

# Figures of Merit for Best-First Probabilistic Chart Parsing

Sharon A. Caraballo and Eugene Charniak  
Brown University  
{sc,ec}@cs.brown.edu

## Abstract

Best-first parsing methods for natural language try to parse efficiently by considering the most likely constituents first. Some figure of merit is needed by which to compare the likelihood of constituents, and the choice of this figure has a substantial impact on the efficiency of the parser. While several parsers described in the literature have used such techniques, there is no published data on their efficacy, much less attempts to judge their relative merits. We propose and evaluate several figures of merit for best-first parsing.

## Introduction

Chart parsing is a commonly-used algorithm for parsing natural language texts. The chart is a data structure which contains all of the constituents which may occur in the sentence being parsed. At any point in the algorithm, there exist constituents which have been proposed but not actually included in a parse. These proposed constituents are stored in a data structure called the keylist. When a constituent is removed from the keylist, the system considers how this constituent can be used to extend its current structural hypothesis. In general this can lead to the creation of new, more encompassing constituents which themselves are then added to the keylist. When we are finished processing one constituent, a new one is chosen to be removed from the keylist, and so on. Traditionally, the keylist is represented as a stack, so that the last item added to the keylist is the next one removed.

Best-first chart parsing is a variation of chart parsing which attempts to find the most likely parses first, by adding constituents to the chart in order of the likelihood that they will appear in a correct parse, rather than simply popping constituents off of a stack. Some figure of merit is assigned to potential constituents, and the constituent maximizing this value is the next to be added to the chart.

In best-first probabilistic chart parsing a probabilistic measure is used. In this paper we consider probabilities primarily based on probabilistic context-free grammars, though in principle other, more complicated schemes could be used.

Ideally, we would like to use as our figure of merit the conditional probability of that constituent, given the entire sentence, in order to choose a constituent that not only appears likely in isolation, but maximizes the likelihood of the sentence as a whole; that is, we would like to pick the constituent that maximizes the following quantity:

$$p(N_{j,k}^i | t_{0,n})$$

where  $t_{0,n}$  is the sequence of the  $n$  tags, or parts of speech, in the sentence (numbered  $t_0, \dots, t_{n-1}$ ), and  $N_{j,k}^i$  is a nonterminal of type  $i$  covering terms  $t_j \dots t_{k-1}$ . However, we cannot calculate this quantity, since in order to do so, we would need to completely parse the sentence. In this paper, we examine the performance of several proposed figures of merit that approximate it in one way or another.

In our experiments, we use only tag sequences for parsing. More accurate probability estimates should be attainable using lexical information.

## Figures of Merit

### Straight $\beta$

It seems reasonable to base a figure of merit on the inside probability  $\beta$  of the constituent. Inside probability is defined as the probability of the words or tags in the constituent given that the constituent is dominated by a particular nonterminal symbol. This seems to be a reasonable basis for comparing constituent probabilities, and has the additional advantage that it is easy to compute during chart parsing.

The inside probability of the constituent  $N_{j,k}^i$  is defined as

$$\beta(N_{j,k}^i) \equiv p(t_{j,k}|N^i)$$

where  $N^i$  represents the  $i$ th nonterminal symbol.

In terms of our earlier discussion, our “ideal” figure of merit can be rewritten as:

$$\begin{aligned} p(N_{j,k}^i|t_{0,n}) &= \frac{p(N_{j,k}^i, t_{0,n})}{p(t_{0,n})} \\ &= \frac{p(N^i, j, k, t_{0,j}, t_j, k, t_k, n)}{p(t_{0,n})} \\ &= \frac{p(t_{0,j}, N_{j,k}^i, t_k, n)p(t_{j,k}|t_{0,j}, N_{j,k}^i, t_k, n)}{p(t_{0,n})}. \end{aligned}$$

