

An Empirical Study on Statistical Disambiguation of Japanese Dependency Structures Using a Lexically Sensitive Language Model

SHIRAI Kiyooki and INUI Kentaro and
TANAKA Hozumi and TOKUNAGA Takenobu

Department of Computer Science, Graduate School of Information Science and Engineering,
Tokyo Institute of Technology

Abstract

We are proposing a new framework of statistical language modeling which integrates lexical association statistics with syntactic preference, while maintaining the modularity of those different statistics types, facilitating both training of the model and analysis of its behavior. In this paper, we report the result of an empirical evaluation of our model, where the model is applied to disambiguation of dependency structures of Japanese sentences. We also discussed the room remained for further improvement based on our error analysis.

1 Introduction

In the statistical parsing literature, it has already been established that statistics of lexical association have real potential for improvement of disambiguation performance. The question is how lexical association statistics should be incorporated into the overall statistical parsing framework. In exploring this issue, we consider the following four basic requirements:

- *Integration of different types of statistics:*
Lexical association statistics should be integrated with other types of statistics that are also expected to be effective in statistical parsing, such as short-term POS n-gram statistics and long-term structural preferences over parse trees.
- *Modularity of statistics types:*
The total score of a parse derivation should be decomposable into factors derived from different types of statistics, which would facilitate analysis of a model's behavior in terms of each statistics type.
- *Probabilistically well-founded semantics:*
The language model used in a statistical parser should have probabilistically well-founded semantics, which would also facilitate the analysis of the model's behavior.

- *Trainability:*

Since incorporation of lexical association statistics would make the model prohibitively complex, the model's complexity should be flexibly controllable depending on the amount of available training data.

However, it seems to be the case that no existing framework of language modeling (Black et al., 1993; Collins, 1996; Li, 1996; Magerman and Marcus, 1991; Magerman, 1995; Resnik, 1992; Schabes, 1992) satisfies these basic requirements simultaneously¹. In this context, we newly designed a framework of statistical language modeling taking all of the above four requirements into account (Inui et al., 1997a; Inui et al., 1997b). This paper reports on the results of our preliminary experiment where our framework was applied to structural disambiguation of Japanese sentences.

In what follows, we first briefly review our framework (Section 2). We next describe the setting of our experiment, including a brief introduction of Japanese dependency structures, the data sets, the baseline of the performance, etc. (Section 3). We then describe the results of the experiment, which was designed to assess the impact of the the incorporation of lexical association statistics (Section 4). We finally discuss the current problems revealed through our error analysis, suggesting some possible solutions (Section 5).

2 Overview of our framework

As with the most statistical parsing frameworks, given an input string A , we rank its parse derivations according to the joint distribution $P(R, W)$, where W is a word sequence candidate for A , and R is a parse derivation candidate for W whose terminal symbols constitute a POS tag sequence L (see Figure 1²). We first decompose $P(R, W)$ into two

¹For further discussion, see (Inui et al., 1997a). This is also the case with recent works such as (Charniak, 1997) and (Collins, 1997) due to the lack of modularity of statistical types.

²Although syntactic structure R is represented as a dependency structure in this figure, our framework does not impose

submodels, the syntactic model $P(R)$ and the lexical model $P(W|R)$:

$$P(R, W) = P(R) \cdot P(W|R) \quad (1)$$

The syntactic model, which is lexically insensitive, reflects both POS n-gram statistics and structural preference, whereas the lexical model reflects lexical association statistics. This division of labor allows for distinct modularity between the syntactic-based statistics and lexically sensitive statistics, while maintaining the probabilistically well-foundedness of the overall model.

