# Context Analysis System for Japanese Text

Hitoshi Isahara     and     Shun Ishizaki

Electrotechnical Laboratory
1-1-4, Umezono, Sakura-mura, Niihari-gun,
Ibaraki, Japan 305

## ABSTRACT

A natural language understanding system is described which extracts contextual information from Japanese texts. It integrates syntactic, semantic and contextual processing serially. The syntactic analyzer obtains rough syntactic structures from the text. The semantic analyzer treats modifying relations inside noun phrases and case relations among verbs and noun phrases. Then, the contextual analyzer obtains contextual information from the semantic structure extracted by the semantic analyzer. Our system understands the context using precoded contextual knowledge on terrorism and plugs the event information in input sentences into the contextual structure.

## 1: Introduction

Despite the advanced state of syntactic analysis research for natural language processing and the many useful results it has produced, there have been few studies involving contextual information, and many problems remain unsolved.

The natural language understanding system described here employs a syntactic analyzer, a semantic analyzer treating modifying relations inside noun phrases and the relations among verbs and phrases, that is, word-level semantics, and a contextual analyzer (Fig. 1). These analyzers operate in a serially integrated fashion. Though humans seem to understand natural language texts using these three analyzers simultaneously, we have made their methodology essentially different from their human counterparts for more efficient computing. Our system uses a context-free grammar parser named Extended-Lingol as a syntactic analyzer to analyze the Japanese sentences and produce parsing trees. From an analysis of these, in turn, it obtains word-level semantic structures expressed in frame-like representations. Finally, it extracts contextual information, using our representation from the semantic structures. We remain far from certain at this stage whether this system represents the best realization of an engineering-based natural language understanding system. Future plans include combining these three processes into one process and bringing the system closer to the human process.

Because our system uses bottom-up analysis first (including syntactic analysis and word-level semantic analysis), it can obtain not only the outline of the input sentences but also their details, as necessary. This method is the best one in situations where the detailed information of texts are quite important, such as Machine-Translation systems and precise question-answering systems. Of course, in this way, we must build up a sizable dictionary of precise word definitions.

In our system, predictive-style processing is not used in syntactic analysis and word-level semantic analysis. But, in the contextual analysis part, predictions from the tree structure of the contextual information are used for instantiation of the contextual structure.

We are now developing a system which can understand newspaper articles through contextual structure (see Fig. 2a). After applying the procedures outlined above, the system obtains
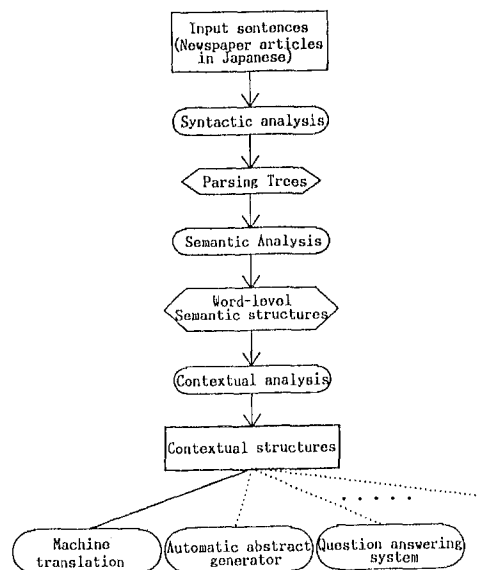


Fig.1  System flow chart of this paper and its applications.



a: Original input (Morning edition of the Asahi Shimbun--July 30, 1983).

THE BOMB KILLS FOUR PEOPLE INCLUDING A JUDGE.
[Rome 29th = correspondent Hirano]

In the morning of the 29th, at Palermo, Sicily in Italy, a parked car exploded, which killed 4 people including a judge who had directed an investigation into Mafia crimes, and injured about 10 people seriously or slightly. This is the fourth murder case on judges at Palermo and is of the largest scale.

Judge Rocco Chinnici, 58, the director of the Palermo preliminary court, police bodyguards and others were murdered. At the moment when the judge left home, the bomb exploded which had been set in the car of Fiatt parked near there. The explosion involved the residents, windows of the apartment and about 10 cars near there.

b: The translation of the example article (a) from Japanese into English.

Fig. 2. An example of newspaper articles

contextual representations expressed as shown in Fig. 3. Some details of the input text are abbreviated in the figure.

