# STRUCTURAL MATCHING OF PARALLEL TEXTS

Yuji Matsumoto

Graduate School of Information Science
Advanced Institute of Science and Technology, Nara
Takayama-cho, Ikoma-shi, Nara 630-01 Japan
matsu@is.aist-nara.ac.jp

Hiroyuki Ishimoto      Takehito Utsuro

Department of Electrical Engineering
Kyoto University
Sakyo-ku, Kyoto 606 Japan
{ishimoto, utsuro}@pine.kuee.kyoto-u.ac.jp

## Abstract

This paper describes a method for finding structural matching between parallel sentences of two languages, (such as Japanese and English). Parallel sentences are analyzed based on unification grammars, and structural matching is performed by making use of a similarity measure of word pairs in the two languages. Syntactic ambiguities are resolved simultaneously in the matching process. The results serve as a useful source for extracting linguistic and lexical knowledge.

## INTRODUCTION

Bilingual (or parallel) texts are useful resources for acquisition of linguistic knowledge as well as for applications such as machine translation. Intensive research has been done for aligning bilingual texts at the sentence level using statistical techniques by measuring sentence lengths in words or in characters (Brown 91), (Gale 91a). Those works are quite successful in that far more than 90% of sentences in bilingual corpora are aligned correctly.

Although such parallel texts are shown to be useful in real applications such as machine translation (Brown 90) and word sense disambiguation (Dagan 91), structured bilingual sentences are undoubtedly more informative and important for future natural language researches. Structured bilingual or multilingual corpora serve as richer sources for extracting linguistic knowledge (Kaji 92), (Klavans 90), (Sadler 91), (Utsuro 92).

Phrase level or word level alignment has also been done by several researchers. The Textual Knowledge Bank Project (Sadler 91) is building monolingual and multilingual text bases structured by linking the elements with grammatical (dependency), referential, and bilingual relations. (Kaji 92) reports a method to obtain phrase level correspondence of parallel texts by coupling phrases of two languages obtained in CKY parsing processes.

This paper presents another method to obtain structural matching of bilingual texts. Sentences in both languages are parsed to produce (disjunctive) feature structures, from which dependency structures are extracted. Ambiguities are represented as disjunction. Then, the two structures are matched to establish a one-to-one correspondence between their substructures. The result of the match is obtained as a set of pairs of minimal corresponding substructures of the dependency structures. Examples of the results are shown in Figures 1, 2 and 3. A dependency structure is represented as a tree, in which ambiguity is specified by a disjunctive node (OR node). Circles in the figure show substructures and bidirectional arrows show corresponding substructures.

Our technique and the results are different from those of other methods mentioned above. (Kaji 92) identifies corresponding phrases and aims at producing translation templates by abstracting those corresponding phrases. In the Bilingual Knowledge Bank (Sadler 91), the correspondence is shown by

23

links between words in two sentences, equating two whole subtrees headed by the words. We prefer the minimal substructure correspondence and the relationship between substructures. Such a minimal substructure stands for the minimal meaningful component in the sentence, which we believe is very useful for our target application of extracting lexical knowledge from bilingual corpora.

# SPECIFICATION OF STRUCTURAL MATCHING PROBLEM

Although the structural matching method shown in this paper is language independent, we deal with parallel texts of Japanese and English. We assume that alignment at the sentence level is already pre-processed manually or by other methods such as those in (Brown 91), (Gale 91a). Throughout this paper, we assume to match simple sentences.[1]

## DEFINITIONS OF DATA STRUCTURES

A pair of Japanese and English sentences are parsed independently into (disjunctive) feature structures. For our present purpose, a part of a feature structure is taken out as a dependency structure consisting of the content words[2] that appear in the original sentence. Ambiguity is represented by disjunctive feature structures (Kasper 87). Since any relation other than modifier-modifyee dependencies is not considered here, path equivalence is not taken into consideration. Both of value disjunction and general disjunction are allowed.

We are currently using LFG-like grammars for both Japanese and English, where the value of the 'pred' label in an f-structure is the content word that is the head of the corresponding c-structure.

We start with the definitions of simplified disjunctive feature structures, and then disjunctive dependency structures, that are extracted from the disjunctive feature structures obtained by the parsing process.

