# Noun Phrase Identification in Dialogue and its Application

Izuru NOGAITO and Hitoshi IIDA

*ATR Interpreting Telephony Research Laboratories*
*Twin 21 MID Tower*
*2-1-61 Shiromi, Higashi-ku, Osaka 540, Japan*

[ABSTRACT]

Noun identifications are studied as 'anaphora'. Noun-noun relationships are ambiguous, as are noun-pronoun relations. Generally, nouns must match more 'antecedent' information than pronouns. But noun's 'antecedent' can be more remote. Therefore, the analysis scope of a noun-noun relationship will be larger than that of a noun-pronoun relationship. This expanded scope of analysis may result in errors. The nearest noun which a satisfies conditions is not always the 'antecedent'. A noun phrase identification model must be more comprehensive to deal with this problem.

In general, noun-noun anaphora are difficult to translate into another language because of the difference in the meaning of words. A machine translation system for dialogue must be able to identify noun phrases.

In this paper, a noun phrase identification model for understanding and translating the dialogue through use of the domain knowledge and a plan recognition model is presented. The presented model determines an area of analysis for dis-ambiguity.

A noun phrase identification model for understanding and translating dialogue through the use of domain knowledge and a plan recognition model is presented.

## 1. Introduction

A model for dialogue interpretation for machine translation is being studied. Previous machine translation techniques were inadequate for translating sentences containing the many ellipses characteristic of spoken dialogue. ' Topic' provides a key to recover the ellipses. If the 'topic' continues to the next sentence, some ellipses will be filled by nouns in previous sentences [Kuno78]. However, in dialogue it is normal to use different words to indicate the 'same' thing/concept. This expressional difference prevents matching and recovery of ellipses in understanding dialogue. A model is needed to match two different expressions which indicate the 'same' thing/concept.

Consider the following dialogue in Japanese with the topic "inquiries regarding an international conference". ('[]' shows word(s) <u>not</u> in Japanese sentences. 'Paper' is general translation of "shorui".)

Example    Dialogue 1

Applicant [d 1-1]:

> *Sankamoushikomiyoushi wo okutte kudasai.*
> registration form                 send     please
> (Could [you] send [me] [a] registration form?)

Secretary [d 1-2]:

> *Hai. Shorui wa sakihodo no jusho ni okureba yoroshiideshou ka ?*
> Ok. paper  given              address to send     shall
> (Ok) (Shall [I] send [a] paper to [the] address [you] have given [me]?)

In the Dialogue 1 example, 'registration form' and 'paper' are identical. 'Paper' in [d1-2] means 'registration form'. However, the words 'registration form' and 'paper' are different, and naive matching, which checks by 'words', fails. Because the 'topic' continuity is missing, 'I' and 'you' cannot be recovered. It is necessary to understand that 'registration form' and 'paper' indicate the 'same' thing.

Noun identifications are studied as 'anaphora'. Noun-noun relationships are ambiguous, as are noun-pronoun relations. Generally, nouns must match more 'antecedent' information than pronouns. But noun's 'antecedent' can be more remote. Therefore, the analysis scope of a noun-noun relationship will be larger than that of a noun-pronoun relationship. This expanded scope of analysis may result in errors. The nearest noun which a satisfies conditions is not always the 'antecedent'. A noun phrase identification model must be more comprehensive to deal with this problem

In general, noun-noun anaphora are difficult to translate into another language because of the difference in the meaning of words. A machine translation system for dialogue must be able to identify noun phrases.

In this paper, a noun phrase[1] identification model for understanding and translating the dialogue through use of the domain knowledge and a plan recognition model is presented. The presented model determines an area of analysis for dis-ambiguity.

## 2. Identification of Noun Phrases

### 2.1 Identification

A noun corresponds to a 'set name' of the thing/concept which the noun indicates. If there are two 'sets' where two nouns correspond, and each set is a subset of the other, then the two nouns are IDENT. If the two nouns are IDENT, two nouns indicate one set. For example, in dialogue 1 'registration form' and 'paper' are IDENT, and indicate the set called "registration form". IDENT merely shows that each thing/concept that the noun phrases indicate belong to the same set. IDENT does not show they are the same element of the same set. For instance, 'paper' in dialogue 1 shows the set of "registration form" but does not show that they are one and the same thing[2].

