

# Parsing Dependencies

## Abstract

We present a new algorithm for transforming dependency parse trees into phrase-structure parse trees. We cast the problem as structured prediction and learn a statistical model. Our algorithm is faster than traditional phrase-structure parsing and achieves accuracy near to the state of the art on English and Chinese benchmarks.

## 1 Introduction

Natural language parsers typically produce phrase-structure (or constituent) trees or dependency trees. These representations capture some of the same syntactic phenomena, and the two can be produced jointly (Carreras et al., 2008; Rush et al., 2010). Yet it appears to be completely unpredictable which will be preferred by a particular application. Both continue to receive the attention of parsing researchers.

Further, it appears to be a historical accident that phrase-structure syntax was used in annotating the Penn Treebank, and that English dependency annotations are largely derived through mechanical, rule-based transformations (reviewed in Section 2). Indeed, despite extensive work on direct-to-dependency parsing algorithms (which we call *d-parsing*), the most accurate dependency parsers for English still involve phrase-structure parsing (which we call *c-parsing*) followed by rule-based extraction of dependencies (Kong and Smith, 2014).

What if dependency annotations had come first? Because d-parsers are generally much faster than c-parsers, we consider an alternate pipeline in Sec-

tion 3: d-parse first, then transform the dependency representation into a phrase-structure tree constrained to be consistent with the dependency parse. This idea was explored by Xia and Palmer (2001) and Xia et al. (2009) using hand-written rules. Instead, we present a data-driven algorithm in Section 4 using the structured prediction framework. The approach can be understood as a pipeline, or as a specially-trained coarse-to-fine decoding algorithm where a d-parser provides “coarse” structure and the second stage refines it (Petrov and Klein, 2007).

Our lexicalized phrase-structure parser is asymptotically faster than parsing with a lexicalized context-free grammar:  $O(n^2)$  plus d-parsing, vs.  $O(n^5)$  worst case runtime in sentence length  $n$ , with the same grammar constant. Experiments in Section 7 show that with simple pruning, our approach achieves linear observable runtime, and that accuracy similar to state-of-the-art phrase-structure parsers without reranking or semisupervised training.

## 2 Background

We begin with the conventional development by first introducing c-parsing and then defining d-parses through a mechanical conversion using head-rules. In the next section, we flip this development and consider the reverse transformation.

### 2.1 CFG Parsing

The phrase-structure trees annotated in the Penn Treebank are derivation trees from a context-free

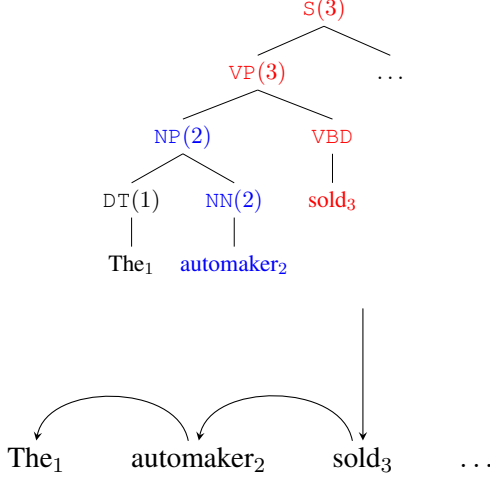


Figure 1: Illustration of c-parse to d-parse conversion with head rules  $\{VP \rightarrow NP VBD^*, NP \rightarrow DT NP^*, \dots\}$ . The c-parse is an ordered tree with fringe  $x_1, \dots, x_n$ . Each vertex is annotated with a non-terminal tag and a derived head index. The blue and red vertices have the words `automaker2` and `sold3` as heads respectively. The vertex  $VP_3$  implies that `automaker2` is a dependent of `sold3`, and that  $d_2 = 3$  in the d-parse.

grammar. Define a binarized<sup>1</sup> context-free grammar (CFG) as a 4-tuple  $(\mathcal{N}, \mathcal{R}, \mathcal{T}, r)$  where:

- $\mathcal{N}$ ; a set of nonterminal symbols, e.g. NP, VP,
- $\mathcal{T}$ ; a set of terminal symbols, consisting of the words in the language,
- $\mathcal{R}$ ; a set of binarized rule productions of the form  $A \rightarrow \beta_1 \beta_2$
- $r \in \mathcal{N}$ ; a distinguished root symbol.

Given an input sentence  $x_1, \dots, x_n$  of terminal symbols from  $\mathcal{T}$ , define the set of c-parses for the sentence as  $\mathcal{Y}(x)$ . This set consists of all binary ordered trees with fringe  $x_1, \dots, x_n$ , internal nodes labeled from  $\mathcal{N}$ , all tree productions  $A \rightarrow \beta$  consisting of members of  $\mathcal{R}$ , and root label  $r$ .

