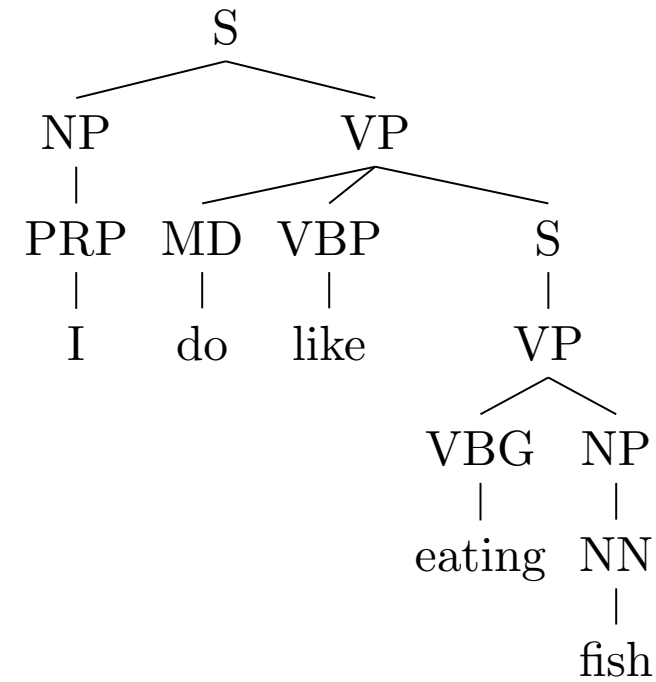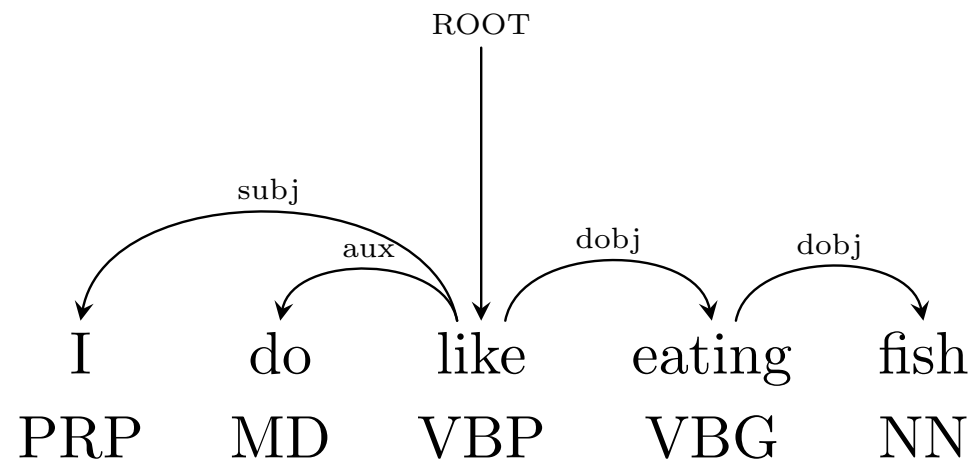# Span-Based Constituency Parsing with Provably Optimal Dynamic Oracles

James Cross and Liang Huang
Oregon State University

EMNLP, Austin, TX
November 2, 2016

**Oregon State**
UNIVERSITY
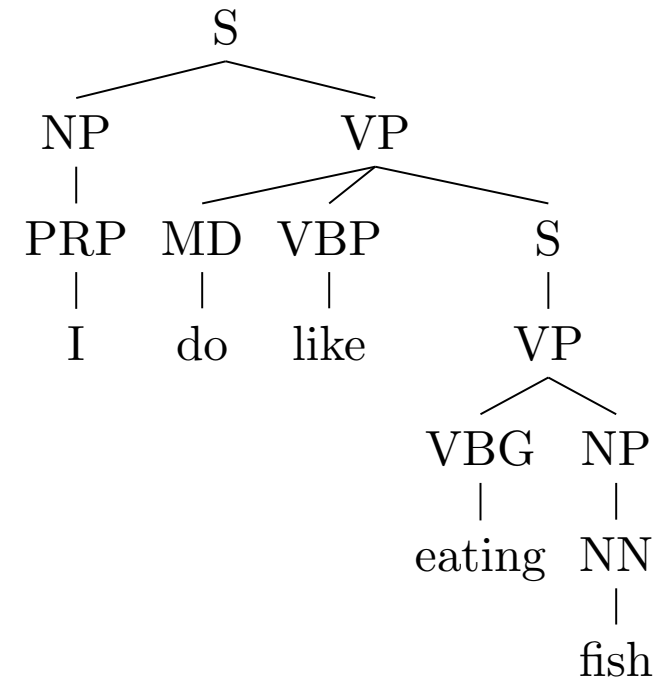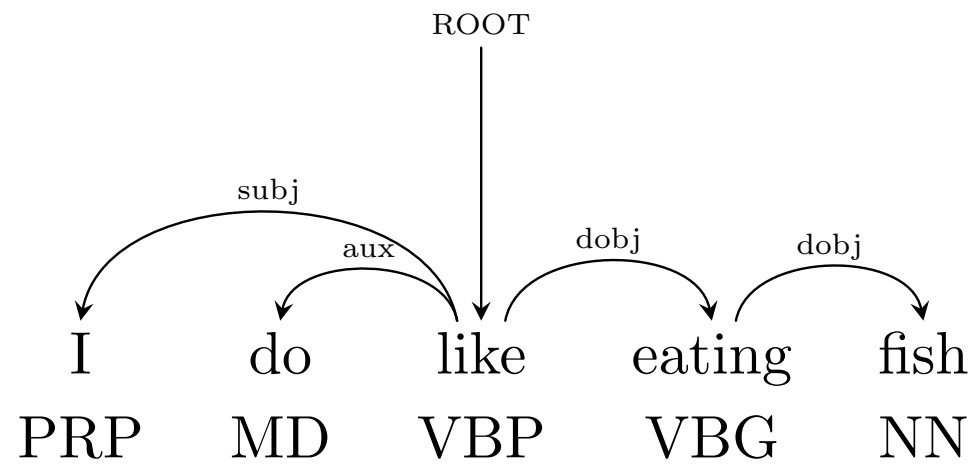
# Dependency vs. Constituency



| | search | UAS |
|---|---|---|
| **Zhang & Nivre 2011** | beam | 92.9 |
| **Chen & Manning 2014** | **greedy** | 91.8 |
| **Zhou et al. (2015)** | beam | 93.3 |
| **Weiss et al. (2015)** | beam | 94.0 |
| **our work (ACL 2016)** | **greedy** | 93.4 |
| **Andor et al. (2016)** | beam | **94.4** |

| | search | F$_1$ |
|---|---|---|
| **Carreras et al. (2008)** | cubic | 91.1 |
| **Shindo et al. (2012)** | cubic | 91.1 |
| **Thang et al. (2015) (A*)** | ~cubic | 91.1 |
| **Watanabe et al. (2015)** | beam | 90.7 |
| **Vinyals et al. (2015) (WSJ)** | beam | 90.5 |