IDENT is the relationship between analyzing 'entry' (generalized noun), analyzed 'entry', and 'element' in our domain knowledge. It is assumed that every noun has an a priori 'element' of domain knowledge in accordance with the noun syntax form.

If entry-i and entry-j are IDENT, entry-j upon analysis is linked to analyzed entry-i, and is linked to element-i which corresponds to entry-i in domain knowledge.

### 2.2 Analysis of Noun Phrase Identification

### 2.2.1 Domain knowledge

Domain knowledge is specified knowledge of "task" and "domain". In this paper, domain knowledge is knowledge on 'international conference'

Domain knowledge consists of two types of knowledge; an 'element' which has some relationship to another 'element'.

'Element' is used to determine noun hierarchy, and corresponds to 'concept' of KL-ONE [Brachman85]. The "rel" of "element" corresponds to "role" in KL-ONE. This "rel" consists of three relationship "is-a", "has-a", and "causal". A classified "is-a" relationship does not include the "generic/individual relation" [Brachrnan83] the "is-a" relationship in

---

[1] In this paper,a 'noun phrase' means a 'simple noun phrase' without a predicate such as a verb or an adjective.

[2] There are more detailed identification analysis such as, definite/indefinite, specific/non-specific, number, etc. IDENT is one step of this more detailed identification analysis
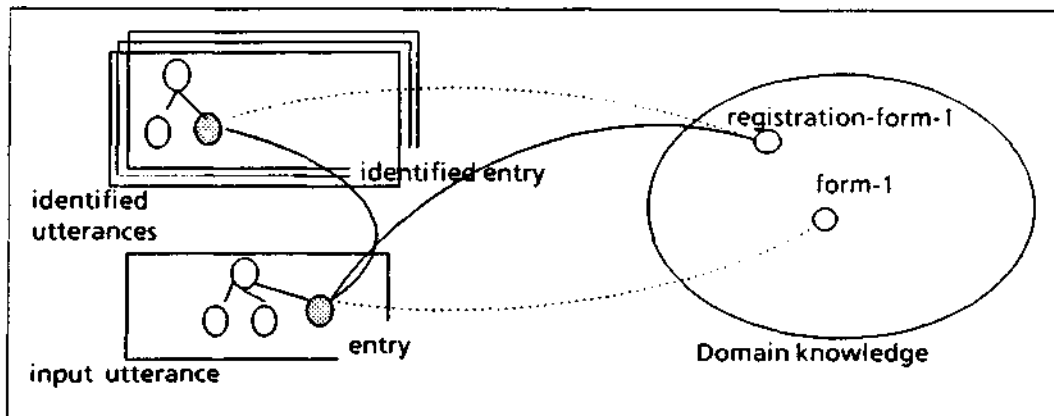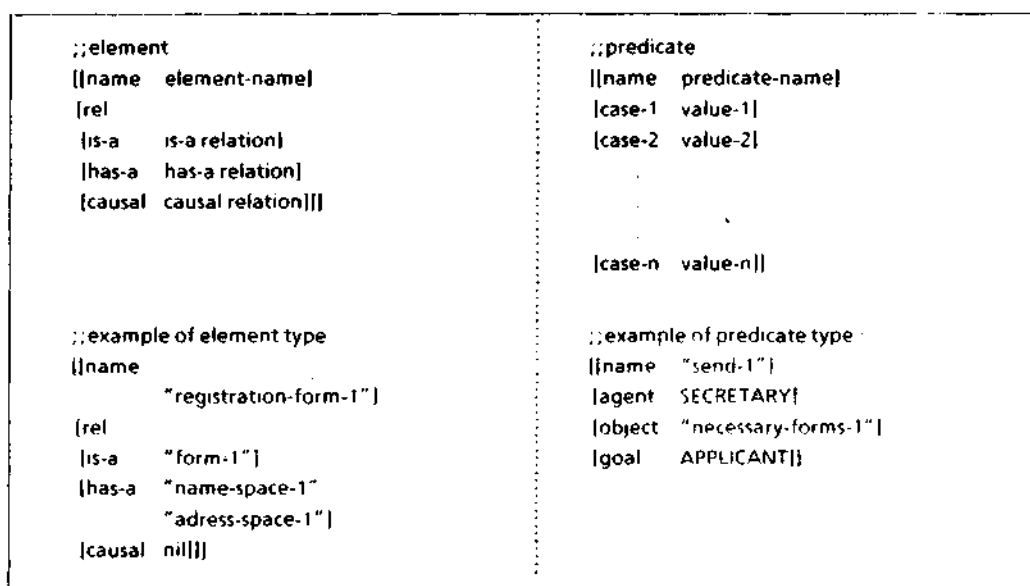
Figure1 Identification of Noun Phrases



Figure 2   Domain Knowledge

"element" gives only "generic/generic relations".  "Is-a","has-a", and  "causal"
relationships in "element" are conventionally defined except as above.

