

# A Fast, Accurate, Non-Projective, Semantically-Enriched Parser

**Stephen Tratz and Eduard Hovy**  
Information Sciences Institute  
University of Southern California  
Marina del Rey, California 90292  
{stratz, hovy}@isi.edu

## Abstract

Dependency parsers are critical components within many NLP systems. However, currently available dependency parsers each exhibit at least one of several weaknesses, including high running time, limited accuracy, vague dependency labels, and lack of non-projectivity support. Furthermore, no commonly used parser provides additional shallow semantic interpretation, such as preposition sense disambiguation and noun compound interpretation. In this paper, we present a new dependency-tree conversion of the Penn Treebank along with its associated fine-grain dependency labels and a fast, accurate parser trained on it. We explain how a non-projective extension to shift-reduce parsing can be incorporated into non-directional easy-first parsing. The parser performs well when evaluated on the standard test section of the Penn Treebank, outperforming several popular open source dependency parsers; it is, to the best of our knowledge, the first dependency parser capable of parsing more than 75 sentences per second at over 93% accuracy.

## 1 Introduction

Parsers are critical components within many natural language processing (NLP) systems, including systems for information extraction, question answering, machine translation, recognition of textual entailment, summarization, and many others. Unfortunately, currently available dependency parsers suffer from at least one of several weaknesses including high running time, limited accuracy, vague dependency labels, and lack of non-projectivity support. Furthermore, few parsers include any sort of

additional semantic interpretation, such as interpretations for prepositions, possessives, or noun compounds.

In this paper, we describe 1) a new dependency conversion (Section 3) of the Penn Treebank (Marcus, et al., 1993) along with the associated dependency label scheme, which is based upon the Stanford parser’s popular scheme (de Marneffe and Manning, 2008), and a fast, accurate dependency parser with non-projectivity support (Section 4) and additional integrated semantic annotation modules for automatic preposition sense disambiguation and noun compound interpretation (Section 5). We show how Nivre’s (2009) swap-based reordering technique for non-projective shift-reduce-style parsing can be integrated into the non-directional easy-first framework of Goldberg and Elhadad (2010) to support non-projectivity, and we report the results of our parsing experiments on the standard test section of the PTB, providing comparisons with several freely available parsers, including Goldberg and Elhadad’s (2010) implementation, MALTPARSER (Nivre et al., 2006), MSTPARSER (McDonald et al., 2005; McDonald and Pereira, 2006), the Charniak (2000) parser, and the Berkeley parser (Petrov et al., 2006; Petrov and Klein, 2007).

The experimental results show that the parser is substantially more accurate than Goldberg and Elhadad’s original implementation, with fairly similar overall speed. Furthermore, the results prove that Stanford-granularity dependency labels can be learned by modern dependency parsing systems when using our Treebank conversion, unlike the Stanford conversion, for which Cer et al. (2010) show that this isn’t the case.

The optional semantic annotation modules also

perform well, with the preposition sense disambiguation module exceeding the accuracy of the previous best reported result for fine-grained preposition sense disambiguation (85.7% vs Hovy et al.'s (2010) 84.8%), the possessives interpretation system achieving over 85% accuracy, and the noun compound interpretation system performing similarly to an earlier version described by Tratz and Hovy (2010) at just over 79% accuracy.

## 2 Background

The NLP community has recently seen a surge of interest in dependency parsing, with several CoNLL shared tasks focusing on it (Buchholz and Marsi, 2006; Nivre et al., 2007). One of the main advantages of dependency parsing is the relative ease with which it can handle non-projectivity<sup>1</sup>. Additionally, since each word is linked directly to its head via a link that, ideally, indicates the syntactic dependency type, there is no difficulty in determining either the syntactic head of a particular word or the syntactic relation type, whereas these issues often arise when dealing with constituent parses<sup>2</sup>.

Unfortunately, most currently available dependency parsers produce relatively vague labels or, in many cases, produce no labels at all. While the Stanford fine-grain dependency scheme (de Marneffe and Manning, 2008) has proven to be popular, recent experiments by Cer et al. (2010) using the Stanford conversion of the Penn Treebank indicate that it is difficult for current dependency parsers to learn. Indeed, the highest scoring parsers trained using the MSTPARSER (McDonald and Pereira, 2006) and MALTPARSER (Nivre et al., 2006) parsing suites achieved only 78.8 and 81.1 labeled attachment  $F_1$ , respectively. This contrasted with the much higher performance obtained using a constituent-to-dependency conversion approach with accurate, but much slower, constituency parsers such as the Charniak and Johnson (2005) and Berkeley (Petrov et al., 2006; Petrov and Klein, 2007) parsers, which achieved 89.1 and 87.9 labeled  $F_1$  scores, respectively.

<sup>1</sup>A tree is non-projective if the sequence of words visited in a left-to-right, depth-first traversal of the sentence's parse tree is different than the actual word order of the sentence.

<sup>2</sup>These latter two issues are not problems for constituent parses with binarized output and functional tags.

Though there are many syntactic parsers that can reconstruct the grammatical structure of a text, there are few, if any, accurate and widely accepted systems that also produce shallow semantic analysis of the text. For example, a parser may indicate that, in the case of 'ice statue', 'ice' modifies 'statue' but will not indicate that 'ice' is the *substance* of the statue. Similarly, a parser will indicate which words a preposition connects but will not give any semantic interpretation (e.g., 'the boy *with* the pirate hat' → *wearing or carrying*, 'wash *with* cold water' → *means*, 'shave *with* the grain' → *in the same direction*). While, in some cases, it may be possible to use the output from a separate system for this purpose, doing so is often difficult in practice due to a wide variety of complications, including programming language differences, alternative data formats, and, sometimes, other parsers.