For a c-parse  $y \in \mathcal{Y}(x)$ , we further associate a span  $\langle i, j \rangle$  with each vertex in the tree. This specifies the sequence  $\{x_i, \dots, x_j\}$  of the sentence covered by this vertex.

## 2.2 Dependency Parsing

Dependency trees provide an alternative, and in some sense simpler, representation of sentence

<sup>1</sup>For notational simplicity we ignore unary rules for this section.

structure. These d-parses can be derived through mechanical transformation from context-free trees. There are several popular transformation in wide use; each provides a different representations of a sentence’s structure (Collins, 2003; De Marneffe and Manning, 2008; Yamada and Matsumoto, 2003). In this work we consider transformations that are defined through local transformations known as head rules.

For a binary CFG, define a collection of head rules as a function  $\mathcal{H} : \mathcal{R} \mapsto \{L, R\}$  mapping each CFG rule to its head preference for its left- or right-child. We use the notation  $A \rightarrow \beta_1^* \beta_2$  and  $A \rightarrow \beta_1 \beta_2^*$  to indicate a left- or right-headed rule respectively.

For each vertex  $v$  of a c-parse let the head of the vertex  $h(v)$  be defined recursively,

1. If the vertex is leaf  $x_m$ , then  $h(v) = m$ .
2. Otherwise,  $h(v)$  matches the head child where  $A \rightarrow \beta_1^* \beta_2$  or  $A \rightarrow \beta_1 \beta_2^*$  is the rule production at this vertex.

The head rules can be used to map a c-parse to a dependency tree. Define a sentence’s dependencies as a sequence  $d_1, \dots, d_n$  where word  $m$  is a dependent of word  $d_m \in \{0, \dots, n\}$  and 0 is a special pseudo-root symbol. Let vertex  $v$  be the first parent vertex of  $x_m$  not headed by  $m$ , i.e.  $h(v) \neq m$ . If this vertex exists, then  $d_m = h(v)$ , otherwise  $d_m = 0$ . Figure ?? illustrates the conversion.

These dependencies  $d_1, \dots, d_n$  can be viewed as a directed tree with arcs  $(d_m, m)$  for all  $m$ . This tree differs from the original c-parse since the words are directly connected. However we can relate the two trees through their spans. Define a dependency span  $\langle m_{\leftarrow}, m_{\rightarrow} \rangle$  as the sequence of words reachable from word  $m$ . By construction this span is the same as the span  $\langle i, j \rangle$  of the top vertex  $v$  with  $h(v) = m$ .

## 3 Parsing Dependencies

Now let us flip this setup. In recent years there has been significant progress in developing efficient direct-to-dependency parsers. These parsers can be trained only on dependency annotations and do not require full phrase-structure trees.<sup>2</sup> In many

<sup>2</sup>For English these parsers are still often trained on trees converted from c-parses; however for other languages d-parse-only treebanks are common.

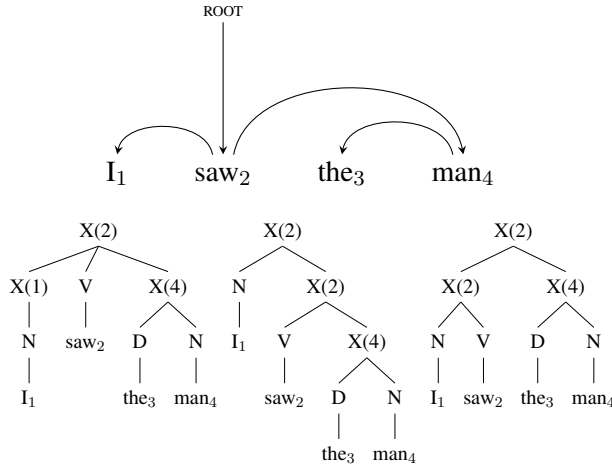


Figure 2: [Adapted from (Collins et al., 1999)] Several simple c-parses that are valid conversion from the original d-parse. The goal of this work is to select the best parse from this set.

ways this setup is preferable since the parser can be trained on the specific dependencies required for a downstream task. However, there are many applications that require or prefer full c-parses as input, and so c-parsers are still widely-used.

With these applications in mind, we consider the problem of converting a d-parser to a c-parser by converting fixed d-parses into c-parses. Since this problem is more challenging than its inverse, we use a structured prediction setup: learn a function to score possible c-parse conversions, and then generate the highest-scoring c-parse possible.

We start by considering the search problem with a given scoring function.