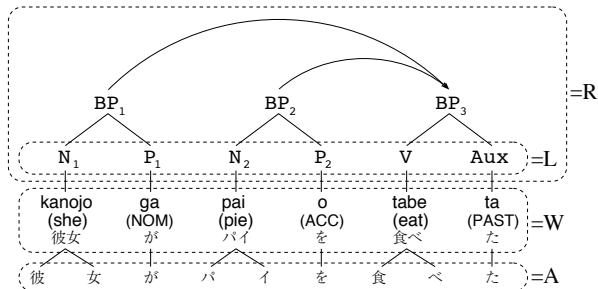


Figure 1. A parse derivation for an input string “彼女がパイを食べた (She ate a pie)”

2.1 The syntactic model

The syntactic model $P(R)$ can be estimated using a wide range of existing syntactic-based language modeling frameworks, from simple PCFG models to more context-sensitive models including those proposed in (Sekine and Grishman, 1995; Magerman and Marcus, 1991; Black et al., 1993). Among these, we, at present, use probabilistic GLR (PGLR) language modeling, which is given by incorporating probabilistic distributions into the GLR parsing framework (Inui et al., 1997c; Sornlertlamvanich et al., 1997)³. The advantages of PGLR modeling are (a) PGLR models are mildly context-sensitive, compared with PCFG models, and (b) PGLR models inherently capture both structural preferences and POS bigram statistics, which meets our integration requirement. For further discussion, see (Inui et al., 1997c).

2.2 The lexical model

The lexical model $P(W|R)$ is the product of the probability of each lexical derivation $l_i \rightarrow w_i$, where $l_i \in L$ ($L \subset R$) is the POS tag of $w_i \in W$:

$$P(W|R) = \prod_i P(w_i|R, w_1, \dots, w_{i-1}) \quad (2)$$

any restriction on the representation of syntactic structures.

³The idea of PGLR modeling was originally proposed by Briscoe and Carroll (Briscoe and Carroll, 1993). Inui et al. proposed a new formalization of PGLR models, resolving the drawbacks of Briscoe and Carroll’s model (Inui et al., 1997c; Sornlertlamvanich et al., 1997).

The key idea for estimating each factor $P(w_i|R, w_1, \dots, w_{i-1})$ (a lexical derivation probability) is in assuming that each lexical derivation depends only on a certain small part of its whole context. We first assume that syntactic structure R in $P(w_i|R, w_1, \dots, w_{i-1})$ can always be reduced to l_i ($\in R$), which allows us to deal with the lexical model separately from the syntactic model. The question then is which subset C of $\{w_1, \dots, w_{i-1}\}$ has the strongest influence on the derivation $l_i \rightarrow w_i$. We refer to a member of such a subset C as a *lexical context* of the derivation $l_i \rightarrow w_i$. As our point of departure, we consider the following two types of lexical associations in identifying lexical contexts:

- the lexical association between a head word (verb, adjective, nominal predicate, etc.) and its subordinated slot-marker,
- the lexical association between a head word coupled with a slot-marker, and the filler of that slot.

Let us illustrate this through the previous example shown in Figure 1. Suppose that the derivation order for W is head-driven, as given below, to guarantee that, for each of the words subordinated by a head word, the context of the derivation of that subordinated word always includes that head word.

$$\begin{aligned} ta \text{ (PAST)} &\rightarrow tabe \text{ (eat)} \rightarrow ga \text{ (NOM)} \rightarrow \\ o \text{ (ACC)} &\rightarrow kanojo \text{ (she)} \rightarrow pai \text{ (pie)} \end{aligned}$$

First, for each lexical item that is not either a slot-marker or a slot-filler, we simply assume that the probability of its derivation can be estimated independently of the derivation of the other words:

$$P(ta|R) \approx P(ta|Aux) \quad (3)$$

$$P(tabe|R, ta) \approx P(tabe|V) \quad (4)$$

Second, we estimate the probability of deriving each slot-marker, e.g. “*ga* (NOM)” and “*o* (ACC)”, by considering not only the dependency between the head word and each of its slot-markers, but also the dependency between slot-markers subordinated by the same head:

$$P(ga|R, tabe, ta) \simeq P(ga|P_1[h(tabe, [P_1, P_2])]) \quad (5)$$

$$P(o|R, ga, tabe, ta) \simeq P(o|P_2[h(tabe, [P_1:ga, P_2])]) \quad (6)$$

where $h(h, [s_1, \dots, s_n])$ is a lexical context denoting a head word h that subordinates the set of slots s_1, \dots, s_n , and $P(w_i|l_i[h(h, [s_1, \dots, s_n])])$ is the probability of a lexical derivation $l_i \rightarrow w_i$, given that w_i functions as a slot-marker of lexical head $h(h, [s_1, \dots, s_n])$.