We apply the usual independence assumption that given a nonterminal, the tag sequence it generates depends only on that nonterminal, giving

$$\begin{aligned} p(N_{j,k}^i|t_{0,n}) &\approx \frac{p(t_{0,j}, N_{j,k}^i, t_k, n)p(t_{j,k}|N_{j,k}^i)}{p(t_{0,n})} \\ &= \frac{p(t_{0,j}, N_{j,k}^i, t_k, n)\beta(N_{j,k}^i)}{p(t_{0,n})}. \end{aligned}$$

The first term in the numerator is just the definition of the outside probability  $\alpha$  of the constituent. Outside probability  $\alpha$  of a constituent  $N_{j,k}^i$  is defined as the probability of that constituent and the rest of the words in the sentence (or rest of the tags in the tag sequence, in our case).

$$\alpha(N_{j,k}^i) \equiv p(t_{0,j}, N_{j,k}^i, t_k, n).$$

We can therefore rewrite our ideal figure of merit as

$$p(N_{j,k}^i|t_{0,n}) \approx \frac{\alpha(N_{j,k}^i)\beta(N_{j,k}^i)}{p(t_{0,n})}.$$

In this equation, we can see that  $\alpha(N_{j,k}^i)$  and  $p(t_{0,n})$  represent the influence of the surrounding words. Thus using  $\beta$  alone assumes that  $\alpha$  and  $p(t_{0,n})$  can be ignored.

We will refer to this figure of merit as **straight  $\beta$** .

### Normalized $\beta$

One side effect from omitting the  $\alpha$  and  $p(t_{0,n})$  terms in the  $\beta$ -only figure above is that inside probability alone tends to prefer shorter constituents to longer ones, as the inside probability of a longer constituent involves the product of

more probabilities. This can result in a “thrashing” effect, where the system parses short constituents, even very low probability ones, while avoiding combining them into longer constituents. To avoid thrashing, typically some technique is used to normalize the inside probability for use as a figure of merit. One approach is to take the geometric mean of the inside probability, to obtain a “per-word” inside probability. (In the “ideal” model, the  $p(t_{0,n})$  term acts as a normalizing factor.)

The per-word inside probability of the constituent  $N_{j,k}^i$  is calculated as

$${}^{k-j}\sqrt{\beta(N_{j,k}^i)}.$$

We will refer to this figure as **normalized  $\beta$** .

### Normalized $\alpha_L\beta$

In the previous section, we showed that our ideal figure of merit can be written as

$$p(N_{j,k}^i|t_{0,n}) \approx \frac{\alpha(N_{j,k}^i)\beta(N_{j,k}^i)}{p(t_{0,n})}.$$

However, the  $\alpha$  term, representing outside probability, cannot be calculated directly during a parse, since we need the full parse of the sentence to compute it. In some of our figures of merit, we use the quantity  $p(N_{j,k}^i, t_{0,j})$ , which is closely related to outside probability. We call this quantity the left outside probability, and denote it  $\alpha_L$ .

The following recursive formula can be used to compute  $\alpha_L$ . Let  $\mathcal{E}_{j,k}^i$  be the set of all completed edges, or rule expansions, in which the nonterminal  $N_{j,k}^i$  appears. For each edge  $e$  in  $\mathcal{E}_{j,k}^i$ , we compute the product of  $\alpha_L$  of the nonterminal appearing on the left-hand side (lhs) of the rule, the probability of the rule itself, and  $\beta$  of each nonterminal  $N_{r,s}^q$  appearing to the left of  $N_{j,k}^i$  in the rule. Then  $\alpha_L(N_{j,k}^i)$  is the sum of these products:

$$\begin{aligned} \alpha_L(N_{j,k}^i) &= \sum_{e \in \mathcal{E}_{j,k}^i} \alpha_L(N_{\text{start}(e), \text{end}(e)}^{\text{lhs}(e)}) p(\text{rule}(e)) \prod_{N_{r,s}^q} \beta(N_{r,s}^q). \end{aligned}$$

This formula can be infinitely recursive, depending on the properties of the grammar. A method for calculating  $\alpha_L$  more efficiently can be derived from the calculations given in (Jelinek and Lafferty, 1991).

