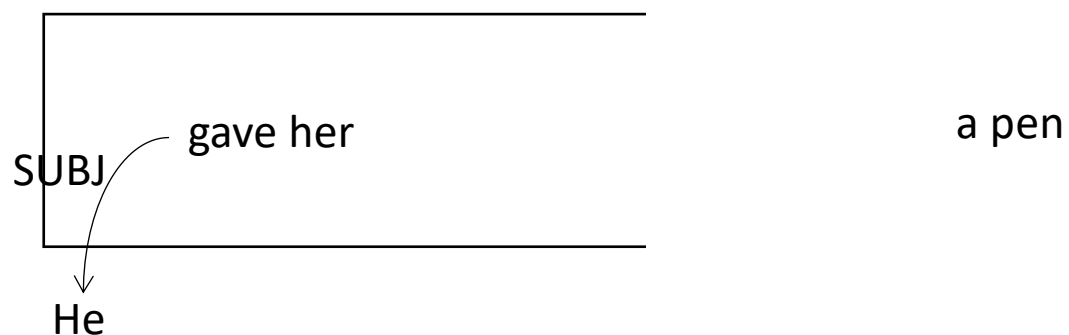


# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: RIGHT-ARC-IOBJ

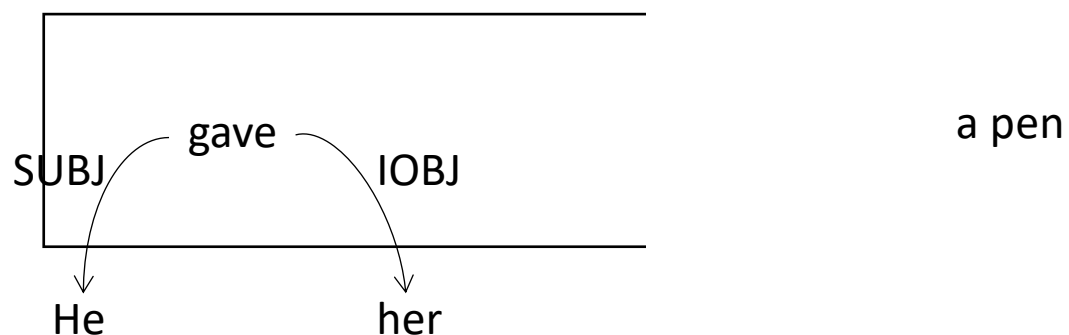


# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: SHIFT

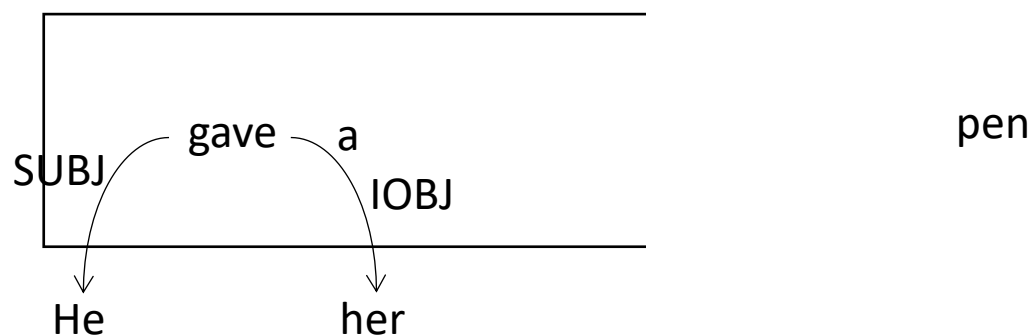


# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: SHIFT

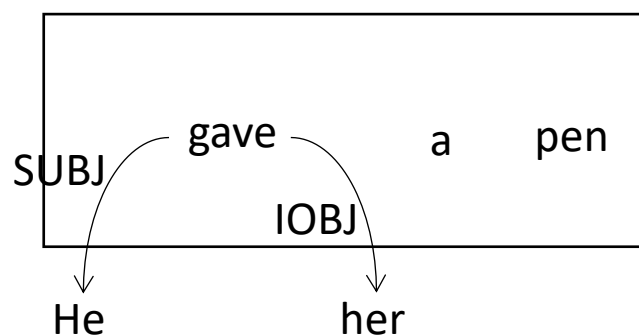


# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: LEFT-ARC-DET

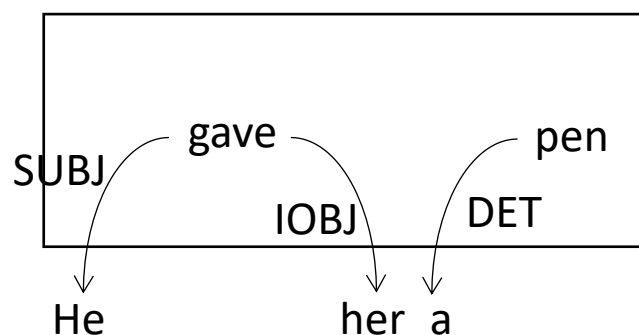


# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: RIGHT-ARC-DOBJ

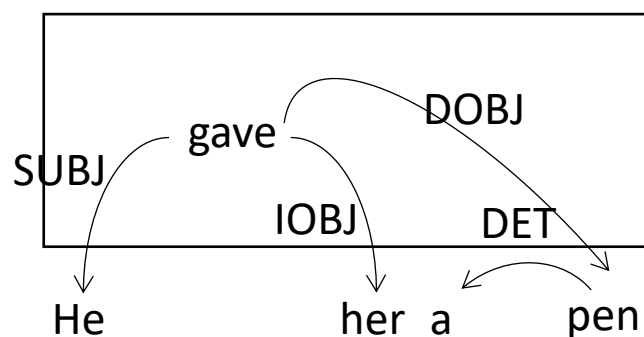


# Local Transition-Based Models

## Arc-standard Projective Dependency

- Example

Next action: END

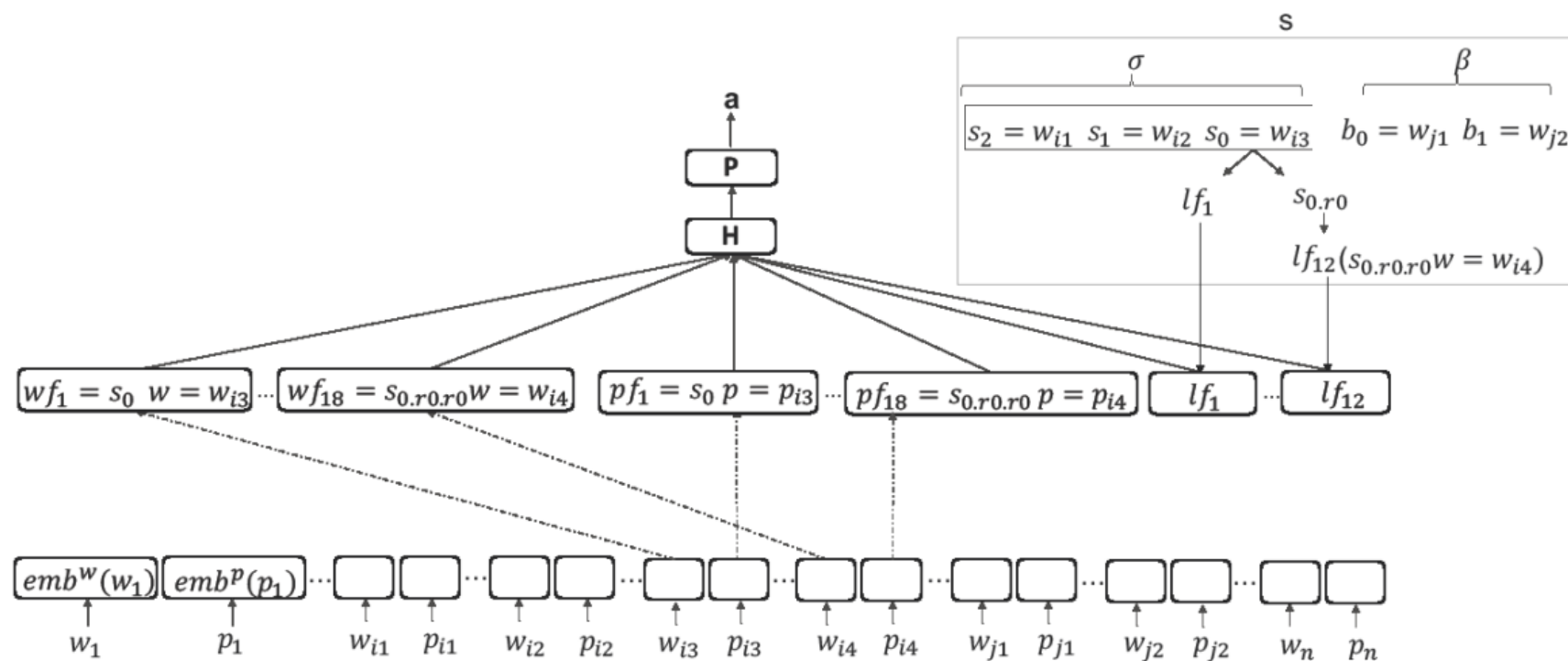




- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - **15.2.1 Model 1**
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization

# Local Transition-Based Models

- Model 1 -- directly takes embedding features.



