# Fast(er) Exact Decoding and Global Training for *TBDPs* via a Minimal Feature Set

Tianze Shi, Liang Huang, Lillian Lee

Cornell & Oregon State –– EMNLP17

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

AntNLP –– Tao Ji

taoji.cs@gmail.com

# Outline

- Transition-based Dependency Parsing

- Three Transition Systems

- Motivation

- A Minimal Feature Set

- Dynamic Programming for TBDPs

- Practical Optimal Algorithms

- Experiments & CoNLL 2017 Shared Task

- Conclusion

# TBDPs

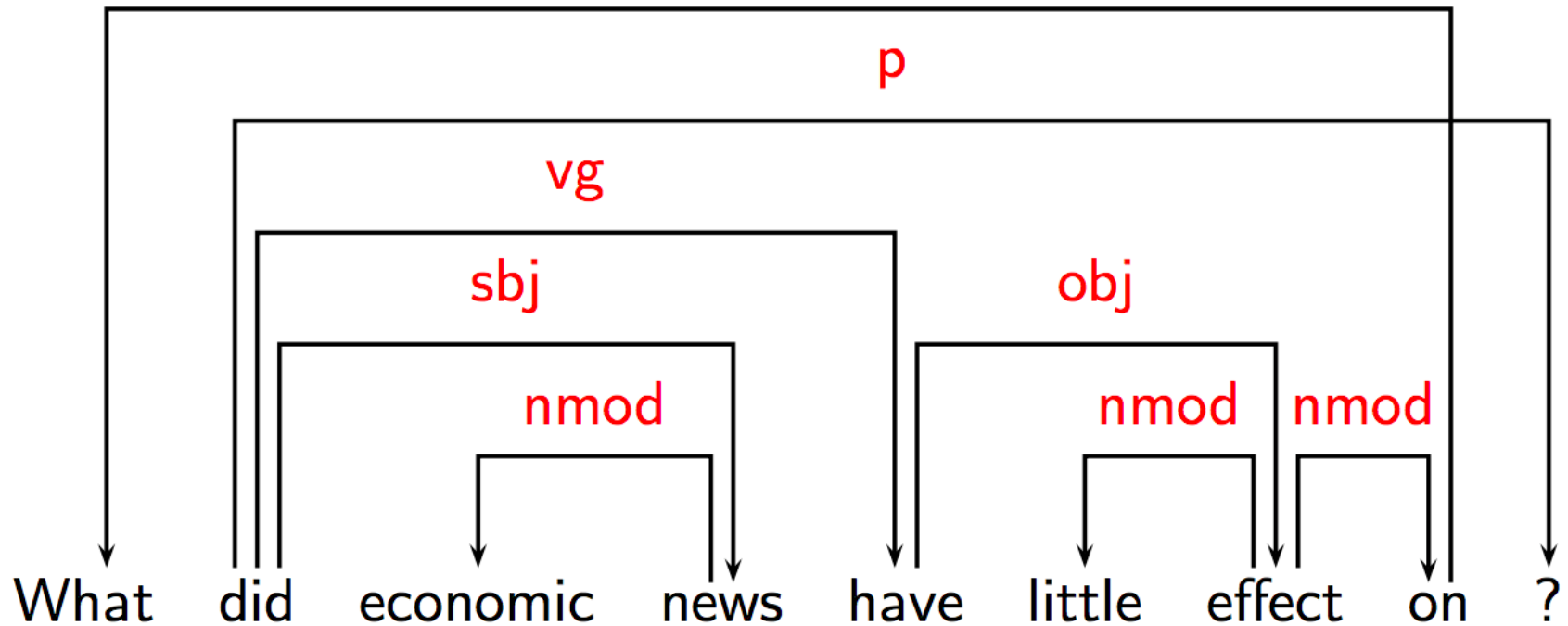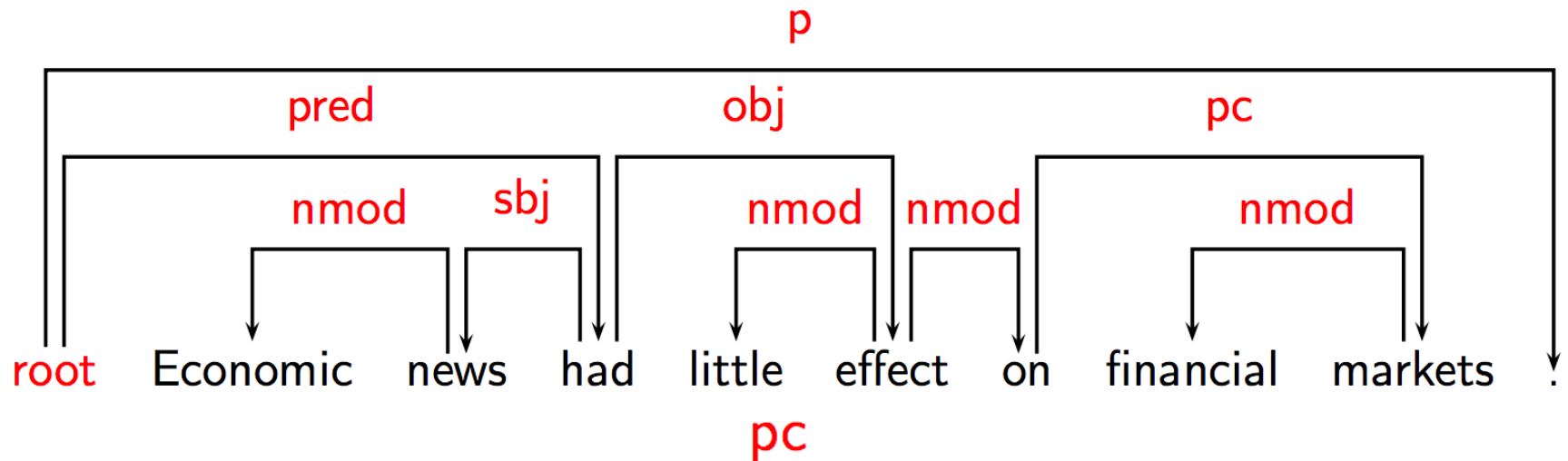| | |
|---|---|
| **Configuration:** | $(S, B, A)$ |
| **Initial:** | $([\ ], [0, 1, \ldots, n], \{\ \})$ |
| **Terminal:** | $(S, [\ ], A)$ |
| **Shift:** | $(S, i|B, A) \Rightarrow (S|i, B, A)$ |
| **Reduce:** | $(S|i, B, A) \Rightarrow (S, B, A)$ |
| **Right-Arc($k$):** | $(S|i, j|B, A) \Rightarrow (S|i|j, B, A \cup \{(i, j, k)\})$ |
| **Left-Arc($k$):** | $(S|i, j|B, A) \Rightarrow (S, j|B, A \cup \{(j, i, k)\})$ |

