

# Isolating Cross-linguistic Parsing Complexity with a Principles-and-Parameters Parser: A Case Study of Japanese and English \*

Sandiway Fong & Robert C. Berwick

NEC Research Institute, 4 Independence Way, Princeton, NJ, 08540 USA, sandiway@research.nec.com  
Rm 838, MIT AI Laboratory, 545 Technology Sq., Cambridge, MA 02139, berwick@ai.mit.edu

## 1 Introduction

As parsing models and linguistic theories have broadened to encompass a wider range of non-English languages, a particularly useful "stress test" is to build a single theory/parser pair that can work for multiple languages, in the best case with minor variation, perhaps restricted to the lexicon. This paper reports on the results of just such a test applied to a fully operational (Prolog) implementation of a so-called principles-and-parameters model of syntax, for the case of Japanese and English. This paper has two basic aims: (1) to show how an implemented model for an entire principles-and-parameters model (essentially all of the linguistic theory in Lasnik & Uriagereka (1988)), see figure 2 for a computer snapshot, leads directly to both a parser for multiple languages and a useful "computational linguistics workbench" in which one can easily experiment with alternative linguistic theoretical formulations of grammatical principles as well as alternative computational strategies; (2) to use this system to uncover sources of parsing complexity in Japanese as opposed to English. In particular, we examine the "null hypothesis" that a single parsing design suffices for efficient processing of both Head-first and Head-final languages, in contrast to approaches that posit, e.g., a right-corner or other mirror-image strategy for parsing Japanese as compared to English (e.g., BUP; Mazuka (1990)). In this case we can confirm computationally and precisely, in accordance with much current psycholinguistic work (Frazier and Raynett (1988); Inoue and J.D. Fodor (1991); Nagai (1991)) that it is *not* the Head-final character of Japanese that results in processing difficulty so much as the possibility of scrambling and free deletion of NPs (so-called "super Pro Drop"). We do this by empirically investigating the effects of 3 possible "optimizations" of the parsing system for Japanese: (1) the use of right-context information via automatic source transformations, using a programming language compiler technique to introduce dummy nonterminals and corresponding semantic actions; (2) modification of the Japanese grammar to put the specifier of CP

(=  $\bar{S}$ ) on the right and so eliminate unnecessary center-embedding; and (3) eliminating of scrambling and NP drop to isolate the separate effects of Head-final (e.g., Verb-final) phrase structure in Japanese.

By explicit construction, the implementation demonstrates that it is possible to build an efficient principle-and-parameters parser for multiple languages, using 25 principles that are expressed in a language quite close in form to that of the original linguistic theory. The English-Japanese differences handled include the basic Subject-Object-Verb (SOV) order of Japanese; free "scrambling" of Japanese noun phrases; topic-comment structure; nonappearance of noun phrases that are discourse recoverable; and lack of *wh*-word movement in Japanese questions. No rule reprogramming is required to accommodate these differences, but changes to only 4 binary switches and a minimally distinct lexicon with different thematic grids in some cases. The parser couples several already-known parsing design strategies to obtain efficient parsing times, e.g., type-checking; multiple-entry canonical LR(1) parsing; and automatic (source-to-source) grammar transformations.<sup>1</sup>

## 2 Principle-based parsing

In a principle-based parser, construction- and language-specific rules are replaced with broader principles that remain invariant aside from parametric variation (see below). The parser works by a (partially interleaved) generate-and-test technique that uses a canonical LR(1) covering grammar (derived from  $\bar{X}$  theory plus the theory of movement) to first build an initial set of tree structures; these structures are then run through a se-

<sup>1</sup>To the best of our knowledge, this system is the first and broadest-coverage of its type to be able to parse Japanese and English by setting just a few parameter switches. Dorr (1987), under the supervision of the second author, developed a conceptually similar scheme to handle English, Spanish, and German. However, Dorr's system did not have the same broad coverage of English; did not handle Japanese; used hand rather than automatic compiling; and was approximately 15 times slower. Gunji's (1987) Japanese unification grammar comes closest to the principle-based model, but requires hand-modification from a set of core principles and does not really accommodate the important Japanese phenomenon of scrambling; see below. Other such systems work only on much smaller parts of English, e.g., Sharp (1985); Wehrli (1987); Crocker (1989); Correa (1988); Johnson, (1989); or are not in fact parsers, but proof-checkers, e.g., Stabler, (1991, forthcoming).

