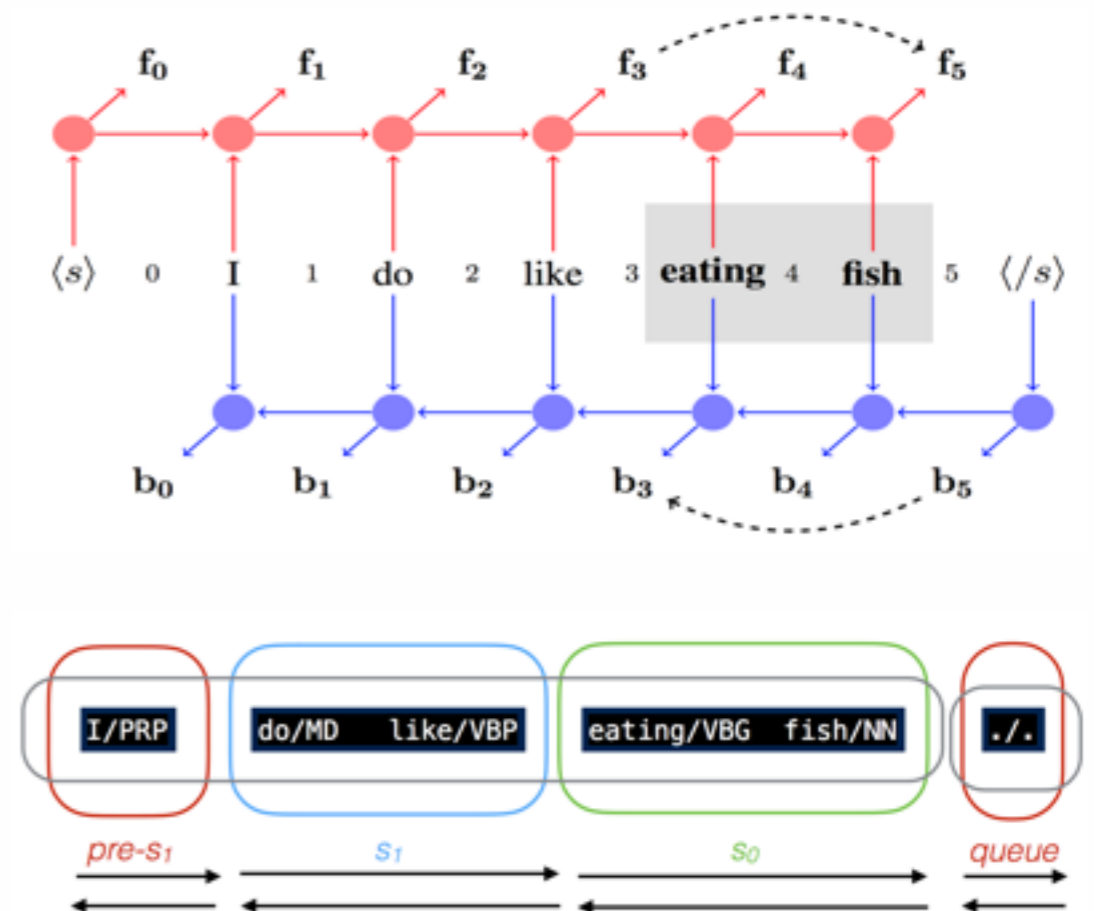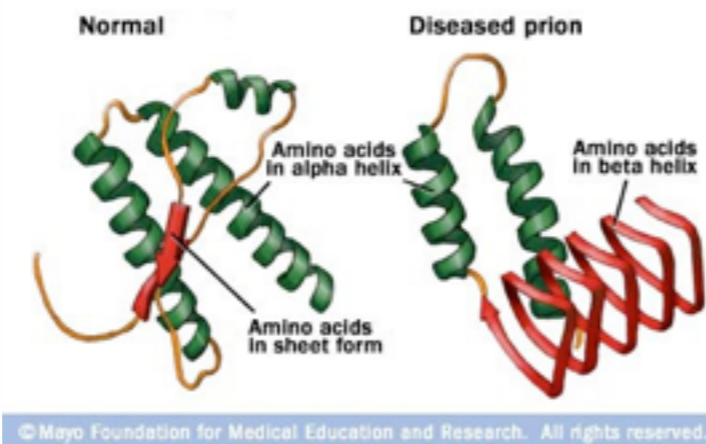# Marrying Dynamic Programming with Recurrent Neural Networks



I eat sushi with tuna from Japan

**Liang Huang**

Oregon State University

Structured Prediction Workshop, EMNLP 2017, Copenhagen, Denmark

# Marrying Dynamic Programming with Recurrent Neural Networks



**Liang Huang**

Oregon State University

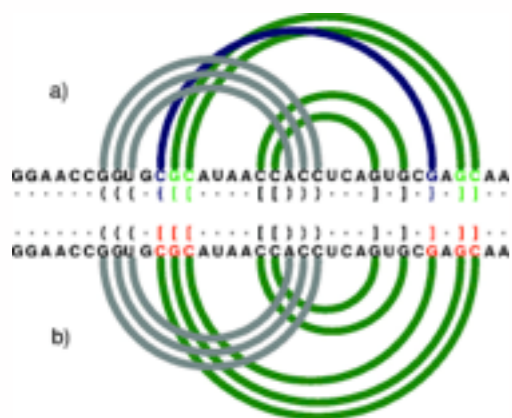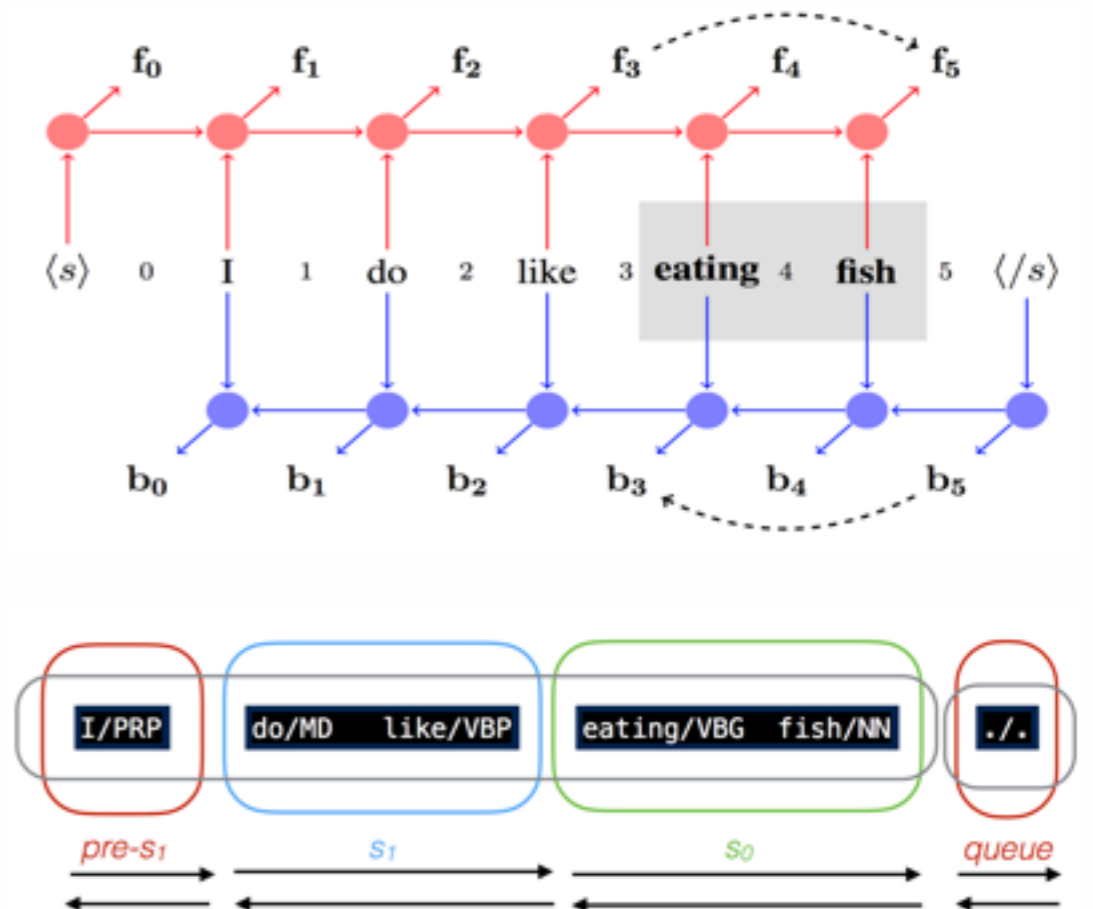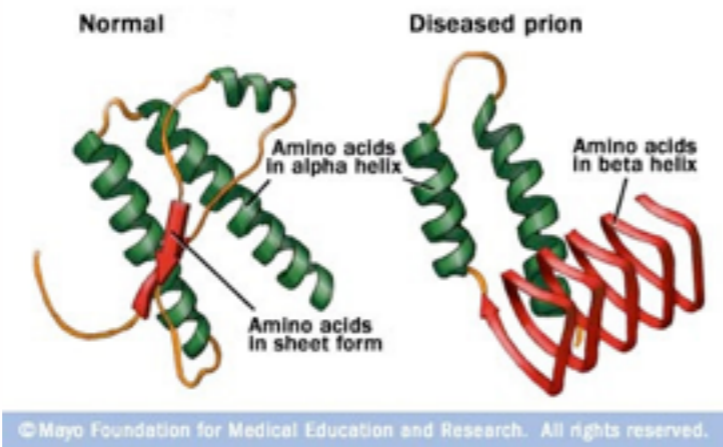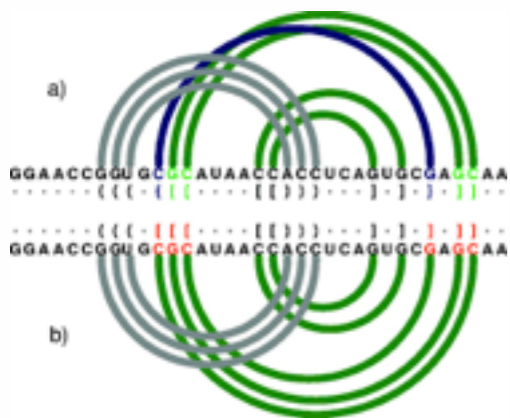Structured Prediction Workshop, EMNLP 2017, Copenhagen, Denmark

# Marrying Dynamic Programming with Recurrent Neural Networks



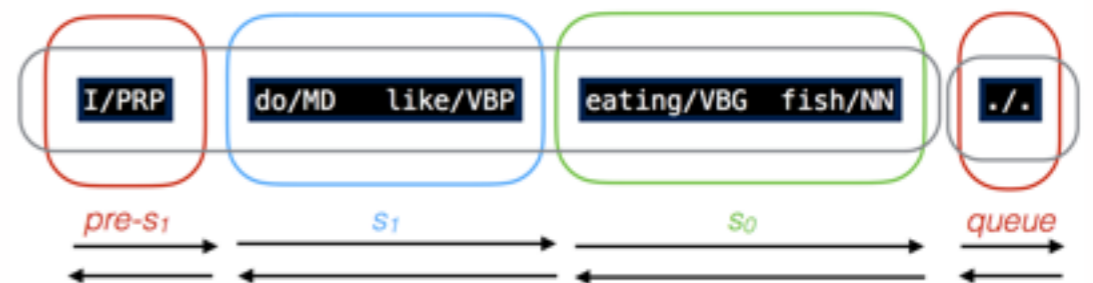I eat sushi with tuna from Japan

**Liang Huang**

Oregon State University

James Cross

Structured Prediction Workshop, EMNLP 2017, Copenhagen, Denmark

# Structured Prediction is Hard!

# Not Easy for Humans Either...



(structural ambiguity :-P)

# Not Even Easy for Nature!

- prion: "misfolded protein"
  - structural ambiguity for the same amino-acid sequence
  - similar to different interpretations under different contexts
  - causes mad-cow diseases etc.



