

PARSER EVALUATION USING TEXTUAL ENTAILMENTS

by

Önder Eker

B.S., Computer Engineering, Boğaziçi University, 2007

Submitted to the Institute for Graduate Studies in
Science and Engineering in partial fulfillment of
the requirements for the degree of
Master of Science

Graduate Program in Computer Engineering
Boğaziçi University

2009

PARSER EVALUATION USING TEXTUAL ENTAILMENTS

APPROVED BY:

Assoc. Prof. Tunga Güngör

(Thesis Supervisor)

Assist. Prof. Deniz Yüret

(Thesis Co-supervisor)

Dr. Ali Vahit Şahiner

Assist. Prof. Alper Şen

Dr. Tamer Şıkoğlu

DATE OF APPROVAL: 14.08.2009

ACKNOWLEDGEMENTS

I would like to thank Assoc. Prof. Tunga Güngör and Assist. Prof. Deniz Yüret for their guidance.

I would like to thank Dr. Ali Vahit Şahiner, Assist. Prof. Alper Şen and Dr. Tamer Şikođlu for kindly accepting to be committee members.

I would like to thank my family for their love, support and encouragement.

The author was supported by TÜBİTAK MS Fellowship BİDEB-2228.

ABSTRACT

PARSER EVALUATION USING TEXTUAL ENTAILMENTS

Syntactic parsing is a basic problem in natural language processing. It can be defined as assigning a structure to a sentence. Two prevalent approaches to parsing are phrase-structure parsing and dependency parsing. A related problem is parser evaluation. PETE is a dependency-based evaluation where the parse is represented as a list of simple sentences, similar to the Recognizing Textual Entailments task. Each entailment focuses on one relation. A priori training of annotators is not required. A program generates entailments from a dependency parse. Phrase-structure parses are converted to dependency parses to generate entailments. Additional entailments are generated for phrase-structure coordinations. Experiments are carried out with a function-tagger. Parsers are evaluated on the set of entailments generated from the Penn Treebank WSJ and Brown test sections. A phrase-structure parser obtained the highest score.

ÖZET

METİNSEL GEREKTİRİMLER İLE AYRIŞTIRICI DEĞERLENDİRMESİ

Sözdizimsel ayrıştırma doğal dil işlemede temel bir problemdir. Cümleye bir yapı atamak olarak tanımlanabilir. En yaygın iki ayrıştırma, öbek yapısı ayrıştırma ve bağımsallık ayrıştırmasıdır. İlgili bir konu ayrıştırıcı değerlendirmesidir. PETE, ayrıştırmanın Metinsel Gerektirimleri Tanıma görevinde olduğu gibi bir dizi basit cümle ile ifade edildiği bağımsallık tabanlı bir değerlendirmedir. Her gerektirim bir bağlantıya odaklanır. Yorumcuların önceden eğitilmesine gerek yoktur. Bir program bağımsallık ayrıştırmasından gerektirimleri üretmektedir. Öbek yapısı ayrıştırmaları gerektirim üretmek için bağımsallık ayrıştırmasına çevrilmiştir. Öbek yapısı eşgüdümleden ek gerektirimler üretilmektedir. Bir işlev etiketçi ile deneyler yapılmıştır. Ayrıştırıcılar Penn Treebank WSJ ve Brown test kısımlarından üretilen gerektirim kümesi üzerinde değerlendirilmiştir. Bir öbek yapısı ayrıştırıcı en yüksek puanı almıştır.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	iii
ABSTRACT	iv
ÖZET	v
LIST OF FIGURES	vii
LIST OF TABLES	viii
LIST OF SYMBOLS/ABBREVIATIONS	ix
1. INTRODUCTION	1
1.1. Phrase-Structure Parsing	1
1.2. Dependency Parsing	2
1.3. Recognizing Textual Entailments	2
2. LITERATURE SURVEY	4
3. GENERATING ENTAILMENTS	8
3.1. Generating Entailments From Dependency Parses	8
3.2. Generating Entailments From Phrase-Structure Parses	9
4. EVALUATION	12
4.1. Evaluation on the WSJ Corpus	12
4.2. Evaluation on the Brown Corpus	14
5. CONCLUSION	17
APPENDIX A: SAMPLE ENTAILMENTS	18
REFERENCES	19

LIST OF FIGURES

Figure 3.1.	Pseudocode for generating entailments from dependency parses . . .	8
Figure 3.2.	Coordination in dependency representation	9
Figure 3.3.	Coordination in phrase-structure representation	10
Figure 3.4.	Generating entailments from phrase-structure parse trees	11

LIST OF TABLES

Table 4.1.	Bracketing scores on the WSJ test section	12
Table 4.2.	LAS and UAS scores on the WSJ test section	12
Table 4.3.	LAS and UAS scores of all parsers on the WSJ test section	13
Table 4.4.	PETE scores on the WSJ test section	14
Table 4.5.	Bracketing scores on the Brown test section	15
Table 4.6.	LAS and UAS scores on the Brown test section	15
Table 4.7.	PETE scores on the Brown test section	16