## 2: Syntactic and semantic analysis[2]

Let us proceed to an explanation of the methodologies adopted by our system, using the newspaper article in Fig. 2a as an example. First, the system analyzed each sentence syntactically, obtaining parsing trees. Next, the system constructs a semantic structure for each phrase. Word meanings in our word dictionary are described in SRL (Semantic Representation Language) which uses frame-like expression as shown in Fig. 4. Each word meaning shares a suitable position in the hierarchy of concepts. SRL enables deep semantic analysis in a flexible way. The formal definition of its syntax and semantics is not stated here. In our system, a word meaning written in the lexical entry using SRL plays an important role in semantic analysis. The interaction between the word meanings is the central issue of the semantic analysis. The modifying relations inside noun phrases and the case relations among verbs and noun phrases are determined in the word-level semantic structure. In Fig. 4, three scenes (explosion, death and injury) are obtained by analyzing the first sentence of the article in Fig. 2a. "Human" is a dummy node that means human beings. Here, the people who died include a judge and some policemen.

There are several types of ambiguity in input text. In syntactic analysis, ambiguity means the existence of several parsing trees. Word-level semantics often specify which should be selected. Here, we should use a kind of prediction. For example, people who are in authority could be a target of terrorism (See Fig. 2a). These constraints are very helpful in eliminating ambiguity, as well as surface syntactic information. Some of this processing is done in an interactive way in our system. Our system asks the user how to specify the relations between events in some decision points. Even after the elimination of ambiguity by the word semantics, there may be unsolved ambiguities. These will be eliminated by contextual analysis with the contextual structure.

## 3: Features of contextual representation

Our contextual structure fits into a tree structure with one root node and a number of leaf nodes. Relations between events in a story are defined in the structure as "scenes", and the relations among our structure are defined by a tree structure. Our structure can share scenes with others.

Leaf nodes with a shared root node have either an "and" or an "or" relationship with each other. The hierarchy shown in Fig. 5 is an example. The node "terrorism involving bomb" has, as in Fig. 5, three leaf nodes (scenes) - "explosion," "damage" and "rescue". Since these seem to occur serially, the relationship among them is an "and" relationship. On the other hand, the root node "terrorist action" in Fig. 5 has several leaf nodes - "terrorism involving bomb", "shooting" and so on. As only one of these usually corresponds to the main topic in newspaper stories, they share an "or" relationship with each other.

Input events are matched not only directly with scenes in the structure, but also with higher concepts in accordance with a predefined tree structure of a concept hierarchy like that in Fig. 6. In other words, the system has a concept thesaurus. So, matching between the scene of the structure and the input events becomes flexible.
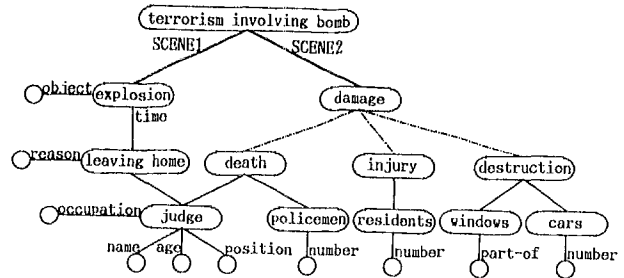


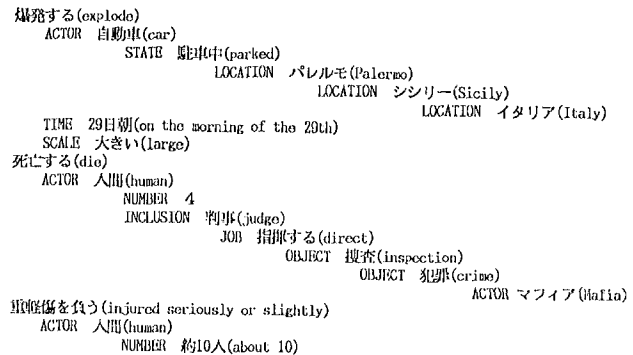Fig. 3. An example of the contextual structure.



Fig. 4. The word-level semantic structure extracted from the first sentence in Fig. 2a.
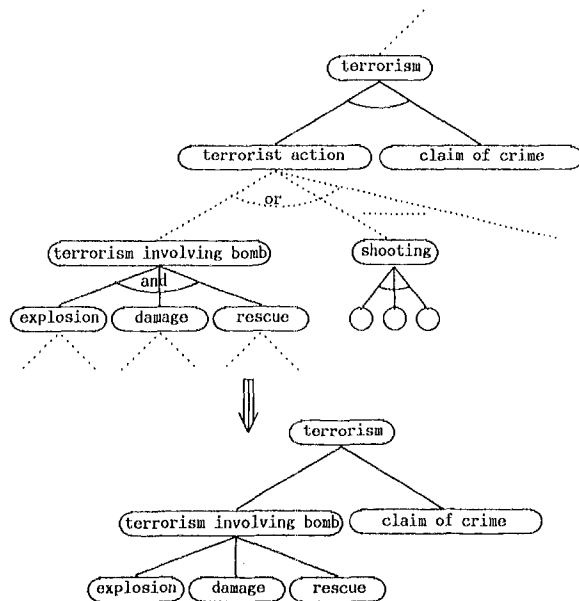


Fig. 5 The contextual structure (upper diagram) and its reorganization (lower diagram).