Finally, we estimate the probability of deriving each slot-filler, e.g. “*kanojo* (she)”, in assuming that the derivation of a slot-filler depends only on its head word and slot:

$$\begin{aligned} P(kanojo|R, ga, o, tabe, ta) &\approx \\ P(kanojo|N[s(tabe, ga)]) &\quad (7) \end{aligned}$$

$$P(\text{pai}|R, \text{kanojo}, \text{ga}, \text{o}, \text{tabe}, \text{ta}) \approx P(\text{pai}|N[\text{s}(\text{tabe}, \text{o})]) \quad (8)$$

where $\text{s}(h, s)$ is a lexical context denoting a slot s of a head word h , and $P(w_i|l_i[\text{s}(h, s)])$ is the probability of a lexical derivation $l_i \rightarrow w_i$ given that w_i functions as a filler of a slot $\text{s}(h, s)$.

Combining equations (3), (4), (5), (6), (7) and (8), we produce (9):

$$P(W|R) \approx P(\text{ta}|Aux) \cdot P(\text{tabe}|V) \cdot P(\text{ga}|P[\text{h}(\text{tabe}, [P, P])]) \cdot P(\text{o}|P[\text{h}(\text{tabe}, [P:\text{ga}, P])]) \cdot P(\text{kanojo}|N[\text{s}(\text{tabe}, \text{ga})]) \cdot P(\text{pai}|N[\text{s}(\text{tabe}, \text{o})]) \quad (9)$$

Generalizing equation (9), we obtain the following equations, which state that the lexical model $P(W|R)$ can be estimated by the context-free lexical derivation model $P_{cf}(W|L)$ and the lexical association model $D(W|R)$:

$$P(W|R) \approx P_{cf}(W|L) \cdot D(W|R) \quad (10)$$

$$P_{cf}(W|L) = \prod_{i=1}^m P(w_i|l_i) \quad (11)$$

$$D(W|R) = \prod_{i=1}^m D(w_i|l_i[c_{w_i}]) \quad (12)$$

where c_{w_i} is the lexical context of w_i , and $D(w_i|l_i[c_{w_i}])$ is what we call a lexical dependency parameter, given by:

$$D(w_i|l_i[c]) = P(w_i|l_i[c]) / P(w_i|l_i) \quad (13)$$

$D(w_i|l_i[c])$ measures the degree of dependency between the lexical derivation $l_i \rightarrow w_i$ and its lexical context c . It is close to one, which means it is negligible, if w_i and c are highly independent. It becomes greater than one if w_i and c are positively correlated, whereas it becomes less than one and close to zero if w_i and c are negatively correlated.

The formulation of the lexical model as in (10) has two significant advantages. First, it localizes lexical association statistics into the lexical association model $D(W|R)$, which allows us to analyze the behavior of lexical association statistics in parsing independently of the other statistics types. Second, as we mention in Section 5, a lexical derivation can be associated with more than one lexical context (multiple lexical contexts) in such a case as a coordinate structure, relative clause, etc., which makes the space of lexical dependency parameters prohibitively large. However, given the definition of lexical dependency parameter (13), the parameter space can be reduced as follows:

$$D(w_i|l_i[C]) \approx \prod_{c \in C} D(w_i|l_i[c]) \quad (14)$$

where C is the set of the lexical contexts associated with the lexical derivation $l_i \rightarrow w_i$ ⁴. Note that

⁴For the proof, see (Inui et al., 1997a).

$P(w_i|l_i[C])$ cannot be decomposed in any similar manner.

The modularity of the lexical model also facilitates parameter estimation. Although the syntactic model ideally requires *fully* bracketed training corpora, training it is expected to be manageable since the model’s parameter space tends to be only a small part of the overall parameter space. The lexical association statistics, on the other hand, may have a much larger parameter space, and thus may require much larger amounts of training data, as compared to the syntactic model. However, since our lexical model can be trained independently of syntactic preference, one can train it using *partially* parsed tagged corpora, which can be produced at a lower cost (i.e. automatically), as well as fully bracketed corpora. In fact, we used both a full-bracketed corpus and a partially parsed corpus in our experiment.

