

EFFICIENT PARSING FOR CCGS WITH GENERALIZED TYPE-RAISED CATEGORIES

Nobo Komagata *

Department of Computer and Information Science

University of Pennsylvania

komagata@linc.cis.upenn.edu

Abstract

A type of ‘non-traditional constituents’ motivates an extended class of Combinatory Categorical Grammars, CCGs with Generalized Type-Raised Categories (CCG-GTRC) involving variables. Although the class of standard CCGs is known to be polynomially parsable, unrestricted use of variables can destroy this essential requirement for a practical parser. This paper argues for polynomial parsability of CCG-GTRC from practical and theoretical points of view. First, we show that an experimental parser runs polynomially in practice on a realistic fragment of Japanese by eliminating spurious ambiguity and excluding genuine ambiguities. Then, we present a worst-case polynomial recognition algorithm for CCG-GTRC by extending the polynomial algorithm for the standard CCGs.

1 Introduction

In Japanese, as in other SOV languages, a sequence of NPs can form a conjunct as exemplified below.

- (1) John-ga Mary-o , Ken-ga Naomi-o tazuneta.
 { John-NOM May-ACC } CONJ { Ken-NOM Naomi-ACC } visited
 “John visited Mary and Ken [visited] Naomi.”

This type of ‘non-traditional constituents’ poses a problem to many grammar formalisms and parsers, including those specifically designed for Japanese, e.g., JPSG [Gunji, 1987] and JLE (based on finite-state syntax) [Kameyama, 1995]. Although these systems could be extended to cover the presented case, such extensions would not generalize to the wide range of non-traditional constituency.

Combinatory Categorical Grammar (CCG) has been proposed to account for non-traditional constituency in various areas of syntax [Ades and Steedman, 1982, Dowty, 1988, Steedman, 1985, Steedman, 1996] and also in the related areas of prosody, information structure, and quantifier scope [Prevost and Steedman, 1993, Hoffman, 1995, Park, 1995]. The mechanisms independently motivated to cover the wide range of non-traditional constituency can also provide an analysis for the NP-NP sequences in (1) as follows:

- (2) John-ga Mary-o
 NP NP
 \Downarrow \Downarrow type raising
 $S/(S \setminus NP)$ $(S \setminus NP)/((S \setminus NP) \setminus NP)$ functional composition (underlined portions are cancelled)
 $S/((S \setminus NP) \setminus NP)$

Informally, the NPs are assigned *higher-order function* categories associated with the basic category NP , and these functions can compose to derive another function category representing the NP-NP sequence. The two instances of such a category can then be coordinated and take the transitive verb category, $(S \setminus NP) \setminus NP$, as the argument to derive the category S .

A similar type of constituents can also be formed of NPs extracted from different levels of embedding as in the following example:

*I am grateful to Mark Steedman for his numerous suggestions and comments. I would also like to thank Jason Eisner, B. Srinivas, David Weir, K. Vijay-Shanker, and the anonymous reviewers for their comments. The research was supported in part by NSF Grant Nos. IRI95-04372, STC-SBR-8920230, ARPA Grant No. N66001-94-C6043, and ARO Grant No. DAAH04-94-G0426.

2 CCGs with Generalized Type-Raised Categories

CCG-GTRC involves the class of *constant categories* (Const) and the class of *Generalized Type-Raised Categories* (GTRC).

A constant (derivable) category c can always be represented as $F|a_n...|a_1$ where F is an atomic *target* category and a_i 's with their directionality are *arguments*. We will use “ A, \dots, Z ” for atomic, constant categories, “ a, \dots, z ” for possibly complex, constant categories, and ‘ $|$ ’ as a meta-variable for directional slashes $\{/, \backslash\}$. Categories are in the “result-leftmost” representation and associate left. Thus we usually write $F|a_n...|a_1$ for $(\dots (F|a_n) \dots |a_1)$. We will call “ $|a_i...|a_j$ ” a *sequence* (of arguments). The length of a sequence is defined as $\left| |a_i...|a_j \right| = i$ while the nil sequence is defined to have the length 0. Thus an atomic constant category is considered as a category with the target category with the nil sequence. We may also use the term ‘sequence’ to represent an ordered set of categories such as “ c_1, \dots, c_2 ” but these two uses can be distinguished by the context. The standard CCGs (CCG-Std) solely utilize the class of Const.

GTRC is a generalization of *Lexical Type-Raised Category* (LTRC) which has the form $\overline{T} \left\langle \left(T \right) a \right\rangle |b_i...|b_1$ associated with a lexical category $a|b_i...|b_1$ where \overline{T} is a variable over categories with the atomic target category T . The target indication may be dropped when it is not crucial or all the atomic categories are allowed for the target. We assume the order-preserving form of LTRC using the following notation. ‘ \langle ’ and ‘ \rangle ’ indicate that either set of slashes in the upper or the lower tier can be chosen but a mixture such as ‘ \langle ’ and ‘ $/$ ’ is prohibited.⁵ GTRC is defined as having the form of $\overline{T} \left\langle \left(T |a_m...|a_2 \right) a_1 \right\rangle |b_n...|b_1$ resulting from compositions of LTRCs where $m \geq 1, n \geq 0$, and the directional constraint is carried over from the involved LTRCs. When the directionality is not critical, we may simply write a GTRC as $\overline{T} | (T |a_m...|a_2 |a_1) |b_n...|b_1$. For $gtrc = \overline{T} | (T |a_m...|a_1) |b_n...|b_1$, we define $|gtrc| = n + 1$, ignoring the underspecified valency of the variable. Note that the introduction of LTRCs in the lexicon is non-recursive and thus does not suffer from the problem of the overgeneration discussed in [Carpenter, 1991].