LIST OF SYMBOLS/ABBREVIATIONS

<i>c</i>	Conjunct
<i>D</i>	Dependency parse tree
<i>e</i>	Entailment
<i>E</i>	List of entailments
<i>h</i>	Head
<i>l</i>	Lowest common ancestor
<i>P</i>	Phrase-structure parse tree
<i>s</i>	Subtree
<i>w</i>	Word
<i>W</i>	List of words
CCG	Combinatory Categorical Grammar
CKY	Cocke-Kasami-Younger
CoNLL	Conference on Computational Natural Language Learning
GR	Grammatical Relation
IE	Information Extraction
LAS	Labeled Attachment Score
LFG	Lexical Functional Grammar
MST	Maximum Spanning Tree
MT	Machine Translation
PCFG	Probabilistic Context-Free Grammar
PETE	Parser Evaluation using Textual Entailments
PRD	Predicate
QA	Question Answering
RTE	Recognizing Textual Entailments
SBJ	Subject
TMP	Temporal
UAS	Unlabeled Attachment Score

1. INTRODUCTION

Parsing is a basic problem in natural language processing. Broadly it can be defined as assigning a structure to a sentence [1]. Different formalisms have been proposed to represent the sentence structure. However, the common goal is to provide semantic analysis for downstream applications.

Recently, data-driven methods in syntactic parsing have gained prominence. Instead of manually specifying grammar rules, they rely on a treebank of syntactically annotated sentences. The most widely used treebank is the Penn Treebank [2].

Creating a treebank is a non-trivial task. Certain conventions should be set forth. Any disagreements arising during the annotation process should be resolved. In the end, the decisions underlying a treebank amount to a significant volume. For example, bracketing manual for the Penn Treebank [3] has more than 300 pages. Annotators go through a learning period, which may take a few months, before becoming fluent.

The following sections give information on phrase-structure and dependency parsing, which are two prevalent approaches to data-driven syntactic parsing. Recognizing Textual Entailments (RTE) challenge is also introduced.

1.1. Phrase-Structure Parsing

Phrase-structure grammars make use of the constituency concept. A constituent is a group of words behaving as a single unit [1]. An example is the noun phrase “the red apple”. A phrase-structure grammar applies rules on the constituents to derive a sentence. A probabilistic context-free grammar (PCFG) has probabilities for the derivational rules, which are calculated by referring to the treebank.

In this framework, parsing a sentence amounts to finding the derivation with the maximum probability. Most of the current parsers use dynamic programming such as

the Cocke-Kasami-Younger (CKY) or chart parsing algorithms to guide the search. Additionally, a lexicalized grammar takes word forms into account for calculating the probabilities. Head percolation rules [4] are used to find the lexical head.

The standard evaluation method for the phrase-structure framework has been the bracketing precision and recall [5]. A constituent assignment in the parser output is deemed correct if the span of words and the label are the same as in the gold standard parse. The crossing brackets score counts the number of parentheses that are only partially contained within a pair of parentheses in the gold standard parse.

1.2. Dependency Parsing

Dependency grammars focus on relations between two words. The syntactic analysis of sentence is represented by a dependency tree, which is a labeled directed graph [6]. The nodes stand for the words in the sentence and the labels show the dependency types. A dependency tree is projective if there is no crossing dependency.

Two dependency parsers evaluated in this thesis are MaltParser [6] and MSTParser [7]. They differ in the definition of the optimization problem since MaltParser carries out a local search and MSTParser a global search [8]. MaltParser applies a classifier to choose the next best step in a transition-based algorithm. MSTParser uses a large-margin learning algorithm to find the highest-scoring maximum spanning tree.

The standard evaluation metrics for dependency parsing are Unlabeled Attachment Score (UAS) and Labeled Attachment Score (LAS). UAS shows the proportion of words with the correct head. LAS additionally requires dependency labels be correct.

1.3. Recognizing Textual Entailments

Recognizing Textual Entailments (RTE) is a series of challenges carried out in 2005 [9], 2006 [10] and 2007 [11]. The task is to recognize whether a hypothesis text can be reasonably inferred from the source text. RTE proposes a generic framework

for researchers working in different areas such as machine translation (MT) evaluation, question answering (QA) and information extraction (IE).

Source text usually contains multiple sentences, while the hypothesis text is a single sentence. Yes and no answers are equally distributed. Considering the various potential applications, data have been collected from different sources. For example, a gold standard human translation is taken as the source text and an automatic translation is taken as the hypothesis text. There are approximately 1500 source–hypothesis text pairs in development and test data.

The submitted systems use a range of resources and methods [11]. WordNet and the Web are frequently used. Some systems implement logical inference. Some systems carry out anaphora resolution and named entity recognition. Machine learning with lexical and syntactic features have proven useful. The highest score obtained in the third challenge is around 80 per cent.

