

Probabilistic Parsing and Psychological Plausibility

Thorsten Brants and Matthew Crocker
Saarland University, Computational Linguistics
D-66041 Saarbrücken, Germany
{brants,crocker}@coli.uni-sb.de

In *Proceedings of the 18th International Conference on Computational Linguistics*,
COLING-2000, July 29 – August 6, 2000, Saarbrücken / Luxembourg / Nancy

Abstract

Given the recent evidence for probabilistic mechanisms in models of human ambiguity resolution, this paper investigates the plausibility of exploiting current wide-coverage, probabilistic parsing techniques to model human linguistic performance. In particular, we investigate the performance of standard stochastic parsers when they are revised to operate incrementally, and with reduced memory resources. We present techniques for ranking and filtering analyses, together with experimental results. Our results confirm that stochastic parsers which adhere to these psychologically motivated constraints achieve good performance. Memory can be reduced down to 1% (compared to exhaustive search) without reducing recall and precision. Additionally, these models exhibit substantially faster performance. Finally, we argue that this general result is likely to hold for more sophisticated, and psycholinguistically plausible, probabilistic parsing models.

1 Introduction

Language engineering and computational psycholinguistics are often viewed as distinct research programmes: engineering solutions aim at practical methods which can achieve good performance, typically paying little attention to linguistic or cognitive modelling. Computational psycholinguistics, on the other hand, is often focussed on detailed modelling of human behaviour for a relatively small number of well-studied constructions. In this paper we suggest that, broadly, the human sentence processing mechanism (HSPM) and current statistical parsing technology can be viewed as having similar objectives: to optimally (i.e. rapidly and accurately) understand the text and utterances

they encounter.

Our aim is to show that large scale probabilistic parsers, when subjected to basic cognitive constraints, can still achieve high levels of parsing accuracy. If successful, this will contribute to a plausible explanation of the fact that people, in general, are also extremely accurate and robust. Such a result would also strengthen existing results showing that related probabilistic mechanisms can explain specific psycholinguistic phenomena.

To investigate this issue, we construct a standard 'baseline' stochastic parser, which mirrors the performance of a similar systems (e.g. (Johnson, 1998)). We then consider an incremental version of the parser, and evaluate the effects of several probabilistic filtering strategies which are used to prune the parser's search space, and thereby reduce memory load.

To assess the generality of our results for more sophisticated probabilistic models, we also conduct experiments using a model in which parent-node information is encoded on the daughters. This increase in contextual information has been shown to improve performance (Johnson, 1998), and the model is also shown to be robust to the incrementality and memory constraints investigated here.

We present the results of parsing performance experiments, showing the accuracy of these systems with respect to both a parsed corpus and the baseline parser. Our experiments suggest that a strictly incremental model, in which memory resources are substantially reduced through filtering, can achieve precision and recall which equals that of 'unrestricted' systems. Furthermore, implementation of these restrictions leads to substantially faster performance. In conclusion, we argue that such broad-coverage probabilistic parsing

models provide a valuable framework for explaining the human capacity to rapidly, accurately, and robustly understand “garden variety” language. This lends further support to psycholinguistic accounts which posit probabilistic ambiguity resolution mechanisms to explain “garden path” phenomena.

It is important to reiterate that our intention here is only to investigate the performance of probabilistic parsers under psycholinguistically motivated constraints. We do not argue for the psychological plausibility of SCFG parsers (or the parent-encoded variant) per se. Our investigation of these models was motivated rather by our desire to obtain a generalizable result for these simple and well-understood models, since obtaining similar results for more sophisticated models (e.g. (Collins, 1996; Ratnaparkhi, 1997)) might have been attributed to special properties of these models. Rather, the current result should be taken as support for the potential scalability and performance of probabilistic psychological models such as those proposed by (Jurafsky, 1996) and (Crocker and Brants, to appear).

2 Psycholinguistic Motivation

Theories of human sentence processing have largely been shaped by the study of pathologies in human language processing behaviour. Most psycholinguistic models seek to explain the *difficulty* people have in comprehending structures that are ambiguous or memory-intensive (see (Crocker, 1999) for a recent overview). While often insightful, this approach diverts attention from the fact that people are in fact extremely accurate and effective in understanding the vast majority of their “linguistic experience”. This observation, combined with the mounting psycholinguistic evidence for statistically-based mechanisms, leads us to investigate the merit of exploiting robust, broad coverage, probabilistic parsing systems as models of human linguistic performance.

