# Fast(er) Exact Decoding and Global Training for Transition-Based Dependency Parsing via a Minimal Feature Set

Tianze Shi*     Liang Huang†     Lillian Lee*

* Cornell University          † Oregon State University

$O(n^3)$

Minimal

$O(n^3)$

~~$O(n^6)$~~

Feature Set

Theoretical

Practical

# Short Version

- Transition-based dependency parsing has an exponentially-large search space

- $O(n^3)$ exact solutions exist ☺

- In practice, however, we needed rich features $\implies O(n^6)$ ☹

- (This work) with bi-LSTMs, now we can do $O(n^3)$! ☺

- And we get state-of-the-art results

# Short Version

- Transition-based dependency parsing has an exponentially-large search space

- $O(n^3)$ exact solutions exist 😃

- In practice, however, we needed rich features $\Longrightarrow O(n^6)$ 😖

- (This work) with bi-LSTMs, now we can do $O(n^3)$! 😃

- And we get state-of-the-art results

3

# Short Version

- Transition-based dependency parsing has an exponentially-large search space

- $O(n^3)$ exact solutions exist 😃

- In practice, however, we needed rich features $\implies O(n^6)$ 😣

- (This work) with bi-LSTMs, now we can do $O(n^3)$! 😃

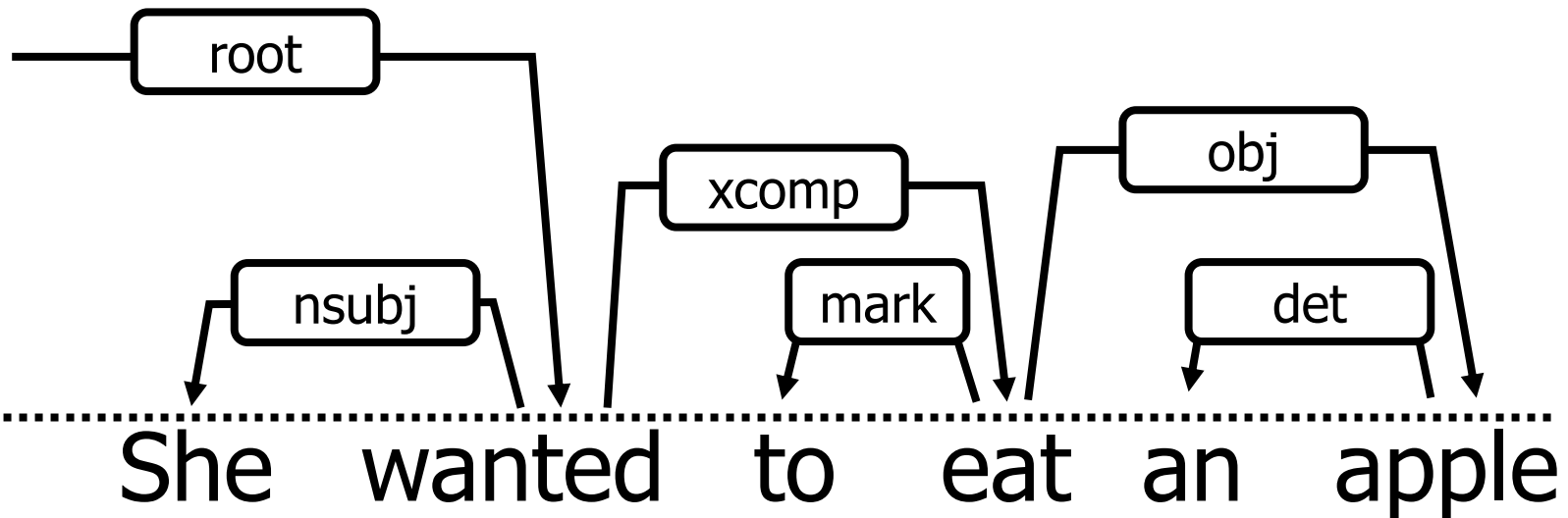- And we get state-of-the-art results

# Short Version

- Transition-based dependency parsing has an exponentially-large search space

- $O(n^3)$ exact solutions exist 😃

- In practice, however, we needed rich features $\implies O(n^6)$ 😣

- (This work) with bi-LSTMs, now we can do $O(n^3)$! 😃

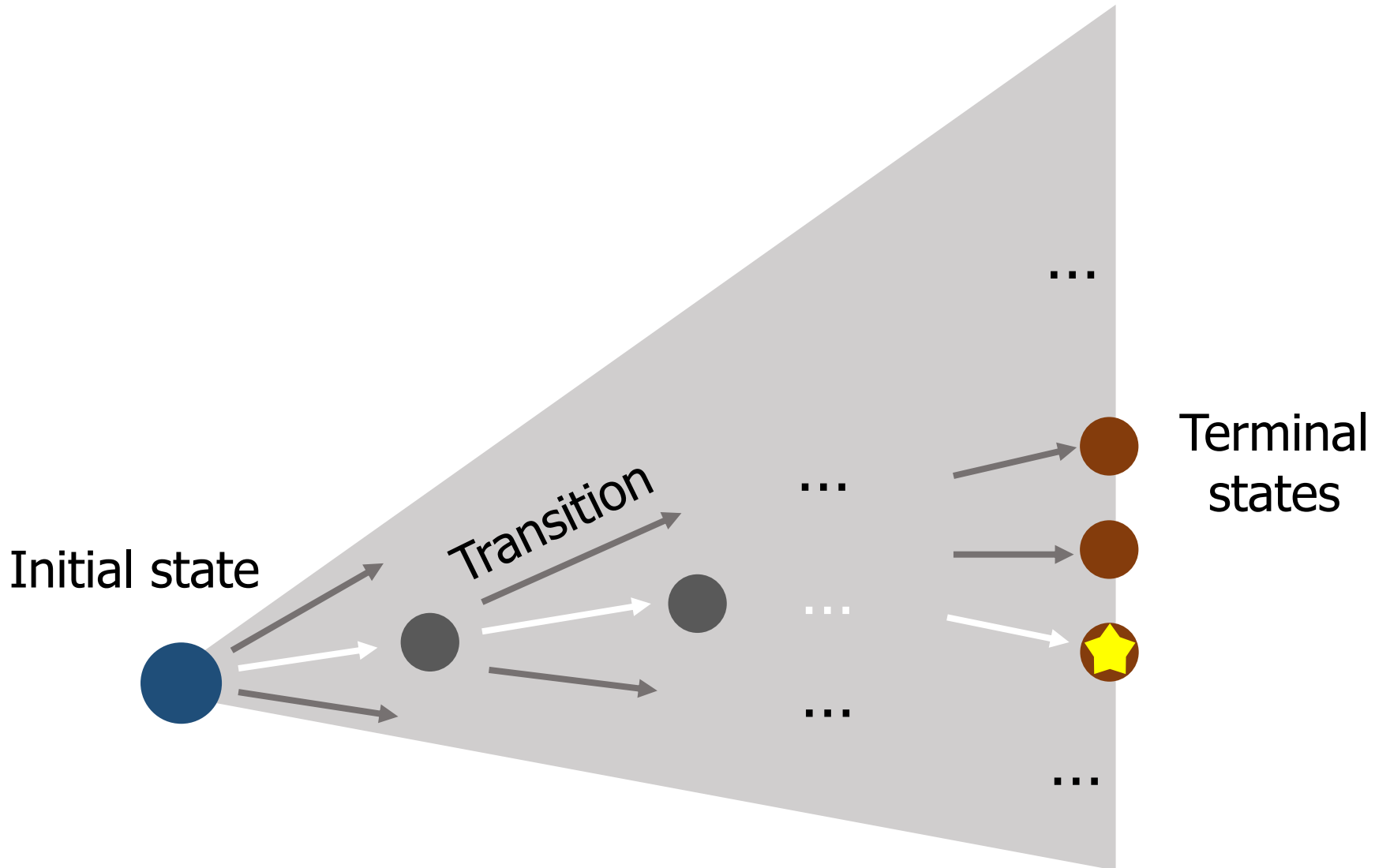- And we get state-of-the-art results

# Short Version

- Transition-based dependency parsing has an exponentially-large search space

- $O(n^3)$ exact solutions exist 😃

- In practice, however, we needed rich features $\implies O(n^6)$ 😣

- (This work) with bi-LSTMs, now we can do $O(n^3)$! 😃

- And we get state-of-the-art results

# Dependency Parsing



**OUTPUT**

**INPUT**    She    wanted    to    eat    an    apple

# Transition-based Dependency Parsing



Initial state

Transition

Terminal states

8

# Transition-based Dependency Parsing



*Goal:*

$$\max \text{ score}(\bullet \to \bullet \to \bullet \to \cdots \to \bullet)$$

$$= \max \sum \text{score}(\bullet \to \bullet)$$

Initial state

Transition

Terminal states
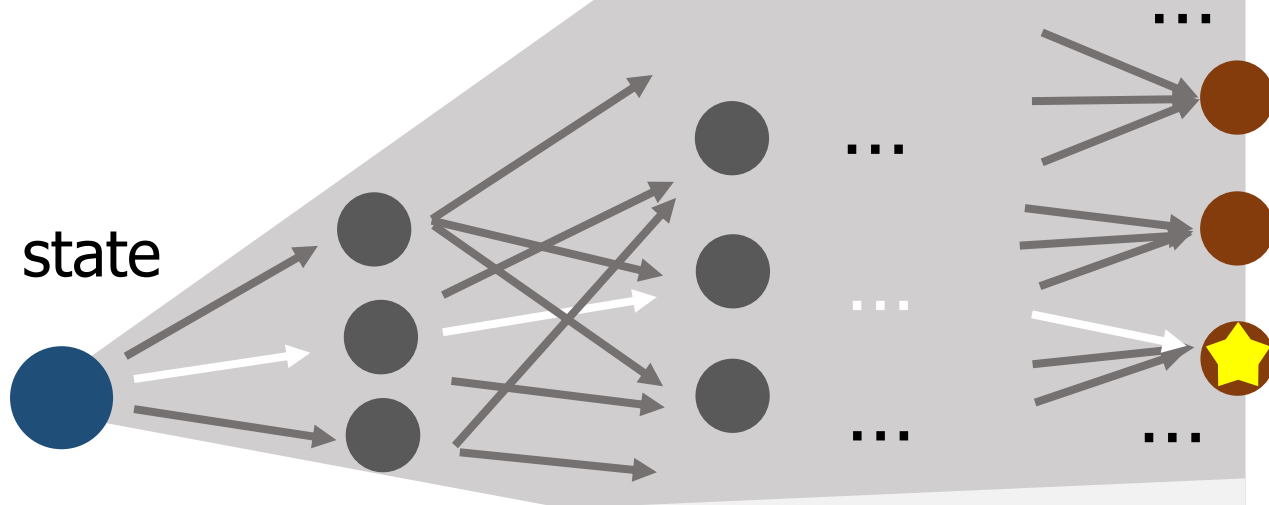
# Exact Decoding with Dynamic Programming