**Definition 1** *Simple* feature structures (FS) *(L is the set of feature labels, and A is the set of atomic values) are defined recursively:*

$NIL$

| | |
|---|---|
| $a$ | *where $a \in A$* |
| $l : \phi$ | *where $l \in L$, $\phi \in$ FS* |
| $\phi \wedge \psi$ | *where $\phi, \psi \in$ FS* |
| $\phi \vee \psi$ | *where $\phi, \psi \in$ FS* |

To define *(Disjunctive) Dependency Structures* as a special case of an FS, we first require the following definitions.

**Definition 2** Top label set *of an FS $\phi$, written as $tl(\phi)$, is defined:*

1. *If $\phi = l : \phi_1$, then $tl(\phi) = \{l\}$.*

2. *If $\phi = \phi_1 \wedge \phi_2$ or $\phi = \phi_1 \vee \phi_2$, then $tl(\phi) = tl(\phi_1) \cup tl(\phi_2)$.*

**Definition 3** *A relation 'sibling' between feature labels in $\phi$ is defined:*

1. *If $\phi = l : \phi_1$, then $l$ and labels in $\phi_1$ are not sibling, and sibling relation holding in $\phi_1$ also holds in $\phi$.*

2. *If $\phi = \phi_1 \wedge \phi_2$, then labels in $tl(\phi_1)$ and labels in $tl(\phi_2)$ are sibling.*

3. *If $\phi = \phi_1 \vee \phi_2$, then labels in $\phi_1$ and labels in $\phi_2$ are not sibling.*

Note that the sibling relation is not an equivalence relation. We refer to a set of feature labels in $\phi$ that are mutually sibling as a *sibling label set* of $\phi$. Now, we are ready to define a dependency structure (DS).

**Definition 4** *A dependency structure $\overline{\psi}$ is an FS that satisfies the following condition:*

Condition: *Every* sibling label set *of $\psi$ includes exactly one 'pred' label.*

The idea behind those are that the value of a 'pred' label is a content word appearing in the original sentence, and that a sibling label set defines the dependency relation between content words. Among the labels in a sibling label set, the values of the labels other than 'pred' are dependent on (i.e., modify) the value of the 'pred' label. A DS can be drawn as a tree structure where the nodes are either a content word or disjunction operator and the edges represent the dependency relation.

**Definition 5** *A substructure of an FS $\phi$ is defined (sub($\phi$) stands for the set of all substructures of $\phi$.):*

1. *$NIL$ and $\phi$ itself are substructures of $\phi$.*

2. *If $\phi = a$ ($a \in A$), then $a$ is a substructure of $\phi$.*

---

[1] Matching of compound sentences are done by cutting them up into simple sentence fragments.

[2] In the present system, nouns, pronouns, verbs, adjectives, and adverbs are regarded as content words.

24

English: She has long hair.
Japanese: 彼女 - の　　髪 - は　　長い
she - GEN　hair - TOP　long
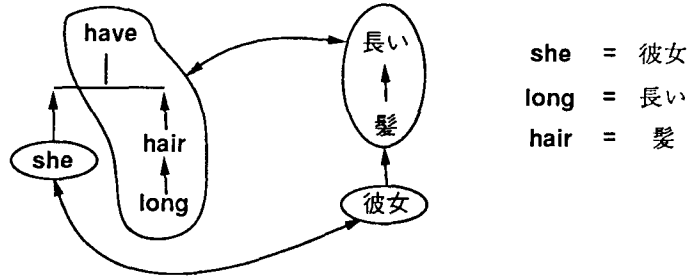
she　=　彼女
long　=　長い
hair　=　髪



Figure 1: Example of structural matching, No.1

English: This child is starving for parental love.
Japanese: この　子 - は　　親 - の　　愛 - に　　飢えている
this　child - TOP　parent - GEN　love - DAT　be-starving
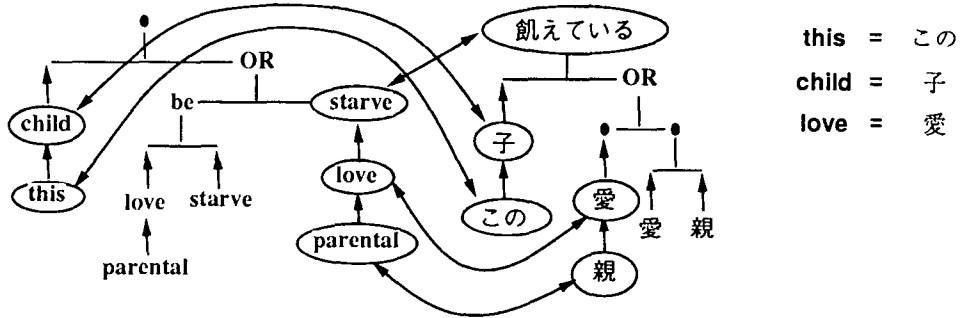
this　=　この
child　=　子
love　=　愛