$\Longleftrightarrow$



Economic news had little effect on financial markets .
adj    noun  verb adj  noun   prep adj      noun

# Projectivity

# Arc-Eager

**Configuration:**  $(S, B, A)$    $[S = \text{Stack}, B = \text{Buffer}, A = \text{Arcs}]$

**Initial:**   $([\,], [0, 1, \ldots, n], \{\,\})$

**Terminal:**   $(S, [\,], A)$

**Shift:**   $(S, i|B, A) \quad \Rightarrow \quad (S|i, B, A)$

**Reduce:**   $(S|i, B, A) \quad \Rightarrow \quad (S, B, A)$    $h(i, A)$

**Right-Arc($k$):**   $(S|i, j|B, A) \quad \Rightarrow \quad (S|i|j, B, A \cup \{(i, j, k)\})$

**Left-Arc($k$):**   $(S|i, j|B, A) \quad \Rightarrow \quad (S, j|B, A \cup \{(j, i, k)\})$    $\neg h(i, A) \wedge i \neq 0$

Notation:   $S|i$ = stack with top $i$ and remainder $S$

$j|B$ = buffer with head $j$ and remainder $B$

$h(i, A)$ = $i$ has a head in $A$

# Arc-Standard

**Configuration:** $(S, B, A)$    [$S =$ Stack, $B =$ Buffer, $A =$ Arcs]

**Initial:** $([\,], [0, 1, \ldots, n], \{\,\})$

**Terminal:** $([0], [\,], A)$

**Shift:** $(S, i|B, A) \Rightarrow (S|i, B, A)$

**Right-Arc($k$):** $(S|i|j, B, A) \Rightarrow (S|i, B, A \cup \{(i, j, k)\})$

**Left-Arc($k$):** $(S|i|j, B, A) \Rightarrow (S|j, B, A \cup \{(j, i, k)\})$    $i \neq 0$

# Arc-Hybrid

$$\text{sh}[(\sigma, b_0|\beta, A)] = (\sigma|b_0, \beta, A)$$

$$\text{re}_{\curvearrowright}[(\sigma|s_1|s_0, \beta, A)] = (\sigma|s_1, \beta, A \cup \{(s_1, s_0)\})$$
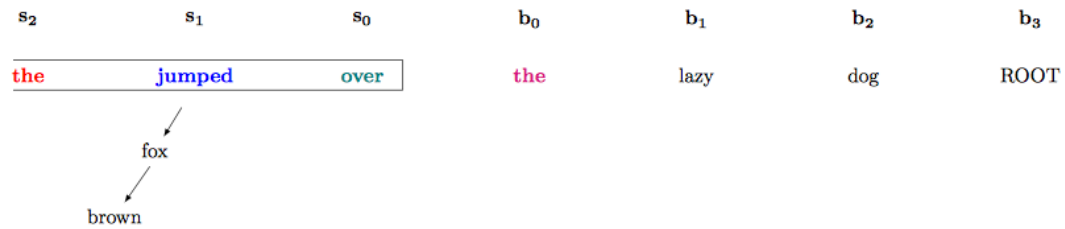
$$\text{re}_{\curvearrowleft}[(\sigma|s_0, b_0|\beta, A)] = (\sigma, b_0|\beta, A \cup \{(b_0, s_0)\})$$
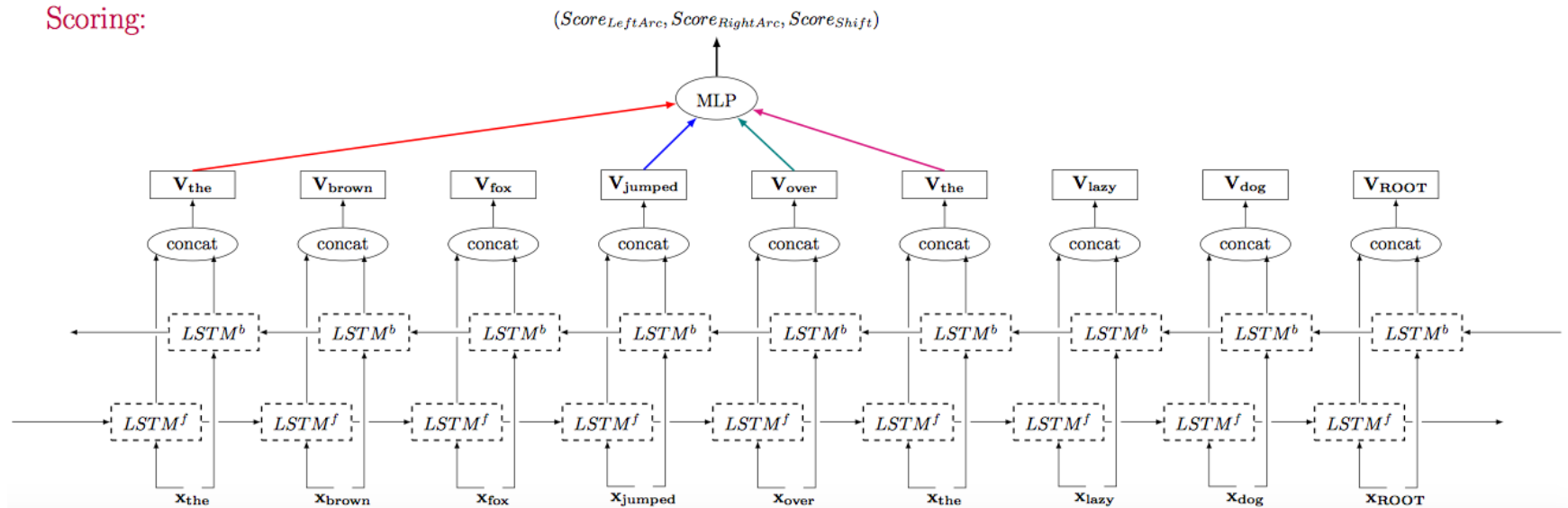
# Introduction