# Local Transition-Based Models

- Model 1 -- MLP classifier

Type	Feature template
word features ( $n_w=18$ )	$s_0w, s_1w, s_2w, b_0w, b_1w, b_2w$ $s_{0.l0}w, s_{0.l1}w, s_{1.l0}w, s_{1.l1}w, s_{0.r0}w, s_{0.r1}w, s_{1.r0}w, s_{1.r1}w$ $s_{0.l0.l0}w, s_{0.r0.l0}w, s_{0.l0.r0}w, s_{0.r0.r0}w$
POS features ( $n_p=18$ )	$s_0p, s_1p, s_2p, b_0p, b_1p, b_2p$ $s_{0.l0}p, s_{0.l1}p, s_{1.l0}p, s_{1.l1}p, s_{0.r0}p, s_{0.r1}p, s_{1.r0}p, s_{1.r1}p$ $s_{0.l0.l0}p, s_{0.r0.l0}p, s_{0.l0.r0}p, s_{0.r0.r0}p$
arc features ( $n_l=12$ )	$s_{0.l0}l, s_{0.l1}l, s_{1.l0}l, s_{1.l1}l, s_{0.r0}l, s_{0.r1}l, s_{1.r0}l, s_{1.r1}l$ $s_{0.l0.l0}l, s_{0.r0.l0}l, s_{0.l0.r0}l, s_{0.r0.r0}l$

$$\mathbf{X}^w = [e; m; b(wf_1) \dots emb(wf_{n_w})]$$

$$\mathbf{X}^p = [e; m; b^p(pf_1) \dots emb^p(pf_{n_p})]$$

$$\mathbf{X}^l = [e; m; b^l(lf_1) \dots emb^l(lf_{n_l})]$$

$$\mathbf{h} = f(\mathbf{W}^w\mathbf{X}^w + \mathbf{W}^p\mathbf{X}^p + \mathbf{W}^l\mathbf{X}^l + \mathbf{b}_h)$$