**Red = Neural**

# Dependency vs. Constituency



| | search | UAS |
|---|---|---|
| **Zhang & Nivre 2011** | beam | 92.9 |
| **Chen & Manning 2014** | **greedy** | 91.8 |
| **Zhou et al. (2015)** | beam | 93.3 |
| **Weiss et al. (2015)** | beam | 94.0 |
| **our work (ACL 2016)** | **greedy** | 93.4 |
| **Andor et al. (2016)** | beam | **94.4** |

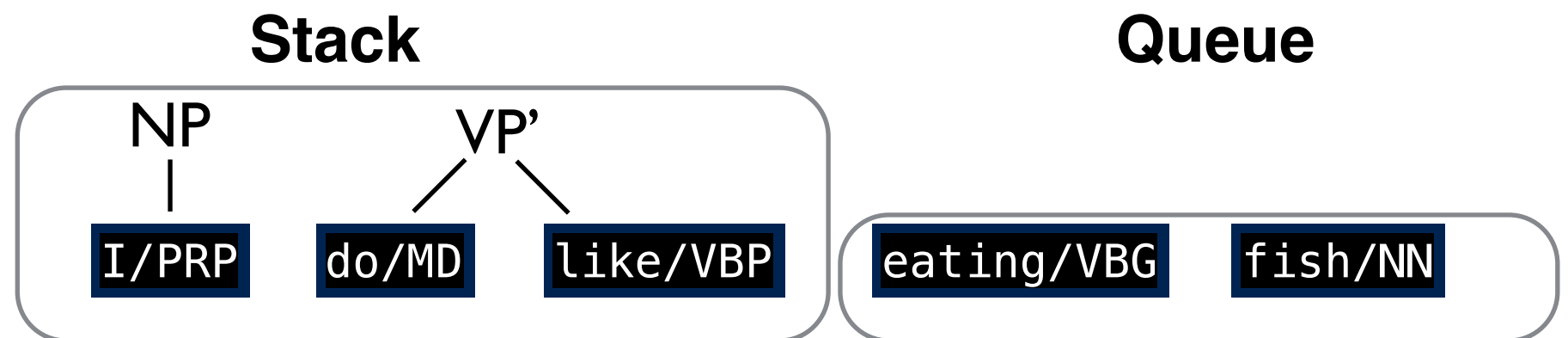| | search | $F_1$ |
|---|---|---|
| **Carreras et al. (2008)** | cubic | 91.1 |
| **Shindo et al. (2012)** | cubic | 91.1 |
| **Thang et al. (2015) (A\*)** | ~cubic | 91.1 |
| **Watanabe et al. (2015)** | beam | 90.7 |
| **Vinyals et al. (2015) (WSJ)** | beam | 90.5 |
| **This Work** | **greedy** | **91.3** |

**Red = Neural**

# Outline

- **Span-Based Constituency Parsing**

- Bi-Directional LSTM Span Features

- Provably Optimal Dynamic Oracle

- Experiments

# Span-Based Parsing

- Previous work uses tree structures on stack

- We simplify to operate directly on sentence spans

- Simple-to-implement linear-time parsing

**Stack**          **Queue**

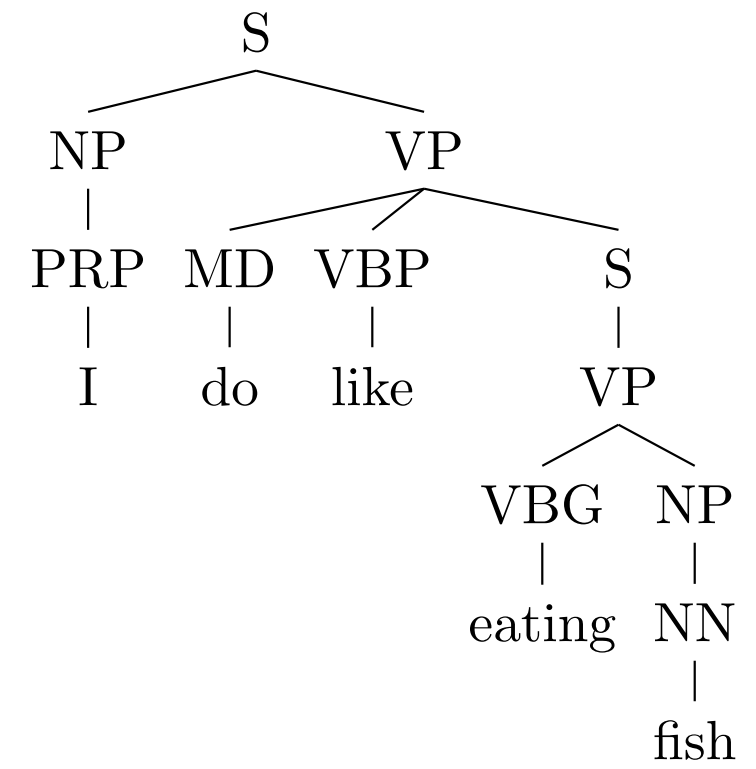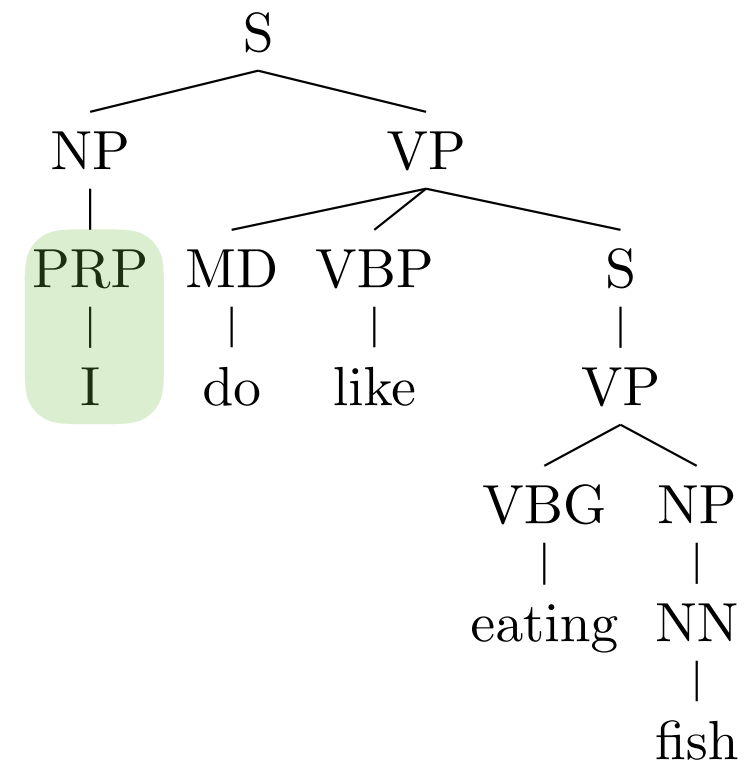*previous work*