Figure 2: Example of structural matching, No.2

English: Japan benefits from free trade.
Japanese: 日本 - は　　自由貿易 - の　　恩恵 - を　　受ける
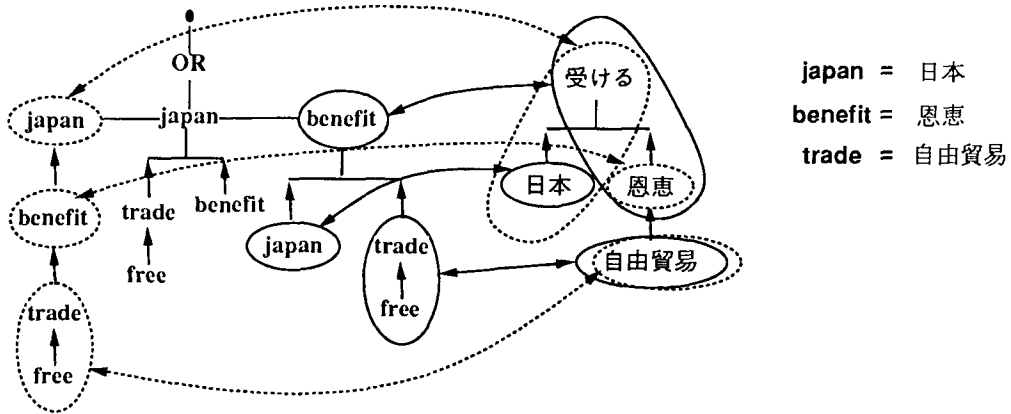Japan - TOP　free-trade - GEN　benefit - ACC　receive

japan　=　日本
benefit　=　恩恵
trade　=　自由貿易



Figure 3: Example of structural matching, No.3

3. If $\phi = l : \phi_1$, then $sub(\phi_1)$ are substructures of $\phi$.

4. If $\phi = \phi_1 \wedge \phi_2$, then for any $\psi_1 \in sub(\phi_1)$ and for any $\psi_2 \in sub(\phi_2)$, $\psi_1 \wedge \psi_2$ is a substructure of $\phi$.

5. If $\phi = \phi_1 \vee \phi_2$, then for for any $\psi_1 \in sub(\phi_1)$ and for any $\psi_2 \in sub(\phi_2)$, $\psi_1 \vee \psi_2$ is a substructure of $\phi$.

The DS derived from an FS is the maximum substructure of the FS that satisfies the condition in Definition 4. The DS is uniquely determined from an FS.

**Definition 6** *A disjunction-free maximal substructure of an FS $\phi$ is called a* complete FS *of $\phi$.*

An FS does not usually have a unique complete FS. This concept is important since the selection of a complete FS corresponds to ambiguity resolution. Naturally, a maximal disjunction-free substructure of a DS $\psi$ is again a DS and is called a *complete DS* of $\psi$.

**Definition 7** *A semi-complete DS of a DS $\psi$ is a substructure of a complete DS of $\psi$ that satisfies the condition in Definition 4.*

Note that a substructure of a DS is not necessarily a DS. This is why the definition requires the condition in Definition 4.

A complete DS $\psi$ can be decomposed into a set of non-overlapping semi-complete DSs. Such a decomposition defines the units of structural matching and plays the key role in our problem.

