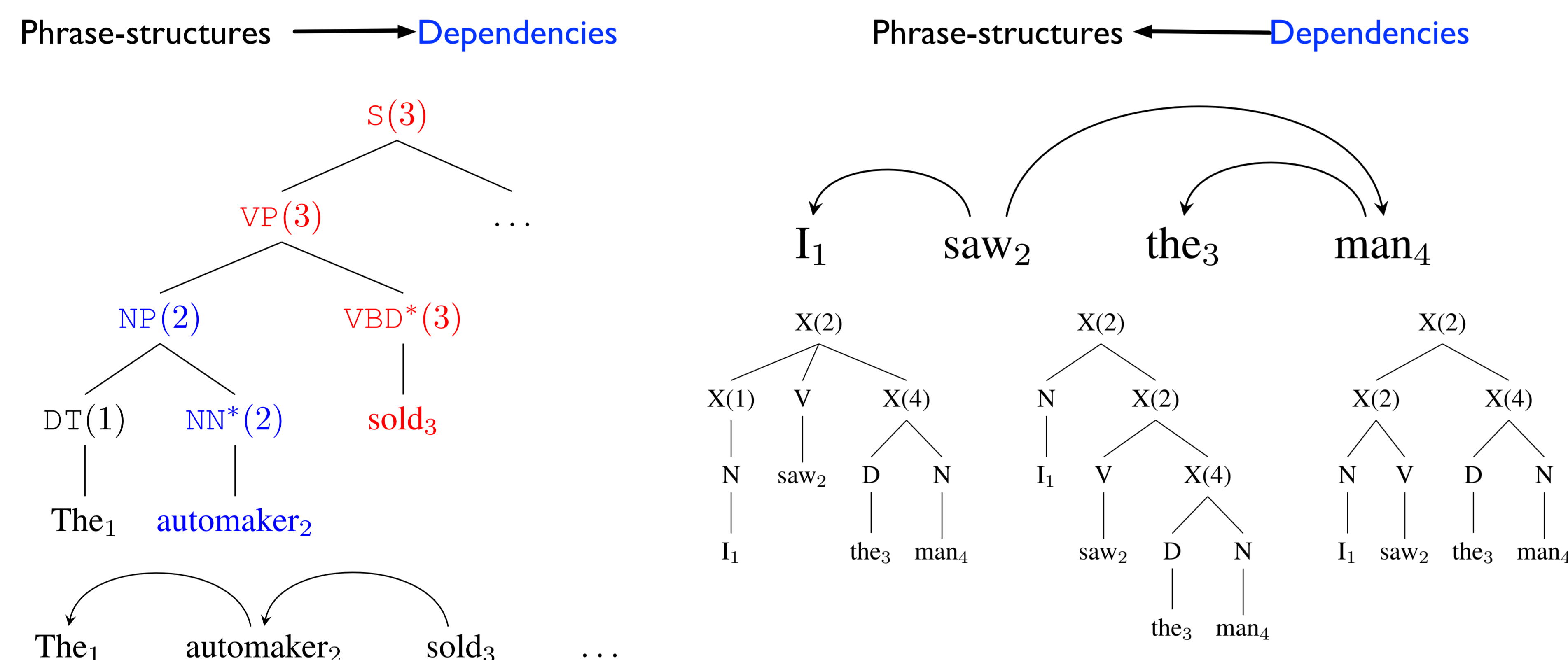


## Contributions

- A **phrase-structure parser (PAD)** achieves **0.4% higher f-score** on the Penn Treebank and **~7x faster** than the Berkeley parser, without reranking or semisupervised training.
- An **linear observable time** algorithm for **transforming dependency parse trees into phrase-structure parse trees**.

## Transformation



• **Various Head-rules** [Collins, 2003; De Marneffe and Manning, 2008; Yamada and Matsumoto, 2003; Johansson and Nugues, 2007].

• Transformation is **deterministic**.

• **Hand-written rules** — Need to make various **decisions**. [Xia and Palmer, 2001; Xia et al., 2009; Collins et al., 1999]

• **Our approach** — **data-driven** algorithm using the **structured prediction** framework

• Transformation is **ambiguous**.