\*This research has been supported by NSF Grant DCR85552543 under a Presidential Young Investigator Award to Professor Robert C. Berwick, and a grant from the Kapor Family Foundation. We would like to thank Howard Lasnik, Alec Marantz, Shigeru Miyagawa, David Pesetsky, and Mamoro Saito for valuable discussions and valiant attempts to tell us about Japanese.

ries of predicates whose conjunction defines the remainder of the constraints the (sentence, phrase structure, LF) triple must satisfy. This is done using familiar machinery from Prolog to output LFs that satisfy all the declarative constraints of the linguistic theory. In practice, a straightforward generate-and-test mechanism is grossly inefficient, since the principles that apply to at the level of surface structure (S-structure) are but a fraction of those that apply in the overall system. The usual problems of lexical and structural ambiguity the the underconstrained nature of the initial  $\bar{X}$  system means that the number of possible S-structures to hypothesize may be huge. To obtain an efficient parser we use a full multiple-entry table with backtracking (as in Tomita, 1986), extending it to a canonical LR(1) parser. The LR machine uses an automatically-built S-structure grammar that folds together enough of the constraints from other principles, parameters, lexical subcategory information offline to produce a 25-fold improvement over the onliue phrase structure recovery procedure originally proposed by Fong and Berwick (1989). Optimizations include extra conditions in action clauses to permit interleaving of other principles (like movement) with structure-building (the 'interleaving' noted by principles marked 'I' in the snapshot in figure 2 below); control structure flexibility in principle ordering; precomputation of the LR transition function; elimination of infinite recursion of empty elements by an additional stack mechanism, and so forth. We exploit the explicit modularity of the principle-based system in way that is impossible in an ordinary rule-based system: we can build a grammar for phrase structure that is small enough to make full, canonical LR(1) parsing usable, unlike large CFGs. The earlier error detection of full LR(1) parsing over LALR methods means that fail as early as possible, to avoid expensive tree constructions that can never participate in final solutions.<sup>2</sup>

### 3 The Japanese parser

We begin with a very simple parameterization of Japanese that will nonetheless be able to cover all the Lasnik and Saito *wh*-questions, scrambling, and so forth; see the table on the next page that follows the example sentences. The important point is that very little additional must be said in order to parse a wide variety of distinctive Japanese sentences; the principles as shown on the righthand side of the computer snapshot do not change.<sup>3</sup>

Consider first the example *wh*-movement sentences found in the linguistics paper *On the Nature of Proper Government* by Lasnik & Saito (1984).<sup>4</sup> These sen-

<sup>2</sup>To provide a rough measure of the machine size for the phrase structure grammar of S-structure for both English and Japanese, the augmented CFG consists of about 74 productions derived from a schema of 30-34 rules. The resulting characteristic finite state automaton (CFSM) consists of 123 states with 550 transitions between the various states. The action table consists of a total of 984 individual (nonerror) entries.

<sup>3</sup>We will scramble only from direct object positions here, even though it is straightforward to scramble from indirect object positions. Informally, we have noted that scrambling from IO greatly increases computation time. A tighter set of constraints on scrambling seems called for.

<sup>4</sup>Average best parsing time for the Japanese sentences shown is 0.37secs/word on a Symbolics 3650 (64K LIPS) ( $\sigma$

tences (listed below) display many familiar typological Japanese-English differences, and cover a rather sophisticated set of differences between English and Japanese: for instance, why (6) is fine in Japanese but not in English; free omission of NPs; "scrambling" of subjects and objects; Verb-final (more generally, Head-final) constituent structure, and no overt movement of *wh*-phrases. We also consider a different set of Japanese sentences (also listed below) designed to illustrate a range of the same phenomena, taken from Hosokawa (1990). We stress that these sentences are designed to illustrate a range of sentence distinctions in Japanese, as well as our investigative method, rather than serve as any complete list of syntactic differences between the two languages (since they are obviously not).<sup>5</sup>

[Lasnik & Saito (1984)]