A simple extension to the normalized  $\beta$  model allows us to estimate the per-word probability of all tags in the sentence through the end of the

constituent under consideration. This allows us to take advantage of information already obtained in a left-right parse. We calculate this quantity as follows:

$$\sqrt[k]{\alpha_L(N_{j,k}^i)\beta(N_{j,k}^i)}.$$

We are again taking the geometric mean to avoid thrashing by compensating for the  $\alpha\beta$  quantity's preference for shorter constituents, as explained in the previous section.

We refer to this figure of merit as **normalized  $\alpha_L\beta$** .

### Trigram estimate

An alternative way to rewrite the "ideal" figure of merit is as follows:

$$\begin{aligned} p(N_{j,k}^i|t_{0,n}) &= \frac{p(N_{j,k}^i, t_{0,n})}{p(t_{0,n})} \\ &= \frac{p(t_{0,j}, t_{k,n})p(N_{j,k}^i|t_{0,j}, t_{k,n})p(t_{j,k}|N_{j,k}^i, t_{0,j}, t_{k,n})}{p(t_{0,j}, t_{k,n})p(t_{j,k}|t_{0,j}, t_{k,n})}. \end{aligned}$$

Once again applying the usual independence assumption that given a nonterminal, the tag sequence it generates depends only on that nonterminal, we can rewrite the figure of merit as follows:

$$p(N_{j,k}^i|t_{0,n}) \approx \frac{p(N_{j,k}^i|t_{0,j}, t_{k,n})\beta(N_{j,k}^i)}{p(t_{j,k}|t_{0,j}, t_{k,n})}.$$

To derive an estimate of this quantity for practical use as a figure of merit, we make some additional independence assumptions. We assume that  $p(N_{j,k}^i|t_{0,j}, t_{k,n}) \approx p(N_{j,k}^i)$ , that is, that the probability of a nonterminal is independent of the tags before and after it in the sentence. We also use a trigram model for the tags themselves, giving  $p(t_{j,k}|t_{0,j}, t_{k,n}) \approx p(t_{j,k}|t_{j-2,j})$ . Then we have:

$$p(N_{j,k}^i|t_{0,n}) \approx \frac{p(N^i)\beta(N_{j,k}^i)}{p(t_{j,k}|t_{j-2,j})}.$$

We can calculate  $\beta(N_{j,k}^i)$  as usual. The  $p(N^i)$  term is estimated from our PCFG as the sum of the counts for all rules having  $N^i$  as their left-hand side, divided by the sum of the counts for all rules. The  $p(t_{j,k}|t_{j-2,j})$  term is just the probability of the tag sequence  $t_j \dots t_{k-1}$  according to a trigram model.<sup>1</sup> (Technically, this is not a trigram model but a tritag model, since we are considering sequences of tags, not words.) We refer to this model as the **trigram estimate**.

<sup>1</sup>Our results show that the  $p(N^i)$  term can be omitted without much effect.

### Prefix estimate

We also derived an estimate of the ideal figure of merit which takes advantage of statistics on the first  $j-1$  tags of the sentence as well as  $t_{j,k}$ . This estimate represents the probability of the constituent in the context of the preceding tags.

$$\begin{aligned} p(N_{j,k}^i|t_{0,n}) &= \frac{p(N_{j,k}^i, t_{0,n})}{p(t_{0,n})} \\ &= \frac{p(t_{k,n})p(N_{j,k}^i, t_{0,j}|t_{k,n})p(t_{j,k}|N_{j,k}^i, t_{0,j}, t_{k,n})}{p(t_{k,n})p(t_{0,k}|t_{k,n})} \\ &= \frac{p(N_{j,k}^i, t_{0,j}|t_{k,n})p(t_{j,k}|N_{j,k}^i, t_{0,j}, t_{k,n})}{p(t_{0,k}|t_{k,n})}. \end{aligned}$$