Another type of knowledge is called a 'predicate'. The 'predicate' corresponds to a
'verb' in an utterance. A 'predicate' has some case-role which fills an 'element'. This
'predicate' gives sentence-type domain knowledge such as "Secretary sends necessary-
forms to applicant", "applicant sends completed-necessary-form to office", etc.

### 2.2.2 Analysis

In the first step, an 'entry' is checked to determine the noun-noun relationship. The
'rel' of "elements" which corresponds to "entries" are matched with "identification
conditions".  In "identification conditions" there are patterns of "rel" which satisfy
"anaphoric relations". In general, this step determines whether one 'entry' is at higher

level in the hierarchy than another 'entry'. If entry-j is at a higher level than entry-i, the set of entry-j includes the set of entry-i.

The second step refers to 'predicate' domain knowledge. The result of the first step is 'context free'. 'Predicate' of domain knowledge is matched with the present utterance. The given candidates (entries) from the first step are substituted for the original entry in the utterance predicate's cases. The substituted utterance is then compared with 'predicate' in the domain knowledge.

Example Dialogue 2
(Continue from a 'topic' of 'accommodation charge for hotel')
Applicant [d 2-1]:

>*Hai. Wakarimashi ta. Tokorode, Kyoto eki kara takushii wo tukka ta* baai, *ryoukin wa ikrura gurai kakari masu ka?*
>(Ok. I see. By the way, How much is fee/fare when I take a taxi from Kyoto Station [to the Hall]?)

>In Japanese, "shukuhaku-ryoukin" (accommodation charge),and "takushii ryoukin" (taxi fare) can be called "ryoukin" (fee/charge/fare).

In dialogue 2, "ryoukin" (fee/charge/fare) does not mean "shukuhaku-ryoukin" (accommodation charge for hotel) immediately before, but rather "takusi ryoukin" (taxi fare). In this case, the 'clue words[3]' "by the way" show 'topic changing'.

The second step analysis does not deal with all dialogues. This model determines the analysis area by using the plan of the speaker and that of the hearer from the "plan recognition model" [Arita88]. The scope of analysis is "local scope" and "given scope" as follows.

Local Scope:　the 'topic' is assumed to continue:
>the last 'topic' is the same 'topic', or there is no 'topic' and no 'clue word' such as "by the way" which show 'topic' changes.

Scope:　　　the 'topic' is assumed to return or continue:
>the 'topic' is the 'sub-topic' of the last 'topic', the 'topic' returns to an old 'topic', or the 'topic' is confirmation of ending the dialogue. This information are from the analysis of the plan recognition model's output.

---------------------------------

[3] Clue words are expressions in utterances to control dialogue  Clue words can be classified into start signals (ex. "at the first", "by the way") and end signal (ex "I see")

2.2.3 Input and Output

An input of this noun identification model is assumed to be the semantic structure output from the unification parser [Kogure88]. input is given to this model for every utterance.

Example    Dialogue 3

   Secretary: [d 4-1]    *Tourokuyoushi wa sude ni o-moti de shou ka.*
                   Do you already have a registration form?

   Applicant: [d 4-2]    *Iie.*
                   No.

   Secretary: [d 4-3]    *Wakarimashita.*
                   I see.

   Secretary: [d 4-4]    *Dewa, Kotira kara youshi wo o-okuri shimasu.*
                   Then I will send you a form.

