

A Discourse Structure Analyzer for Japanese Text*

K. Sumita, K. Ono, T. Chino, T. Ukita, and S. Amano

Toshiba Corp. R&D Center
Komukai-Toshiba-cho 1, Saiwai-ku, Kawasaki 210, Japan
sumita@isl.rdc.toshiba.co.jp

Abstract

This paper presents a practical procedure for analyzing discourse structures for Japanese text, where the structures are represented by binary trees. In order to construct discourse structures for Japanese argumentative articles, the procedure uses local *thinking-flow* restrictions, *segmentation rules*, and *topic flow* preference. The thinking-flow restrictions restrict the consecutive combination of relationships detected by connective expressions. Whereas the thinking-flow restrictions restrict the discourse structures locally, the segmentation rules constrain them globally, based on rhetorical dependencies between distant sentences. In addition, the topic flow preference, which is the information concerning the linkage of topic expressions and normal noun phrases, chooses preferable structures. Using these restrictions, the procedure can recognize the scope of relationships between blocks of sentences, which no other discourse structure analysis methods can handle. The procedure has been applied to 18 Japanese articles, different from the data used for algorithm development. Results show that this approach is promising for extracting discourse information.

1 Introduction

A computational theory for analyzing linguistic discourse structure and its practical procedure are necessary to develop machine systems dealing with plural sentences; e.g., systems for text summarization and for knowledge extraction from a text corpus.

Hobbs developed a theory in which he arranged three kinds of relationships between sentences from the text coherency viewpoint [Hobbs 1979]. Grosz and Sidner proposed a theory which accounted for interactions between three notions on discourse: linguistic structure, intention, and attention [Grosz and Sidner 1986]. Litman and Allen described a model in which a discourse structure of conversation was built by recognizing a participant's plans [Litman and Allen 1987]. These theories all de-

pend on extra-linguistic knowledge, the accumulation of which presents a problem in the realization of a practical analyzer. The authors aim to build a practical analyzer which dispenses with such extra-linguistic knowledge dependent on topic areas of articles to be analyzed.

Mann and Thompson proposed a linguistic structure of text describing relationships between sentences and their relative importance [Mann and Thompson 1987]. However, no method for extracting the relationships from superficial linguistic expressions was described in their paper. Cohen proposed a framework for analyzing the structure of argumentative discourse [Cohen 1987], yet did not provide a concrete identification procedure for 'evidence' relationships between sentences, where no linguistic clues indicate the relationships. Also, since only relationships between successive sentences were considered, the scope which the relationships cover cannot be analyzed, even if explicit connectives are detected.

This paper discusses a practical procedure for analyzing the discourse structure of Japanese text. The authors present a machine analyzer for extracting such structure, the main component of which is a structure analysis using thinking-flow restrictions for processing of argumentative documents. These restrictions, which examine possible sequences of relationships extracted from connective expressions in sentences, indicate which sentences should be grouped together to define the discourse structure.

2 Discourse structure of Japanese text

2.1 Discourse structure

This paper focuses on analyzing discourse structure, representing relationships between sentences. In text, various rhetorical patterns are used to clarify the principle of argument. Among them, connective expressions, which state inter-sentence relationships, are the most significant. They can be divided into the categories described in Table 1.

Here, connective expressions include not only normal connectives such as "therefore", but also idiomatic

*This work was supported by ICOT (Institute for New Generation Computer Technology), and was carried out as a part of the Fifth Generation Computer Systems research.

expressions stating relations to the other part of text such as "in addition" and "here ... is described." The authors extracted 800 connective expressions from a preliminary analysis of more than 1,000 sentences in several argumentative articles [Ono *et al.* 1989]. Then, connective relationships were classified into 18 categories as shown in Table 1. Using these relationships, linguistic structures of articles are captured.

Table 1 is the current version of the relationship categories. The number of relationship categories necessary and sufficient to represent discourse structures must be determined through further experimentation. New categories will be formed as need becomes apparent; likewise, categories found to overlap in function will be merged. Final categorization can only be fixed after extensive analysis.

Sentences of similar content may be grouped together into a block. Just as each sentence in a block serves specific roles, e.g., "serial", "parallel", and "contrast", each block in text serves a similar function. Thus, the discourse structure must be able to represent hierarchical structures as well as individual relationships between sentences. In this paper, a discourse structure is represented as a binary tree whose terminal nodes are sentences; sub-trees correspond to local blocks of sentences in text.

