Analysis Grammar of Japanese in the Mu-Project

- A Procedural Approach to Analysis Grammar -

Jun-ichi TSUJII, Jun-ichi NAKAMURA and Makoto NAGAO

Department of Electrical Engineering
Kyoto University
Kyoto, JAPAN

Abstract

Analysis grammar of Japanese in the Mu-project
is presented. It is emphasized that rules
expressing constraints on single linguistic
structures and rules for selecting the most
preferable readings are completely different in
nature, and that rules for selecting preferale
readings should be utilized in analysis grammars of
practical MT systems. It is also claimed that
procedural control is essential in integrating such
rules into a unified grammar. Some sample rules
are given to make the points of discussion clear
and concrete.

1. Introduction

The Mu-Project is a Japanese national project
supported by grants from the Special Coordination
Funds for Promoting Science & Technology of
STA(Science and Technology Agency), which aims to
develop Japanese-English and English-Japanese
machine translation systems. We currently restrict
the domain of translation to abstracts of
scientific and technological papers. The systems
are based on the transfer approach[1], and consist
of three phases: analysis, transfer and generation.
In this paper, we focus on the analysis grammar of
Japanese in the Japanese-English system. The
grammar has been developed by using GRADE which is
a programming language specially designed for this
project[2]. The grammar now consists of about 900
GRADE rules. The experiments so far show that the
grammar works very well and is comprehensive enough
to treat various linguistic phenomena in abstracts.
In this paper we will discuss some of the basic
design principles of the grammar together with its
detailed construction. Some examples of grammar
rules and analysis results will be shown to make
the points of our discussion clear and concrete.


2. Procedural Grammar

There has been a prominent tendency in recent
computational linguistics to re-evaluate CFG and
use it directly or augment it to analyze
sentences[3,4,5]. In these systems(frameworks),
CFG rules independently describe constraints on
single linguistic structures, and a universal rule
application mechanism automatically produces a set
of possible structures which satisfy the given
constraints. It is well-known, however, that such
sets of possible structures often become
unmanageably large.

Because two separate rules such as

NP -----> NP PREP-P
VP -----> VP PREP-P

are usually prepared in CFG grammars in order to
analyze noun and verb phrases modified by
prepositional phrases, CFG grammars provide two
syntactic analyses for

She was given flowers by her uncle.

Furthermore, the ambiguity of the sentence is
doubled by the lexical ambiguity of "by", which can
be read as either a locative or an agentive
preposition. Since the two syntactic structures
are recognized by completely independent rules and
the semantic interpretations of "by" are given by
independent processes in the later stages,it is
difficult to compare these four readings during the
analysis to give a preference to one of these four
readings.

A rule such as

"If a sentence is passive and there is a
"by"-prepositional phrase, it is often the case
that the prepositional phrase fills the deep
agentive case. (try this analysis first)"

seems reasonable and quite useful for choosing the
most preferable interpretation, but it cannot be
expressed by refining the ordinary CFG rules. This
kind of rule is quite different in nature from a
CFG rule. It is not a rule of constraint on a
single linguistic structure(in fact, the above four
readings are all linguistically possible), but it
is a "heuristic" rule concerned with preference of
readings, which compares several alternative
analysis paths and chooses the most feasible one.
Human translaters (or humans in general) have many

such preference rules based on various sorts of cue such as morphological forms of words, collocations of words, text styles, word semantics, etc. These heuristic rules are quite useful not only for increasing efficiency but also for preventing proliferation of analysis results. As Wilks[6] pointed out, we cannot use semantic information as constraints on single linguistic structures, but just as preference cues to choose the most feasible interpretations among linguistically possible interpretations. We claim that many sorts of preference cues other than semantic ones exist in real texts which cannot be captured by CFG rules. We will show in this paper that, by utilizing various sorts of preference cues, our analysis grammar of Japanese can work almost deterministically to give the most preferable interpretation as the first output, without any extensive semantic processing (note that even "semantic" processing cannot disambiguate the above sentence. The four readings are semantically possible. It requires deep understanding of contexts or situations, which we cannot expect in a practical MT system).