Example input of utterance[d 4-4] in Dialogue 4 is as follows.

```
?X04-01[[HEAD [[POS V]                        ;;;part-of-speech
             [CTYPE MASU]                     ;;;conjugation-type
             [CFORM SENF]                     ;;;conjugation-form
             [MODL [[POLT +]]]]]              ;;;modal
       [SUBCAT END]
       [SLASH  []]
       [SEM [[RELN OKURU-1]                   ;;;SEND-1
           [AGEN ?X04-02[[LABEL *SPEAKER*]]]  ;;;agent
           [RECP ?X04-03[[LABEL *HEARER*]]]   ;;;recipient
           [OBJE ?X04-04[[PARM ?X04-05]       ;;;object
                 [RESTR [[RELN YOUSHI-1]      ;;;FORM-1
                       [OBJE ?X04-05]]]]]
           [INFMANN [[RESTR [[RELN DEWA-1]    ;;;restriction, THEN-1
                         [OBJE ?X04-06]]]
                    [PARM ?X04-06]]]]]
       [PRAG [[SPEAKER ?X04-02]               ;;;pragmatics
            [HEARER ?X04-03]
                                ]]]]]]]]
       [WH []]]]]]
```

This structure is based on the output of the parser [Kogure88], but there are two expansions of the feature structure.

The figures which begin with '?' show feature structure token identity (e.g. ?X01-17). The figures added to a PREFIX shows the utterance number for identification in the dialogue. And a feature for CLUE WORDS is added. If there is a CLUE WORDS in an utterance, the next structure is added.

```
[CLUE-WORDS [[RELN+ TOKORODE-1 ]]]   :;;By the way
```

The feature structures in bold face in the example above , namely SEMANTIC[SEM], TOPIC, and CLUE-WORDS, are analyzed in Noun Identification. In this example, *SPEAKER* and *HEARER* are speaker and hearer of the utterance.

The Noun Identification Model deals with this input and IDENTIFIED UTTERANCES. IDENTIFIED UTTERANCES is the queued previous output of the Noun Identification Model. The output of the Noun Identification Model consists of the above input structure and the N-IDENT structure which shows noun identification. The N-IDENT structure of utterance [D4-4] is as follows,

```
[N-IDENT [TOPIC+ [SCOPE LOCAL]              ;;;topic and analysis area
                [RELN+ []]                  ;;;relation
                [OBJE+ []]]                 ;;;object
           [ENTRIES
            [[RELN+ YOUSHI-1]               ;;;FORM-1
             [OBJE+ ?X04-04]
             [IDENT+ [[RELN+ TOUROKUYOUSHI-1] ;;;REGISTRATION-FORM-1
                     [OBJE+ ?X01-02]]]]] ]   ;;;in utterance [04-1]
                                     ]]]]]
```

TOPIC + shows the analysis area and TOPIC in that utterance. TOPIC is obtained from SEMANTIC structure. An analysis area SCOPE is obtained from previous N-IDENT in IDENTIFIED UTTERANCE. There is no TOPIC in this example utterance. LOCAL is selected for analysis area.

ENTRIES shows noun-noun identification, in this example YOUSHI in [D4-4] is identified with TOUROKUYOUSHI in [D4-1].

The output of [D4-1] is as follows. In this output, the bold face part is the antecedent of YOUSHI-1. The N-IDENT part shows that TOUROKUYOUSHI-1 is TOPIC and there is no antecedent of TOUROKUYOUSHI-1.

```
?X01-01[[HEAD [[POS V]
            [TOPIC ?X01-02[[RESTR [[RELN TOUROKUYOUSHI-1] ;;;TOPIC
                                   [OBJE ?X01-03]]]  ;;;REGISTRATION-FORM-1
                    [PARM ?X01-03]]]           :;; parameter
           [CFORM SENF]                 ;;;conjugation-form. sentence-final
           [MODL [[SFP-1 KA]]]          :;;modal
           [CTYPE NONC]]]               ;;;conjugation-type, nonconjugate
      [SUBCAT END]
      [SLASH  []]
      [SEM [[RELN S-REQUEST]                          ;;;surface request
            [AGEN ?X01-04[]]                          ;;;agent
            [RECP ?X01-05[]]                          ;;;recipient
            [OBJE [[RELN ?X01-06 INFORMIF]
                   [AGEN ?X01-05]
                   [RECP ?X01-04]
                   [OBJE [[RELN ?X01-07 MOTSU-1]               ;;; HAVE-1
                          [OBJE ?X01-02]
                          [MANN ?X01-08[RESTR [[RELN SUDENI-1]    ;;; ALREADY-1
```

```
                                          [OBJE ?X01-09]
                                          [PARM ?X01-09]]]]]]
        [PRAG [[SPEAKER ?X01-04]
        [HEARER  ?X01-05]
        [RESTR [[FIRST
          ...
[N-IDENT [TOPIC+ [SCOPE+ LOCAL]                ;;;topic and analysis area
               [RELN+ TOUROKUYOUSHI-1]    ;;:REGISTRATION-FORM-1
               [OBJE+ ?X01-02]]
          [ENTRIES
             [[RELN+ TOUROKUYOUSHI-1]          ;;;REGISTRATION-FORM-1
              [OBJE+ ?X01-02]]]
             [IDENT+ []]]]]]                    ;;;NO ANTECEDENT
                                          ]]]  ]]
```