## 3 Dependency Conversion

### 3.1 Relations and Structure

Most recent English dependency parsers produce one of three sets of dependency types: unlabeled, some variant of the coarse labels used by the CoNLL dependency parsing shared-tasks (Buchholz and Marsi, 2006; Nivre et al., 2007) (e.g., ADV, NMOD, PMOD), or Stanford's dependency labels (de Marneffe and Manning, 2008). Unlabeled dependencies are clearly too impoverished for many tasks. Similarly, the coarse labels of the CoNLL tasks are not very specific; for example, the same relation, NMOD, is used for determiners, adjectives, nouns, participle modifiers, relative clauses, etc. that modify nouns. In contrast, the Stanford relations provide a more reasonable level of granularity.

Our dependency relation scheme is similar to Stanford's basic scheme but has several differences. It introduces several new relations including *ccinit* "initial coordinating conjunction", *cleft* "cleft clause", *combo* "combined term", *extr* "extraposed element", *infmark* "infinitive marker 'to'", *objcomp* "object complement", *postloc* "post-modifying location", *sccomp* "clausal complement of 'so'", *vch* "verbal chain" and *whadvmod* "wh- adverbial modifier". The *nsubjpass*, *csubjpass*, and *auxpass* relations of Stanford's are left out because adding them up front makes learning more difficult and the fact

<i>abbrev</i>	abbreviation	<i>csubjpass</i>	clausal subject (passive)	<i>pobj</i>	prepositional object
<i>acompl</i>	adjectival complement	<i>det</i>	determiner	<i>poss</i>	possessive
<i>advcl</i>	adverbial clause	<i>dobj</i>	direct object	<i>possessive</i>	possessive marker
<i>advmod</i>	adverbial modifier	<b>extr</b>	extraposed element	<b>postloc</b>	post-modifying location
<i>agent</i>	'by' agent	<i>expl</i>	'there' expletive	<i>preconj</i>	pre conjunct
<i>amod</i>	adjectival modifier	<b>infmark</b>	infinitive marker ('to')	<i>predet</i>	predeterminer
<i>appos</i>	appositive	<i>infmod</i>	infinite modifier	<i>prep</i>	preposition
<i>attr</i>	attributive	<i>iobj</i>	indirect object	<i>prt</i>	particle
<i>aux</i>	auxiliary	<i>mark</i>	subordinate clause marker	<i>punct</i>	punctuation
<i>auxpass</i>	auxiliary (passive)	<i>measure</i>	measure modifier	<i>purpcl</i>	purpose clause
<b>cleft</b>	cleft clause	<i>neg</i>	negative	<i>quantmod</i>	quantifier modifier
<i>cc</i>	coordination	<i>nn</i>	noun compound	<i>rcmod</i>	relative clause
<b>ccinit</b>	initial CC	<i>nsubj</i>	nominal subject	<i>rel</i>	relative
<i>ccomp</i>	clausal complement	<i>nsubjpass</i>	nominal subject (passive)	<b>sccomp</b>	clausal complement of 'so'
<b>combo</b>	combination term	<i>num</i>	numeric modifier	<i>tmod</i>	temporal modifier
<i>compl</i>	complementizer	<i>number</i>	compound number	<b>vch</b>	verbal chain
<i>conj</i>	conjunction	<b>objcomp</b>	object complement	<b>whadvmod</b>	wh- adverbial
<i>cop</i>	copula complement	<i>parataxis</i>	parataxis	<i>xcomp</i>	clausal complement w/o subj
<i>csubj</i>	clausal subject	<i>partmod</i>	participle modifier		

Table 1: Dependency scheme with differences versus basic Stanford dependencies highlighted. Bold indicates the relation does not exist in the Stanford scheme. Italics indicate the relation appears in Stanford's scheme but not ours.

that a *nsubj*, *csubj*, or *aux* is passive can easily be determined from the final tree. Stanford's *aux* dependencies are replaced using verbal chain (*vch*) links; conversion of these to Stanford-style *aux* dependencies is also trivial as a post-processing step.<sup>3</sup> The *attr* dependency is excluded because it is redundant with the *cop* relation due to different handling of copula, and the dependency scheme does not have an *abbrev* label because this information is not provided by the Penn Treebank. The dependency scheme with differences with Stanford highlighted is presented in Table 1.

In addition to using a slightly different set of dependency names, a handful of relations, notably *cop*, *conj*, and *cc*, are treated in a different manner. These differences are illustrated by Figure 1. The Stanford scheme's treatment of copula may be one reason why dependency parsers have trouble learning and applying it. Normally, the head of the clause is a verb, but, under Stanford's scheme, if the verb happens to be a copula, the complement of the copula (*cop*) is treated as the head of the clause instead.

<sup>3</sup>The parsing system includes an optional script that can convert *vch* arcs into *aux* and *auxpass* and the subject relations into *csubjpass* and *nsubjpass*.

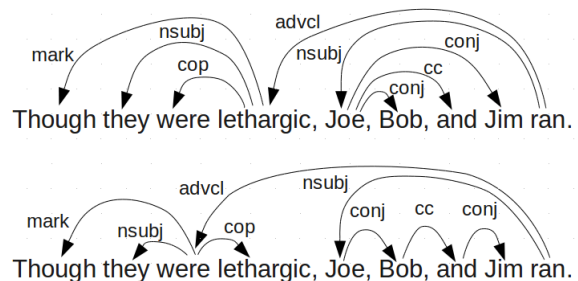


Figure 1: Example comparing Stanford's (top) handling of copula and coordinating conjunctions with ours (bottom).