The view that human language processing can be viewed as an optimally adapted system, within a probabilistic framework, is advanced by (Chater et al., 1998), while (Jurafsky, 1996) has proposed a specific probabilistic parsing model of human sentence processing. In work on human lexical category dis-

ambiguation, (Crocker and Corley, to appear), have demonstrated that a standard (incremental) HMM-based part-of-speech tagger models the finding from a range of psycholinguistic experiments. In related research, (Crocker and Brants, 1999) present evidence that an incremental stochastic parser based on Cascaded Markov Models (Brants, 1999) can account for a range of experimentally observed local ambiguity preferences. These include NP/S complement ambiguities, reduced relative clauses, noun-verb category ambiguities, and ‘that’-ambiguities (where ‘that’ can be either a complementizer or a determiner) (Crocker and Brants, to appear).

Crucially, however, there are differences between the classes of mechanisms which are psychologically plausible, and those which prevail in current language technology. We suggest that two of the most important differences concern *incrementality*, and *memory resources*. There is overwhelming experimental evidence that people construct connected (i.e. semantically interpretable) analyses for each initial substring of an utterance, as it is encountered. That is, processing takes place incrementally, from left to right, on a word by word basis.

Secondly, it is universally accepted that people can at most consider a relatively small number of competing analyses (indeed, some would argue that number is one, i.e. processing is strictly serial). In contrast, many existing stochastic parsers are “unrestricted”, in that they are optimised for accuracy, and ignore such psychologically motivated constraints. Thus the appropriateness of using broad-coverage probabilistic parsers to model the high level of human performance is contingent upon being able to maintain these levels of accuracy when the constraints of incrementality and resource limitations are imposed.

3 Incremental Stochastic Context-Free Parsing

The following assumes that the reader is familiar with stochastic context-free grammars (SCFG) and stochastic chart-parsing techniques. A good introduction can be found, e.g., in (Manning and Schütze, 1999). We use standard abbreviations for terminal nodes, non-terminal nodes, rules and probabilities.

This paper investigates stochastic context-free parsing based on a grammar that is derived from a treebank, starting with part-of-speech tags as terminals. The grammar is derived by collecting all rules $X \rightarrow \alpha$ that occur in the treebank and their frequencies f . The probability of a rule is set to

$$P(X \rightarrow \alpha) = \frac{f(X \rightarrow \alpha)}{\sum_{\beta} f(X \rightarrow \beta)} \quad (1)$$

For a description of treebank grammars see (Charniak, 1996). The grammar does not contain ϵ -rules, otherwise there is no restriction on the rules. In particular, we do not require Chomsky-Normal-Form.

In addition to the rules that correspond to structures in the corpus, we add a new start symbol ROOT to the grammar and rules $\text{ROOT} \rightarrow X$ for all non-terminals X together with probabilities derived from the root nodes in the corpus¹.

For parsing these grammars, we rely upon a standard bottom-up chart-parsing technique with a modification for incremental parsing, i.e., for each word, all edges are processed and possibly pruned before proceeding to the next word. The outline of the algorithm is as follows.

A chart entry E consists of a start and end position i and j , a dotted rule $X \rightarrow \alpha.\gamma$, the inside probability $\beta(X_{i,j})$ that X generates the terminal string from position i to j , and information about the most probable inside structure. If the dot of the dotted rule is at the rightmost position, the corresponding edge is an *inactive* edge. If the dot is at any other position, it is an *active* edge. Inactive edges represent recognized hypothetical constituents, while active edges represent prefixes of hypothetical constituents.

The i th terminal node t_i that enters the chart generates an inactive edge for the span $(i-1, i)$. Based on this, new active and inactive edges are generated according to the standard algorithm. Since we are interested in the most probable parse, the chart can be minimized in the following way while still performing an exhaustive search. If there is more than one edge that covers a span (i, j) having the same non-terminal symbol on the left-hand side of the dotted rule,

only the one with the highest inside probability is kept in the chart. The others cannot contribute to the most probable parse.