These categories can be combined by combinatory rule schemata. Rules of (forward) “generalized functional composition” have the following form:

$$(5) \quad \begin{array}{ccc} x/y & \blacktriangleright^k & y|z_k...|z_1 \quad \longrightarrow \quad x|z_k...|z_1 \\ \text{functor category} & & \text{input category} \quad \quad \quad \text{result category} \end{array}$$

The integer ‘ k ’ in this schema is bounded by k_{max} , specific to the grammar, as in CCG-Std. Rules of functional application, “ $x/y \blacktriangleright^0 y \rightarrow x$ ”, can be considered as a special case of (5) where the sequence z_i 's is nil. The index k may be dropped when no confusion occurs. We say “ $x/y \blacktriangleright y|z_k...|z_1$ derives $x|z_k...|z_1$ ”, and “ $x|z_k...|z_1$ generates the string of nonterminals ‘ $x/y, y|z_k...|z_1$ ’ or the string of terminals ‘ ab ’ ” where the terminals a and b are associated with x/y and $y|z_k...|z_1$, respectively. The case with backward rules involving ‘ \blacktriangleleft ’ is analogous.

Inclusion of GTRCs calls for thorough examination of each combinatory case depending on the involved category classes. A summary is given in Table 1. The categorial representations for a few relevant cases will be shown below, but the details are left to the accompanying technical report [Komagata, 1997a].

The three cases indicated by ‘*’ in Table 1 introduce categories which are neither Const nor GTRC due to the residual variables. For example, observe Case 3c below.

(6) Combination type of “Const \blacktriangleright GTRC” where $k > |\text{input}|$ (and $k > 2$):

$$a/b \blacktriangleright^k \overbrace{T_0 | T_1 | (T_0 | T_1 | c_m \dots | c_1) | d_{k-2} \dots | d_1}^k \longrightarrow a | T_1 | (b | T_1 | c_m \dots | c_1) | d_{k-2} \dots | d_1$$

↑
residual

This is an unintended, accidental use of functional composition. The closure of the system must be maintained by excluding these cases.⁷ We are now in the position to define the class of CCG-GTRC:

Definition 1 A CCG-GTRC is a six tuple $(V_N, V_T, S, \mathcal{T}, f, R)$ where

- V_N is a finite set of nonterminals (atomic categories)
- V_T is a finite set of terminals (lexical items, written as a, \dots, z)

⁵For a related discussion, see [Steedman, 1991].

⁶ \overline{T} could also be decomposed into $T_0 | T_k \dots | T_1$ for a larger k but all these share the same characteristics with the above scheme.

⁷This condition is particularly important for implementation since the residual variables can behave beyond our imagination.

Case	Functor cat		Input cat		Result cat		
	Class	Outer seq	Class	$k \leq \text{input} $	Class	Residual variable	Unbounded const argument
1	Const	-	Const	\leq	Const	no	no
2a	GTRC	yes	Const	\leq	GTRC	no	no
2b	GTRC	no	Const	\leq	Const	no	no
3a	Const	-	GTRC	$<$	Const	no	no
3b	Const	-	GTRC	$=$	Const	no	possible
* 3c	Const	-	GTRC	$>$	neither	yes	possible
4a	GTRC	yes	GTRC	$<$	GTRC	no	no
4b	GTRC	yes	GTRC	$=$	GTRC	no	possible
* 4c	GTRC	yes	GTRC	$>$	neither	yes	possible
4di	GTRC	no	GTRC	$<$	GTRC	no	no
4dii	GTRC	no	GTRC	$<$	Const	no	no
4e	GTRC	no	GTRC	$=$	GTRC	no	no
* 4f	GTRC	no	GTRC	$>$	neither	yes	possible

Table 1: Combinatory Cases for CCG-GTRC

- S is a distinguished member of V_N
- \mathcal{T} is a countable set of variables⁸
- f is a function that maps elements of V_T to finite subsets of “Const \cup LTRC”
- R is a finite set of rule instances summarized in Table 1 (except for those with ‘*’).⁹

3 Progress Towards a Practical Parser for CCG-GTRC

This section investigates the performance of the experimental parser and demonstrates that it runs polynomially in practice. Since any practical parser for CCGs (or categorical grammars in general) requires some kind of spurious ambiguity elimination, we start with a discussion about spurious ambiguity elimination techniques.

3.1 Spurious Ambiguity Elimination

Let us first define different types of ambiguities referred to in this paper:

- (7) *a.* Categorical ambiguity: Availability of multiple categories (lexical/derivational), e.g., $\{NP, S \setminus NP\}$.
- b.* Spurious ambiguity: Multiple derivations of the same category which are semantically *equivalent* [Wittenburg, 1986], e.g., “John visited Bill.” has two derivations of S (left and right branching) in CCG with the identical semantics ‘(visited bill john)’.
- c.* Genuine ambiguity:
 - (i) Lexico-semantic ambiguity: Multiple semantic assignments to a single lexical category:
 - (ii) Attachment ambiguity: Multiple derivations of the same category with *distinct* semantics. E.g., PP attachment.

Note that by ‘semantics’, we assume predicate-argument structure (a kind of logical form) as the input to the semantic module.

Since spurious ambiguity is often accused of as the source of inefficient parsing, CCG parsers must implement some means of spurious ambiguity elimination. We review three classes of approaches: (i) syntactic, (ii) semantic, and (iii) those which do not belong to the previous two.

The syntactic approaches eliminate ‘spurious derivations’ which are not ‘the normal form’. Hendriks [1993] and König [1994] worked on Lambek calculus which does not have functional composition as a primitive rule. Hepple and Morrill

⁸Each instance of GTRC must be assigned a new variable when the GTRC is instantiated at a particular string position in order to avoid unintended variable binding.

⁹Due to the introduction of GTRC, the rule instances may involve variables even at the first argument of the functor category and at the input category.

[1989] cover a subset of the current formalism but do not have the mixed instances of function composition nor type raising. Eisner [1996] covers a wider range of CCGs but the case including type raising remains to be shown correct. Labeled deduction [Morrill, 1994] has a means to interpret semantics syntactically but normal form deduction must be adjusted.