$$\mathbf{o} = \mathbf{W}^o\mathbf{h} + \mathbf{b}^o$$

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

- Model 1
  - Training

$$D = \{(W_i, T_i)\}_{i=1}^N$$

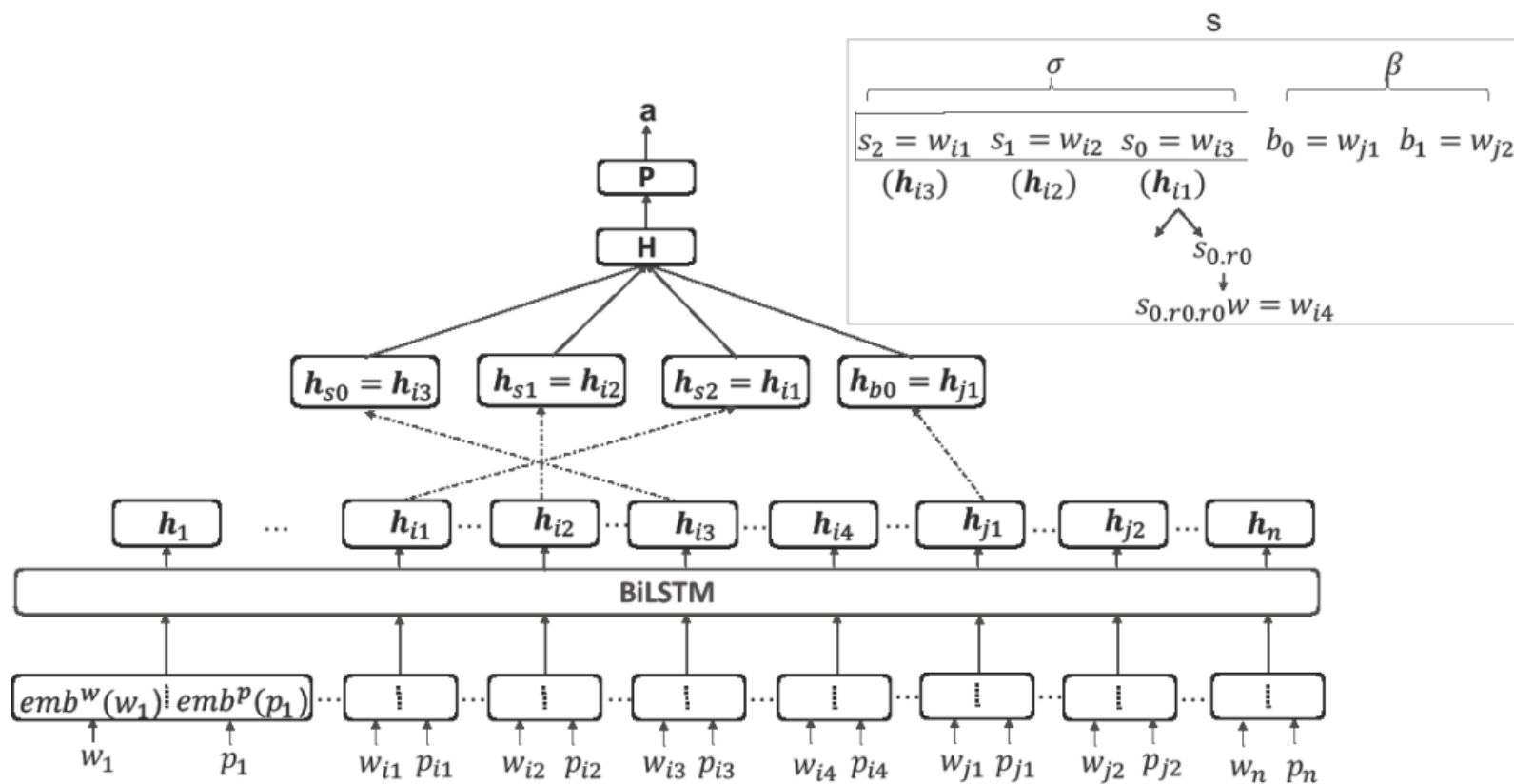
$$T_i = \langle (s_0^i, a_1^i), (s_1^i, a_2^i), \dots, (s_{|W_i|-2}^i, a_{|W_i|-1}^i) \rangle$$

$$L = - \sum_{i=1}^N \sum_{j=1}^{2|W_i|-1} \log P(a_j^i | s_{j-1}^i)$$

- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - **15.2.2 Model 2**
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization

# Local Transition-Based Models

- Model 2 – enrich input sequence features
  - Sequence encoding



# Local Transition-Based Models

- Model 2 – reduce stack feature sources.

$$\mathbf{x}_i = \text{emb}(w_i) \oplus \text{emb}^p(t_i)$$

$$\mathbf{H}_{1:n} = \text{BiLSTM}(\mathbf{X}_{1:n})$$

$$\mathbf{h} = \mathbf{h}_{s_0} \oplus \mathbf{h}_{s_1} \oplus \mathbf{h}_{s_2} \oplus \mathbf{h}_0$$

$$\mathbf{o} = \mathbf{W}^o \tanh(\mathbf{W}^h \mathbf{h} + \mathbf{b}^h) + \mathbf{b}^o$$