Figure 1 shows a paragraph from an article titled "a zero-crossing rate which estimates the frequency of a speech signal," where underlined words indicate connective expressions. Figure 2 shows its discourse structure. Extension relationships are set to sentences without any explicit connective expressions. Although the fourth and fifth sentences are clearly the exemplification of the first three sentences, the sixth is not. Thus, the first five can be grouped into a block.

Discourse structure can be represented by a formula. The discourse structure in Figure 2 corresponds to the following formula.

$$[[[1 \langle EX \rangle [2 \langle EX \rangle 3]] \langle EG \rangle [4 \langle EX \rangle 5]] \langle SR \rangle 6].$$

2.2 Local constraint for consecutive relationships

For analyzing discourse structure, a local constraint on consecutive relationships between blocks of sentences is introduced. The example shown in Figures 1 and 2 suggests that the sequence of connective relationships can limit the accepted discourse structures to those most accurately representative of original argumentative text. Consider the sequence [P <EG> Q <SR> R], where P, Q, R are arbitrary (blocks of) sentences. The premise of R is obviously not only Q but both P and Q. Since the argument in P and Q is considered to close locally, the two should be grouped into a block. This is a local constraint on natural argumentation.

Table 1: Connective relationships.

RELATION	EXAMPLES and EXPLANATION
serial connection <SR>	だから (thus, therefore), よって (then) <i>dakara yotte</i>
negative connection <NG>	だが (but), しかし (though) <i>daga shikashi</i>
reason <RS>	なぜなら (because), <i>nazenara</i> その訳は (the reason is ...) <i>sono wake wa</i>
parallel <PA>	同時に (at the same time), <i>doujini</i> さらに (in addition) <i>sarani</i>
contrast <CT>	一方 (however), 反面 (on the contrary) <i>ippou hanmen</i>
exemplification <EG>	例えば (for example), <i>tatoeba</i> ... 等である (and so on) ... <i>nado dearu</i>
repetition <RP>	というのは (in other words), <i>toiuowa</i> それは (it is ...) <i>sore wa</i>
supplementation <SP>	もちろん (of course) <i>mochiron</i>
rephrase <RH>	つまり, すなわち (that is ...) <i>tsumari sunawachi</i>
summarization <SM>	結局 (after all), まとめると (in sum) <i>kekkyoku matomeruto</i>
extension <EX>	これは (this is) <i>kore wa</i>
definition <DF>	ここで ... とする (... is defined as ...) <i>koko de ... to suru</i>
rhetorical question <RQ>	なぜ ... なのだろうか (Why is it ...) <i>naze ... nanodarouka</i>
direction <DI>	ここでは ... を述べる <i>kokode wa ... wo noberu</i> (here ... is described)
reference <RF>	図Xに ... を述べる (Fig.X shows ...) <i>zu X ni ... wo noberu</i>
topic shift <TS>	さて, ところで (well, now) <i>sate tokorode</i>
background <BG>	従来 (hitherto) <i>juurai</i>
enumeration <EN>	第一に (in the first place), <i>dai 1 ni</i> 第二に (in the second place) <i>dai 2 ni</i>

- 1 : In the context of discrete-time signals, zero-crossing is said to occur if successive samples have different algebraic signs.
- 2 : The rate at which zero crossings occur is a simple measure of the frequency content of a signal.
- 3 : This is particularly true of narrow band signals.
- 4 : For example, a sinusoidal signal of frequency F_0 , sampled at a rate F_s , has F_s/F_0 samples per cycle of the sine wave.
- 5 : Each cycle has two zero crossings so that the long-term average rate of zero-crossings is $Z = 2F_0/F_s$
- 6 : Thus, the average zero-crossing rate gives a reasonable way to estimate the frequency of a sine wave.

Figure 1: Text example 1.



Figure 2: Discourse structure for the text example 1.
 (This structure can be represented as the form
 [[[1 <EX> [2 <EX> 3]] <EG> [4 <EX> 5]] <SR> 6].)

Thinking-flow is defined by a sequence of connective relationships and the way in which the sequence fits into the allowable structure. The authors have investigated all 324 (18×18) pairs of connective relationships and derived possible local structures for thinking-flow restrictions. The pairs of connective relationships can be represented by (r_1, r_2) , where the relations r_1 and r_2 are arbitrary connective relationships. They can be classified into the following four major groups.