In order to integrate heuristic rules based on various levels of cues into a unified analysis grammar, we have developed a programming langauage, GRADE. GRADE provides us with the following facilities.

- Explicit Control of Rule Applications : Heuristic rules can be ordered according to their strength(See 4-2).

- Multiple Relation Representation : Various levels of information including morphological, syntactic, semantic, logical etc. are expressed in a single annotated tree and can be manipulated at any time during the analysis. This is required not only because many heuristic rules are based on heterogeneous levels of cues, but also because the analysis grammar should perform semantic/logical interpretation of sentences at the same time and the rules for these phases should be written in the same framework as syntactic analysis rules (See 4-2, 4-4).

- Lexicon Driven Processing : We can write heuristic rules specific to a single or a limited number of words such as rules concerned with collocations among words. These rules are strong in the sense that they almost always succeed. They are stored in the lexicon and invoked at appropriate times during the analysis without decreasing efficiency (See 4-1).

- Explicit Definition of Analysis Strategies : The whole analysis phase can be divided into steps. This makes the whole grammar efficient, natural and easy to read. Furthermore, strategic consideration plays an essential role in preventing undesirable interpretations from being generated (See 4-3).

3 Organization of Grammar

In this section, we will give the organization of the grammar necessary for understanding the discussion in the following sections. The main components of the grammar are as follows.

(1) Post-Morphological Analysis
(2) Determination of Scopes
(3) Analysis of Simple Noun Phrases
(4) Analysis of Simple Sentences
(5) Analysis of Embedded Sentences (Relative Clauses)
(6) Analysis of Relationships of Sentences
(7) Analysis of Outer Cases
(8) Contextual Processing (Processing of Omitted case elements, Interpretation of 'Ha' , etc.)
(9) Reduction of Structures for Transfer Phase

Each component consists of from 60 to 120 GRADE rules.

47 morpho-syntactic categories are provided for Japanese analysis, each of which has its own lexical description format. 12,000 lexical entries have already been prepared according to the formats. In this classification, Japanese nouns are categorized into 8 sub-classes according to their morpho-syntactic behaviour, and 53 semantic markers are used to characterize their semantic behaviour. Each verb has a set of case frame descriptions (CFD) which correspond to different usages of the verb. A CFD gives mapping rules between surface case markers (SCM - postpositional case particles are used as SCM's in Japanese) and their deep case interpretations (DCI - 33 deep cases are used). DCI of an SCM often depends on verbs so that the mapping rules are given to CFD's of individual verbs. A CFD also gives a normal collocation between the verb and SCM's(postpositonal case particles). Detailed lexical descriptions are given and discussed in another paper[7].

The analysis results are dependency trees which show the semantic relationships among input words.
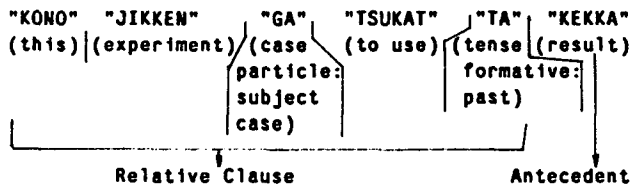
4. Typical Steps of Analysis Grammar

In the following, we will take some sample rules to illustrate our points of discussion.

4-1 Relative Clauses

Relative clause constructions in Japanese express several different relationships between modifying clauses (relative clauses) and their antecedents. Some relative clause constructions

cannot be translated as relative clauses in English. We classified Japanese relative clauses into the following four types, according to the relationships between clauses and their antecedents.

(1) Type 1 : Gaps in Cases

One of the case elements of the relative clause is deleted and the antecedent fills the gap.

(2) Type 2 : Gaps in Case Elements

The antecedent modifies a case element in the clause. That is, a gap exists in a noun phrase in the clause.