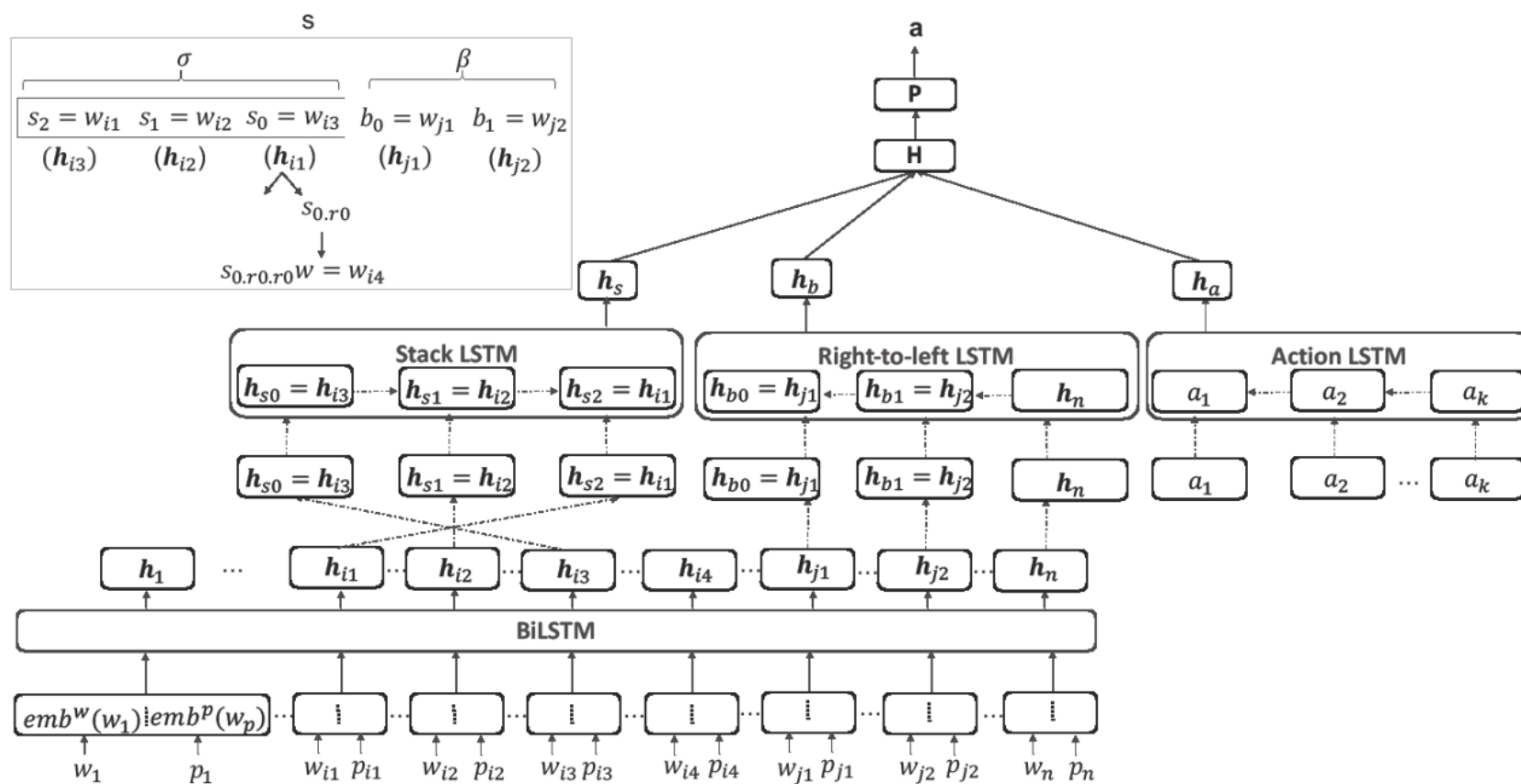
$$\mathbf{p} = \text{softmax}(\mathbf{o})$$

- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - **15.2.3 Model 3**
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization



# Local Transition-Based Models

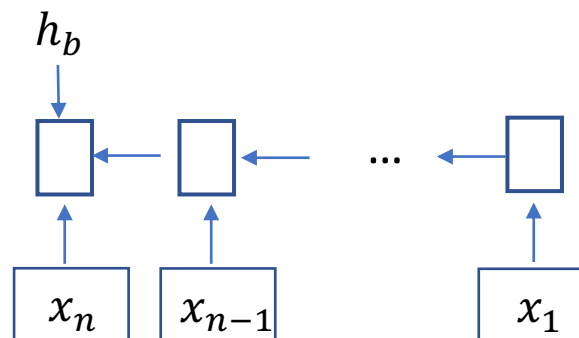
- Model 3 –better represents the state



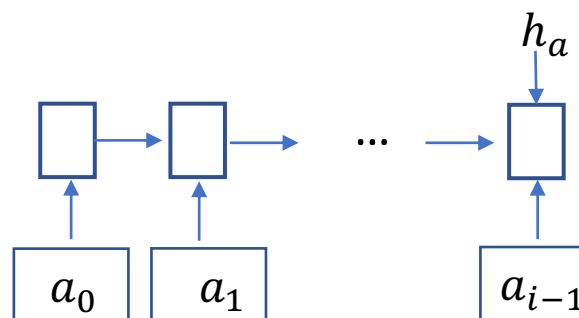
# Local Transition-Based Models

- Model 3

- Buffer LSTM

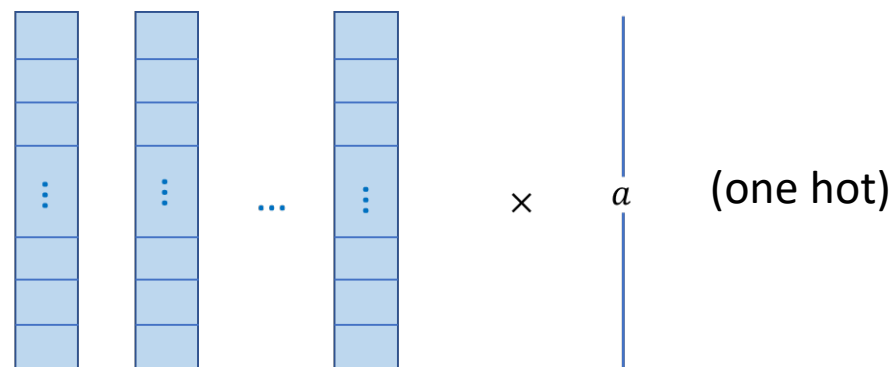


- Action LSTM



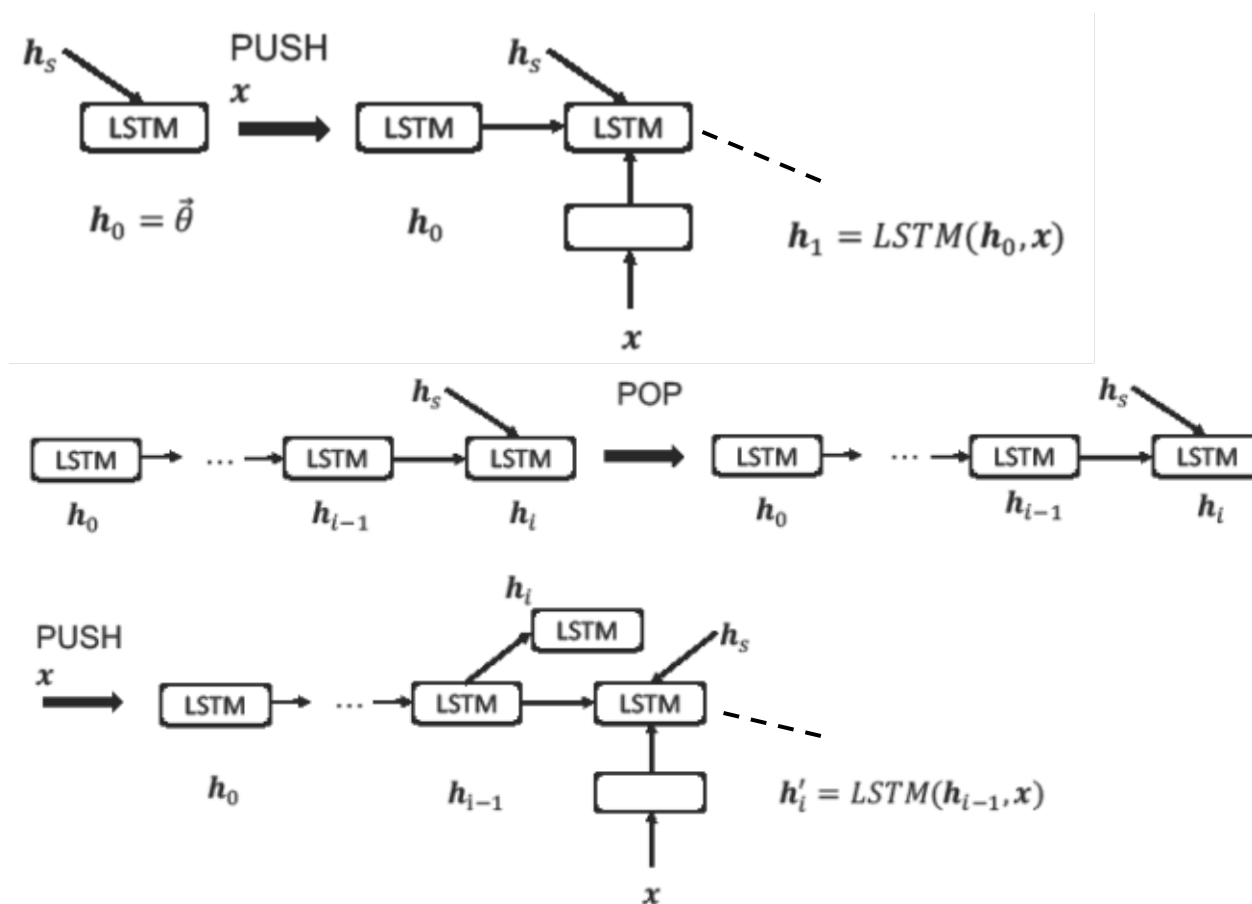
- action embedding

$$emb(a) = W$$



# Local Transition-Based Models

- Model 3
  - Stack LSTM



# Local Transition-Based Models

- Model 3
  - Further reduces stack features

$$\mathbf{h} = \mathbf{h}_s \oplus \mathbf{h}_b \oplus \mathbf{h}_a$$

$$\mathbf{o} = MLP(\mathbf{h})$$

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

- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - **15.3.1 Neural CRF**
  - 15.3.2 Neural Transition-Based Models with Global Normalization

- Neural CRF
  - CRF

$$P(T_{1:n}|W_{1:n}) = \frac{\exp(\vec{\theta} \cdot \vec{\phi}(T_{1:n}, W_{1:n}))}{\sum_{T'_{1:n}} \exp(\vec{\theta} \cdot \vec{\phi}(T'_{1:n}, W_{1:n}))}$$

$$\vec{\theta} \cdot \vec{\phi}(T_{1:n}, W_{1:n}) = \sum_{i=1}^n \vec{\theta} \cdot \vec{\phi}(t_i, t_{i-1}, W_{1:n})$$

$$P(T_{1:n}|W_{1:n}) = \frac{\exp\left(\sum_{i=1}^n \vec{\theta} \cdot \vec{\phi}(t_i, t_{i-1}, W_{1:n})\right)}{\sum_{T'_{1:n}} \exp\left(\sum_{i=1}^n \vec{\theta} \cdot \vec{\phi}(t'_i, t'_{i-1}, W_{1:n})\right)}$$

- Neural CRF – replace  $\vec{\theta} \cdot \vec{\phi}$  with a neural score.