- (1) POP-type : permitting $[[P \ r_1 \ Q] \ r_2 \ R]$
 (eliminating $[P \ r_1 \ [Q \ r_2 \ R]]$)
 ex. $[[P \ <EG> \ Q] \ <SR> \ R]$,
 <EG> : exemplification,
 <SR> : serial.
- (2) PUSH-type : permitting $[P \ r_1 \ [Q \ r_2 \ R]]$
 ex. $[P \ <RS> \ [Q \ <SR> \ R]]$,
 <RS> : reason.

- (3) NEUTRAL-type : permitting both (1) and (2)
 ex. $[[P \ <PA> \ Q] \ <EG> \ R]$,
 $[P \ <PA> \ [Q \ <EG> \ R]]$,
 <PA> : parallel.
- (4) NON-type : permitting non-structure
 $[P \ r_1 \ Q \ r_2 \ R]$
 ex. $[P \ <PA> \ Q \ <PA> \ R]$.

The relationship sequence of POP-type means that the local structure for the first two blocks should be popped up, because the local argument is closed. On the other hand, the relationship sequence of PUSH-type means that the local structure should be pushed down.

The relationship sequence of NON-type permits non-structure, which is of the form $[P \ r_1 \ Q \ r_2 \ R]$. Therefore, to be exact, the discourse structure which contains the sequence of this type is not a binary tree.

The thinking-flow restrictions can be used to eliminate structures expressing unnatural argumentative extensions, by examining their local structures. Although the thinking-flow restrictions define local constraints on relationships to neighbors, the scope of relationships is analyzed by recursively checking all local structures of a discourse structure.

2.3 Distant dependencies

The greater part of text can be appropriately analyzed, using the above local constraints on connective relationships to neighbors, if the relationships are extracted correctly. However, in real text, there are rhetorical dependencies concerning distant sentences, which cannot be detected by examining only the normal relationships to neighbors. Two kinds of linguistic clues to distant dependencies must be considered in the realization of a precise discourse analyzer: rhetorical expressions which cover distant sentences, and referential relations of words, in particular, *topics*.

2.3.1 Rhetorical expressions stating global structure

First, rhetorical expressions which relate to an entire article play an important role. Examples are :

- "...? ...? The reason is, ...",
 "... as follows. ... (TENSE=present).
 ... (TENSE=present).",
 "... is not an exceptional case."

Consider the text example in Figure 3, in which unnecessary words are omitted for expositional clarity. In this text the rhetorical expressions which relate to the entire paragraph affect its discourse structure. The expressions "first" and "second" in the last two sentences correspond to the expression "two pieces" in the first sentence; the second and the third sentences, therefore, can be said to be connected by parallel relationship, as they have similar relations with the first sentence. Thus, the discourse structure in Figure 4 is a natural representation.

While, in real text, there is a wide variety of rhetorical expressions of this type, those that are often used in argumentative articles can be determined through analysis. A robust discourse analysis system must detect these rhetorical expressions to restrict discourse structures.

2.3.2 Topic flow

The other significant phenomenon concerning the distant dependencies is *reference*. While English uses pronouns and definite noun phrases in reference, in Japanese, a phrase that is identical to or a part of the original noun phrase is used when referring to some other part of the text. By analyzing the appearance of the same expressions, a restriction or a preference for building discourse structures can be determined. However, the same expressions tend to scatter in a text, and it is difficult to determine the referent for a reference without task-dependent knowledge [Sumita *et al.* 1991]. The author's aim is to create a system not dependent on such extralinguistic knowledge; the reappearance of certain expressions is used as a preference for structure determination.

Figure 5 shows a text example in Japanese, where the underlined words are the same expressions. Note that many underlined words are followed by the character "は (*wa*)". This character is a postpositional particle topicalizing the preceding noun in a sentence.

A *topic* of a sentence is an object indicating what the sentence is about; it can localize the reader's attention in the area that the object relates to. In contrast to topic processing for English (cf. [Schank 1977], [Sidner 1983]), we can use a linguistic device to extract topics for Japanese; some postpositional words are said to indicate a topic of a sentence [Nagano 1986].