For an inactive edge spanning i to j and representing the rule $X \rightarrow Y^1 \dots Y^k$, the inside probability β_I is set to

$$\beta_I(X_{i,j}) = P(X \rightarrow Y_1 \dots Y_k) \prod_{l=1}^k \beta_I(Y_{i_l, j_l}^l) \quad (2)$$

where i_l and j_l mark the start and end position of Y^l , having $i = i_1$ and $j = j_1$. The inside probability for an active edge β_A with the dot after the k th symbol of the right-hand side is set to

$$\beta_A(X_{i,j}) = \prod_{l=1}^k \beta_I(Y_{i_l, j_l}^k) \quad (3)$$

We do not use the probability of the rule at this point. This allows us to combine all edges with the same span and the dot at the same position but with different symbols on the left-hand side. Introducing a distinguished left-hand side only for inactive edges significantly reduces the number of active edges in the chart. This goes one step further than implicitly right-binarizing the grammar; not only suffixes of right-hand sides are joined, but also the corresponding left-hand sides.

4 Memory Restrictions

We investigate the elimination (pruning) of edges from the chart in our incremental parsing scheme. After processing a word and before proceeding to the next word during incremental parsing, low ranked edges are removed. This is equivalent to imposing memory restrictions on the processing system.

The original algorithm keeps one edge in the chart for each combination of span (start and end position) and non-terminal symbol (for inactive edges) or right-hand side prefixes of dotted rules (for active edges). With pruning, we restrict the number of edges allowed per span. The limitation can be expressed in two ways:

1. *Variable beam.* Select a threshold $\theta \geq 1$. Edge e is removed, if its probability is p_e , the best probability for the span is p_1 , and

$$p_e < \frac{p_1}{\theta}. \quad (4)$$

¹The ROOT node is used internally for parsing; it is neither emitted nor counted for recall and precision.

2. *Fixed beam.* Select a maximum number of edges per span m . An edge e is removed, if its probability is not in the first m highest probabilities for edges with the same span.

We performed experiments using both types of beams. Fixed beams yielded consistently better results than variable beams when plotting chart size vs. F-score. Therefore, the following results are reported for fixed beams.

We compare and rank edges covering the same span only, and we rank active and inactive edges separately. This is in contrast to (Charniak et al., 1998) who rank all edges. They use normalization in order to account for different spans since in general, edges for longer spans involve more multiplications of probabilities, yielding lower probabilities. Charniak *et al.*'s normalization value is calculated by a different probability model than the inside probabilities of the edges. So, in addition to the normalization for different span lengths, they need a normalization constant that accounts for the different probability models.

This investigation is based on a much simpler ranking formula. We use what can be described as the unigram probability of a non-terminal node, i.e., the *a priori* probability of the corresponding non-terminal symbol(s) times the inside probability. Thus, for an inactive edge $\langle i, j, X \rightarrow \alpha, \beta_I(X_{i,j}) \rangle$, we use the probability

$$\begin{aligned} P_{RI}(X_{i,j}) &= P(X) \cdot P(t_i \dots t_{j-1} | X) \quad (5) \\ &= P(X) \cdot \beta_I(X_{i,j}) \quad (6) \end{aligned}$$

for ranking. This is the probability of the node and its yield being present in a parse. The higher this value, the better is this node. β_I is the inside probability for inactive edges as given in equation 2, $P(X)$ is the *a priori* probability for non-terminal X , (as estimated from the frequency in the training corpus) and P_{RI} is the probability of the edge for the non-terminal X spanning positions i to j that is used for ranking.