We again make the independence assumption that  $p(t_{j,k}|N_{j,k}^i, t_{0,j}, t_{k,n}) \approx \beta(N_{j,k}^i)$ . Additionally, we assume that  $p(N_{j,k}^i, t_{0,j})$  and  $p(t_{0,k})$  are independent of  $p(t_{k,n})$ , giving

$$p(N_{j,k}^i|t_{0,n}) \approx \frac{p(N_{j,k}^i, t_{0,j})\beta(N_{j,k}^i)}{p(t_{0,k})}.$$

The denominator,  $p(t_{0,k})$ , is once again calculated from a tritag model. The  $p(N_{j,k}^i, t_{0,j})$  term is just  $\alpha_L$ , defined above in the discussion of the normalized  $\alpha_L\beta$  model. Thus this figure of merit can be written as

$$\frac{\alpha_L(N_{j,k}^i)\beta(N_{j,k}^i)}{p(t_{0,k})}.$$

We will refer to this as the **prefix estimate**.

## The Experiment

We used as our grammar a probabilistic context-free grammar learned from the Brown corpus (see (Francis and Kučera, 1982), Carroll and Charniak (1992a) and (1992b), and (Charniak and Carroll, 1994)). We parsed 500 sentences of length 3 to 30 (including punctuation) from the Penn Treebank Wall Street Journal corpus using a best-first parsing method and each of the following estimates for  $p(N_{j,k}^i|t_{0,n})$  as the figure of merit:

1. straight  $\beta$
2. normalized  $\beta$
3. normalized  $\alpha_L\beta$
4. trigram estimate
5. prefix estimate

The probability  $p(N^i)$  in the trigram estimate was determined from the same training data from which our grammar was learned initially. Our trigram probabilities for the trigram and prefix estimates were learned from this data as well, using the deleted interpolation method for smoothing.

For each figure of merit, we compared the performance of best-first parsing using that figure of merit to exhaustive parsing. By exhaustive parsing, we mean continuing to parse until there are no more constituents available to be added to the chart. We parse exhaustively to determine the total probability of a sentence, that is, the sum of the probabilities of all parses found for that sentence.

We then computed several quantities for best-first parsing with each figure of merit at the point where the best-first parsing method has found parses contributing at least 95% of the probability mass of the sentence.

## Results

The chart below presents the following measures for each figure of merit:

1. %E: The percentage of edges, or rule expansions, in the exhaustive parse that have been used by the best-first parse to get 95% of the probability mass. Edge creation is generally considered the best measure of CFG parser effort.
2. %non-0 E: The percentage of nonzero-length edges used by the best-first parse to get 95%. Zero-length edges are required by our parser as a book-keeping measure, and as such are virtually un-eliminable. We anticipated that removing them from consideration would highlight the “true” differences in the figures of merit.
3. %popped: The percentage of constituents in the exhaustive parse that were used by the best-first parse to get 95% of the probability mass.

Figure of Merit	%E	%non-0 E	%popped
straight $\beta$	97.6	97.5	93.8
normalized $\beta$	34.7	31.6	61.5
normalized $\alpha_L\beta$	39.7	36.4	57.3
trigram estimate	25.2	21.7	44.3
prefix estimate	21.8	17.4	38.3

The statistics converged to their final values quickly. The edge-count percentages were generally within .01 of their final values after processing only 200 sentences, so the results were quite stable by the end of our 500-sentence test corpus.

We gathered statistics for each sentence length from 3 to 30. Sentence length was limited to a maximum of 30 because of the huge number of edges that are generated in doing a full parse of

long sentences; using this grammar, sentences in this length range have produced up to 130,000 edges. Figure 1 shows a graph of %non-0 E, that is, the percent of nonzero-length edges needed to get 95% of the probability mass, for each sentence length.

We also measured the total CPU time (in seconds) needed to get 95% of the probability mass for each of the 500 sentences. The results are presented in the following chart:

Figure of Merit	CPU time
straight $\beta$	3966
normalized $\beta$	1631
normalized $\alpha_L\beta$	68660
trigram estimate	1547
prefix estimate	26520

Figure 2 shows the average CPU time to get 95% of the probability mass for each estimate and each sentence length. Each estimate averaged below 1 second on sentences of fewer than 7 words. (The  $y$ -axis has been restricted so that the normalized  $\beta$  and trigram estimates can be better compared.)