In this paper, topic information is used for preference judgment of discourse structures, but not as an element of the structures. To simplify explanation, let us denote a topic of the sentence Q by T^Q , and a case where T^Q refers to a word in the previous sentence P by $T^Q \Rightarrow P$. In the case of the text shown in Figure 5, $T^2 \Rightarrow 1$, $T^3 \Rightarrow 2$, and $T^4 \Rightarrow 3$ hold. If a topic in a sentence refers to a word in the previous sentence, it is regarded as an elaboration of the earlier sentence. Thus, these sentences must be kept close together in their discourse structure; the structure depicted in Figure 6 is appropriate for this text.

In addition, relative importance of relationship connecting sentences in text must be considered for the topic flow analysis. Connective relationships can be classified into three categories according to their relative importance: left-hand, right-hand, and neutral type. For example, the exemplification relationship is a left-hand type; i.e., for $[P \langle EG \rangle Q]$, P strongly relates to the global flow of argumentation beyond the outside of this block, and in this sense P is more important than Q . In contrast, the serial relationship is a right-hand type, and the parallel relationship is a neutral type.

Consider the structure $[[P \ r1 \ Q] \ r2 \ R]$, where 'r1' is a left-hand type relationship, and 'r2' can be any relationship. If $T^R \Rightarrow P$, the above structure is natural, even if there is the same word as T^R in Q . However, if $T^R \Rightarrow Q$, this structure is unnatural, in the sense of coherency. In this case, the structure $[P \ r1 \ [Q \ r2 \ R]]$ is preferable to $[[P \ r1 \ Q] \ r2 \ R]$.

On the contrary, in the case where 'r1' is a right-hand type, $[[P \ r1 \ Q] \ r2 \ R]$ is a natural structure, even if $T^R \Rightarrow Q$. In short, the naturalness of a discourse structure closely depends on the appearance position of topics and their referents, and the relative importance of the referred nodes.

- 1 : Two pieces of X are relevant.
 2 : First, ...
 3 : Second, ...

Figure 3: Text example 2 (X is a noun phrase.)

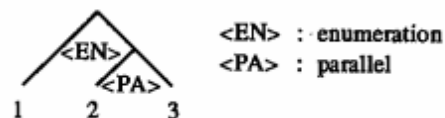


Figure 4: Discourse structure for the text example 2.

- 1 : AはBとCからなる。
 A wa B to C kara naru
 A consists of B and C.
 2 : Cは... DとEに分けられる。
 C wa ... D to E ni wakerareru
 C is divided into D and E.
 3 : Dは... Fを持つ。
 D wa ... F wo motsu
 D has ... F.
 4 : Fは...
 F wa ...
 F is ...

Figure 5: Text example 3 (A - F are noun phrases.)

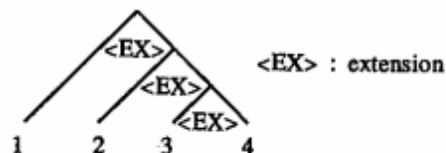


Figure 6: Discourse structure for the text example 3.

3 Discourse structure analyzer

3.1 System configuration

Figure 7 shows the discourse structure analyzer, which consists of five parts: pre-processing, segmentation, candidate generation, candidate reduction and preference judgement. If input text consists of multiple paragraphs or multiple sections, every section or every paragraph in the text is analyzed individually. Figure 8 outlines the input/output data of each stage for a paragraph. The outline of each stage of the discourse structure analyzer is described in the following sections.

3.1.1 Pre-processing

In this stage, input sentences are analyzed, character strings are divided into words, and the dependency structure for each sentence is constructed. The stage consists of the following sub-processes :

- (1) Extracting the text of an article from chapters or sections.
- (2) Accomplishing morphological and syntactic analysis.
- (3) Extracting topic expressions and the reappearance of the targeted expression.
- (4) Detecting connective relationships and constructing their sequence.

In Step (1), the title of an article is eliminated, and the body is extracted. Next, in Step (2), sentences in the body of the article, extracted in Step (1), are morphologically and syntactically analyzed. In Step (3), topic expressions are extracted, according to a table of topic denotation expressions. The following are examples of topic expressions.

"... wa" (as for ...),
 "... niwa" (in ...),
 "... dewa" (in ...),
 "... noitewa" (in ...).