### 3.1 Parsing Algorithm

First consider the standard problem of predicting the best c-parse under a CFG with head rules. Assume that we are given a binary CFG defining a set of valid c-parses  $\mathcal{Y}(x)$ . The parsing problem is to find the highest-scoring parse in this set, i.e.

$$\hat{y} \leftarrow \arg \max_{y \in \mathcal{Y}(x)} s(y; x)$$

where  $s$  is a scoring function that factors over rule productions.

This problem is known as lexicalized context-free parsing. It can be solved using the standard lexicalized CKY parsing algorithm. This algorithm is defined by the productions in Figure ?? . The productions are of the form

$$\frac{(\langle i, k \rangle, m, \beta_1) \quad (\langle k+1, j \rangle, h, \beta_2)}{(\langle i, j \rangle, h, A)}$$

for all rules  $A \rightarrow \beta_1 \beta_2^* \in \mathcal{R}$  and spans  $i \leq k < j$ . This particular production indicates that rule  $A \rightarrow \beta_1 \beta_2^*$  was applied at a vertex covering  $\langle i, j \rangle$  to produce two vertices covering  $\langle i, k \rangle$  and  $\langle k+1, j \rangle$ , and that the new head is index  $h$  which is modified by index  $m$ .

The highest-scoring parse can be found by bottom-up dynamic programming over these productions. The standard lexicalized CKY algorithm requires  $O(n^5 |\mathcal{R}|)$  time. While simple to implement, this algorithm is intractable to run without heavy pruning or further assumptions on the scoring function.

Now, consider the same setup, but constrained to parses that will be converted to a given d-parse,  $d_1, \dots, d_n$ . Define this set as  $\mathcal{Y}(x, d)$ .

$$\arg \max_{y \in \mathcal{Y}(x, d)} s(y; x, d)$$

This new problem has a nice property. For any word  $x_m$  the span  $\langle i, j \rangle$  of the highest  $v$  with  $h(v) = n$  is the same as span  $\langle m_{\leftarrow}, m_{\rightarrow} \rangle$  in the dependency parse. Since we have the dependency parse these spans can be efficiently pre-computed.

This property greatly limits the search space of the parsing problem. Instead of searching over all possible spans  $\langle i, j \rangle$  of each modifier, we simply pre-compute the span  $\langle m_{\leftarrow}, m_{\rightarrow} \rangle$ . Figure ?? shows the new algorithm. While the productions are the same as the original algorithm, there are many fewer of them. There is one production for each index and each possible modifier, leading to an algorithm with  $O(n^2 |\mathcal{R}|)$  running time.

### 3.2 Pruning

In addition to constraining the number of c-parses possible, the d-parse also provides valuable information about the labeling and structure of the c-parse. We can use this information to provide further pruning of the search space. We experimented with two heuristic pruning methods based on this information

- We observe that in training the part-of-speech of the head word  $x_h$  greatly limits the possible rules  $A \rightarrow \beta_1 \beta_2$  available. To exploit this

**Premise:**

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

**Rules:**

For  $i \leq h \leq k < m \leq j$ , and rule  $A \rightarrow \beta_1^* \beta_2$ ,

$$\frac{(\langle i, k \rangle, h, \beta_1) \quad (\langle k+1, j \rangle, m, \beta_2)}{(\langle i, j \rangle, h, A)}$$

For  $i \leq m \leq k < h \leq j$ , rule  $A \rightarrow \beta_1 \beta_2^*$ ,

$$\frac{(\langle i, k \rangle, m, \beta_1) \quad (\langle k+1, j \rangle, h, \beta_2)}{(\langle i, j \rangle, h, A)}$$

**Goal:**

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

**Premise:**

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

**Rules:**

For all  $i < m, h = d_m$  and rule  $A \rightarrow \beta_1^* \beta_2$ ,

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

For all  $m < j, h = d_m$  and rule  $A \rightarrow \beta_1 \beta_2^*$ ,

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

**Goal:**

$$(\langle 1, n \rangle, m, r) \text{ for any } m \text{ s.t. } d_m = 0$$

Figure 3: The two algorithms written as deductive parsers. Starting from the *premise*, any valid application of *rules* that leads to a *goal* is a valid parse. (a) Lexicalized CKY algorithm for CFG parsing with head rules. For this algorithm there are  $O(n^5|\mathcal{R}|)$  rules where  $n$  is the length of the sentence. (b) The constrained CKY parsing algorithm for  $\mathcal{Y}(x, d)$ . The algorithm is nearly identical except that many of the free indices are now fixed to the dependency parse. Finding the optimal parse is now  $O(n^2|\mathcal{R}|)$ .

property we build tables  $\mathcal{R}_t$  for each part-of-speech tag  $t$  and limit the search to rules seen for the current head tag. This reduces the effect of  $|\mathcal{R}|$  at runtime.