### 3.2 Conversion Process

A three-step process is used to convert the Penn Treebank (Marcus, et al., 1993) from constituent parses into dependency trees labeled according to the dependency scheme presented in the prior section. The first step is to apply the noun phrase structure patch created by Vadas and Curran (2007), which adds structure to the otherwise flat noun phrases (NPs) of the Penn Treebank (e.g., '(metal soup pot cover)' would become '(metal (soup pot) cover)'). The second step is to apply a version of Johansson and Nugues' (2007) constituent-to-dependency converter with some head-finding rule modifications; these rules, with changes highlighted

(WH)?NP NX NML NAC	<u>FW NML NN*</u> JJR \$ # CD FW <u>QP</u> JJ NAC JJS PRP ADJP RB[SR] VBG DT WP
	RB NP-ε S SBAR UCP PP SINV SBARQ SQ UH VP NP VB VBP
ADJP JJP	NNS QP NN \$ # JJ VBN VBG (AD J)JP ADVP JJR NP NML JJS DT FW RBR RBS SBAR RB
ADVP	<b>RB RB JJ JJR</b> RBS FW ADVP TO CD IN NP NML JJS NN
PRN	S* <b>VP NN*</b>  NX NML NP W* PP IN ADJP JJ ADVP RB NAC <b>VP INTJ</b>
QP	\$ # NNS NN <b>CD</b> JJ RB DT NCD QP IN <u>CC</u> JJR JJS
SBARQ	SQ S SBARQ SINV FRAG
SQ	VBZ VBD VBP VB MD *-PRD SQ VP <b>FRAG X</b>
UCP	[ <b>QNV</b> P S* UCP NML PR[NT] RR C NX NAC FRAG INTJ AD[JV] P LST WH* X
VP	VBD AUX VBN MD VBZ VB VBG VBP VP <b>POS</b> *-PRD ADJP JJ NN NNS NP NML
WHADJP	CC <b>JJ</b> WRB ADJP
WHADVP	<u>CC</u> WRB RB
X	[ <b>QNV</b> P S* UCP NML PR[NT] RR C NX NAC FRAG INTJ AD[JV] P LST WH* X CONJP
LST	LS : DT NN SYM

Figure 2: Modified head-finding rules. Underline indicates that the search is performed in a left-to-right fashion instead of the default right-to-left order. NML and JJP are both products of Vadas and Curran’s (2007) patch. Bold indicates an added or moved element; for the original rules, see the paper by Johansson and Nugues (2007).

in bold, are provided in Figure 2. Finally, an additional script makes additional changes and converts the intermediate output into the dependency scheme.

This dependency conversion has several advantages to it. Using the modified head-finding rules for Johansson and Nugues’ (2007) converter results in fewer buggy trees than were present in the CoNLL shared tasks, including fewer trees in which words are headed by punctuation marks. For sections 2–21, there are far fewer generic *dep/DEP* relations (2,765) than with the Stanford conversion (34,134) or the CoNLL 2008 shared task conversion (23,811). Also, the additional conversion script contains various rules for correcting part-of-speech (POS) errors using the syntactic structure as well as additional rules for some specific word forms, mostly common words with inconsistent taggings. Many of these changes cover part-of-speech problems discussed by Manning (2011), including VBD/VBN, VBZ/NNS, NNP/NNPS, and IN/WDT/DT issues. In total, the script changes over 9,500 part-of-speech tags, with the most common change being to change preposition tags (IN) into adverb tags (RB) for cases where there is no prepositional complement/object. The top fifteen of these changes are presented in Table 2. The conversion script contains a variety of additional rules for modifying the parse structure and fixing erroneous trees as well, including cases where one or more POS tags were incorrect and, as such, the initial dependency parse was flawed. Quick manual inspections of the changes suggested that the

vast majority are accurate.

In the final output from the conversion, the number of sentences with one or more words dependent on non-projective arcs in sections 2–21 is 3,245—about 8.1% of the dataset. About 1.3% of this, or 556 of sentences, is due to the secondary conversion script, with sentences containing approximate currency amounts (e.g., *about \$ 10*) comprising the bulk of difference. For these, the quantifying text (e.g., *about, over, nearly*), is linked to the number following the currency symbol instead of to the currency symbol as it was in the CoNLL 2008 task.

Original	New	# of changes
IN	RB	1128
JJ	NN	787
VBD	VBN	601
RB	IN	462
VBN	VBD	441
NN	JJ	409
NNPS	NNP	405
IN	WDT	388
VBG	NN	223
DT	IN	220
RB	JJ	214
VB	VBP	184
NN	NNS	169
RB	NN	157
NNS	VBZ	148

Table 2: Top 15 part-of-speech tag changes performed by the conversion script.

## 4 Parser

### 4.1 Algorithm

The parsing approach is based upon the *non-directional easy-first* algorithm recently presented by Goldberg and Elhadad (2010). Their original algorithm behaves as follows. For a sentence of length  $n$ , the algorithm performs a total of  $n$  steps. In each step, one of the unattached tokens is added as a child to one of its current neighbors and is then removed from the list of unprocessed tokens. When only one token remains unprocessed, it is designated as the root. Provided that only a constant number of potential attachments need to be re-evaluated after each step, which is the case if one restricts the context for feature generation to a constant number of neighboring tokens, the algorithm can be implemented to run in  $O(n \log n)$ . However, since only  $O(n)$  dot products must be calculated by the parser and these have a large constant associated with them, the running time will rival  $O(n)$  parsers for any reasonable  $n$ , and, thus, a naive  $O(n^2)$  implementation will be nearly as fast as a priority queue implementation in practice.<sup>4</sup>

The algorithm has a couple potential advantages over standard shift-reduce style parsing algorithms. The first advantage is that performing easy actions first may make the originally difficult decisions easier. The second advantage is that performing parse actions in a more flexible order than left-to-right/right-to-left shift-reduce parsing reduces the chance of error propagation.

Unfortunately, the original algorithm does not support non-projective trees. To extend the algorithm to support non-projective trees, we introduce *move-right* and *move-left* operations similar to the stack-to-buffer *swaps* proposed by Nivre (2009) for shift-reduce style parsing. Thus, instead of attaching a token to one of its neighbors at each step, the algorithm may instead decide to *move* a token past one of its neighbors. Provided that no node is allowed to be moved past a token in such a way that a previous *move* operation is undone, there can be at most  $O(n^2)$  *moves* and the overall worst-case complexity becomes  $O(n^2 \log n)$ . While theoretically slower, this has a limited impact upon actual parsing times