*Goal:*

$$\max \; \text{score}( \bullet \rightarrow \bullet \rightarrow \bullet \rightarrow \cdots \rightarrow \bullet )$$

$$= \max \; \sum \text{score}( \bullet \rightarrow \bullet )$$

Exponential to polynomial

Initial state

... Terminal states

(Huang and Sagae, 2010; Kuhlmann, Gómez-Rodríguez and Satta, 2011)
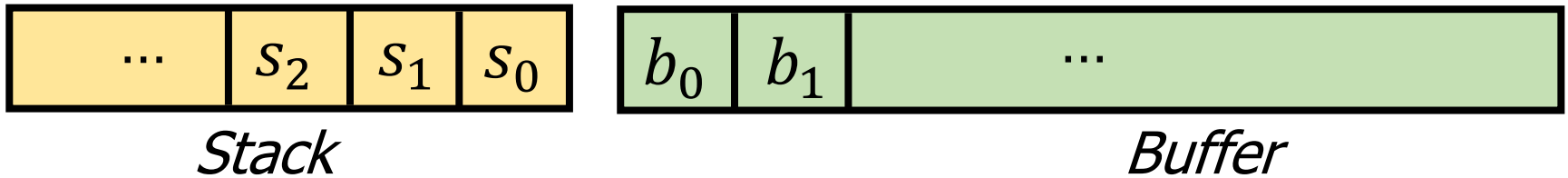
10

# Transition Systems

| | DP Complexity | # Action Types |
|---|---|---|
| Arc-standard | $O(n^4)$ | **3** |
| Arc-eager | $\boldsymbol{O(n^3)}$ | 4 |
| Arc-hybrid | $\boldsymbol{O(n^3)}$ | **3** |

In our paper

Presentational convenience

# Arc-hybrid Transition System

**State**

| ... | $s_2$ | $s_1$ | $s_0$ |
|---|---|---|---|

| $b_0$ | $b_1$ | ... |
|---|---|---|

*Stack*    *Buffer*

**Initial State**

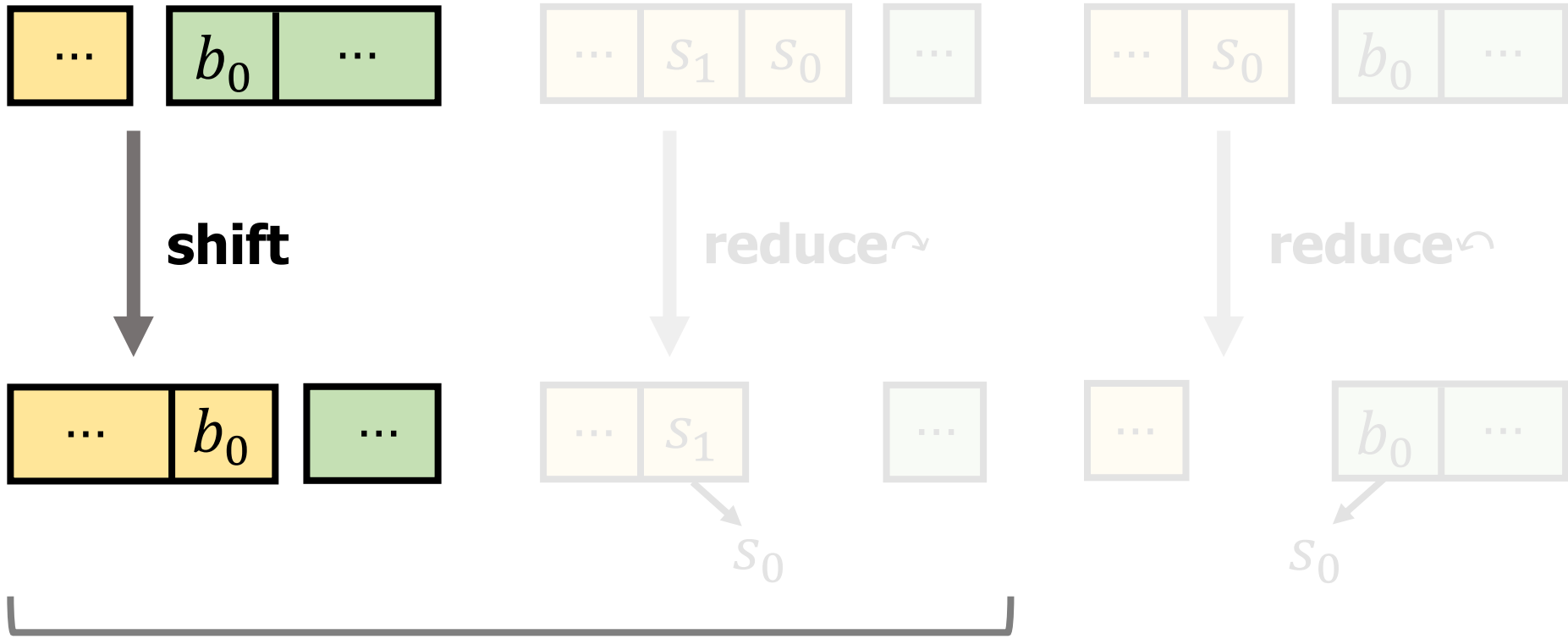| | ROOT | She | wanted | ... |
|---|---|---|---|---|

**Terminal State**

| ROOT | |
|---|---|

(Yamada and Matsumoto, 2003)
(Gómez-Rodríguez et al., 2008)
(Kuhlmann et al., 2011)

12

# Arc-hybrid Transition System

**Transitions**

shift

reduce

reduce

Same as arc-standard

# Arc-hybrid Transition System

**Transitions**

shift

reduce $\curvearrowright$

reduce $\curvearrowleft$

$s_0$

$b_0$

$s_0$
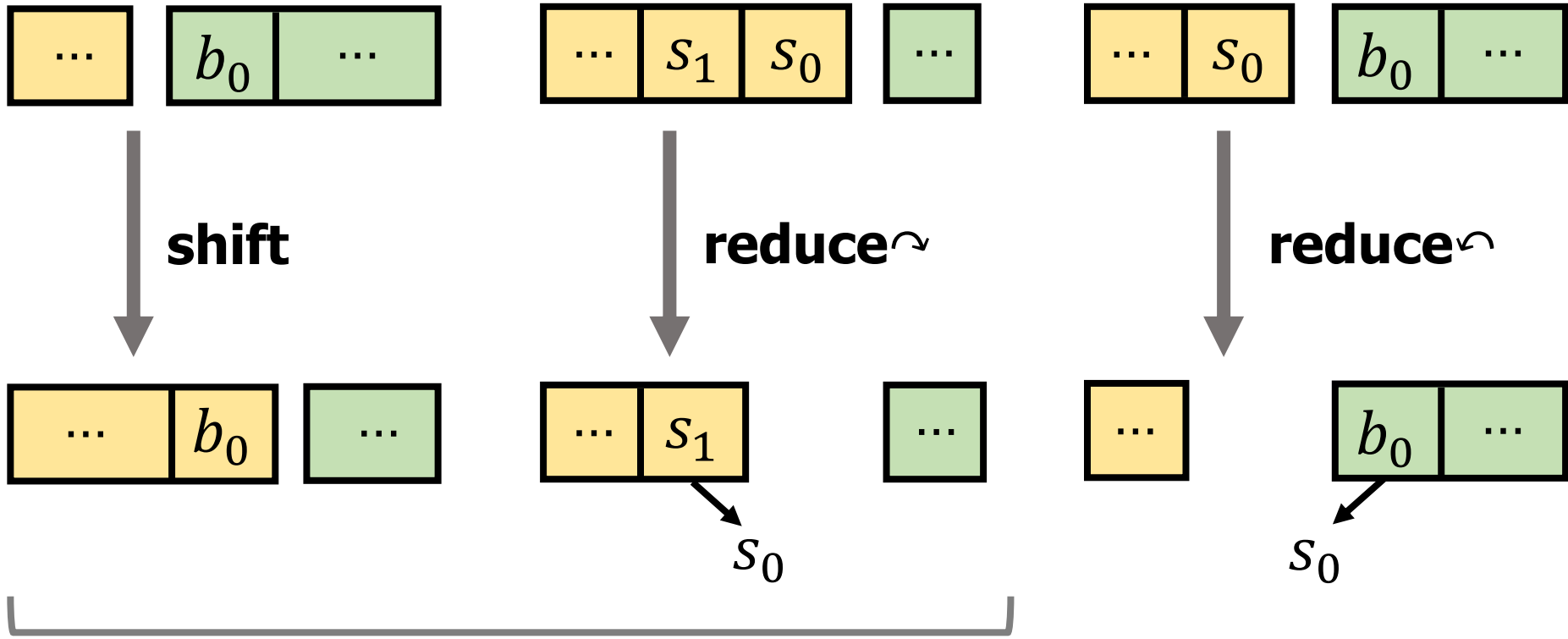
Same as arc-standard

# Arc-hybrid Transition System

**Transitions**



Same as arc-standard

# Arc-hybrid Transition System

| | Stack | Buffer |
|---|---|---|



*initial*   | | ROOT | She | wanted | to | eat | an | apple |

**shift**    ROOT | She | wanted | to | eat | an | apple

**shift**    ROOT | She | wanted | to | eat | an | apple

**reduce**    ROOT | wanted | to | eat | an | apple

She

**shift**    ROOT | wanted | to | eat | an | apple

**shift**    ROOT | wanted | to | eat | an | apple

16

# Arc-hybrid Transition System

| | *Stack* | | *Buffer* | | | | | |
|---|---|---|---|---|---|---|---|---|

*initial*

|  | ROOT | She | wanted | to | eat | an | apple |
|---|---|---|---|---|---|---|---|