## Previous work

The literature shows many implementations of best-first parsing, but none of the previous work shares our goal of explicitly comparing figures of merit.

Bobrow (1990) and Chitrao and Grishman (1990) introduced statistical agenda-based parsing techniques. Chitrao and Grishman implemented a best-first probabilistic parser and noted the parser’s tendency to prefer shorter constituents. They proposed a heuristic solution of penalizing shorter constituents by a fixed amount per word.

Miller and Fox (1994) compare the performance of parsers using three different types of grammars, and show that a probabilistic context-free grammar using inside probability (unnormalized) as a figure of merit outperforms both a context-free grammar and a context-dependent grammar.

Kochman and Kupin (1991) propose a figure of merit closely related to our prefix estimate. They do not actually incorporate this figure into a best-first parser.

Magerman and Marcus (1991) use the geometric mean to compute a figure of merit that is independent of constituent length. Magerman and Weir (1992) use a similar model with a different parsing algorithm.

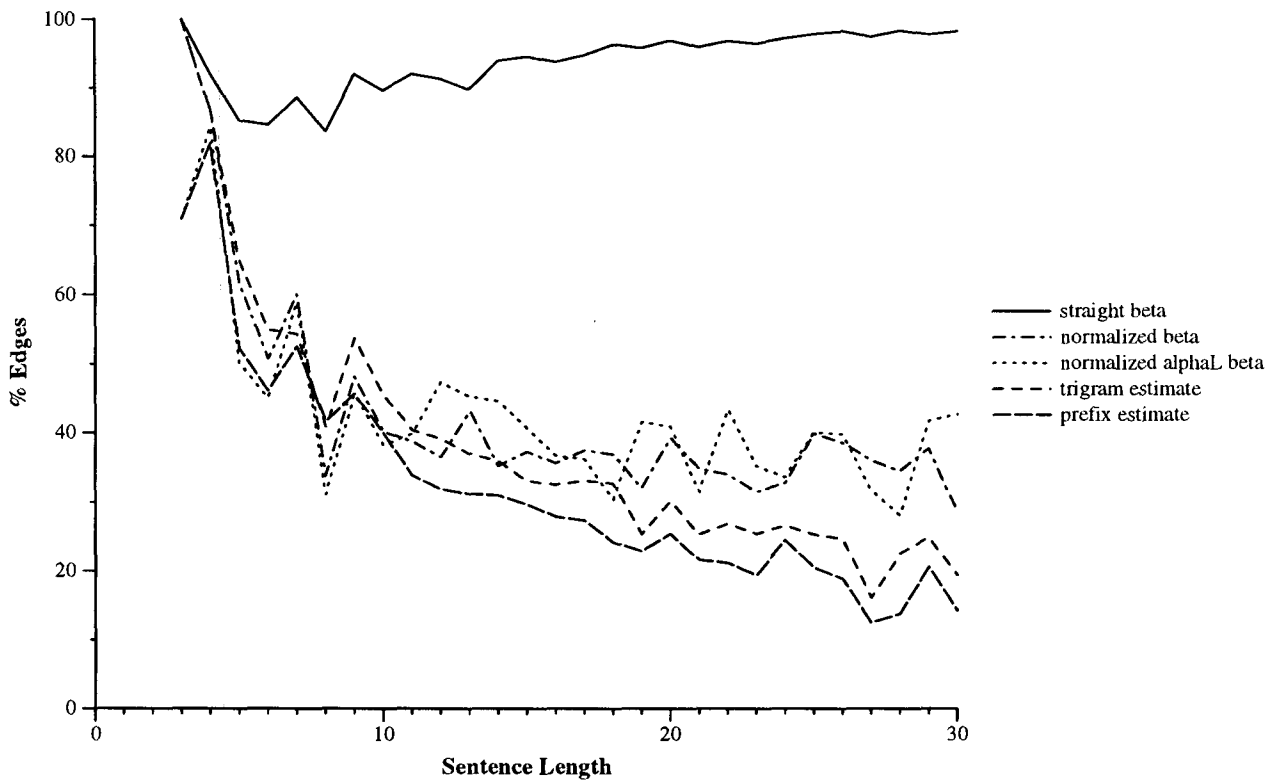


Figure 1: Nonzero-length edges for 95% of Probability Mass

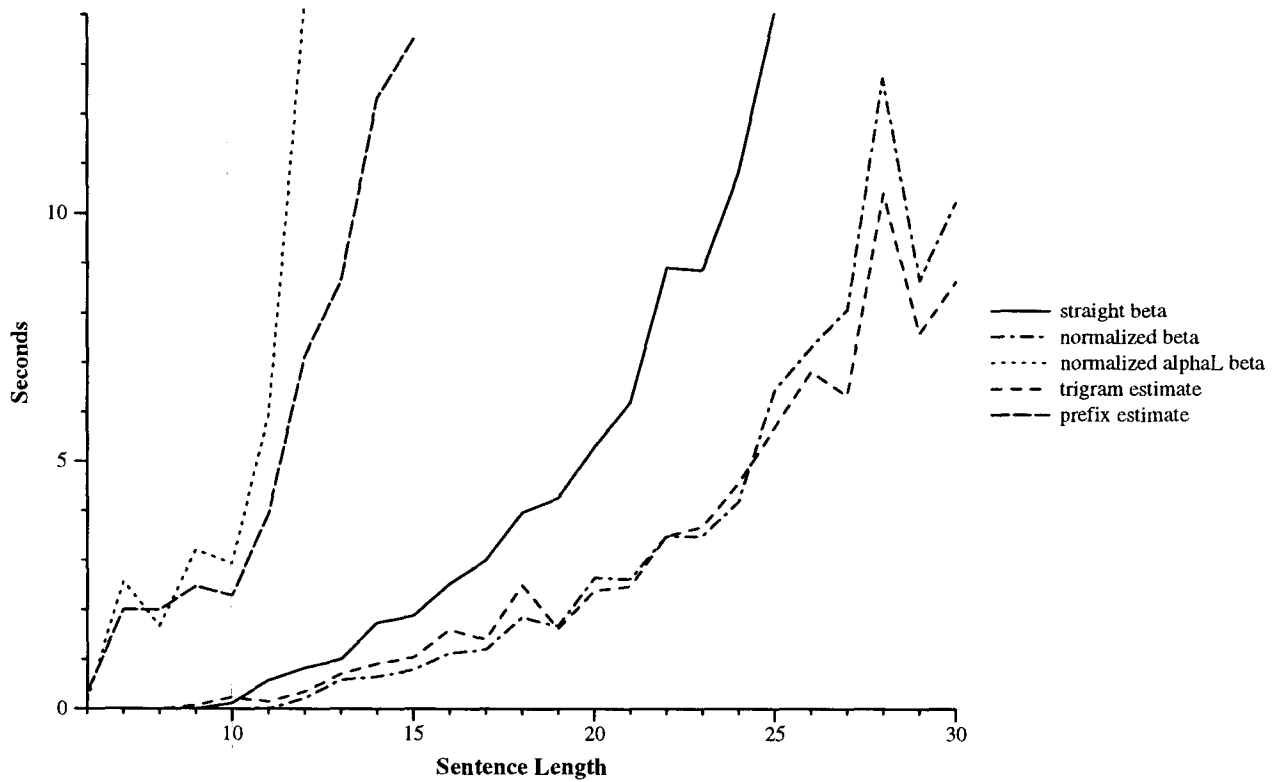


Figure 2: Average CPU Time for 95% of Probability Mass

## Conclusions

From the edge count statistics, it is clear that straight  $\beta$  is a poor figure of merit. Figure 1 also demonstrates that its performance generally worsens as sentence length increases.