NP          VP'

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

**Stack**          **Queue**

*our work*

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

4

| | | |
|---|---|---|
| **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 = {}

5

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

current brackets    t = {}

*Shift*

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       2         3           4        5

current brackets   $t = \{\}$

*Shift*

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

0          1       2         3           4        5

*Label-NP*   $t = \{_0NP_1\}$

| 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 | |
|---|---|---|---|---|---|---|
| 0 | | 1 | 2 | 3 | 4 | 5 |

current brackets    $t = \{\}$

*Shift*

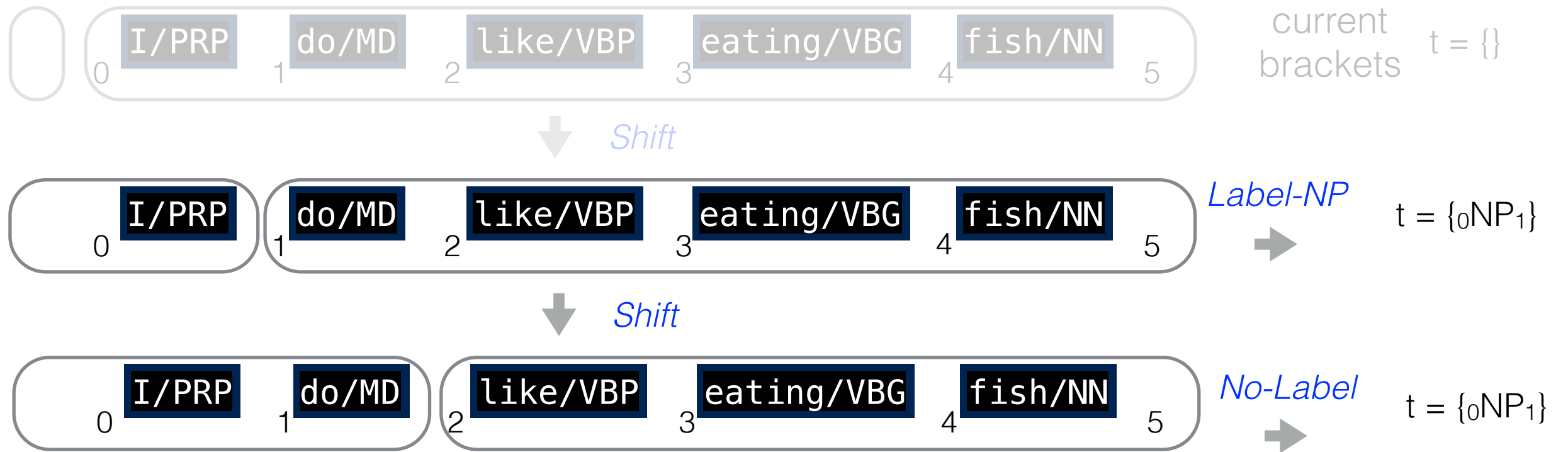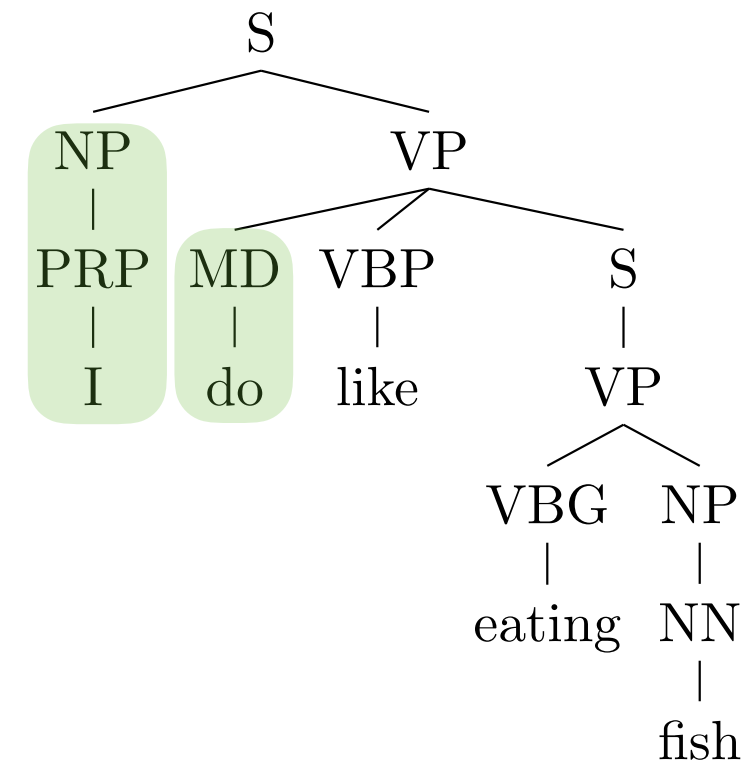| I/PRP | do/MD | like/VBP | eating/VBG | fish/NN |
|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 |

