

New Figures of Merit for Best-First Probabilistic Chart Parsing

Sharon A. Caraballo*
Brown University

Eugene Charniak*
Brown University

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 little 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, and we identify an easily computable figure of merit that provides excellent performance on various measures and two different grammars.

1. Introduction

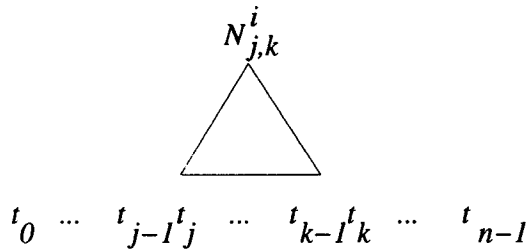
Chart parsing is a commonly used algorithm for parsing natural language texts. The chart is a data structure that contains all of the constituents for which subtrees have been found, that is, constituents for which a derivation has been found and which may therefore appear in some complete parse of the sentence. The agenda is a structure that stores a list of constituents for which a derivation has been found but which have not yet been combined with other constituents. Initially, the agenda contains the terminal symbols from the sentence to be parsed. A constituent is removed from the agenda and added to the chart, and the system considers how this constituent can be used to extend its current structural hypothesis by combining with other constituents in the chart according to the grammar rules. (We will often refer to these expansions of rules as “edges”.) In general this can lead to the creation of new, more encompassing constituents, which themselves are then added to the agenda. When one constituent has been processed, a new one is chosen to be removed from the agenda, and so on. Traditionally, the agenda is represented as a stack, so that the last item added to the agenda is the next one removed. Chart parsing is described extensively in the literature; for one such discussion see Section 1.4 of Charniak (1993).

Best-first probabilistic chart parsing is a variation of chart parsing that 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 probabilistic figure of merit is assigned to the constituents on the agenda, and the constituent maximizing this value is the next to be added to the chart.

In this paper we consider probabilities primarily based on probabilistic context-free grammars, though in principle, other, more complicated schemes could be used.

The purpose of this work is to compare how well several figures of merit select

* Computer Science Department, Box 1910, Brown University, Providence, RI 02912. E-mail: {sc, ec}@cs.brown.edu

**Figure 1**

Constituent $N_{j,k}^i$ in a sentence $t_{0,n}$.

constituents to be moved from the agenda to the chart. 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 is most likely given 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 $t_j \dots t_{k-1}$. (See Figure 1.)

In our experiments, we use only tag sequences (as given in the test data) for parsing. More accurate probability estimates should be attainable using lexical information in future experiments, as more detail usually leads to better statistics, but lexicalized figures of merit are beyond the scope of the research described here.

Note that our “ideal” figure is simply a heuristic, since there is no guarantee that a constituent that scores well on this measure will appear in the correct parse of a sentence. For example, there may be a very large number of low-probability derivations of $N_{j,k}^i$, which are combined here to give a high value, but a parse of the sentence can only include one of these derivations, making it unlikely that $N_{j,k}^i$ appears in the most probable parse of the sentence. On the other hand, there is no reason to believe that such cases are common in practice.

We cannot calculate $p(N_{j,k}^i | t_{0,n})$, 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, using two different grammars. We identify a figure of merit that gives superior results on all of our performance measures and on both grammars.

Section 2 of this paper describes the method we used to determine the effectiveness of figures of merit, that is, to compare how well they choose constituents to be moved from the agenda to the chart. Section 2.1 explains the experiment, Section 2.2 describes the measures we used to compare the performance of the figures of merit, and Section 2.3 describes a model we used to represent the performance of a traditional parser using a simple stack as an agenda.

In Section 3, we describe and compare three simple and easily computable figures of merit based on inside probability. Sections 3.1 through 3.3 describe each figure in detail, and Section 3.4 presents the results of an experiment comparing these three figures. Sections 4 and 5 have a similar structure to Section 3, with Section 4 evaluating two figures of merit using statistics on the left-side context of the constituent, and

Section 5 evaluating three additional figures of merit using statistics on the context on both sides of the constituent. Section 6 contains a table summarizing the results from Sections 3, 4, and 5.

In Section 7, we use another grammar in the experiment, to verify that our results are not an artifact of the grammar used for parsing. Section 8 describes previous work in this area, and Section 9 presents our conclusions and recommendations.

There are also three appendices to this paper. Appendix A gives our method for computing inside probability estimates while maintaining parser speed. Appendix B explains how we obtained our boundary statistics used in Section 5. Appendix C presents data comparing the parsing accuracy obtained by each of our parsers as the number of edges they create increases.

2. Comparing Figures of Merit

2.1 The Experiment

We used as our first grammar a probabilistic context-free grammar learned from the Brown corpus (see Francis and Kučera [1982] for a description of the Brown Corpus, and Carroll and Charniak [1992a, 1992b], and Charniak and Carroll [1994] for grammar and training details). This grammar contains about 5,000 rules using 32 terminal and nonterminal symbols. We parsed 500 sentences of length 3 to 30 (including punctuation) from the Penn Treebank *Wall Street Journal* corpus (Marcus, Santorini, and Marcinkiewicz 1993) using a best-first parsing method and various estimates for $p(N_{j,k}^i | t_{0,n})$ as the figure of merit.

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. The 95% figure is simply a convenience; see Appendix C for a discussion of speed versus accuracy.

2.2 Measures Used

We compared the figures of merit using the following measures:

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 a good measure of CFG parser effort, since it is independent of platform and implementation.
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 bookkeeping measure, and, as such, virtually cannot be eliminated. 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. This measure was included to confirm that a figure of merit that is efficient in terms of edge creation is also efficient in terms of constituent creation.

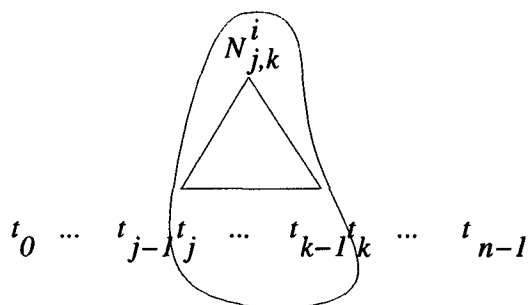


Figure 2

β includes only words within the constituent.

4. CPU time: The total CPU time (in seconds) needed to get 95% of the probability mass for all of the 500 sentences.

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.

2.3 The "Stack" Model

As a basis for comparison, we measured the CPU time for a non-best-first version of the parser to completely parse all 500 sentences. The CPU time needed by this version of the parser was 4,882 seconds. For a best-first version of the parser to be useful, it must be able to find the most probable parse (or a reasonably good parse, depending on the application) in less than this amount of time. Here, for the best-first parsers, we will use for convenience the time needed to get 95% of the sentence's total probability mass.

3. Simple Figures of Merit

3.1 Straight β

It seems reasonable to base a figure of merit on the inside probability β 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; see Figure 2. 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. Appendix A gives details of our on-line computation of β .

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.

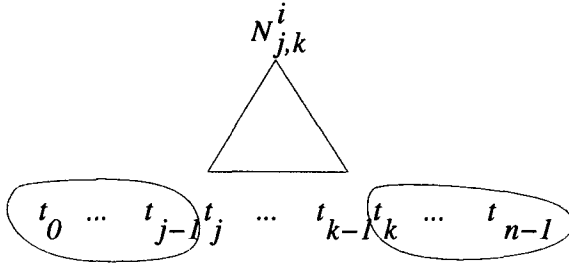


Figure 3
 α includes the entire context of the constituent.