- Plug their minimal feature set into the dynamic-programming framework.

- Produce the first implementation of worst-case $O(n^3)$ exact decoders for arc-hybrid and arc-eager transition systems.

- With their minimal features, we also present $O(n^3)$ global training methods.

- Their achieve the best UAS reported (to our knowledge) on the CTB and the "second-best-in-class" result on the PTB.

- Their had the top average performance on the four surprise languages and on the small treebank subset.

# BIST Parser

# Minimal Feature Set

| Features | Arc-standard | Arc-hybrid | Arc-eager |
|---|---|---|---|
| $\{\overset{\rightarrow\leftarrow}{s}_2, \overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.95_{\pm0.12}$ | $94.08_{\pm0.13}$ | $93.92_{\pm0.04}$ |
| $\{\overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $94.13_{\pm0.06}$ | $94.08_{\pm0.05}$ | $93.91_{\pm0.07}$ |
| $\{\overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $54.47_{\pm0.36}$ | $94.03_{\pm0.12}$ | $93.92_{\pm0.07}$ |
| $\{\overset{\rightarrow\leftarrow}{b}_0\}$ | $47.11_{\pm0.44}$ | $52.39_{\pm0.23}$ | $79.15_{\pm0.06}$ |

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

# Dynamic Programming

- If beam search reduces search errors, why not exact inference?
- Dynamic programming for transition-based parsers:
    - Using a graph-structured stack [Huang and Sagae 2010]
    - Using push-computations [Kuhlmann et al. 2011]
- Adds constraints on feature representations

## Features

Overlapping Subproblems

Optimal Substructure

# Deduction System for Arc-Eager Parsing

**Items:**   $[i^b, j] \Leftrightarrow (S, i|B, A) \Rightarrow^* (S|i, j|B', A')$

$$b = \begin{cases} 1 & \text{if } [\![ h(i) \in A' ]\!] \\ 0 & \text{otherwise} \end{cases}$$

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

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

**Rules:**

Shift:   $[i^b, j] \Rightarrow [j^0, j+1]$

Reduce:   $[i^b, m] \wedge [m^1, j] \Rightarrow [i^b, j]$

Right-Arc:   $[i^b, j] \Rightarrow [j^1, j+1]$

Left-Arc:   $[i^b, m] \wedge [m^0, j] \Rightarrow [i^b, j]$

[Kuhlmann et al. 2011]