Karttunen [1986] proposed the following semantic method. For each new derivation, discard it if its semantics is *equivalent to* (or mutually subsumes) that of some entry with the same category already in the cell.¹⁰ This directly enforces the definition of spurious ambiguity and does not depend on the syntax. Note that ‘equivalence’ depends on the semantic representation [Thompson, 1991]. For the case where the semantics is represented in λ -calculus, equivalence is not computable [Paulson, 1991]. For the case of feature structure, equivalence is defined as alphabetic variants and characterized by the isomorphism between the structures [Carpenter, 1992]. Our case corresponds to the latter although the equivalence check on predicate-argument structures are term unification rather than graph unification for the feature structure.

Among the third type, Pareschi and Steedman [1987] employed a control-based mechanism but the published algorithm has been shown to be incomplete [Hepple, 1987]. Wittenburg and Wall [1991] use a compilation scheme but the crucial notion of flexible constituency is compromised.

To be precise, the definition of spurious ambiguity for syntactic methods is not the same as our definition given above. Since syntactic methods are insensitive to semantics, they cannot distinguish genuine (attachment) and spurious ambiguities and eliminate both of these. With semantic equivalence check, we can adjust the way the genuine ambiguity is handled. While equivalence check in its original form is sensitive to genuine ambiguity, equivalence check applied only to category (ignoring semantics) is insensitive to genuine ambiguity much like syntactic methods. We refer to these two cases as (i) *category+semantics* and (ii) *category-only*. Both cases are explored in the experiment in the following subsection.

Among the presented methods, the only proposal immediately compatible with the present formalism is Karttunen’s semantic equivalence check. Let us review some arguments against this approach. Eisner [1996] argued that a sequence of categories exemplified by “ $X/X \dots X \dots X \backslash X$ ” can slow down the parser exponentially. Note that we assume that X/X and $X \backslash X$ are ‘modifiers’ of X with distinct semantics. But this is an instance of genuine ambiguity, a problem to any parser. Note that for syntactic methods, this appears as spurious ambiguity and the semantic processing is simply left for a later stage. Wittenburg [1987] objected to the cost of equivalence check. But an equivalence check (for our semantics) is inherently easier than the general case of subsumption check, the latter requiring the *occurs check* for soundness [Pereira and Shieber, 1987]. Hepple and Morrill [1989] have raised another objection. While syntactic methods detect spurious ambiguities before deriving a result, equivalence checking needs to compare the derived result with every entry in the current table cell. However, the cost associated with the semantic method depends on how many genuinely ambiguous entries are in the cell but not on the number of spuriously-ambiguous entries (they are eliminated as soon as they are derived and do not accumulate). This does not introduce additional complexity specific to spurious ambiguity check. Further, once semantics is involved, it is not possible to distinguish between spurious and genuine ambiguities unless we actually check the involved semantics.

3.2 Experiment

We now look at the results of a pilot experiment done on Sun Ultra E4000 2x167MHz Ultraspacs with 320MB memory running SunOS 5.5.1. The program (100KB approx., about a half is the grammar) was written in Sicstus Prolog Ver. 3 and CPU time was measured by Sicstus’ built-in predicate `statistics`.

We parsed 22 contiguous sentences (6 paragraphs) in Japanese in a section arbitrarily chosen from “Anata-no Byouin-no Kusuri (Your Hospital Drugs) 1996” by Tadami Kumazawa and Ko-ichi Ushiro. Japanese is chosen because the current work is a part of a larger project of medical database interface using this language. The romanized sentences are partially-segmented to the word level but the verb inflection and suffixes are considered as a part of the word. The average number of words in a sentence is 20 and the longest sentence contains 41 words. The sentences are realistically difficult, and include complex clauses (relative and complement), coordination (up to 4 conjuncts), nominal/verbal modifications (adjectives/adverbs), scrambling, and verb argument dropping.

The parser is based on a CKY algorithm [Aho and Ullman, 1972] equipped with Karttunen’s equivalence check for spurious ambiguity elimination but without the worst-case polynomial algorithm introduced in the next section.¹¹ LTRCs are assigned to words by lexical rules and GTRCs are restricted to unidirectional forms. Coordination is handled by

¹⁰Shieber [1986] contains a detailed discussion of subsumption.

¹¹Earlier applications of a CKY-style algorithm to CCG parsing include [Pareschi and Steedman, 1987].

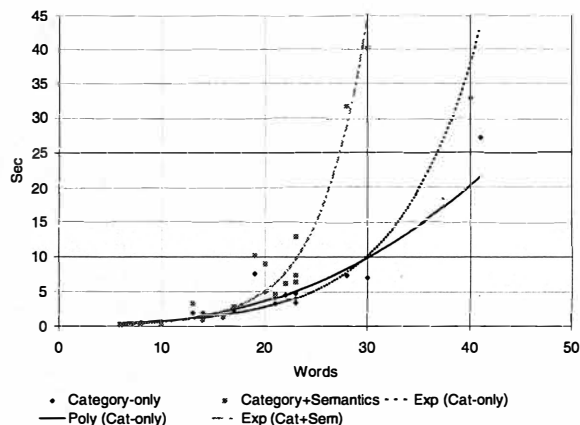


Figure 1: Basic Data Set (linear scale)