Normal | Diseased prion
Amino acids in alpha helix
Amino acids in beta helix
Amino acids in sheet form

Scrapie | Mad cow disease | Various

Spontaneous misfolding
Spontaneous mutation in PrP gene

?? Another species
?? Spontaneous

Feeding cows sheep meat/bones

Tasty steak

Vaccines Cannibalism

Cannibalism

Inheritance of mutation
Iatrogenic transmission
Cannibalism

# Case Study: Parsing and Folding

- both problems have exponentially large search space

  - both can be modeled by grammars (context-free & above)

- question 1: how to search for the highest-scoring structure?

- question 2: how to make gold structure score the highest?

I eat sushi with tuna from Japan

# Solutions to Search and Learning

- question 1: how to search for the highest-scoring structure?
  - answer: dynamic programming to factor search space

- question 2: how to make gold structure score the highest?
  - answer: neural nets to automate feature engineering

- But do DP and neural nets like each other??

I eat sushi with tuna from Japan

6

# Solutions to Search and Learning

- question 1: how to search for the highest-scoring structure?

  - answer: dynamic programming to factor search space

- question 2: how to make gold structure score the highest?

  - answer: neural nets to automate feature engineering

- But do DP and neural nets like each other??

# In this talk...

- Background

- Dynamic Programming for Incremental Parsing

- Features: from sparse to neural to recurrent neural nets

- Bidirectional RNNs: minimal features; no tree structures!

  - dependency parsing (Kiperwaser+Goldberg, 2016, Cross+Huang, 2016a)

  - span-based constituency parsing (Cross+Huang, 2016b)

- Marrying DP & RNNs *(mostly not my work!)*

  - transition-based dependency parsing (Shi et al, EMNLP 2017)

  - minimal span-based constituency parsing (Stern et al, ACL 2017)

# Spectrum: Neural Incremental Parsing

constituency
dependency

*bottom-up*

*edge-factored*
*(McDonald+ 05a)*

DP incremental parsing
(Huang+Sagae 10, Kuhlmann+ 11)

Feedforward NNs
*(Chen + Manning 14)*

Stack LSTM

*(Dyer+ 15)*

RNNG
*(Dyer+ 16)*

biRNN dependency
*(Kiperwaser+Goldberg 16;
Cross+Huang 16a)*

*biRNN graph-based
dependency
(Kiperwaser+Goldberg 16;
Wang+Chang 16)*

biRNN span-based
constituency
*(Cross+Huang 16b)*

*minimal span-based
constituency
(Stern+ ACL 17)*

minimal dependency
*(Shi+ EMNLP 17)*

*all tree info
(summarize output **y**)*

*minimal or no tree info
(summarize input **x**)*

DP impossible     enables slow DP          enables fast DP      fastest DP: $O(n^3)_8$

# Spectrum: Neural Incremental Parsing

constituency
dependency

*bottom-up*

*edge-factored*
*(McDonald+ 05a)*

DP incremental parsing
(Huang+Sagae 10, Kuhlmann+ 11)

Feedforward NNs
*(Chen + Manning 14)*

Stack LSTM

*(Dyer+ 15)*

RNNG
*(Dyer+ 16)*

biRNN dependency
*(Kiperwaser+Goldberg 16;
Cross+Huang 16a)*

*biRNN graph-based
dependency
(Kiperwaser+Goldberg 16;
Wang+Chang 16)*

biRNN span-based
constituency
*(Cross+Huang 16b)*

*minimal span-based
constituency
(Stern+ ACL 17)*

minimal dependency
*(Shi+ EMNLP 17)*

*all tree info
(summarize output* **y**)

*minimal or no tree info
(summarize input* **x**)

DP impossible          enables slow DP                    enables fast DP          fastest DP: $O(n^3)$ 8

# Incremental Parsing with Dynamic Programming

(Huang & Sagae, ACL 2010*; Kuhlmann et al., ACL 2011; Mi & Huang, ACL 2015)

*best paper nominee

# Incremental Parsing with Dynamic Programming

(Huang & Sagae, ACL 2010[*]; Kuhlmann et al., ACL 2011; Mi & Huang, ACL 2015)

[*] *best paper nominee*

I eat sushi with tuna from Japan in a restaurant

action             stack                queue

I eat sushi with tuna from Japan in a restaurant

|   | action | stack | queue |
|---|--------|-------|-------|
| 0 | - | | I eat sushi ... |

# Incremental Parsing (Shift-Reduce)

I eat sushi with tuna from Japan in a restaurant

➡

| | action | stack | queue |
|---|---|---|---|
| | | | I eat sushi ... |
| | | | eat sushi with ... |
| 0 | - | | |
| 1 | shift | I | |

# Incremental Parsing (Shift-Reduce)

I eat sushi with tuna from Japan in a restaurant

→

| | action | stack | queue |
|---|---|---|---|
| | | | I eat sushi ... |
| 0 | - | | eat sushi with ... |
| 1 | shift | I | sushi with tuna ... |
| 2 | shift | I eat | |

# Incremental Parsing (Shift-Reduce)

I eat sushi with tuna from Japan in a restaurant

→

| action | stack | queue |
|--------|-------|-------|

| | action | stack |
|---|---|---|
| 0 | - | [ ] |
| 1 | shift | I |
| 2 | shift | I eat |
| 3 | l-reduce | eat / I |

| queue |
|-------|
| I eat sushi ... |
| eat sushi with ... |
| sushi with tuna ... |
| sushi with tuna ... |

# Incremental Parsing (Shift-Reduce)

I eat sushi with tuna from Japan in a restaurant

| | action | stack | | queue |
|---|---|---|---|---|
| 0 | - | | | I eat sushi ... |
| 1 | shift | I | | eat sushi with ... |
| 2 | shift | I eat | | sushi with tuna ... |
| 3 | l-reduce | eat | | sushi with tuna ... |
| 4 | shift | eat sushi | | with tuna from ... |

# Incremental Parsing (Shift-Reduce)

I eat sushi with tuna from Japan in a restaurant

| action | | stack | | queue |
|--------|--|-------|--|-------|

0   -   (empty)

1   shift   I

2   shift   I eat

3   l-reduce   eat → I

4   shift   eat sushi → I

5a   r-reduce   eat → I   sushi

**queue:**

| |
|---|
| I eat sushi ... |
| eat sushi with ... |
| sushi with tuna ... |
| sushi with tuna ... |
| with tuna from ... |
| with tuna from ... |

I eat sushi with tuna from Japan in a restaurant

| | action | stack | queue |
|---|---|---|---|
| 0 | - | | I eat sushi ... |
| 1 | shift | I | eat sushi with ... |
| 2 | shift | I eat | sushi with tuna ... |
| 3 | l-reduce | eat | sushi with tuna ... |
| 4 | shift | eat sushi | with tuna from ... |
| 5a | r-reduce | eat sushi | with tuna from ... |
| 5b | shift | eat sushi with | tuna from Japan ... |

# Incremental Parsing (Shift-Reduce)

I eat sushi with tuna from Japan in a restaurant

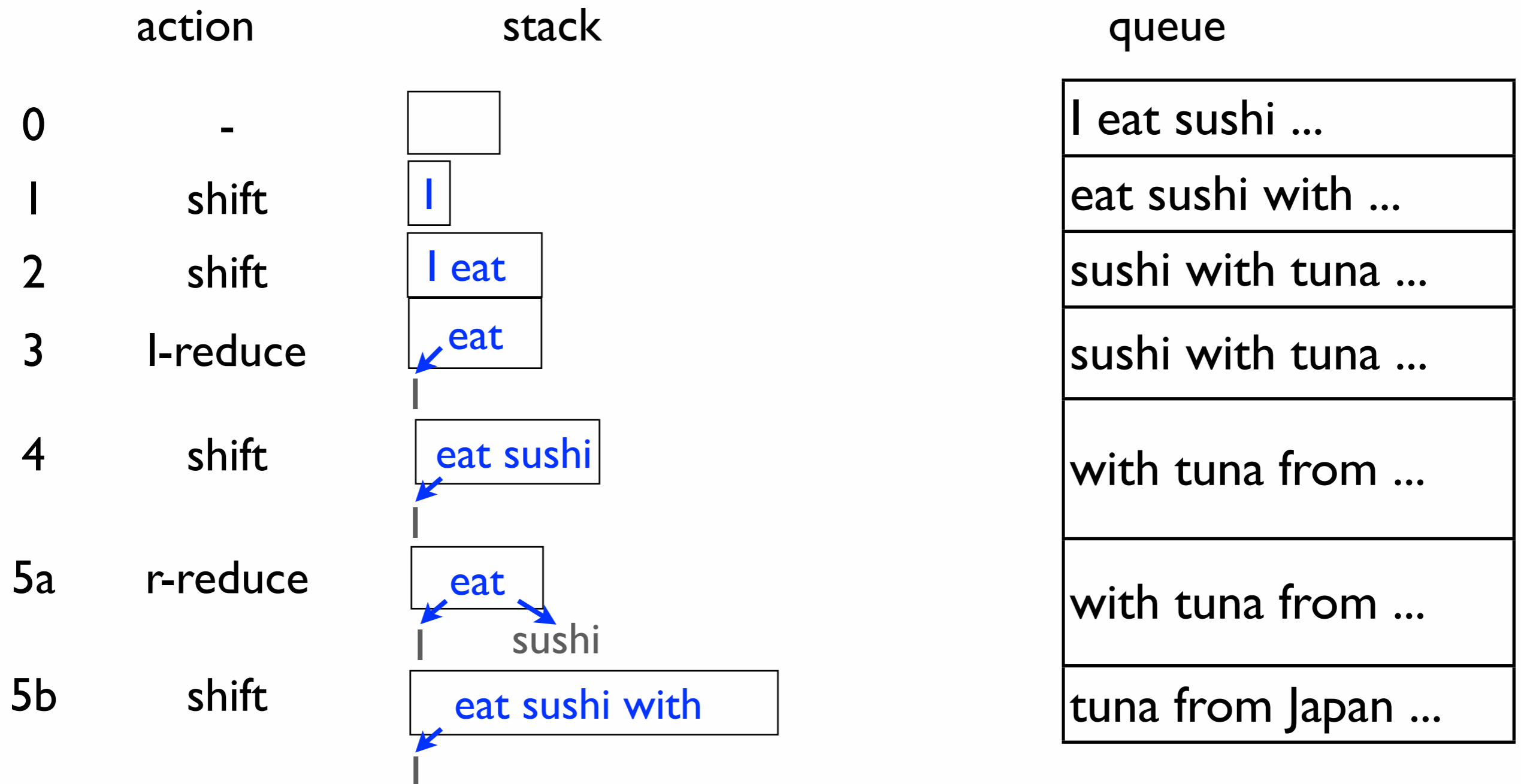| | action | stack | | queue |
|---|---|---|---|---|
| 0 | - | | | I eat sushi ... |
| 1 | shift | I | | eat sushi with ... |
| 2 | shift | I eat | | sushi with tuna ... |
| 3 | l-reduce | eat | | sushi with tuna ... |
| 4 | shift | eat sushi | | with tuna from ... |
| 5a | r-reduce | eat  sushi | | with tuna from ... |
| 5b | shift | eat sushi with | | tuna from Japan ... |

**shift-reduce conflict**

# Greedy Search

- each state => three new states (shift, l-reduce, r-reduce)

- greedy search:  always pick the best next state

  - "best" is defined by a score learned from data

sh

l-re

r-re

# Greedy Search

- each state => three new states (shift, l-reduce, r-reduce)

- greedy search: always pick the best next state

  - "best" is defined by a score learned from data

# Beam Search

- each state => three new states (shift, l-reduce, r-reduce)