- A word  $x_h$  with  $L$  dependents to its left and  $R$  dependents to its right may at any point combine with a left or right dependent. There are therefore  $L \times R$  ( $O(n^2)$ ) possible orderings of these dependents.

However, we can often predict with very high accuracy which side will come next. To make this prediction, we estimate a distribution

$$p(\text{side} | \text{tag}(x_h), \text{tag}(x_l), \text{tag}(x_r))$$

where  $x_l$  is the next left dependent,  $x_r$  is the next right dependent, and  $\text{tag}()$  is the part-of-speech tag of the word. This reduces the effect of the  $n^2$  at runtime.

Table ?? shows a comparison of these pruning methods. In addition to speed, an important question whether it is even possible to recover c-parses under these constraints. These experiments show that there is a small drop from pruning, but that there is

Model	Comp.	Speed	Oracle
Full[< 20]	$O(n^5 \mathcal{R} )$	0.25	100.0
Dep	$O(n^2 \mathcal{R} )$	76.2	92.9
Dep+Prune1	$O(n^2 \mathcal{R} )$	375.9	92.8
Dep+Prune	$O(n^2 \mathcal{R} )$	487.3	92.6

Table 1: Comparison of three parsing setups: Full[< 20] is the standard full lexicalized grammar limited to sentence of length less than 20 words, LCFG(dep) is limited to the dependency skeleton, and LCFG(prune) is the pruning described in Section 3.2. *Oracle* is the oracle f-score on the development data (described in Section 7). *Speed* is the efficiency on development data in sentences per second.

a sizable drop in oracle accuracy with the d-parse constraints. Still this upper-bound on accuracy is high enough to make it possible to still recover high-accuracy parses. We return to this discussion in Section ??.

### 3.3 Binarization

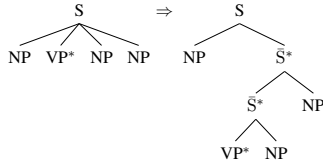
To this point, we have assumed that the CFG and head rules are binary; however the standard treebank grammars have arbitrarily large rules. In order to apply the algorithm, we need to binarize the grammar and preserve the head relations.

Consider a non-binarized rule of the form  $A \rightarrow \beta_1 \dots \beta_m$  with head rule  $\mathcal{H}(A \rightarrow \beta_1 \dots \beta_m) = k$ . Relative to the head  $\beta_k$  the rule has left-side  $\beta_1 \dots \beta_{k-1}$  and right-side  $\beta_{k+1} \dots \beta_m$ .

We replace this rule with binary rules that consume each side independently as a simple chain. The transformation introduces new rules<sup>3</sup>

- $\bar{A} \rightarrow \beta_i \bar{A}^*$  for  $i \in \{2, \dots, k\}$
- $\bar{A} \rightarrow \bar{A}^* \beta_i$  for  $i \in \{k, \dots, m\}$
- $A \rightarrow \beta_1 \bar{A}^*$

As an example consider the transformation of a rule with four children:



Note that each rule contains at most two nonterminals from the original grammar, so the size of the new binarized rule set is bounded by  $O(|\mathcal{N}|^2)$ . The small size of the grammar further helps the speed of the algorithm.

We experimented with binarization that use horizontal and vertical markovization to include additional context of the tree as in unlexicalized approaches (Klein and Manning, 2003), but we found that this increased the size of the grammar without leading to improvements in accuracy.

## 4 Structured Prediction

While the parser only requires a d-parse at prediction time, it uses annotated c-parses to learn the transformation. Define this scoring function  $s$  as

$$s(y; x, d, \theta) = \theta^\top f(x, d, y)$$

where  $\theta$  is a parameter vector and  $f(x, d, y)$  is a feature function that maps parse productions to sparse feature vectors.

To learn the parameter for scoring possible c-parses we use a standard structured prediction setup. In this section, we describe the features of the model and then the learning algorithm.

<sup>3</sup>These rules are slightly modified when  $k = 1$

### 4.1 Features

In theory, the scoring function for this system could be directly adapted from existing c-parsers. However, the existing structure of the dependency parse limits the number of decisions that need to be made, and allows for a smaller set of features.