(3) Type 3 : Apposition

The clause describes the content of the antecedent as the English "that"-clause in 'the idea that the earth is round'.

(4) Type 4 : Partial Apposition

The antecedent and the clause are related by certain semantic/pragmatic relationships. The relative clause of this type doesn't have any gaps. This type cannot be translated directly into English relative clauses. We have to interpolate in English appropriate phrases or clauses which are implicit in Japanese, in order to express the semantic/pragmatic relationships between the antecedents and relative clauses explicitly. In other words, gaps exist in the interpolated phrases or clauses.

Because the above four types of relative clauses have the same surface forms in Japanese

```
    --------- (verb ) (noun).
   |_____|         |
   Relative Clause   Antecedent
```

careful processing is required to distinguish them (note that the 'antecedents' -modified nouns- are located after the relative clauses in Japanese). A sophisticated analysis procedure has already been developed, which fully utilizes various levels of heuristic cues as follows.

(Rule 1) There are a limited number of nouns which are often used as antecedents of Type 3 clauses.

(Rule 2) When nouns with certain semantic markers appear in the relative clauses and those nouns are followed by one of specific postpositional case particles, there is a high possibility that the relative clauses are Type 2. In the following example, the word "SHORISOKUDO"(processing speed) has the semantic marker AO (attribute).

[ex-1] [Type 2]

"SHORISOKUDO"        "GA"      "HAYAI" "KEISANKI"
(processing speed) (case     (high)  (computer)
                   particle:
                   subject
                   case)

Relative Clause              Antecedent

-->(English Translation)
    A computer whose processing speed is high

(Rule 3) Nouns such as "MOKUTEKI"(purpose), "GEN_IN"(reason), "SHUDAN"(method) etc. express deep case relationships by themselves, and, when these nouns appear as antecedents, it is often the case that they fill the gaps of the corresponding deep cases in the relative clauses.

[ex-2] [Type 1]

"KONO" "SOUCHI" "O"    "TSUKAT" "TA"      "MOKUTEKI"
(this) (device) (case (to use)(tense     (purpose)
                particle:     formative:
                object        past)
                case)

Relative Clause              Antecedent

--> (English Translation)

    The purpose for which (someone) used this device
    The purpose of using this device

(Rule 4) There is a limited number of nouns which are often used as antecedents in Type 4 relative clauses. Each of such nouns requires a specific phrase or clause to be interpolated in English.

[ex-3] [Type 4]

"KONO" "SOUCHI" "O"    "TSUKAT"  "TA"      "KEKKA"
(this) (device) (case (to use) (tense     (result)
                particle:      formative:
                object         past)
                case)

Relative Clause              Antecedent

--> (English Translation)

    The result which was obtained by using this device

In the above example, the clause "the result which someone obtained (the result : gap)" is ommited in Japanese, which relates the antecedent "KEKKA"(result) and the relative clause "KONO SOUCHI O TSUKAT_TA"(someone used this device).

269

A set of lexical rules is defined for
"KEKKA"(result), which basically works as follows :
it examines first whether the deep object case has
already been filled by a noun phrase in the
relative clause. If so, the relative clause is
taken as type 4 and an appropriate phrase is
interpolated as in [ex-3]. If not, the relative
clause is taken as type 1 as in the following
example where the noun "KEKKA" (result) fills the
gap of object case in the relative clause.

[ex-4] [Type 1]

```
"KONO"  "JIKKEN"    "GA"   "TSUKAT"  "TA"  "KEKKA"
(this)|(experiment)/(case\  (to use)/(tense|(result)
                   /particle:\        formative:|
                   |subject  |        past)    |
                   |case)    |                 |
     |_____|_____|  |
                      |                         |
              Relative Clause              Antecedent
```

-->(English Translation)

    The result which this experiment used


    Such lexical rules are invoked at the beginning of
the relative clause analysis by a rule in the main
flow of processing. The noun "KEKKA" (result) is
given a mark as a lexical property which indicates
the noun has special rules to be invoked when it
appears as an antecedent of a relative clause. All
the nouns which require special treatments in the
relative clause analysis are given the same marker.
The rule in the main flow only checks this mark and
invokes the lexical rules defined in the lexicon.