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



**Inference Rules**

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



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



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



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



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



(c) Arc-eager

# Deduction System

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

**Inference Rules**

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

**re** $\dfrac{[i, h_1, k] \quad [k, h_2, j]}{[i, h_1, j]}$ $\qquad h_1 \frown h_2$

**re** $\dfrac{[i, h_1, k] \quad [k, h_2, j]}{[i, h_2, j]}$ $\qquad h_1 \frown h_2$

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

(a) Arc-standard

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

**Inference Rules**

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

**re** $\dfrac{[k, i] \quad [i, j]}{[k, j]}$ $\qquad k \frown i$

**re** $\dfrac{[k, i] \quad [i, j]}{[k, j]}$ $\qquad i \frown j$

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

(b) Arc-hybrid

# Exact Decoding

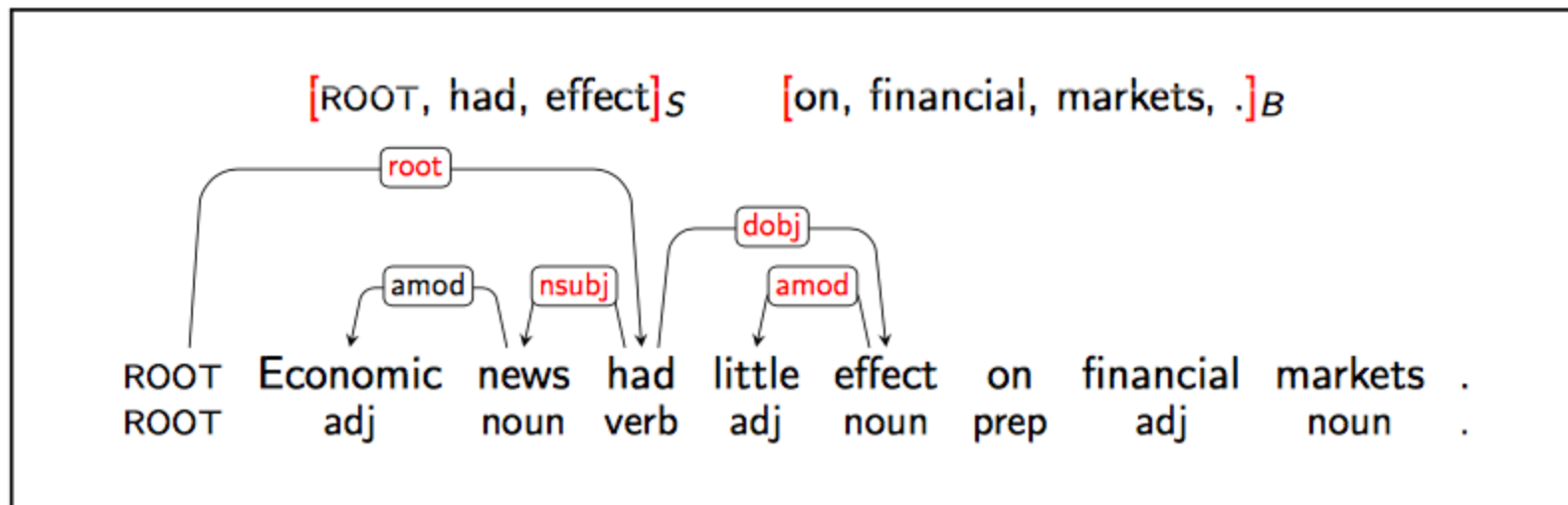$$\frac{[i^b, j] : v}{[j^0, j+1] : 0} \text{(sh)} \qquad \frac{[k^b, i] : v_1 \qquad [i^0, j] : v_2}{[k^b, j] : v_1 + v_2 + \Delta} \text{(re} \curvearrowleft )$$