**Definition 8** *A set of semi-complete DS of a DS $\psi$, $D = \{\psi_1, \cdots \psi_n\}$, is called a decomposition of $\psi$, iff every $\psi_i$ in the set contains at least one occurrence of 'pred' feature label, and every content word at the 'pred' feature label appearing in $\psi$ is contained in exactly one $\psi_i$.*

**Definition 9** *The reduced DS of a DS $\psi$ with respect to a decomposition $D = \{\psi_1, \cdots \psi_n\}$ is constructed as follows:*

1. *$\psi_i$ is transformed to a DS, 'pred : $S_i$', where $S_i$ is the set of all content words appearing in $\psi_i$. This DS is referred to as $red(\psi_i)$.*

2. *If there is a direct dependency relation between two content words $w_1$ and $w_2$ that are in $\psi_i$ and $\psi_j$ ($i \neq j$), then the dependency relation is allotted between $\psi_i$ and $\psi_j$.*

Although this definition should be described precisely, we leave it with this more intuitive description. Examples of dependency structures and reduced dependency structures are found in Figures 1, 2 and 3, where the decompositions are indicated by circles.

It is not difficult to show that the reduced DS satisfies the condition of Definition 4.

## STRUCTURAL MATCHING OF BILINGUAL DEPENDENCY STRUCTURES

Structural matching problem of bilingual sentences is now defined formally.

Parsing parallel English and Japanese sentences results in feature structures, from which dependency structures are derived by removing unrelated features.

Assume that $\psi_E$ and $\psi_J$ are dependency structures of English and Japanese sentences. The structural matching is to find *the most plausible one-to-one mapping* between a decomposition of a complete DS of $\psi_E$ and a decomposition of a complete DS of $\psi_J$, provided that the reduced DS of $\psi_E$ and the reduced DS of $\psi_J$ w.r.t. the decompositions are isomorphic over the dependency relation. The isomorphism imposes a natural one-to-one correspondence on the dependency relations between the reduced DSs.

Generally, the mapping need not always be one-to-one, i.e., all elements in a decomposition need not map into another decomposition. When the mapping is not one-to-one, we assume that dummy nodes are inserted in the dependency structures so that the mapping naturally extends to be one-to-one.

When the decompositions of parallel sentences have such an isomorphic one-to-one mapping, we assume that there are systematic methods to compute *similarity* between corresponding elements in the decompositions and to compute *similarity* between the corresponding dependency relations[3].

We write the function defining the former similarity as $f$, and that of the latter as $g$. Then, $f$ is a function over semi-complete DSs derived from English and Japanese parallel sentences into a real number, and $g$ is a function over feature label sets

---

[3] in the case of similarity between dependency relations, the original feature labels are taken into account.

of English and Japanese into a real number.

**Definition 10** *Given dependency structures, $DS_1$ and $DS_2$, of two languages, the structural matching problem is to find an isomorphic one-to-one mapping $m$ between decompositions of $DS_1$ and $DS_2$ that maximizes the sum of the values of similarity functions, $f$ and $g$.*

*That is, the problem is to find the function $m$ that maximizes*

$$\sum_d (f(d, m(d)) + \sum_l g(l, m(l)))$$

*where $d$ varies over semi-complete DS of $DS_1$ and $l$ varies over feature labels in $DS_1$.*

The similarity functions can be defined in various ways. We assume some similarity measure between Japanese and English words. For instance, we assume that the similarity function $f$ satisfies the following principles:

1. $f$ is a simple function defined by the similarity measure between content words of two languages.
2. Fine-grained decompositions get larger similarity measure than coarse-grained decompositions.
3. Dummy nodes should give some negative value to $f$.

The first principle is to simplify the complexity of the structural matching algorithm. The second is to obtain detailed structural matching between parallel sentences and to avoid trivial results, e.g., the whole DSs are matched. The third is to avoid the introduction of dummy nodes when it is possible.

The function $g$ should be defined according to the language pair. Although feature labels represent grammatical relation between content words or phrases and may provide useful information for measuring similarity, we do not use the information at our current stage. The reason is that we found it difficult to have a clear view on the relationship between feature labels of English and Japanese and on the meaning of feature labels between semi-complete dependency structures.

## STRUCTURAL MATCHING ALGORITHM

The structural matching of two dependency structures are combinatorially difficult problem. We apply the branch-and-bound method to solve the problem.