## Algorithm

### Rules:

For all  $h, m \in \mathcal{R}(h)$ , rule  $A \rightarrow \beta_1^* \beta_2$ ,  
and  $i \in \{m'_{\leftarrow} : m' \in \mathcal{L}(h)\} \cup \{h\}$ ,

$$\frac{\langle\langle i, m_{\leftarrow} - 1 \rangle, h, \beta_1 \rangle \quad \langle\langle m_{\leftarrow}, m_{\rightarrow} \rangle, m, \beta_2 \rangle}{\langle\langle i, m_{\rightarrow} \rangle, h, A \rangle}$$

For all  $h, m \in \mathcal{L}(h)$ , rule  $A \rightarrow \beta_1 \beta_2^*$ ,  
and  $j \in \{m'_{\rightarrow} : m' \in \mathcal{R}(h)\} \cup \{h\}$ ,

$$\frac{\langle\langle m_{\leftarrow}, m_{\rightarrow} \rangle, m, \beta_1 \rangle \quad \langle\langle m_{\rightarrow} + 1, j \rangle, h, \beta_2 \rangle}{\langle\langle m_{\leftarrow}, j \rangle, h, A \rangle}$$

### Premise:

$$\langle\langle i, i \rangle, i, A \rangle \quad \forall i \in \{1 \dots n\}, A \in \mathcal{N}$$

### Goal:

$$\langle\langle 1, n \rangle, m, r \rangle \text{ for any } m \in \mathcal{R}(0)$$

• Note: Binarization of the phrase structure tree should be done **with respect to the head**, so that each use of a rule implies a dependency arc.

## Learning

$$\min_{\theta} \sum_{i=1}^D \ell(x^i, d^i, y^i, \theta) + \lambda \|\theta\|_1$$

$$\ell(x, d, y, \theta) = -s(y) + \max_{y' \in \mathcal{Y}(x, d)} (s(y') + \Delta(y, y'))$$

For a production $\frac{\langle\langle i, k \rangle, m, \beta_1 \rangle \quad \langle\langle k+1, j \rangle, h, \beta_2 \rangle}{\langle\langle i, j \rangle, h, A \rangle}$	
<b>Nonterm Features</b>	<b>Rule Features</b>
$(A, \beta_1)$ $(A, \beta_1, \text{tag}(m))$	(rule)
$(A, \beta_2)$ $(A, \beta_2, \text{tag}(h))$	(rule, $x_h, \text{tag}(m)$ )
<b>Span Features</b>	(rule, $\text{tag}(h), x_m$ )
$(\text{rule}, x_i)$ $(\text{rule}, x_{i-1})$	(rule, $\text{tag}(h), \text{tag}(m)$ )
$(\text{rule}, x_j)$ $(\text{rule}, x_{j+1})$	(rule, $x_h$ )
$(\text{rule}, x_k)$ $(\text{rule}, x_{k+1})$	(rule, $\text{tag}(h)$ )
$(\text{rule}, \text{bin}(j-i))$	(rule, $x_m$ )
	(rule, $\text{tag}(m)$ )

## Objective

• Parameters are estimated using a structural support vector machine [Taskar et al., 2004]

•  $\Delta(y, y')$  is a hamming loss where  $\mathcal{Y}$  is an indicator for production rules firing over pairs of adjacent spans