**shift**

| ROOT | | She | wanted | to | eat | an | apple |
|---|---|---|---|---|---|---|---|

shift

| ROOT | She | | wanted | to | eat | an | apple |
|---|---|---|---|---|---|---|---|

reduce

| ROOT | | | wanted | to | eat | an | apple |
|---|---|---|---|---|---|---|---|

She

shift

| ROOT | wanted | | to | eat | an | apple |
|---|---|---|---|---|---|---|

shift

| ROOT | wanted | to | | eat | an | apple |
|---|---|---|---|---|---|---|

17

# Arc-hybrid Transition System

|            | Stack | Buffer |
|------------|-------|--------|

*initial*

**shift**

**shift**

reduce

shift

shift



18

# Arc-hybrid Transition System

*Stack*                                    *Buffer*

*initial*    [ ] | ROOT | She | wanted | to | eat | an | apple |

**shift**    | ROOT |     | She | wanted | to | eat | an | apple |

**shift**    | ROOT | She |     | wanted | to | eat | an | apple |

**reduce↶**  | ROOT |     | wanted | to | eat | an | apple |

                                    She

**shift**    | ROOT | wanted |     | to | eat | an | apple |

**shift**    | ROOT | wanted | to |     | eat | an | apple |

19

# Arc-hybrid Transition System

| | Stack | | Buffer |
|---|---|---|---|

*initial*

| | | ROOT | She | wanted | to | eat | an | apple |

**shift**

| ROOT | | | She | wanted | to | eat | an | apple |

**shift**

| ROOT | She | | | wanted | to | eat | an | apple |

**reduce** ↶

| ROOT | | | | wanted | to | eat | an | apple |

She

**shift**

| ROOT | wanted | | | | to | eat | an | apple |

**shift**

| ROOT | wanted | to | | | | eat | an | apple |

20

# Arc-hybrid Transition System

| | Stack | Buffer |
|---|---|---|

*initial*

| | |
|---|---|

| ROOT | She | wanted | to | eat | an | apple |
|---|---|---|---|---|---|---|

**shift**

| ROOT |
|---|

| She | wanted | to | eat | an | apple |
|---|---|---|---|---|---|

**shift**

| ROOT | She |
|---|---|

| wanted | to | eat | an | apple |
|---|---|---|---|---|

**reduce**↶

| ROOT |
|---|

| wanted | to | eat | an | apple |
|---|---|---|---|---|

She

**shift**

| ROOT | wanted |
|---|---|

| to | eat | an | apple |
|---|---|---|---|

**shift**

| ROOT | wanted | to |
|---|---|---|

| eat | an | apple |
|---|---|---|

21

# Arc-hybrid Transition System

*Stack*                                                    *Buffer*

| ROOT | wanted | to |

| eat | an | apple |

reduce

| ROOT | wanted |

| eat | an | apple |

shift

to

| ROOT | wanted | eat |

| an | apple |

shift

| ROOT | wanted | eat | an |

| apple |

reduce

| ROOT | wanted | eat |

| apple |

shift

an

| ROOT | wanted | eat | apple |

# Arc-hybrid Transition System

*Stack*                                    *Buffer*

| ROOT | wanted | to |

| eat | an | apple |

**reduce** ↶

| ROOT | wanted |

| eat | an | apple |

to

shift

| ROOT | wanted | eat |

| an | apple |

shift

| ROOT | wanted | eat | an |

| apple |

reduce ↶

| ROOT | wanted | eat |

| apple |

an

shift

| ROOT | wanted | eat | apple |

| |

23

# Arc-hybrid Transition System

*Stack*                                    *Buffer*

| ROOT | wanted | to |          | eat | an | apple |

**reduce**⌒

| ROOT | wanted |          | eat | an | apple |
                                        ↓
**shift**                                to

| ROOT | wanted | eat |          | an | apple |

shift

| ROOT | wanted | eat | an |          | apple |

reduce⌒

| ROOT | wanted | eat |          | apple |
                                        ↓
shift                                    an

| ROOT | wanted | eat | apple |          |

# Arc-hybrid Transition System

# Arc-hybrid Transition System

# Arc-hybrid Transition System



*Stack*                          *Buffer*

**reduce**↶

**shift**

**shift**

**reduce**↶

**shift**

27

# Arc-hybrid Transition System

*Stack*                                                    *Buffer*

| ROOT | wanted | eat | apple |
|------|--------|-----|-------|

**reduce**⤸

| ROOT | wanted | eat |
|------|--------|-----|

apple

**reduce**⤸

| ROOT | wanted |
|------|--------|

eat

**reduce**⤸

| ROOT |
|------|

*terminal*

wanted

28

# Arc-hybrid Transition System

*Stack*                                                                 *Buffer*

| ROOT | wanted | eat | apple |

**reduce↷**

| ROOT | wanted | eat |

apple

reduce↷

| ROOT | wanted |

eat

reduce↶

| ROOT |

wanted

29

# Arc-hybrid Transition System

*Stack*                                      *Buffer*

| ROOT | wanted | eat | apple |
|------|--------|-----|-------|

**reduce**↷

| ROOT | wanted | eat |
|------|--------|-----|

apple

**reduce**↷

| ROOT | wanted |
|------|--------|

eat

reduce↷

terminal   ROOT

wanted

30

# Arc-hybrid Transition System

*Stack*                                                    *Buffer*

| ROOT | wanted | eat | apple |

**reduce**↷

| ROOT | wanted | eat |

apple

**reduce**↷

| ROOT | wanted |

eat

**reduce**↷

*terminal* | ROOT |

wanted

31

# Dynamic Programming for Arc-hybrid

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 $(n)$ |
|---|---|---|---|---|---|---|---|
| | ROOT | She | wanted | to | eat | an | apple |

*Stack*                                *Buffer*

- Deduction Item

| ... | $i$ | | $j$ | ... |
|---|---|---|---|---|

$$[i, j]$$

- Goal

| 0 | | $n + 1$ |
|---|---|---|

$$[0, n + 1]$$

32

# Dynamic Programming for Arc-hybrid

$$\text{shift} \quad \frac{[i,j]}{[j,j+1]}$$

# Dynamic Programming for Arc-hybrid

**reduce**↶    $\dfrac{[k,j]}{[?,j]}$



**reduce**↶

34

# Dynamic Programming for Arc-hybrid

**reduce**$\curvearrowright$   $\dfrac{[k,j]}{\textcolor{red}{[i,j]}}$



**reduce**$\curvearrowright$

# Dynamic Programming for Arc-hybrid

**reduce**↶   $\dfrac{[i,k]\,[k,j]}{[i,j]}$

# Dynamic Programming for Arc-hybrid

In Kuhlmann, Gómez-Rodríguez and Satta (2011)'s notation

$$* \quad [i, k]$$

$$\text{reduce} \quad \frac{[i,k]\,[k,j]}{[i,j]}$$

$$* \quad [k, j]$$

$$* \quad [i, j]$$

**reduce**

37

# Dynamic Programming for Arc-hybrid



**reduce** $\frac{[i,k]\,[k,j]}{[i,j]}$

# Dynamic Programming for Arc-hybrid

**shift**    $$\frac{[i,j]}{[j,j+1]}$$

**Goal:**    $[0, n+1]$

**reduce**⌢    $$\frac{[i,k]\,[k,j]}{[i,j]}\quad k \frown j$$

$$O(n^3)$$

**reduce**⌢    $$\frac{[i,k]\,[k,j]}{[i,j]}\quad i \frown k$$

39

# Time Complexity in Practice

- Complexity depends on feature representation!

- Typical feature representation:
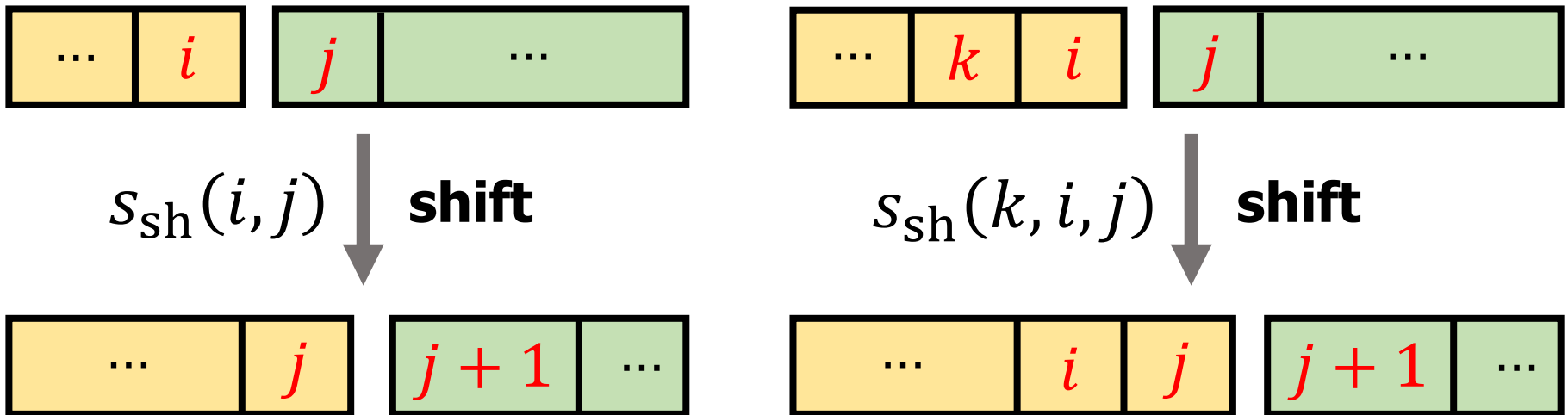  - Feature templates look at specific _positions_ in the stack and in the buffer

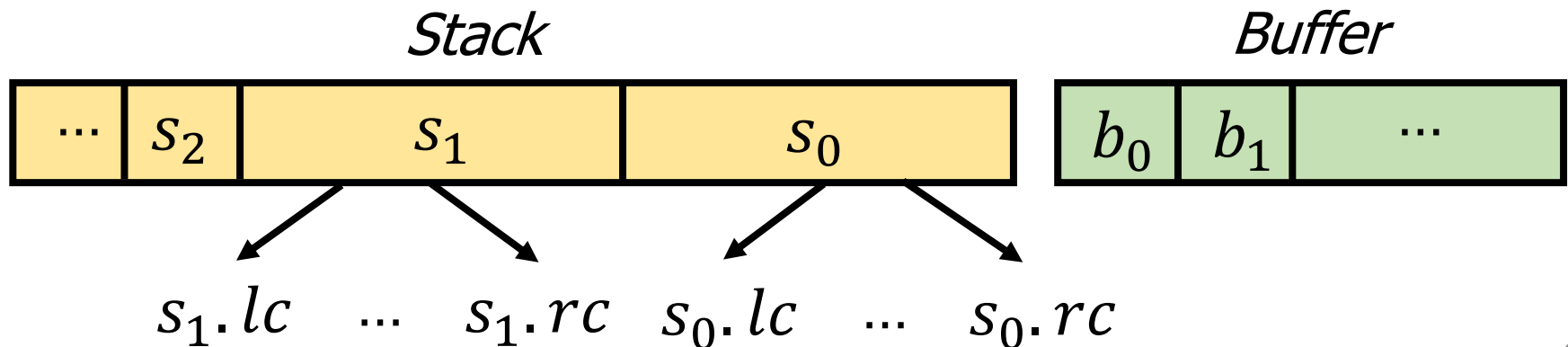# Time Complexity in Practice

- Compare the following features

$$\boxed{\cdots \mid s_0} \quad \boxed{b_0 \mid \cdots} \qquad\qquad \boxed{\cdots \mid s_1 \mid s_0} \quad \boxed{b_0 \mid \cdots}$$

- Time complexities are different!!!

$$\boxed{\cdots \mid i} \quad \boxed{j \mid \cdots} \qquad\qquad \boxed{\cdots \mid k \mid i} \quad \boxed{j \mid \cdots}$$

$$s_{\mathrm{sh}}(i,j) \quad \textbf{shift} \qquad\qquad s_{\mathrm{sh}}(k,i,j) \quad \textbf{shift}$$

$$\boxed{\cdots \mid j} \quad \boxed{j+1 \mid \cdots} \qquad\qquad \boxed{\cdots \mid i \mid j} \quad \boxed{j+1 \mid \cdots}$$

41

# Time Complexity in Practice

- Complexity depends on feature representation!