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

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

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 α of the constituent. Outside probability α 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); see Figure 3.

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

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 β alone assumes that α and $p(t_{0,n})$ can be ignored.

We will refer to this figure of merit as **straight β** .

3.2 Normalized β

One side effect of omitting the α and $p(t_{0,n})$ terms in the straight β 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 as noted in Chitrao and Grishman (1990), where the system parses short constituents, even very low-probability ones, while avoiding combining them into longer constituents. To avoid thrashing, 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:

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

We will refer to this figure as **normalized β** .

3.3 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 $p(t_{j,k} | t_{0,j}, t_{k,n}) \approx p(t_{j,k} | t_{j-2}, t_{j-1})$. Then we have:

$$p(N_{j,k}^i | t_{0,n}) \approx \frac{p(N^i)\beta(N_{j,k}^i)}{p(t_{j,k} | t_{j-2}, t_{j-1})}.$$

We can calculate $\beta(N_{j,k}^i)$ as usual.

The $p(N^i)$ term is estimated from our PCFG and the training data from which the grammar was learned. We estimate $p(N^i)$ 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.¹

The $p(t_{j,k} | t_{j-2}, t_{j-1})$ term is just the probability of the tag sequence $t_j \dots t_{k-1}$ according to a trigram model. (Technically, this is not a trigram model but a tritag model, since we are considering sequences of tags, not words.) Our tritag probabilities $p(t_a | t_{a-2}, t_{a-1})$ were learned from the training data used for the grammar, using

¹ Our results show that the $p(N^i)$ term can be omitted from this figure of merit without much effect.

Table 1
Results for the β estimates.

Figure of Merit	%E	%non-0 E	%popped	CPU Time
straight β	97.6	97.5	93.8	3,966
normalized β	34.7	31.6	61.5	1,631
trigram estimate	25.2	21.7	44.3	1,547
"stack"	—	—	—	4,882

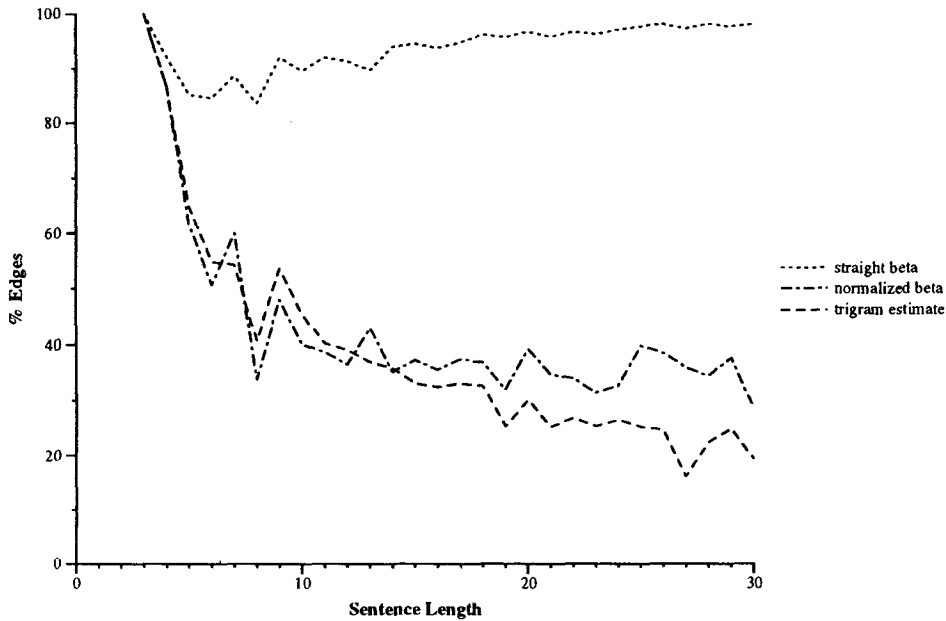


Figure 4
Nonzero-length edges for 95% of the probability mass for the β estimates.

the deleted interpolation method for smoothing. Our figure of merit uses:

$$p(t_{j,k} | t_{j-2}, t_{j-1}) \approx \prod_{a=j}^{k-1} p(t_a | t_{a-2}, t_{a-1})$$

We refer to this figure of merit as the **trigram estimate**.

3.4 Results

The results for the three figures of merit introduced in the last section according to the measurements given in Section 2.2 are shown in Table 1 (the time to fully parse using the "stack" model is included for easy reference).

Figure 4 expands the %non-0 E data to show the percent of nonzero-length edges needed to get 95% of the probability mass for each sentence length.

Straight β performs quite poorly on this measure. In order to find 95% of the probability mass for a sentence, a parser using this figure of merit typically needs to do over 90% of the work. On the other hand, normalized β and the trigram estimate both result in substantial savings of work. However, while these two models produce

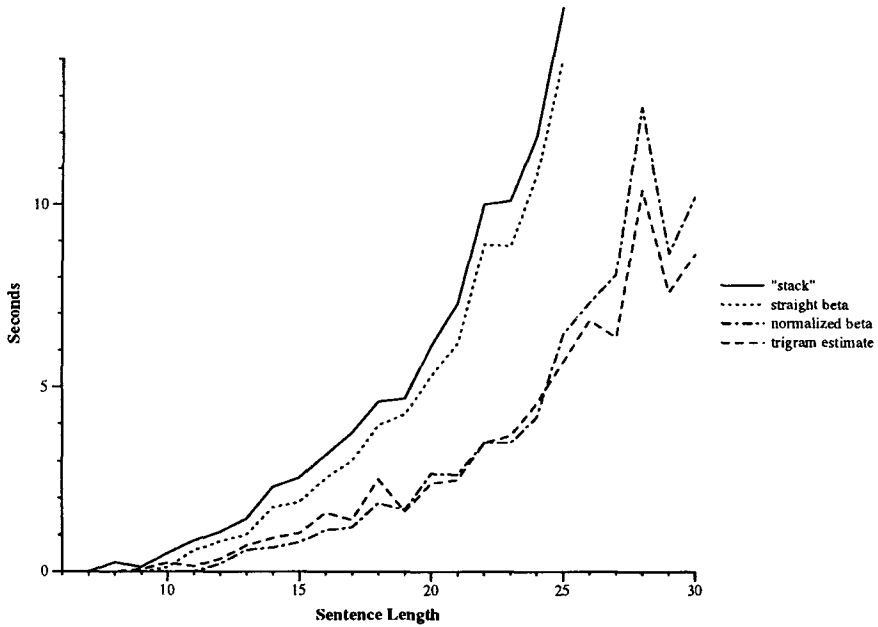


Figure 5
Average CPU time for 95% of the probability mass for the β estimates.

near-equivalent performance for short sentences, for longer sentences, with length greater than about 15 words, the trigram estimate gains a clear advantage. In fact, the performance of normalized β appears to level off in this range, while the amount of work done using the trigram estimate shows a continuing downward trend.

Figure 5 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 β and trigram estimates can be better compared).