3 A preliminary experiment

Let us first briefly describe some fundamental features of Japanese syntax. A Japanese sentence can be analyzed as a sequence of so-called *bunsetsu* phrases (BPs, hereafter) as illustrated in Figure 1. A BP is a chunk of words consisting of a content word (noun, verb, adjective, etc.) accompanied by some function word(s) (postposition, auxiliary, etc.). For example, the BP “*kanojo-ga*” (BP_1) in Figure 1 consists of the noun “*kanojo* (she)” followed by the postposition “*ga* (NOM)”, which functions as a case-marker. The BP “*tabe-ta*” (BP_3), on the other hand, consists of the verb “*tabe* (eat)” followed by the auxiliary “*ta* (PAST)”.

Given a sequence of BPs, one can recognize dependency relations between them as illustrated in Figure 1. In Japanese, if BP_i precedes BP_j , and BP_i and BP_j are in a dependency relation, then BP_i is always the modifier of BP_j , and we say “ BP_i modifies BP_j .” For example, in Figure 1, both BP_1 and BP_2 modify BP_3 .

For the preliminary evaluation of our model, we restricted our focus only on the model’s performance for structural disambiguation excluding morphological disambiguation. Thus, the task of the parser was restricted to determination of the dependency structure of an input sentence, which is given together with the specification of word segments, their POS tags, and the boundaries between BPs.

In developing the grammar used by our PGLR parser, we first established a categorization of BPs based on the POS of their constituents: postpositional BPs, verbal BPs, nominal predicative BPs, etc. We then developed a modification constraint matrix that describes which BP category can modify which BP category, based on examples collected from the EDR Japanese corpus (EDR, 1995). We finally transformed this matrix into a CFG; for in-

stance, the constraint that a BP of category C_i can modify a BP of category C_j can be transformed into context-free rules such as $\langle \bar{C}_j \rightarrow C_i C_j \rangle$, $\langle \bar{C}_j \rightarrow \bar{C}_i C_j \rangle$, etc., where \bar{X} denotes a nonterminal symbol.

For the text data, we used roughly 22,000 sentences⁵ collected from the EDR bracketed corpus for training the syntactic model, and the whole EDR corpus and the RWC POS-tagged corpus (Real World Computing Partnership, 1995) for training the lexical model. For testing, we used 784 sentences collected from the EDR corpus with the average sentence length being 6.4 BPs. The data sets used for training and testing are mutually exclusive. The grammar used by our probabilistic GLR parser was a CFG automatically acquired from the training sentences, consisting of 590 context-free rules containing 39 nonterminal symbols and 30 terminal symbols (i.e. BP categories).

The baseline of the disambiguation performance was assessed by way of a naive strategy which selects the nearest possible modifiee (similarly to the right association principle in English) under the non-crossing constraint. The performance of this naive strategy was 60.4% in BP-based accuracy, where BP-based accuracy is the ratio of the number of the BPs whose modifiee is correctly identified to the total number of BPs (excluding the two rightmost BPs for each sentence).

The contribution of the syntactic model $P(R)$ to structural disambiguation can be assessed by removing the lexical association model, namely:

$$P(R, W) \approx P(R) \cdot \prod_i P(w_i | l_i) \quad (15)$$

This lexically insensitive model achieved 74.3% in BP-based accuracy, 13.9 points above the baseline.

4 The contribution of the lexical model

In our experiment, we considered as the major factors of the lexical model (a) $D(p|P[h(h, [s_1, \dots, s_n])])$, the dependencies between slot markers and their lexical head, and (b) $D(n|N[s(v, p)])$, the dependencies between case fillers and their head verb coupled with the corresponding case markers.

$D(p|P[h(h, [s_1, \dots, s_n])])$ can be computed from $P(p^n | P^n[h(h, [])])$, the distribution of n postpositions (case markers) given that all of them are subordinated by a single lexical head h . We trained this distribution using 150,000 instances of p^n - $\{verb, adjective, nominal_predicate\}$ collocation collected from the EDR full-bracketed corpus. For parameter estimation, we used the maximum entropy

⁵We collected only sentences associated with a complete binary tree. The test set was collected according to the same criterion.

estimation technique (Berger et al., 1996; Ratnaparkhi et al., 1994). For further details of this estimation process, see (Shirai et al., 1997).