We model our features after two bare-bones parsing systems. The first set is the basic arc-factored features used by McDonald (2006). These features include combinations of:

- rule and top non-terminal
- modifier word and part-of-speech
- head word and part-of-speech

The second set of features is modeled after the span features described in the X-Bar style parser of Hall et al. (2014). These include conjunction of the rule with:

- first and last word of current span.
- preceding and following word of current span
- adjacent words at split of current span
- length of the span

The full feature set is shown in Figure ???. After training, there are a total of around 2 million non-zero features. For efficiency, we use lossy feature hashing. We found this had no impact on parsing accuracy but made the parsing significantly faster.

### 4.2 Training

The parameters  $\theta$  are estimated using a structured SVM setup. Given a set of gold-annotated parse examples,  $(x^1, y^1), \dots, (x^D, y^D)$ , and paired dependency structures  $d^1 \dots d^D$  induced from the head rules, we estimate the parameters to minimize the regularized empirical risk

$$\min_{\theta} \sum_{i=1}^D \max\{0, \ell(x^i, d^i, y^i, \theta)\} + \frac{\lambda}{2} \|\theta\|_1$$

where we define  $\ell$  as

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

For a production  $\frac{(\langle i, k \rangle, m, \beta_1) \quad (\langle k + 1, j \rangle, h, \beta_2)}{(\langle i, j \rangle, h, A)}$

Nonterm Features	Rule Features
$(A, \beta_1)$	(rule)
$(A, \beta_2)$	$(\text{rule}, x_h, \text{tag}(m))$
$(A, \beta_1, \text{tag}(m))$	$(\text{rule}, \text{tag}(h), x_m)$
$(A, \beta_2, \text{tag}(h))$	$(\text{rule}, \text{tag}(h), \text{tag}(m))$
Span Features	$(\text{rule}, x_h)$
$(\text{rule}, x_i)$	$(\text{rule}, \text{tag}(h))$
$(\text{rule}, x_j)$	$(\text{rule}, x_m)$
$(\text{rule}, x_{i-1})$	$(\text{rule}, \text{tag}(m))$
$(\text{rule}, x_{j+1})$	
$(\text{rule}, x_k)$	
$(\text{rule}, x_{k+1})$	
$(\text{rule}, \text{bin}(j - i))$	

Figure 4: The feature templates used in the function  $f(x, d, y)$ . For the span features, the symbol rule is expanded into both  $A \rightarrow B \ C$  and  $A$ . The function  $\text{tag}(i)$  gives the part-of-speech tag of word  $x_i$ . The function  $\text{bin}(i)$  partitions a span length into one of 10 bins.

and where  $\Delta$  is a problem specific cost-function. In experiments, we use a hamming loss  $\Delta(y, y') = ||y - y'||$  where  $y$  is an indicator of rule productions.

The objective is optimized using Adagrad (Duchi et al., 2011). The gradient calculation requires computing a loss-augmented argmax for each training example which is done using the algorithm of Figure ??.

## 5 Related Work

The problem of converting dependency to phrase-structured trees has been studied previously from the perspective of building multi-representational treebanks. Xia and Palmer (2001) and Xia et al. (2009) develop a rule-based system for the converting human-annotated dependency parses. Our work differs in that we learn a data-driven structured prediction model that is also able to handle automatically predicted input. **[mention the table in the exp section where the comparison number show our approach better here? -lpk]**

There has been successful work combining dependency and phrase-structure parsing. Carreras et al. (2008) build a high-accuracy parser that uses a dependency parsing model both for pruning and within a richer lexicalized parser. Similarly Rush

et al. (2010) use dual decomposition to combine a dependency parser with a simple phrase-structure model. This work differs in that we remove the dependency parsing model entirely. By starting with a complete dependency parse structure we can build a much faster and smaller model that is still able to generate highly phrase-structure parses.

**[only fixing the dependency structure entirely before parsing sounds more like a step backward than a step further... i am not sure if this is a good way to say that... i think there is a important difference here, which is, when combining the phrase-structure and dep parsing, when ppl do this, the binarization does not imply dep arcs, so even you can push info from dep to phrase structure parsing, you can't directly pruning so much like we do here...right? -lpk]**

Finally there have also been several papers that use ideas from dependency parsing to simplify and speed up phrase-structure prediction. Zhu et al. (2013) build a high-accuracy phrase-structure parser using a transition-based system. Hall et al. (2014) use a stripped down parser based on a simple X-bar grammar and a small set of lexicalized features.

## 6 Methods

We ran a series of experiments to assess the accuracy, efficiency, and applicability of the parser to several tasks. These experiments use the following setup.

For English experiments we use the standard Penn Treebank (PTB) experimental setup (Marcus et al., 1993). Training is done on section 2-21, development on section 22, and test of section 23.