- beam search: always keep top-$b$ states

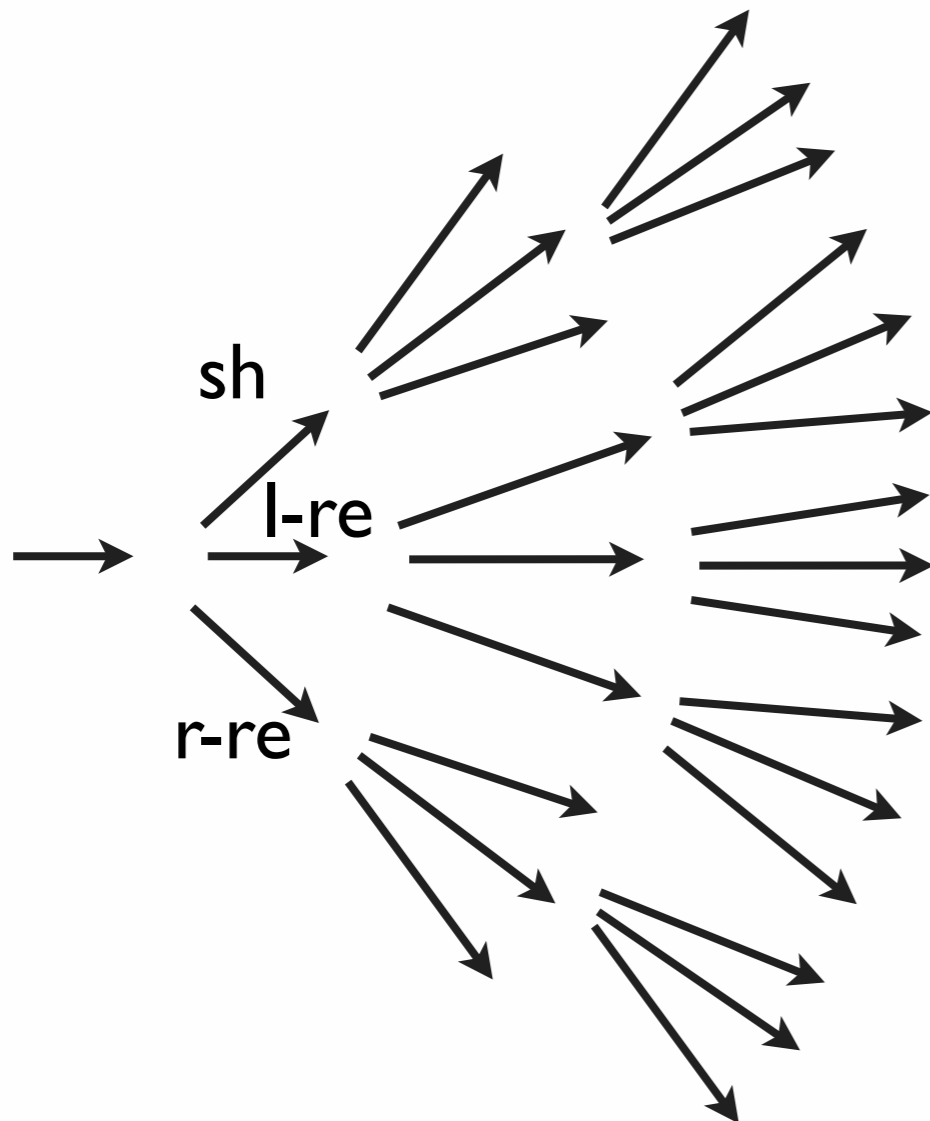  - still just a tiny fraction of the whole search space

# Beam Search

- each state => three new states (shift, l-reduce, r-reduce)

- beam search: always keep top-*b* states

  - still just a tiny fraction of the whole search space



psycholinguistic evidence: parallelism (Fodor et al, 1974; Gibson, 1991)

# Dynamic Programming

- each state => three new states (shift, l-reduce, r-reduce)

- key idea of DP: share common subproblems

  - merge equivalent states => polynomial space

# Dynamic Programming

- each state => three new states (shift, l-reduce, r-reduce)

- key idea of DP: share common subproblems

  - merge equivalent states => polynomial space

# Dynamic Programming

- each state => three new states (shift, l-reduce, r-reduce)

- key idea of DP: share common subproblems

  - merge equivalent states => polynomial space

# Dynamic Programming

- each state => three new states (shift, l-reduce, r-reduce)

- key idea of DP: share common subproblems

  - merge equivalent states => polynomial space

  each DP state corresponds to exponentially many non-DP states

graph-structured stack
(Tomita, 1986)

# Dynamic Programming

- each state => three new states (shift, l-reduce, r-reduce)

- key idea of DP: share common subproblems

  - merge equivalent states => polynomial space

each DP state corresponds to exponentially many non-DP states

# Dynamic Programming

- each state => three new states (shift, l-reduce, r-reduce)

- key idea of DP: share common subproblems

  - merge equivalent states => polynomial space

each DP state corresponds to exponentially many non-DP states



graph-structured stack
(Tomita, 1986)

(Huang and Sagae, 2010)

# Merging (Ambiguity Packing)

- two states are equivalent if they agree on features

  - because same features guarantee same cost

  - example: if we only care about the last 2 words on stack



(Huang and Sagae, 2010)

# Merging (Ambiguity Packing)

- two states are equivalent if they agree on features
  - because same features guarantee same cost
  - example: if we only care about the last 2 words on stack



two equivalent classes

I sushi  … I sushi

I eat sushi  … eat sushi

eat sushi

(Huang and Sagae, 2010)

# Merging (Ambiguity Packing)

- two states are equivalent if they agree on features

  - because same features guarantee same cost

  - example: if we only care about the last 2 words on stack

two equivalent classes

I  sushi        ... I sushi

I  eat  sushi        ... eat sushi

eat  sushi

psycholinguistic evidence
(eye-tracking experiments):

**delayed disambiguation**

John and Mary had 2 papers
John and Mary had 2 papers

Frazier and Rayner (1990), Frazier (1999)

# Merging (Ambiguity Packing)

- two states are equivalent if they agree on features

  - because same features guarantee same cost

  - example: if we only care about the last 2 words on stack



two equivalent classes

I sushi

... I sushi

I eat sushi

eat sushi

... eat sushi

psycholinguistic evidence
(eye-tracking experiments):

**delayed disambiguation**

John and Mary had 2 papers each
John and Mary had 2 papers together

Frazier and Rayner (1990), Frazier (1999)

# Result: linear-time, DP, and accurate!

- very fast linear-time dynamic programming parser

- explores *exponentially* many trees (and outputs forest)

- state-of-the-art parsing accuracy on English & Chinese

# Result: linear-time, DP, and accurate!

- very fast linear-time dynamic programming parser

- explores *exponentially* many trees (and outputs forest)

- state-of-the-art parsing accuracy on English & Chinese

# Result: linear-time, DP, and accurate!

- very fast linear-time dynamic programming parser

- explores *exponentially* many trees (and outputs forest)

- state-of-the-art parsing accuracy on English & Chinese

# In this talk...

- Background

- Dynamic Programming for Incremental Parsing

- Features: from sparse to neural to recurrent neural nets

- Bidirectional RNNs: minimal features; no tree structures!

  - dependency parsing (Kiperwaser+Goldberg, 2016, Cross+Huang, 2016a)

  - span-based constituency parsing (Cross+Huang, 2016b)

- Marrying DP & RNNs *(mostly not my work!)*

  - minimal span-based constituency parsing (Stern et al, ACL 2017)

  - transition-based dependency parsing (Shi et al, EMNLP 2017)

# Sparse Features



- score each action using features $\mathbf{f}$ and weights $\mathbf{w}$

  - features are drawn from a *local window*

    - abstraction (or signature) of a state -- this inspires DP!

  - weights trained by structured perceptron (Collins 02)

# Sparse Features

← stack   queue →

... feed cats    in the garden ...

nearby

← stack   queue →

...   $s_2$   $s_1$   $s_0$   $q_0$ $q_1$ ...

- score each action using features **f** and weights **w**
  - features are drawn from a *local window*
    - abstraction (or signature) of a state -- this inspires DP!
  - weights trained by structured perceptron (Collins 02)

# Sparse Features

← stack   queue →

... feed cats | in the garden ...

nearby

features:
$(s_0.w, s_0.rc, q_0, ...) = (\text{cats}, \text{nearby}, \text{in}, ...)$

← stack   queue →

... $s_2$   $s_1$   $s_0$   $q_0$ $q_1$ ...

- score each action using features **f** and weights **w**

  - features are drawn from a *local window*

    - abstraction (or signature) of a state -- this inspires DP!

  - weights trained by structured perceptron (Collins 02)

# From Sparse to Neural to RNN

# From Sparse to Neural to RNN

- neural nets can automate feature engineering :-)

- but early neural work (e.g., Chen+Manning 14) still use lots of manually designed atomic features on the stack

# From Sparse to Neural to RNN

$$\cdots \mid s_2 \mid s_1 \mid s_0 \mid \qquad \mid b_0 \mid b_1 \mid b_2 \mid \cdots$$

(Chen+Manning 2014)

$$s_1.lc_i \quad \cdots \quad s_1.rc_i \quad s_0.lc_i \quad \cdots \quad s_0.rc_i$$

$$s_1.lc_0.lc_0 \qquad s_1.rc_0.rc_0 \quad s_0.lc_0.lc_0 \qquad s_0.rc_0.rc_0$$

- neural nets can automate feature engineering :-)

- but early neural work (e.g., Chen+Manning 14) still use lots of manually designed atomic features on the stack

- can we automate even more?

  - option 1: summarize the whole stack (part of **y**) using RNNs => stack LSTM / RNNG (Dyer+ 15, 16)

  - option 2: summarize the whole input (**x**) using RNNs => biLSTM dependency parsing (Kiperwaser+Goldberg 16, Cross+Huang 16a) biLSTM constituency parsing (Cross+Huang 16b)

# From Sparse to Neural to RNN

- neural nets can automate feature engineering :-)

- but early neural work (e.g., Chen+Manning 14) still use lots of manually designed atomic features on the stack

- can we automate even more?

  - option 1: summarize the whole stack (part of **y**) using RNNs => stack LSTM / RNNG (Dyer+ 15, 16)    rules out DP! :(

  - option 2: summarize the whole input (**x**) using RNNs => biLSTM dependency parsing (Kiperwaser+Goldberg 16, Cross+Huang 16a) biLSTM constituency parsing (Cross+Huang 16b)    enables DP! :)

# Spectrum: Neural Incremental Parsing

constituency
dependency

*bottom-up*

*edge-factored*
*(McDonald+ 05a)*

DP incremental parsing
(Huang+Sagae 10, Kuhlmann+ 11)

Feedforward NNs
*(Chen + Manning 14)*

Stack LSTM

*(Dyer+ 15)*

RNNG
*(Dyer+ 16)*

biRNN dependency
*(Kiperwaser+Goldberg 16;*
*Cross+Huang 16a)*

*biRNN graph-based*
*dependency*
*(Kiperwaser+Goldberg 16;*
*Wang+Chang 16)*

biRNN span-based
constituency
*(Cross+Huang 16b)*

*minimal span-based*
*constituency*
*(Stern+ ACL 17)*

minimal dependency
*(Shi+ EMNLP 17)*

*all tree info*
*(summarize output **y**)*

*minimal or no tree info*
*(summarize input **x**)*

DP impossible    enables slow DP        enables fast DP      fastest DP: $O(n^3)$ 23

# In this talk...

- Background

- Dynamic Programming for Incremental Parsing

- Interlude: NN Features: from feedforward to recurrent

- **Bidirectional RNNs: minimal features; no tree structures!**

  - dependency parsing (Kiperwaser+Goldberg, 2016, Cross+Huang, 2016a)

  - span-based constituency parsing (Cross+Huang, 2016b)

- **Marrying DP & RNNs** *(mostly not my work!)*

  - minimal span-based constituency parsing (Stern et al, ACL 2017)