- Typical feature representation:
  - Feature templates look at specific _positions_ in the stack and in the buffer

- Best-known complexity in practice: $O(n^6)$
  (Huang and Sagae, 2010)



*Stack*      *Buffer*

| ... | $s_2$ | $s_1$ | $s_0$ | | $b_0$ | $b_1$ | ... |

$s_1.lc$   ...   $s_1.rc$   $s_0.lc$   ...   $s_0.rc$

# Best-known Time Complexities (recap)

$$O(n^3)$$

Theoretical

Gap:

Feature representation

$$O(n^6)$$

Practical

# In Practice, Instead of Exact Decoding …

- Greedy search (Nivre, 2003, 2004, 2008; Chen and Manning, 2014)

- Beam search (Zhang and Clark, 2011; Weiss et al.,2015)

- Best-first search (Sagae and Lavie, 2006; Sagae and Tsujii, 2007; Zhao et al., 2013)

- Dynamic oracles (Goldberg and Nivre, 2012, 2013)

- "Global" normalization on the beam (Zhou et al., 2015; Andor et al., 2016)

- Reinforcement learning (Lê and Fokkens, 2017)

- Learning to search (Daumé III and Marcu, 2005; Chang et al., 2016; Wiseman and Rush, 2016)

- …

## How Many Positional Features Do We Need?

Non-neural (manual engineering)

☞ Chen and Manning (2014)

*Stack*                                                *Buffer*

| ... | $s_2$ | $s_1$ | $s_0$ | | $b_0$ | $b_1$ | $b_2$ | ... |

$s_1.lc_i$    ...    $s_1.rc_i$   $s_0.lc_i$    ...    $s_0.rc_i$

$s_1.lc_0.lc_0$       $s_1.rc_0.rc_0$   $s_0.lc_0.lc_0$      $s_0.rc_0.rc_0$

# How Many Positional Features Do We Need?

Non-neural (manual engineering)

☞ Chen and Manning (2014)

**Stack LSTM**
☞ Dyer et al. (2015)

**Bi-LSTM**
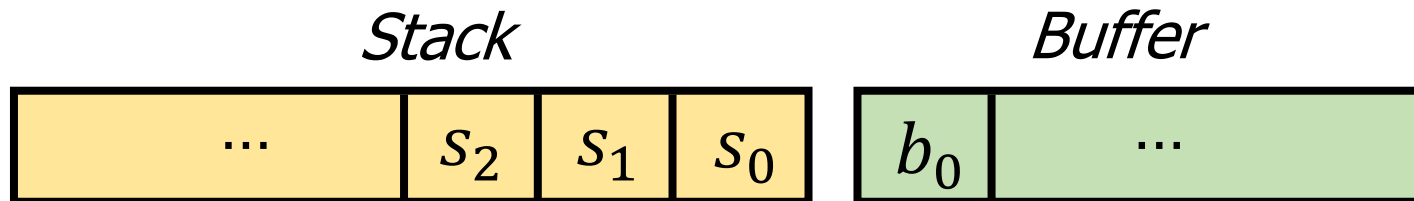☞ Kiperwasser and Goldberg (2016)

☞ Cross and Huang (2016)

*Stack*

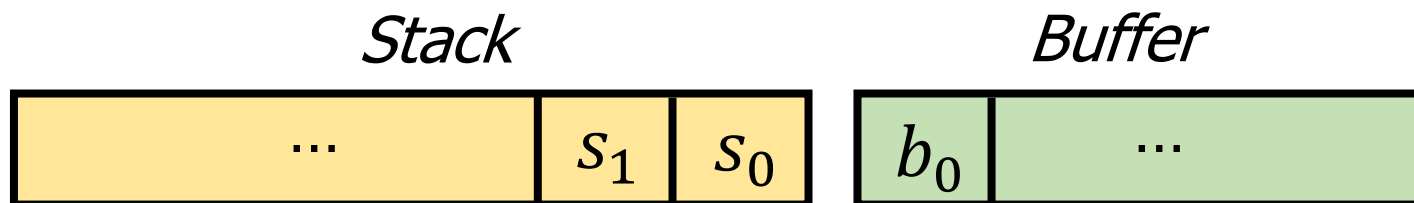| $s_{|\sigma|-1}$ | ... | $s_2$ | $s_1$ | $s_0$ |
|---|---|---|---|---|

...

# How Many Positional Features Do We Need?

- Bi-LSTMs give compact feature representations
  (Kiperwasser and Goldberg, 2016; Cross and Huang, 2016)

- Features used in Kiperwasser and Goldberg (2016)

*Stack*                                    *Buffer*

| ... | $s_2$ | $s_1$ | $s_0$ | | $b_0$ | ... |

- Features used in Cross and Huang (2016)

*Stack*                                    *Buffer*

| ... | $s_1$ | $s_0$ | | $b_0$ | ... |

# How Many Positional Features Do We Need?

Non-neural (manual engineering)

☞ Chen and Manning (2014)

**Stack LSTM**
☞ Dyer et al. (2015)

**Bi-LSTM**
☞ Kiperwasser and Goldberg (2016)

☞ Cross and Huang (2016)

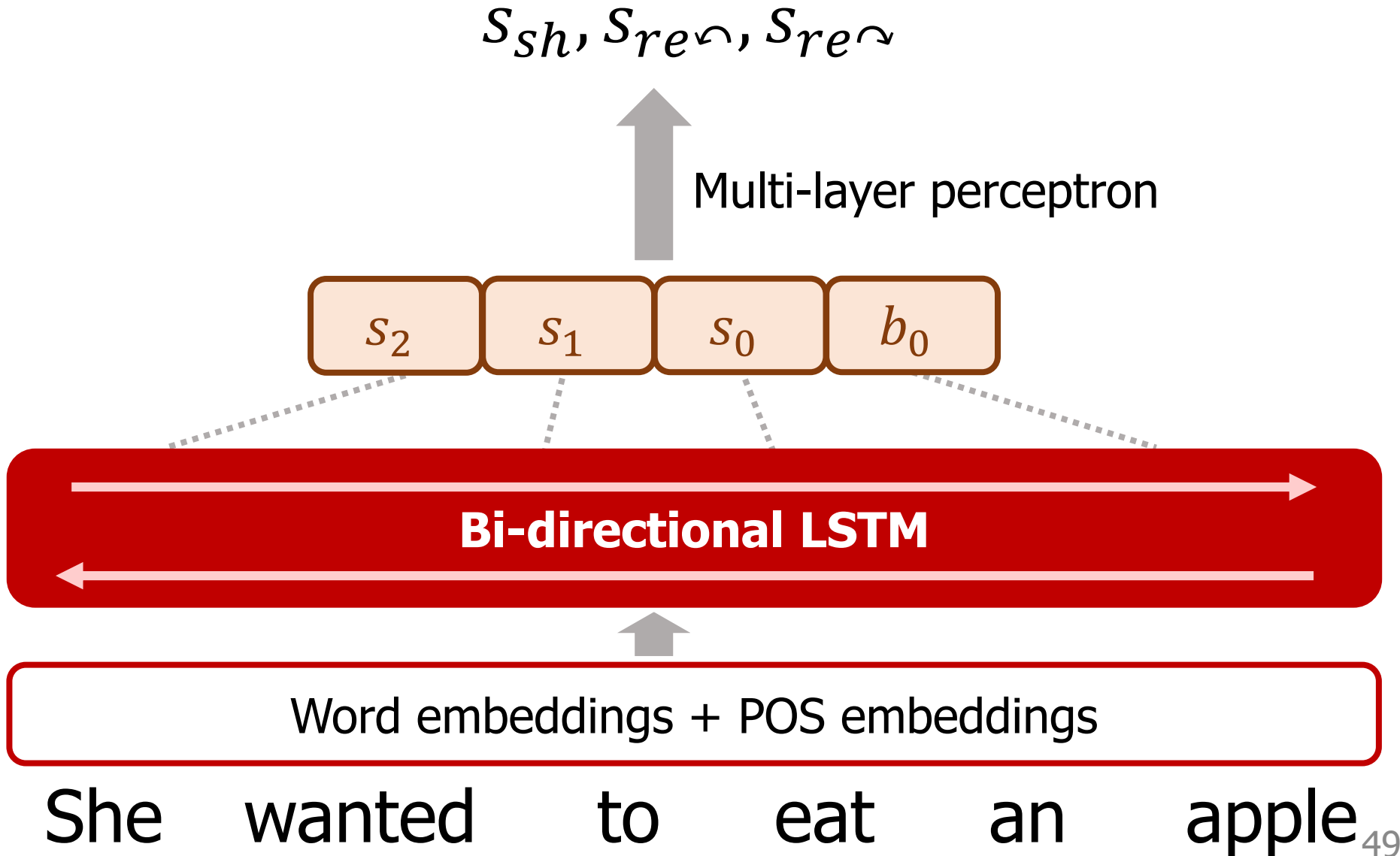Summarizing trees on stack                    *Summarizing input*
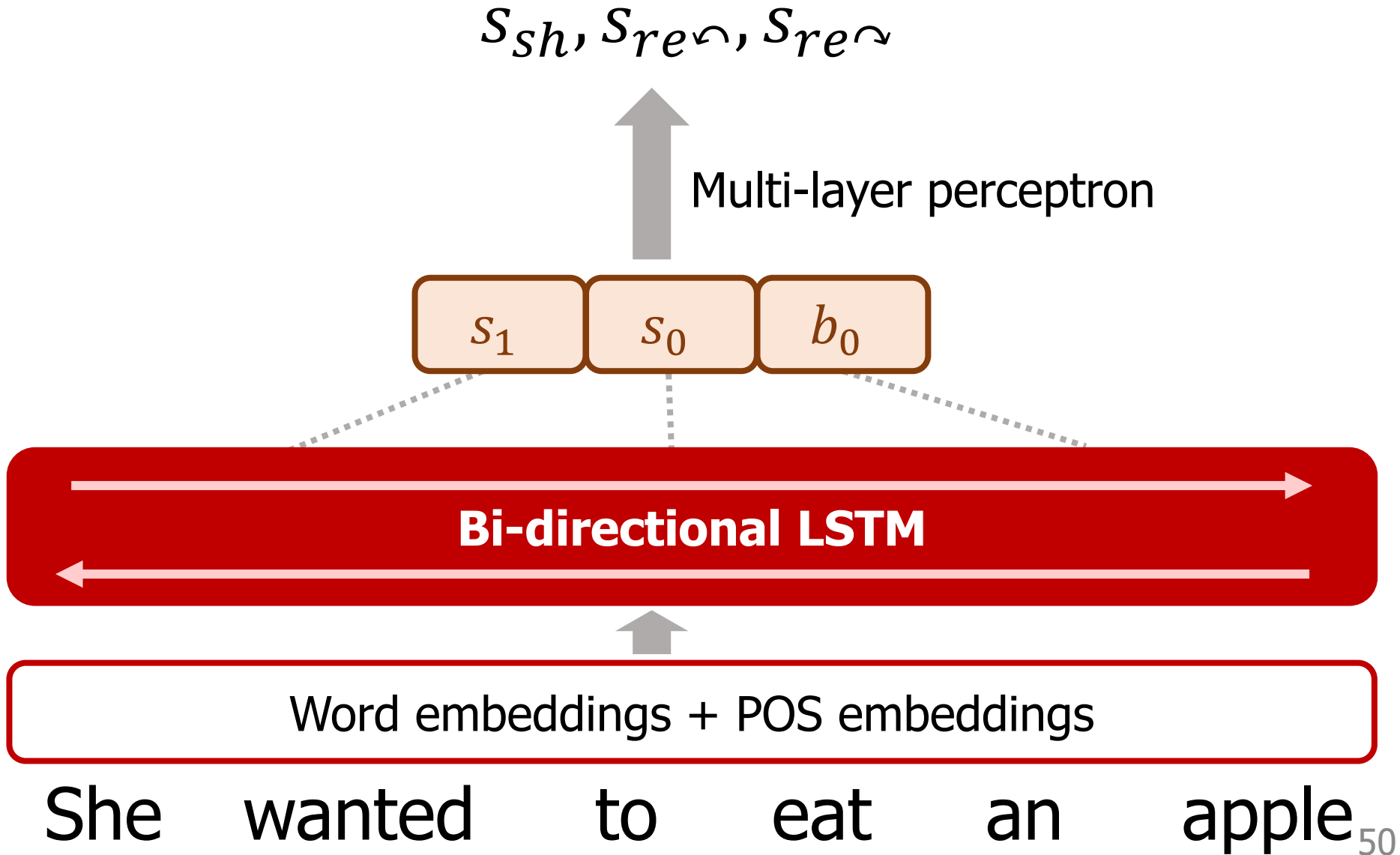
Exponential DP          Enables Slow DP          Enables Fast DP