$D(n|N[s(v, p)])$ was trained using 6.7 million instances of *noun-postposition-verb* collocation collected from both the EDR and RWC corpora. For parameter estimation, we used 115 non-hierarchical semantic noun classes derived from the NTT semantic dictionary (Ikehara et al., 1997) to reduce the parameter space:

$$D(n|N[s(v, p)]) \approx \frac{\sum_{c_n} P(c_n | N[s(v, p)]) \cdot P(n | c_n)}{P(n | N)} \quad (16)$$

$P(c_n | N[s(v, p)])$ was estimated using a simple back-off smoothing technique: for any given lexical verb v and postposition p , if the frequency of $s(v, p)$ is less than a certain threshold λ (in our experiment, $\lambda = 100$), then $P(c_n | N[s(v, p)])$ was approximated to be $P(c_n | N[s(c_v, p)])$ where c_v is a class of v whose frequency is more than λ . Obviously, the more abstract the class chosen for c_v is, the closer to one $D(n|N[s(v, h)])$ becomes.

Table 1 summarizes the results of the experiment. “+filler” denotes the setting where dependency parameters for slot fillers are considered, but not those for slot markers. “+marker” denotes the setting where dependency parameters for slot markers are considered, but not those for slot fillers. “+both” denotes the setting where the both types of parameters are considered. As shown in the table, we achieved a 12.3 point gain in BP-based accuracy over the figure for the syntactic model of 74.3%, by incorporating lexical association statistics.

Table 1. The contribution of the lexical model

	accuracy
base line	60.4 %
syntactic model	74.3 %
+filler	75.7 %
+marker	86.5 %
+both	86.6 %

5 Error analysis

In the test set, there were 457 BPs whose modifiee was not correctly identified by the system. Among these errors, we particularly explored 169 errors that were associated with postpositional BPs functioning as a case of either a verb, adjective, or nominal predicate, since, for lexical association statistics in the lexical model, we took only the dependencies between cases (i.e. case markers and case fillers) and their heads into account. In this exploration, we identified three major error types: (a) errors associated with a coordinate clause (43 cases), (b) errors associated with relative clauses (34 cases), (c) errors associated with compound predicates (28 cases).

5.1 Coordinate structures

Figure 2 illustrates a typical error associated with a coordinate clause. The sentence in this figure has at least three alternative interpretations in terms of which BP is modified by the left-most BP “*kanojo-wa* (she-TOP)”: (a) “*tabe-ta* (eat-PAST)”, (b) “*dekake-ta* (leave-PAST)”, (c) both “*tabe-ta* (eat-PAST)” and “*dekake-ta* (leave-PAST)”. Among these alternatives, the most reasonable interpretation is obviously (c), where the two predicative BPs constitute a coordinate structure.

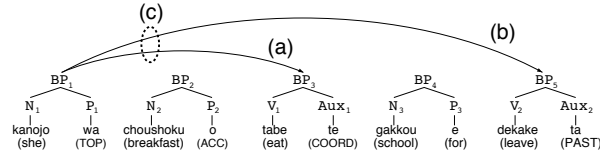


Figure 2. An example sentence containing a coordinate structure: “She ate breakfast and left for school”

In our experiment, however, neither the training data nor the test data indicates such coordinate structures. Thus, in the above sentence, for example, the system was required to choose one of two alternatives (a) and (b), where (b) is the preferred candidate according to the structural policy underlying our corpora. However, this choice is not really meaningful. Furthermore, the system systematically prefers (a), the wrong choice, since (i) the syntactic model tends to prefer shorter-distance modification relations (similarly to the right association principle in English), and (ii) the lexical model is expected to support both candidates because both $D(kanojo|N[s(tabe, wa)])$ in (a) and $D(kanojo|N[s(dekake, wa)])$ in (b) should be high. This problem makes the performance of our model lower than what it should be.