- (2) Watashi-wa Taro-ga nani-o katta ka shitte iru  
'I know what John bought'
- (6) Kimi-wa dare-ni Taro-ga naze kubi-ni natta tte itta no  
'To whom did you say that John was fired why'
- (32) \*Meari-wa Taro-ga nani-o katta ka do ka shiranai  
'Mary does not know whether or not John bought what'
- (37a) Taro-wa naze kubi-ni natta no  
'Why was John fired'
- (37b) Biru-wa Taro-ga naze kubi-ni natta tte itta no  
'Why did Bill say that John was fired'
- (39a) Taro-ga nani-o te-ni ireta koto-o sonnani okotteru no  
'What are you so angry about the fact that Taro obtained'
- (39b) \*Taro-ga naze sore-o te-ni ireta koto-o sonnani okotteru no  
'Why are you so angry about the fact that Taro obtained it'
- (41a) Hanoko-ga Taro-ga nani-o te-ni ireta tte itta koto-o sonnani okotteru no  
'What are you so angry about the fact that Hanoko said that Taro obtained'
- (41b) \*Hanoko-ga Taro-ga naze sore-o te-ni ireta tte itta koto-o sonnani okotteru no  
'Why are you so angry about the fact that Hanoko said that Taro obtained it'
- (60) Kimi-wa nani-o doko-de katta no  
'Where did you buy what'
- (63) Kimi-wa nani-o sagashiteru no  
'Why are you looking for what'

### Complement/noncomplement asymmetry, scrambling and unexpected parses

To see how the parser handles one Japanese example (see the actual computer output in figure 1 or figure 2), consider (39a) (and the corresponding illicit (39b)), where a complement *wh* but not a noncomplement *wh* can be extracted from a complex NP: (a) Taro-ga nani-o te-ni ireta koto-o sonnani okotteru no; (b) \*Taro-ga naze sore-o te-ni ireta koto-o 'What/\*Why are you so angry about the fact that Taro obtained'

This example illustrates several Japanese typological differences with English. The subject of the matrix clause (= you) has been omitted. *Nani* ('what') and *te* ('and') have been scrambled; the direct object

= 1.52sec,  $n=100$ ). Parsing time on a Sun Sparcstation 2 is approximately an order of magnitude faster.

<sup>5</sup>E.g., the double- $\sigma$  constraint; case-overwriting, passive and causative constructions, etc. all remain to be fully implemented.

(marked -o) now appearing in front of the indirect object *te*. Phrase structure is Ilead final. Our relaxation of the Case Adjacency parameter and the rule that allows adjunction of NP' to VP, plus transmission of Case to the scrambled NP' will let this analysis through. The LF for this sentence should be something along the lines of: for what *x*, *pro* is so angry about [the fact that Taro obtained *x*]

In this example *pro* denotes the understood subject of *okotteru* ("be angry"). The LF's actually returned by the parser are shown in the snapshot in figure 1.<sup>6</sup>

[Hosokawa (1990)]

- (1b)' Gengogaku-no gakusei-ga tiizu-o tabeta  
linguistics-gen student-nom cheese-acc eat-past  
'A student of linguistics ate cheese'
- (2b)' Nagai kami-no gakusei-ga tiizu-o tabeta  
long hair-gen student-nom cheese-acc eat-past  
'A long haired student ate cheese'
- (3b) Taro-ga hon-o katta  
John-nom book-acc buy-past  
'John bought a book'
- (4b) Taro-ga Hanoko-ni hon-o ageta  
John-nom Mary-dat book-acc give-past  
'John bought Mary a book'
- (5b) Taro-ga hon-o table-no ue-ni oita  
John book-acc table-gen top-dat (top of table)  
put-past  
'John put the book on the table'
- (6b) Taro-wa gakkoo-ni itta  
John-top school-dat go-past  
'John went to school'
- (15b) Watashi-wa taro-ga nani-o katta ka shiranai  
I-top John-nom what-acc bought Q know-not  
'I don't know what John bought'
- (17b) Taro-wa Chomsky-no Barriers-o yomimashita ka  
John-top Chomsky-gen Barriers-acc read-past Q  
'Did John read Chomsky's Barriers'
- (18b) Hanoko-wa  
Biru-ga Chomsky-no Barriers-o yonda ka do ka  
shiranai  
Mary-top Bill-nom Chomsky-gen Barriers-acc  
read Q know-not  
'Mary does not know whether or not Bill read  
Chomsky's...'