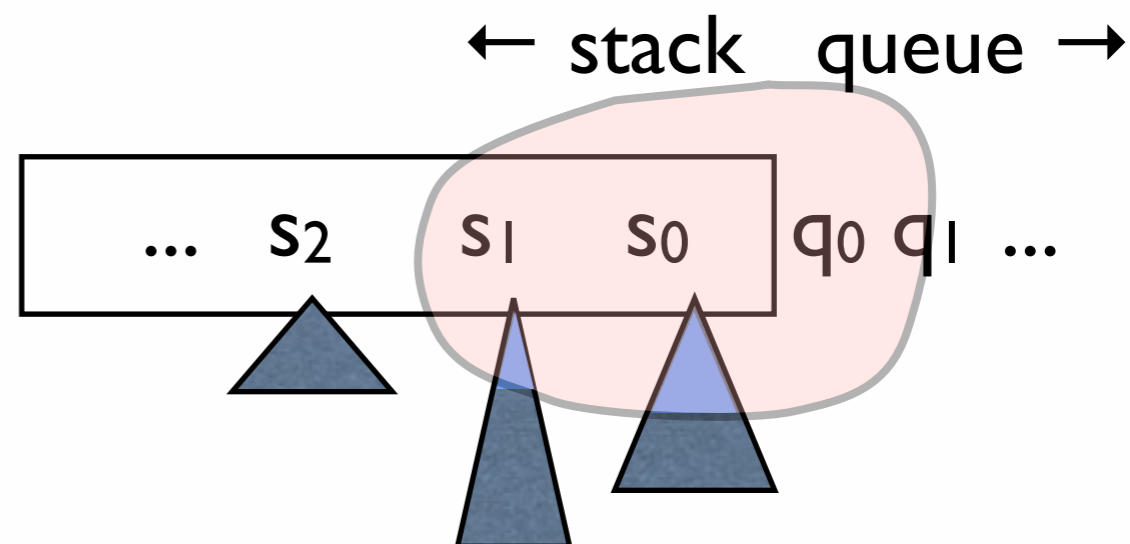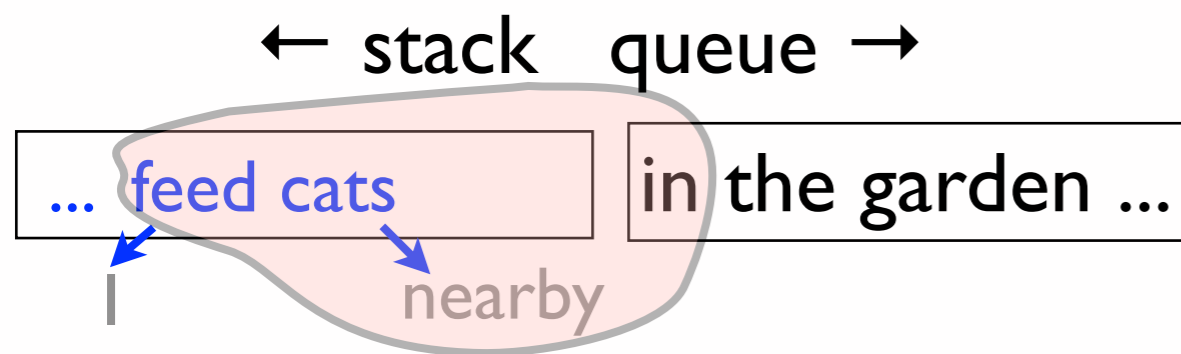  - transition-based dependency parsing (Shi et al, EMNLP 2017)

# biRNN for Dependency Parsing

- several parallel efforts in 2016 used biLSTM features

  - Kiperwaser+Goldberg 2016: four positional feats; arc-eager

  - Cross+Huang ACL 2016: three positional feats; arc-standard

  - Wang+Chang 2016: two positional feats; graph-based

- all inspired by sparse edge-factored model (McDonald+05)

  - use positions to summarize the input **x**, not the output **y**!

  - => $O(n^3)$ DP, e.g. graph-based, but also incremental!

(Kiperwaser and Goldberg 2016)

(Cross and Huang, ACL 2016)

these developments lead to state-of-the-art in dependency parsing

# Span-Based Constituency Parsing

- previous work uses tree structures on stack

- we simplify to operate directly on sentence spans

- simple-to-implement linear-time parsing

*previous work*

**Stack**

NP    VP'

I/PRP    do/MD    like/VBP

**Queue**

eating/VBG    fish/NN

*our work*

**Stack**

I/PRP    do/MD    like/VBP

0    1    3

**Queue**

eating/VBG    fish/NN

4    5

| Structural (even step) | Shift |
| Structural (even step) | Combine |
| Label (odd step) | Label-*X* |
| Label (odd step) | No-Label |

S
├─ NP
│  └─ PRP
│     └─ I
└─ VP
   ├─ MD — do
   ├─ VBP — like
   └─ S
      └─ VP
         ├─ VBG — eating
         └─ NP
            └─ NN — fish

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1       2          3            4        5

current brackets    t = {}

| Structural (even step) | Shift |
| | Combine |
| Label (odd step) | Label-X |
| | No-Label |

S
NP        VP
PRP  MD  VBP        S
I    do  like       VP
                VBG    NP
                eating  NN
                        fish

I/PRP  do/MD  like/VBP  eating/VBG  fish/NN
0      1      2         3           4       5

current brackets    t = {}

⇩ *Shift*

I/PRP    do/MD  like/VBP  eating/VBG  fish/NN
0        1      2         3           4       5

| Structural (even step) | Shift |
| --- | --- |
| | Combine |
| Label (odd step) | Label-*X* |
| | No-Label |

S
├─ NP
│   └─ PRP
│       └─ I
└─ VP
    ├─ MD
    │   └─ do
    ├─ VBP
    │   └─ like
    └─ S
        └─ VP
            ├─ VBG
            │   └─ eating
            └─ NP
                └─ NN
                    └─ fish

```
( )  [ I/PRP    do/MD    like/VBP    eating/VBG    fish/NN ]
   0         1         2           3             4         5
```

current brackets    $t = \{\}$

$\Downarrow$ *Shift*

```
( I/PRP )  [ do/MD    like/VBP    eating/VBG    fish/NN ]
 0        1         2           3             4         5
```

*Label-NP*    $t = \{_0 NP_1\}$

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-*X* |
| | No-Label |

S
├── NP
│   └── PRP
│       └── I
└── VP
    ├── MD
    │   └── do
    ├── VBP
    │   └── like
    └── S
        └── VP
            ├── VBG
            │   └── eating
            └── NP
                └── NN
                    └── fish

| | I/PRP | do/MD | like/VBP | eating/VBG | fish/NN | | current brackets | $t = \{\}$ |
|---|---|---|---|---|---|---|---|---|

*Shift*

| I/PRP | do/MD | like/VBP | eating/VBG | fish/NN | *Label-NP* | $t = \{_0 NP_1\}$ |
|---|---|---|---|---|---|---|

*Shift*

| I/PRP | do/MD | like/VBP | eating/VBG | fish/NN |
|---|---|---|---|---|

Liang Huang (Oregon State)      (Cross and Huang, EMNLP 2016)      27

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-*X* |
| | No-Label |

S
├── NP
│    └── PRP
│         └── I
└── VP
     ├── MD
     │    └── do
     ├── VBP
     │    └── like
     └── S
          └── VP
               ├── VBG
               │    └── eating
               └── NP
                    └── NN
                         └── fish



```
  I/PRP   do/MD   like/VBP   eating/VBG   fish/NN        current      t = {}
0       1       2          3            4          5     brackets
```

⇓ *Shift*

```
  I/PRP   do/MD   like/VBP   eating/VBG   fish/NN      Label-NP     t = {₀NP₁}
0       1       2          3            4          5
```

$t = \{_0 NP_1\}$

⇓ *Shift*

```
  I/PRP   do/MD   like/VBP   eating/VBG   fish/NN      No-Label     t = {₀NP₁}
0       1       2          3            4          5
```

$t = \{_0 NP_1\}$

| Structural (even step) | Shift |
| --- | --- |
| | Combine |
| Label (odd step) | Label-*X* |
| | No-Label |

```
          S
        /   \
      NP      VP
      |      / \
     PRP   MD  VBP    S
      |    |    |     |
      I    do  like   VP
                      / \
                    VBG  NP
                     |   |
                  eating NN
                         |
                        fish
```

current brackets

$[\ ]_0$ I/PRP $_1$ do/MD $_2$ like/VBP $_3$ eating/VBG $_4$ fish/NN $_5$    $t = \{\}$

⇩ *Shift*

$_0$ I/PRP $_1$ do/MD $_2$ like/VBP $_3$ eating/VBG $_4$ fish/NN $_5$    *Label-NP*    $t = \{_0 NP_1\}$

⇩ *Shift*

$_0$ I/PRP do/MD $_1$  $_2$ like/VBP $_3$ eating/VBG $_4$ fish/NN $_5$    *No-Label*    $t = \{_0 NP_1\}$

⇩ *Shift*

$_0$ I/PRP do/MD $_1$ like/VBP $_2$  $_3$ eating/VBG $_4$ fish/NN $_5$

| | | |
|---|---|---|
| **Structural (even step)** | Shift | |
| | Combine | |
| **Label (odd step)** | Label-*X* | |
| | No-Label | |

S
NP VP
PRP MD VBP S
I do like VP
VBG NP
eating NN
fish



| | | | |
|---|---|---|---|
| 0 I/PRP 1 do/MD 2 like/VBP 3 eating/VBG 4 fish/NN 5 | | current brackets | $t = \{\}$ |

*Shift*

| | | |
|---|---|---|
| 0 I/PRP 1 | do/MD 2 like/VBP 3 eating/VBG 4 fish/NN 5 | *Label-NP* | $t = \{_0NP_1\}$ |

*Shift*

| | | |
|---|---|---|
| 0 I/PRP 1 do/MD 2 | like/VBP 3 eating/VBG 4 fish/NN 5 | *No-Label* | $t = \{_0NP_1\}$ |

*Shift*

| | | |
|---|---|---|
| 0 I/PRP 1 do/MD 2 like/VBP 3 | eating/VBG 4 fish/NN 5 | *No-Label* | $t = \{_0NP_1\}$ |

(Cross and Huang, EMNLP 2016)

| | |
|---|---|
| **Structural (even step)** | Shift |
| | Combine |
| **Label (odd step)** | Label-$X$ |
| | No-Label |

S
NP    VP
PRP  MD  VBP   S
I    do  like  VP
VBG   NP
eating  NN
fish

| I/PRP | do/MD | like/VBP | | eating/VBG | fish/NN | |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | |

$t = \{_0 NP_1\}$

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-$X$ |
| | No-Label |



$t = \{_0 NP_1\}$

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-$X$ |
| | No-Label |

S
├ NP
│ └ PRP
│    └ I
└ VP
   ├ MD VBP
   │  do  like
   └ S
      └ VP
         ├ VBG
         │  eating
         └ NP
            └ NN
               └ fish

| I/PRP | do/MD | like/VBP | | eating/VBG | fish/NN | | $t = \{_0NP_1\}$ |

*Combine*

| I/PRP | do/MD        like/VBP | | eating/VBG | fish/NN | | *No-Label* | $t = \{_0NP_1\}$ |

(Cross and Huang, EMNLP 2016)

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-$X$ |
| | No-Label |



S
NP → PRP → I
VP
MD → do
VBP → like
S
VP
VBG → eating
NP → NN → fish

| | I/PRP | do/MD | like/VBP | | eating/VBG | fish/NN | | $t = \{_0NP_1\}$ |
| 0 | | 1 | | 2 | | 3 | | 4 | 5 |

⇩ *Combine*

| | I/PRP | do/MD      like/VBP | | eating/VBG | fish/NN | |
| 0 | | 1 | | 3 | | 4 | 5 |

*No-Label* ⇨   $t = \{_0NP_1\}$

⇩ *Shift*

| | I/PRP | do/MD      like/VBP | | eating/VBG | | fish/NN | |
| 0 | | 1 | | 3 | | 4 | 5 |

(Cross and Huang, EMNLP 2016)

| Structural (even step) | Shift |
| --- | --- |
| | Combine |
| Label (odd step) | Label-X |
| | No-Label |

S
NP    VP
PRP    MD  VBP      S
I     do   like     VP
                VBG   NP
              eating   NN
                      fish

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1       2          3            4        5

$t = \{_0 NP_1\}$

*Combine*

I/PRP   do/MD        like/VBP   eating/VBG   fish/NN
0       1                       3            4        5

*No-Label*   $t = \{_0 NP_1\}$

*Shift*

I/PRP   do/MD        like/VBP   eating/VBG   fish/NN
0       1                       3            4        5

*No-Label*   $t = \{_0 NP_1\}$

| Structural (even step) | Shift |
| --- | --- |
| | Combine |
| Label (odd step) | Label-*X* |
| | No-Label |