- Neural CRF
  - A simple version

$$\mathbf{X}_{i-2:i+2} = emb(w_{i-2}) \oplus emb(w_{i-1}) \oplus emb(w_i) \oplus emb(w_{i+1}) \oplus emb(w_{i+2})$$

$$\mathbf{h}_i = f(W_{1:n}, i) = \tanh(\mathbf{W}^x \mathbf{X}_{i-2:i+2} + \mathbf{b}^x)$$

$$P(T_{1:n} | W_{1:n}) = \frac{\exp\left(\sum_{i=1}^n (\mathbf{U}(t_i) \mathbf{h}_i + b(t_i, t_{i-1}))\right)}{\sum_{T_1^n} \left( \exp\left(\sum_{i=1}^n (\mathbf{U}(t'_i) \mathbf{h}_i + b(t'_i, t'_{i-1}))\right) \right)}$$

- Neural CRF

- Training  $D = \{(W_i, T_i)\}_{i=1}^N$

$$\begin{aligned} L(W_i, T_i, \theta) &= -\frac{1}{N} \sum_{i=1}^N \log P(T_i | W_i) \\ &= -\frac{1}{N} \left( \sum_{i=1}^N \left( \sum_{j=1}^{|W_i|} (\mathbf{U}(t_j^i) \mathbf{h}_j^i + b(t_j^i, t_{j-1}^i)) \right) \right) \\ &\quad - \log \sum_{T'} \left( \exp \left( \sum_{j=1}^{|W_i|} (\mathbf{U}(t'_j) \mathbf{h}_j + b(t'_j, t'_{j-1})) \right) \right) \end{aligned}$$

- Use SGD



# Global Structured Models

- Neural CRF
  - A simple version
  - Training

$$\begin{aligned}\frac{\partial L(W_i, T_i, \Theta)}{\partial U(l_k)} &= -\left(\sum_{j=1}^{|W_i|} \mathbf{h}_j^i \delta(t_j^i = l_k)\right. \\ &\quad \left.- \sum_{T'} \frac{\exp\left(\sum_{j=1}^{|W_i|} (U(t'_j) \mathbf{h}_j^i + b(t'_j, t'_{j-1}))\right)}{\sum_{T''} \left(\exp\left(\sum_{j=1}^{|W_i|} (U(t''_j) \mathbf{h}_j^i + b(t''_j, t''_{j-1}))\right)\right)} \sum_{j=1}^{|W_i|} \mathbf{h}_j^i \delta(t'_j = l_k)\right) \\ &= -\left(\sum_{j=1}^{|W_i|} \mathbf{h}_j^i \delta(t_j^i = l_k) - \sum_{T'} P(T' | w_i) \mathbf{h}_j^i \delta(t'_j = l_k)\right) \\ &= -\sum_{j=1}^{|W_i|} \left(\mathbf{h}_j^i \delta(t_j^i = l_k) - \sum_{T'} (P(T' | w_i) \mathbf{h}_j^i \delta(t'_j = l_k))\right) \\ &= -\sum_{j=1}^{|W_i|} \left(\mathbf{h}_j^i \delta(t_j^i = l_k) - \mathbb{E}_{T' \sim P(T' | W_i)} \mathbf{h}_j^i \delta(t'_j = l_k)\right) \\ \mathbb{E}_{T' \sim P(T' | W_i)} \mathbf{h}_j^i \delta(t'_j = l_k) &= \mathbb{E}_{t'_j \sim P(t'_j | W_i)} \mathbf{h}_j^i \delta(t'_j = l_k)\end{aligned}$$

- Neural CRF
  - A simple version
    - Training

$$\begin{aligned}\frac{\partial L(W_i, T_i, \Theta)}{\partial \mathbf{h}_j^i} &= \mathbf{U}(t_j^i) - \frac{\exp\left(\sum_{j=1}^{|W_i|} (\mathbf{U}(t_j^i) \mathbf{h}_j^i + b(t_j^i, t_{j-1}^i))\right)}{\sum_{T'} \left(\exp\left(\sum_{j=1}^{|W_i|} (\mathbf{U}(t'_j) \mathbf{h}_j^i + b(t'_j, t'_{j-1}))\right)\right)} \mathbf{U}(t'_j) \\&= \mathbf{U}(t_j^i) - \sum_{T'} P(T' | W_i) \mathbf{U}(t'_j) \\&= \mathbf{U}(t_j^i) - \mathbb{E}_{T' \sim P(T' | W_i)} \mathbf{U}(t'_j) \\&= \mathbf{U}(t_j^i) - \mathbb{E}_{t'_j \sim P(t'_j | W_i)} \mathbf{U}(t'_j)\end{aligned}$$

- Neural CRF
  - A more complex version BiLSTM CRF.

$$\mathbf{H}_{1:n} = BiLSTM(\mathbf{X}_{1:n})$$

$$P(T_{1:n} | W_{1:n}) = \frac{\exp\left(\sum_{i=1}^n (\mathbf{U}(t_i)\mathbf{h}_i + b(t_i, t_{i-1}))\right)}{\sum_{T'} \left(\exp\left(\sum_{i=1}^n (\mathbf{U}(t'_i)\mathbf{h}_i + b(t'_i, t'_{i-1}))\right)\right)}$$

(add a sequence encoder)

- Gradient calculation remains similar.

# Global Structured Models

- Neural CRF
  - More variants?
  - Consider label embedding

- Neural CRF
  - Tree CRF

$$P(T|W) = \frac{\exp\left(\vec{\theta} \cdot \vec{\phi}(W, T)\right)}{\sum_{T' \in \text{Gen}(W)} \exp\left(\vec{\theta} \cdot \vec{\phi}(W, T')\right)}$$

$$\vec{\phi}(W, T) = \sum_{r \in T} \vec{\phi}(W, r)$$

$$P(T|W) = \frac{\exp\left(\sum_{r \in T} \vec{\theta} \cdot \vec{\phi}(W, r)\right)}{\sum_{T' \in \text{Gen}(W)} \exp\left(\sum_{r' \in T'} \vec{\theta} \cdot \vec{\phi}(W, r')\right)}$$

- Neural version – replace  $\vec{\theta} \cdot \vec{\phi}$  with a neural score.