The parametric differences that we need to accommodate all these differences between English and Japanese are quite few:

<sup>6</sup>We will not have room to describe in detail the derivation of these LF's. But, it should be noted that the derivation sequence is quite complex. Note, for example, that *nani* ('what') undergoes movement at two levels of phrase structure in order to get to the specifier position of the matrix Complementizer: [CP nani[IP Taro[NP[CP pro[VP t' [VP te t ireta]] koto]. . ]]

Furthermore, the LF trace *t'* violates the so-called empty category principle unless it is deleted (as indicated by [] in the snapshot), under the present theory. The lack of *wh*-movement at S-structure in Japanese, and its presence in English, interacts with these constraints to bar examples like (6) in English; see Lasnik & Saito.

English and Japanese parameter settings		
Parameter	English	Japanese
Spec order	specInitial. specFinal :- \+ specInitial.	specInitial. specFinal :- \+ specInitial.
*Head order	headInitial. headFinal :- \+ headInitial.	headFinal. headInitial :- \+ headFinal.
Agreement Bounding	agr(weak). boundingNode(i2). boundingNode(np).	agr(weak). boundingNode(i2). boundingNode(np).
*Case Adjacency *WA in Syntax *Pro-Drop	caseAdjacency. whInSyntax. :- no proDrop.	:- no caseAdjacency. :-no whInSyntax. proDrop.

As one can see from the figure, the system does correctly recover the right LF, as the last one in snapshot. However, it also (surprisingly) discovers three *additional* LF's, illustrating the power of the system to uncover alternative interpretations that a proper theory of context would have the job of ruling out. Ignoring indices, they all have the same form: for what *x*, Taro is so angry about [the fact that *pro* obtained *x*]

Here the embedded subject Taro has been interchanged with the matrix subject *pro*. It turns out that the sentence happens to be ambiguous with respect to the two basic interpretations.<sup>7</sup> For completeness, here are the three variants of that correspond to the first three LF's reported by the parser. S. Miyagawa (p.c.) informs us that the last two, given proper context, are in fact possible. These include: (1) *pro* is coreferent with *koto* ("fact");<sup>8</sup> i.e., for what *x*, Taro is so angry about [the fact that the fact obtained *x*]; (2) *pro* is coreferent with *taro*: for what *x*, Taro is so angry about [the fact that Taro obtained *x*]; and (3) *pro* is free in the sentence: for what *x*, Taro is so angry about [the fact that (someone else) obtained *x*].<sup>9</sup>

#### 4 Parsing Japanese: the computational effects of scrambling, *pro*-drop, and phrase structure

Next we turn to the investigation of the *computational* differences between the two languages that we have explored, and show how to use the system in an exploratory mode to discover complexity differences between English and Japanese. In the discussion that follows, we shall need to draw on comparisons between the complexity of different parses. While this is a delicate matter, there are two obvious metrics to use in comparing this parser's complexity. The first is the total number of principle operations used to analyze a sentence -- the number of S-structures, chain formations, indexings, the case filter and other constraint applications, etc. We can treat these individually and as a whole to give an account of the entire "search space" the parser moves through to discover analyses. However, this is

<sup>7</sup>This was pointed out by D. Pesetsky, and confirmed by M. Saito. However, presumably the use of *wa* rather than *ga* and intonational pauses could be exploited as a surface cue to rule out more generally ambiguity in this example and others like it. See Fong and Berwick (1989) for a discussion of how to integrate surface cues into the principle-based system.

<sup>8</sup>This interpretation can be eliminated by imposing selectional restrictions on the possible "agents" of *okotteru* (let us say that they must be animate).

<sup>9</sup>Having a parsing system that can recover *all* such linguistic alternatives is of interest in its own right, both to verify and correct the linguistic theory, as well as ensure that no possibilities are overlooked by human interpreters.





## Optimization 2: Spec of CP on the right

A second obvious strategy is to remove the center-embedding itself. Here there is a grammatical move we can make. Evidently, in Japanese the only elements that appear in Spec of CP are put there by LF movement. Thus, these elements can never be visible in this position on the surface. If this is so, then there is really nothing to prevent us from placing just the Spec of CP on the right, rather than the left. This is an example of the "testbed" property of the system; this change takes two lines of code. Given this change, the resulting structures will have their Operators on the right, rather than the left, and will *not* be center-embedded. In addition, in this test the parser will *not* take advantage of right-hand information, thus eliminating this as a possible source of speedup.

Parsing complexity is reduced by this move, by a factor of just about one-half, if one considers either LR state transitions or principle operations; not as good as the first optimization; see below for some representative results. Also, with the most deeply center-embedded sentence the total number of principle operations actually is worse than in the base case. Evidently we have not located the source of the parser's problems in center-embedding alone.