*Label-NP*    $t = \{_0NP_1\}$

*Shift*

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

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 | 2 | 3 | 4 | 5 |

current brackets    $t = \{\}$

*Shift*

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

*Label-NP*    $t = \{_0NP_1\}$

*Shift*

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

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

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

current brackets

$t = \{\}$

*Shift*

*Label-NP*

$t = \{_0\text{NP}_1\}$

*Shift*

*No-Label*

$t = \{_0\text{NP}_1\}$

*Shift*

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        2        3            4        5

current brackets    $t = \{\}$

*Shift*

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

0        1        2        3            4        5

*Label-NP*    $t = \{_0NP_1\}$

*Shift*

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

0        1        2        3            4        5

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

*Shift*

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

0        1        2        3            4        5

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

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

$t = \{_0 \text{NP}_1\}$

6

| 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

6

| | |
|---|---|
| **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

$t = \{_0 NP_1\}$

*Combine*

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

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

| **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

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

6

| 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

| 0 I/PRP | 1 do/MD | 2 like/VBP | | 3 eating/VBG | 4 fish/NN | 5 |  $t = \{_0\text{NP}_1\}$

*Combine*

| 0 I/PRP | 1 do/MD   like/VBP | | 3 eating/VBG | 4 fish/NN | 5 |  *No-Label*  $t = \{_0\text{NP}_1\}$

*Shift*

| 0 I/PRP | 1 do/MD   like/VBP | 3 eating/VBG | | 4 fish/NN | 5 |  *No-Label*  $t = \{_0\text{NP}_1\}$

6

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

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

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

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

*Shift*

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

6

| 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
$t = \{_0NP_1\}$

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

*Shift*

I/PRP   do/MD   like/VBP   eating/VBG   fish/NN
0   1   3   4   5
*Label-NP*   $t = \{_0NP_1, _4NP_5\}$

6

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

S
NP — VP
PRP: I
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

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

7

| | |
|---|---|
| **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
```

$t = \{_0NP_1, _4NP_5\}$