A relevant observation is that syntax by itself has less than 50 per cent applicability for RTE [12]. In this study human annotators assume that gold standard syntactic parse and a general thesaurus are available and try to find out the upper bound using these resources only. Using only syntax enables to answer 34 per cent of the questions. The availability of a thesaurus raises this number to 48 per cent. By contrast, PETE solely focuses on syntax.

2. LITERATURE SURVEY

PETE is a dependency-based parser evaluation. Instead of specifying a dependency type, words in the relation are represented as a simple sentence, similar to a hypothesis sentence in RTE. This chapter reviews some of the previous work in parser evaluation, with an emphasis on dependency-based methods.

A survey on parser evaluation is given by [13]. This paper summarizes the state-of-the-art up to 1998. Evaluation methods can be classified as non-corpus and corpus-based. Corpus-based methods are further divided based on whether the corpus is annotated or unannotated. Coverage methods measure the percentage of sentences in an unannotated corpus that receive at least one parse. A downside is that a parser returning trivial parses for every sentence would still score high in this method. Average Parse Base method calculates the geometric mean of the number of analyses divided by the number of tokens in a sentence. It shows the amount of ambiguity in the parser grammar, plotted against the sentence length. Its disadvantage is that a low coverage parser would perform well in this measure if the parses are relatively unambiguous. Entropy/Perplexity method applies a probabilistic language model on unannotated corpus and finds out how much a parser captures regularities and decreases ambiguity. Part-of-speech Assignment Accuracy has the advantage that there is already a large amount of part-of-speech tagged corpus. However, many parsers take pre-tagged corpus as input and applicability of this method is low.

This paper also proposes grammatical relations (GR) for parser evaluation [13]. GRs are arranged in a LFG-like notation. A question mark is placed for unspecified information. All possible derivations out of a parse are computed. Evaluation is done by measuring precision and recall on the set of manually annotated GRs. Basically, a grammatical relation shows the syntactic dependency between a head and a dependent. GRs are organized hierarchically. For evaluation SUSANNE corpus was converted into GR scheme by first automatic processing and then manual inspection. On average there are around four GRs per sentence.

Lin [14] converts phrase structures into dependency representations. The proposed evaluation is illustrated using MiniPar, which follows the Minimalist Program. Dependency relations are asymmetric binary relations between a head and a modifier. To represent the dependency tree, a tuple is written for each word. Precision and recall metrics are calculated by comparing tuples in gold standard and response sentences. Dependency labels are ignored.

Gaizauskas et al. [15] propose flatter annotation structures that would have relatively more consensus across grammatical theories. A two-step approach is proposed. First, problematic items are deleted. Second, transformational rules are used to convert annotations into canonical forms. Relations that are common to different parsers are considered. Due to omissions, sentence structures are flatter. Authors acknowledge that divergences may still exist, therefore an additional post-processing step on the parser output may be required. An advantage is that it is easier to develop a corpus with flatter structures. Recall and conformance metrics are used for evaluation. Within this framework the precision metric is unsuitable because the gold standard annotation is minimal. The conformance metric is defined as the proportion of gold standard constituents that are not crossed by any constituent in the response.

De Marneffe et al. [16] extract dependency relations from phrase structure parses. Stanford dependency (SD) relations are similar to GRs [13]. Tree regular expressions are used to produce typed dependencies. Dependencies are extracted and then dependency types are determined. Semantic heads are retrieved instead of syntactic heads. Furthermore, some words such as prepositions and conjunctions are collapsed. Additional links may be inserted for better semantic interpretability. They perform a small-scale evaluation between MiniPar, Link and Stanford parsers using 10 sentences from the Brown corpus. However, such an evaluation is difficult because their dependency representations are different.

Similarly, Miyao et al. [17] convert parser output to GR and SD representations. Conversion is done as a post processing step. Conversion from the Penn Treebank is approximate and problematic. Heuristic approaches are used to overcome problems.

Although SD and GR look similar, conversion between them is not trivial either. 560 sentences, annotated in both GR and SD schemes, from the Penn Treebank test section are used in the experiments. It is estimated that 20 per cent of conversions are problematic.

Clark and Curran [18] evaluate a CCG parser using GRs. Each of 425 lexical categories is mapped to a GR. A post-processing step handles the remaining discrepancies. Gold standard CCGBank data have been used for the development of transformation rules. Evaluation is done on 560 test sentences from the CCGBank version of Penn Treebank. The upper bound for a score that can be obtained in the evaluation is 84 per cent.

In order to get around conversion problems, Tam et al. [19] propose that linguistic features be collected from the parser output and compared to the gold standard. Annotators list the linguistic phenomena for the gold standard set. The most salient phenomenon is taken. Additionally, annotators write the most likely error that parsers would make for each sentence. A recognizer lists the linguistic phenomena in the parser output. Points are given if the recognizer list contains the correct phenomenon or if it does not contain the incorrect phenomenon. A small scale evaluation consisting of 10 sentences is carried out. Its disadvantage is that the linguistic phenomenon may not be recognized even if the parser output is correct.