Note that while straight β does perform better than the "stack" model in CPU time, the two models approach equivalent performance as sentence length increases, which is what would be expected from the edge count measures. The other two models provide a real time savings over the "stack" model, as can be seen from Figure 5 and from the total CPU times given earlier. Through most of the length range, the CPU time needed by the normalized β and the trigram estimate is quite close, but at the upper end of the range we can see better performance by the trigram estimate. (This improvement comes later than in the edge count statistics because of the small additional amount of overhead work needed to use the trigram estimate.)

4. Figures Involving Left Outside Probability

4.1 Normalized $\alpha_L\beta$

Earlier, 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 α term, representing outside probability, cannot be calculated di-

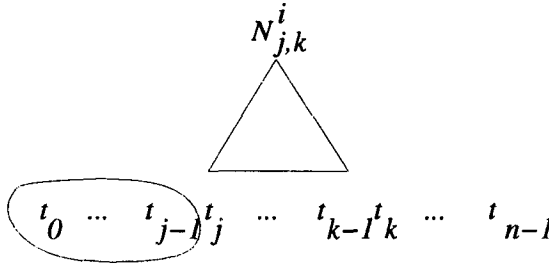


Figure 6
Left outside context.

rectly 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 α_L (see Figure 6).

The following recursive formula can be used to compute α_L . Let $\mathcal{E}_{j,k}^i$ be the set of all 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 α_L of the nonterminal appearing on the left-hand side (lhs) of the rule, the probability of the rule itself, and β 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:

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

Given a complete parse of the sentence, the formula above gives an exact value for α_L . During parsing, the set $\mathcal{E}_{j,k}^i$ is not complete, and so the formula gives an approximation of α_L .

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

A simple extension to the normalized β 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 $\alpha_L \beta$ quantity's preference for shorter constituents, as explained in the previous section.

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

4.2 Prefix Estimate

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

Table 2
Results for the $\alpha_L\beta$ estimates.

Figure of Merit	%E	%non-0 E	%popped	CPU Time
normalized $\alpha_L\beta$	39.7	36.4	57.3	68,660
prefix estimate	21.8	17.4	38.3	26,520

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 α_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**.

4.3 Results

The results for the figures of merit introduced in the previous section according to the measurements given in Section 2.2 are shown in Table 2.

Figure 7 shows a graph of %non-0 E for each sentence length for the two α_L models and the related β models.

Figure 7 illustrates two main points. First, the deterioration of the performance of the geometric-mean-based models with sentence length can be seen clearly. Second, when we consider only the two conditional-probability models, we can see that the additional information obtained from context in the prefix estimate gives a substantial improvement in this measure as compared to the trigram estimate.

However, the CPU time needed to compute the α_L term exceeds the time saved by processing fewer edges. Note that using this estimate, the parser took over 26,000 seconds to get 95% of the probability mass, while the "stack" model can exhaustively parse the test data in less than 5,000 seconds. Figure 8 shows the average CPU time for each sentence length.

While chart parsing and calculations of β can be done in $O(n^3)$ time (see Appendix A), we have been unable to find an algorithm to compute the α_L terms faster

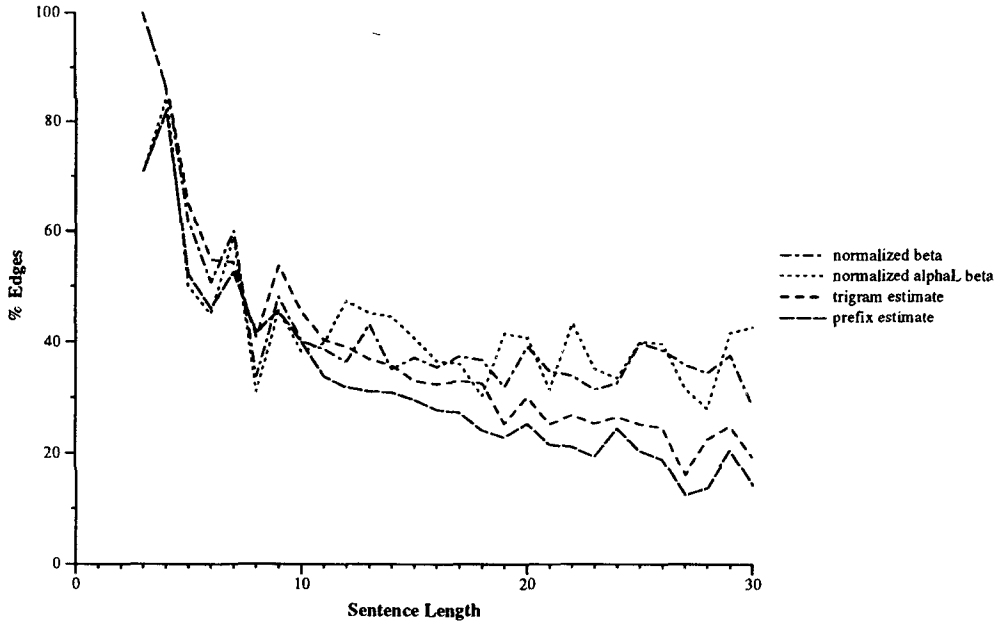


Figure 7
Nonzero-length edges for 95% of the probability mass for the $\alpha_L\beta$ estimates.

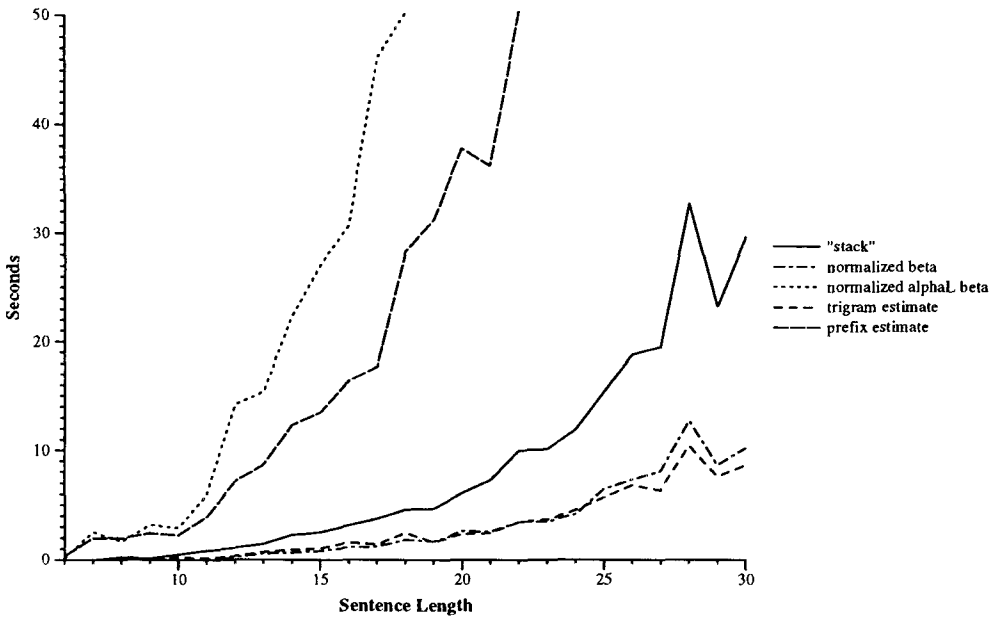


Figure 8
Average CPU time for 95% of the probability mass for the $\alpha_L\beta$ estimates.