48

# Model Architecture

$$S_{sh}, S_{re}\frown, S_{re}\frown$$

Multi-layer perceptron

| $s_2$ | $s_1$ | $s_0$ | $b_0$ |

**Bi-directional LSTM**

Word embeddings + POS embeddings

She    wanted    to    eat    an    apple

# Model Architecture

$$S_{sh}, S_{re\curvearrowleft}, S_{re\curvearrowright}$$
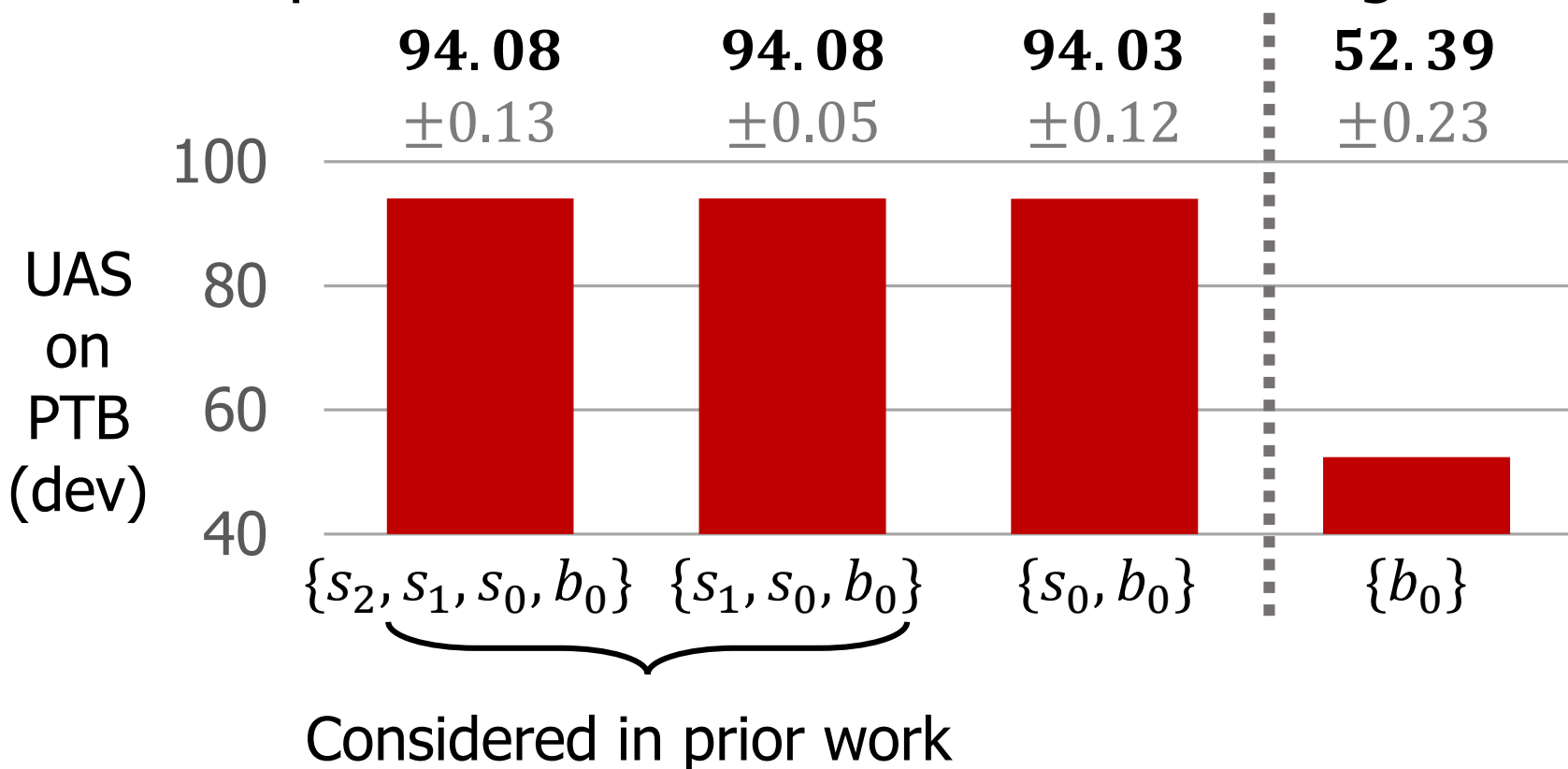
Multi-layer perceptron

| $s_1$ | $s_0$ | $b_0$ |

**Bi-directional LSTM**

Word embeddings + POS embeddings

She    wanted    to    eat    an    apple

# Model Architecture

$$s_{sh}, s_{re\curvearrowleft}, s_{re\curvearrowright}$$

Multi-layer perceptron

$s_0$   $b_0$

**Bi-directional LSTM**

Word embeddings + POS embeddings

She    wanted    to    eat    an    apple

# Model Architecture

$$S_{sh}, S_{re\frown}, S_{re\frown}$$

Multi-layer perceptron

$b_0$

**Bi-directional LSTM**

Word embeddings + POS embeddings

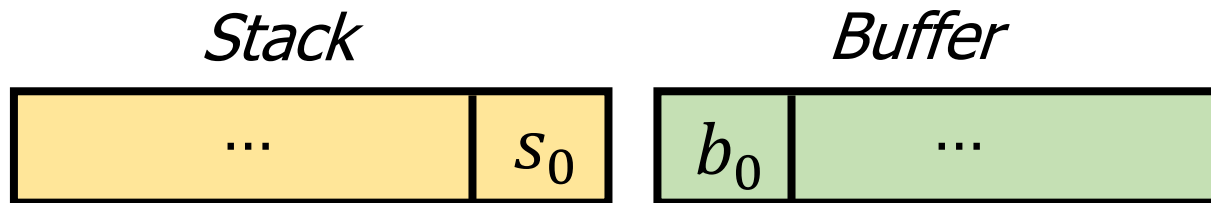She    wanted    to    eat    an    apple

# How Many Positional Features Do We Need?

- We answer the question empirically

  … experimented with <span style="color:red">greedy decoding</span>
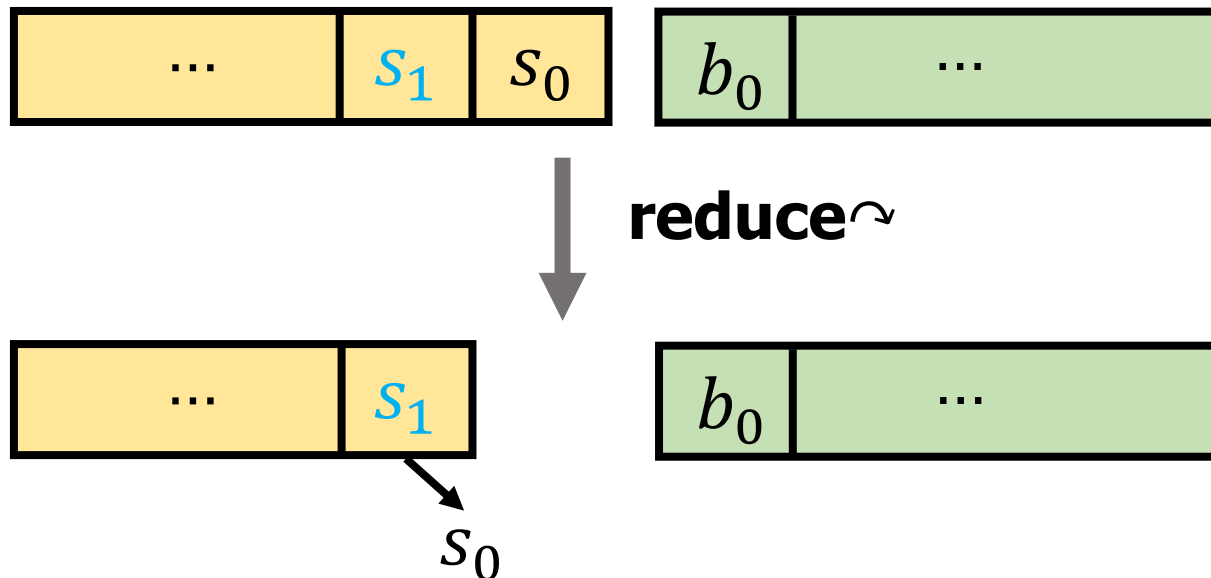
- Two positional feature vectors are enough!