in practice, especially for languages with relatively fixed word order such as English.<sup>5</sup> Though Goldberg and Elhadad's (2010) original implementation only supports unlabeled dependencies, the algorithm itself is in no way limited in this regard, and it is simple enough to add labeled dependency support by treating each dependency label as a specific type of *attach* operation (e.g., *attach\_as\_nsubj*), which is the method used by this implementation. Pseudocode for the non-directional easy-first algorithm with non-projective support is given in Algorithm 1.

```
input :  $w_1 \dots w_n$ , #the sentence
         $m$ , #the model
         $k$ , #the context width
         $actions$ , #the list of parse actions
         $\phi$ , #the feature generator
output:  $tree$  #a collection of dependency arcs
 $words = copyOf(s)$ ;
 $stale = copyOf(s)$ ;
 $cache$ ; #cache of action scores
while  $|words| > 1$  do
    for  $w \in stale$  do
        for  $act \in actions$  do
             $cache[w,act] = score(act, \phi(w, \dots), m)$ ;
             $stale.remove(w)$ ;
         $best = \arg \max_{a \in actions \& valid(a), w \in words} cache[w, a]$ 
    if  $isMove(best)$  then
         $i = words.index(getTokenToMove(best))$ ;
         $words.move(i, isMoveLeft(best) ? -1 : 1)$ ;
    else
         $arc = createArc(best)$ ;
         $tree.add(arc)$ ;
         $i = words.index(getChild(arc))$ ;
         $words.remove(i)$ ;
    for  $x \in -k, \dots, k$  do
         $stale.add(words.get(index+x))$ ;
return  $tree$ 
```

**Algorithm 1:** Modified version of Goldberg and Elhadad's (2010) Easy-First Algorithm with non-projective support.

<sup>4</sup>See Goldberg and Elhadad (2010) for more explanation.

<sup>5</sup>See Nivre (2009) for more information on the effect of re-ordering operations on parse time.

## 4.2 Features

One of the key aspects of the parser is the complex set of features used. The feature set is based off the features used by Goldberg and Elhadad (2010) but has a significant number of extensions. Various feature templates are specifically designed to produce features that help with several syntactic issues including preposition attachment, coordination, adverbial clauses, clausal complements, and relative clauses. Unfortunately, there is insufficient space in this paper to describe them all here. However, a list of feature templates will be provided with the parser download.

Several of the feature templates use unsupervised word clusters created with the Brown et al. (1992) hierarchical clustering algorithm. The use of this algorithm was inspired by Koo et al. (2008), who used the top branches of the cluster hierarchy as features. However, unlike Koo et al.'s (2008) parser, the fine-grained cluster identifiers are used instead of just the top 4-6 branches of the cluster hierarchy. The 175 word clusters utilized by the parser were created from the New York Times corpus (Sandhaus, 2008). Some examples from the clusters are presented in Figure 3. The ideal number of such clusters was not thoroughly investigated.

while where when although despite unless unlike ...  
why what whom whatever whoever whomever whence ...  
based died involved runs ended lived charged born ...  
them him me us himself themselves herself myself ...  
really just almost nearly simply quite fully virtually ...  
know think thought feel believe knew felt hope mean ...  
into through on onto atop astride Saturday/Early thru ...  
Ms. Mr. Dr. Mrs. Judge Miss Professor Officer Colonel ...  
John President David J. St. Robert Michael James George ...  
wife own husband brother sister grandfather beloved ...  
often now once recently sometimes clearly apparently ...  
everyone it everybody somebody anybody nobody hers ...  
around over under among near behind outside across ...  
Clinton Bush Johnson Smith Brown Williams King ...  
children companies women people men things students ...

Figure 3: High frequency examples from 15 of the Brown clusters.

## 4.3 Training

The parsing model is trained using a variant of the structured perceptron training algorithm used in the original Goldberg and Elhadad (2010) implementa-

tion. The general idea of the algorithm is to iterate over the sentences and, whenever the model predicts an incorrect action, update the model weights. Following Goldberg and Elhadad, parameter averaging is used to reduce overfitting.

Our implementation varies slightly from that of Goldberg and Elhadad (2010). The difference is that, at any particular step for a given sentence, the algorithm continues to update the weight vector as long as *any* invalid action is scored higher than *any* valid action, not just the highest scoring valid action; unfortunately, this change significantly slowed down the training process. In early experiments, this change produced a slight improvement in accuracy though it also slowed training significantly. In later experiments using additional feature templates, this change ceased to have any notable impact on the overall accuracy, but it was kept anyway.<sup>6</sup>

The oracle used to determine whether a *move* operation should be considered legal during the training phase is similar to Nivre et al.'s (2009) improved oracle based upon *maximal projective subcomponents*. As an additional restriction, during training, *move* actions were only considered valid either if no other action was valid or if the token to be moved already had all its children attached and moving it caused it to be adjacent to its parent. This fits with Nivre et al.'s (2009) intuition that it is best to delay word reordering as long as possible.

## 4.4 Speed Enhancements

To enhance the speed for practical use, the parser uses constraints based upon the part-of-speech tags of the adjacent word pairs to eliminate invalid dependencies from even being evaluated. A relation is only considered between a pair of words if such a relation was observed in the training data between a pair of words with the same parts-of-speech (with the exception of the generic *dep* dependency, which is permitted between any POS tag pair). Early experiments utilizing similar constraints showed an improvement in parsing speed of about 16% with no significant impact on accuracy, regardless of whether the constraints were enforced during training.

<sup>6</sup>See Goldberg and Elhadad (2010) for more description of the general training procedure.