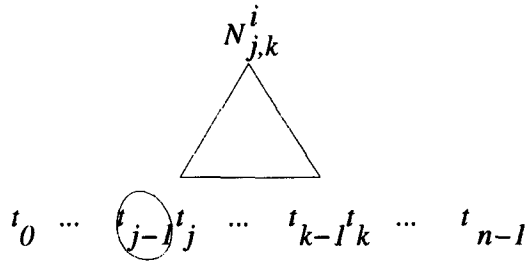


Figure 9
Left boundary context.

than $O(n^5)$. When a constituent is removed from the agenda, it only affects the β values of its ancestors in the parse trees; however, α_L values are propagated to all of the constituent’s siblings to the right and all of its descendants. Recomputing the α_L terms when a constituent is removed from the agenda can be done in $O(n^3)$ time, and since there are $O(n^2)$ possible constituents, the total time needed to compute the α_L terms in this manner is $O(n^5)$.

5. Figures Using Boundary Statistics

5.1 Left Boundary Trigram Estimate

Although the α_L -based models seem impractical, the edge-count and constituent-count statistics show that contextual information is useful. We can derive an estimate similar to the prefix estimate but containing a much simpler model of the context 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:

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

As usual, we use a trigram model for the tags, giving $p(t_{j,k} | t_{0,j}, t_{k,n}) \approx p(t_{j,k} | t_{j-2}, t_{j-1})$.

Now, we assume that $p(N_{j,k}^i | t_{0,j}, t_{k,n}) \approx p(N_{j,k}^i | t_{j-1})$, that is, that the probability of a nonterminal is dependent on the tag immediately before it in the sentence (see Figure 9). Then we have:

$$p(N_{j,k}^i | t_{0,n}) \approx \frac{p(N_{j,k}^i | t_{j-1})\beta(N_{j,k}^i)}{p(t_{j,k} | t_{j-2}, t_{j-1})}.$$

We can calculate $\beta(N_{j,k}^i)$ and the tritag probabilities as usual. The $p(N_{j,k}^i | t_{j-1})$ probabilities are estimated from our training data by parsing the training data and

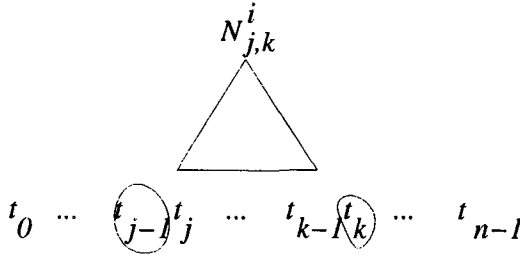


Figure 10
Boundary context.

counting the occurrences of the nonterminal and the tag weighted by their probability in the parse. (Further details are provided in Appendix B.)

We will refer to this figure as the **left boundary trigram estimate**.

5.2 Boundary Trigram Estimate

We can derive a similar estimate using context on both sides of the constituent as follows:

$$\begin{aligned}
 & p(N_{j,k}^i | t_{0,n}) \\
 &= \frac{p(N_{j,k}^i, t_{0,n})}{p(t_{0,n})} \\
 &= \frac{p(t_{0,j})p(N_{j,k}^i | t_{0,j})p(t_{j,k} | N_{j,k}^i, t_{0,j})p(t_k | t_{0,j}, N_{j,k}^i, t_{j,k})p(t_{k+1,n} | t_{0,j}, N_{j,k}^i, t_{j,k}, t_k)}{p(t_{0,j})p(t_{j,k} | t_{0,j})p(t_k | t_{0,k})p(t_{k+1,n} | t_{0,k+1})} \\
 &= \frac{p(N_{j,k}^i | t_{0,j})p(t_{j,k} | N_{j,k}^i, t_{0,j})p(t_k | t_{0,j}, N_{j,k}^i, t_{j,k})p(t_{k+1,n} | t_{0,k+1}, N_{j,k}^i)}{p(t_{j,k} | t_{0,j})p(t_k | t_{0,k})p(t_{k+1,n} | t_{0,k+1})}
 \end{aligned}$$

Once again applying the usual independence assumption that given a nonterminal, the tag sequence it generates depends only on that nonterminal and also assuming that the probability of $t_{k+1,n}$ depends only on the previous tags, we can rewrite the figure of merit as follows:

$$p(N_{j,k}^i | t_{0,n}) \approx \frac{p(N_{j,k}^i | t_{0,j})\beta(N_{j,k}^i)p(t_k | t_{0,k}, N_{j,k}^i)}{p(t_{j,k+1} | t_{0,j})}$$

Now we add some new independence assumptions. We assume that the probability of the nonterminal depends only on the immediately preceding tag, and that the probability of the tag immediately following the nonterminal depends only on the nonterminal (see Figure 10), giving:

$$p(N_{j,k}^i | t_{0,n}) \approx \frac{p(N_{j,k}^i | t_{j-1})\beta(N_{j,k}^i)p(t_k | N_{j,k}^i)}{p(t_{j,k+1} | t_{0,j})}$$

As usual, we use a trigram model for the tags, giving $p(t_{j,k} | t_{0,j}, t_{k,n}) \approx p(t_{j,k} | t_{j-2}, t_{j-1})$. Then we have:

$$p(N_{j,k}^i | t_{0,n}) \approx \frac{p(N_{j,k}^i | t_{j-1})\beta(N_{j,k}^i)p(t_k | N_{j,k}^i)}{p(t_{j,k+1} | t_{j-2}, t_{j-1})}$$

We can calculate $\beta(N_{j,k}^i)$ and the tritrag probabilities as usual. The $p(N_{j,k}^i | t_{j-1})$ and $p(t_k | N_{j,k}^i)$ probabilities are estimated from our training data by parsing the training data and counting the occurrences of the nonterminal and the tag weighted by their probability in the parse.² Again, see Appendix B for details of how these estimates were obtained.

We will refer to this figure as the **boundary trigram estimate**.

5.3 Boundary Statistics Only

We also wished to examine whether contextual information by itself is sufficient as a figure of merit. We can derive an estimate based only on easily computable contextual information as follows:

$$\begin{aligned} p(N_{j,k}^i | t_{0,n}) &= \frac{p(N_{j,k}^i, t_{0,n})}{p(t_{0,n})} \\ &= \frac{p(t_{0,j})p(N_{j,k}^i | t_{0,j})p(t_{j,k} | N_{j,k}^i, t_{0,j})p(t_k | t_{0,j}, N_{j,k}^i, t_{j,k})p(t_{k+1,n} | t_{0,j}, N_{j,k}^i, t_{j,k}, t_k)}{p(t_{0,j})p(t_{j,k} | t_{0,j})p(t_k | t_{0,k})p(t_{k+1,n} | t_{0,k+1})} \\ &= \frac{p(N_{j,k}^i | t_{0,j})p(t_{j,k} | N_{j,k}^i, t_{0,j})p(t_k | t_{0,j}, N_{j,k}^i, t_{j,k})p(t_{k+1,n} | t_{0,k+1}, N_{j,k}^i)}{p(t_{j,k} | t_{0,j})p(t_k | t_{0,k})p(t_{k+1,n} | t_{0,k+1})}. \end{aligned}$$

Most of the independence assumptions we make are the same as in the boundary trigram estimate. We assume that the probability of the nonterminal depends only on the previous tag, that the probability of the immediately following tag depends only on the nonterminal, and that the probability of the tags following that depend only on the previous tags. However, we make one independence assumption that differs from all of our previous estimates. Rather than assuming that the probability of the tags within the constituent depends on the nonterminal, giving an inside probability term, we assume that the probability of these tags depends only on the previous tags. Then we have

$$\begin{aligned} p(N_{j,k}^i | t_{0,n}) &\approx \frac{p(N_{j,k}^i | t_{0,j})p(t_{j,k} | t_{0,j})p(t_k | N_{j,k}^i)p(t_{k+1,n} | t_{0,k+1})}{p(t_{j,k} | t_{0,j})p(t_k | t_{0,k})p(t_{k+1,n} | t_{0,k+1})} \\ &= \frac{p(N_{j,k}^i | t_{0,j})p(t_k | N_{j,k}^i)}{p(t_k | t_{0,k})}. \end{aligned}$$

In the denominator, we take $p(t_k | t_{0,k}) \approx p(t_k)$, giving:

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

which is simply the product of the two boundary statistics described in the previous section.