- Neural CRF
  - Neural version using label embedding parameterization

$$P(T|W) = \frac{\exp\left(\sum_{r \in T} f(W, r)\right)}{\sum_{T' \in \text{Gen}(W)} \exp\left(\sum_{r' \in T'} f(W, r')\right)}$$

$$f(W, r) = f(W, c, \rightarrow, c_1, c_2, bb'e\vec{\theta}) = \vec{\tau}(W, b, b', e)^T \mathbf{W}^f \vec{\gamma}(c \rightarrow c_1 c_2)$$

$$\vec{\tau}(W, b, b', e) = \text{ReLU}\left(\mathbf{W}^w(\mathbf{h}_b \oplus \mathbf{h}_{b'-1} \oplus \mathbf{h}_e)\right)$$

$$\begin{aligned} \vec{\gamma}(c \rightarrow c_1 c_2) \\ = \text{ReLU}\left(\mathbf{W}^a \text{emb}^s(c) + \mathbf{W}^l \text{emb}^s(c_1) + \mathbf{W}^r \text{emb}^s(c_2) + \mathbf{b}\right) \end{aligned}$$

# Global Structured Models

- Neural CRF

- Training  $D = \{(W_i, T_i)\}_{i=1}^N$

$$L = -\frac{1}{N} \sum_{i=1}^N \log P(T_i | W_i) = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp\left(\sum_{r \in T_i} \vec{\tau}(r)^T \mathbf{W}^f \vec{\gamma}(r)\right)}{\sum_{T'_i \in \text{Gen}(W_i)} \exp\left(\sum_{r' \in T'_i} \vec{\tau}(r')^T \mathbf{W}^f \vec{\gamma}(r')\right)}$$

- Use SGD

- Neural CRF
  - Tree CRF
    - Training

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{W}^f} &= - \left( \sum_{r \in T_i} \vec{\tau}(r) \vec{\gamma}(r)^T - \sum_{T'_i \in \text{Gen}(W_i)} P(T'_i | W_i) \sum_{r' \in T'_i} \vec{\tau}(r') \vec{\gamma}(r')^T \right) \\ &= - \left( \sum_{r \in T_i} \vec{\tau}(r) \vec{\gamma}(r)^T - \mathbb{E}_{T'_i \sim P(T'_i | W_i)} \sum_{r' \in T'_i} \vec{\tau}(r') \vec{\gamma}(r')^T \right)\end{aligned}$$

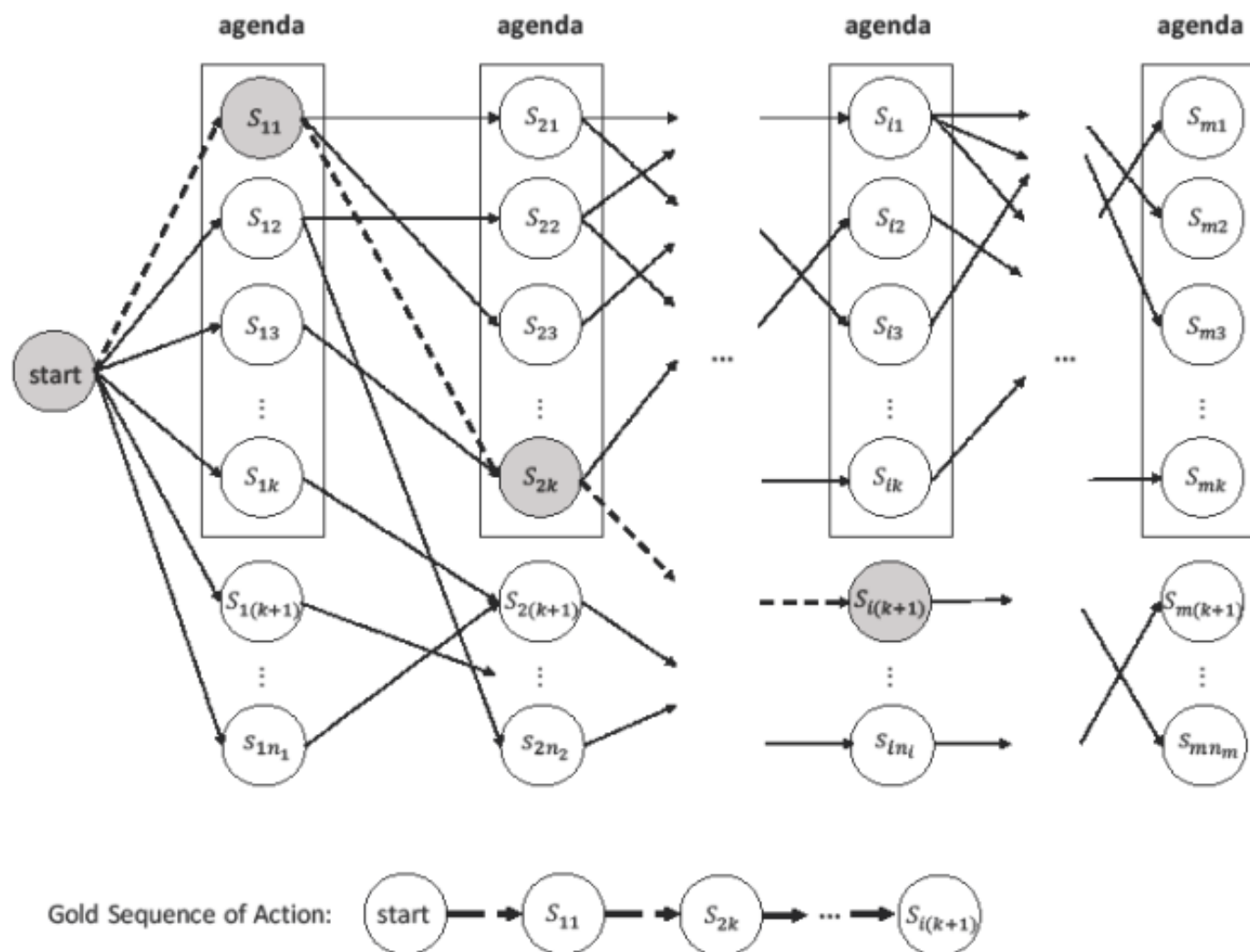
$$\begin{aligned}\frac{\partial L}{\partial \vec{\tau}(r)} &= - \left( \sum_{r \in T_i} \mathbf{W}^f \vec{\gamma}(r) - \sum_{T'_i \in \text{Gen}(W_i)} P(T'_i | W_i) \sum_{r' \in T'_i} \mathbf{W}^f \vec{\gamma}(r') \right) \\ &= - \left( \sum_{r \in T_i} \mathbf{W}^f \vec{\gamma}(r) - \mathbb{E}_{T'_i \sim P(T'_i | W_i)} \sum_{r' \in T'_i} \mathbf{W}^f \vec{\gamma}(r') \right)\end{aligned}$$