**Configuration**



$[\text{ROOT, had, effect}]_S \qquad [\text{on, financial, markets, .}]_B$

# Global Training

$$\max_{\mathbf{t}} \left( F(\mathbf{t}) + cost(\mathbf{t}^{\mathrm{gold}}, \mathbf{t}) - F(\mathbf{t}^{\mathrm{gold}}) \right)$$
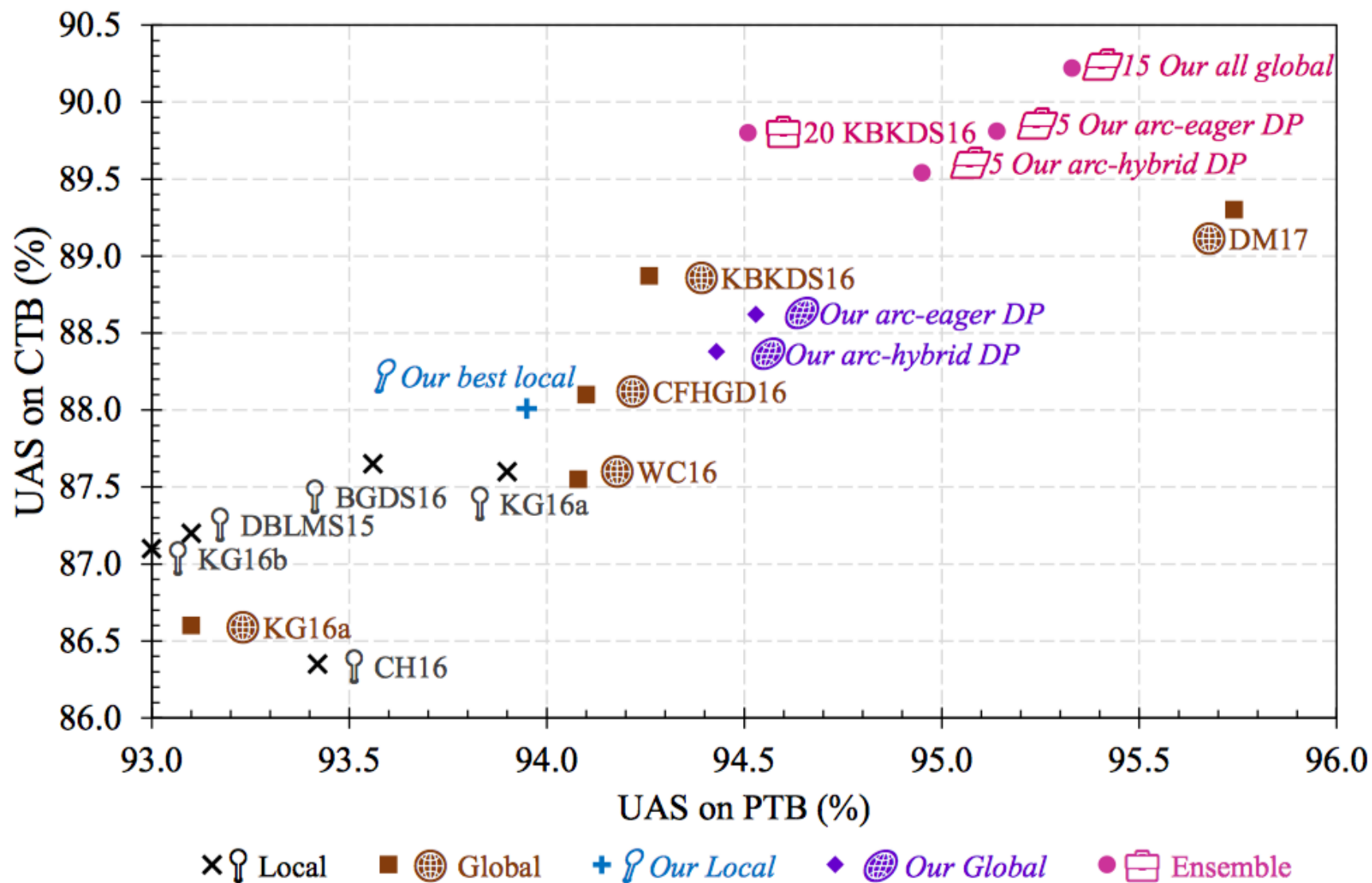
$$\frac{[k^b, i] : v_1 \quad [i^0, j] : v_2}{[k^b, j] : v_1 + v_2 + \Delta'} (\mathbf{re}_\curvearrowleft)$$

$$\text{where } \Delta' = \Delta + \mathbf{1} \left( head(w_i) \neq w_j \right).$$