(Rule 5) Only the cases marked by postpositional
case particles 'GA', 'WO' and 'NI' can be deleted
in Type 1 relative clauses, when the antecedents
are ordinary nouns. Gaps in Type 1 relative clauses
can have other surface case marks, only when the
antecedents are special nouns such as described in
Rule (3).


4-2 Conjuncted Noun Phrases

    Conjuncted noun phrases often appear in
abstracts of scientific and technological papers.
It is important to analyze them correctly,
especially to determine scopes of conjunctions
correctly, because they often lead to proliferation
of analysis results. The particle "TO" plays
almost the same role as the English "and" to
conjunct noun phrases. There are several heuristic
rules based on various levels of information to
determine the scopes.


<Scope Decision Rules of Conjuncted Noun Phrases
by Particle 'TO'>

(Rule 1) Since particle "TO" is also used as a case
particle, if it appears in the position:
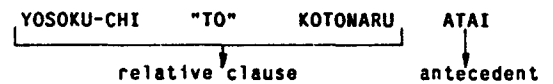
    Noun  'TO'  verb      Noun,
    Noun  'TO'  adjective  Noun,


there are two possible interpretations, one in
which "TO" is a case particle and 'noun TO
adjective(verb)' forms a relative clause that
modifies the second noun, and the other one in
which "TO" is a conjunctive particle to form a
conjuncted noun phrase. However, it is very likely
that the particle 'TO' is not a conjunctive
particle but a post-positional case particle, if
the adjective (verb) is one of adjectives (verbs)
which require case elements with surface case mark
'TO' and there are no extra words between "TO" and
the adjective (verb). In the following example,
"KOTONARU(to be different)" is an adjective which
is often collocated with a noun phrase followed by
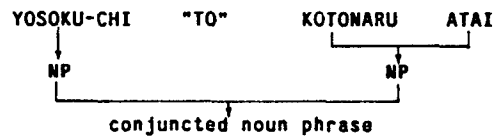case particle "TO".

[ex-5]
    YOSOKU-CHI    "TO"    KOTONARU    ATAI
    (predicted value)   (to be different)  (value)

    [dominant interpretation]

```
    YOSOKU-CHI    "TO"    KOTONARU  ATAI
   |_____|  |
                   |                     |
            relative clause         antecedent
```

    = the value which is different from the
      predicted value

    [less dominant interpretation]

```
    YOSOKU-CHI    "TO"     KOTONARU    ATAI
       |                      |_____|
       |                           |
       NP                          NP
   |_____|
                   |
          conjuncted noun phrase
```

    = the predicted value and the different value


(Rule 2) If two 'TO' particles appear in the
position:

    Noun-1  'TO'  .......... Noun-2  'TO'  'NO' NOUN-3

the right boundary of the scope of the conjuction
is almost always Noun-2. The second 'TO' plays a
role of a delimiter which delimits the right
boundary of the conjunction. This 'TO' is
optional, but in real texts one often places it to
make the scope unambiguous, especially when the
second conjunct is a long noun phrase and the scope
is highly ambiguous without it. Because the second
'TO' can be interpreted as a case particle (not as
a delimiter of the conjunction) and 'NO' following
a case particle turns the preceding phrase to a

270

modifier of a noun, an interpretation in which "NOUN-2 TO NO" is taken as a modifier of NOUN-3 and NOUN-3 is taken as the head noun of the second conjunt is also linguistically possible. However, in most cases, when two 'TO' particles appear in the above position, the second 'TO' is just a delimiter of the scope(see [ex-6]).