In Step (4), a connective expression is detected based on an expression table consisting of a word and its part of speech for individual connective relationships. In this step, *connection sequence*, a sequence of sentence identifiers and connective relationships, is acquired. For example, a connection sequence is of the form

[1 <EN> 2 <EX> 3 <EX> 4 <EN> 5 <SR> 6],

as is shown as the final result in Figure 8.

3.1.2 Segmentation

In this stage, rhetorical expressions between distant sentences, which define discourse structure, are detected. They form restrictions on segmentation of text.

This stage is implemented as a rule-based proce-

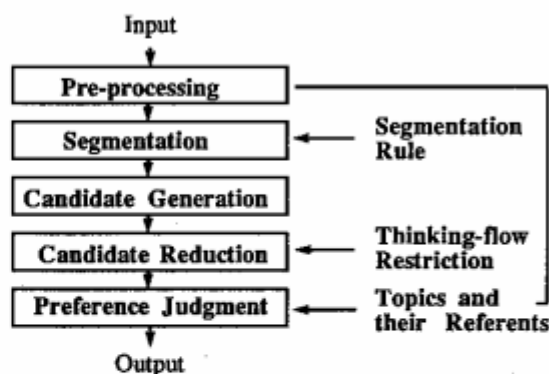


Figure 7: System overview.

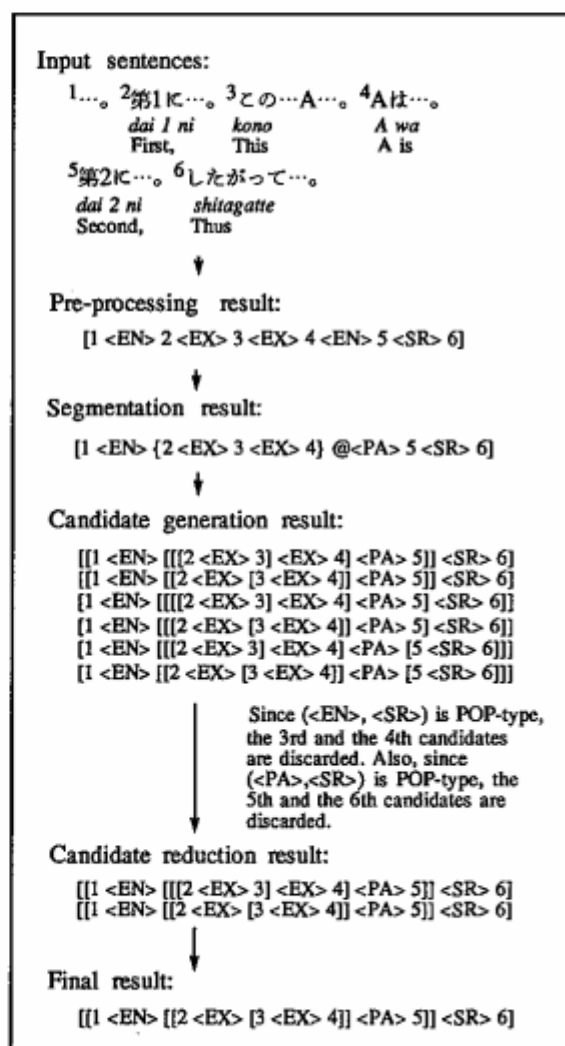


Figure 8: Output example of each process.

ture [Ono *et al.* 1991]. *If-then* rules, called *segmentation rules*, have been formulated in advance. The *if*-part of a segmentation rule corresponds to linguistic surface patterns to detect inter-sentence rhetorical expressions, e.g. "as follows. ... First ... Second ...". The *then*-part represents a connection sequence embedded with control operators discussed below. Also, the *then*-part can indicate an exchange of connective relationships.

There are three kinds of control operators. They are '{' and '}', '(' and ')', and '@'. Sentences enclosed by '{' and '}' must be grouped together as a block of sentences. Operators '(' and ')' are similar to '{' and '}'. They can be used singly, while the operators '{' and '}' must be used in pairs. The operator '@' means that the position must not be a boundary of a sentence block.

Figure 9 shows examples of the segmentation rules. The first example means that if the Nth sentence includes expression "tashikani" (of course), and the Mth sentence includes expression "shikashi" (though), then from N+1st to M-1st sentences must be grouped together.

For the input sentences and the connection sequence in Figure 8, the second rule is activated. The connection sequence is then converted into

[1 <EN> {2 <EX> 3 <EX> 4}@<PA> 5 <SR> 6].