System	Arc Accuracy		Perfect Sentences		Non-Proj Arcs	
	Labeled	Unlabeled	Labeled	Unlabeled	Labeled	Unlabeled
THIS WORK	<b>92.1 (93.3)</b>	<b>93.7 (94.3)</b>	<b>38.4 (42.5)</b>	<b>46.2 (48.5)</b>	66.5 (69.7)	<b>69.3 (71.7)</b>
THIS WORK <sub>no clusters</sub>	91.8 (93.1)	93.4 (94.1)	38.2 (42.3)	45.5 (47.3)	<b>67.3 (70.9)</b>	<b>69.3 (72.5)</b>
THIS WORK <sub>moves disabled</sub>	91.7 (92.9)	93.3 (93.9)	37.1 (40.8)	44.2 (46.2)	21.1 (21.1)	22.7 (21.9)
NON-DIR EASY FIRST	*	91.2 (92.0)	*	37.8 (39.4)	*	15.1 (16.3)
EISNER <sup>†</sup> <sub>MST</sub>	90.9 (92.2)	92.8 (93.5)	32.1 (35.6)	40.6 (42.3)	62.5 (65.3)	63.7 (66.9)
CHU-LIU-EDMONDS <sub>MST</sub>	90.0 (91.2)	91.8 (92.5)	28.4 (31.3)	35.0 (36.4)	62.9 (65.3)	64.1 (66.5)
ARC-EAGER <sub>Malt</sub>	89.8 (91.1)	91.3 (92.1)	31.6 (34.2)	37.4 (38.5)	19.5 (19.5)	20.3 (19.9)
ARC-STANDARD <sub>Malt</sub>	88.3 (89.5)	89.7 (90.4)	31.4 (34.1)	36.1 (37.3)	13.1 (12.0)	13.9 (12.7)
STACK-EAGER <sub>Malt</sub>	90.0 (91.2)	91.5 (92.3)	34.5 (37.5)	40.4 (41.9)	51.8 (53.8)	53.8 (55.4)
STACK-LAZY <sub>Malt</sub>	90.4 (91.7)	91.9 (92.8)	34.8 (37.7)	40.6 (42.5)	61.8 (63.3)	63.3 (65.3)
CHARNIAK <sup>‡</sup>	*	93.2	*	43.5	*	32.3
BERKELEY <sup>‡</sup>	*	93.3	*	43.6	*	34.3

Table 3: Parsing results for section 23 of the Penn Treebank (punctuation excluded). Results in parentheses were produced using gold POS tags. <sup>†</sup>Eisner (1996) algorithm with non-projective rewriting and second order features. <sup>‡</sup>Results not directly comparable; see text. \*Labeled dependencies not available/comparable.

#### 4.5 Evaluation

The following split of the Penn Treebank (Marcus, et al., 1993) was used for the experiments: sections 2–21 for training, 22 for development, and 23 for testing.

For part-of-speech (POS) tagging, we used an in-house SVM-based POS tagger modeled after the work of Giménez and Márquez (2004)<sup>7</sup>. The training data was tagged in a 10-fold fashion; each fold was tagged using a tagger trained from the nine remaining folds. The development and test sections were tagged by an instance of the tagger trained using the entire training set. The full details of the POS tagger are outside the scope of this paper; it is included with the parser download.

The final parser was trained for 31 iterations, which is the point at which its performance on the development set peaked. One test run was performed with non-projectivity support disabled in order to get some idea of the impact of the *move* operations on the parser’s overall performance; also, since the parsers used for comparison had no access to the unsupervised word clusters, an additional instance of the parser was trained with every word treated as belonging to the same cluster so as to facilitate a more fair comparison.

Seven different dependency parsing models were

<sup>7</sup>97.42% accuracy on traditional POS evaluation (Penn Treebank WSJ sections 22-24).

trained for comparison using the following open source parsing packages: Goldberg and Elhadad’s (2010)’s *non-directional easy-first* parser, MALTPARSER (Nivre et al., 2006), and MSTPARSER (McDonald and Pereira, 2006)<sup>8</sup>. The model trained using Goldberg and Elhadad’s (2010) easy-first parser serves as something of a baseline. The four MALTPARSER parsing models used the *arc-eager*, *arc-standard*, *stack-eager*, and *stack-lazy* algorithms. One of the MSTPARSER models used the Chu-Liu-Edmonds maximum spanning tree approach, and the other used the Eisner (1996) algorithm with second order features and a non-projective rewriting post-processing step.

Unfortunately, it is not possible to directly compare the parser’s accuracy with most popular constituent parsers such as the Charniak (2000) and Berkeley (Petrov et al., 2006; Petrov and Klein, 2007) parsers<sup>9</sup> both because they do not produce functional tags for subjects, direct objects, etc., which are required for the final script of the constituent-to-dependency conversion routine, and because they determine part-of-speech tags in conjunction with the parsing. However, it is possible to compute approximate unlabeled accuracy scores by training the constituent parsers on the NP-patched (Vadas and Curran, 2007) version of the data and then running the test output through just the first conversion script—that is, the modified version of Johansson and Nugues’ (2007) converter.

The results of the experiment are given in Table 3, including accuracy for individual arcs, non-projective arcs only, and full sentence match. Punctuation is excluded in all the result computations. To determine whether an arc is non-projective, the following heuristic was used. Traverse the sentence in a depth-first search, starting from the imaginary root node and pursuing child arcs in order of increasing absolute distance from their parent. Whenever an arc being traversed is found to cross a previously traversed arc, mark it as non-projective and continue. To evaluate the impact of part-of-speech tagging error, results for parsing using the gold standard part-of-speech tags are also included.

We also measured the speed of the parser on the various sentences in the test collection. For reasonable sentence lengths, the parser scales quite well. The scatterplot depicting the relation between sentence length and parsing time is presented in Figure 5.

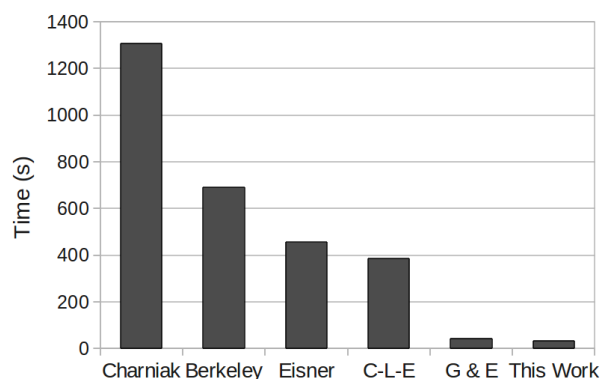


Figure 4: Parse times for Penn Treebank section 23 for the parsers on a PC with a 2.4Ghz Q6600 processor and 8GB RAM. MALTPARSER ran substantially slower than the others, perhaps due to its use of polynomial kernels, and isn't shown. (C-L-E - Chu-Liu-Edmonds, G&E - Goldberg and Elhadad (2010)).