[ex-6]

```
YOSOKU-CHI TO JIKKEN     DE NO JISSOKU-CHI TO NO SA
(predicted| (experiment)(case| (actual value)   |
 value)   |             |particle:            (difference)
                         |place)          |
```

[dominant interpretation]

```
YOSOKU-CHI TO JIKKEN DE NO JISSOKU-CHI TO NO SA
    |          |                    |
    NP         NP
    |_____|
         |
     Conjuncted NP
         |_____
              NP
```

= the difference between the predicted value and the actual value in the experiment

[less dominant interpretations]

(A)
```
YOSOKU-CHI TO JIKKEN DE NO JISSOKU-CHI TO NO SA
    |          |
    NP         NP
    |_____|
         |
     Conjuncted NP
```

= the difference with the actual value in the predicted value and the experiment

(B)
```
YOSOKU-CHI TO JIKKEN DE NO JISSOKU-CHI TO NO SA
    |          |_____|
    NP                           NP
    |_____|
              |
          Conjuncted NP
```

= the predicted value and the difference with the actual value in the experiment

(Rule 3) If a special noun which is often collocated with conjunctive noun phrases appear in the position:

Noun-1 'TO' ........ Noun-2 'NO'<special-noun>,

the right boundary of the conjunction is almost always Noun-2. Such special nouns are marked in the lexicon. In the following example, "KANKEI" is such a special noun.

[ex-7]
```
JISSOKU-CHI "TO" RIRON-DE E-TA YOSOKU-CHI NO KANKEI
(actual value)|  |(theory) (to (predicted (relation-
                          obtain) value)   ship)
                                              ||
                                         special noun
```

[dominant interpretation]

```
JISSOKU-CHI   "TO" ....... YOSOKU-CHI NO KANKEI
    |                          |
                      (relative antecedent
                       clause)
                          |_____|
    |                          |
    NP                        NP
    |_____|
              |
          conjuncted NP
```

= the relationship between the actual value and the predicted value obtained by the theory

[less dominant interpretations]
(A)
```
JISSOKU-CHI "TO" RIRON-DE ...YOSOKU-CHI NO KANKEI
    |            |                        |
    NP          NP
    |_____|
         |
    conjuncted NP
         |_____|
    relative clause   antecedent
```

= the relationship of the predicted value which was obtained by the actual value and the theory

(B)
```
JISSOKU-CHI "TO" .......... YOSOKU-CHI NO KANKEI
    |                          |_____|
    NP                         NP
    |_____|
              |
          conjuncted NP
```

= the actual value and the relationship of the predicted value which was obtained by the theory

(Rule 4) In

Noun-1 'TO' ...... Noun-2,

if Noun-1 and Noun-2 are the same nouns, the right boundary of the conjunction is almost always Noun-2.

(Rule 5) In

Noun-1 'TO' ....... Noun-2,

if Noun-1 and Noun-2 are not exactly the same but nouns with the same morphemes, the right boundary

271

is often Noun-2. In [ex-7] above, both of the head
nouns of the conjuncts, JISSOKU-CHI(actual value)
and YOSOKU-CHI(predicted value), have the same
morpheme "CHI" (which meams "value"). Thus, this
rule can correctly determine the scope, even if the
special word "KANKEI"(relationship) does not exist.

(Rule 6) If some special words (like 'SONO'
'SORE-NO' etc. which roughly correspond to 'the',
'its' in English) appear in the position:

```
|Phrases which|Noun-1 'TO' <special word> Noun-2,
|modify noun  |
|phrases      |
```

the modifiers preceding Noun-1 modify only Noun-1
but not the whole conjuncted noun phrase.

(Rule 7) In

```
...... Noun-1 'TO' ............ Noun-2,
```

if Noun-1 and Noun-2 belong to the same specific
semantic categories, like action nouns, abstract
nouns etc. the right boundary is often Noun-2.

(Rule 8) In most conjuncted noun phrases, the
structures of conjuncts are well-balanced.
Therefore, if a relative clause precedes the first
conjunct and the length of the second conjunct (the
number of words between 'TO' and Noun-2) is short
like

```
[Relative Clause] Noun-1 'TO' ........ Noun-2
                              |length of the|
                              |2nd conjunct |
```

the relative clause modifies both conjuncts, that
is, the antecedent of the relative clause is the
whole conjuncted phrase.