This structure directs the next stage to generate discourse structure candidates whose second, third and fourth sentences are grouped into a block.

At present, approximately 100 rules are available in the system.

3.1.3 Candidate generation

All possible discourse structures, described by binary-trees which do not violate segmentation restrictions, are generated as discourse structure candidates. The generation is performed in a bottom-up manner of sentence parsing by the CYK algorithm. After the generation of sub-trees for blocks directed by segmentation restrictions, the whole trees are generated based on these sub-trees. In case of the example in Figure 8, only 6 candidates are generated, while 42 binary trees would be produced without the segmentation rules.

3.1.4 Candidate reduction

Local structures of generated structure candidates are checked by inspecting thinking-flow restrictions. The candidates including a local structure violating the restrictions are discarded. Only legal candidates are passed on to the next stage.

In order to show the effectiveness of the thinking-flow restrictions, consider the following connection sequence.

[1 <EX> 2 <EG> 3 <PA> 4 <SR> 5].

Figure 10 shows discourse structure candidates for the

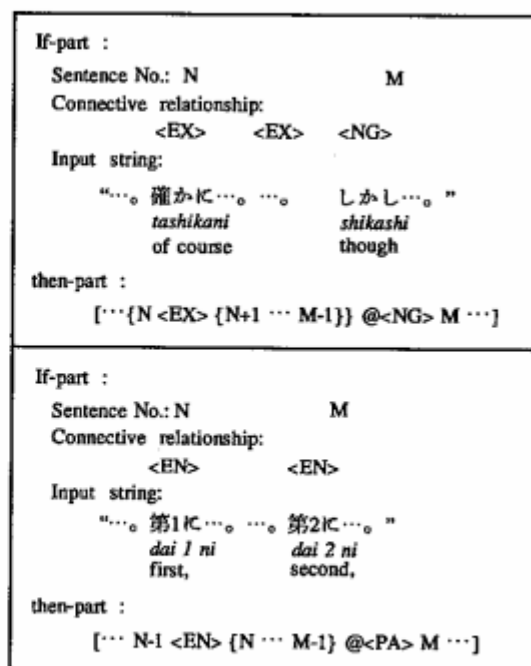


Figure 9: Segmentation rule example.

above sequence. There are 14 binary tree possibilities.

The candidates violating the thinking-flow restrictions are eliminated. For example, the first structure is discarded, because it contains the local structure [2 <EG> [3 <PA> 4]], and the pair (<EG>, <PA>) is POP-type. For the same reason, the seventh structure is also eliminated. This local structure would also be discarded after the exemplification relationship ("<EG> [[3 <PA>]"). As a result of elimination through thinking-flow restrictions, 11 candidates can be discarded, and the second, the fourth and the tenth structures remain.

In the above example, in the case outlined in Figure 8, structure candidates unnatural from the viewpoint of thinking-flow are discarded. Since the third through sixth candidates violate thinking-flow restrictions, the candidates are reduced to two structures.

The thinking-flow restrictions are represented in the system as a table of the applicable pairs of consecutive relationships and their acceptable local structures.

3.1.5 Preference judgment

The final result of discourse analysis is the structure with the lowest *penalty score*, a value associated with topic-referent relationships.

A penalty is set against each arc of path on a discourse structure, which leads from a sentence containing a topic to a sentence referred to by the topic. The concrete arc of a discourse structure, on which a penalty is imposed, is either an arc to or from an unimportant node or an arc to an equally important node. For example, for the structure [[P <EG> Q] <EX> R] where $T^R \Rightarrow Q$,