#### 4.5.1 Results Discussion

The parser achieves 92.1% labeled and 93.7% unlabeled accuracy on the evaluation, a solid result and about 2.5% higher than the original easy-first implementation of Goldberg and Elhadad (2010). Furthermore, the parser processed the entire test section in

<sup>8</sup>Versions 1.4.1, 0.4.3b, and 0.2, respectively

<sup>9</sup>Versions 1.1 and 05Aug16, respectively

just over 30 seconds—a rate of over 75 sentences per second, substantially faster than most of the other parsers.

Not surprisingly, the results for non-projective arcs are substantially lower than the results for all arcs, and the systems that are designed to handle them outperformed the strictly projective parsers in this regard.

The negative effect of part-of-speech tagging error appears to impact the different parsers about the same amount, with a loss of .6% to .8% in unlabeled accuracy and 1.1% to 1.3% in labeled accuracy.

The 93.2% and 93.3% accuracy scores achieved by the Charniak and Berkeley parsers are not too different from the 93.7% result, but, of course, it is important to remember that these scores are not directly comparable.

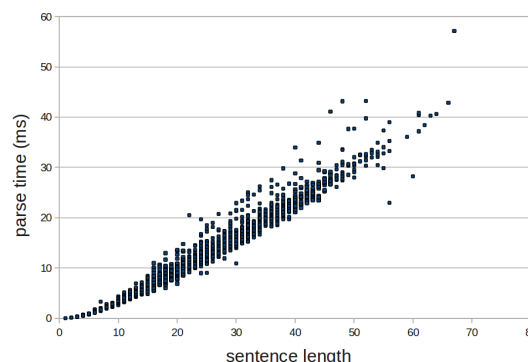


Figure 5: Sentence length versus parse time. Median times for five runs over section 23.

## 5 Shallow Semantic Annotation

To create a more informative parse, the parser includes four optional modules, a preposition sense disambiguation (PSD) system, a work-in-progress 's-possessive interpretation system, a noun compound interpretation system, and a PropBank-based semantic role labeling system<sup>10</sup>. Taken together, these integrated modules enable the parsing system to produce substantially more informative output than a traditional parser.

**Preposition Sense Disambiguation** The PSD system is a newer version of the system described

<sup>10</sup>Lack of space prohibits a sufficiently thorough discussion of these individual components and their evaluations, but additional information will be available with the system download.



by Tratz and Hovy (2009) and Hovy et al. (2010); it achieves 85.7% accuracy on the SemEval-2007 fine-grain PSD task (Litkowski and Hargraves, 2007), which is a statistically significant ( $p \leq 0.05$ ; upper-tailed z test) increase over the previous best reported result for this dataset, Hovy et al.’s (2010) 84.8%.

**Noun Compound Interpretation** The noun compound interpretation system is a newer version of the system described by Tratz and Hovy (2010) with similar accuracy (79.6% vs 79.3% using 10-fold cross-validation<sup>11</sup>).

**Possessives Interpretation** The possessive interpretation system assigns interpretations to ’s possessives (e.g., John’s arm  $\rightarrow$  PART-OF, Mowgli’s capture  $\rightarrow$  PATIENT/THEME). The current system achieves over 85.0% accuracy, but it is important to note that the annotation scheme, automatic classifier, and dataset are all still under active development.

**PropBank SRL** The PropBank-based semantic role labeling system achieves 86.8 combined  $F_1$  measure for automatically-generated parse trees calculated over both predicate disambiguation and argument/adjunct classification (89.5  $F_1$  on predicate disambiguation, 85.6  $F_1$  on argument and adjuncts corresponding to dependency links, and 86.8  $F_1$ ); this score is not directly comparable to any previous work due to some differences, including differences in both the parse tree conversion and the PropBank conversion. The most similar work is that of the CoNLL shared task work (Surdeanu et al., 2008; Hajič et al., 2009).

## 6 Related Work

**Non-projectivity.** There are two main approaches used in recent NLP literature for handling non-projectivity in parse trees. The first is to use an algorithm, like the one presented in this paper, that has inherent support for non-projective trees. Examples of this include the Chu-Liu-Edmonds’ approach for maximum spanning tree (MST) parsing (McDonald et al., 2005) and Nivre’s (2009) swap-based reordering method for shift-reduce parsing. The second approach is to create an initial projective parse and then apply transformations to intro-

duce non-projectivity into it. Examples of this include McDonald and Pereira’s (2006) rewriting of projective trees produced by the Eisner (1996) algorithm, and Nivre and Nilsson’s (2005) pseudo-projective approach that creates projective trees with specially marked arcs that are later transformed into non-projective dependencies.

**Descriptive dependency labels.** While most recent dependency parsing research has used either vague labels, such as those of the CoNLL shared tasks, or no labels at all, some descriptive dependency label schemes exist. By far the most prominent of these is the Stanford typed dependency scheme (de Marneffe and Manning, 2008). Another descriptive scheme that exists, but which is less widely used in the NLP community, is the one used by Tapanainen and Järvinen’s parser (1997). Unfortunately, the Stanford dependency conversion of the Penn Treebank has proven difficult to learn for current dependency parsers (Cer et al., 2010), and there is no publicly available dependency conversion according to Tapanainen and Järvinen’s scheme.

**Faster parsing.** While the fastest reasonable parsing algorithms are the  $O(n)$  shift-reduce algorithms, such as Nivre’s (2003) algorithm and an expected linear time dynamic programming approach presented by Huang and Sagae (2010), a few other fast alternatives exist. Goldberg and Elhadad’s (2010) easy-first algorithm is one such example. Another example, is Roark and Hollingshead’s (2009) work that uses chart constraints to achieve linear time complexity for constituency parsing.

**Effective features for parsing.** A variety of work has investigated the use of more informative features for parsing. This includes work that integrates second and even third order features (McDonald et al., 2006; Carreras, 2007; Koo and Collins, 2010). Also, some work has incorporated unsupervised word clusters as features, including that of Koo et al. (2008) and Suzuki et al. (2009), who utilized unsupervised word clusters created using the Brown et al. (1992) hierarchical clustering algorithm.