```
┌─────────────────────────────────────────────────────────────┐  ┐
│  │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     │  │
└─────────────────────────────────────────────────────────────┘  ┘
```

7

| 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*

*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

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

*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

7

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

| | | |
|---|---|---|
| **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* ➡ $t = \{_0NP_1, _4NP_5, _3S_5, _3VP_5\}$

⬇ *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 — 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* → $t = \{_0NP_1, \, _4NP_5, \\ _3S_5, \, _3VP_5\}$

⬇ *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

*Label-S* → $t = \{_0NP_1, \, _4NP_5, \\ _3S_5, \, _3VP_5, \\ _1VP_5, \, _0S_5\}$

7

# Advantages of Span-Based System

- Linear-time and fixed number of steps (well-suited for beam search)

- Separates prediction of structure and labels

- Predicts rules of arbitrary arity with no binarization

**Stack**                    **Queue**

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

0       1                        3            4           5

# Advantages of Span-Based System

- Linear-time and fixed number of steps (well-suited for beam search)

- Separates prediction of structure and labels

- Predicts rules of arbitrary arity with no binarization

**Stack**                                    **Queue**

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

Q: How to decide which action to take?
What features represent spans?

# Outline

- Span-Based Constituency Parsing

- **Bi-Directional LSTM Span Features**

- Provably Optimal Dynamic Oracle

- Experiments

# 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: $b_3$ - $b_5$

# Span Features for Structure Action

to predict:
*Combine*



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

*pre-s₁*            *s₁*                    *s₀*                  *queue*

4 bi-LSTM span features

**(no tree-structure information used)**

# Span Features for Label Action

to predict:
*Label-VP*

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

*pre-s₀*                            *s₀*                         *queue*

3 bi-LSTM span features

**(no tree-structure information used)**

# Training Scheme: Local

- Every parser state is paired with a correct action

- Separate multilayer perceptron for each action type

- Baseline training scheme (static oracle) uses canonical order with short-stack preference

# Training Scheme: Local

- Every parser state is paired with a correct action

- Separate multilayer perceptron for each action type