Considered in prior work

53

# How Many Positional Features Do We Need?

- Our minimal feature set

*Stack*                          *Buffer*

| ... | $s_0$ |

| $b_0$ | ... |

- Counter-intuitive, but works for *greedy decoding*

| ... | $s_1$ | $s_0$ |

| $b_0$ | ... |

**reduce**↷

| ... | $s_1$ |

| $b_0$ | ... |

$s_0$

54

# How Many Positional Features Do We Need?

- Our minimal feature set

| *Stack* | | *Buffer* | |
|---|---|---|---|
| ... | $s_0$ | $b_0$ | ... |

- Counter-intuitive, but works for *greedy decoding*

- The bare deduction items already contain enough information to extract features for DP

- Leads to the first $O(n^3)$ implementation of global decoders!

# How Many Positional Features Do We Need?

Non-neural (manual engineering)

☞ Chen and Manning (2014)

**Stack LSTM**
☞ Dyer et al. (2015)

**Bi-LSTM**
☞ Kiperwasser and Goldberg (2016)

☞ Cross and Huang (2016)

☞ *Our work*

*Summarizing trees on stack*          *Summarizing input*

Exponential DP          Enables Slow DP          Enables Fast DP          *Fast(er) DP*          56

# Best-known Time Complexities (recap)

$$O(n^3)$$

Theoretical

Gap:

Feature representation

$$O(n^6)$$

Practical
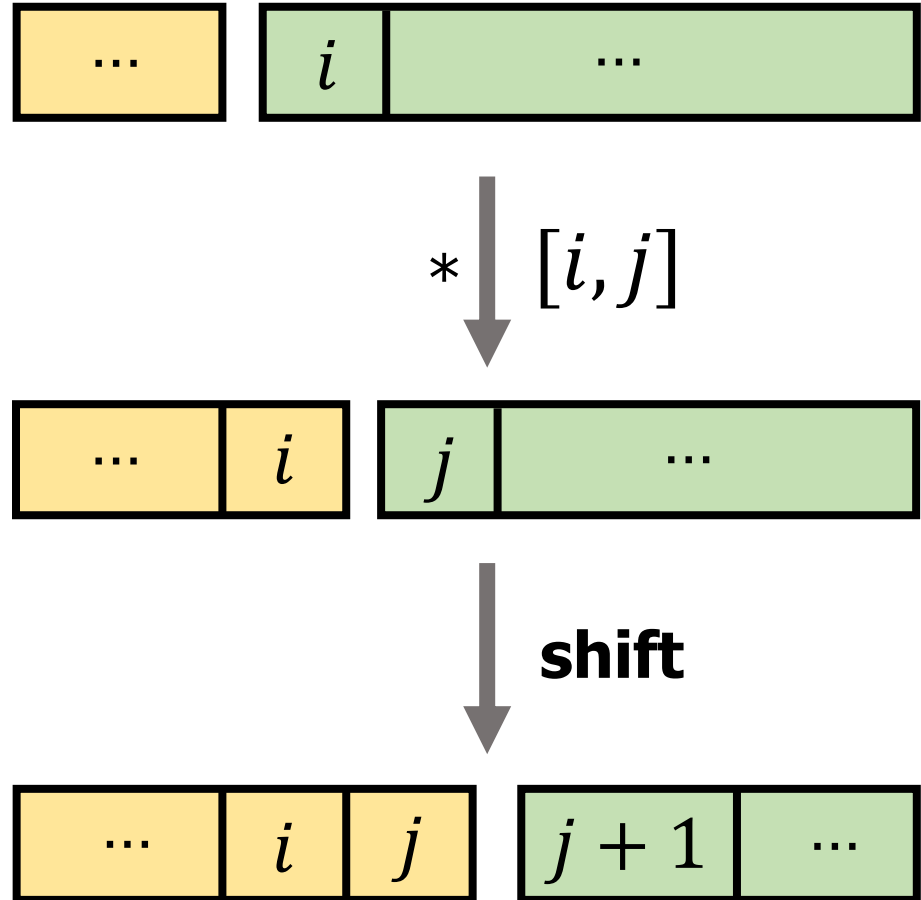
57

# Our contribution

$O(n^3)$

Minimal

Feature Set

$O(n^3)$

$O(n^6)$

Theoretical
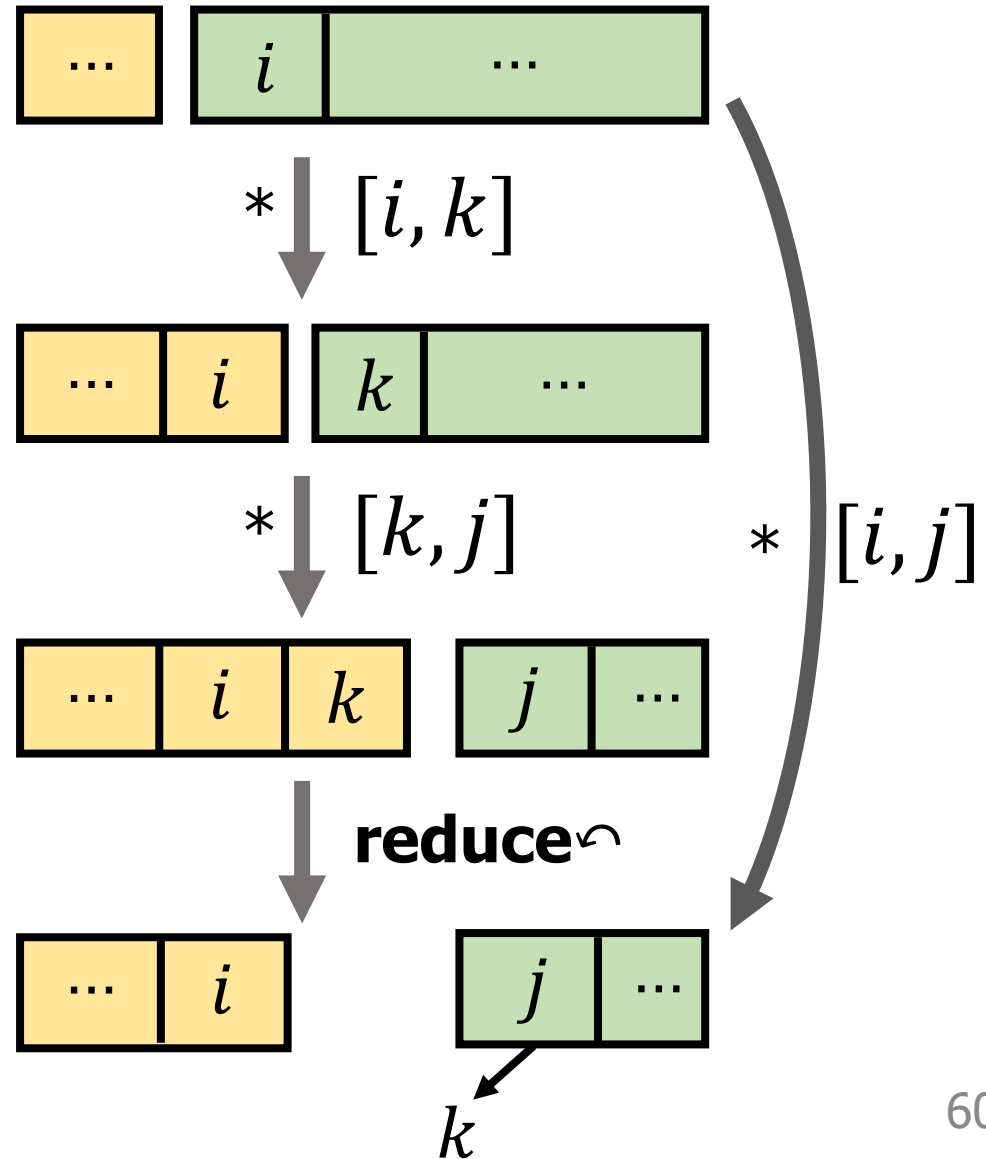
Practical

# Decoding

Score of the sub-sequence

**shift** $\dfrac{[i,j]: v}{[j, j+1]: 0}$

$* \quad [i, j]$

# Decoding



$$\textbf{reduce} \curvearrowright \frac{[i,k]:v_1 \quad [k,j]:v_2}{[i,j]:v_1 + v_2 + \Delta}$$

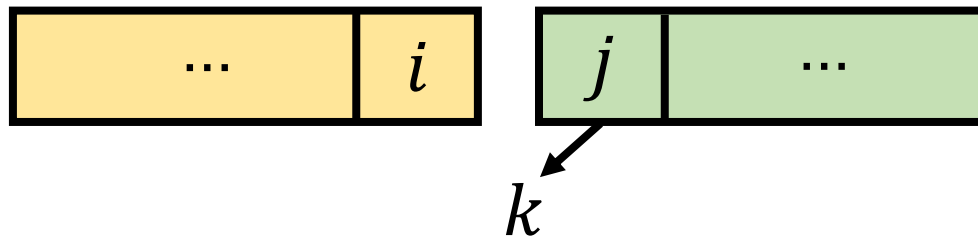$$\Delta = s_{\mathrm{sh}}(i,k) + s_{\mathrm{re}\curvearrowright}(k,j)$$