For an active edge $\langle i, j, X \rightarrow Y^1 \dots Y^k, Y^{k+1} \dots Y^m, \beta_A(Y_{i_1, j_1}^1 \dots Y_{i_k, j_k}^k) \rangle$ (the dot is after the k th symbol of the right-hand side) we use:

$$\begin{aligned} P_{RA}(Y_{i_1, j_1}^1 \dots Y_{i_k, j_k}^k) & \quad (7) \\ &= P(Y^1 \dots Y^k) \cdot P(t_i \dots t_{j-1} | Y^1 \dots Y^k) \quad (8) \end{aligned}$$

$$= P(Y^1 \dots Y^k) \cdot \beta_A(Y_{i_1, j_1}^1 \dots Y_{i_k, j_k}^k) \quad (9)$$

$P(Y^1 \dots Y^k)$ can be read off the corpus. It is the *a priori* probability that the right-hand side of a production has the prefix $Y^1 \dots Y^k$, which is estimated by

$$\frac{f(Y^1 \dots Y^k \text{ is prefix})}{N} \quad (10)$$

where N is the total number of productions in the corpus², $i = i_1, j = j_k$ and β_A is the inside probability of the prefix.

5 Experiments

5.1 Data

We use sections 2 – 21 of the Wall Street Journal part of the Penn Treebank (Marcus et al., 1993) to generate a treebank grammar. Traces, functional tags and other tag extensions that do not mark syntactic category are removed before training³. No other modifications are made. For testing, we use the 1578 sentences of length 40 or less of section 22. The input to the parser is the sequence of part-of-speech tags.

5.2 Evaluation

For evaluation, we use the parseval measures and report labeled F-score (the harmonic mean of labeled recall and labeled precision). Reporting the F-score makes our results comparable to those of other previous experiments using the same data sets. As a measure of the amount of work done by the parser, we report the size of the chart. The number of active and inactive edges that enter the chart is given for the exhaustive search, not counting those hypothetical edges that are replaced or rejected because there is an alternative edge with higher probability⁴. For pruned search, we give the percentage of edges required.

5.3 Fixed Beam

For our experiments, we define the beam by a maximum number of edges per span. Beams for active and inactive edges are set separately. The beams run from 2 to 12, and we test all

²Here, we use proper prefixes, i.e., all prefixes not including the last element.

³As an example, PP-TMP=3 is replaced by PP.

⁴The size of the chart is comparable to the “number of edges popped” as given in (Charniak et al., 1998).

Results with Original and Parent Encoding

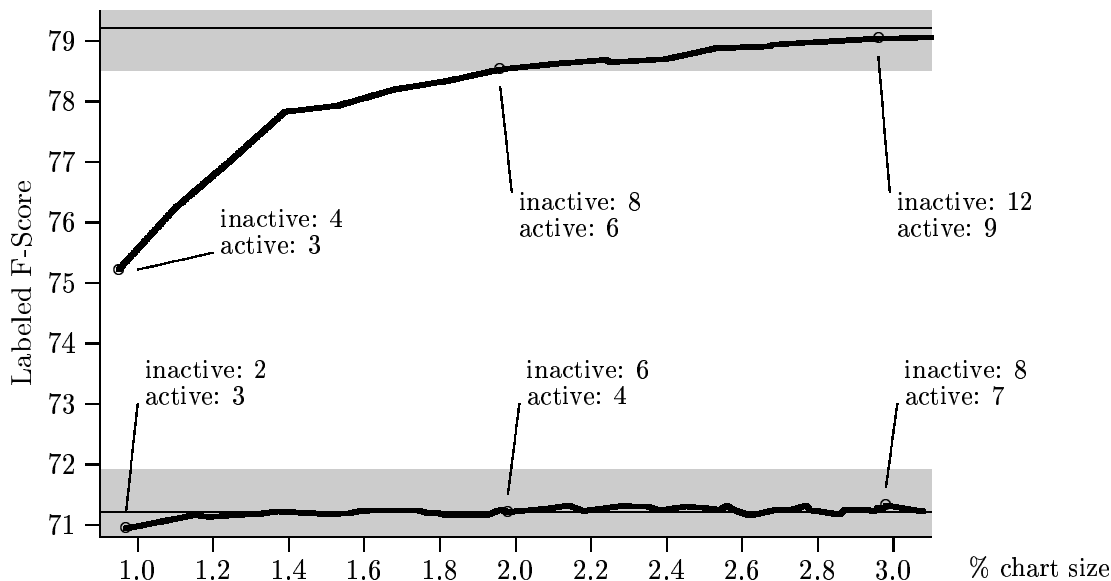


Figure 1: Experimental results for incremental parsing and pruning. The figure shows the percentage of edges relative to exhaustive search and the F-score achieved with this chart size. Exhaustive search yielded 71.21% for the original encoding and 79.28% for the parent encoding. Results in the grey areas are equivalent with a confidence degree of $\alpha = 0.99$.