- Baseline training scheme (static oracle) uses canonical order with short-stack preference

# Outline

- Span-Based Constituency Parsing

- Bi-Directional LSTM Span Features

- **Provably Optimal Dynamic Oracle**

- Experiments

# Dynamic Oracle: Motivation

- Static oracle training assumes all correct actions

- What to do after decoding mistakes?



*gold path*

# Dynamic Oracle: Motivation

- Static oracle training assumes all correct actions

- What to do after decoding mistakes?

- Need a way to decide best action in arbitrary state: **Dynamic Oracle** (everywhere-defined optimal policy)

*gold path*

*best action?*
**Dynamic Oracle!**

# Dynamic Oracle: Example

# Dynamic Oracle: Example

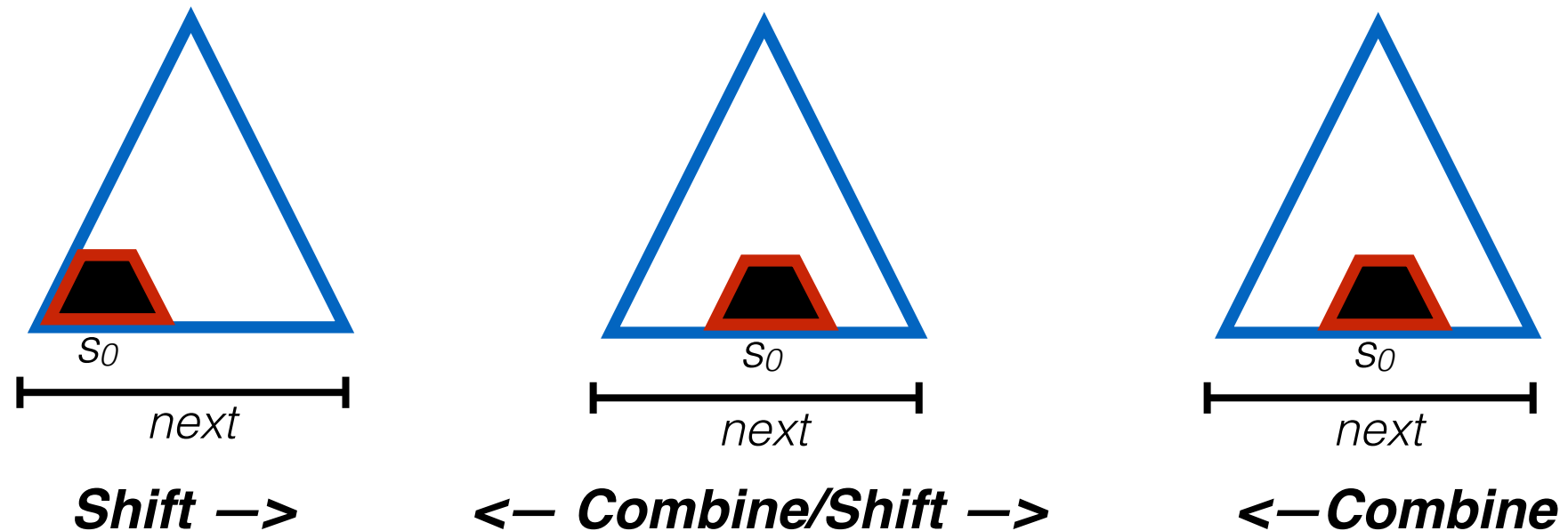# Dynamic Oracle: Example

S

NP     VP     .

PRP   MD   VBP     S

VP

VBG   NP

NN

I   do   like   eating   fish   .

$s_1$      $s_0$

smallest reachable gold bracket incl. $s_0$

16

# Dynamic Oracle: Example

S

NP · · · VP · · · .

PRP · · · MD · · VBP · · · S

VP

VBG · · · NP

NN

I · do · like · eating · fish · .

*S₁* · · · *S₀*

smallest reachable
gold bracket incl. *s₀*

# Dynamic Oracle: Example

S

NP VP .

PRP MD VBP S

VP

VBG NP

NN

I do like eating fish .

$s_1$ $s_0$

smallest reachable
gold bracket incl. $s_0$

$s_0$

next reachable

# Dynamic Oracle: Example

S

NP       VP       .

PRP    MD    VBP     S

VP

VBG    NP

NN

I     do    like    eating     fish    .

$S_1$        $S_0$

smallest reachable
gold bracket incl. $s_0$



$s_0$

*next reachable*

**Dynamic Oracle:
Shift or Combine**

16

# Dynamic Oracle: Example

# Dynamic Oracle: Example

S

NP

PRP

VP

MD  VBP

S

VP

VBG  NP

NN

I  do  like  eating  fish  .

$s_1$  $s_0$

smallest reachable
gold bracket incl. $s_0$



$s_0$

*next reachable*

# Dynamic Oracle: Example

S

NP       VP       .

PRP   MD   VBP      S

VP

VBG    NP

NN

I     do    like   eating    fish   .

$s_1$          $s_0$

smallest reachable
gold bracket incl. $s_0$



$s_0$

*next reachable*

**Dynamic Oracle:
Combine**

# Dynamic Oracle: Example



S

NP VP .

PRP MD VBP S

VP

VBG NP

NN

I do like eating fish .

$s_1$ $s_0$

# Dynamic Oracle: Example

# Dynamic Oracle: Example



smallest reachable
gold bracket incl. $s_0$

# Dynamic Oracle: Example



smallest reachable
gold bracket incl. $s_0$

**Dynamic Oracle:
Shift**

# Dynamic Oracle: Full Definition



**Shift —>**  **<— Combine/Shift —>**  **<—Combine**
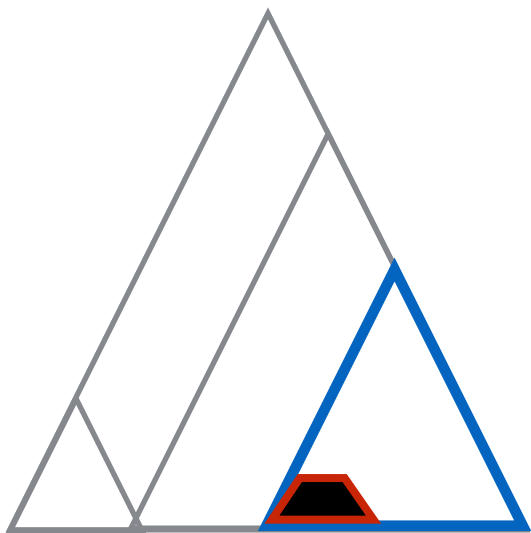
- Structure actions depend on next <u>reachable</u> bracket in gold tree
- All non-bracket label states —> *No-Label*
- All gold-bracket label states —> Correct label(s)

# Dynamic Oracle: Full Definition



**Shift —>**    **<— Combine/Shift —>**    **<—Combine**

- Structure actions depend on next <u>reachable</u> bracket in gold tree
- All non-bracket label states —> *No-Label*