**Semantically-enriched output.** The 2008 and 2009 CoNLL shared tasks (Surdeanu et al., 2008; Hajič et al., 2009), which required participants to build systems capable of both syntactic parsing and Semantic Role Labeling (SRL) (Gildea and Jurafsky, 2002), are the most notable attempts to encour-

<sup>11</sup>These accuracy figures are higher than what should be expected for unseen datasets; see Tratz and Hovy (2010) for more detail.

age the development of parsers with additional semantic annotation. These tasks relied upon PropBank (2005) and NomBank (2004) for the semantic roles. A variety of other systems have focused on FrameNet-based (1998) SRL instead, including those that participated in the SemEval-2007 Task 19 (Baker et al., 2007) and work by Das et al. (2010).

## 7 Conclusion

In this paper, we have described a new high-quality dependency tree conversion of the Penn Treebank (Marcus, et al., 1993) along with its labeled dependency scheme and presented a parser that is fast, accurate, supports non-projective trees and provides rich output, including not only informative dependency labels similar to Stanford’s but also additional semantic annotation for prepositions, possessives, and noun compound relations. We showed how the easy-first algorithm of Goldberg and Elhadad (Goldberg and Elhadad, 2010) can be extended to support non-projective trees by adding *move* actions similar to Nivre’s (2009) swap-based reordering for shift-reduce parsing and evaluated our parser on the standard test section of the Penn Treebank, comparing with several other freely available parsers.

The Penn Treebank conversion process fixes a number of buggy trees and part-of-speech tags and produces dependency trees with a relatively small percentage of generic *dep* dependencies. The experimental results show that dependency parsers can generally produce Stanford-granularity labels with high accuracy when using the new dependency conversion of the Penn Treebank, something which, according to the findings of Cer et al. (2010), does not appear to be the case when training and testing dependency parsers on the Stanford conversion.

The parser achieves high labeled and unlabeled accuracy in the evaluation, 92.1% and 93.7%, respectively. The 93.7% result represents a 2.5% increase over the accuracy of Goldberg and Elhadad’s (2010) implementation. Also, the parser proves to be quite fast, processing section 23 of the Penn Treebank in just over 30 seconds (a rate of over 75 sentences per second).

The parsing system is capable of not only producing fine-grained dependency relations, but can also produce shallow semantic annotations for preposi-

tions, possessives, and noun compounds by using several optional integrated modules. The preposition sense disambiguation (PSD) module achieves 85.7% accuracy on the SemEval-2007 PSD task, exceeding the previous best published result of 84.8% by a statistically significant margin, the possessives module is over 85% accurate, the noun compound interpretation module achieves 79.6% accuracy on Tratz and Hovy’s (2010) dataset. The PropBank SRL module achieves 89.5  $F_1$  on predicate disambiguation and 85.6  $F_1$  on argument and adjuncts corresponding to dependency links, for an overall  $F_1$  of 86.8. Combined with the core parser, these modules allow the system to produce a substantially more informative textual analysis than a standard parser.

## 8 Future Work

There are a variety of ways to extend and improve upon this work. We would like to change our handling of coordinating conjunctions to treat the coordinating conjunction as the head because this has fewer ambiguities than the current approach and also add the ability to produce traces for WH- words. It would also be interesting to examine the impact on final parsing accuracy of the various differences between our dependency conversion and Stanford’s.

To aid future NLP research work, the code, including the treebank converter, part-of-speech tagger, parser, and semantic annotation add-ons, will be made publicly available for download via <http://www.isi.edu>.

## Acknowledgements

We would like to thank Richard Johansson for providing us with the code for the *penncconverter* constituent-to-dependency converter. We would also like to thank Dirk Hovy and Anselmo Peñas for many valuable discussions and suggestions.

## References

- Collin Baker, and Michael Ellsworth and Katrin Erk. 2007. SemEval’07 task 19: Frame Semantic Structure Extraction. In *Proc. of the 4th International Workshop on Semantic Evaluations*
- Collin Baker, Charles J. Fillmore and John B. Lowe. 1998. The Berkeley FrameNet Project. In *Proc. of the*