Rimell and Clark [20] analyze the decisions made by several GR-like evaluation schemes. A decision involves constructions such as subject of verb, direct object of verb, passive voice. Another decision is choosing words that enter a particular construction. A related decision is choosing words to represent the construction. Usually a few words that are deemed important represent a construction that may span more words. Lastly a decision involves the choice of arguments. Empirical study is needed to show which decisions are better in different situations.

Miyao et al. [21] give a task-oriented parser evaluation. They compare recent parsers from different frameworks in a practical application. In total eight depen-

dependency, phrase-structure and deep parsers are evaluated. Another goal is to analyze the domain-adaptability of these data-driven parsers. The task-oriented method uses parser outputs in a machine learning classifier to detect protein–protein interactions. Parse tree indicates a interaction if there is a close syntactic relationship between lexical items. All parsers achieved a higher score than the baseline of bag-of-words. Parser accuracies are similar. Accuracies increase after training with domain specific data. Parsing times vary significantly. Dependency parsing is the fastest, phrase-structure parsing is the slowest and deep parsing is in between. As for the efficacy of the representation, CoNLL format seems to be better than the Penn Treebank format. An ensemble of parsers results in higher accuracy values.

3. GENERATING ENTAILMENTS

In PETE, entailments are designed as minimal sentences, each focusing on one relation. Entailments always include the head and the modifier, plus auxiliary words as necessary for grammaticality. This chapter describes how entailments are generated from dependency parses and phrase-structure parses.

3.1. Generating Entailments From Dependency Parses

Given a dependency parse, the program scans the sentence from left to right. For each words, a decision is made as to whether an entailment should be generated. Words such as determiners are skipped. Pseudocode of the entailment generation algorithm is listed in 3.1.

```

Input:  $D$  a dependency parse tree
Output:  $E$  a list of entailments
for all word  $w$  in  $D$  do
  if  $w$  is not a skip word then
     $h \leftarrow$  head word of  $w$ 
    if  $h$  is a verb then
       $e \leftarrow$  “subject( $h$ ) verb( $h$ ) object( $h$ ) adverb( $h$ )”
    else if  $w$  is appositive then
       $e \leftarrow$  “ $h$  is  $w$ ”
    else
       $e \leftarrow$  “there is  $w$   $h$ ”
    end if
     $E \leftarrow E \cup \{e\}$ 
  end if
end for

```

Figure 3.1. Pseudocode for generating entailments from dependency parses

An entailment is generated with the current word as the target word. A template is selected based on the target word or its head. If the head word is a verb, a simple sentence is formed including the head and the target words. If the target word is an appositive to the head word, a predicative sentence is formed. In the case where the target and head words are inside a noun phrase, an existential construct is formed. In all cases entailments are reorderings of the source sentence words. Few extra words, such as copulae, are added.

3.2. Generating Entailments From Phrase-Structure Parses

There is not a separate program for generating entailments from phrase-structure parses, nor there is any reason to do so. Having a program mimic another program's output is not a trivial task. That would cause unnecessary differences in the generated entailments, and eventually lost points in the evaluation. Instead, phrase-structure parses are converted to dependency parses using the standard conversion program [22] and the dependency program of Section 3.1 is used to generate entailments.

Some information is lost in the conversion. Coordination ambiguity arises due to Melcuk style analysis of coordination [23] in dependency parses. Modification scope can be determined in phrase-structure parses by looking at the attachment level. However, in dependency parses a dependent is always linked to the first conjunct, irrespective of whether it modifies only its immediate head or all conjuncts.

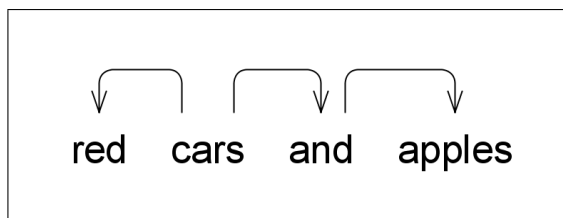


Figure 3.2. Coordination in dependency representation

Figure 3.2 illustrates the coordination ambiguity problem. According to this dependency parse, cars are red. However, it is unclear whether apples are red. Although there are alternative coordination styles that avoid coordination ambiguity, Melcuk

style is preferred because links to a conjunction word are harder to learn for data-driven dependency parsers.