121 combinations of these beams for active and inactive edges. Each setting results in a particular average size of the chart and an F-score, which are reported in the following section.

5.4 Experimental Results

The results of our 121 test runs with different settings for active and inactive beams are given in figure 1. The diagram shows chart sizes vs. labeled F-scores. It sorts chart sizes across different settings of the beams. If several beam settings result in equivalent chart sizes, the diagram contains the one yielding the highest F-score.

The main finding is that we can reduce the size of the chart to between 1% and 3% of the size required for exhaustive search without affecting the results. Only very small beams degrade performance⁵. The effect occurs for both models despite the simple ranking formula. This significantly reduces memory requirements

⁵Given the amount of test data (26,322 non-terminal nodes), results within a range of around 0.7% are equivalent with a confidence degree of $\alpha = 99\%$.

(given as size of the chart) and increases parsing speed.

Exhaustive search yields an F-Score of 71.21% when using the original Penn Treebank encoding. Only around 1% the edges are required to yield equivalent results with incremental processing and pruning after each word is added to the chart. This result is, among other settings, obtained by a fixed beam of 2 for inactive edges and 3 for active edges⁶

For the parent encoding, exhaustive search yields an F-Score of 79.28%. Only between 2 and 3% of the edges are required to yield an equivalent result with incremental processing and pruning. As an example, the point at size = 3.0% F-score = 79.1% is generated by the beam setting of 12 for inactive and 9 for active edges. The parent encoding yields around 8% higher F-scores but it also imposes a higher absolute and relative memory load on the process. The higher degree of parallelism in the inactive

⁶Using variable beams, we would need 1.95% of the chart entries to achieve an equivalent F-score.

chart stems from the parent hypothesis in each node. In terms of pure node categories, the average number of parallel nodes at this point is 3.5⁷.

Exhaustive search for the base encoding needs in average 140,000 edges per sentence, for the parent encoding 200,000 edges; equivalent results for the base encoding can be achieved with around 1% of these edges, equivalent results for the parent encoding need between 2 and 3%.

The lower number of edges significantly increases parsing speed. Using exhaustive search for the base model, the parser processes 3.0 tokens per second (measured on a Pentium III 500; no serious efforts of optimization have gone into the parser). With a chart size of 1%, speed is 630 tokens/second. This is a factor of 210 without decreasing accuracy. Speed for the parent model is 0.5 tokens/second (exhaustive) and 111 tokens/seconds (3.0% chart size), yielding an improvement by factor 220.

6 Related Work

Probably mostly related to the work reported here are (Charniak et al., 1998) and (Roark and Johnson, 1999). Both report on significantly improved parsing efficiency by selecting only a subset of edges for processing. There are three main differences to our approach. One is that they use a ranking for best-first search while we immediately prune hypotheses. They need to store a large number edges because it is not known in advance how many of the edges will be used until a parse is found. The second difference is that we proceed strictly incrementally without look-ahead. (Charniak et al., 1998) use a non-incremental procedure, (Roark and Johnson, 1999) use a look-ahead of one word. Thirdly, we use a much simpler ranking formula.

Additionally, (Charniak et al., 1998) and (Roark and Johnson, 1999) do not use the original Penntree encoding for the context-free structures. Before training and parsing, they change/remove some of the productions and introduce new part-of-speech tags for auxiliaries. The exact effect of these modifications is unknown, and it is unclear if these affect compa-

⁷For the active chart, parallelism cannot be given for different nodes types since active edges are introduced for right-hand side prefixes, collapsing all possible left-hand sides.

rability to our results.

The heavy restrictions in our method (immediate pruning, no look-ahead, very simple ranking formula) have consequences on the accuracy. Using right context and sorting instead of pruning yields roughly 2% higher results (compared to our base encoding⁸). But our work shows that even with these massive restrictions, the chart size can be reduced to 1% without a decrease in accuracy when compared to exhaustive search.