These heuristic rules are based on different
levels of information (some are based on surface
lexical items, some are based on morphemes of
words, some on semantic information) and may lead
to different decisions about scopes. However, we
can distinguish strong heuristic rules (i.e. rules
which almost always give correct scopes when they
are applied) from others. In fact, there exists
some ordering of heuristic rules according to their
strength. Rules (1), (2), (3), (4) and (6), for
example, almost always succeed, and rules like (7)
and (8) often lead to wrong decisions. Rules like
(7) and (8) should be treated as default rules
which are applied only when the other stronger
rules cannot decide the scopes. We can define in
GRADE an arbitrary ordering of rule applications.
This capability of controlling the sequences of
rule applications is essential in integrating
heuristic rules based on heterogeneous levels of
information into a unified set of rules.

Note that most of these rules cannot be
naturally expressed by ordinary CFG rules. Rule
(2), for example, is a rule which blocks the
application of the ordinary CFG rule such as

NP ---> NP <case-particle> NO N

when the <case-particle> is 'TO' and a conjunctive
particle 'TO' precedes this sequence of words.

4-3 Determination of Scopes

Scopes of conjuncted noun phrases often
overlap with scopes of relative clauses, which
makes the problem of scope determination more
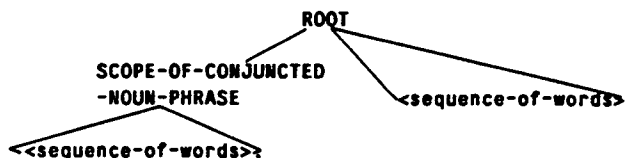complicated. For the surface sequence of phrases
like

NP-1 'TO' NP-2 <case-particle> ..... <verb> NP-3

there are two possible relationships between the
scopes of conjuncted noun phrase and the relative
clause like

(1) NP-1 'TO' NP-2 <case-particle> .... <verb> NP-3
```
      |
      v
   conjuncted
   noun phrase
                        Relative Clause        Antecedent

                              NP
```
(2)NP-2 'TO' NP-2 <case-particle> ..... <verb> NP-3
```
                              Relative Clause    Antecedent

                                  N,P
           Conjuncted Noun Phrase
```
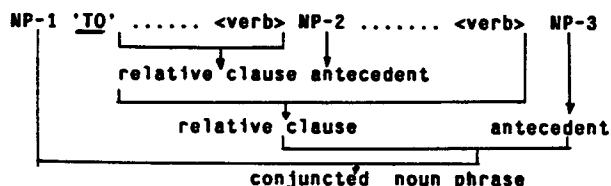
This ambiguity together with genuine ambiguities in
scopes of conjuncted noun phrases in 4-2 produces
combinatorial interpretations in CFG grammars, most
of which are linguistically possible but
practically unthinkable. It is not only
inefficient but also almost impossible to compare
such an enormous number of linguistically possible
structures after they have been generated. In our
analysis grammar, a set of scope decision rules are
applied in the early stages of processing in order
to block the generation of combinatorial
interpretations. In fact, the structure (2) in
which a relative clause exists within the scope of
a conjuncted noun phrase is relatively rare in real
texts, especially when the relative clause is
rather long. Such constructions with long relative
clauses are a kind of garden path sentence.
Therefore, unless strong heuristic rules like (2),
(3) and (4) in 4-2 suggest the structure (2), the
structure (1) is adopted as the first choice (Note
that, in [ex-7] in 4-2, the strong heuristic
rule[rule (3)] suggests the structure (2)). Since