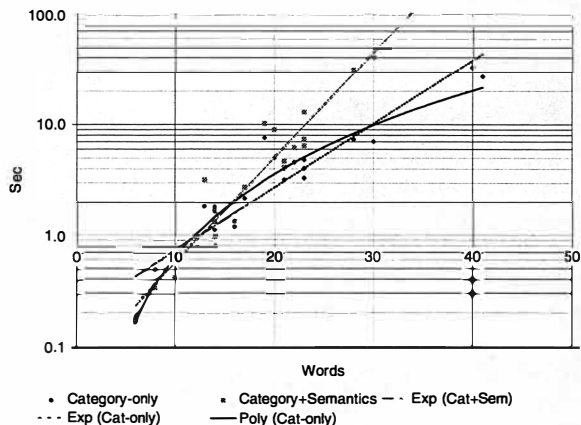


Figure 2: Basic Data Set (log scale)

special trinomial rules [Steedman, 1996] with a few categorial features added to limit the coordination involving multiple constituents only to the left-branching structure. Verb argument dropping is handled by lexical rules which change the verb valency. Morphological analysis is a complete substring match and the results are dynamically ‘asserted’ among the code. About 200 lexical entries are asserted after parsing the 22 sentences. Currently, morphological analysis takes about 0.2 seconds per word on average and needs improvement. The output of a parse is an enumeration of the final result categories associated with the features and the semantics as seen below.

```
(8) itumo 95mmHg o koeru baai_wa tiryou ga hituyoudesu.
    always num(95mmHg) -ACC exceed in_case treatment [-NOM,-CONJv] necessary
    Cat:
    SS: s
    Fs: []
    PA: (in_case always((exceed num(95mmHg) $1)) (necessary treatment))
    Cat:
    SS: s
    Fs: []
    PA: always((in_case (exceed num(95mmHg) $2) (necessary treatment)))
    CPU time: 280 ms Elapsed: 320 ms Words: 8 Solutions: 2
```

The unresolved pronoun is shown as ‘\$n’ where $n \in \mathbb{N}$. The ambiguity regarding adverbial modification is left unresolved. The implementation has a simplified treatment of quantifiers and scope ambiguity too is left unresolved.

We consider the following two cases: (i) category-only and (ii) category+semantics. As we have discussed in the previous subsection, the application of equivalence check to the category-only case not only eliminates spurious ambiguities but also provide a result without genuine ambiguities.

Let us start with the analysis of the category-only case.¹² This case corresponds to the situation involving syntactic methods and also the polynomial algorithms introduced in the next section. The results are shown in Figure 1 (linear scale) and 2 (log scale). Both exponential ($y = 0.1963 \times 1.14^n$) and polynomial ($y = 0.002 \times n^{2.496}$) regression lines calculated by Microsoft Excel are provided. Although it is often easier to fit either an exponential or a polynomial curve on a log-scale graph, the data do not seem to be enough for such a conclusion. To see how the experiment might extend to the case with words longer than 45 words, we parsed pseudo-long sentences. That is, some of the test sentences are conjoined to form long sentences. Although these are semi-fabricated data, most long sentences are in fact the results of coordination. Thus natural data are expected to behave similarly rather than differently from our pseudo-long sentences. The results are shown in Figures 3 and 4. The polynomial curves ($y = 0.0017 \times n^{2.5616}$) seem to represent the data better than the exponential curve ($y = 0.4375 \times 1.09^n$), especially on the log-scale graph. Since the data is sparse, we do not attempt to obtain a significant statistic analysis for these and simply eye-fit the data. With these qualifications, we conclude that the performance appears no worse than n^3 . The result also shows that categorial ambiguities still present in the parses

¹²Category-only is the case also corresponding to the spurious ambiguity check of syntactic methods.

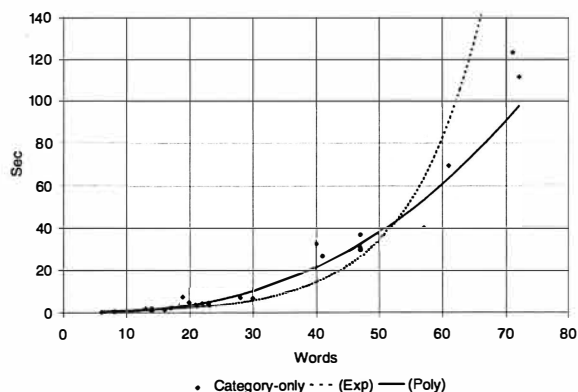


Figure 3: Extended Data Set (linear scale)

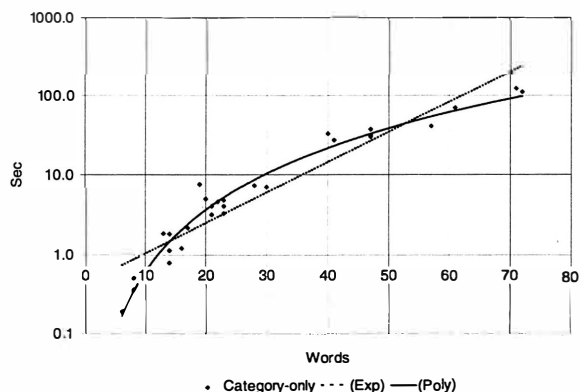


Figure 4: Extended Data Set (log scale)

are in practice within this bound.

A few remarks are in order. We compared our results with the following experiment to see how the figures stand. Tanaka and Ueki [1995] report that the LR-based syntactic analysis of a 19-word sentence in Japanese took 3.240sec.¹³ The range of CPU times for our sentences with 19-23 words is between 3 to 8sec (category-only case). The performance of our parser seems to be within a comparable range.

Another point is that the effect of spurious ambiguity check is immediate. Without the check, only the sentences with 10 or fewer words were parsed. Under this condition, the maximum number of cell entries easily exceeds 300 for longer sentences, which resulted in out-of-space errors. We thus confirmed that the exponential effect of spurious ambiguity is well controlled by semantic equivalence check.

Since one of the advantages of CCG parsers is the ability to derive semantics along with syntactic structure, the results of the category+semantics case is of special interest. The situation naturally looks quite different. The exponential regression line $t = 0.0638 \times 1.24^n$ (Figures 1 and 2) seems to fit the data closely. In fact, the two longest sentences with 40 and 41 words results in out-of-space errors. Since spurious ambiguities are eliminated by equivalence check and categorial ambiguity is only polynomial as shown in the category-only experiment, the exponential slow down is due exclusively to genuine ambiguity. Genuine ambiguity is a major problem for our parser as it is for any parser.