Complexity for Spec on the right

Sentence	LR trans.	Total ops
ce1	122	32
ce2	4930	97
ce3	209,980	721
ce4	16,290,667	12605

## Optimization 3: Factoring out the effects of scrambling and *pro-drop*

While it appears that Head-final information helps the most, we must also remember that part of the complexity of Japanese is the result of free scrambling and *pro-drop*. To factor apart these effects, we ran a series of computer experiments on a quasi-Japanese grammar, J\*, that was just like Japanese except scrambling and *pro-drop* were barred. The changes were again simple to make: one change was automatic, just turning off a parameter value, while the second involved 3 lines of hand-coding in the X schemas to force the system to look for a lexical NP in DO (and IO) positions further, we did *not* optimize for right-hand information (so that the Head-final character was left intact). Of course, we now can no longer parse sentences with scrambled direct objects.

The table below shows the results. This was the *best* optimization of all. Without scrambling, and hence no movement at all compared to English, the Head-final quasi-Japanese was for the most part parsed 5-10 times more efficiently than English, and at worst (for the triply-embedded sentence) with three times fewer LR transitions and only about 30% more principle operations than English. Thus, this was even more efficient than the righthand information optimized Japanese parser. (The first column gives the number of LR transitions and the second gives the total number of principle operations for this "no scramble/drop" version, while the last two columns give the same information for English.)

No scrambling/drop vs. English

Sentence	LR trans.	No. ops	Eng. LR	Eng. ops
ce1	93	32	745	109
ce2	274	88	2431	168
ce3	1241	445	4979	558
ce4	8241	3719	21,074	2874

As before, with a short sentence, there is little difference between optimization methods, but over a range of sentences and with longer sentences, the no-scramble or *pro-drop* optimization works better than any other. Evidently, given the framework of assumptions we have made, the Head-final character of Japanese does not hurt the most; rather, it is scrambling and *pro-drop* that does, since if we remove these latter two effects we get the biggest improvement in parsing efficiency. We can confirm this by looking at the LR transitions for the other sentences (1b)-(18b) across methods, summarizing our tests. We can summarize the three experiments across sentences in figure 3.

Summary of complexity across tests

Sentence	Unopt.	Opt.	Spec-Final	No Scramble/drop
1b	1196	730	602	216
2b	1894	790	957	298
3b	345	289	185	103
4b	601	422	307*	149
5b	1742	1051	878*	370
6b	526	377	267	138
15b	32,955	19,998	11,205	1681
17b	1287	1789	685	272
18b	210,036	84,727	43,745	5306

## 6 Conclusions

Given our limited set of test sentences, our results must be correspondingly tentative. Nonetheless, we can draw several initial conclusions:

- One *can* parse Japanese by parametrically varying a grammar, much as expected. The limits of the method are theory-bound: we can accommodate just as much as we understand about Japanese syntax.
- Attempting to parse more than one language with the same grammar and parser can quickly reveal what is wrong with one's theory for either language. In our case, we discovered omissions in the implementation relating to Case transmission, the Wh-Comp Requirement, and trace deletion, among other items.
- A single parser suffices for very distinct languages. The grammar is parameterized, but not the parser, confirming much recent other research in Japanese sentence processing cited in the introduction. Japanese at first appears much more complex to parse than corresponding English sentences. We suggest, tentatively, that complexity is introduced by scrambling and omission of NPs, rather than Head-final properties. Unoptimized, the system is too slow. Some efficiency is obtained if one can "lift" information from the right for use in parsing with an LR machine. From a heuristic standpoint, this suggests that strategies limiting what may appear in a scrambled position or dropped in a certain context will aid such an LR-based device more than switching to a parser based presumably geared for a different branching direction.
- The principle-based system affords a new and generally straightforward way to precisely explore different grammatical theories, structural assumptions, and parsing methods and their computational consequences