245

## 4: Contextual structure selection process

Now we have implemented two selection methods for the selection of the contextual structure, a "two-event method" and a "title-based method". First, we will explain the "two-event method".

In the "two-event method", titles are not processed by the system for selection. In sentence processing, after two events are obtained, the system begins a search for a structure involving these two events as their scenes. The use of two events helps decrease the number of possible structures during the search. As mentioned previously, selection of suitable structures and scenes can be accomplished flexibly with the concept thesaurus.

After developing the "two-event method", we began to implement the "title-based method". In the case of newspaper articles, titles have important information for the selection of suitable contextual structures. If there is a special word (noun or verb) in the title, contextual representation indicated by that word is selected. In this way, the system can almost always select suitable structures. Newspaper titles should be written so that readers can get enough information for the selection of the topic from its title only. The correct selection rate of our "title-based method" is shown in Table 1. Derivatives point to their original words, and, through them, derivatives can select suitable structure.

Within our experience, there are no differences in the correct selection rates between these two methods. In our system, at present, we use the "title-based method" because of its similarity to human behaviour.

## 5: Contextual analysis

Once a promising structure is discovered, scenes corresponding to the input events are selected in the following manner: if an event in the input sentence matches one of the scenes already activated in the system, it identifies the event with that scene.

For example, from the article shown in fig. 2a our system extracts three events - "explosion", "murder (death)" and "injury". The contextual structure of "terrorism involving bomb" is then selected using its title, and the contextual analysis begins. In the contextual analysis, first, "explosion" matches directly the first scene (scene1) in the structure, "explosion", and this event is plugged into scene1. Next, "murder" is checked comparing with each scene. Here, there is no scene which directly matches "murder" but one of the higher concept of "murder" is "damage". So the system identifies "murder" with scene2, "damage" and plugs that event into scene2. In these cases there are no events already plugged into the selected scene, so the system can easily plug the events into the scenes. When processing the third event, "injury", it is quite important to determine whether this event is the same or different from "murder". (Here, "injury" also has the higher concept, "damage".)

The determination whether this part of the input sentence is giving absolutely new information about a new event now being introduced or more precise information about the event already described is accomplished in the following manner. When the input event is identified with one scene in the contextual structure, the contextual analyzer begins to search for an event already plugged into the selected scene that has the same concept (or a higher concept) as the input event. If there is no conflict between the values of the attributes (for example, ACTOR, OBJECT, TIME) in the input event and the event found by the search, the
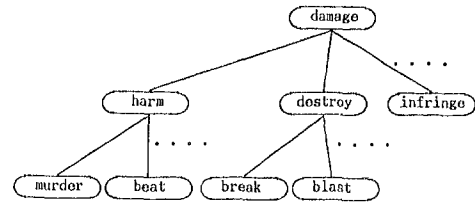


Fig. 6. Concept thesaurus.

Table 1. Topic selection by the "title-based method".

| Topic | Successes | | Failures |
|---|---|---|---|
| | By title | By subtitle | |
| Kidnapping | 2 4 | 1 | 0 |
| Explosion | 7 6 | 1 2 | 2 |
| Shooting | 2 2 | 1 | 0 |
| Attack | 1 2 | 1 | 3 |
| Highjack | 1 6 | 6 | 0 |
| Total | 1 5 0 | 2 1 | 5 |

information obtained by the input event is treated as a detail of a predescribed event. If this is not the case, the input event is treated as a parallel event of the events in that scene.

## 6: Conclusion

In the above sections we have proposed a system for extracting contextual information from natural language texts using a contextual representation structure as a knowledge structure. Our structure has proven itself useful for expressing contents of Japanese newspaper articles. Though we propounded the method used in our system to understand natural language texts in every field, some of its specifications such as the treatment of titles are oriented toward special features of newspaper articles.

At present, the applications of this system are restricted to stories dealing with terrorism. For these limits to be extended, the number of the contextual structures must be increased and the concept thesaurus scale enlarged. We believe that the natural language understanding system described in this paper is flexible enough to allow for such extension. Computer facilities must, of course, also be taken into account.

As our system is still in the development stage, some parts are not yet complete. Our dictionary is still rather small. For these reasons, the scope of this paper has been limited to processing ability for a restricted category of newspaper articles.

Reference
1: Schank, R. C., "Dynamic Memory", Cambridge Univ. Press (1983)
2: Tanaka, H., "A Semantic Processing System for Natural Language Understanding", (in Japanese) Researches of the Electrotechnical Laboratory, No.797 (1979)
3: Lytinen, S. L. and Schank, R. C., "Representation and Translation", Text 2:1/3 Yale Research Report (1982)
4: Ishizaki, S., Isahara, H. and K. Handa, "Natural Language Processing System with Deductive Learning Mechanism", International Symposium on Language and Artificial Intelligence, March 16-21, 1986, Kyoto