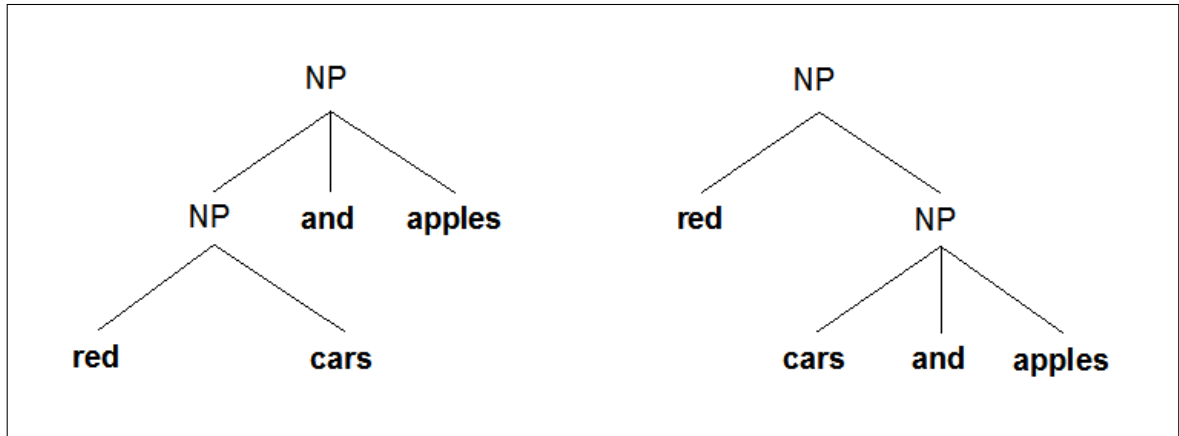


Figure 3.3. Coordination in phrase-structure representation

Figure 3.3 shows two simplified coordination examples in phrase structure representation with different attachment levels. In the left parse tree, red modifies only cars. On the right one, red modifies both cars and apples. After conversion, these two parses map to a single dependency parse, effectively resulting in an ambiguity.

Another point to consider for the conversion is the output format. The Penn Treebank [3] uses function tags such as SBJ (subject) and TMP (temporal) to give additional syntactic and semantic information. Empty nodes show non-local dependencies. Most phrase-structure parsers output a skeletal syntactic structure without function tags and empty nodes. Entailment generation program requires dependency relations, which must be output by the conversion program. The conversion program gets the dependency relations from the function tags of the parse input in turn. To generate entailments correctly either phrase-structure parses should be augmented with function tags or the dependency relations of the conversion output should be corrected.

Two heuristics are used to correct dependency relations. Subject heuristic tries to assign SBJ tag to a nominal dependent of a verb. Predicate heuristic tries to assign PRD tag to a nearest right dependent of a predicative verb. Both heuristics are a few lines of code. They are always applied for a fair evaluation of phrase-structure parsers.

Blaheta’s function-tagger [24] is used. It is based on a relatively simple statistical method and is the first one to recover function tags. It takes skeletal parser output as input and appends function tags to some node labels. It is explicitly stated when the function-tagger is used in the experiments.

Figure 3.4 gives the algorithm that generates the additional entailments lost due to the coordination ambiguity after conversion. Subtrees of conjuncts are replaced to generate a new entailment. If a modifier modifies both conjuncts, it is at a higher level in the parse tree and it is not discarded with the first conjunct. Thus a new entailment showing the modification of the second conjunct is generated.

```

Input:  $D$  a dependency parse tree
          $P$  a phrase structure parse tree
          $E_D$  a list of entailments generated from  $D$ 
Output:  $E$  a list of entailments
for all word  $w$  in  $D$  do
  if  $\exists$  entailment  $e$  for  $w$  in  $E_D$  then
     $W \leftarrow$  words in  $e$ 
     $l \leftarrow$  lowest common ancestor of  $W$  in  $P$ 
    Check for coordination recursively starting at  $l$ 
    if  $\exists$  node  $n$  having coordination then
       $s_i \leftarrow$  subtree of the conjunct  $c_i \in W$ 
       $s_o \leftarrow$  subtree of the conjunct  $c_o \notin W$ 
       $W' \leftarrow$  remove  $s_i$  add  $s_o$ 
       $e' \leftarrow$  entailment with words  $W'$ 
       $E_P \leftarrow E_P \cup \{e'\}$ 
    end if
  end if
end for
 $E \leftarrow E_D \cup E_P$ 

```

Figure 3.4. Generating entailments from phrase-structure parse trees

4. EVALUATION

State-of-the-art phrase-structure and dependency parsers are evaluated. Three phrase-structure parsers are Charniak parser [25], Bikel’s implementation [26] of Collins parser [27], and Stanford parser [28]. Collins parser runs Collins Model 2. Stanford parser is unlexicalized, other two parsers are lexicalized. Pre-trained models with basic configurations are used.

Dependency parsers evaluated in this study are MSTParser [7] and MaltParser [6]. They were trained with the dependency data converted from Penn Treebank training section using the LTH converter [22]. Both parsers run in projective mode.

4.1. Evaluation on the WSJ Corpus

Table 4.1 shows the bracketing scores on the WSJ section 23. There are 2399 sentences in the test set. Sentence lengths are not considered in the reported results.