We refer to this estimate as **boundary statistics only**.

² Actually, in our implementation, the $p(t_k)$ in the denominator is included in the following-tag statistic, for which we use $\frac{p(t_k | N_{j,k}^i)}{p(t_k)}$. Then at run time we only use the trigram probabilities for $t_{0,k}$.

Table 3
Results for the boundary estimates.

Figure of Merit	%E	%non-0 E	%popped	CPU Time
boundary statistics only	53.2	50.8	59.6	2,759
left boundary trigram estimate	22.1	18.4	39.6	1,700
boundary trigram estimate	18.2	13.9	31.2	1,111

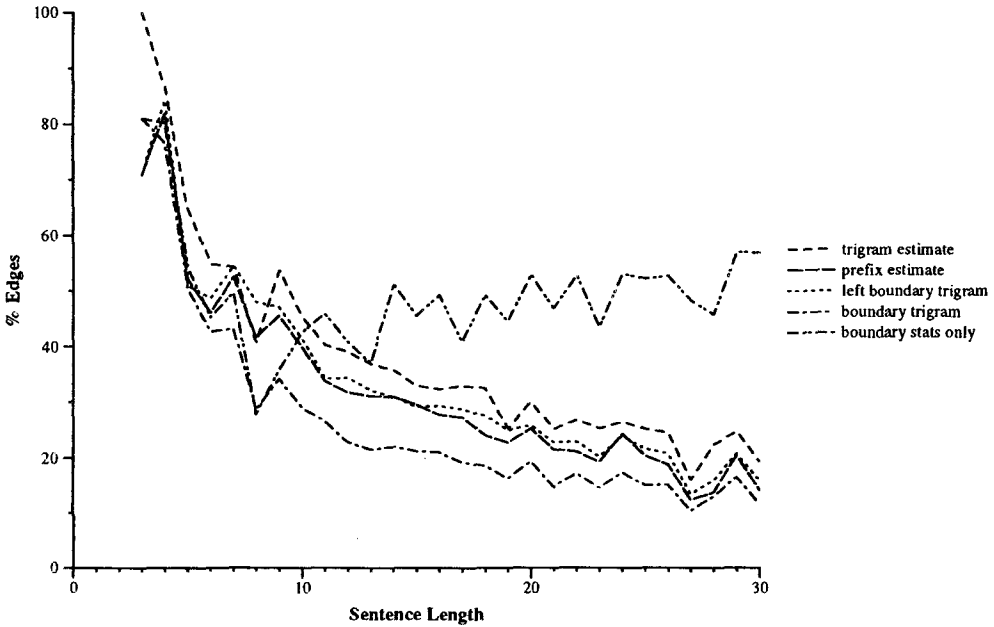


Figure 11
Nonzero-length edges for 95% of the probability mass for the boundary estimates.

5.4 Results

The results for the figures of merit introduced in the previous section according to the measurements given in Section 2.2 are shown in Table 3.

Figure 11 shows a graph of %non-0 E for each sentence length for the boundary models and the trigram and prefix estimates. This graph shows that the contextual information gained from using α_L in the prefix estimate is almost completely included in just the previous tag, as illustrated by the left boundary trigram estimate. Adding right contextual information in the boundary trigram estimate gives us the best performance on this measure of any of our figures of merit.

We can consider the left boundary trigram estimate to be an approximation of the prefix estimate, where the effect of the left context is approximated by the effect of the single tag to the left. Similarly, the boundary trigram estimate is an approximation to an estimate involving the full context, i.e., an estimate involving the outside probability α . However, the parser cannot compute the outside probability of a constituent during a parse, and so in order to use context on both sides of the constituent, we need to use something like our boundary statistics. Our results suggest that a single tag before or after the constituent can be used as a reasonable approximation to the full context on

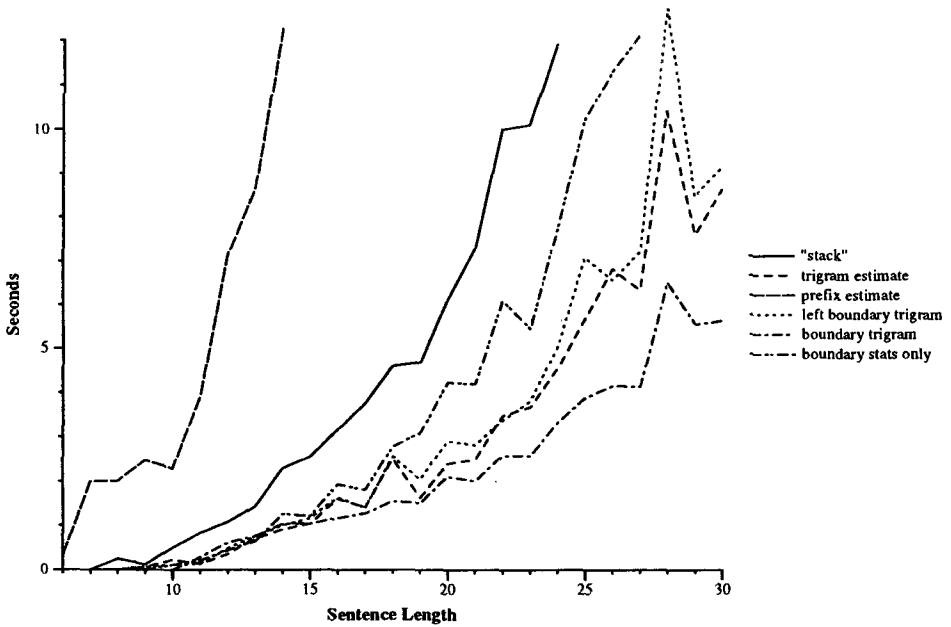


Figure 12
Average CPU time for 95% of the probability mass for the boundary estimates.

that side of the constituent. Figure 12 shows the average CPU time for each sentence length.

Since the boundary trigram estimate has none of the overhead associated with the prefix estimate, it is the best performer in terms of CPU time as well. We can also see that using just the boundary statistics, which can be precomputed and require no extra processing during parsing, still results in a substantial improvement over the non-best-first “stack” model.

As another method of comparison between the two best-performing estimates, the context-dependent boundary trigram model and the context-independent trigram model, we compared the number of edges needed to find the first parse for average-length sentences. The average length of a sentence in our test data is about 22 words. Figure 13 shows the percentage of sentences of length 18 through 26 for which a parse could be found within 2,500 edges. For this experiment, we used a separate test set from the *Wall Street Journal* corpus, containing approximately 570 sentences in the desired length range. This measure also shows a real advantage of the boundary trigram estimate over the trigram estimate.