# Experiments

| Model | Training | Features | PTB | | CTB | |
|---|---|---|---|---|---|---|
| | | | UAS (%) | UEM (%) | UAS (%) | UEM (%) |
| Arc-standard | Local | $\{\overset{\rightarrow\leftarrow}{s}_2, \overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.95_{\pm0.12}$ | $52.29_{\pm0.66}$ | $88.01_{\pm0.26}$ | $36.87_{\pm0.53}$ |
| Arc-hybrid | Local | $\{\overset{\rightarrow\leftarrow}{s}_2, \overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.89_{\pm0.10}$ | $50.82_{\pm0.75}$ | $87.87_{\pm0.17}$ | $35.47_{\pm0.48}$ |
| | Local | $\{\overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.80_{\pm0.12}$ | $49.66_{\pm0.43}$ | $87.78_{\pm0.09}$ | $35.09_{\pm0.40}$ |
| | Global | $\{\overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $94.43_{\pm0.08}$ | $53.03_{\pm0.71}$ | $88.38_{\pm0.11}$ | $36.59_{\pm0.27}$ |
| Arc-eager | Local | $\{\overset{\rightarrow\leftarrow}{s}_2, \overset{\rightarrow\leftarrow}{s}_1, \overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.80_{\pm0.12}$ | $49.66_{\pm0.43}$ | $87.49_{\pm0.20}$ | $33.15_{\pm0.72}$ |
| | Local | $\{\overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $93.77_{\pm0.08}$ | $49.71_{\pm0.24}$ | $87.33_{\pm0.11}$ | $34.17_{\pm0.41}$ |
| | Global | $\{\overset{\rightarrow\leftarrow}{s}_0, \overset{\rightarrow\leftarrow}{b}_0\}$ | $\mathbf{94.53}_{\pm0.05}$ | $53.77_{\pm0.46}$ | $\mathbf{88.62}_{\pm0.09}$ | $\mathbf{37.75}_{\pm0.87}$ |
| Edge-factored | Global | $\{\overset{\rightarrow\leftarrow}{h}, \overset{\rightarrow\leftarrow}{m}\}$ | $94.50_{\pm0.13}$ | $\mathbf{53.86}_{\pm0.78}$ | $88.25_{\pm0.12}$ | $36.42_{\pm0.52}$ |

# Experiments

# CoNLL 2017 Shared Task

# CoNLL 2017 Shared Task

Following Dozat and Manning (2017), we use a deep bi-affine scoring function:

$$\text{score}^{\text{MST}}(h, m) = v_h^\intercal U v_m + b_h \cdot v_h + b_m \cdot v_m + b$$

where

$$v_h = \text{MLP}^{\text{MST-head}}(\overset{\rightarrow\leftarrow}{h})$$

$$v_m = \text{MLP}^{\text{MST-mod}}(\overset{\rightarrow\leftarrow}{m})$$

# CoNLL 2017 Shared Task

|  | UAS F1 | LAS F1 | Official Ranking |
|---|---|---|---|
| Big Treebanks | 85.16 | 79.85 | 2 |
| Small Treebanks | 70.59 | 61.49 | 1 |
| PUD Treebanks | 80.17 | 71.49 | 2 |
| Surprise Languages | 58.40 | 47.54 | 1 |
| Overall | 80.35 | 75.00 | 2 |

# CoNLL 2017 Shared Task

| Target | Source | UAS F1 | LAS F1 | Official Ranking |
|--------|--------|--------|--------|------------------|
| bxr | hi | 50.79 | 31.98 | 2 |
| hsb | cs | 69.45 | 61.70 | 1 |
| kmr | fa | 54.51 | 47.53 | 1 |
| sme | fi | 58.85 | 48.96 | 1 |
| Average | | 58.40 | 47.54 | 1 |

# CoNLL 2017 Shared Task

# Thank you!

# Q&A