7 Conclusions

A central challenge in computational psycholinguistics is to explain how it is that people are so accurate and robust in processing language. Given the substantial psycholinguistic evidence for statistical cognitive mechanisms, our objective in this paper was to assess the plausibility of using wide-coverage probabilistic parsers to model human linguistic performance. In particular, we set out to investigate the effects of imposing incremental processing and significant memory limitations on such parsers.

The central finding of our experiments is that incremental parsing with massive (97% – 99%) pruning of the search space does not impair the accuracy of stochastic context-free parsers. This basic finding was robust across different settings of the beams and for the original Penn Treebank encoding as well as the parent encoding. We did however, observe significantly reduced memory and time requirements when using combined active/inactive edge filtering. To our knowledge, this is the first investigation on tree-bank grammars that systematically varies the beam for pruning.

Our aim in this paper is not to challenge state-of-the-art parsing accuracy results. For our experiments we used a purely context-free stochastic parser combined with a very simple pruning scheme based on simple “unigram” probabilities, and no use of right context. We do, however suggest that our result should apply to richer, more sophisticated probabilistic

⁸Comparison of results is not straight-forward since (Roark and Johnson, 1999) report accuracies only for those sentences for which a parse tree was generated (between 93 and 98% of the sentences), while our parser (except for very small beams) generates parses for virtually all sentences, hence we report accuracies for all sentences.

models, e.g. when adding word statistics to the model (Charniak, 1997).

We therefore conclude that wide-coverage, probabilistic parsers do not suffer impaired accuracy when subject to strict cognitive memory limitations and incremental processing. Furthermore, parse times are substantially reduced. This suggests that it may be fruitful to pursue the use of these models within computational psycholinguistics, where it is necessary to explain not only the relatively rare 'pathologies' of the human parser, but also its more frequently observed accuracy and robustness.

References

- Thorsten Brants. 1999. Cascaded Markov models. In *Proceedings of 9th Conference of the European Chapter of the Association for Computational Linguistics EACL-99*, Bergen, Norway.
- Eugene Charniak, Sharon Goldwater, and Mark Johnson. 1998. Edge-based best-first chart parsing. In *Proceedings of the Sixth Workshop on Very Large Corpora (WVLC-98)*, Montreal, Kanada.
- Eugene Charniak. 1996. Tree-bank grammars. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 1031–1036, Menlo Park: AAAI Press/MIT Press.
- Eugene Charniak. 1997. Statistical parsing with a context-free grammar and word statistics. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence*, pages 1031–1036, Menlo Park: AAAI Press/MIT Press.
- Nicholas Chater, Matthew Crocker, and Martin Pickering. 1998. The rational analysis of inquiry: The case for parsign. In Chater and Oaksford, editors, *Rational Models of Cognition*. Oxford University Press.
- Michael Collins. 1996. A new statistical parser based on bigram lexical dependencies. In *Proceedings of ACL-96*, Santa Cruz, CA, USA.
- Matthew Crocker and Thorsten Brants. 1999. Incremental probabilistic models of human linguistic performance. In *The 5th Conference on Architectures and Mechanisms for Language Processing*, Edinburgh, U.K.
- Matthew Crocker and Thorsten Brants. to appear. Wide coverage probabilistic sentence processing. *Journal of Psycholinguistic Research*, November 2000.
- Matthew Crocker and Steffan Corley. to appear. Modular architectures and statistical mechanisms: The case from lexical category disambiguation. In Merlo and Stevenson, editors, *The Lexical Basis of Sentence Processing*. John Benjamins.
- Matthew Crocker. 1999. Mechanisms for sentence processing. In Garrod and Pickering, editors, *Language Processing*. Psychology Press, London, UK.
- Mark Johnson. 1998. PCFG models of linguistic tree representations. *Computational Linguistics*, 24(4):613–632.
- Daniel Jurafsky. 1996. A probabilistic model of lexical and syntactic access and disambiguation. *Cognitive Science*, 20:137–194.
- Christopher Manning and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, Massachusetts.
- Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- Adwait Ratnaparkhi. 1997. A linear observed time statistical parser based on maximum entropy models. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing EMNLP-97*, Providence, RI.
- Brian Roark and Mark Johnson. 1999. Efficient probabilistic top-down and left-corner parsing. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics ACL-99*, Maryland.