$t = \{_0NP_1\}$

*Combine*

$t = \{_0NP_1\}$

*No-Label*

*Shift*

*No-Label*   $t = \{_0NP_1\}$

*Shift*

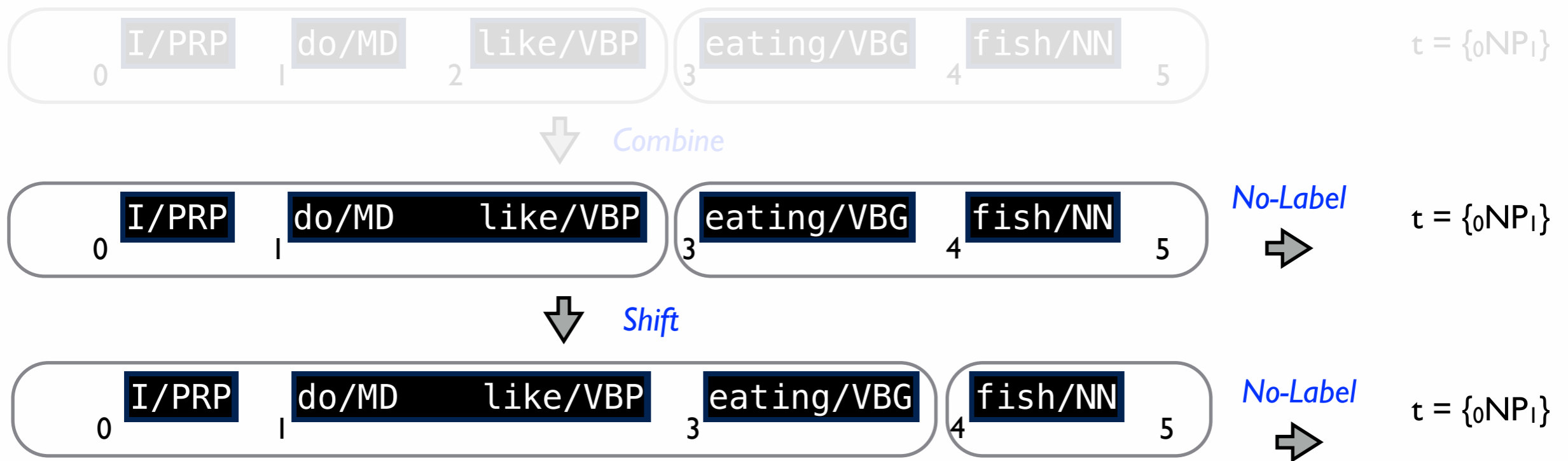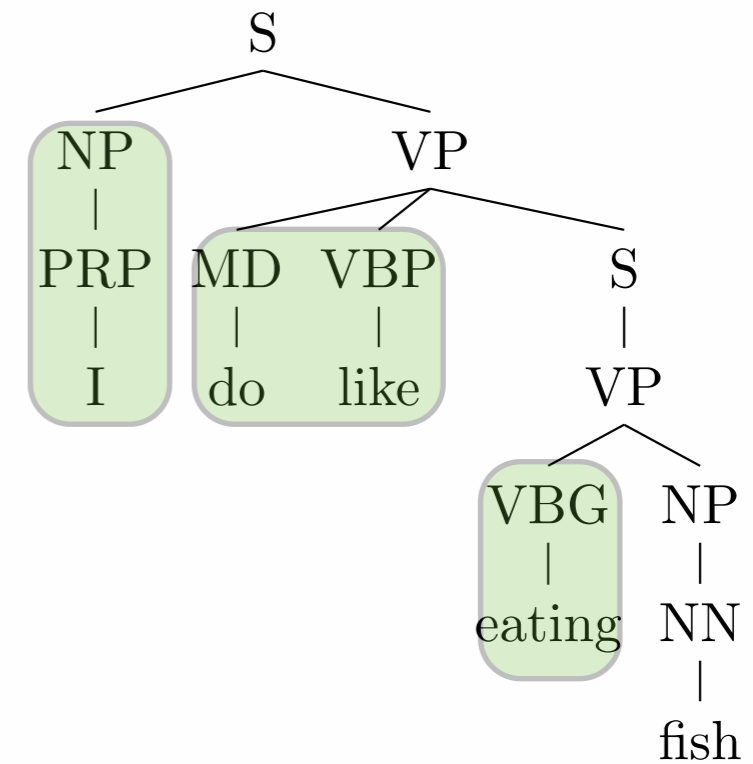| | | |
|---|---|---|
| **Structural (even step)** | Shift | |
| | Combine | |
| **Label (odd step)** | Label-*X* | |
| | No-Label | |



S

NP     VP

PRP   MD   VBP    S

I    do    like    VP

VBG   NP

eating   NN

fish

| | | |
|---|---|---|
| I/PRP   do/MD   like/VBP    eating/VBG   fish/NN | | $t = \{_0NP_1\}$ |
| 0    1    2    3    4    5 | | |

*Combine*

I/PRP   do/MD    like/VBP    eating/VBG   fish/NN     *No-Label*    $t = \{_0NP_1\}$

0    1     3    4    5

*Shift*

I/PRP   do/MD    like/VBP   eating/VBG    fish/NN     *No-Label*    $t = \{_0NP_1\}$

0    1      3    4    5

*Shift*

I/PRP   do/MD    like/VBP   eating/VBG   fish/NN     *Label-NP*    $t = \{_0NP_1, _4NP_5\}$

0    1      3    4    5

(Cross and Huang, EMNLP 2016)

| Structural (even step) | Shift |
|---|---|
| | Combine |
| **Label (odd step)** | Label-$X$ |
| | No-Label |

S
├── NP
│   └── PRP
│       └── I
└── VP
    ├── MD — do
    ├── VBP — like
    └── S
        └── VP
            ├── VBG — eating
            └── NP
                └── NN — fish

```
  I/PRP      do/MD      like/VBP      eating/VBG      fish/NN
0         1                        3             4              5
```

$t = \{_0 NP_1, \, _4 NP_5\}$

| | |
|---|---|
| **Structural (even step)** | Shift |
| | Combine |
| **Label (odd step)** | Label-$X$ |
| | No-Label |

S
NP — PRP — I
VP
MD — do
VBP — like
S
VP
VBG — eating
NP — NN — fish

$t = \{_0 NP_1, {_4}NP_5\}$

I/PRP do/MD like/VBP eating/VBG fish/NN
0 1 3 4 5

*Combine*

I/PRP do/MD like/VBP eating/VBG fish/NN
0 1 3 5

(Cross and Huang, EMNLP 2016)

| | |
|---|---|
| **Structural (even step)** | Shift |
| | Combine |
| **Label (odd step)** | Label-$X$ |
| | No-Label |

S
├─ NP
│  └─ PRP
│     └─ I
└─ VP
   ├─ MD → do
   ├─ VBP → like
   └─ S
      └─ VP
         ├─ VBG → eating
         └─ NP
            └─ NN → fish

| I/PRP | do/MD    like/VBP | eating/VBG | fish/NN |
| 0 | 1 | 3 | 4 | 5 |

$t = \{_0NP_1, _4NP_5\}$

⬇ *Combine*

| I/PRP | do/MD    like/VBP | eating/VBG    fish/NN |
| 0 | 1 | 3 | 5 |

*Label-S-VP* ⇨

$t = \{_0NP_1, _4NP_5, _3S_5, _3VP_5\}$

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-$X$ |
| | No-Label |

S
NP → PRP → I
VP
MD → do
VBP → like
S → VP
VBG → eating
NP → NN → fish

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1                  3            4        5

$t = \{_0NP_1, {_4}NP_5\}$

*Combine*

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1                  3                     5

*Label-S-VP*

$t = \{_0NP_1, {_4}NP_5, {_3}S_5, {_3}VP_5\}$

*Combine*

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1                                    5

(Cross and Huang, EMNLP 2016)

| | |
|---|---|
| **Structural (even step)** | Shift |
| | Combine |
| **Label (odd step)** | Label-*X* |
| | No-Label |

S
NP — VP
PRP — MD VBP — S
I — do like — VP
VBG — NP
eating — NN
fish

I/PRP  do/MD  like/VBP  eating/VBG  fish/NN
0      1              3                5

$t = \{_0NP_1, _4NP_5\}$

⬇ *Combine*

I/PRP  do/MD  like/VBP  eating/VBG  fish/NN
0      1              3                5

*Label-S-VP* ➡  $t = \{_0NP_1, _4NP_5, _3S_5, _3VP_5\}$

⬇ *Combine*

I/PRP  do/MD  like/VBP  eating/VBG  fish/NN
0      1                               5

*Label-VP* ➡  $t = \{_0NP_1, _4NP_5, _3S_5, _3VP_5, _1VP_5\}$

| Structural (even step) | Shift |
|---|---|
| | Combine |
| Label (odd step) | Label-*X* |
| | No-Label |

```
          S
        /   \
      NP      VP
      |      / | \
     PRP   MD VBP    S
      |    |   |     |
      I   do like   VP
                   /  \
                 VBG   NP
                  |    |
               eating  NN
                       |
                      fish
```

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1                3            4        5

$t = \{_0NP_1, {}_4NP_5\}$

⬇ *Combine*

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1                3                     5

*Label-S-VP* ➡

$t = \{_0NP_1, {}_4NP_5, {}_3S_5, {}_3VP_5\}$

⬇ *Combine*

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0       1                                      5

*Label-VP* ➡

$t = \{_0NP_1, {}_4NP_5, {}_3S_5, {}_3VP_5, {}_1VP_5\}$

⬇ *Combine*

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0                                              5

| | |
|---|---|
| **Structural (even step)** | Shift |
| | Combine |
| **Label (odd step)** | Label-*X* |
| | No-Label |

S
— NP  VP
— PRP  MD  VBP  S
— I  do  like  VP
— VBG  NP
— eating  NN
— fish

$t = \{_0NP_1, {}_4NP_5\}$

*Combine*

$t = \{_0NP_1, {}_4NP_5, {}_3S_5, {}_3VP_5\}$

*Label-S-VP*

*Combine*

*Label-VP*

$t = \{_0NP_1, {}_4NP_5, {}_3S_5, {}_3VP_5, {}_1VP_5\}$

*Combine*

*Label-S*

$t = \{_0NP_1, {}_4NP_5, {}_3S_5, {}_3VP_5, {}_1VP_5, {}_0S_5\}$

I/PRP  do/MD  like/VBP  eating/VBG  fish/NN

# Bi-LSTM Span Features



- Sentence segment "eating fish" represented by two vectors:
  - Forward component: $f_5 - f_3$ (Wang and Chang, ACL 2016)
  - Backward component: $b_3 - b_5$