Table 4.1. Bracketing scores of phrase-structure parsers on the WSJ test section

System	Recall	Precision	F-measure
Charniak	89.57	89.91	89.74
Collins	88.13	88.26	88.19
Stanford	85.09	86.52	85.80

Table 4.2. LAS and UAS scores of dependency parsers on the WSJ test section

System	UAS	LAS
MST	92.0	88.7
Malt	89.8	86.8

Similarly, LAS and UAS scores are listed in Table 4.2.

Dependency-based evaluation of all parsers is given in Table 4.3. Phrase-structure parser outputs are converted to dependency representation using LTH converter. For dependency parsers, scores are the same as in Table 4.2. Charniak parser is best at getting the unlabeled dependency tree correct, as shown in the UAS column. The large gap between UAS and LAS scores show that phrase-structure parsers need to incorporate function tags in their outputs. The gap is 8.3 for Charniak parser and 3.3 for MSTParser. Blaheta’s function-tagger does a good job by reducing the gap to less than 4.0. Table 4.3 is not a fair comparison since coordination information in phrase-structure parses is not considered.

Table 4.3. LAS and UAS scores on the WSJ test section

System	UAS	LAS	UAS–LAS
Charniak + F-Tags	93.2	89.6	3.6
Charniak	93.0	84.7	8.3
MST	92.0	88.7	3.3
Collins + F-Tags	91.6	87.7	3.9
Collins	91.4	83.1	8.3
Stanford + F-Tags	90.3	86.5	3.8
Stanford	90.2	82.1	8.1
Malt	89.8	86.8	3.0

PETE scores are given in Table 4.4. The gold standard set of entailments are generated by applying the algorithm in Figure 3.1 to the gold standard dependency parse and the algorithm in Figure 3.4 to the gold standard phrase-structure parse. The first three lines indicate upper bounds. Recall upper bound for dependency parsers is less than 100 per cent because they miss the additional coordination entailments generated from the gold standard phrase-structure parse. On the other hand, phrase-structure parsers lose points due to their impoverished output format. For second and third rows, function tags and empty nodes are removed from the gold standard phrase-structure parse. Imperfect conversion determines the upper bound in this case. The function-tagger increases the score by four per cent. In PETE evaluation, a full

phrase-structure parser [29], not included in this study, has 100 per cent upper bound.

The recall is more relevant than f-measure. The decision of when to return an answer is generally systematically determined. Dependency parsers do not return an answer for the additional coordination entailments. Consequently their precision scores are higher than their recall scores. F-measures are relevant for parsers using the same framework.

Two phrase-structure parsers got the highest scores in PETE evaluation. The function-tagger improves the results, although its effect is less pronounced than in LAS scores. Charniak parser without the function-tagger performs better than the best performing dependency parser.

Table 4.4. PETE scores on the WSJ test section

System	Recall	Precision	F-measure
Gold Dep.	98.12	100.00	99.05
Gold Phrase + F-Tags	96.40	96.26	96.33
Gold Phrase	92.42	92.02	92.22
Charniak + F-Tags	83.34	83.35	83.34
Collins + F-Tags	81.60	82.16	81.88
Charniak	80.57	80.33	80.45
MST	79.24	81.17	80.19
Collins	78.99	79.35	79.17
Malt	77.66	80.82	79.21
Stanford + F-Tags	77.58	77.63	77.60
Stanford	75.10	75.02	75.06

4.2. Evaluation on the Brown Corpus

Whereas WSJ corpus of the Penn Treebank includes only financial texts, Brown corpus includes various genres such as fiction and natural sciences. The rationale for

evaluation on Brown corpus is to see how much parsers overfit to the WSJ data and how well they perform in a new domain. There are 425 sentences in the test set, the same sentences as in the CoNLL 2008 shared task data set [30].

Table 4.5 gives the bracketing scores. While the ranking is the same as in Table 4.1 the gap between Stanford parser and Collins parser has shrunk.

Table 4.5. Bracketing scores of phrase-structure parsers on the Brown test section

System	Recall	Precision	F-measure
Charniak	85.74	85.90	85.82
Collins	83.57	83.91	83.74
Stanford	83.00	83.01	83.00

Attachment scores of dependency parsers are given in Table 4.6. While the scores are lower than in Table 4.2, the ranking has not changed.

Table 4.6. LAS and UAS scores of dependency parsers on the Brown test section

System	UAS	LAS
MST	88.2	81.9
Malt	85.2	79.2

PETE scores are given in Table 4.7. Coordination ambiguity has a slightly larger role in the Brown corpus since the upper bound of recall for dependency parsing is 96.75 per cent compared to 98.12 per cent. There are some notable differences in the ranking. Charniak parser performs even better. Its score is still higher than other parsers without using the function-tagger. Stanford parser has risen in the ranking, probably its unlexicalized model helps avoid overfitting. Dependency parsers seem to be worst affected by the overfitting problem since the decline in scores cannot be explained by the decreased upper bound. Their rich feature models may make them more susceptible to overfitting.