As a potential solution to genuine ambiguity problem, Dörre [1997] recently reported a polynomial algorithm to build shared semantic structures from shared syntactic structures. Unfortunately, this technique may not be directly applicable to our parser because semantic equivalence check is now required to traverse the shared structures. It is not clear if the traversal can be done in polynomial time. This concern is shared by the situation of applying structure sharing technique to conceptual dependency [Bröker et al., 1994].

4 Worst-Case Polynomial Recognition Algorithm

Although the implemented parser runs efficiently (the category-only case), its worst-case performance is still exponential due to categorial ambiguity. This section presents a worst-case polynomial recognition algorithm for a subclass of CCG-GTRC (Poly-GTRC) by extending the polynomial algorithm of Vijay-Shanker and Weir [1990] for CCG-Std (Poly-Std). We will observe below that the crucial property of CCG-Std employed by Poly-Std can be extended to the subclass of CCG-GTRC with an additional condition. Let us start with a brief review of the intuition behind Poly-Std and then move on to Poly-GTRC. Note that Poly-Std has the second stage of structure building but we concentrate on the more critical part of recognition.

¹³The word count is based on our criteria. Other details are ignored for now.

4.1 Polynomial Algorithm for CCG-Std

First, observe the following properties of CCG categories:

- (9) *a.* The length of a category in a cell can grow proportionally to the input size.
- b.* The number of categories in a cell may grow exponentially to the input size.

For example, consider a lexicon $f = \{(a, S/NP/S), (a, S/PP/S)\}$. Then, for the input “ $\overbrace{a \dots a}^n$ ”, the top CKY-cell includes

2^n combinations of categories like $S \underbrace{\left\{ \begin{matrix} /NP \\ /PP \end{matrix} \right\} \dots \left\{ \begin{matrix} /NP \\ /PP \end{matrix} \right\}}_n /S$ derived by functional composition. Thus, we have exponential worst-case performance with respect to the input size.

The idea behind Poly-Std is to store categories as if they were some kind of linked list. Informally, a long category $F|a_n \dots |a_2|a_1$ is stored as ‘ F this portion is *linked* in a cell $[p, q]$ with *index* $|a_2|$ $|a_1$ ’. A crucial point here is that the instances of target category F and arguments a_2 and a_1 are *bounded*. We will come back to this point in the next subsection. The pair $[p, q]$ can be represented as a n^2 matrix. Thus by setting up n^2 subcells in each CKY-table cell, we can represent a category in a finite manner.

The effectiveness of such a representation comes from the fact that CCG rule application does not depend on the entire category. Namely, in order to verify “ $F \underbrace{|a_n \dots |a_2|}_{\star} |a_1| \blacktriangleright^k \underbrace{b_0|b_k \dots |b_1|}_{\star} \rightarrow F |a_n \dots |a_2| \underbrace{|b_k \dots |b_1|}_{\star}$ ”, the sequence marked by ‘ \star ’

does not need to be examined. Thus, for the functor category, we only need to check F and a_1 available in the current cell. In addition, since b_0 must be unified with a bounded a_1 , and k is also bounded by k_{max} , the entire input category is bounded and thus can be stored in the current cell. Therefore, the proposed representation does not slow down this type of process. When the result category exceeds a certain limit, we leave the excessive portion right in the original cell and set up a link to it.

One complication is that when an argument (e.g., a_1 in the above example) is canceled, we may have to restore a portion of the category from the linked cell (as the ‘index’ for the cell is required). We need to scan the linked cells and find the categories with the same index from n^2 subcells. Even though there may be multiple such categories, all of them can be restored in one of n^2 subcells associated with the result category. This case dominates the computational complexity but can be done in $O(n^3)$. Since this is inside i, j, k of CKY-style loop, the overall complexity is $O(n^6)$.

4.2 Polynomial Algorithm for CCG-GTRC

We first note that there are cases where a crucial property of CCG-Std cannot be maintained in CCG-GTRC. The property is that arguments of derived categories are bounded. Although there might be a polynomial algorithm for CCG-GTRC which does not depend on this property, we pursue a straightforward extension of Poly-Std with an additional condition on the rules.¹⁴

The problematic cases are *3b* and *4b* in Table 1. As we have motivated in Introduction, the inner sequence of GTRC must be allowed to grow unboundedly long (otherwise, the resulting grammar becomes equivalent to CCG-Std). Then an argument can grow unboundedly long as exemplified by Case *3b* below.

- (10) Combination type of “Const \blacktriangleright GTRC” where $k = |\text{input}|$ (and $k \geq 1$):

$$a/\underline{b} \blacktriangleright^k \underline{\text{T}} \overbrace{\text{T} |c_m \dots |c_1|}^k |d_{k-1} \dots |d_1 \longrightarrow a|(b |c_m \dots |c_1|) |d_{k-1} \dots |d_1$$

$\underbrace{\hspace{10em}}_{\text{unbounded}}$

But we do not want to exclude Cases *3b* and *4b* entirely since a case involving a sentential adverb may be possible as follows: “ $S/S \blacktriangleright \text{T} |(\text{T}|A|B)$ ”. The present approach is to place the following condition on the rule application:

- (11) **Bounded Argument Condition:** The input GTRCs in the cases *3b* and *4b* must be instantiated.

This is in a sense analogous to placing a bound k_{max} on functional composition instantiating the input category. By this condition, we have the following property:

- (12) The set of arguments of constant categories and the set of arguments of the inner and outer sequences of GTRC are all finite.

¹⁴The length of the argument is still bounded by $O(n)$ in CCG-GTRC since the only source of unboundedness is GTRC inner sequences. If every argument can be represented in some finite manner with link information similar to the one used for the Poly-Std, polynomial recognition might be possible.

That is, by considering the inner sequence at the same level as the outer sequence of a GTRC or the arguments of a constant category, we can maintain the property analogous to the one found for CCG-Std. Note that the Bound Argument Condition is not implemented in our parser. The inner sequence of GTRC does not grow unboundedly in practice and thus the arguments of categories do not grow unboundedly either.