# Structural & Label Actions

## Structural Action: 4 spans

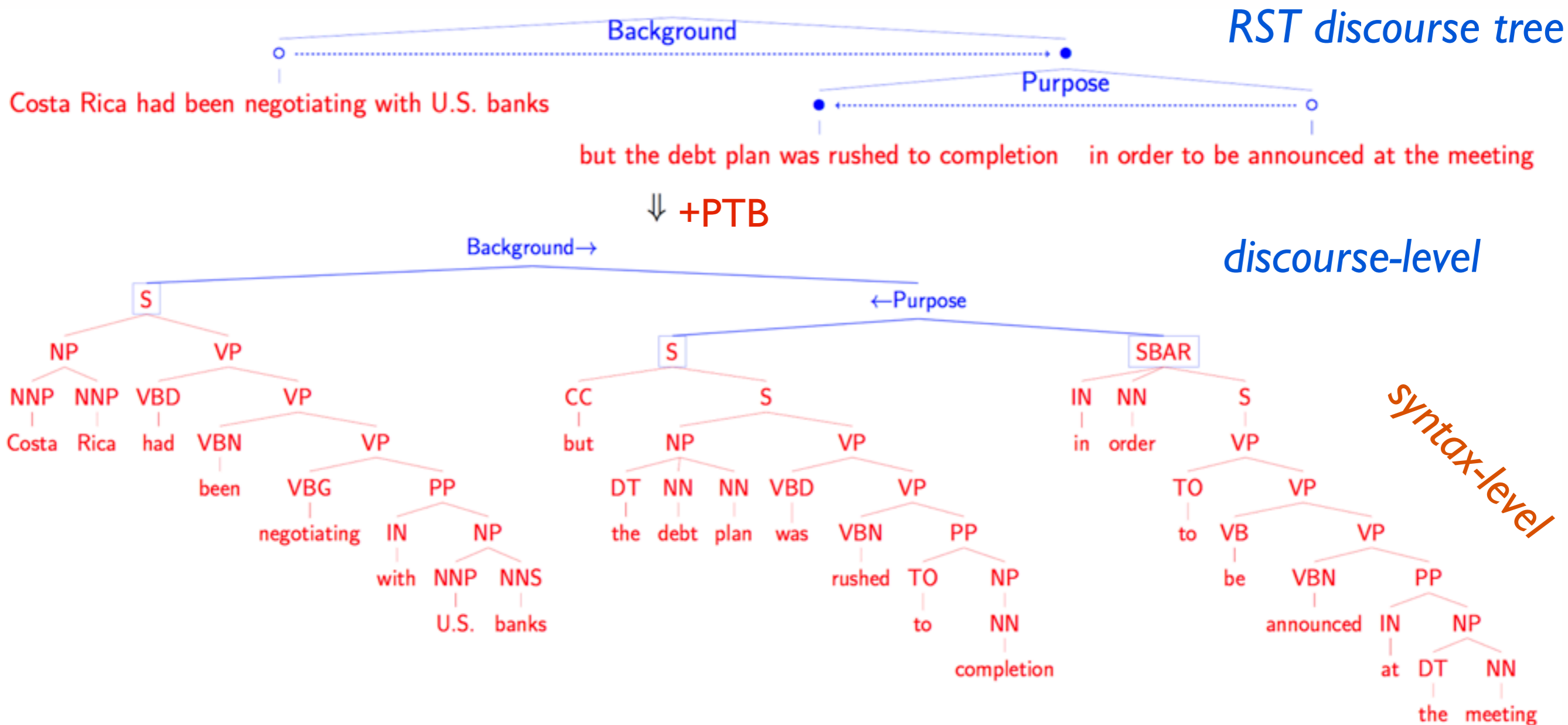`I/PRP`    `do/MD    like/VBP`    `eating/VBG   fish/NN`    `./.`

*pre-s$_1$*        *s$_1$*        *s$_0$*        *queue*

## Label Action: 3 spans

`I/PRP`    `do/MD    like/VBP    eating/VBG   fish/NN`    `./.`

*pre-s$_0$*        *s$_0$*        *queue*

# Results on Penn Treebank

| Parser | Search | Recall | Prec. | $F_1$ |
|---|---|---|---|---|
| Carreras et al. (2008) | cubic | 90.7 | 91.4 | 91.1 |
| Shindo et al. (2012) | cubic | | | 91.1 |
| Thang et al. (2015) | ~cubic | | | 91.1 |
| Watanabe et al. (2015) | beam | | | 90.7 |
| **Static Oracle** | **greedy** | 90.7 | 91.4 | 91.0 |
| **Dynamic + Exploration** | **greedy** | 90.5 | 92.1 | **91.3** |

- state of the art despite simple system with greedy actions and small embeddings trained from scratch

- first neural constituency parser to outperform sparse features

# Extension: Joint Syntax-Discourse Parsing

- extend span-based parsing to discourse parsing

  - end-to-end, joint syntactic and discourse parsing



*RST discourse tree*

*discourse-level*

*syntax-level*

⇓ +PTB

(Kai and Huang, EMNLP 2017)

# In this talk...

- Background

- Dynamic Programming for Incremental Parsing

- Interlude: NN Features: from feedforward to recurrent

- Bidirectional RNNs: minimal features; no tree structures!

  - dependency parsing (Kiperwaser+Goldberg, 2016, Cross+Huang, 2016a)

  - span-based constituency parsing (Cross+Huang, 2016b)

- Marrying DP & RNNs *(mostly not my work!)*

  - minimal span-based constituency parsing (Stern et al, ACL 2017)

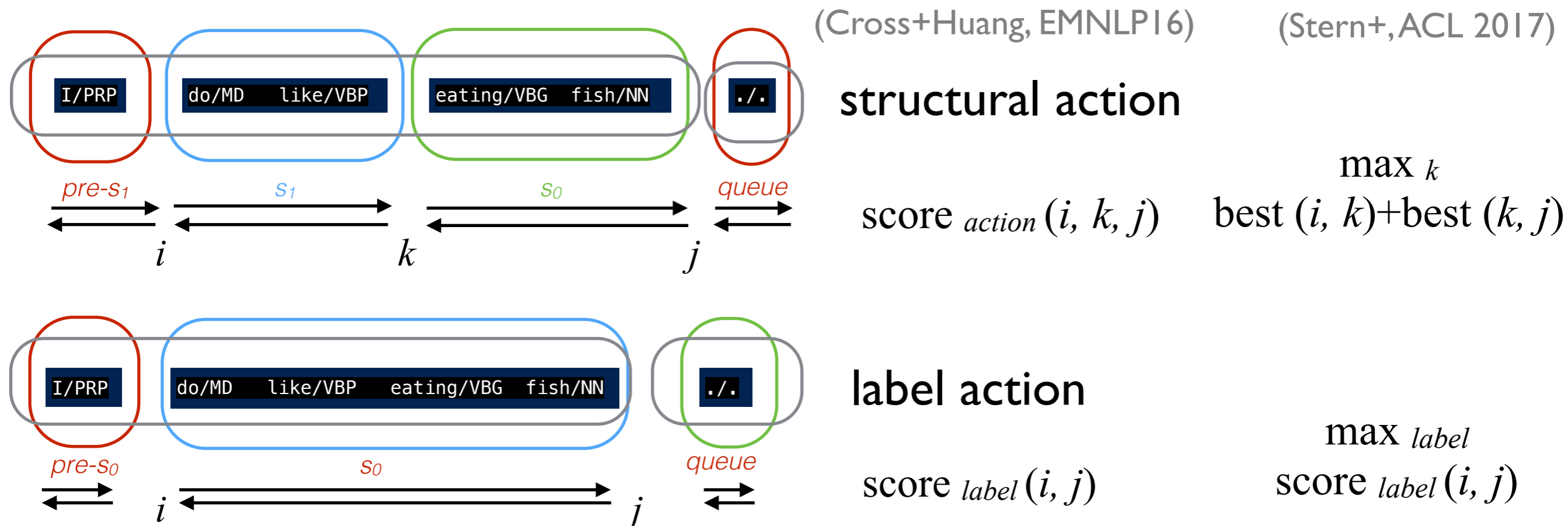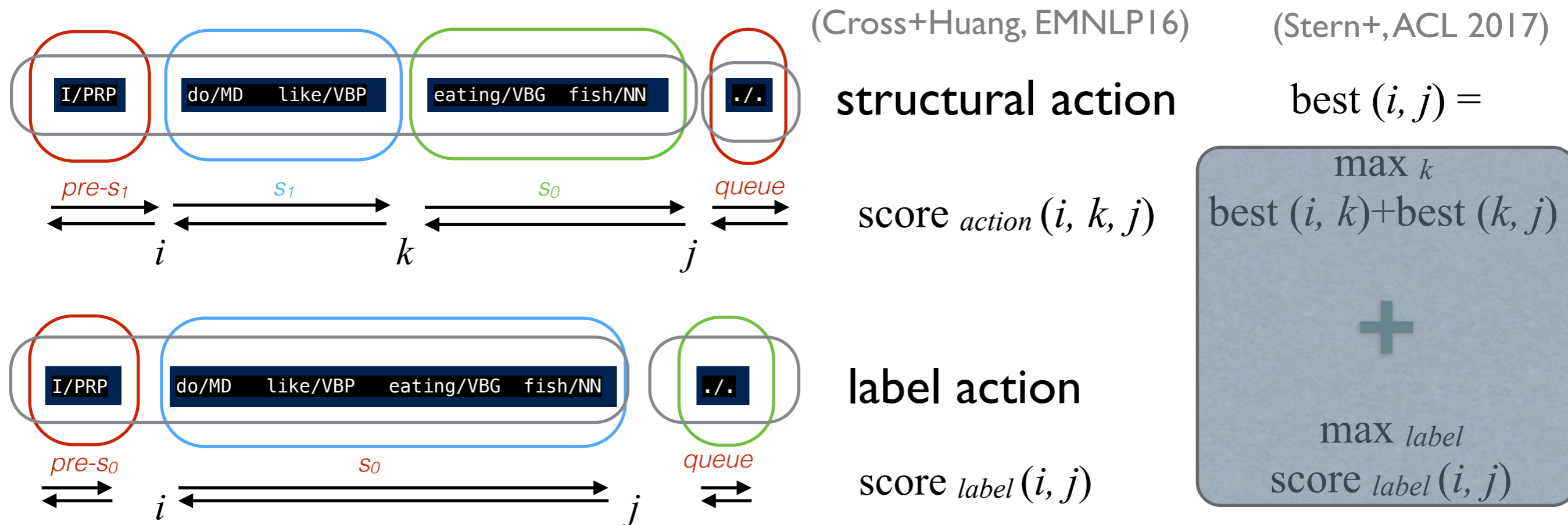  - transition-based dependency parsing (Shi et al, EMNLP 2017)

# Minimal Span-based Const. Parsing

- chart-based bottom-up parsing instead of incremental

  - an even simpler score formulation

  - $O(n^3)$ exact DP (CKY) instead of greedy search

  - global loss-augmented training instead of local training

# Minimal Span-based Const. Parsing

- chart-based bottom-up parsing instead of incremental

  - an even simpler score formulation

  - $O(n^3)$ exact DP (CKY) instead of greedy search

  - global loss-augmented training instead of local training

(Cross+Huang, EMNLP16)

# Minimal Span-based Const. Parsing

- chart-based bottom-up parsing instead of incremental

  - an even simpler score formulation

  - $O(n^3)$ exact DP (CKY) instead of greedy search

  - global loss-augmented training instead of local training

(Cross+Huang, EMNLP16)



structural action

$$\text{score}_{action}(i, k, j)$$

label action

$$\text{score}_{label}(i, j)$$

# Minimal Span-based Const. Parsing

- chart-based bottom-up parsing instead of incremental

  - an even simpler score formulation

  - $O(n^3)$ exact DP (CKY) instead of greedy search

  - global loss-augmented training instead of local training



(Cross+Huang, EMNLP16)  (Stern+, ACL 2017)

structural action

$$\max_k \text{score}_{action}(i, k, j) \quad \text{best}(i, k) + \text{best}(k, j)$$

label action

$$\max_{label} \text{score}_{label}(i, j) \quad \text{score}_{label}(i, j)$$

35

# Minimal Span-based Const. Parsing

- chart-based bottom-up parsing instead of incremental
  - an even simpler score formulation
  - $O(n^3)$ exact DP (CKY) instead of greedy search
  - global loss-augmented training instead of local training



(Cross+Huang, EMNLP16)     (Stern+, ACL 2017)

structural action

$$\text{score}_{action}(i, k, j)$$

label action

$$\text{score}_{label}(i, j)$$

$$\text{best}(i, j) =$$

$$\max_{k} \text{best}(i, k) + \text{best}(k, j)$$

$$+$$

$$\max_{label} \text{score}_{label}(i, j)$$

35

# Global Training & Loss-Augmented Decoding

want   $s_{\text{tree}}(T^*) > s_{\text{tree}}(T)$   for all   $T \neq T^*$

and larger margin for worse trees:   $s_{\text{tree}}(T^*) \geq \Delta(T, T^*) + s_{\text{tree}}(T)$

loss-augmented decoding in training (find the most-violated tree, i.e., a *bad tree* with *good score*)

$$\hat{T} = \max_T [\underbrace{\Delta(T, T^*)}_{\textit{bad tree}} + \underbrace{s_{\text{tree}}(T)}_{\textit{good score}}]$$

loss-augmented decoding for Hamming loss (approximating F1):
    simply replace score $_{label}(i, j)$

        with score $_{label}(i, j) + \mathbf{1}(label \neq label^*_{ij})$

                            gold tree label for span $(i, j)$

                            (could be "nolabel")

Liang Huang (Oregon State)