6. Results Summary

Table 4 summarizes the results obtained for each figure of merit.

7. Comparing Figures of Merit Using a Treebank Grammar

7.1 Background

To verify that our results are not an artifact of the particular grammar we chose for testing, we also tested using a treebank grammar introduced in Charniak (1996). This

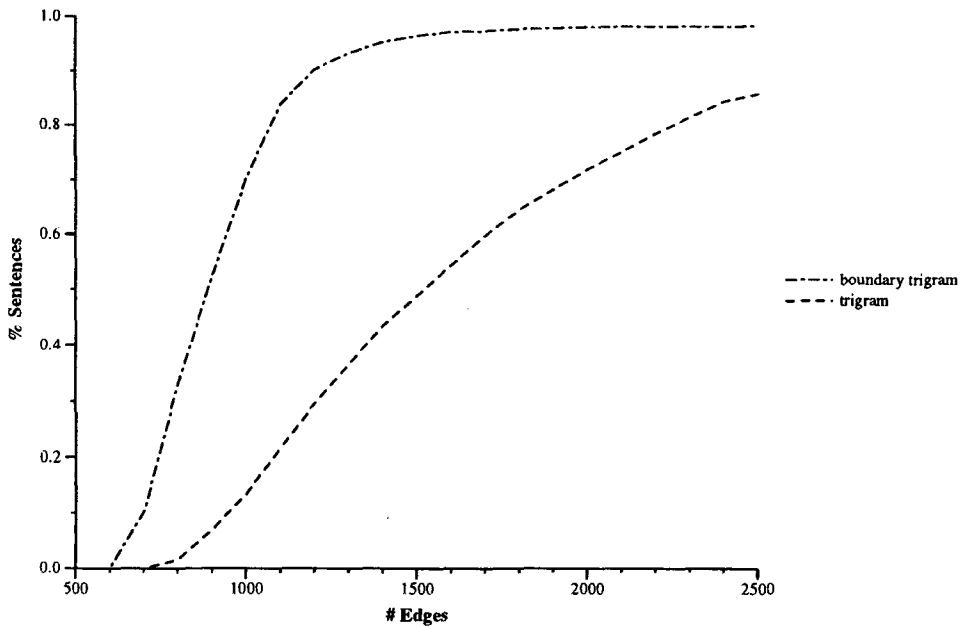


Figure 13
% of the 18- to 26-word sentences finding a parse in a fixed number of edges.

Table 4
Results for all figures of merit.

Figure of Merit	%E	%non-0 E	%popped	CPU Time
"stack" model				4,882
straight β	97.6	97.5	93.8	3,966
normalized β	34.7	31.6	61.5	1,631
trigram estimate	25.2	21.7	44.3	1,547
normalized $\alpha_L\beta$	39.7	36.4	57.3	68,660
prefix estimate	21.8	17.4	38.3	26,520
boundary statistics only	53.2	50.8	59.6	2,759
left boundary trigram estimate	22.1	18.4	39.6	1,700
boundary trigram estimate	18.2	13.9	31.2	1,111

grammar was trained in a straightforward way by reading the grammar directly (with minor modifications) from a portion of the Penn Treebank *Wall Street Journal* data comprised of about 300,000 words. The boundary statistics were counted directly from the training data as well. The treebank grammar is much larger and more ambiguous than our original grammar, containing about 16,000 rules and 78 terminal and nonterminal symbols, and it was impractical to parse sentences to exhaustion using our existing hardware, so the figures based on 95% of the probability mass could not be computed. We were able to use this grammar to compare the number of edges needed to find the first parse using the trigram and boundary trigram estimates.

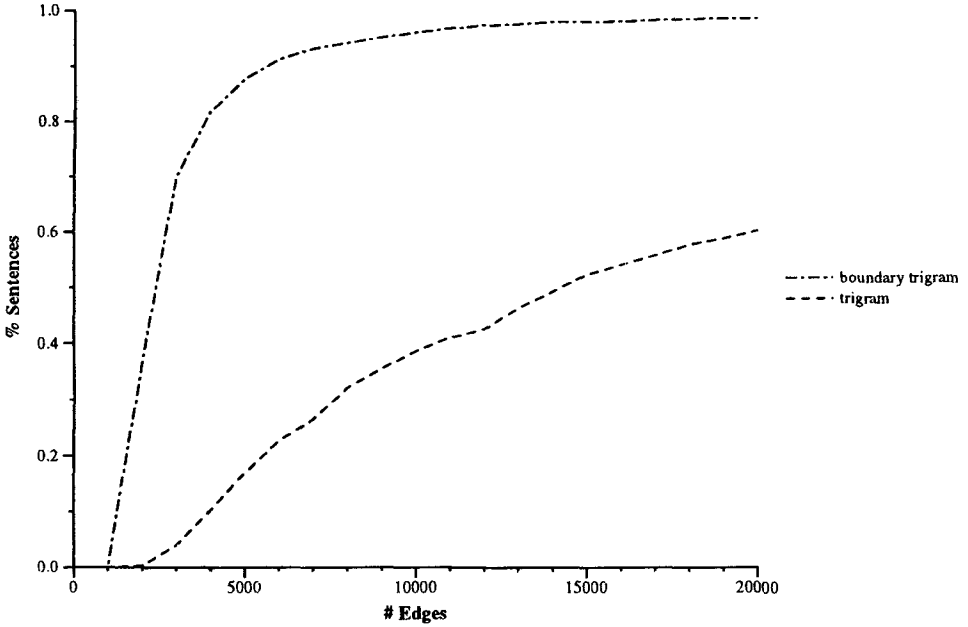


Figure 14
 % of the 18- to 26-word sentences finding a parse in a fixed number of edges for a treebank grammar.

7.2 Results

Figure 14 shows the percentage of sentences of length 18 through 26 for which a parse could be found within 20,000 edges. Again, we used a test set of approximately 570 sentences of the appropriate length from the *Wall Street Journal* corpus. Although the x-axis covers a much wider range than in Figure 13, the relationship between the two estimates is quite similar.

8. Previous Work

In an earlier version of this paper (Caraballo and Charniak 1996), we presented the results for several of these models using our original grammar. The treebank grammar was introduced in Charniak (1996), and the parser in that paper is a best-first parser using the boundary trigram figure of merit.

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.

9. Conclusions

We have presented and evaluated several figures of merit for best-first parsing. The best performer according to all of our measures was the parser using the boundary trigram estimate as a figure of merit, and this result holds for two different grammars. 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. (As mentioned earlier, the efficient online computation of β is described in Appendix A.) We strongly recommend this figure of merit as the basis for best-first statistical parsers.

The measurements presented here almost certainly underestimate the true benefits of this model. 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 improvement would be even more dramatic for sentences longer than 30 words.

Appendix A: Efficient On-Line Computation of β

We compute estimates of the inside probability β for each proposed constituent incrementally as new constituents are added to the chart. Initially, β is set to 1 for each terminal symbol, since our input is given as a stream of tags, which are our terminals. When a new proposed constituent is added to the agenda, its β estimate is set to its current inside probability according to the constituents already in the chart. However, as more constituents are added to the chart, we may find a new way to build up a proposed constituent, i.e., additional evidence for that proposed constituent, so we need to update the β for the proposed constituent (and also for affected constituents already in the chart, since these may in turn affect other proposed constituents).