272

the result of such a decision is explicitly expressed in the tree:

```
                    ROOT
              ___/    \___
  SCOPE-OF-CONJUNCTED
  -NOUN-PHRASE           <sequence-of-words>
    /    \
<<sequence-of-words>>
```
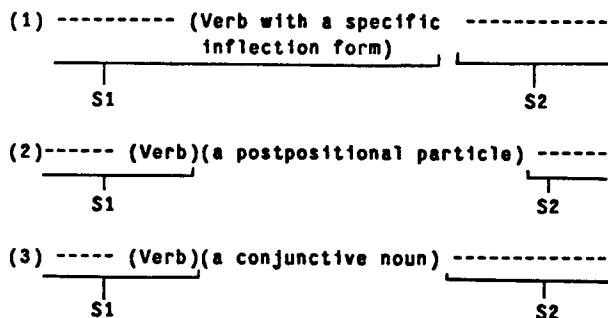
and the grammar rules in the later stages of processing work on this structure, the other interpretations of scopes will not be tried unless the first choice fails at a later stage for some reason or alternative interpretations are explicitly requested by a human operator. Note that a structure like

```
NP-1 'TO' ...... <verb> NP-2 ....... <verb> NP-3
    |        |                  |
    relative clause antecedent
    |
          relative clause              antecedent
    |
               conjuncted noun phrase
```

which is linguistically possible but extremely rare in real texts, is naturally blocked.


4-4 Sentence Relationships and Outer Case Analysis

Corresponding to English sub-ordinators and co-ordinators like 'although', 'in order to', 'and' etc., we have several different syntactic constructions as follows.

```
(1) ---------- (Verb with a specific ------------
               inflection form)
    _____|  |_____
         |                        |
         S1                       S2
```

```
(2)------ (Verb)(a postpositional particle) ------
    _____|                    |_____
         |                                  |
         S1                                 S2
```

```
(3) ----- (Verb)(a conjunctive noun) -------------
    _____|                   |_____
         |                              |
         S1                             S2
```

(1) roughly corresponds to English co-ordinate constructions, and (2) and (3) to English sub-ordinate constructions. However, the correspondence between the forms of Japanese and English sentence connections is not so straightforward. Some postpositional particles in (2), for example, are used to express several different semantic relationships between sentences, and therefore, should be translated into different sub-ordinators in English according to the semantic relationships. The postpositional particle 'TAME' expresses either 'purpose-action' relationships or

'cause-effect' relationships. In order to disambiguate the semantic relationships expressed by 'TAME', a set of lexical rules is defined in the dictionary of 'TAME'. The rules are roughly as follows.

(1) If S1 expresses a completed action or a stative assertion, the relationship is 'cause-effect'.

(2) If S1 expresses neither a completed event nor a stative assertion and S2 expresses a controllable action, the relationship is 'purpose-action'.

[ex-8]
(A) S1: TOKYO-NI      IT-   TEITA    TAME
        (Tokyo)       (to go) (aspect
                              formative)

    S2: KAIGI-NI  SHUSSEKI  DEKINAKA-   TA
        (meeting) (to attend) (cannot)(tense format-
                                       ive : past)

        S1: completed action
            (the aspect formative "TEITA" means
             completion of an action)

    ---> [cause-effect]
         • Because I was in Tokyo, I couldn't
           attend the meeting.

(B) S1: TOKYO-NI      IKU        TAME
        (Tokyo)       (to go)

    S2: KAIGI-NI  SHUSSEKI   DEKINAI
        (meeting) (to attend) (cannot)

        S1: neither a completed action nor
            a stative assertion
        S2: "whether I can attend the meeting
            or not" is not controllable.

    ---> [cause-effect]
         • Because I go to Tokyo, I cannot attend
           the meeting.

(C) S1: TOKYO-NI      IKU        TAME
        (Tokyo)       (to go)

    S2: KIPPU-O    KAT-         TA
        (ticket)   (to buy) (tense formative: past)

        S1: neither a completed action nor
            a stative assertion
        S2: volitional action

    ---> [purpose-action]
         • In order to go to Tokyo, I bought a
           ticket.