# Penn Treebank Results

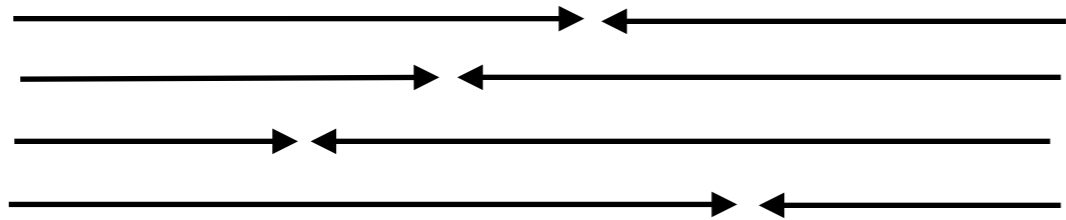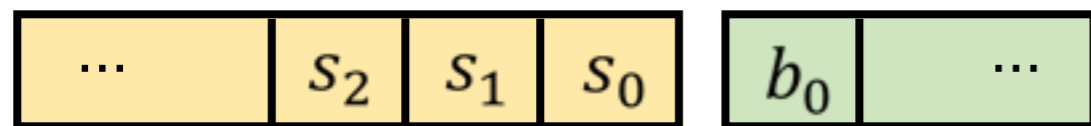| Parser | F1 Score |
|---|---|
| Hall et al. (2014) | 89.2 |
| Vinyals et al. (2015) | 88.3 |
| Cross and Huang (2016b) | 91.3 |
| Dyer et al. (2016) corrected | 91.7 |
| Liu and Zhang (2017) | 91.7 |
| Chart Parser | 91.7 |
| +refinement | 91.8 |

Liang Huang (Oregon State)
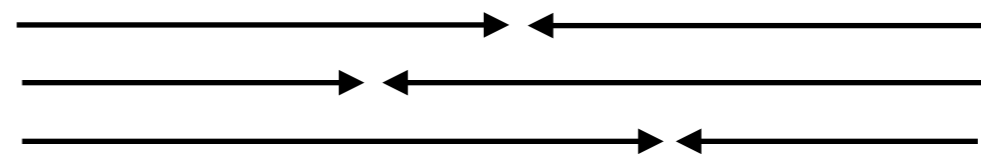
(Stern+, ACL 2017)

(Kiperwaser and Goldberg 2016)
arc-eager

(Cross and Huang, ACL 2016)
arc-standard

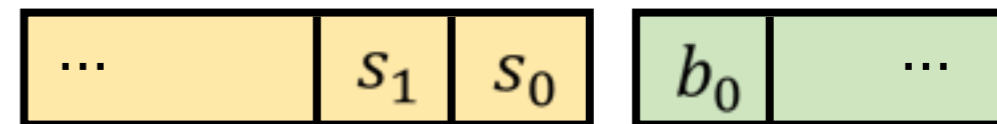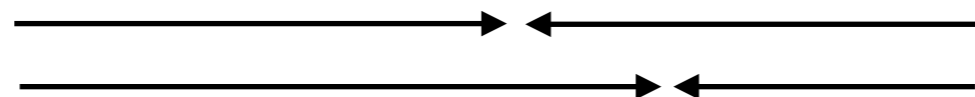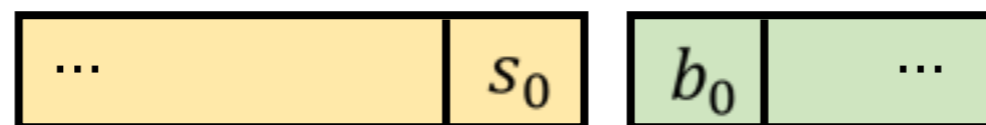# Minimal Feats for Incremental Dep. Parsing



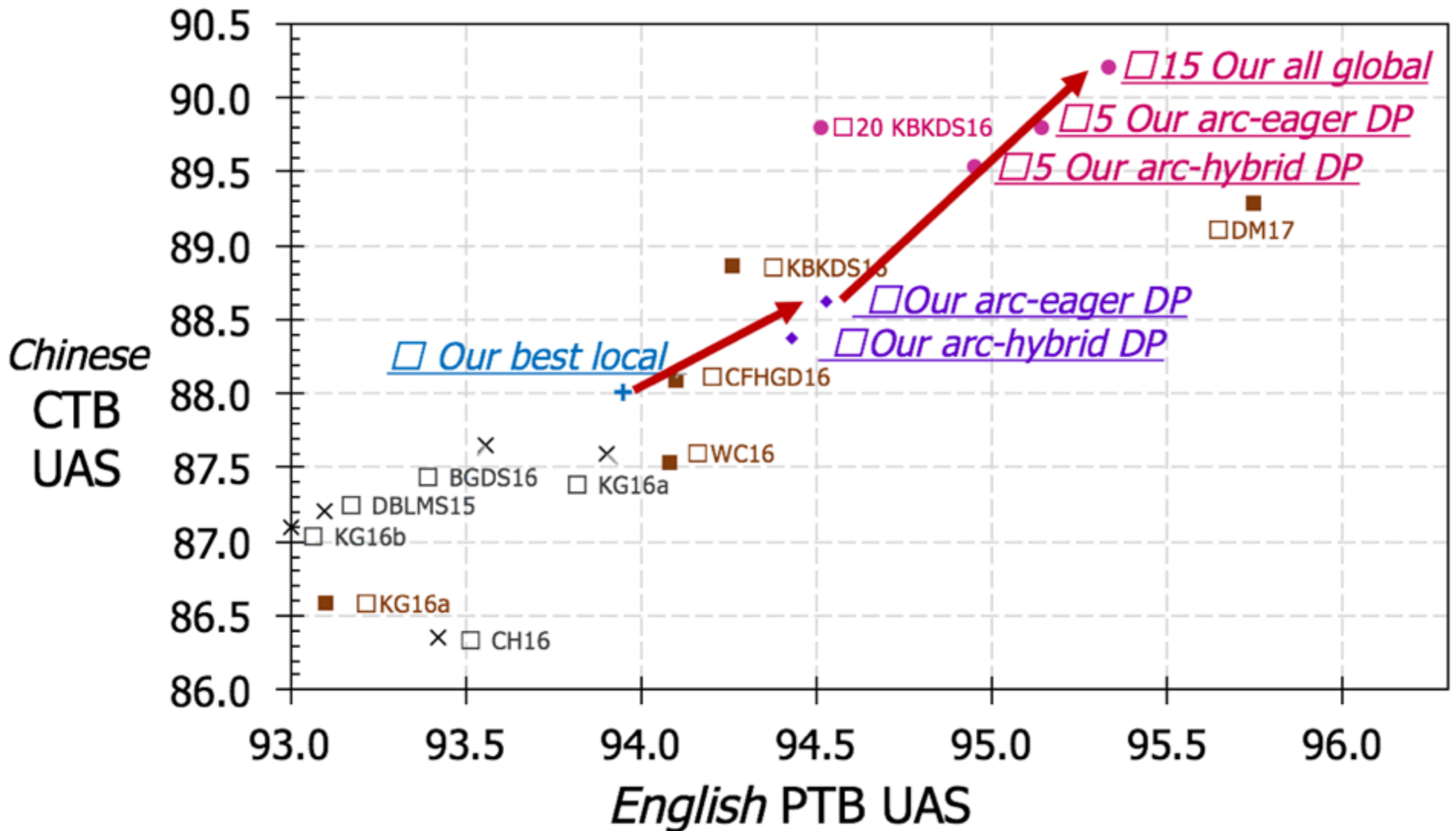(Kiperwaser and Goldberg 2016)
arc-eager

(Cross and Huang, ACL 2016)
arc-standard

(Shi, Huang, Lee, EMNLP 2017)
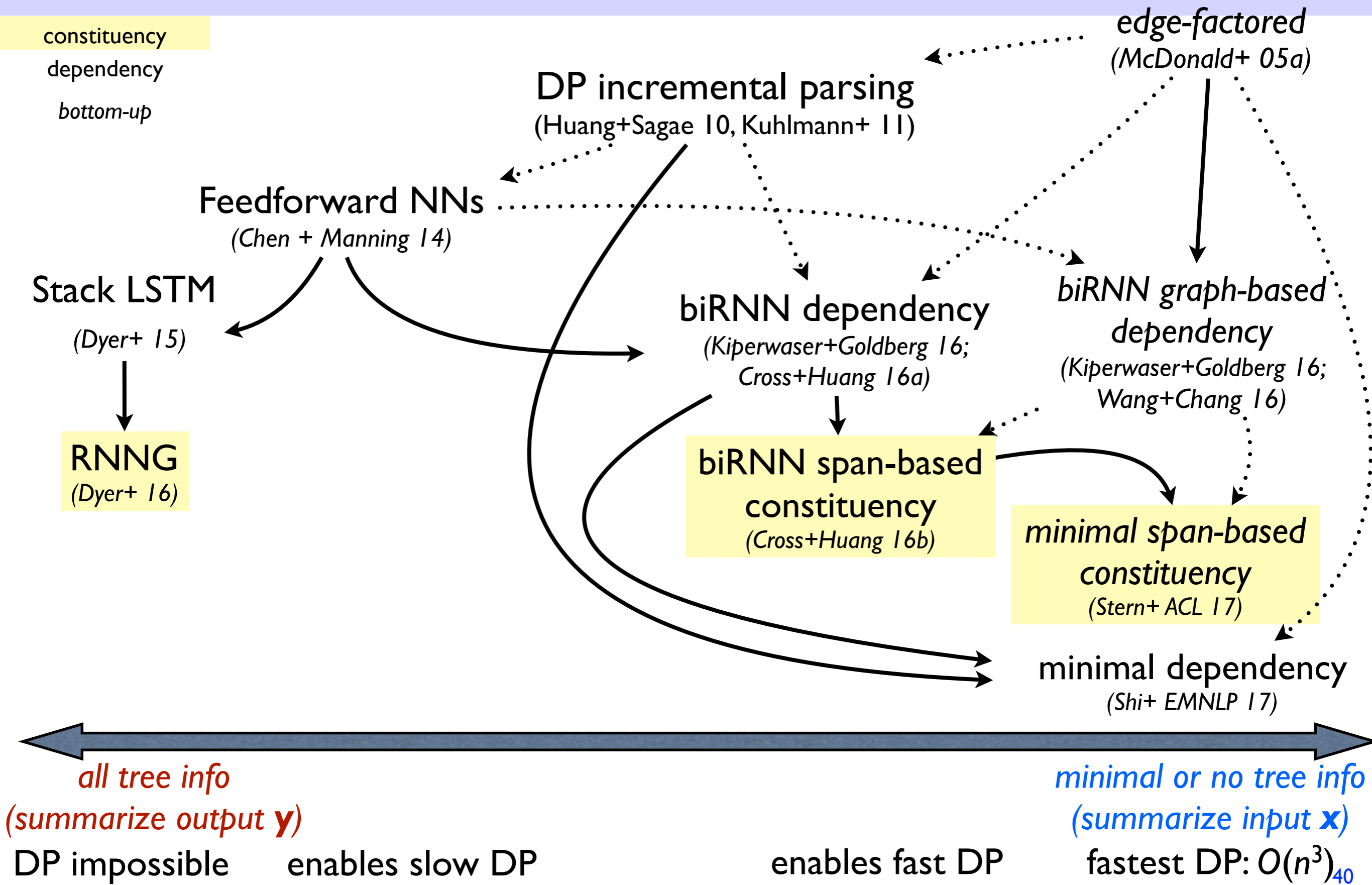Saturday talk!
arc-hybrid and arc-eager