The branch-and-bound algorithm is a top-down depth-first backtracking algorithm for search problems. It looks for the answers with the BEST score. In each new step, it estimates the maximum value of the expected scores along the current path and compares it with the currently known best score. The maximum expected score is usually calculated by a simplified problem that guarantees to give a value not less than the best score attainable along the current path. If the maximum expectation is less than the currently known best score, it means that there is no chance to find better answers by pursuing the path. Then, it gives up the current path and backtracks to try remaining paths.

We regard a dependency structure as a tree structure that includes disjunction (OR nodes), and call a content word and a dependency relation as a node and an edge, respectively. Then a semi-complete dependency structure corresponds to a connected subgraph in the tree.

The matching of two dependency trees starts from the top nodes and the matching process goes along edges of the trees. During the matching process, three types of nondeterminism arise:

1. Selection of top-most subgraphs in both of the trees (i.e., selection of a semi-complete DS)
2. Selection of edges in both of the trees to decide the correspondence of dependency relations
3. Selection of one of the disjuncts at an 'OR' node

While the matching is done top-down, the exact score of the matched subgraphs is calculated using the similarity function $f$.[4] When the matching process proceeds to the selection of the second type, it selects an edge in each of the dependency trees. The maximum expected score of matching the subtrees under the selected edges are calculated from the sets of content words in the subtrees. The calculation method of the maximum expected score is defined in some relation with the similarity function $f$.

Suppose $h$ is the function that gives the maximum expected score of two subgraphs. Also, suppose $B$ and $P$ be the currently known best score

---

[4] We do not take into account the similarity measure between dependency relations as stated in the preceding section.

27

and the total score of the already matched subgraphs, respectively. If $s$ and $t$ are the subgraphs under the selected edges and $s'$ and $t'$ are the whole remaining subgraphs, the matching under $s$ and $t$ will be undertaken further only when the following inequation holds:

$$P + h(s,t) + h(s',t') \geq B$$

Any selection of edges that does not satisfy this inequality cannot provide better matching than the currently known best ones.

All of the three types of nondeterminism are simply treated as the nondeterminism in the algorithm.

The syntactic ambiguities in the dependency structures are resolved spontaneously when the matching with the best score is obtained.

## EXPERIMENTS

We have tested the structural matching algorithm with 82 pairs of sample sentences randomly selected from a Japanese-English dictionary.

We used a machine readable Japanese-English dictionary (Shimizu 79) and Roget's thesaurus (Roget 11) to measure the similarity of pairs of content words, which are used to define the function $f$.

### Similarity of word pairs

Given a pair of Japanese and English sentences, we take two methods to measure the similarity between Japanese and English content words appearing in the sentences.

For each Japanese content word $w_J$ appearing in the Japanese sentence, we can find a set of translatable English words from the Japanese-English dictionary. When the Japanese word is a polysemous word, we select an English word from each polysemous entry. Let $C_{EJ}$ be the set of such translatable English words of $w_J$. Suppose $C_E$ is the set of contents words in the English sentence. The translatable pairs of $w_J$, $Tp(w_J)$, is defined as follows:

$$Tp(w_j) = \{ \langle w_J, w_E \rangle \mid w_E \in C_{EJ} \cap C_E \}$$

We use Roget's thesaurus to measure similarity of other word pairs. Roget's thesaurus is regarded as a tree structure where words are allocated at the leaves of the tree: For each Japanese content word $w_J$ appearing in the Japanese sentence, we can define the set of translatable English words of $w_J$, $C_{EJ}$. From each English word in the set, the minimum distance to each of the English content words

appearing in the English sentence is measured.[5] This minimum distance defines the similarity between pairs of Japanese and English words.

We decided to use this similarity only for estimating dissimilarity between Japanese and English word pairs. We set a predetermined threshold distance. If the minimal distance exceeds the threshold, the exceeded distance is counted as the negative similarity.