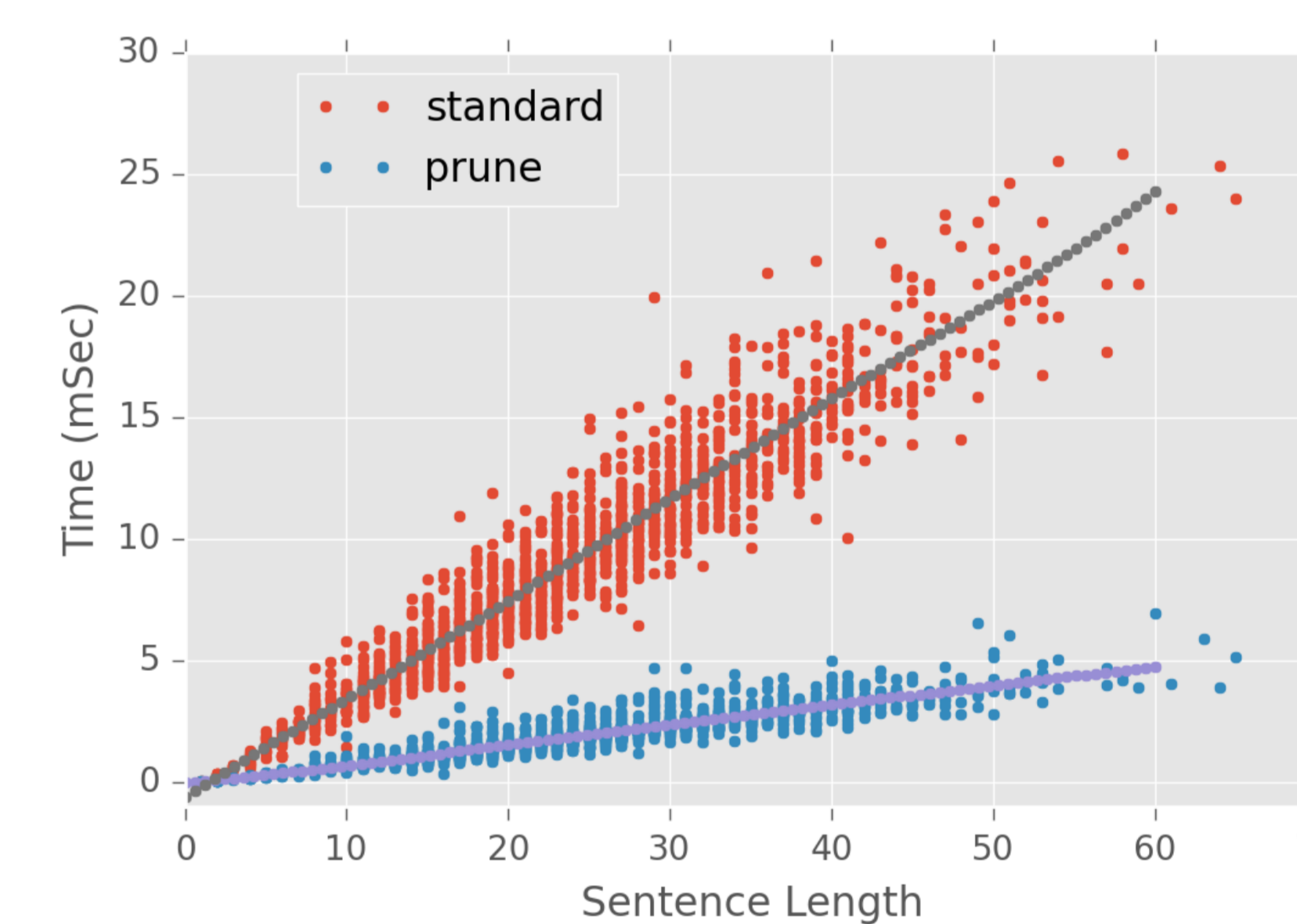
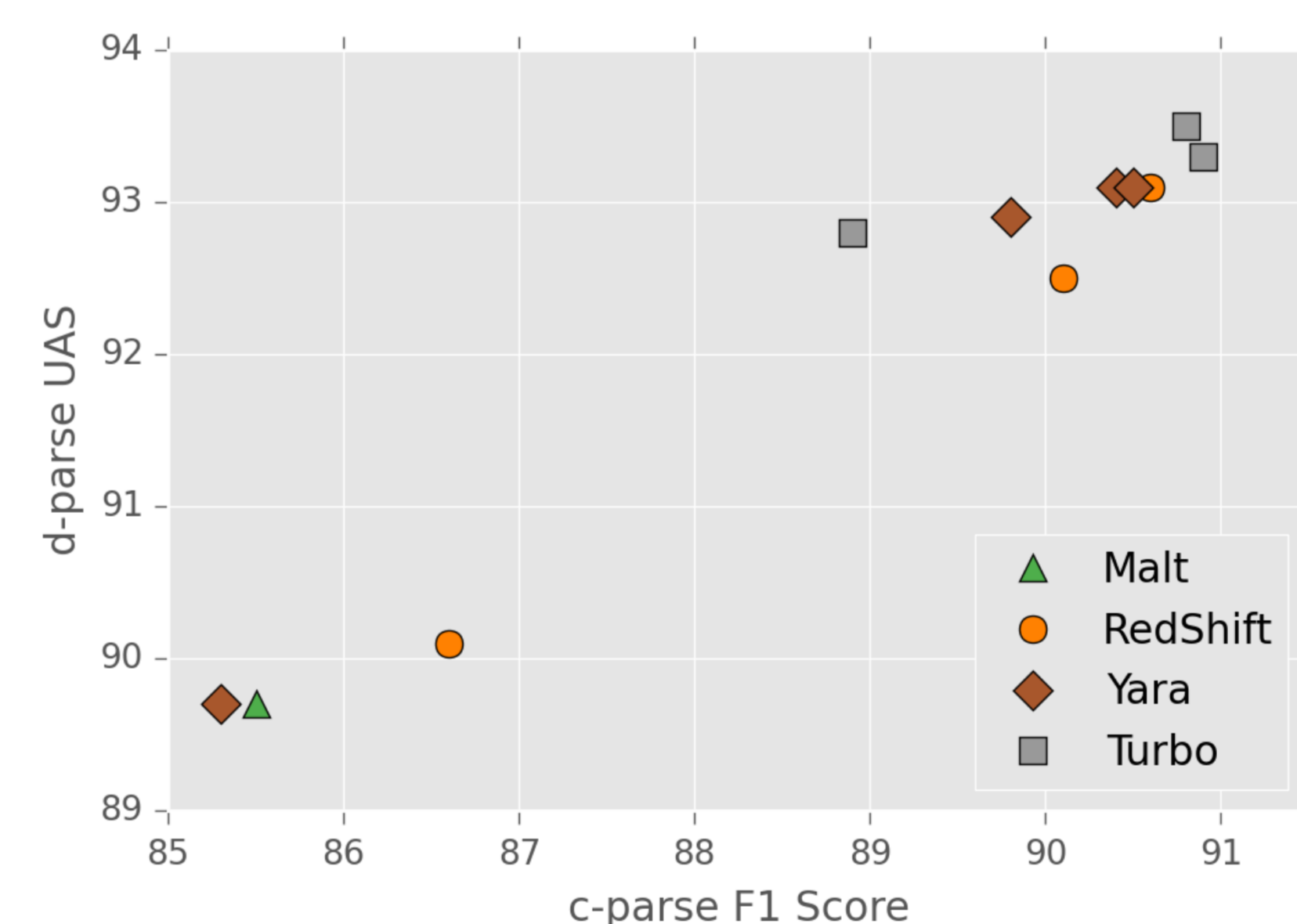
• optimized using AdaGrad of [Duchi et al., 2011]

## Features

• **arc-factored features** [McDonald (2006)]

• **span features** used in the X-bar-style parser of [Hall et al. (2014)]

## Experiments



• The effect of d-parsing accuracy (PTB §22) on **PAD**

• The runtime of our transforming algorithm

PTB 23		
Model	F1	Sent./s.
Charniak (2000)	89.5	-
Stanford PCFG (2003)	85.5	5.3
Petrov (2007)	90.1	8.6
Zhu (2013)	90.3	39.0
Carreras (008)	91.1	-
CJ Reranking (2005)	91.5	4.3
Stanford RNN (2013)	90.0	2.8
<b>PAD</b>	90.6	34.3
<b>PAD (Pruned)</b>	90.5	58.6

CTB 5	
Model	F1
Charniak (2000)	80.8
Bikel (2004)	80.6
Petrov (2007)	83.3
Zhu (2013)	83.2
<b>PAD</b>	82.4

• Accuracy and Speed on the English Penn Treebank (PTB) and the Chinese Penn Treebank (CTB)