- 15.1 Local Graph-Based Models
  - 15.1.1 Sequence labelling
  - 15.1.2 Dependency Parsing
  - 15.1.3 Constituent Parsing
  - 15.1.4 Comparison with Linear Models
- 15.2 Local Transition-Based Models
  - 15.2.1 Model 1
  - 15.2.2 Model 2
  - 15.2.3 Model 3
- 15.3 Global Structured Models
  - 15.3.1 Neural CRF
  - 15.3.2 Neural Transition-Based Models with Global Normalization

# Global Structured Models

- Neural Transition-Based Models with Global Normalization



- Neural Transition-Based Models with Global Normalization
- Linear version

$$score(A_{1:i}) = \sum_{j=1}^i score(a_j) = \sum_{j=1}^i \vec{\theta} \cdot \vec{\phi}(s_{j-1} \cdot a_j)$$

- Neural version

$$score(A_{1:i}) = \sum_{j=1}^i score(a_j) = \sum_{j=1}^i \mathbf{o}_j[a_j]$$

$$A_{1:i} = a_1, a_2, \dots, a_i$$

- Neural Transition-Based Models with Global Normalization
  - Training

$$P(A_{1:i}|W_{1:n}) = \text{softmax}(A_{1:i}, A'_{1:i} \in \text{gen}(W_{1:n}, i)) = \frac{\exp\left(\sum_{j=1}^i \mathbf{o}_j[a_j]\right)}{\sum_{A'_{1:i}} \exp\left(\sum_{j=1}^i \mathbf{o}'_j[a_j]\right)}$$

$$L = -\log P(A_{1:i}|W_{1:n})$$
$$= -\log \frac{\exp\left(\sum_{j=1}^i \mathbf{o}_j[a_j]\right)}{\sum_{A'_{1:i}} \exp\left(\sum_{j=1}^i \mathbf{o}'_j[a_j]\right)}$$

$$= -\log \frac{\exp\left(\sum_{j=1}^i \mathbf{o}_j[a_j]\right)}{Z}$$

$$= \log Z - \sum_{j=1}^i \mathbf{o}_j[a_j]$$

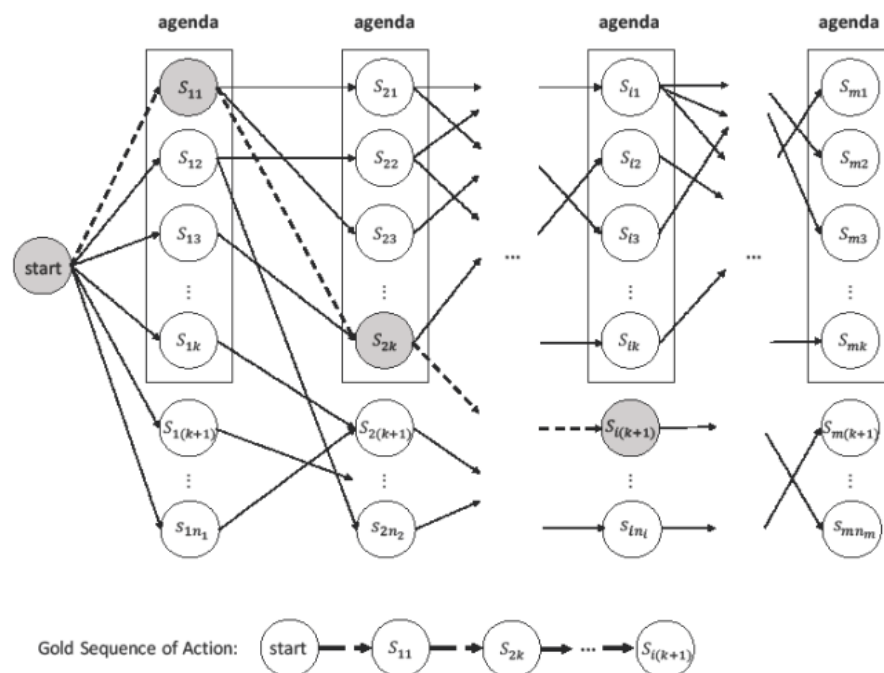
$$\text{where } Z = \sum_{A'_{1:i}} \exp\left(\sum_{j=1}^i \mathbf{o}_j[a'_j]\right)$$

- Neural Transition-Based Models with Global Normalization

- $A'_{1:i}$  is exponential to  $i$

- Contrastive estimation

$$Z'(x_i, \theta) = \sum_{A'_{1:i} \in [S_{i,1}, \dots, S_{i,n_i}]} \exp \left( \sum_{j=1}^i \mathbf{o}_j [a'_j] \right)$$



# Summary

- Local graph-based neural models
- Local transition-based neural models
- Global graph-based neural models
- Global transition-based neural models