Table 4.7. PETE scores on the Brown test section

System	Recall	Precision	F-measure
Gold Dep.	96.75	100.00	98.35
Gold Phrase + F-Tags	94.79	94.89	94.84
Gold Phrase	89.82	89.26	89.54
Charniak + F-Tags	74.89	74.64	74.76
Charniak	72.24	71.60	71.92
Collins + F-Tags	71.38	72.68	72.02
Collins	68.76	69.68	69.22
Stanford + F-Tags	68.23	68.61	68.42
MST	67.66	70.98	69.28
Stanford	66.60	66.78	66.69
Malt	65.87	71.44	68.54

5. CONCLUSION

PETE has been proposed as a cross-framework parser evaluation scheme. State-of-the-art phrase-structure and dependency parsers have been evaluated using PETE. A program generates entailments from dependency parses. Phrase-structure parses are converted to dependency parses to generate entailments. Additional entailments are generated for unambiguous coordinations in phrase-structure parses. Charniak parser has ranked first in the evaluation.

Unlike other dependency-based evaluations, PETE is designed with annotation in mind. Traditionally, annotators go through a training period to acquire the intricacies of an annotation scheme. Limited number of annotators are usually based in the same institution. Consequently, annotation is time consuming and expensive. PETE offers a new annotation model. A priori training is not required. Many people can participate in the annotation process over the Web. Although annotation should not take significantly more time than the comprehension time of a simple sentence, trial annotations are required to compare the speed of annotation.

Evaluation methods are expected to provide a clear assessment of current technologies and directions for future research. The detrimental effect of impoverished output format of phrase-structure parsers have been noted. However, the standard bracketing evaluation does not create incentives for a richer output. Therefore most state-of-the-art phrase-structure parsers output only a skeletal parse tree. PETE clearly shows that incorporating functions tags and empty nodes is beneficial.

APPENDIX A: SAMPLE ENTAILMENTS

- Ms. Haag plays Elianti.
 - There is Ms. Haag.
 - Ms. Haag plays X.
 - X plays Elianti.
- Rolls-Royce Motor Cars Inc. said it expects its U.S. sales to remain steady at about 1,200 cars in 1990.
 - There are Rolls-Royce Motor Cars Inc.
 - Rolls-Royce Motor Cars Inc. said X.
 - It expects X.
 - X said X expects X.
 - There are its U.S. sales.
 - X expects its sales to remain steady.
 - X expects X to remain steady.
 - X are to remain steady.
 - X are to remain X at about 1,200 cars.
 - X are to remain X in 1990.
- Companies would be compelled to publish in annual proxy statements the names of insiders who fail to file reports on time.
 - Companies would be compelled to publish X.
 - There are annual proxy statements.
 - X would be compelled to publish X in statements.
 - X would be compelled to publish the names.
 - There are the names of insiders.
 - Insiders fail to file X.
 - X fail to file reports.
 - X fail to file X on time.

REFERENCES

1. Jurafsky, D. and J. Martin, *Speech and Language Processing*, Prentice Hall, 2008.
2. Marcus, M., G. Kim, M. Marcinkiewicz, R. MacIntyre, A. Bies, M. Ferguson, K. Katz, and B. Schasberger, “The Penn Treebank: Annotating Predicate Argument Structure”, *ARPA Human Language Technology Workshop*, 1994.
3. Bies, A., M. Ferguson, K. Katz, and R. MacIntyre, “Bracketing Guidelines for Treebank II Style Penn Treebank Project”, Technical report, University of Pennsylvania, 1995.
4. Collins, M., *Head-Driven Statistical Models for Natural Language Parsing*, Ph.D. thesis, University of Pennsylvania, 1999.
5. Abney, S., S. Flickenger, C. Gdaniec, C. Grishman, P. Harrison, D. Hindle, R. Ingria, F. Jelinek, J. Klavans, and M. Liberman, “Procedure for Quantitatively Comparing the Syntactic Coverage of English Grammars”, *Proceedings of the Workshop on Speech and Natural Language*, pp. 306–311, Association for Computational Linguistics Morristown, NJ, USA, 1991.
6. Nivre, J., J. Hall, J. Nilsson, A. Chanev, G. Eryigit, S. Kübler, S. Marinov, and E. Marsi, “MaltParser: A language-independent system for data-driven dependency parsing”, *Natural Language Engineering*, Vol. 13, No. 2, pp. 95–135, 2007.
7. McDonald, R., K. Crammer, and F. Pereira, “Online Large-Margin Training of Dependency Parsers”, *Proceedings of ACL*, pp. 91–98, Association for Computational Linguistics, Ann Arbor, Michigan, June 2005.
8. McDonald, R., J. Nivre, and S. Sweden, “Characterizing the Errors of Data-Driven Dependency Parsing Models”, *Proceedings of EMNLP-CoNLL*, pp. 122–131, 2007.