The presented Noun Identification Model is implemented on a rule based expert system for verifying rules. Using data of simulated dialogue, noun identification rules have been verifying in detail in a while dialogue.

Examples of rules are as follows,

IF:: There is an ENTRY in the present utterance.
    THEN:: ENTRY is copied to N-IDENT.

IF:: There is an ENTRY in IDENTIFIED UTTERANCES. And
    ENTRY in IDENTIFIED UTTERANCES and present ENTRY are satisfied with
    IDENTIFICATION CONDITIONS.
        Then::  ENTRY in IDENTIFIED UTTERANCES and present ENTRY are in WEAK-
            IDENT.

IF:: Corresponding ELEMENTS in Domain Knowledge are in IS-A relationship.
        THEN:: ENTRY in IDENTIFIED UTTERANCES and present ENTRY are satisfied
            with IDENTIFICATION CONDITIONS.


IF:: Corresponding ELEMENTS in Domain Knowledge are in CAUSAL relationship.
        THEN:: ENTRY in IDENTIFIED UTTERANCES and present ENTRY are satisfied
        with IDENTIFICATION CONDITIONS.

IF:: There is a TOPIC in the present utterance.
        THEN:: TOPIC is copied to N-IDENT.


3. Application

3.1 Translation and generation

A machine translation system for dialogue requires the generation of natural dialogue in a target language.

A noun's meaning may be different in another language. Identification of nouns is needed. "Youshi" is usually correspond with paper, sheet or form
The meaning of "youshi" and "form" are shown below.

YOUSHI

*:a paper which is used for some purpose.*

Form

*:a printed paper divided by lines into separate parts,in each of which* answers to *questions must be written down.*

In translating from the Japanese word "youshi" to English, "youshi" includes no information about "lines" and so on in the current utterance, so a machine translation system must understand what "youshi" is, and whether "youshi" has "lines" or not. If "form" is given as a wrong translation of "youshi" which has no lines, A hearer cannot find the form that 'form' indicates. An erroneous translation does not keep an "anaphoric relation" equivalent between Japanese and English, and thus makes a translated sentence difficult to understand.

In identifying "youshi", a system can obtain 'more detailed' information from an identified 'element' such a "moushikomiyoushi" which has "lines" And this information gives a correct translation. "Youshi" can be transferred into "registration form" [see figure 3] The transferred word is linked to "form" by IS-A relationship. And "form", which is one of translation candidates, can be a translation of "youshi". "Registration form" and "form" can be in an anaphoric relationship.

Hence, in the translation of nouns in anaphoric relationship, these nouns need to be in anaphoric relationship in the target language. For an anaphoric relationship in the target language, the noun identification model is able to understand the anaphoric relationship in the source language, and a system can choose a translation to maintain the identification relationship in the target language.