- All gold-bracket label states —> Correct label(s)
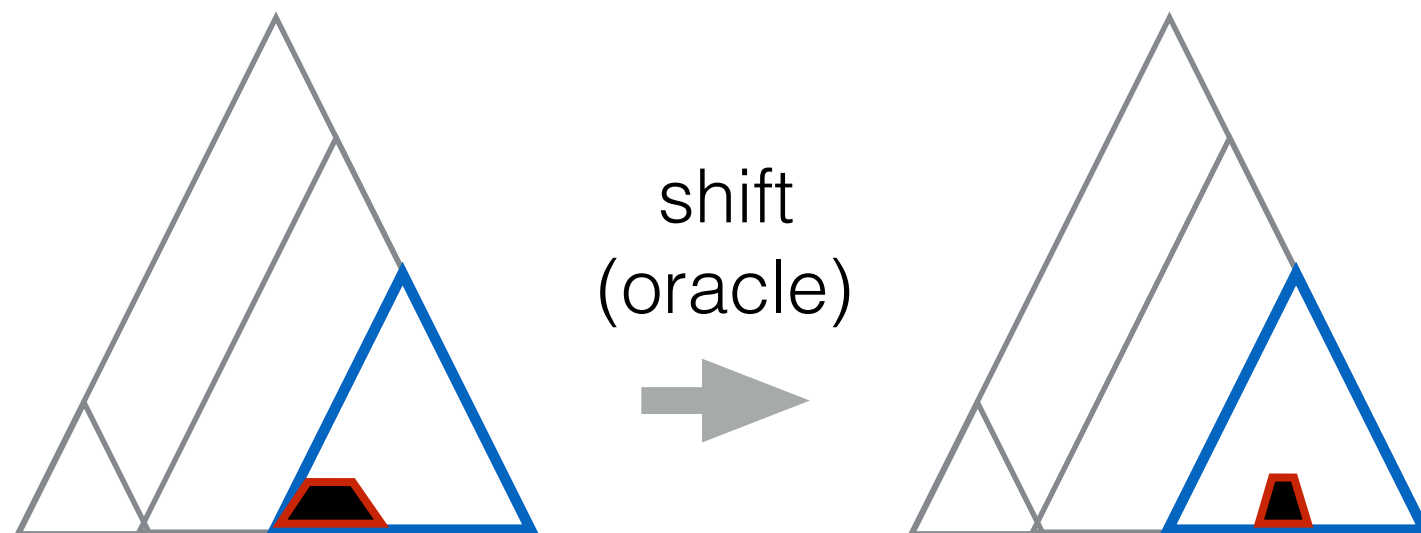
Gold Brackets

Current Brackets

# Dynamic Oracle: Full Definition



**Shift —>**　　　**<— Combine/Shift —>**　　　**<—Combine**

- Structure actions depend on next <u>reachable</u> bracket in gold tree
- All non-bracket label states —> *No-Label*
- All gold-bracket label states —> Correct label(s)



Gold Brackets

Current Brackets

Reachable

# Dynamic Oracle: Full Definition



**Shift —>**   **<— Combine/Shift —>**   **<—Combine**

- Structure actions depend on next <u>reachable</u> bracket in gold tree
- All non-bracket label states —> *No-Label*
- All gold-bracket label states —> Correct label(s)



Gold Brackets
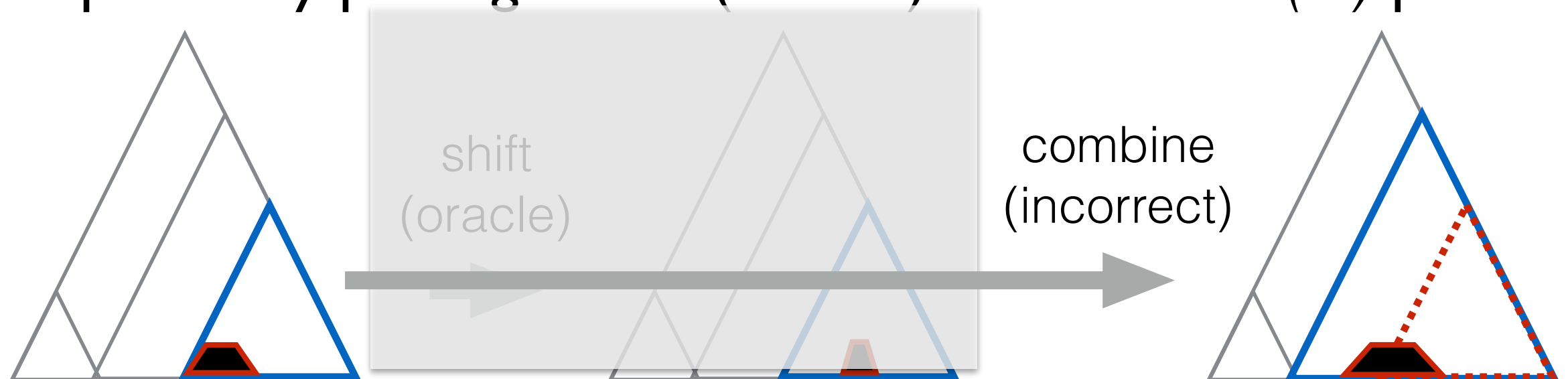
Current Brackets
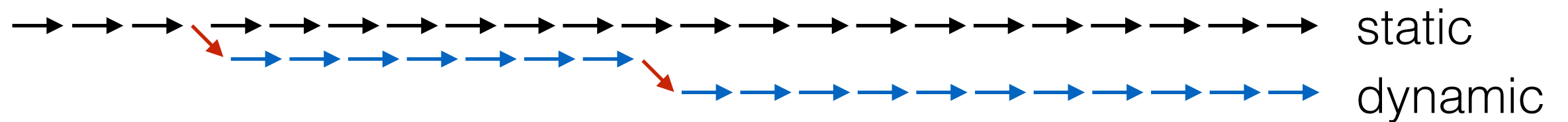
**optimal tree**

Reachable

# Dynamic Oracle: Optimality/Complexity

- First provably optimal oracle for constituency parsing (optimal in both precision and recall)

- After each action next reachable may (or may not) be updated by tracing parent link in gold tree

- Also $O(n)$ steps, thus amortized $O(1)$ time

- Dependency parsing oracle (arc-std): worst case $O(n^3)$ per step

# Dynamic Oracle: Optimality/Complexity

- First provably optimal oracle for constituency parsing (optimal in both precision and recall)

- After each action next reachable may (or may not) be updated by tracing parent link in gold tree

- Also $O(n)$ steps, thus amortized $O(1)$ time

- Dependency parsing oracle (arc-std): worst case $O(n^3)$ per step
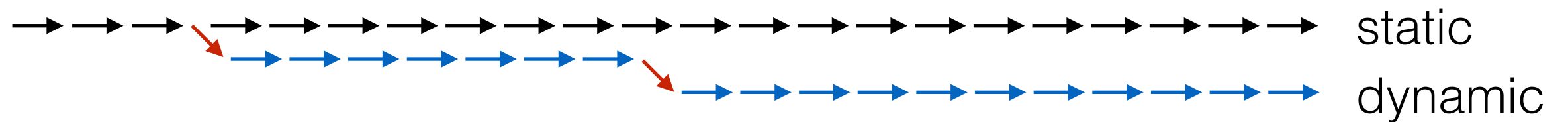
shift
(oracle)

# Dynamic Oracle: Optimality/Complexity

- First provably optimal oracle for constituency parsing (optimal in both precision and recall)

- After each action next reachable may (or may not) be updated by tracing parent link in gold tree

- Also $O(n)$ steps, thus amortized $O(1)$ time

- Dependency parsing oracle (arc-std): worst case $O(n^3)$ per step



shift
(oracle)

combine
(incorrect)

# Training with Dynamic Oracle

$$\rightarrow \rightarrow \rightarrow \ \ \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \text{ static}$$

| (scores on PTB 22) | Recall | Prec. | $F_1$ |
|---|---|---|---|
| Static Oracle | 91.34 | 91.43 | 91.38 |

# Training with Dynamic Oracle

- Basic dynamic oracle: follow current model

static

dynamic

| (scores on PTB 22) | Recall | Prec. | $F_1$ |
|---|---|---|---|
| Static Oracle | 91.34 | 91.43 | 91.38 |

# Training with Dynamic Oracle

- Basic dynamic oracle: follow current model

static

dynamic

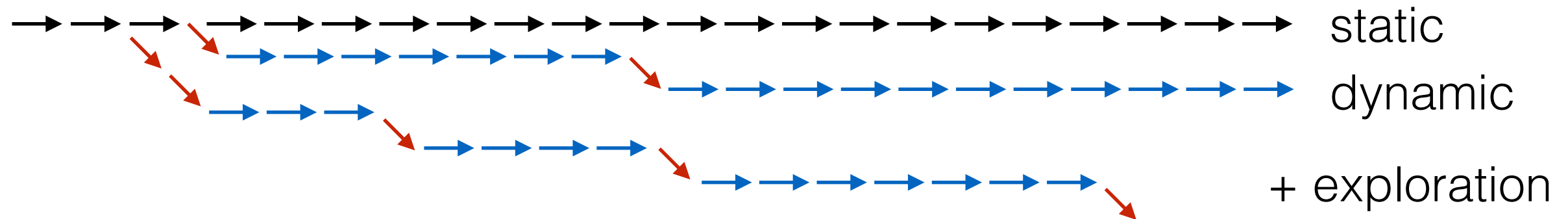| (scores on PTB 22) | Recall | Prec. | $F_1$ |
|---|---|---|---|
| Static Oracle | 91.34 | 91.43 | 91.38 |
| Dynamic Oracle | 91.14 | 91.61 | 91.38 |