9. Dagan, I., O. Glickman, and B. Magnini, “The PASCAL Recognising Textual Entailment Challenge”, Quinonero-Candela, J., I. Dagan, B. Magnini, and F. d’Alché Buc (editors), *Machine Learning Challenges. Lecture Notes in Computer Science*, Vol. 3944, pp. 177–190, Springer, 2006.
10. Bar-Haim, R., I. Dagan, B. Dolan, L. Ferro, D. Giampiccolo, B. Magnini, and I. Szpektor, “The Second PASCAL Recognising Textual Entailment Challenge”, *Proceedings of the Second PASCAL Challenges Workshop on Recognising Textual Entailment*, pp. 1–9, 2006.
11. Giampiccolo, D., B. Magnini, I. Dagan, and B. Dolan, “The Third PASCAL Recognizing Textual Entailment Challenge”, *Proceedings of the ACL-PASCAL Workshop on Textual Entailment and Paraphrasing*, pp. 1–9, 2007.
12. Vanderwende, L. and W. Dolan, “What Syntax Can Contribute in the Entailment Task”, *Lecture Notes In Computer Science*, Vol. 3944, p. 205, 2006.
13. Carroll, J., T. Briscoe, and A. Sanfilippo, “Parser Evaluation: a Survey and a New Proposal”, *Proceedings of LREC*, pp. 447–454, 1998.
14. Lin, D., “Dependency-Based Evaluation of MINIPAR”, *Proceedings of the LREC Workshop on Evaluation of Parsing Systems*, Granada, Spain, 1998.
15. Gaizauskas, R., M. Hepple, and C. Huyck, “A Scheme for Comparative Evaluation of Diverse Parsing Systems”, *Proceedings of LREC*, pp. 143–149, 1998.
16. de Marneffe, M., B. MacCartney, and C. Manning, “Generating Typed Dependency Parses From Phrase Structure Parses”, *LREC*, 2006.
17. Miyao, Y., K. Sagae, and J. Tsujii, “Towards Framework-Independent Evaluation of Deep Linguistic Parsers”, King, T. H. and E. M. Bender (editors), *Proceedings of the Grammar Engineering Across Frameworks Workshop (GEAF-07)*, pp. 238–258, CSLI, Stanford, CA, 2007.

18. Clark, S. and J. Curran, “Formalism-Independent Parser Evaluation with CCG and DepBank”, *Proceedings of ACL*, pp. 248–255, Association for Computational Linguistics, Prague, Czech Republic, June 2007.
19. Tam, W. L., Y. Sato, Y. Miyao, and J. Tsujii, “Parser Evaluation Across Frameworks without Format Conversion”, *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, pp. 29–35, Coling 2008 Organizing Committee, Manchester, UK, August 2008.
20. Rimell, L. and S. Clark, “Constructing a Parser Evaluation Scheme”, *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, pp. 44–50, Manchester, UK, August 2008.
21. Miyao, Y., R. Sætre, K. Sagae, T. Matsuzaki, and J. Tsujii, “Task-Oriented Evaluation of Syntactic Parsers and Their Representations”, *Proceedings of ACL-HLT*, pp. 46–54, 2008.
22. Johansson, R. and P. Nugues, “Extended Constituent-to-dependency Conversion for English”, *Proceedings of NODALIDA 2007*, Tartu, Estonia, May 25-26 2007.
23. Melcuk, I. A., *Dependency Syntax: Theory and Practice*, State University Press of New York, 1988.
24. Blaheta, D., *Function Tagging*, Ph.D. thesis, Brown University, 2003.
25. Charniak, E., “A Maximum-Entropy-Inspired Parser”, *Proceedings of NAACL*, pp. 132–139, 2000.
26. Bikel, D., “Design of a Multi-lingual, Parallel-processing Statistical Parsing Engine”, *Proceedings of HLT*, pp. 178–182, Morgan Kaufmann Publishers Inc. San Francisco, CA, USA, 2002.
27. Collins, M., “Head-Driven Statistical Models for Natural Language Parsing”, *Computational Linguistics*, Vol. 29, No. 4, pp. 589–637, 2003.

28. Klein, D. and C. Manning, “Accurate Unlexicalized Parsing”, *Proceedings of ACL*, pp. 423–430, Association for Computational Linguistics, 2003.
29. Gabbard, R., S. Kulick, and M. Marcus, “Fully Parsing The Penn Treebank”, *Proceedings of HLT-NAACL*, pp. 184–191, 2006.
30. Surdeanu, M., R. Johansson, A. Meyers, L. Màrquez, and J. Nivre, “The CoNLL-2008 Shared Task on Joint Parsing of Syntactic and Semantic Dependencies”, *Proceedings of CoNLL*, 2008.