- 17th international conference on Computational linguistics*
- Adam L. Berger, Vincent J. Della Pietra, and Stephen A. Della Pietra. 1996. A maximum entropy approach to natural language processing. In *Computational Linguistics* 22(1):39–71
- Peter F. Brown, Peter V. deSouza, Robert L. Mercer, Vincent J. Della Pietra, and Jenifer C. Lai. 1992. Class-Based n-gram Models of Natural Language. *Computational Linguistics* 18(4):467–479.
- Sabine Buchholz and Erwin Marsi. 2006. CoNLL-X shared task on multilingual dependency parsing. In *Proc. of CoNLL 2006*.
- Xavier Carreras. 2007. Experiments with a Higher-Order Projective Dependency Parser. In *Proc. of the CoNLL Shared Task Session of EMNLP-CoNLL 2007*.
- Daniel Cer, Marie-Catherine de Marneffe, Daniel Jurafsky, and Christopher D. Manning. 2010. Parsing to Stanford Dependencies: Trade-offs between speed and accuracy. In *Proc. of LREC 2010*.
- Eugene Charniak. 2000. A Maximum-Entropy-Inspired Parser. In *Proc. of NAACL 2000*.
- Eugene Charniak and Mark Johnson. 2005. Coarse-to-fine-grained *n*-best parsing and discriminative reranking. In *Proc. of ACL 2005*.
- Michael A. Covington. 2001. A Fundamental Algorithm for Dependency Parsing. In *Proc. of the 39th Annual ACM Southeast Conference*.
- Dipanjan Das, Nathan Schneider, Desai Chen, and Noah A. Smith. 2010. Probabilistic Frame-Semantic Parsing. In *Proc. of HLT-NAACL 2010*.
- Marie-Catherine de Marneffe and Christopher D. Manning. 2008. The Stanford typed dependencies representation. In *COLING Workshop on Cross-framework and Cross-domain Parser Evaluation*.
- Jason Eisner. 1996. Three New Probabilistic Models for Dependency Parsing: An Exploration. In *Proc. of COLING 1996*.
- Christiane Fellbaum. 1998. WordNet: An Electronic Lexical Database. MIT Press.
- Daniel Gildea and Daniel Jurafsky. 2002. Automatic labeling of semantic roles. *Computational Linguistics* 28(3):245–288.
- Jesús Giménez and Lluís Márquez. 2004. SVMTool: A General POS Tagger Generator Based on Support Vector Machines. In *Proc. of LREC 2004*.
- Yoav Goldberg and Michael Elhadad. 2010. An Efficient Algorithm for Easy-First Non-Directional Dependency Parsing. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the ACL*.
- Yoav Goldberg and Michael Elhadad. 2009. The CoNLL-2009 Shared Task: Syntactic and Semantic Dependencies in Multiple Languages. In *Proc. of the Thirteenth Conference on Computational Natural Language Learning: Shared Task*.
- Dirk Hovy, Stephen Tratz, and Eduard Hovy. 2010. What’s in a Preposition?—Dimensions of Sense Disambiguation for an Interesting Word Class. In *Proc. of COLING 2010*.
- Liang Huang and Kenji Sagae. 2010. Dynamic Programming for Linear-Time Shift-Reduce Parsing. In *Proc. of ACL 2010*.
- Richard Johansson and Pierre Nugues. 2007. Extended constituent-to-dependency conversion for english. In *Proc. of NODALIDA*.
- Terry Koo, Xavier Carreras, and Michael Collins. 2008. Simple Semi-supervised Dependency Parsing. In *Proc. of ACL 2008*.
- Terry Koo and Michael Collins. 2010. Efficient Third-order Dependency Parsers. In *Proc. of ACL 2010*.
- Ken Litkowski and Orin Hargraves. 2007. SemEval-2007 Task 06: Word-Sense Disambiguation of Prepositions. In *Proc. of the 4th International Workshop on Semantic Evaluations*.
- Christopher D. Manning. 2011. Part-of-Speech Tagging from 97% to 100%: Is It Time for Some Linguistics? In *Proc. of the 12th International Conference on Intelligent Text Processing and Computational Linguistics (CICLING 2011)*.
- 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.
- Ryan McDonald, Fernando Pereira, Kiril Ribarov, and Jan Hajič. 2005. Non-Projective Dependency Parsing Using Spanning Tree Algorithms. In *Proc. of HLT-EMNLP 2005*.
- Ryan McDonald and Fernando Pereira. 2006. Online Learning of Approximate Dependency Parsing Algorithms. In *Proc. of EACL 2006*.
- Adam Meyers, Ruth Reeves, Catherine Macleod, Rachel Szekely, Veronika Zielinska, Brian Young and Ralph Grishman. 2004. The NomBank Project: An Interim Report. In *Proc. of the NAACL/HLT Workshop on Frontiers in Corpus Annotation*.
- Joakim Nivre. 2009. Non-Projective Dependency Parsing in Expected Linear Time. In *Proc. of the 47th Annual Meeting of the ACL and the 4th IJCNLP of the AFNLP*.
- Joakim Nivre. 2003. An Efficient Algorithm for Projective Dependency Parsing. In *Proc. of the 8th International Workshop on Parsing Technologies (IWPT)*.
- Joakim Nivre, Marco Kuhlmann, and Johan Hall. 2009. An Improved Oracle for Dependency Parsing with Online Reordering. In *Proc. of the 11th International Conference on Parsing Technologies (IWPT)*.
- Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. 2007. The CoNLL 2007 shared task on dependency parsing. In *Proc. of EMNLP-CoNLL 2007*.

- Joakim Nivre, Johan Hall, and Jens Nilsson. 2006. Malt-Parser: A Data-Driven Parser-Generator for Dependency Parsing. In *Proc. of LREC 2006*.
- Joakim Nivre and Jens Nilsson. 2005. Pseudo-projective dependency parsing. In *Proc. of ACL-2005*.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The Proposition Bank: An Annotated Corpus of Semantic Roles. In *Computational Linguistics*. 31(1):71–106.
- Slav Petrov and Dan Klein. 2007. Improved Inference for Unlexicalized Parsing. In *Proc. of HLT-NAACL 2007*.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In *Proc. of COLING-ACL 2006*.
- Brian Roark and Kristy Hollingshead. 2009. Linear complexity context-free parsing pipelines via chart constraints. In *Proc. of HLT-NAACL*.
- Evan Sandhaus. 2008. The New York Times Annotated Corpus. Linguistic Data Consortium, Philadelphia.
- Mihai Surdeanu, Richard Johansson, Adam Meyers, Lluís Màrquez, and Joakim Nivre. 2008. The CoNLL-2008 shared task on joint parsing of syntactic and semantic dependencies. In *Proc. of the Twelfth Conference on Computational Natural Language Learning*.
- Jun Suzuki, Hideki Isozaki, Xavier Carrerras, and Michael Collins. 2009. An Empirical Study of Semi-supervised Structured Conditional Models for Dependency Parsing. In *Proc. of EMNLP*.
- Pasi Tapanainen and Timo Järvinen. 1997. A non-projective dependency parser. In *Proc. of the fifth conference on applied natural language processing*.
- Stephen Tratz and Dirk Hovy. 2009. Disambiguation of Preposition Sense using Linguistically Motivated Features. In *Proc. of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Student Research Workshop and Doctoral Consortium*.
- Stephen Tratz and Eduard Hovy. 2010. A Taxonomy, Dataset, and Classifier for Automatic Noun Compound Interpretation. In *Proc. of ACL 2010*.
- David Vadas and James R. Curran. 2007. Adding Noun Phrase Structure to the Penn Treebank. In *Proc. of ACL 2007*.
- Hiroyasu Yamada and Yuji Matsumoto. 2003. Statistical Dependency Analysis With Support Vector Machines. In *Proc. of 8th International Workshop on Parsing Technologies (IWPT)*.