*3.2* Matching to internal data

In a machine translation, there are some types of internal data. A variation in noun expressions also causes failure of naive matching to internal data. Many noun phrases cannot be matched to the internal data which the noun phrases indicate. Thus many noun phrases are difficult to link to "more detailed" and "more specific" internal data. In domain dependent understanding, these nouns are often keywords of understanding. "Accommodation charge" will show that an utterance is about "a hotel reservation", but using "charge" in place of "accommodation charge" will not yield "a hotel reservation". Identification of a noun phrase gives "detailed" information from linked 'elements' and 'entries'.
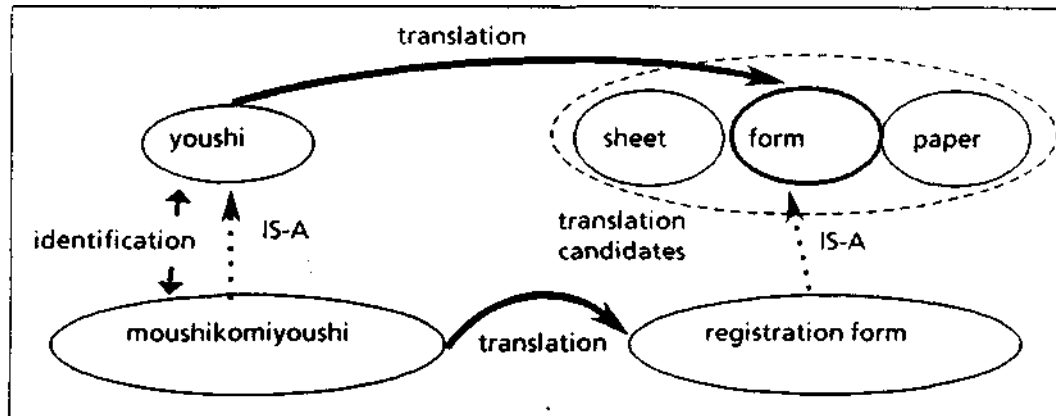
Figure 3 Translation of Anaphoric relationship

3.3 Interpreting Definite Noun Phrases

Japanese has no 'definite' or 'indefinite' marker such as 'the', 'a', or 'an'. Whether a noun phrase is "definite" or "indefinite" must be determined from a noun without a marker. The naive definition of 'definite noun phrase' is 'thing/concept which can be identified by both speaker and hearer'. A "known" word satisfies this naive definition, but a "known word" does not mean the "same" word. The same words of "form" are used in the dialogue, but one will be for "registration form", another may be a new "form" for "bank transfer form". In this case, the second "form" is not definite. If a noun is identified with previous noun, then the noun will be definite.

4. Conclusion

A noun phrase identification model for understanding and translating dialogue through the use of domain knowledge and a plan recognition model is presented.

A model for understanding noun-pronoun relationship using 'Focus Space' [Sidner83] has been presented. This model needs noun-noun identification for obtaining 'Focus Space'. This means that the model for understanding noun-pronoun relations cannot deal with noun-noun relationship.

Noun-noun anaphoric relationships are sometimes ambiguous, as are noun-pronoun relationships in understanding. When 'antecedents' of noun-noun relationship in dialogue are more remote, this expanded analysis scope may result in error.

The presented model determines an area of analysis for eliminating ambiguity.

In simulated dialogue [Arita87], there are many usages of a noun which is a part of another noun. Most nouns in this relationship are also in "anaphoric relations", e.g. HOTEL SHUKUHAKU RYOUKIN / SHUKUHAKU RYOUKIN [hotel accommodation charge/ accommodation charge]. Generally speaking, a candidate which satisfies "Identification Conditions" is also a part of another noun. A candidate is often in an anaphoric

relationship. Using this pragmatic knowledge will narrow the choice of the antecedent in this Noun Identification Model.

## Acknowledgment

## References

|Arita87|    Arita, H. et al : Media dependent conversation manners, 1987, WGNL Meeting Report 61-5, Information Processing Society of Japan

[Arita88|    Arita, H. et al.: Construction of a Discourse Structure Using a Plan Recognition Model, 1988, Annual Meeting Report 5T-3, Information Processing Society of Japan

|Brachman83|    Brachman, R: What IS-A is and isn't: Analysis of Taxonomic Links in  Semantic Networks,"COMPUTER", Oct 1983, IEEE

|Brachman85|    Brachman, R: An Overview of the KL-ONE Knowledge Representation  System, COGNITIVE SCIENCE.9 1985

[Kogure88|    Kogure, K. et al.: A method of analyzing Japanese speech act types. The  2nd International Conference on Theoretical and Methodological Issues in  Machine Translation of Natural Languages, Jun. 1988

|Kuno78|    Kuno, S.: Danwa no Bunpou [Discourse Grammar], 1978, Taisyu-shoten

|Sidner83l    Sidner, C : Focusing in the Comprehension of Definite Anaphora,"  Computational Model of Discourse",1983, MIT PRESS