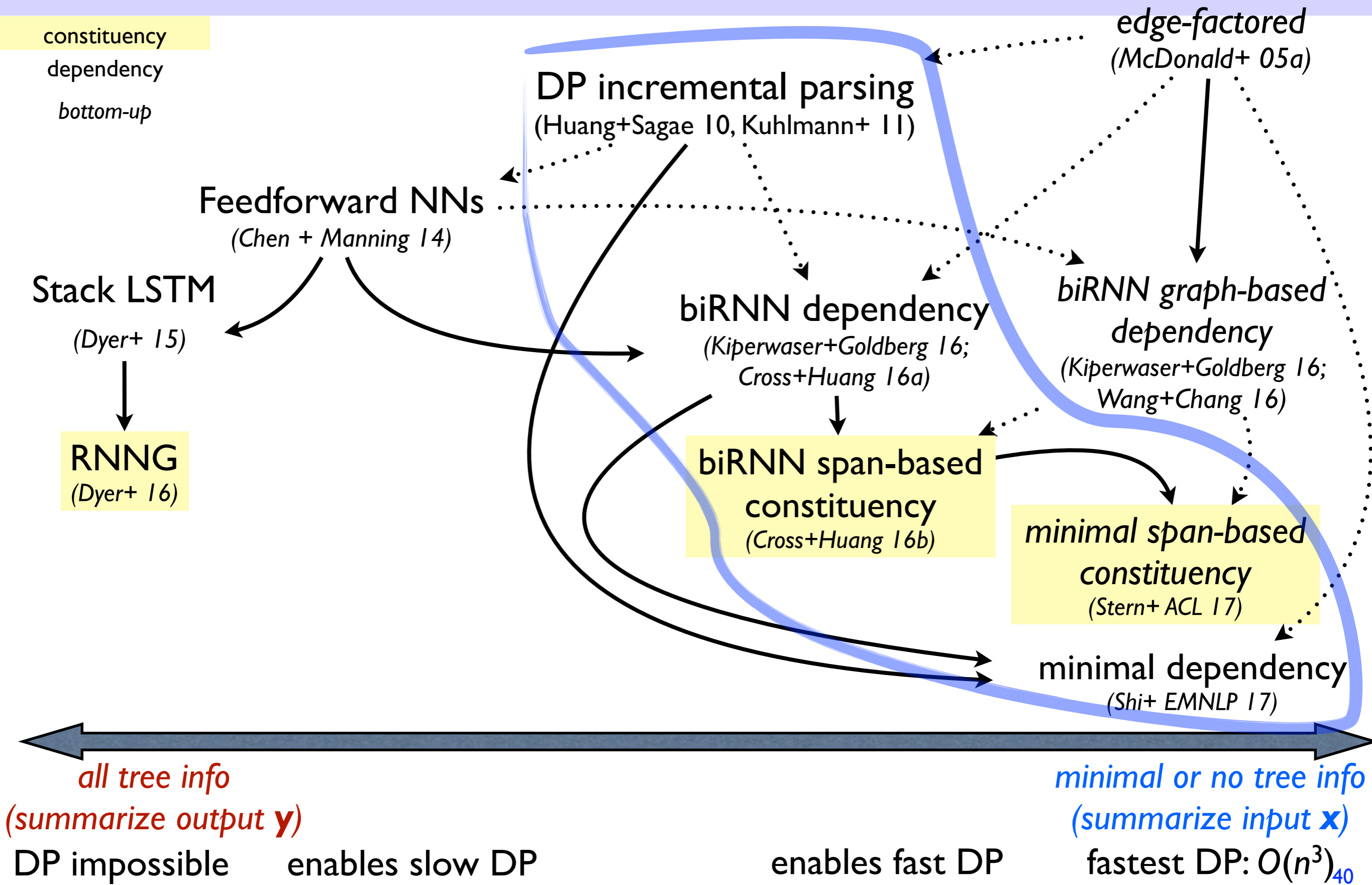works for both greedy and $O(n^3)$ DP

# Spectrum: Neural Incremental Parsing

constituency
dependency

*bottom-up*

DP incremental parsing
*(Huang+Sagae 10, Kuhlmann+ 11)*

*edge-factored*
*(McDonald+ 05a)*

Feedforward NNs
*(Chen + Manning 14)*

Stack LSTM

*(Dyer+ 15)*

RNNG
*(Dyer+ 16)*

biRNN dependency
*(Kiperwaser+Goldberg 16;
Cross+Huang 16a)*

*biRNN graph-based
dependency
(Kiperwaser+Goldberg 16;
Wang+Chang 16)*

biRNN span-based
constituency
*(Cross+Huang 16b)*

*minimal span-based
constituency
(Stern+ ACL 17)*

minimal dependency
*(Shi+ EMNLP 17)*

*all tree info
(summarize output **y**)*

*minimal or no tree info
(summarize input **x**)*

DP impossible     enables slow DP         enables fast DP     fastest DP: $O(n^3)$

# Spectrum: Neural Incremental Parsing

constituency
dependency

*bottom-up*

*edge-factored*
*(McDonald+ 05a)*

DP incremental parsing
*(Huang+Sagae 10, Kuhlmann+ 11)*

Feedforward NNs
*(Chen + Manning 14)*

Stack LSTM

*(Dyer+ 15)*

RNNG
*(Dyer+ 16)*

biRNN dependency
*(Kiperwaser+Goldberg 16;
Cross+Huang 16a)*

biRNN graph-based
dependency
*(Kiperwaser+Goldberg 16;
Wang+Chang 16)*

biRNN span-based
constituency
*(Cross+Huang 16b)*

*minimal span-based
constituency*
*(Stern+ ACL 17)*

minimal dependency
*(Shi+ EMNLP 17)*

*all tree info
(summarize output **y**)*

*minimal or no tree info
(summarize input **x**)*

DP impossible          enables slow DP                    enables fast DP          fastest DP: $O(n^3)$ 40

# Conclusions and Limitations

# Conclusions and Limitations

- DP and RNNs can indeed be married, if done creatively

  - biRNN summarizing input **x** and not output structure **y**

  - this allows efficient DP with exact search

  - combine with global learning (loss-augmented decoding)

# Conclusions and Limitations

- DP and RNNs can indeed be married, if done creatively

  - biRNN summarizing input **x** and not output structure **y**

  - this allows efficient DP with exact search

  - combine with global learning (loss-augmented decoding)

- but exact DP is still too slow

  - future work: linear-time beam search DP with biRNNs

# Conclusions and Limitations

- DP and RNNs can indeed be married, if done creatively

  - biRNN summarizing input **x** and not output structure **y**

  - this allows efficient DP with exact search

  - combine with global learning (loss-augmented decoding)

- but exact DP is still too slow

  - future work: linear-time beam search DP with biRNNs

- what if we want strictly incremental parsing? no biRNN...

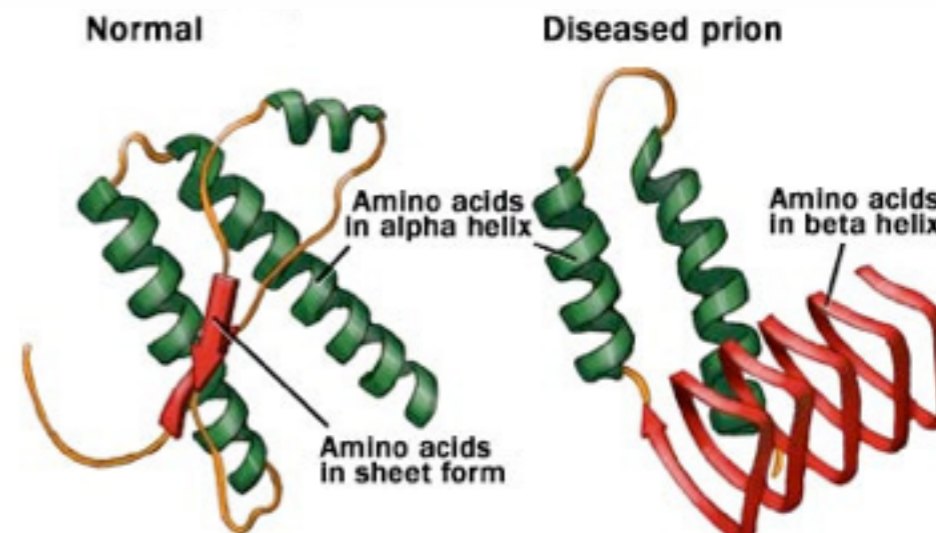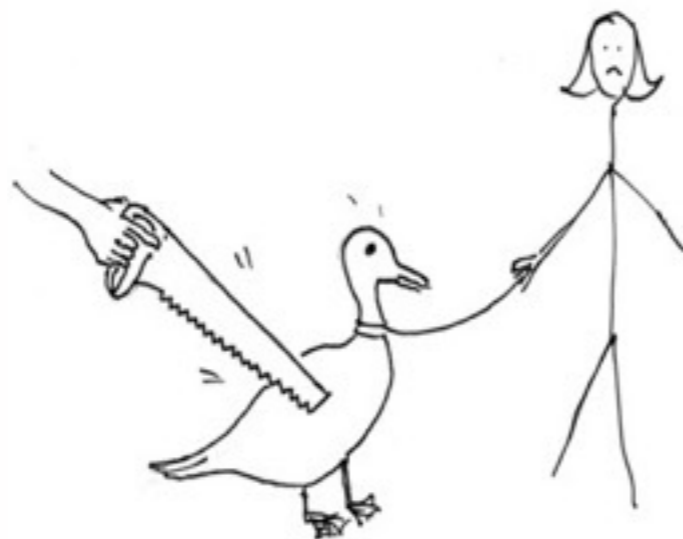  - DP search could compensate for loss of lookahead

# Conclusions and Limitations

- DP and RNNs can indeed be married, if done creatively
  - biRNN summarizing input **x** and not output structure **y**
  - this allows efficient DP with exact search
  - combine with global learning (loss-augmented decoding)
- but exact DP is still too slow
  - future work: linear-time beam search DP with biRNNs
- what if we want strictly incremental parsing? no biRNN...
  - DP search could compensate for loss of lookahead
- what about translation? we do need to model **y** directly...

# 非常 感谢！
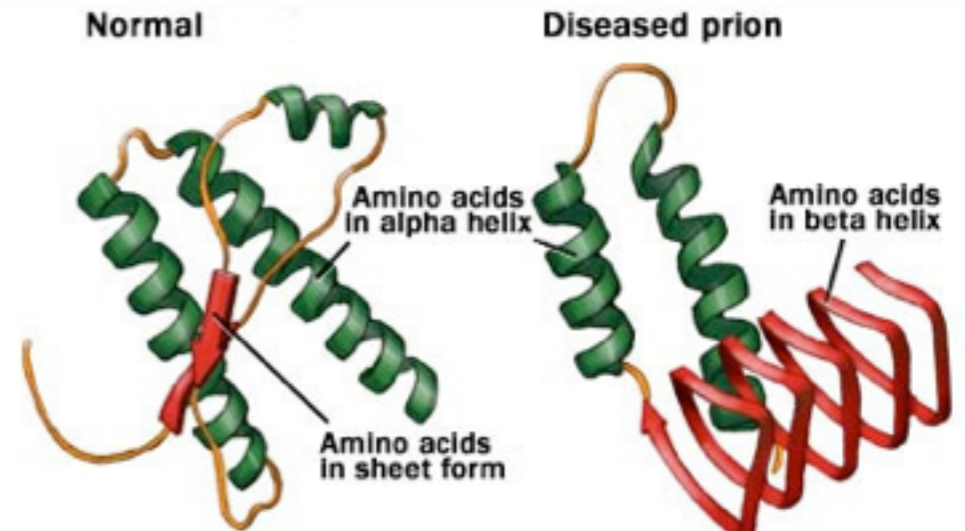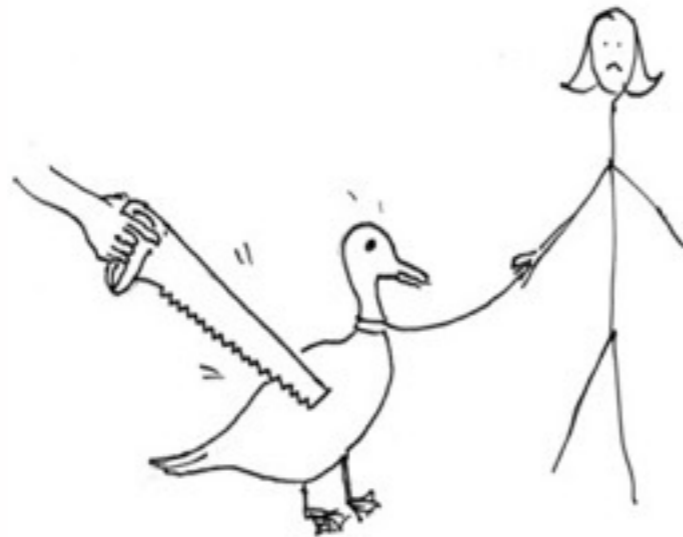
fēi cháng　gǎn xiè

James Cross