For Chinese experiments, we use version 5.1 of the Penn Chinese Treebank 5.1 (CTB) (Xue et al., 2005). We followed previous work and used 001-270 and 440-1151 for training, articles 301-325 as development, and articles 271-300 as test.

Part-of-speech tagging is performed for all models using TurboTagger (Martins et al., 2013). Unless otherwise noted, the input d-parsing is done using RedShift implementation<sup>4</sup> of the Zhang-Nivre Parser (Zhang and Nivre, 2011).

Prior to training the d-parser, the train sections are automatically tagged using 10-fold jackknifing.

<sup>4</sup><https://github.com/syllog1sm/redshift>

The d-parser is trained to produce Collins head rules (Collins, 2003).

For our model the “gold” d-parses  $d^1, \dots, d^D$  used at training come from the gold c-parse converted with the Collins head rules<sup>5</sup>. Since the trained model is independent of the d-parser, for all experiments with other d-parsers we used the same model.

Evaluation for phrase-structure parses is performed using the `evalb`<sup>6</sup> script using the standard setup. We report F1-Score as well as recall and precision. For dependency parsing we use unlabeled accuracy score (UAS).

We implemented the grammar binarization, head rules, and pruning tables in Python, and the parser, features, and training in C++. The core run-time decoding algorithm is self contained and requires less than 400 lines of code. Both are publicly available.<sup>7</sup> Experiments are performed on a Lenovo ThinkCentre desktop computer with 32GB of memory and Core i7-3770 3.4GHz 8M cache CPU.

## 7 Experiments

We ran experiments to assess the accuracy of the method, its run-time efficiency, the amount of phrase-structure data required, and the effect of dependency accuracy.

**Accuracy** Our first experiments, shown in Table ??, give the accuracy of the phrase-structure trees produced by the parser. For these experiments we treat our system and the Zhang-Nivre parser as an independently trained, but complete end-to-end c-parser. This setup is limited in the that the parser cannot change the dependencies given to it. Despite this limitation, the results show that the parser is comparable in accuracy to many widely-used systems, and is significantly faster. The only parser competitive in both speed and accuracy is the Zhu[13], a fast shift-reduce phrase-structure parser.

...

**Dependencies** Our next set of experiments look at the effect of different input dependency trees. We

<sup>5</sup>We experimented with using jackknifed dependency parses  $d'$  at training time with oracle tree structures, i.e.  $\arg \min_{y' \in \mathcal{Y}(x, d')} \Delta(y, y')$ , but found that this did not improve performance.

<sup>6</sup><http://nlp.cs.nyu.edu/evalb/>

<sup>7</sup>Withheld for review

Model	PTB	
	F1	Speed
Charniak[00]	89.5	-
Stanford PCFG[03]	85.5	5.3
Petrov[07]	90.1	8.6
Zhu[13]	90.3	39.0
Carreras[08]	91.1	-
CJ Reranking[05]	91.5	4.3
Stanford RNN[13]	90.0	2.8
PARPAR	90.4	58.6

Table 2: Accuracy results on the Penn Treebank and Chinese Treebank datasets. Comparisons are to state-of-the-art non-reranking phrase-structure parsers including: Charniak[00] (Charniak, 2000), Stanford PCFG[03] (Klein and Manning, 2003), Petrov[07] (Petrov and Klein, 2007), Carreras[08] (Carreras et al., 2008), Zhu[13] (Zhu et al., 2013), as well as semi-supervised parsers including : Stanford Recursive neural network[13] (Socher et al., 2013) and CJ Reranking (Charniak and Johnson, 2005).

ran the same trained converter on seven different dependency inputs of varying quality measured by unlabeled accuracy score(UAS) and speed.

The experiments are shown in Table ??. The corresponding figure shows that there is a direct correlation between the UAS of the parser and the final F1, but that small changes in dependency accuracy can have rather large effects. For instance MaltParser scores 89.7 in UAS but the conversion produces only a 85.5 F1 score.

To upper bound of the accuracy of this system we consider an oracle version of the parser. For a gold parse  $y$  and predicted dependencies  $\hat{d}$ , define the oracle parse  $y'$  as

$$y' = \arg \min_{y' \in \mathcal{Y}(x, \hat{d})} \Delta(y, y')$$

Table ?? shows that the change in oracle score tracks the F1 scores seen for each parser.

**Speed** In Section 3 we considered the theoretical complexity of the parsing model, and gave the main speed results in Table ??. Despite having a quadratic theoretical complexity the practical run-time was quite fast.