# Training

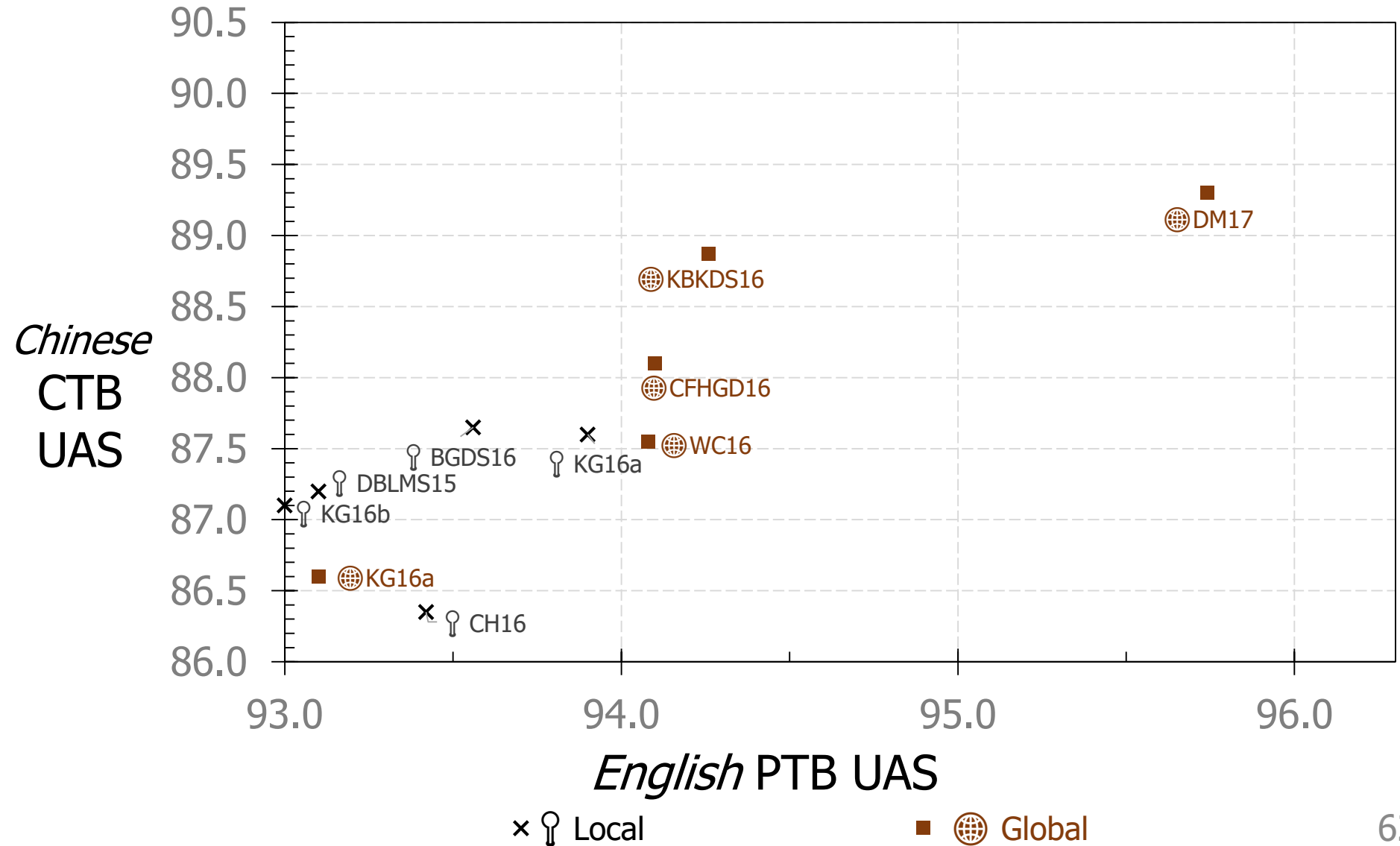- Separate incorrect from correct by a margin

$\max$ score(⬤▸⬤▸ ⋯ ▸⬤) - score(⬤▸⬤▸ ⋯ ▸⭐) + cost(⬤▸⬤▸ ⋯ ▸⬤ / ⬤▸⬤▸ ⋯ ▸⭐)

- Cost-augmented decoding (Taskar et al., 2005)

$$\textbf{reduce}\,\curvearrowright \quad \frac{[i,k]\!:v_1 \qquad [k,j]\!:v_2}{[i,j]\!:v_1 + v_2 + s_{\mathrm{sh}}(i,k) + s_{\mathrm{re}\curvearrowright}(k,j) + \textbf{1}(\mathrm{head}_{\star}(k) \neq j)}$$

$$\boxed{\ \cdots \mid i\ } \quad \boxed{\ j \mid \cdots\ }$$

$k$

61

# Comparing with State-of-the-art



Plot with x-axis "*English* PTB UAS" ranging from 93.0 to 96.0, and y-axis "*Chinese* CTB UAS" ranging from 86.0 to 90.5.

Data points:
- DM17 (Global): ~95.7, 89.3 / label at ~95.6, 89.1
- KBKDS16 (Global): ~94.25, 88.9
- CFHGD16 (Global): ~94.1, 88.1
- WC16 (Global): ~94.1, 87.55
- KG16a (Local): ~93.9, 87.6
- BGDS16 (Local): ~93.55, 87.65
- DBLMS15 (Local): ~93.1, 87.2
- KG16b (Local): ~93.0, 87.1
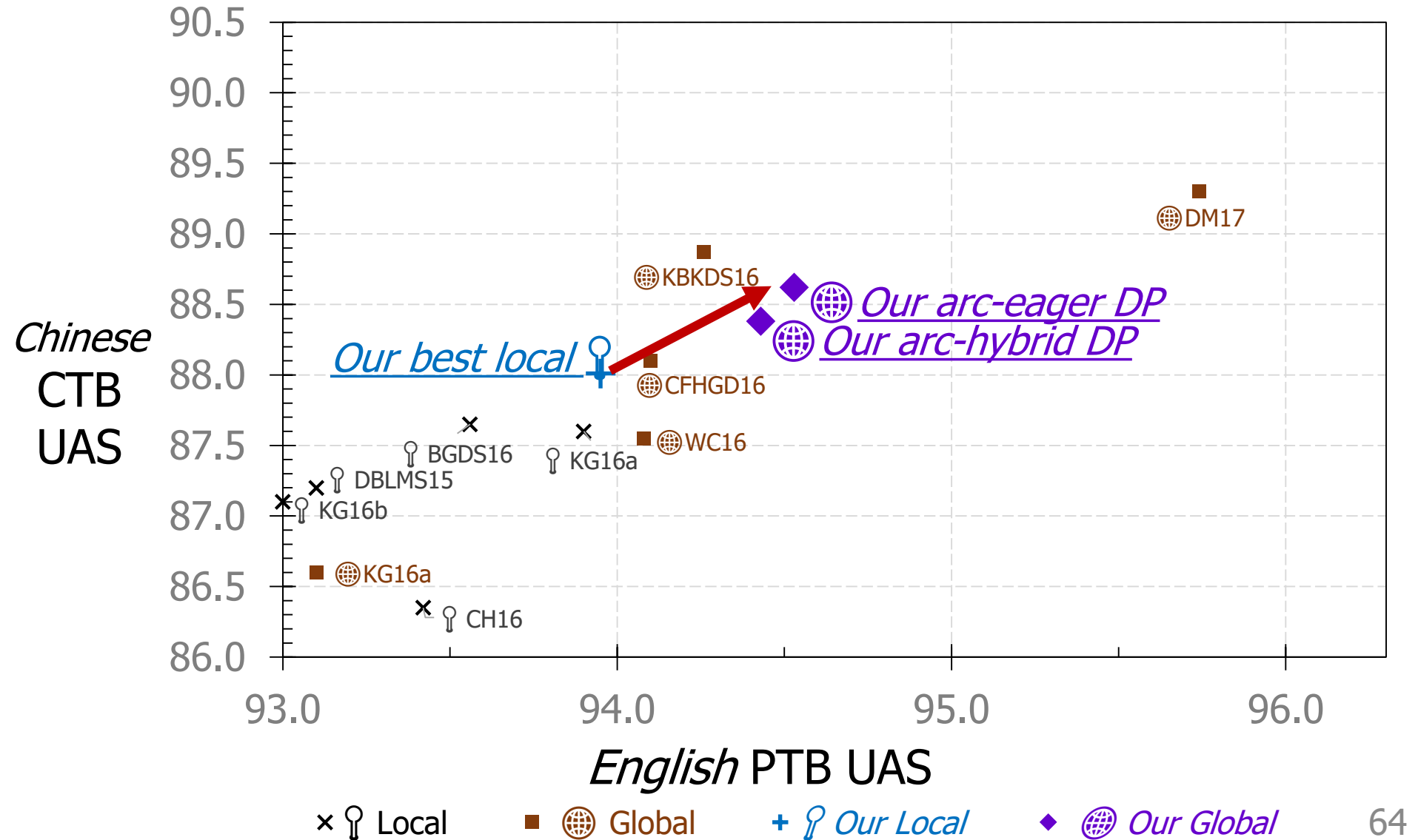- KG16a (Global): ~93.1, 86.6
- CH16 (Local): ~93.45, 86.35