# Training with Dynamic Oracle

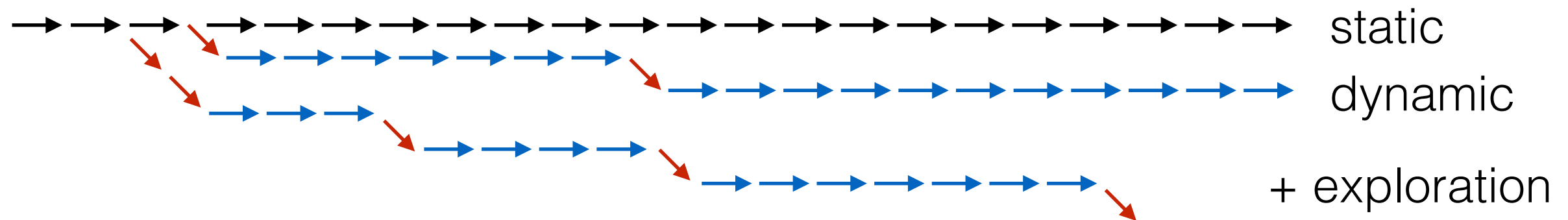- Basic dynamic oracle: follow current model

- Problem: overfits training data, making fewer mistakes than test



| (scores on PTB 22) | Recall | Prec. | $F_1$ |
|---|---|---|---|
| Static Oracle | 91.34 | 91.43 | 91.38 |
| Dynamic Oracle | 91.14 | 91.61 | 91.38 |

# Training with Dynamic Oracle

- Basic dynamic oracle: follow current model

- Problem: overfits training data, making fewer mistakes than test

- Exploration: sample from softmax distribution (Ballesteros et al., 2016) to encourage more mistakes
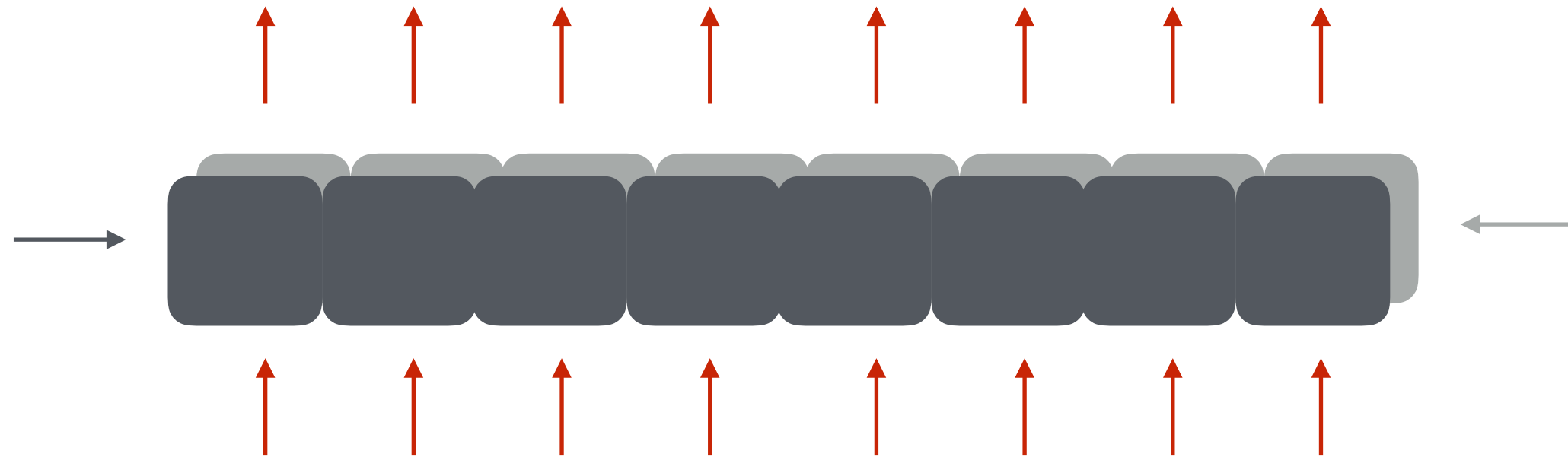


| (scores on PTB 22) | Recall | Prec. | $F_1$ |
|---|---|---|---|
| Static Oracle | 91.34 | 91.43 | 91.38 |
| Dynamic Oracle | 91.14 | 91.61 | 91.38 |

# Training with Dynamic Oracle

- Basic dynamic oracle: follow current model

- Problem: overfits training data, making fewer mistakes than test

- Exploration: sample from softmax distribution (Ballesteros et al., 2016) to encourage more mistakes



| (scores on PTB 22) | Recall | Prec. | $F_1$ |
|---|---|---|---|
| Static Oracle | 91.34 | 91.43 | 91.38 |
| Dynamic Oracle | 91.14 | 91.61 | 91.38 |
| Dynamic + Exploration | 91.07 | 92.22 | **91.64** |

# Outline

- Span-Based Constituency Parsing

- Bi-Directional LSTM Span Features

- Provably Optimal Dynamic Oracle
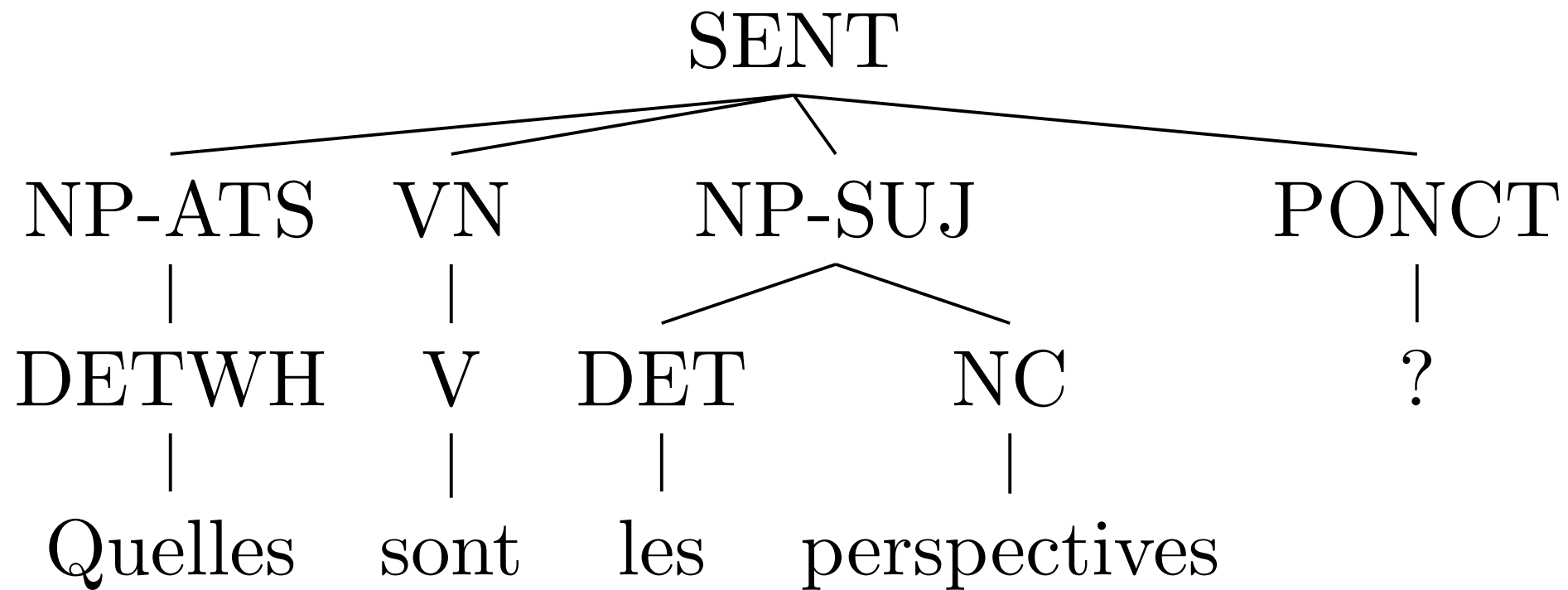
- **Experiments**

# Architecture



- 50-dim word and 20-dim tag embeddings

- No pre-training

- Each LSTM layer 200 units each direction

- 200 ReLU units for each of structure and label predictors

# 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

# Parsing Morphologically Rich Languages



lemma = perspective
coarse_POS = N
gender = feminine
number = plural
subcategory = common

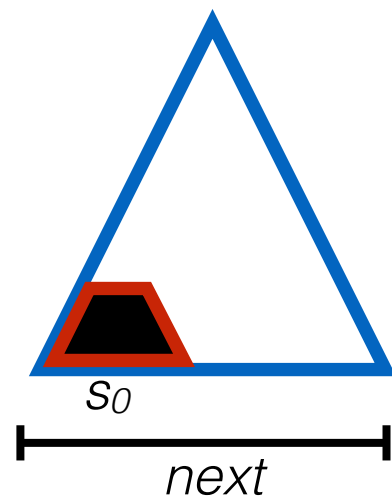# Results on French Treebank

- Morphological feature embeddings (10 dim. each)

- Additional input to recurrent network

- For French, we used SPMRL 2014 predicted features

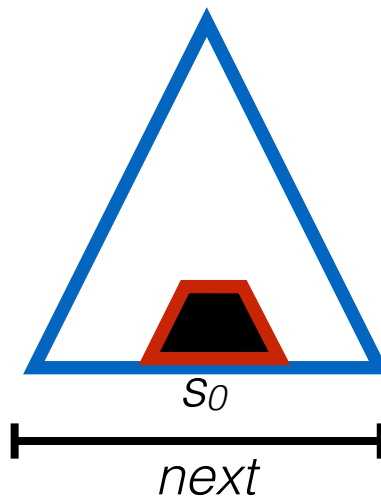| Parser | Recall | Prec. | F$_1$ |
|---|---|---|---|
| Björkelund et al. (2014) | | | 82.53 |
| **Static Oracle** | 83.50 | 82.87 | 83.18 |
| **Dynamic + Exploration** | 81.90 | 84.77 | **83.31** |

# Summary

- Simple, easy-to-implement span-based parsing system

- No tree/label information in features (good candidate for dynamic programming)

- Linear time parsing with greedy decoding

- No pre-trained embeddings, small architecture, and minimal hyper-parameter tuning (trained on CPU)

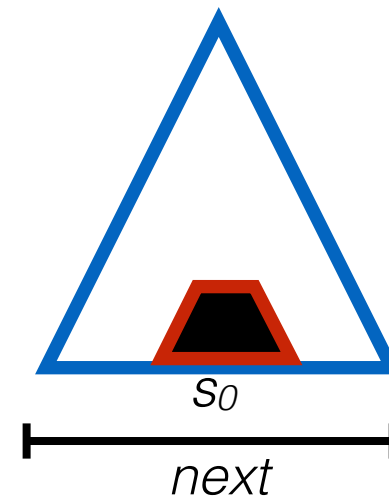- First optimal dynamic oracle for constituency parsing

# Thank You!



**Shift —>**          **<— Combine/Shift —>**          **<—Combine**

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

*pre-$s_1$*          *$s_1$*          *$s_0$*          *queue*

Oregon State
UNIVERSITY