1: [[1 <EX> [2 <EG> [3 <PA> 4]]] <SR> 5] : NG ((<EG>,<PA>):POP, "<EG> [3 <PA>" :NG)
2: [[1 <EX> [[2 <EG> 3] <PA> 4]] <SR> 5]
3: [[[1 <EX> 2] <EG> [3 <PA> 4]] <SR> 5] : NG ((<EG>,<PA>):POP, "<EG> [3 <PA>" :NG)
4: [[[1 <EX> [2 <EG> 3]] <PA> 4] <SR> 5]
5: [[[[1 <EX> 2] <EG> 3] <PA> 4] <SR> 5] : NG ((<EX>,<EG>):PUSH, "<EX> 2] <EG>" :NG)
6: [1 <EX> [2 <EG> [3 <PA> [4 <SR> 5]]]] : NG ((<PA>,<SR>):POP, "<PA> [4 <SR>" :NG)
7: [1 <EX> [2 <EG> [[3 <PA> 4] <SR> 5]]] : NG ((<EG>,<PA>):POP, "<EG> [[3 <PA>" :NG)
8: [1 <EX> [[2 <EG> 3] <PA> [4 <SR> 5]]] : NG ((<PA>,<SR>):POP, "<PA> [4 <SR>" :NG)
9: [1 <EX> [[2 <EG> [3 <PA> 4]] <SR> 5]] : NG ((<EG>,<PA>):POP, "<EG> [3 <PA>" :NG)
10: [1 <EX> [[[2 <EG> 3] <PA> 4] <SR> 5]]]
11: [[1 <EX> 2] <EG> [3 <PA> [4 <SR> 5]]] : NG ((<PA>,<SR>):POP, "<PA> [4 <SR>" :NG)
12: [[[1 <EX> 2] <EG> [[3 <PA> 4] <SR> 5]]] : NG ((<EG>,<PA>):POP, "<EG> [[3 <PA>" :NG)
13: [[[1 <EX> [2 <EG> 3]] <PA> [4 <SR> 5]]] : NG ((<PA>,<SR>):POP, "<PA> [4 <SR>" :NG)
14: [[[[1 <EX> 2] <EG> 3] <PA> [4 <SR> 5]]] : NG ((<PA>,<SR>):POP, "<PA> [4 <SR>" :NG)

Figure 10: Discourse structure candidates.

a penalty is imposed on the arc from the parent node of P and Q to Q because the left node in an exemplification relationship is unimportant.

The penalty of a discourse structure is defined as a sum of penalties for all paths concerning all topics in the paragraph. By selecting the structure candidate with the lowest penalty, the most coherent discourse structure is obtained.

Of the two surviving structures of the candidate reduction process in Figure 8, the second structure is preferable. The difference is the structural relationship between the second and fourth sentences: the local structure for the first candidate is [[2 <EX> 3] <EX> 4], and that for the second candidate is [2 <EX> [3 <EX> 4]]. Since $T^4 \Rightarrow 3$, a penalty is imposed on the first structure, but not on the second structure. As a result, the second structure candidate is chosen.

While every paragraph can be analyzed respectively, a chapter or a section containing multiple paragraphs is analyzed in an analysis manner similar to that of a paragraph. In case of a discourse structure for a chapter or a section, paragraphs rather than sentences are used

as the terminal nodes of the structure. The connective relationship expressed in the first sentence of each paragraph is used for making the connection sequence. After structure candidates are generated based on the connection sequence, candidates unnatural from the viewpoint of thinking-flow are eliminated. Since every paragraph is analyzed into a discourse structure, each node of the discourse structure for a section also forms the discourse structure for the corresponding paragraph.

3.2 Experiment

To evaluate the discourse structure analyzer, 18 journal articles, different from the data used for algorithm development or rule extraction, have been analyzed. The journal used is "Toshiba Review", which publishes short technical papers of three or four pages. An experiment has been carried out on every paragraph. Correct discourse structure for every paragraph was made manually in advance. The system's performance was evaluated by comparing the correct human-produced structures and the structures analyzed by the system,

Table 2 shows analysis results. There are a total of 554 paragraphs. Nearly 50% of them consist of only one sentence and are excluded from consideration. For 114 paragraphs consisting of more than three sentences, a correct analysis was produced for approximately seventy-four percent.

There were 15 errors for all of the processed paragraphs. Most of the errors are due to incorrect detection of relationships (60%), or incorrect candidate reduction (27%). For the former, the procedure failed to detect explicit connective expressions because of insufficient dictionary data, which can be improved by refining the dictionary data. Most of the latter type of errors occur in a paragraph in which the first or last sentence refers to information outside of the paragraph by such phrases as "as shown above" or "as follows." This suggests that the procedure should also take into account relationships to

Table 2: Analysis results

paragraph size (number of sentences)	correct* (unique)	correct* (other candidate)	incorrect*	Total*
1	-	-	-	293
2	-	-	-	147
3	53	8	6	67
4	12	5	7	24
5	7	1	2	10
6	3	0	0	3
7	5	0	0	6
8	2	0	0	2
9	2	0	0	2
Total	84	14	15	114 ⁺ (554)

* Numbers indicate counts of paragraphs, except for the paragraph size.