Legend: × / Local       ■ / ⊕ Global

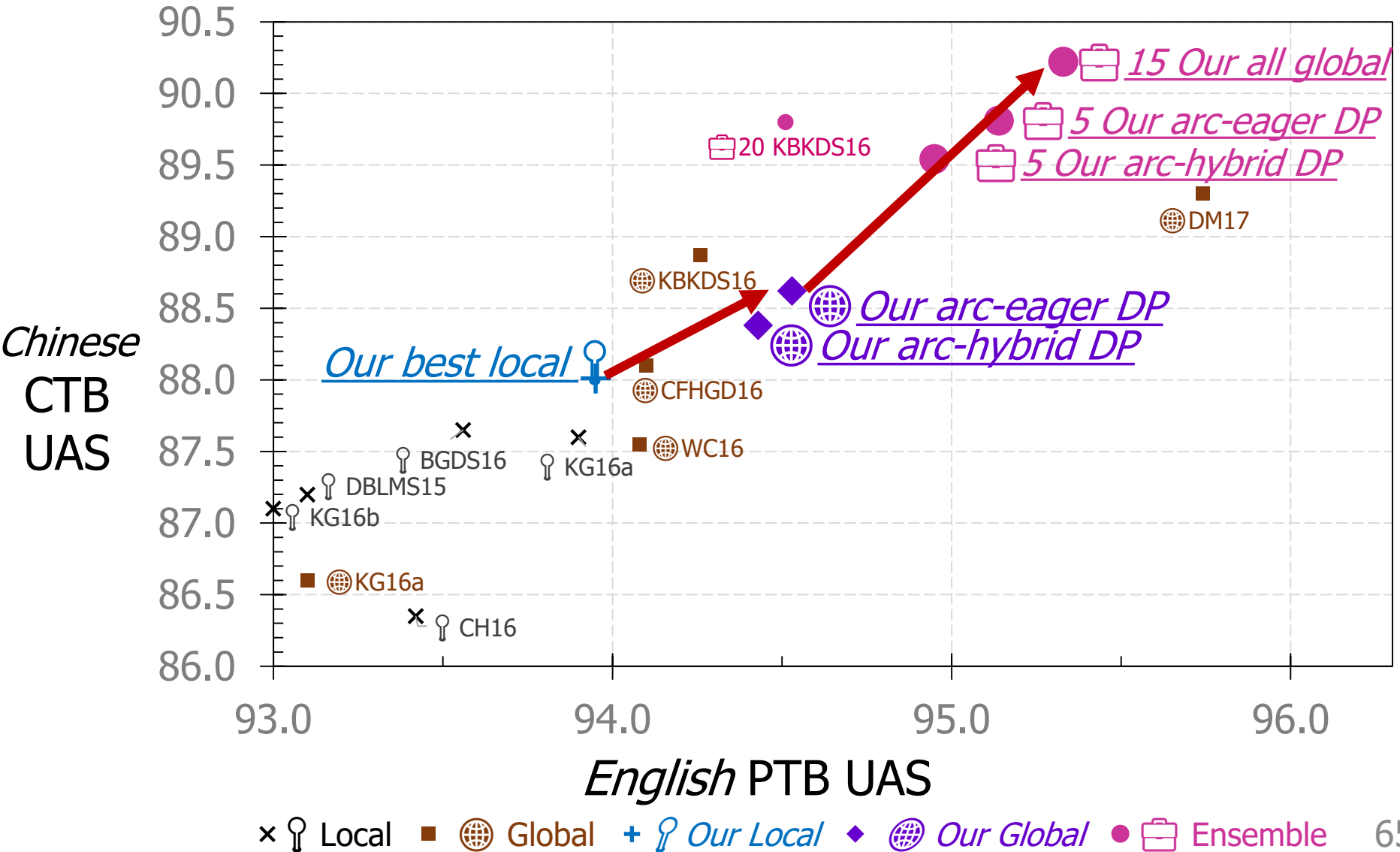# Comparing with State-of-the-art

# Comparing with State-of-the-art
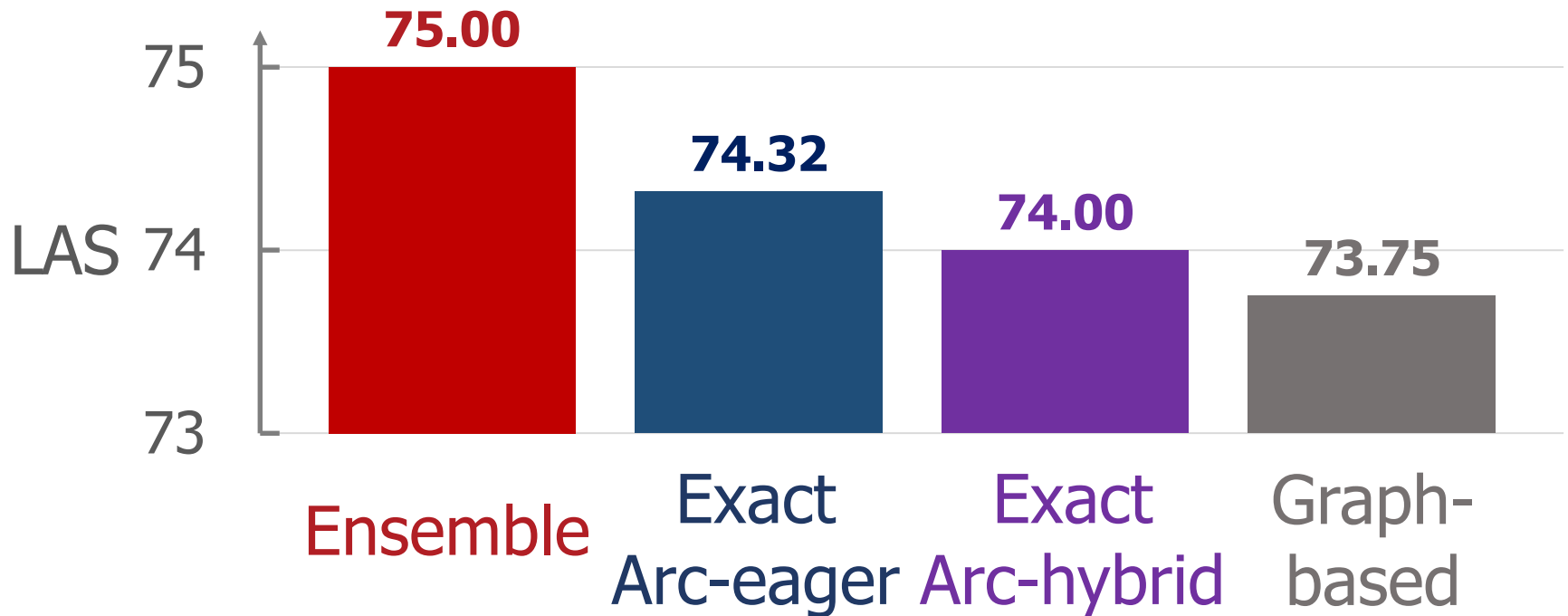
# Comparing with State-of-the-art

# Results – CoNLL'17 Shared Task

- Macro-average of 81 treebanks in 49 languages
- 2[nd]–highest overall performance



(Shi, Wu, Chen and Cheng, 2017; Zeman et al., 2017)  66

# Conclusion

- Bi-LSTM feature set is minimal yet highly effective

- First $O(n^3)$ implementation of exact decoders

- Global training and decoding gave high performance

# More in Our Paper

- Description and analysis of three transition systems (arc-standard, arc-hybrid, arc-eager)

- CKY-style representations of the deduction systems



- Theoretical analysis of the global methods
  - Arc-eager models can "simulate" arc-hybrid models
  - Arc-eager models can "simulate" edge-factored models

# Fast(er) Exact Decoding and Global Training for Transition-Based Dependency Parsing via a Minimal Feature Set

https://github.com/tzshi/dp-parser-emnlp17

Tianze Shi*    Liang Huang†    Lillian Lee*

* Cornell University    † Oregon State University

**Axiom** $[0, 1]$

**Inference Rules**

sh $\quad \dfrac{[i,j]}{[j, j+1]} \qquad j \leqslant n$

re$_{\curvearrowright}$ $\quad \dfrac{[k,i] \quad [i,j]}{[k,j]} \qquad k \overset{\curvearrowright}{} i$

re$_{\curvearrowleft}$ $\quad \dfrac{[k,i] \quad [i,j]}{[k,j]} \qquad i \overset{\curvearrowleft}{} j$

**Goal** $\quad [0, n+1]$

(b) Arc-hybrid

**Axiom** $[0^0, 1]$

**Inference Rules**

sh $\quad \dfrac{[i^b, j]}{[j^0, j+1]} \qquad j \leqslant n$

ra $\quad \dfrac{[i^b, j]}{[j^1, j+1]} \qquad \begin{array}{l} i \overset{\curvearrowright}{} j \\ j \leqslant n \end{array}$

re$_{\curvearrowleft}$ $\quad \dfrac{[k^b,i] \quad [i^0,j]}{[k^b,j]} \qquad i \overset{\curvearrowleft}{} j$

re $\quad \dfrac{[k^b,i] \quad [i^1,j]}{[k^b,j]}$

**Goal** $\quad [0^0, n+1]$

(c) Arc-eager

70

**Axiom** $[0, 1]$

**Inference Rules**

sh $\dfrac{[i,j]}{[j,j+1]}$    $j \leqslant n$

re$\curvearrowright$ $\dfrac{[k,i] \quad [i,j]}{[k,j]}$    $k \curvearrowright i$

re$\curvearrowleft$ $\dfrac{[k,i] \quad [i,j]}{[k,j]}$    $i \curvearrowleft j$

**Goal** $[0, n+1]$

(b) Arc-hybrid

**Axioms**    $0 \leqslant i, j \leqslant n$

**Inference Rules**

right-attach

right-reduce    $i \curvearrowright j$

left-attach

left-reduce    $i \curvearrowleft j$

**Goal**

(d) Edge-factored graph-based parsing.

71

| Features | Arc-standard | Arc-hybrid | Arc-eager |
|---|---|---|---|
| $\{\overset{\rightarrow\leftarrow}{s}_2, \overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.95_{\pm 0.12}$ | $94.08_{\pm 0.13}$ | $93.92_{\pm 0.04}$ |
| $\{\overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $94.13_{\pm 0.06}$ | $94.08_{\pm 0.05}$ | $93.91_{\pm 0.07}$ |
| $\{\overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $54.47_{\pm 0.36}$ | $94.03_{\pm 0.12}$ | $93.92_{\pm 0.07}$ |
| $\{\overset{\rightarrow\leftarrow}{b}_0\}$ | $47.11_{\pm 0.44}$ | $52.39_{\pm 0.23}$ | $79.15_{\pm 0.06}$ |

| Min positions | Arc-standard | Arc-hybrid | Arc-eager |
|---|---|---|---|
| K&G 2016a | - | 4 | - |
| C&H 2016a | 3 | - | - |
| our work | 3 | **2** | **2** |

# Results with Arc-eager and Arc-standard

| Model | Training | Features | PTB | | CTB | |
|---|---|---|---|---|---|---|
| | | | UAS (%) | UEM (%) | UAS (%) | UEM (%) |
| Arc-standard | Local | $\{\overset{\rightharpoonup}{s}_2, \overset{\rightharpoonup}{s}_1, \overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $93.95_{\pm 0.12}$ | $52.29_{\pm 0.66}$ | $88.01_{\pm 0.26}$ | $36.87_{\pm 0.53}$ |
| Arc-hybrid | Local | $\{\overset{\rightharpoonup}{s}_2, \overset{\rightharpoonup}{s}_1, \overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $93.89_{\pm 0.10}$ | $50.82_{\pm 0.75}$ | $87.87_{\pm 0.17}$ | $35.47_{\pm 0.48}$ |
| | Local | $\{\overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $93.80_{\pm 0.12}$ | $49.66_{\pm 0.43}$ | $87.78_{\pm 0.09}$ | $35.09_{\pm 0.40}$ |
| | Global | $\{\overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $94.43_{\pm 0.08}$ | $53.03_{\pm 0.71}$ | $88.38_{\pm 0.11}$ | $36.59_{\pm 0.27}$ |
| Arc-eager | Local | $\{\overset{\rightharpoonup}{s}_2, \overset{\rightharpoonup}{s}_1, \overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $93.80_{\pm 0.12}$ | $49.66_{\pm 0.43}$ | $87.49_{\pm 0.20}$ | $33.15_{\pm 0.72}$ |
| | Local | $\{\overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $93.77_{\pm 0.08}$ | $49.71_{\pm 0.24}$ | $87.33_{\pm 0.11}$ | $34.17_{\pm 0.41}$ |
| | Global | $\{\overset{\rightharpoonup}{s}_0, \overset{\rightharpoonup}{b}_0\}$ | $\mathbf{94.53}_{\pm 0.05}$ | $53.77_{\pm 0.46}$ | $\mathbf{88.62}_{\pm 0.09}$ | $\mathbf{37.75}_{\pm 0.87}$ |
| Edge-factored | Global | $\{\overset{\rightharpoonup}{h}, \overset{\rightharpoonup}{m}\}$ | $94.50_{\pm 0.13}$ | $\mathbf{53.86}_{\pm 0.78}$ | $88.25_{\pm 0.12}$ | $36.42_{\pm 0.52}$ |