Obviously, the first step to resolving this problem is to enhance our corpora and grammar to enable the parser to generate the third interpretation, i.e. to explicitly generate a coordinate structure such as (c) if needed. Once such a setting is established, we then need to consider the lexical contexts of each of the constituents modifying a coordinate structure, such as “*kanojo-wa* (she-TOP)” in the above sentence. In interpretation (c), since “*kanojo-wa* (she-TOP)” modifies both predicative BPs, it is reasonable to associate it with two lexical contexts, $s(tabe, wa)$ and $s(dekake, wa)$. As mentioned in Section 2, our framework allows us to deal with such multiple lexical contexts, namely:

$$\begin{aligned} & D(kanojo|N[s(tabe, wa), s(dekake, wa)]) \\ & \approx D(kanojo|N[s(tabe, wa)]) \cdot \\ & D(kanojo|N[s(dekake, wa)]) \end{aligned} \quad (17)$$

5.2 Treatment of coreference

One may have already noticed that the issue discussed above can be generalized as an issue associated with the treatment of coreference in dependency

structures. Namely, if a prepositional BP is coreferred to by more than one clause as a participant, a naive treatment of this coreference relation could require the parser to make a meaningless choice: which clause subordinates that BP. This problem in the treatment of coreference is considered to cause a significant proportion of errors associated with relative/adverbial clauses or compound predicates. Such errors are expected to be resolvable through an extension of the model, as discussed in Section 5.1.

Let us briefly look at another example in Figure 3, where the matrix clause and relative clause corefer to the leftmost BP “*kanojo-wa* (she-TOP)”, i.e. interpretation (c). Without any refined treatment of this coreference relation, the parser would be required to make a meaningless choice between (a) and (b).

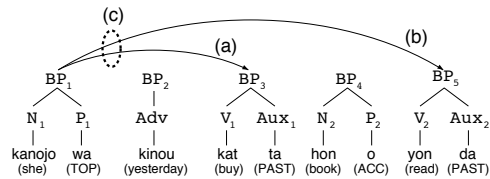


Figure 3. An example sentence containing a relative clause: “She read the book which she bought yesterday”

5.3 Dependency between slot fillers

According to the results summarized in Table 1, the contribution of the dependency between case fillers and their heads seems to be negligibly small. We can enumerate several possible reasons including: (a) most of the sentences used for testing were fairly short, and thus not ambiguous enough to need the help of these types of dependencies, and (b) the estimation of these types of dependency parameters was not sufficiently sophisticated.

In addition to these reasons, we also found that the lack of the consideration of dependency between case markers was also problematic in some cases; there are particular patterns where dependency between case fillers seems to be highly significant. For example, in the clause “*kanojo-wa* (she-TOP) *isha-ni* (doctor-DAT) *nat-ta* (become-PAST)” (she became a doctor), the distribution of the filler of the “*wa* (TOP)” slot is considered to be highly dependent on the filler of the “*ni* (DAT)” slot, “*isha* (doctor)”, since its distribution would be markedly different if “*isha* (doctor)” was replaced with “*mizu* (water)”. Similar patterns include, for example, “*A-wo* (ACC) *B-ni* (DAT) *suru* (make)”, where *A* and *B* are highly dependent, and “*A-ga* (NOM) *B-wo* (ACC) *suru* (do)”, where noun *B* indicating an action strongly influences the distribution of *A*.

In our framework, this type of problem can be treated by means of controlling the choice of lexical contexts. We are now conducting another experiment in which the dependencies between case

fillers are additionally considered in particular patterns. Note that the refinement of our model in this manner illustrates that the modularity of lexical association statistics facilitates rule-based control in choosing the locations where lexical association is considered. This rule-based control allows us to incorporate qualitative knowledge such as linguistic insights and heuristics newly obtained from experiments based on the model.

6 Conclusion

In this paper, we first presented a new framework of language modeling for statistical parsing, which incorporates lexical association statistics while maintaining modularity. We then reported on the results of our preliminary evaluation of the model's performance, showing that both the syntactic and lexical models made a considerable contribution to structural disambiguation, and that the division of labor between those two models thus seemed to be working well to date. We also discussed the fact that room remains for further improvement, suggesting that, when considering lexical association, we need to carefully deal with structures including coreference relations; constituents in coreference relations need to be associated with multiple lexical contexts, which can be treated by introducing the notion of lexical dependency parameters.