These updates can be quite expensive in terms of CPU time. However, many of the updates are quite small, and do not affect the relative ordering of the proposed constituents on the agenda. Instead of propagating every change to β , then, we only want to propagate those changes that we expect to have an effect on this ordering. What we have done is to have each constituent store not only its β value, but also an increment. Increases to the inside probability are added not to β itself, but to this increment, until the increment exceeds some threshold. Experimentally we have found that we can avoid propagating increments until they exceed 1% of the current value of β with very little effect on the parser's selection of constituents from the agenda.

This thresholding on the propagation of β allows us to update the β values on line while still keeping the performance of the parser as $O(n^3)$ empirically.

Appendix B: Estimation of Boundary Statistics

Our figures of merit incorporating boundary statistics use the figures $p(N_{j,k}^i | t_{j-1})$ to represent the effect of the left context and $\frac{p(t_k | N_{j,k}^i)}{p(t_k)}$ to represent the effect of the right context. For our experiments with the first grammar, which was learned from training data taken from the Brown corpus, we estimated these statistics from the same training data.

First, we parsed the training data according to our grammar. (It was necessary to do this, rather than using the hand-annotated parses of the training data, because our grammar does not use the same set of nonterminals as the corpus; see Carroll and Charniak [1992a, 1992b] and Charniak and Carroll [1994] for details.) Since we use the tags as our input, the probability of a nonterminal appearing with a particular previous tag is the same as the probability of that nonterminal appearing in any sentence containing that tag.

We can then count the probability-weighted occurrences of a nonterminal given the previous tag as follows:

$$\begin{aligned} C(N_{j,k}^i, t_{j-1}) &= \sum_{w_{0,n} \text{ containing } t_{j-1}} p(N_{j,k}^i | w_{0,n}) \\ &= \frac{\alpha(N_{j,k}^i) \beta(N_{j,k}^i)}{p(w_{0,n})} \end{aligned}$$

That is, for each sentence that contains the previous tag t_{j-1} , we increment our count by the probability of the nonterminal $N_{j,k}^i$ occurring immediately following t_{j-1} in that sentence.

Since we have a complete parse, the inside and outside probabilities and the sentence probability can be easily computed. We can also obtain the count $C(t_{j-1})$ simply by counting the number of sentences in which that tag appears in position $j - 1$. We then obtain the conditional probability for the left boundary statistic as follows:

$$p(N_{j,k}^i | t_{j-1}) = \frac{C(N_{j,k}^i, t_{j-1})}{C(t_{j-1})}$$

The right boundary statistic is computed in the corresponding way.

For the experiment using the treebank grammar, these statistics were obtained by counting directly from the *Wall Street Journal* treebank corpus, just as the grammar rules and trigram statistics were.

Appendix C: Speed vs. Accuracy

As an additional verification of our results, we gathered data on speed versus accuracy. For this experiment, we used the probabilistic context-free grammar learned from the Brown corpus and the average-length test sentences described in Section 5.4. For each figure of merit, we computed the average precision and recall of the best parse found as compared to the number of edges created. We computed unlabeled precision and recall only, since our grammar uses a different set of nonterminals from those used in the test data.

Precision is defined as the percentage of the constituents proposed by our parser that are actually correct according to the treebank. For each edge count, we measured the precision of the best parse of each sentence found within that number of edges. Figure 15 is a graph of the average precision for the β figures of merit from Section 3, plotted against edge counts.

The fluctuations at the low edge counts are due to the small amount of data at this level. At a low edge count, very few sentences have actually been parsed, and since these sentences tend to be short and simple, the parses are likely to be correct. The sentences that could not be parsed do not contribute to the measurement of precision. As more sentences are parsed, precision settles at about 47%, the highest precision attainable by our particular test grammar, and remains there as edge counts increase.

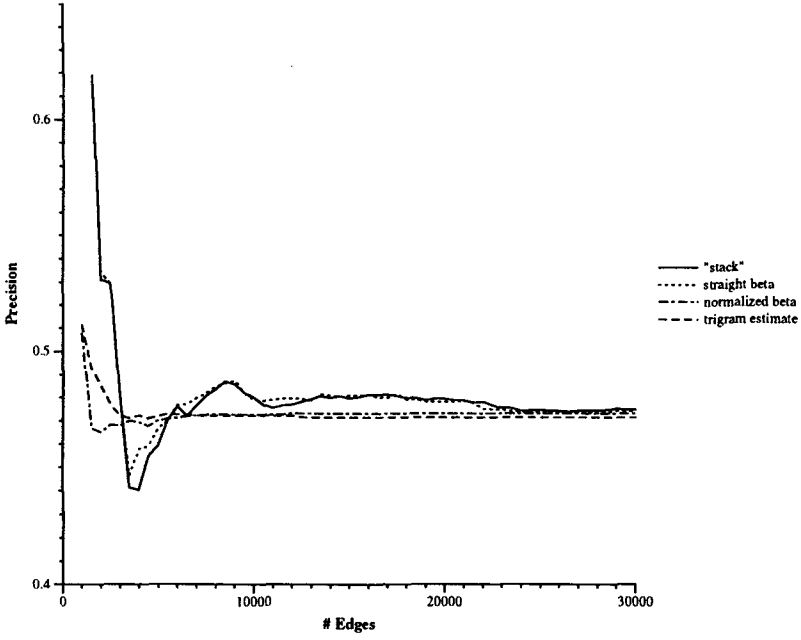


Figure 15
Precision of the best parse found in a fixed number of edges for the β estimates.

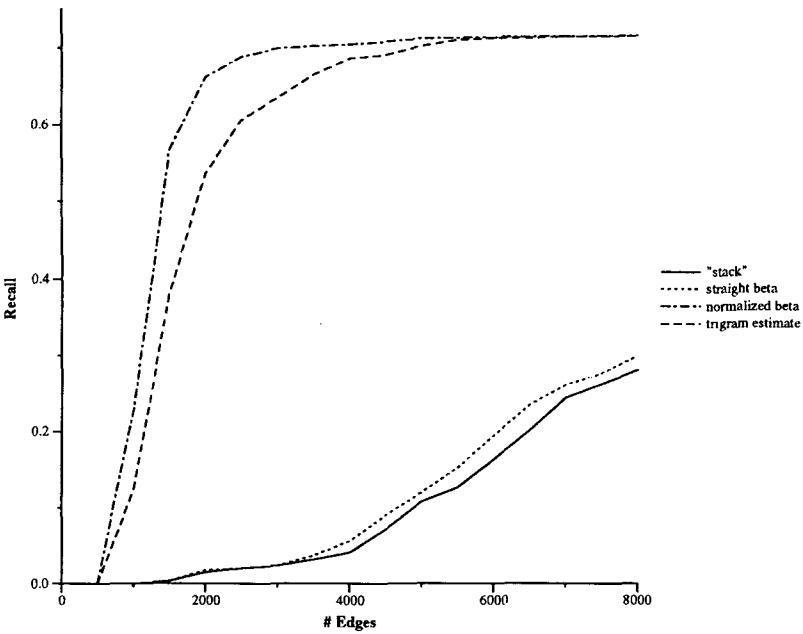


Figure 16
Recall of the best parse found in a fixed number of edges for the β estimates.