Now, consider the subclass of CCG-GTRC constrained by the Bounded Argument Condition. In the rest of this section, we will concentrate on this subclass.

Poly-GTRC is an extension to Poly-Std. The basic organization of the algorithm is analogous to Poly-Std. We use the same $n^2 \times n^2$ CKY-style table and a similar representation for constant categories. But we need to deal with GTRCs in polynomial time as well. First, let us examine two representative cases of rule applications since this reveals the necessary conditions for polynomial parsing.

The inner sequence of GTRC can grow as a result of Case 4e shown below:

$$(13) \quad \begin{array}{c} \uparrow \\ \text{T} / \underline{(\text{T}|a_m \dots |a_2 \backslash a_1)} \blacktriangleright^k \underline{\text{U}} | (\text{U}|c_p \dots |c_1) | d_{k-1} \dots | d_1 \longrightarrow \text{T} | (\text{T}|a_m \dots |a_1 | c_p \dots |c_1) | d_{k-1} \dots | d_1 \quad (k \geq 1) \end{array}$$

The only information needed to determine if the rule is applicable is the directionality of the slash indicated by ‘ \uparrow ’. Thus we do not actually need to know the inner sequence of the functor or input categories. The inner sequence of the result GTRC can thus be represented as two links to the functor and input categories. This link information virtually encodes a kind of grammar for deriving the inner sequence and is thus considered as an application of structure sharing [Billot and Lang, 1989, Dymetman, 1997]. The outer sequence can be represented in a way similar to the argument of constant category. Although there may be exponentially-many GTRCs associated with each CKY cell, the number of cell entries is bounded by the link destinations of the inner sequence and the finite representation for the outer sequence.

Next, consider Case 2b.

$$(14) \quad \text{T} / \underline{(\text{T}|a_m \dots |a_2 \backslash a_1)} \blacktriangleright^k \begin{array}{c} c \\ \parallel \\ c_0 | c_m \dots | c_1 \end{array} | d_k \dots | d_1 \longrightarrow c_0 | d_k \dots | d_1$$

We need to show that the unification process of the underlined portions can be done in polynomial time. As the first approximation, consider this process as an iteration of (backward) functional application of the form “ $a_i \blacktriangleleft c_0 | c_m \dots | c_i$ ” for $i = 1$ to m where a_i and c_i are canceled. But recall that in general, we only store a finite portion of both the functor and the input categories in the current cell and the remaining information must be restored through the links. The restoration of the information could cost exponential time since there may be multiple links to lower locations at any point. Therefore it is crucial that we proceed from $i = 1$ to m so that no enumeration of all the instances of c_i, \dots, c_1 and a_i, \dots, a_1 in (14) is actually generated. The traversal of the link from c_1 and a_1 may introduce sets of categories C_i and \mathcal{A}_i for each position of $i \geq 2$, as schematically shown below.

$$(15) \quad \begin{array}{ccc} C_m & \dots & C_2 \\ \updownarrow & & \updownarrow \\ \mathcal{A}_m & \dots & \mathcal{A}_2 \end{array} \longleftarrow \begin{array}{c} \{c_1\} \\ \updownarrow \\ \{a_1\} \end{array}$$

Note that each set C_i and \mathcal{A}_i are bounded. This is the crucial point we needed the Bounded Argument Condition. Now, suppose that an element in C_i is canceled with some elements in \mathcal{A}_i . We can proceed to the next set C_{i+1} where the elements in C_{i+1} are obtained by traversing the links from the canceled elements in C_i . Once we move from C_i to C_{i+1} , the history of cancellation can be forgotten, as in the case of iterative functional application in Poly-Std. Thus even though we have potentially exponential instances of c_i, \dots, c_1 , the traversal of this side can be done step-by-step without suffering the exponential effect.

The traversal of \mathcal{A}_i ’s is more challenging. The availability of a_i for cancellation with some c_i depends on the history of the cancellation of a_{i-1}, \dots, a_1 . Actually, it depends on the *tree structure* exactly encoded by the structure sharing technique. The ‘GTRC recovery algorithm’ will be introduced below to handle this situation in polynomial time.

The other cases are included in the technical report along with a more detailed description and a worked example of recovering GTRCs [Komagata, 1997a]. Through the examination, we conclude that the polynomial parsability of CCG-GTRC depends on recovery of GTRCs. We present the polynomial GTRC recovery algorithm in Figure 5.

The GTRC recovery algorithm takes advantage of the encoded shared structure, and utilizes an additional n^2 GTRC recovery table to restore possibly ambiguous GTRC derivations in polynomial time. The first stage (*table setup*) is to represent the derivational structure available in the CKY table in a slightly different way. Suppose the following partial CKY table starting from a GTRC in question:

Initialization:

- Create an n^2 GTRC recovery table, R

Table setup (Stage 1):

- For each cell (top-down) $O(n^2)$
 - For each entry (depending on the midpoint) $O(n)$
 - Restore the derivation info from CKY table and store the children in the appropriate cells $O(n^2)$

Recovery (Stage 2):

- For each cell in the bottom row (right-to-left) $O(n)$
 - For each entry $O(n)$
 - If there is a matching category in the target category set
 - Mark the current entry as ‘success’
 - Otherwise
 - Mark the current entry as ‘fail’
 - Do **status percolation**

Status percolation (subprocedure):

- For each cell (bottom-up) $O(n^2)$
 - For each entry $O(n)$
 - For each parent $O(n^2)$
 - If the parent is marked as ‘fail’
 - Mark the current entry as ‘fail’
 - Otherwise
 - If the current entry is the right branch *and* marked as ‘fail’ *and* all the right branch siblings are marked as ‘fail’,
 - Mark the parent as ‘fail’
 - If the current entry is the left branch *and* marked as ‘success’
 - Mark the parent as ‘success’