The similarity of two words $w_1$ and $w_2$ appearing in the given pair of sentences, $sim(\langle w_1, w_2 \rangle)$, is defined as follows:

$sim(\langle w_1, w_2 \rangle) =$
$$\begin{cases} 6 & \langle w_1, w_2 \rangle \in Tp(w_1) \text{ or } \langle w_2, w_1 \rangle \in Tp(w_2) \\ -k & \langle w_1, w_2 \rangle \notin Tp(w_1) \text{ and } \langle w_2, w_1 \rangle \notin Tp(w_2) \\ & \text{and the distance between } w_1 \text{ and } w_2 \\ & \text{exceeds the threshold by } k. \\ 0 & \text{otherwise} \end{cases}$$

### Similarity of semi-complete DSs

The similarity between corresponding semi-complete DSs is defined based on the similarity between the content words. Suppose that $s$ and $t$ are semi-complete DSs to be matched, and that $V_s$ and $V_t$ are the sets of content words in $s$ and $t$. Let $A$ be the less larger set of $V_s$ and $V_t$ and $B$ be the other ($\mid A \mid \leq \mid B \mid$). For each injection $p$ from $A$ into $B$, the set of word pairs $D$ derived from $p$ can be defined as follows.

$$D = \left\{ \langle a, p(a) \rangle \mid a \in A \right\}$$

Now, we define the similarity function $f$ over Japanese and English semi-complete DSs to give the maximum value to the following expression for all possible injections:

$$f(s,t) = \max_p \left\{ \sum_{d \in D} sim(d) \right\} \times 0.95^{\mid V_s \mid + \mid V_t \mid - 2}$$

The summation gives the maximum sum of the similarity of the content words in $s$ and $t$. 0.95 is the penalty when the semi-complete DSs with more than one content words are used in the matching.

Figures 1, 2 and 3 shows the results of the structural matching algorithm, in which the translatable pairs obtained from the Japanese-English dictionary are shown by the equations.

---

[5] The distance between words is the length of the shortest path in the thesaurus tree.

28

Table 1: Results of experiments

| Parsing Japanese and English sentences | | |
|---|---|---|
| Number of sentences | 82 | |
| Parse failure | 23 | |
| Parsable | 59 | |
| **Correct parsability** | | |
| Correct parse | 53 | 89.8% (53/59) |
| Incorrect parse | 6 | 10.2% (6/59) |
| **The match with the best score includes** | | |
| Correct matching | 47 | 89% (47/53) |
| no correct matching | 6 | 11% (6/53) |
| Single correct matching | 34 | 64% (34/53) |

### Results of the experiments

We used 82 pairs of Japanese and English sentences appearing in a Japanese-English dictionary. The results were checked and examined in detail by hand. Some of the sentences are not parsable because of the limited coverage of our current grammars. Although 59 pairs of them are parsable, 6 out of them do not include correct parse results.

The structural matching algorithm with the setting described above is applied to the 53 pairs. The cases where the correct matching is not included in the best rated answers are 6 out of them. The remaining 47 pairs include the correct matching, of which 31 pairs result in the correct matching uniquely. Table 1 summarizes the results.

## EVALUATION AND DISCUSSION

Although the number of sentences used in the experiments is small, the result shows that about two third of the pairs give the unique matching, in which every syntactic ambiguity is resolved.

The cases where no correct matching was obtained needs be examined. Some sentences contain an idiomatic expression that has completely different syntactic structures from the sentence structure of the other. Such an expression will no way be matched correctly except that the whole structures are matched intact. Other cases are caused by complex sentences that include an embedded sentence. When the verbs at the roots of the dependency trees are irrelevant, extraordinary matchings are produced. We intend not to use our method to match complex or compound sentences as a whole. We will rather use our method to find structural matching between simple sentences or verb phrases of two languages.

The matching problem of complex sentences are regarded as a different problem though the similar technique is usable. We think that the scores of matched phrases will help to identify the corresponding phrases when we match complex sentences.

Taking the sources of other errors into consideration, possible improvements are:

1. Enhancement of English and Japanese grammars for wider coverage and lower error rate.
2. Introduction of more precise similarity measurement of content words.
3. Utilization of grammatical information:
   - Feature labels, for estimating matching plausibility of dependency relations
   - Part of speech, for measuring matching plausibility of content words
   - Other grammatical information: mood, voice, etc.

The first two improvements are undoubtedly important. As for the similarity measurement of content words, completely different approaches such as statistical methods may be useful to get good translatable pairs (Brown 90), (Gale 91).

Various grammatical information is kept in the feature descriptions produced in the parsing process. However, we should be very prudent in using it. Since English and Japanese are grammatically quite different, some grammatical relation may not be preserved between them. In Figure 3, solid arrows and circles show the correct matching. While 'benefit' matches with the structure consisting of ' 恩恵 ' and ' 受ける ', their dependent words 'trade' and ' 自由貿易 ' modify them as a verb modifier and as a noun modifier, the grammatical relation of which are quite different.