The best performance in terms of edge counts of the figures we tested was the model which used the most information available from the sentence, the prefix model. However, so far, the additional running time needed for the computation of  $\alpha_L$  terms has exceeded the time saved by processing fewer edges, as is made clear in the CPU time statistics, where these two models perform substantially worse than even the straight  $\beta$  figure.

While chart parsing and calculations of  $\beta$  can be done in  $O(n^3)$  time, we have been unable to find an algorithm to compute the  $\alpha_L$  terms faster than  $O(n^5)$ . When a constituent is removed from the keylist, it only affects the  $\beta$  values of its ancestors in the parse trees; however,  $\alpha_L$  values are propagated to all of the constituent's siblings to the right and all of its descendants. Recomputing the  $\alpha_L$  terms when a constituent is removed from the keylist can be done in  $O(n^3)$  time, and since there are  $O(n^2)$  possible constituents, the total time needed to compute the  $\alpha_L$  terms in this manner is  $O(n^5)$ .

The best performer in running time was the parser using the trigram estimate as a figure of merit. This figure has the additional advantage that it can be easily incorporated into existing best-first parsers using a figure of merit based on inside probability. From the CPU time statistics, it can be seen that the running time begins to show a real improvement over the normalized  $\beta$  model on sentences of length 25 or greater, and the trend suggests that the improvement would be greater for longer sentences.

It is also interesting to note that while the models using figures of merit normalized by the geometric mean performed similarly to the other models on shorter sentences, the superior performance of the other models becomes more pronounced as sentence length increases. From Figure 1, we can see that the models using the geometric mean appear to level off with respect to an exhaustive parse when used to parse sentences of length greater than about 15. The other two estimates seem to continue improving with greater sentence length. In fact, the measurements presented here almost certainly underestimate the true benefits of the better models. We restricted sentence length to a maximum of 30 words, in order to keep the number of edges in the exhaustive parse to a practical size; however, since the percentage of edges needed by the best-first parse decreases with increasing sentence length, we assume that the im-

provement would be even more dramatic for sentences longer than 30 words.

## References

- [1990] Robert J. Bobrow. 1990. Statistical agenda parsing. In *DARPA Speech and Language Workshop*, pages 222–224.
- [1992a] Glenn Carroll and Eugene Charniak. 1992a. Learning probabilistic dependency grammars from labeled text. In *Working Notes*, Fall Symposium Series, pages 25–32. AAAI.
- [1992b] Glenn Carroll and Eugene Charniak. 1992b. Two experiments on learning probabilistic dependency grammars from corpora. In *Workshop Notes, Statistically-Based NLP Techniques*, pages 1–13. AAAI.
- [1994] Eugene Charniak and Glenn Carroll. 1994. Context-sensitive statistics for improved grammatical language models. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, pages 728–733.
- [1990] Mahesh V. Chitrao and Ralph Grishman. 1990. Statistical parsing of messages. In *DARPA Speech and Language Workshop*, pages 263–266.
- [1982] W. Nelson Francis and Henry Kučera. 1982. *Frequency Analysis of English Usage: Lexicon and Grammar*. Houghton Mifflin.
- [1991] Frederick Jelinek and John D. Lafferty. 1991. Computation of the probability of initial substring generation by stochastic context-free grammars. *Computational Linguistics*, 17:315–323.
- [1991] Fred Kochman and Joseph Kupin. 1991. Calculating the probability of a partial parse of a sentence. In *DARPA Speech and Language Workshop*, pages 237–240.
- [1991] David M. Magerman and Mitchell P. Marcus. 1991. Parsing the voyager domain using pearl. In *DARPA Speech and Language Workshop*, pages 231–236.
- [1992] David M. Magerman and Carl Weir. 1992. Efficiency, robustness and accuracy in picky chart parsing. In *Proceedings of the 30th ACL Conference*, pages 40–47.
- [1994] Scott Miller and Heidi Fox. 1994. Automatic grammar acquisition. In *Proceedings of the Human Language Technology Workshop*, pages 268–271.