Figure 5: GTRC Recovery Algorithm

(16) Partial CKY table:

5	$T/(T \dots B)_{[1,2] \blacktriangleright [3,5]}$				
4					
3			$T/(T \dots B)_{[3,3] \blacktriangleleft [4,5]}$		
2	$T/(T \setminus A)_{[1,1] \blacktriangleright [2,2]}$			$T \setminus (T/C)_{[4,4] \blacktriangleleft [5,5]}$	
1	$T/(T \setminus A)/D$	D	$T/(T \setminus B)$	E	$T \setminus (T/C) \setminus E$
	1	2	3	4	5

$T/(T \dots B)_{[1,2] \blacktriangleright [3,5]}$ represents the derivation “ $T/(T \setminus A)_{[1,1] \blacktriangleright [2,2]} \blacktriangleright T/(T \dots B)_{[3,3] \blacktriangleleft [4,5]} \rightarrow T/(T \dots B)$ ” at the designated string positions. A GTRC recovery table can be used to store the same derivational structure with the bottom row corresponding to *the order of the inner sequence* of the GTRC rather than the string position. This is the order to process the inner sequence for $C_i \text{-} \mathcal{A}_i$ comparison shown in (15). Since GTRC recovery process only concerns with the inner sequence of the GTRCs, the recovery table may have a dimension smaller than the corresponding portion of the CKY table as seen in the following example:

(17) GTRC recovery table:

3	$T/(T \dots B)_{[1,1] \blacktriangleright [2,3]}$		
2		$T/(T \dots B)_{[2,2] \blacktriangleright [3,3]}$	
1	$T/(T \setminus A)$	$T \setminus (T/C)$	$T/(T \setminus B)$
	3	2	1

The categorial ambiguities originally aligned at string positions are now aligned in the order of processing.

In the second stage (*recovery stage*), the comparison with the target categories is done while the above-mentioned dependency among LTRCs in the bottom row is checked. The comparison proceeds from right to left in the bottom row. The decision on the cancellation of the argument under consideration, a_i , depends on (i) if it is unifiable with some target category (in C_i) *and* (ii) if the corresponding sequence to the right of a_i was successfully canceled. This latter condition

can be checked by observing the status of the first right branch from the current position since all the processes up to that point must have been completed. For the later processing (for the positions to the left), the success/failure status of the current category must also be percolated to the relevant higher nodes (*status percolation*). The total complexity turns out to be a rather daunting $O(n^{10}) = O(\underbrace{n^3}_{i,j,k} \times \underbrace{n^2}_{\text{recovery}} \times \underbrace{n^5}_{\text{percolation}})$.

The extremely-high cost of GTRC recovery process seems to be related to the following conjecture: the generative power of this subclass of CCG-GTRC is greater than CCG-Std. A more restricted version of CCG-GTRC has been shown to be weakly equivalent to CCG-Std [Komagata, 1997b, Komagata, 1997c]. This subclass restricts the directionality of GTRC to the unidirectional case and only deals with GTRCs without outer sequence, thus allowing only the GTRCs of the form $T \langle (T) a_m \dots \rangle a_1$. In this case, the inner sequence and the string positions agree on the number and the order of the arguments and no GTRC recovery algorithm is needed. Realistically, most languages seem to restrict the use of GTRCs in some sense. For example, the Japanese grammar used in our experiment restricts GTRCs to the unidirectional variety but still with outer sequence.

5 Conclusion

The practical and the theoretical polynomial results are due to distinct factors. The former comes from a practical bound on the number of cell entries after spurious ambiguity elimination, showing CFG-like behavior. The latter (for both Poly-Std and Poly-GTRC) is achieved by efficiently representing and processing the potentially exponentially-many entries in a cell. This is possible even with the presence of spurious/genuine ambiguities. But it turns out that what the polynomial algorithms do is to eliminate a possibility which in practice rarely occurs. The additional cost for Poly-Std/GTRC of managing n^2 subcells and links to cover all the cases of exponential factors including spurious/genuine ambiguities is thus considered as overkill for the practical case. Although it may be possible to add spurious ambiguity check to Poly-Std/GTRC, we are better off with a simple CKY-style parser with equivalence check, without the overhead of the Poly-Std/GTRC.

The above conclusion naturally remains qualified by the small scale of the experiment reported here. But, the test sentences are reasonably representative and relatively challenging. They vary in sentence length and complexity and span the space we may typically encounter. With additional data, it is reasonable to expect that the missing points will be filled and statistic significance will be obtained. It is also reasonable to believe that the experiment with pseudo-long sentences characterizes the kind of complexity that will be found in natural data.

Few grammars or parsers designed for a configurational language such as English perform well on realistic fragments of a language like Japanese. One of the reasons may be abundance of non-traditional constituents found in SOV languages. Our parser, in contrast, demonstrates a reasonable performance although the grammar is basically an instantiation of the same framework applied for English [Steedman, 1996]. This strongly suggests that a framework which can analyze non-traditional constituents as a part of the competence property has an advantage in dealing with cross-linguistic data.

References

- Anthony Ades and Mark J. Steedman. 1982. On the Order of Words. *Linguistics and Philosophy*, 4.
- Alfred V. Aho and J. D. Ullman. 1972. *The Theory of Parsing, Translation, and Compiling. Vol. 1: Parsing*. Prentice-Hall.
- Tilman Becker, Aravind Joshi, and Owen Rambow. 1991. Long-Distance Scrambling and Tree Adjoining Grammars. In *EACL5*, 1991.
- Sylvie Billot and Bernard Lang. 1989. The Structure of Shared Forests in Ambiguous Parsing. In *ACL27*, 1989.
- Norbert Bröker, Udo Hahn, and Susanne Schacht. 1994. Concurrent Lexicalized Dependency Parsing: The ParseTalk Model. In *COLING-94*, 1994.
- Bob Carpenter. 1991. The Generative Power of Categorical Grammars and Head-Driven Phrase Structure Grammars with Lexical Rules. *Computational Linguistics*, 17.
- Robert L. Carpenter. 1992. *The Logic of Typed Feature Structures*. Cambridge University Press.
- Jochen Dörre. 1997. Efficient Construction of Underspecified Semantics under Massive Ambiguity. In *ACL35/EACL8*, 1997.
- David Dowty. 1988. Type Raising, Functional Composition, and Non-Constituent Conjunction. In Richard Oehrle, Emmon Bach, and Deirdre Wheeler, editors, *Categorical Grammars and Natural Language Structures*. D. Reidel.