We therefore considered the empirical complexity of the model by measuring the time spent of individual sentences. Figure 5 shows a graph of the speed

Model	UAS	F1	Speed	Oracle
MALTPARSER	89.7	85.5	240.7	87.8
RS-K1	90.1	86.6	233.9	87.6
RS-K4	92.5	90.1	151.3	91.5
RS-K16	93.1	90.6	58.6	92.4
TP-BASIC	92.8	88.9	132.8	90.8
TP-STANDARD	93.3	90.9	27.2	92.6
TP-FULL	93.5	90.8	13.2	92.9

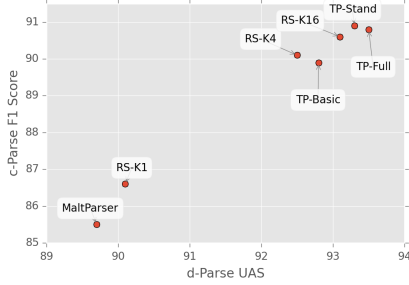


Table 3: Comparison of the effect of downstream dependency prediction. Experiments are run on the development section with different input dependencies. *Oracle* is the oracle F1 on the development data. *Speed* is the efficiency of the *full-system* (c-parse+d-parse) in sentences per second. Inputs include three versions of TurboParser (Martins et al., 2013) (Basic, Standard Full), MaltParser (Nivre et al., 2006), and the RedShift implementation of the Zhang-Nivre parser (Zhang and Nivre, 2011) with beam size  $k \in \{1, 4, 16\}$ .

of sentences of varying length for both the full algorithm and for the algorithm with pruning. In both cases the run-time is clearly linear in practice.

**Conversion** Previous work on this problem has looked at converting dependency trees to phrase-structure trees using linguistic rules (Xia and Palmer, 2001; Xia et al., 2009). This work is targeted towards the development of treebanks, particularly converting dependency treebanks to phrase-structure treebanks. For this application, it is useful to convert gold trees as opposed to predicted trees.

To compare to this work, we train our parser with gold tags and run on gold dependency trees in development. Table 5 give the results for this task.

**Analysis** Finally we consider an internal error analysis of the parser. For this analysis, we group each binary rule production selected by the parser by three properties: is its dependency  $(h, m)$  correct?, is its span  $\langle i, j \rangle$  correct? is its split  $k$  correct? The first is fully determined by the input d-parse, the second are partially determined by the parser.

Table ?? gives the breakdown. As expected the

Dep $(h, m)$	Class		Results	
	Span $\langle i, j \rangle$	Split $k$	Total %	Acc $A$
+	+	+	86.1	97.9
-	+	+	1.0	69.3
+	+	-	2.1	83.3
-	+	-	1.3	85.9
+	-	-	4.5	0.0
-	-	-	4.7	0.0

Table 4: Error analysis of binary CFG rules. Rules used are split into classes based on correct (+) identification of dependency  $(h, m)$ , span  $\langle i, j \rangle$ , and split  $k$ . *Total* is the relative size of the class. *Acc* is accuracy of span nonterminal identification.

Model	Dev		
	Prec	Rec	F1
Xia[09]	88.1	90.7	89.4
PARPAR(Sec19)	95.9	95.9	95.9
PARPAR	97.5	97.8	97.7

Table 5: Comparison with the rule-based system of Xia et al. (2009). Results are from PTB development section 22 using gold tags and gold dependencies. Xia[09] report results from training on only on Section 19, but include a note that further data had little effect. For comparison we report result on complete training as well as just Sec. 19.

conversion is almost always correct when it has this information. Many of the errors come from cases where the dependency was fully incorrect, or there was a span issue either from the dependency parser or from the conversion itself.

## 8 Conclusion

With recent advances in statistical dependency parsing, state-of-the-art parsers have reached the comparable dependency accuracy as the best phrase structure parsers. However, these parser cannot be directly used in applications that require phrase-structure prediction. In this work we have described a simple parsing algorithm and structured prediction system for this comparison, and show that it can produce phrase-structure parses at comparable accuracy to state-of-the-art system.

## References

Xavier Carreras, Michael Collins, and Terry Koo. 2008. Tag, dynamic programming, and the perceptron for efficient, feature-rich parsing. In *Proceedings of the*