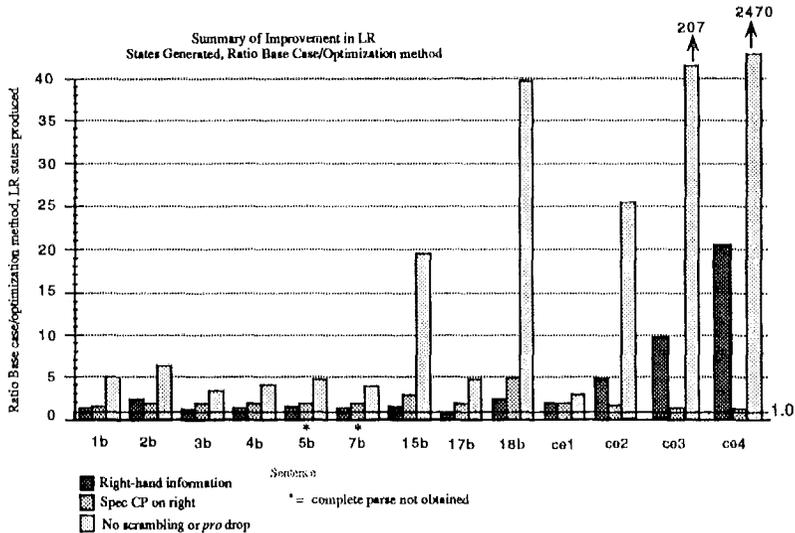


Figure 3: A bar graph showing the improvement in total LR transitions when parsing Japanese examples 1b-18b, and ce1-ce4, compared against the original base case unoptimized parser, across the 3 experiments described here. The horizontal line drawn at 1.0 indicates improvement over the base case.

in a precise way, without extensive hand coding. All of the experiments we tried took no more than a few lines of modification. Of course, the difficult part is to come up with a universal set of principles in the first place—so that in fact, English looks just about like Japanese, and vice-versa.

## 7 References

- Baltin, M.R., and A.S. Kroch (eds.), 1989. *Alternative Conceptions of Phrase Structure*, The University of Chicago Press.
- Chomsky, N.A., 1986. *Knowledge of Language: Its Nature, Origin, and Use*. Prager.
- Correa, N., 1988. Syntactic Analysis of English with respect to Government-binding Grammar, Ph.D. dissertation, Syracuse University.
- Crocker, M.W., 1989. *A Principle-Based System for Syntactic Analysis*, (m.s.).
- Dorr, B.J., 1987. UNITRAN: A Principle-Based Approach to Machine Translation. S.M. thesis. MIT Department of Electrical Engineering and Computer Science.
- Fong, S. & R.C. Berwick, 1989. The computational implementation of principle-based parsers. *International Workshop on Parsing Technologies*, Carnegie Mellon University, in M. Tomita (ed) *Current Issues in Parsing Technologies*, Kluwer.
- Fong, S., 1991. Computational Properties of Principle-based Grammatical Theories. Ph.D., dissertation, MIT Department of Computer Science and Electrical Engineering.
- Frazier, L., and K. Rayner, 1988. Parameterizing the language processing system: Left- vs. right-branching within and across languages. In J. Hawkins (ed.) *Explaining Language Universals*, Basil Blackwell, Oxford, pp. 247-279.
- Hosokawa, H. 1991. Syntactic differences between English and Japanese. *Georgetown Journal of Languages and Linguistics*, 1:4, 401-414.
- Inoue, N. and Fodor, J.D., 1991. Information-paced processing of Japanese sentences. Paper presented at the International Workshop on Japanese Syntactic Processing, Duke University.
- Johnson, M., 1989. Use of the Knowledge of Language, *Journal of Psycholinguistic Research*. 18(1).
- Kuno, S., 1973. *The Structure of the Japanese Language*, Cambridge, MA: MIT Press.
- Lasnik, H. & M. Saito, 1984. On the nature of proper government. *Linguistic Inquiry*, 15:2.
- Lasnik, H. & J. Uriagereka, 1988. *A Course in GB Syntax: Lectures on Binding and Empty Categories*, Cambridge, MA: MIT Press.
- Mazuka, R., 1990. Processing of empty categories in Japanese. Manuscript, Duke University.
- Nagai, N., 1991. Paper presented at the International Conference on Japanese Syntactic Processing, Duke University.
- Sharp, R.M., 1985. A Model of Grammar Based on Principles of Government and Binding Theory. M.S. thesis. Department of Computer Science. University of British Columbia.
- Stabler, E.P., Jr., 1991 forthcoming. *The Logical Approach to Syntax: Foundations, Specifications and Implementations of Theories of Government and Binding*, Cambridge, MA: MIT Press.
- Tomita, M., 1986. *Efficient Parsing for Natural Language: A Fast Algorithm for Practical Systems*. Kluwer.
- Wehrli, E., 1986. *A Government-Binding Parser for French*. Working Paper No. 48. University of Geneva.