Note that whether S1 expresses a completed action or not is determined in the preceding phases

273

by using rules which utilize aspectual features of verbs described in the dictionary and aspect formatives following the verbs (The classification of Japanese verbs based on their aspectual features and related topics are discussed in [8]). We have already written rules (some of which are heuristic ones) for 57 postpositional particles for conjuctions of sentences like 'TAME'.

Postpositional particles for cases, which follow noun phrases and express case relationships, are also very ambiguous in the sense that they express several different deep cases. While the interpretation of inner case elements are directly given in the verb dictionary as the form of mapping between surface case particles and their deep case interpretations, the outer case elements should be semantically interpreted by referring to semantic categories of noun phrases and properties of verbs. Lexical rules for 62 case particles have also been implemented and tested.

5 Conclusions

Analysis Grammar of Japanese in the Mu-project is discussed in this paper. By integrating various levels of heuristic information, the grammar can work very efficiently to produce the most natural and preferable reading as the first output result, without any extensive semantic processings.

The concept of procedural grammars was originally proposed by Winograd[9] and independently persued by other research groups[10]. However, their claims have not been well appreciated by other researchers (or even by themselves). One often argues against procedural grammars, saying that: the linguistic facts Winograd's grammar captures can also be expressed by ATN, and the expressive power of ATN is equivalent with that of the augmented CFG. Therefore, procedural grammars have no advantages over the augmented CFG. They just make the whole grammars complicated and hard to maintain.

The above argument, however, misses an important point and confuses procedural grammar with the representation of grammars in the form of programs (as shown in Winograd[9]). We showed in this paper that: the rules which give structural constraints on final analysis results and the rules which choose the most preferable linguistic structures (or the rules which block "garden path" structures) are different in nature. In order to integrate the latter type of rules in a unified analysis grammar, it is essential to control the sequence of rule applications explicitly and introduce strategic knowledge into grammar organizations. Furthermore, introduction of control specifications doesn't necessarily lead to the grammar in the form of programs. Our grammar writing system GRADE allows us a rule based specification of grammar, and the grammar developed by using GRADE is easy to maintain.

We also discuss the usefulness of lexicon driven processing in treating idiosyncratic phenomena in natural languages. Lexicon driven prcessing is extremely useful in the transfer phase of machine translation systems, because the transfer of lexical items (selection of appropriate target lexical items) is highly dependent on each lexical item[11].

The current version of our analysis grammar works quite well on 1,000 sample sentences in real abstracts without any pre-editing.

References

[1] B.Vauquois: La Traduction Automatique a Grenoble, Documents de Linguistique Quantitative, No. 24, Paris, Dunod, 1975
[2] J.Nakamura et.al.: Grammar Writing System (GRADE) of Mu-Machine Translation Project and its Characteristics, Proc. of COLING 84, 1984
[3] J.Slocum: A Status Report on the LRC Machine Translation System, Working Paper LRC-82-3, Linguistic Research Center, Univ. of Texas, 1982
[4] F.Pereira et.al.: Definite Clause GRammars of Natural Language Analysis, Artificial Intelligence, Vol. 13, 1980
[5] G.Gazdar: Phrase Structure Grammars and Natural Languages, Proc. of 8th IJCAI, 1983
[6] Y.Wilks: Preference Semantics, in The Formal Semantics of Natural Language (ed: E.L.Keenan), Cambridge University Press, 1975
[7] Y.Sakamoto et.al.: Lexicon Features for Japanese Syntactic Analysis in Mu-Project-JE, Proc. of COLING 84, 1984
[8] J.Tsujii: The Transfer Phase in an English-Japanese Translation System, Proc. of COLING 82, 1982
[9] T.Winograd: Understanding Natural Language, Academic Press, 1975
[10] C.Boitet et.al.: Recent Developments in Russian-French Machine Translation at Grenoble, Linguistics, Vol. 19, 1981
[11] M.Nagao, et.al.: Dealing with Incompleteness of Linguistic Knowledge on Language Translation, Proc. of COLING 84, 1984