- Marc Dymetman. 1997. Charts, Interaction-Free Grammars, and the Compact Representation of Ambiguity. In *IJCAI-97*, 1997.
- Jason Eisner. 1996. Efficient Normal-Form Parsing for Combinatory Categorical Grammar. In *ACL 34*, 1996.
- Takao Gunji. 1987. *Japanese Phrase Structure Grammar: A Unification-Based Approach*. D. Reidel.
- Herman Hendriks. 1993. *Studied Flexibility: Categories and Types in Syntax and Semantics*. ILLC Dissertation Series.
- Mark Hepple and Glyn Morrill. 1989. Parsing and Derivational Equivalence. In *EACL 4*, 1989.
- Mark Hepple. 1987. Methods for Parsing Combinatory Grammars and the Spurious Ambiguity Problem.
- Beryl Hoffman. 1993. The Formal Consequences of Using Variables in CCG Categories. In *ACL31*, 1993.
- Beryl Hoffman. 1995. *The Computational Analysis of the Syntax and Interpretation of "Free" Word Order in Turkish*. PhD thesis, University of Pennsylvania.
- Aravind K. Joshi, Tilman Becker, and Owen Rambow. 1994. Complexity of Scrambling: A New Twist to the Competence - Performance Distinction. In *3e Colloque International sur les grammaires d' Arbres Adjoints*, 1994.
- Megumi Kameyama. 1995. The Syntax and Semantics of the Japanese Language Engine. In Reiko Mazuka and Noriko Nagai, editors, *Japanese Sentence Processing*. Lawrence Erlbaum.
- Lauri Karttunen. 1986. Radical Lexicalism. Technical report, CSLI.
- Nobo Komagata. 1997a. Efficient Parsing for CCGs with Generalized Type-Raised Categories. Technical report, University of Pennsylvania.
- Nobo Komagata. 1997b. Generative Power of CCGs with Generalized Type-Raised Categories. In *ACL35/EACL8 (Student Session)*, 1997b.
- Nobo Komagata. 1997c. Generative Power of CCGs with Generalized Type-Raised Categories. Technical report, University of Pennsylvania.
- Esther König. 1994. A Hypothetical Reasoning Algorithm for Linguistic Analysis. *J. Logic and Computation*, 4(1):1-19.
- George Miller and Noam Chomsky. 1963. Finitary Models of Language Users. In R.R. Luce et al., editors, *Handbook of Mathematical Psychology, Vol. II*. John Wiley and Sons.
- Glyn V. Morrill. 1994. *Type logical grammar: categorial logic of signs*. Kluwer.
- Remo Pareschi and Mark Steedman. 1987. A Lazy Way to Chart-Parse with Categorical Grammars. In *ACL25*, 1987.
- Jong C. Park. 1995. Quantifier Scopepe and Constituency. In *ACL33*, 1995.
- Lawrence C. Paulson. 1991. *ML for the Working Programmer*. Cambridge University Press.
- Fernando C.N. Pereira and Stuart M. Shieber. 1987. *Prolog and Natural-Language Analysis*. CSLI.
- Scott Prevost and Mark Steedman. 1993. Generating Contextually Appropriate Intonation. In *EACL6*, 1993.
- Stuart M. Shieber. 1986. *An Introduction to Unification-Based Approaches to Grammar*. CSLI.
- Mark J. Steedman. 1985. Dependency and Coordination in the Grammar of Dutch and English. *Language*, 61:523-56.
- Mark Steedman. 1991. Type-Raising and Directionality in Combinatory Grammar. In *ACL29*, 1991.
- Mark Steedman. 1996. *Surface Structure and Interpretation*. MIT Press.
- Hozumi Tanaka and Masahiro Ueki. 1995. Japanese Parsing System using EDR Dictionary. <http://www.icot.or.jp/AITEC/PUBLICATIONS/Itaku/95/catalogue18-E.html>.
- Simon Thompson. 1991. *Type Theory and Functional Programming*. Addison-Wesley.
- K. Vijay-Shanker and David J. Weir. 1990. Polynomial Time Parsing of Combinatory Categorical Grammars. In *ACL28*, 1990.
- K. Vijay-Shanker and David J. Weir. 1991. Polynomial Parsing of Extensions of Context-Free Grammars. In Masaru Tomita, editor, *Current Issues in Parsing Technology*. Kluwer.
- K. Vijay-Shanker and David J. Weir. 1993. Parsing Constrained Grammar Formalisms. *Computational Linguistics*, 19(4).
- K. Vijay-Shanker and D. J. Weir. 1994. The Equivalence of Four Extensions of Context-Free Grammars. *Mathematical Systems Theory*, 27:511-546.
- David Weir. 1988. *Characterizing Mildly Context-Sensitive Grammar Formalisms*. PhD thesis, University of Pennsylvania.
- Kent Wittenburg and Robert E. Wall. 1991. Parsing with Categorical Grammar in Predictive Normal Form. In Masaru Tomita, editor, *Current Issues in Parsing Technology*. Kluwer.
- Kent Barrows Wittenburg. 1986. *Natural Language Parsing with Combinatory Categorical Grammar in a Graph-Unification-Based Formalism*. PhD thesis, University of Texas, at Austin.
- Kent Wittenburg. 1987. Predictive Combinators: A Method for Efficient Processing of Combinatory Categorical Grammars. In *ACL25*, 1987.