This example highlights another interesting point. Dotted arrows and circles show another matching with the same highest score. In this case, 'japan' is taken as a verb. This rather strange interpretation insists that 'japan' matches with ' 日本 ' and ' 受ける '. Since 'japan' as a verb has little semantic relation with ' 日本 ' as a country, discrimination of part-of-speech seems to be useful. On the other hand, the correspondence between 'benefit' and ' 恩恵 ' is found in their noun entry in the dictionary. Since 'benefit' is used as a verb in the

sentence, taking part-of-speech into consideration may jeopardize the correct matching, either. The fact that the verb and noun usages of 'benefit' bear common concept implies that more precise similarity measurement will solve this particular problem. Since the interpretations of the sample English sentences are in different mood, imperative and declarative, the mood of a sentence is also useful to remove irrelevant interpretations.

## CONCLUSIONS

The structural matching problem of parallel texts is formally defined and our current implementation and experiments are introduced. Although the research is at the preliminary stage and has a very simple setting, the experiments have shown a number of interesting results. The method is easily enhanced by improving the grammars and by incorporating more accurate similarity measurement. Number of other researches of building translation dictionaries and of determining similarity relationship between words are useful to improve our method.

To extract useful information from bilingual corpora, structural matching is inevitable for language pairs like English and Japanese that have quite different linguistic structure. Incidentally, we have found that this dissimilarity plays an important role in resolving syntactic ambiguities since the sources of ambiguities in English and Japanese sentences are in many cases do not coincide (Utsuro 92). We are currently working on extracting verbal case frames of Japanese from the results of structural matching of a Japanese-English corpus (Utsuro 93). The same technique is naturally applicable to acquire verbal case frames of English as well. Another application we are envisaging is to extract translation pattern from the results of structural matching.

We plan to work on possible improvements discussed in the preceding section, and will make large scale experiments using translated newspaper articles, based on the phrase matching strategy.

## ACKNOWLEDGMENTS

## REFERENCES

Brown, P.F., et al., A Statistical Approach to Machine Translation, *Computational Linguistics*, Vol.16, No.2, pp.79-85, 1990.

Brown, P.F., Lai, J.C. and Mercer, R.L., Aligning Sentences in Parallel Corpora, *ACL-91*, pp.169-176, 1991.

Dagan, I., Itai, A. and Schwall, U., Two Languages are More Informative than One, *ACL-91*, pp.130-137, 1991a.

Gale. W.A. and Church, K.W., A Program for Aligning Sentences in Bilingual Corpora, *ACL-91*, pp.177-184, 1991b.

Gale. W.A. and Church, K.W., Identifying Word Correspondences in Parallel Texts, *'91 DARPA Speech and Natural Language Workshop*, pp.152-157, 1991.

Kaji, H., Kida, Y., and Morimoto, Y., Learning Translation Templates from Bilingual Text, *COLING-92*, pp.672-678, 1992.

Kasper, R., A Unification Method for Disjunctive Feature Descriptions, *ACL-87*, pp.235-242, 1987.

Klavans, J. and Tzoukermann, E., The BICORD System: Combining Lexical Information from Bilingual Corpora and Machine Readable Dictionaries, *COLING-90*, pp.174-179, 1990.

Miller, G.A., et al., Five Papers on WordNet, *Cognitive Science Laboratory, Princeton University, CSL Report 43*, July 1990.

Roget, S.R., Roget's Thesaurus, *Crowell Co.*, 1911.

Sadler, V., The Textual Knowledge Bank: Design, Construction, Applications, *Proc. International Workshop on Fundamental Research for the Future Generation of Natural Language Processing (FGNLP)*, pp.17-32, Kyoto, Japan, 1991.

Shimizu, M., et al. (ed.), Japanese-English Dictionary, *Kodansha*, 1979.

Utsuro, T., Matsumoto, Y., and Nagao, M., Lexical Knowledge Acquisition from Bilingual Corpora, *COLING-92*, pp.581-587, 1992.

Utsuro, T., Matsumoto, Y., and Nagao, M., Verbal Case Frame Acquisition from Bilingual Corpora, to appear *IJCAI-93*, 1993.