+ Total number of paragraphs consisting of more than 3 sentences.

neighboring paragraphs.

In the segmentation stage segmentation rules were activated for 35 paragraphs, with 85% of the rules correctly used; 65% have contributed to structure determination for itemized parts of text, and 20% to relationship determination. In addition, the preference judgment stage has increased the accuracy of the analysis by 3%. Except for the effects of these contributions, correct relationships have been detected in 73 paragraphs, and correct results have been obtained for 55 paragraphs. Thus, if correct connective relationships are detected, 73% of discourse structures can be appropriately analyzed using thinking-flow restrictions only.

4 Concluding remarks

A practical analyzer has been described for building discourse structures for Japanese argumentative or explanatory articles. To analyze structures, three types of knowledge are used: thinking-flow restrictions, segmentation rules, and topic-flow preference. They represent relative constraints between connective relationships or structural restrictions spanning a paragraph, as opposed to the relative importance between consecutive sentences on which other discourse structure analysis researchers depend. Using linguistic knowledge, global structures or the scope of relationships can be determined appropriately.

In addition, the above knowledge on which the procedure is based is detected from superficial linguistic clues independent of topic areas in analyzed articles. The authors are convinced that the method is effective for any articles whose aim is persuasion or assertion.

It should be noted that the relative importance of sentences can be evaluated, using the extracted discourse structure. For example, a left-hand node of a structure linked by exemplification relationship is more important than the right-hand node, as discussed in Section 2.3.2. By a recursive application of relative importance judgment from the top node of discourse structure analyzed from a paragraph, the key-sentence in the paragraph can be extracted.

In addition to the key-sentence extraction shown above, the extracted structure can be a promising clue to other various natural language processes, such as topic estimation and knowledge extraction. The authors intend to polish up the presented restrictions and rules, and refine the procedure toward these natural language processes.

References

- [Cohen 1987] Cohen, R.: "Analyzing the Structure of Argumentative Discourse", *Computational Linguistics*, Vol.13, 1987, pp.11-24.
- [Grosz and Sidner 1986] Grosz, B.J. and Sidner, C.L.: "Attention, Intentions and the Structure of Discourse", *Computational Linguistics*, Vol.12, 1986, pp.175-204.
- [Hobbs 1979] Hobbs, J.R.: "Coherence and Coreference", *Cognitive Science*, Vol.3, 1979, pp.67-90.
- [Litman and Allen 1987] Litman, D.J. and Allen, J.F.: "A Plan Recognition Model for Subdialogues in Conversations", *Cognitive Science*, Vol.11, 1987, pp.163-200.
- [Mann and Thompson 1987] Mann, W.C. and Thompson, S.A.: "Rhetorical Structure Theory: A Framework for the Analysis of Texts", *USC/Information Science Institute Research Report RR-87-190*, 1987.
- [Nagano 1986] Nagano, K.: *Bunshouron Sousetsu — Bunpouron-teki Kousatsu— (An Introduction to Theory of Texts —Syntactic Consideration—)*, Asakusa Shoten, 1986, (in Japanese).
- [Ono et al. 1989] Ono, K., Ukita, T., and Amano, S.: "An Analysis of Rhetorical Structure", *IPS Japan Technical Report NL 70-2*, 1989, (in Japanese).
- [Ono et al. 1991] Ono, K., Sumita, K., Ukita, T., and Amano, S.: "Text Segmentation and Discourse Analysis", *Proc. IPS Japan '91 October*, 4E-2, 1991, (in Japanese).
- [Schank 1977] Schank, R.C.: "Rules and Topics in Conversation", *Cognitive Science*, Vol.1, 1977, pp.421-441.
- [Sidner 1983] Sidner, C.L.: "Focusing in Comprehension of Definite Anaphora", M.Brady and R.C.Berwick (Eds.), *Computational Models of Discourse*, MIT Press, 1983, pp.267-330.
- [Sumita et al. 1991] Sumita, K., Ukita, T., and Amano, S.: "Disambiguation in Natural Language Interpretation Based on Amount of Information", *IEICE Trans.*, Vol.E74, No.6, 1991, pp.1735-1746.