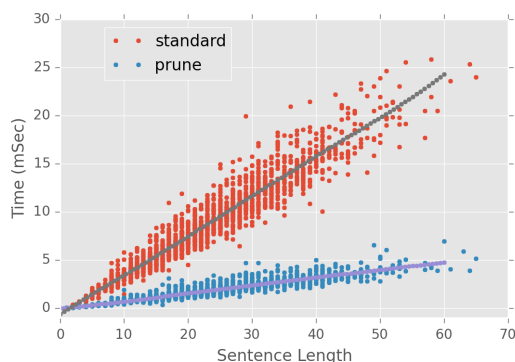
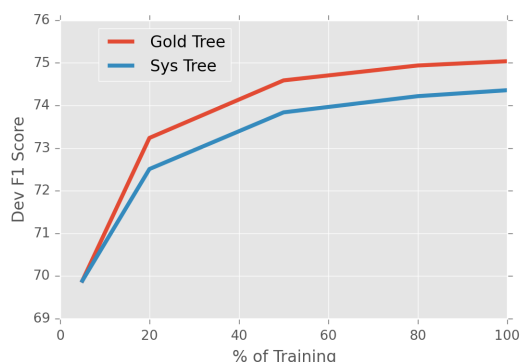


Figure 5: Empirical complexity of the parser on sentences of varying length, with and without pruning. Despite a worst-case quadratic complexity, in practice the runtime is linear. Pruning significantly lowers the slope of the run-time.



*Twelfth Conference on Computational Natural Language Learning*, pages 9–16. Association for Computational Linguistics.

Eugene Charniak and Mark Johnson. 2005. Coarse-to-fine n-best parsing and maxent discriminative reranking. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 173–180. Association for Computational Linguistics.

Eugene Charniak. 2000. A maximum-entropy-inspired parser. In *Proceedings of the 1st North American chapter of the Association for Computational Linguistics conference*, pages 132–139. Association for Computational Linguistics.

Michael Collins, Lance Ramshaw, Jan Hajič, and Christoph Tillmann. 1999. A statistical parser for czech. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 505–512. Association for Computational Linguistics.

Michael Collins. 2003. Head-driven statistical models for natural language parsing. *Computational linguistics*, 29(4):589–637.

Marie-Catherine De Marneffe and Christopher D Manning. 2008. The stanford typed dependencies representation. In *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, pages 1–8. Association for Computational Linguistics.

John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159.

David Hall, Greg Durrett, and Dan Klein. 2014. Less grammar, more features. In *ACL*.

Dan Klein and Christopher D Manning. 2003. Accurate unlexicalized parsing. In *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics-Volume 1*, pages 423–430. Association for Computational Linguistics.

Lingpeng Kong and Noah A Smith. 2014. An empirical comparison of parsing methods for stanford dependencies. *arXiv preprint arXiv:1404.4314*.

Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330.

André FT Martins, Miguel Almeida, and Noah A Smith. 2013. Turning on the turbo: Fast third-order non-projective turbo parsers. In *ACL (2)*, pages 617–622.

Ryan McDonald. 2006. *Discriminative learning and spanning tree algorithms for dependency parsing*. Ph.D. thesis, University of Pennsylvania.

- Joakim Nivre, Johan Hall, and Jens Nilsson. 2006. Malt-parser: A data-driven parser-generator for dependency parsing. In *Proceedings of LREC*, volume 6, pages 2216–2219.
- Slav Petrov and Dan Klein. 2007. Improved inference for unlexicalized parsing. In *HLT-NAACL*, pages 404–411. Citeseer.
- Alexander M Rush, David Sontag, Michael Collins, and Tommi Jaakkola. 2010. On dual decomposition and linear programming relaxations for natural language processing. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 1–11. Association for Computational Linguistics.
- Richard Socher, John Bauer, Christopher D Manning, and Andrew Y Ng. 2013. Parsing with compositional vector grammars. In *In Proceedings of the ACL conference*.
- Fei Xia and Martha Palmer. 2001. Converting dependency structures to phrase structures. In *Proceedings of the first international conference on Human language technology research*, pages 1–5. Association for Computational Linguistics.
- Fei Xia, Owen Rambow, Rajesh Bhatt, Martha Palmer, and Dipti Misra Sharma. 2009. Towards a multi-representational treebank. In *The 7th International Workshop on Treebanks and Linguistic Theories. Groningen, Netherlands*.
- Naiwen Xue, Fei Xia, Fu-Dong Chiou, and Martha Palmer. 2005. The penn chinese treebank: Phrase structure annotation of a large corpus. *Natural language engineering*, 11(02):207–238.
- Hiroyasu Yamada and Yuji Matsumoto. 2003. Statistical dependency analysis with support vector machines. In *Proceedings of IWPT*, volume 3.
- Yue Zhang and Joakim Nivre. 2011. Transition-based dependency parsing with rich non-local features. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, pages 188–193. Association for Computational Linguistics.
- Muhua Zhu, Yue Zhang, Wenliang Chen, Min Zhang, and Jingbo Zhu. 2013. Fast and accurate shift-reduce constituent parsing. In *ACL (1)*, pages 434–443.