Many issues remain unclear. First, most of the sentences used so far in testing are quite short, and thus not highly ambiguous, which made difficult the empirical evaluation of the model's quality. As such, we need to explore the model's behavior when it is applied to longer sentences. Second, we also need to conduct experiments on the combination of the morphological and syntactic disambiguation tasks, which our framework intrinsically is designed for. Third, empirical comparison with other lexically sensitive models is also strongly required. One interesting issue is whether the division of labor between the syntactic and lexical models presented in this paper works well language-independently, or conversely, whether the existing models designed for English are equally applicable to languages like Japanese.

Acknowledgements

The authors would like to thank the staff of NTT for making available their considerable electronic resources.

References

- A. L. Berger, S. A. Della Pietra, and V. J. Della Pietra. 1996. A maximum entropy approach to natural language processing. *Computational Linguistics*, 22(1):39–71.
- E. Black, F. Jelinek, J. Lafferty, D. M. Magerman, R. Mercer, and S. Roukos. 1993. Towards history-based grammars: Using richer models for probabilistic parsing. In *Proceedings of the ACL*, pages 31–37.
- T. Briscoe and J. Carroll. 1993. Generalized probabilistic LR parsing of natural language (corpora) with unification-based grammars. *Computational Linguistics*, 19(1):25–59.
- E. Charniak. 1997. Statistical parsing with a context-free grammar and word statistics. In *Proceedings of the AAAI*.
- M. Collins. 1996. A new statistical parser based on bigram lexical dependencies. In *Proceedings of the ACL*.
- M. Collins. 1997. Three generative, lexicalised models for statistical parsing. In *Proceedings of the ACL*.
- EDR. 1995. The EDR electronic dictionary technical guide (second edition). Technical Report TR-045, Japan Electronic Dictionary Research Institute.
- S. Ikehara, M. Miyazaki, S. Shirai, A. Yokoo, H. Nakaiwa, K. Ogura, Y. Ooyama, and Y. Hayashi. 1997. *A Japanese Lexicon*. Iwanami Shoten.
- K. Inui, K. Shirai, H. Tanaka, and T. Tokunaga. 1997a. Integrated probabilistic language modeling for statistical parsing. Technical Report TR97-0005, Dept. of Computer Science, Tokyo Institute of Technology. <ftp://ftp.cs.titech.ac.jp/lab/tanaka/papers/97/inui97b.ps.gz>.
- K. Inui, K. Shirai, T. Tokunaga, and H. Tanaka. 1997b. Integration of statistical techniques for parsing. In *summary collection of the IJCAI'97 poster session*.
- K. Inui, V. Sornlertlamvanich, H. Tanaka, and T. Tokunaga. 1997c. A new formalization of probabilistic GLR parsing. In *Proceedings of the IWPT*.
- H. Li. 1996. A probabilistic disambiguation method based on psycholinguistic principles. In *Proceedings of WVLC-4*.
- D. M. Magerman and M. Marcus. 1991. Pearl: A probabilistic chart parser. In *Proceedings of the EACL*, pages 15–20.
- D. M. Magerman. 1995. Statistical decision-tree models for parsing. In *Proceedings of the ACL*, pages 276–283.
- A. Ratnaparkhi, J. Reyrner, and S. Roukos. 1994. A maximum entropy model for prepositional phrase attachment. In *Proceedings of the Human Language Technology Workshop*, pages 250–255.
- Real World Computing Partnership. 1995. RWC text database. <http://www.rwcp.or.jp/wswg.html>.
- P. Resnik. 1992. Probabilistic tree-adjointing grammar as a framework for statistical natural language processing. In *Proceedings of the COLING*, pages 418–424.
- Y. Schabes. 1992. Stochastic lexicalized tree-adjointing grammars. In *Proceedings of the COLING*, pages 425–432.
- S. Sekine and R. Grishman. 1995. A corpus-based probabilistic grammar with only two non-terminals. In *Proceedings of the IWPT*.
- K. Shirai, K. Inui, T. Tokunaga, and H. Tanaka. 1997. Learning dependencies between case frames using maximum entropy method. In *Proceedings of Annual Meeting of the Japan Association for Natural Language Processing*. (In Japanese).
- V. Sornlertlamvanich, K. Inui, K. Shirai, H. Tanaka, T. Tokunaga, and T. Takezawa. 1997. Empirical evaluation of probabilistic glr parsing. In *Proceedings of the NLPRS*.