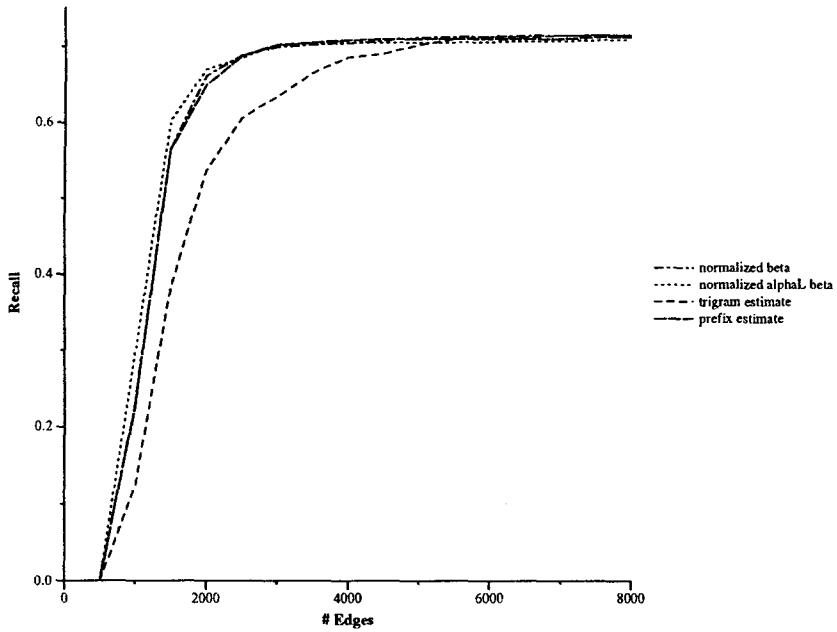


Figure 17
 Recall of the best parse found in a fixed number of edges for the $\alpha_L\beta$ estimates.

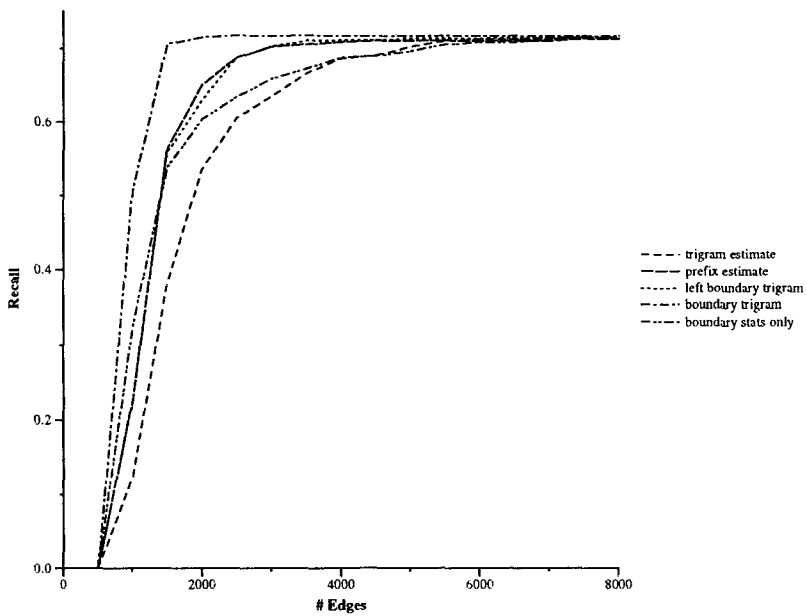


Figure 18
 Recall of the best parse found in a fixed number of edges for the boundary estimates.

This level of precision is independent of the figure of merit used, so measurement of precision does not help evaluate our figures of merit.

A much more useful measure is recall. Recall is defined as the percentage of constituents in the treebank test data that are found by our parser. Again, we measured the recall of the best parse of each sentence found within each number of edges. Figure 16 shows the results for the figures of merit from Section 3.

Straight beta clearly shows little or no improvement over the “stack” parser using no figure of merit at all. The other figures of merit increase quickly to about 64%, the maximum recall attainable with our test grammar. The “stack” parser and the one using straight beta, on the other hand, do not reach this maximum level until about 50,000 edges. We have no explanation for the relatively poor performance of the parser using the trigram estimate compared to the other best-first parsers, as shown in Figures 16, 17, and 18. Figure 17 shows the recall values for the $\alpha_L\beta$ figures of merit from Section 4, and Figure 18 shows recall for the boundary figures of merit from Section 5. Since precision is not a useful measure, we have not included precision data for these figures of merit.

These data confirm that the parser using the boundary trigram figure of merit performs better than any of the others. Recall using this figure of merit is consistently higher than any of the others at low edge counts, and it reaches the maximum value in fewer than 2,000 edges, with the nearest competitors approaching the maximum at about 3,000 edges.

Acknowledgments

The authors are very grateful to Heidi Fox for obtaining the speed vs. accuracy data discussed in Appendix C. We also wish to thank the anonymous reviewers for their comments and suggestions. This research was supported in part by NSF grant IRI-9319516 and by ONR grant N0014-96-1-0549.

References

- Bobrow, Robert J. 1990. Statistical agenda parsing. In *DARPA Speech and Language Workshop*, pages 222–224.
- Caraballo, Sharon and Eugene Charniak. 1996. Figures of merit for best-first probabilistic chart parsing. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 127–132.
- Carroll, Glenn and Eugene Charniak. 1992a. Learning probabilistic dependency grammars from labeled text. In *Working Notes, Fall Symposium Series*, pages 25–32. AAAI.
- Carroll, Glenn and Eugene Charniak. 1992b. Two experiments on learning probabilistic dependency grammars from corpora. In *Workshop Notes, Statistically-Based NLP Techniques*, pages 1–13. AAAI.
- Charniak, Eugene. 1993. *Statistical Language Learning*. MIT Press.
- Charniak, Eugene. 1996. Tree-bank grammars. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 1031–1036. AAAI.
- Charniak, Eugene 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.
- Chitrao, Mahesh V. and Ralph Grishman. 1990. Statistical parsing of messages. In *DARPA Speech and Language Workshop*, pages 263–266.
- Francis, W. Nelson and Henry Kučera. 1982. *Frequency Analysis of English Usage: Lexicon and Grammar*. Houghton Mifflin.
- Jelinek, Frederick and John D. Lafferty. 1991. Computation of the probability of initial substring generation by stochastic context-free grammars. *Computational Linguistics*, 17:315–323.
- Kochman, Fred and Joseph Kupin. 1991. Calculating the probability of a partial parse of sentence. In *DARPA Speech and Language Workshop*, pages 237–240.
- Magerman, David M. and Mitchell P. Marcus. 1991. Parsing the Voyager domain using Pearl. In *DARPA Speech and Language Workshop*, pages 231–236.
- Magerman, David M. and Carl Weir. 1992. Efficiency, robustness and accuracy in Picky chart parsing. In *Proceedings of the 30th Annual Meeting, Association for Computational Linguistics*, pages 40–47.

- Linguistics.
- Marcus, Mitchell P., Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn treebank. *Computational Linguistics*, 19:313–330.
- Miller, Scott and Heidi Fox. 1994. Automatic grammar acquisition. In *Proceedings of the Human Language Technology Workshop*, pages 268–271.