# Verb Paraphrase based on Case Frame Alignment

**Nobuhiro Kaji[†], Daisuke Kawahara[†], Sadao Kurohash[†‡] and Satoshi Sato[*]**

[†]Graduate School of Information Science
and Tech. University of Tokyo
Hongo, Bunkyo-ku,
Tokyo 113-8656, Japan

[‡]PRESTO,
Japan Science and Technology
Corporation (JST)

[*]Graduate School of Informatics
Kyoto University
Yoshida-honmachi, Sakyo-ku,
Kyoto, 606-8501, Japnan

{kaji,kawahara,kuro} @kc.t.u-tokyo.ac.jp
sato@pine.kuee.kyoto-u.ac.jp

## Abstract

This paper describes a method of translating a predicate-argument structure of a verb into that of an equivalent verb, which is a core component of the dictionary-based paraphrasing. Our method grasps several usages of a headword and those of the def-heads as a form of their case frames and aligns those case frames, which means the acquisition of word sense disambiguation rules and the detection of the appropriate equivalent and case marker transformation.

## 1 Introduction

We are conducting a research of automatic paraphrasing using an ordinary dictionary, replacing words in an input with the appropriate expressions in their definitions of a dictionary. Since a dictionary explains a headword by simpler words/expressions, the dictionary-based paraphrasing translates an input sentence into a simpler sentence. This paper describes a core component of the dictionary-based paraphrasing, that is, a method of translating a predicate-argument structure of a verb into that of a simpler equivalent verb. For example, our method can translate "overlook the mistake" into "not notice the mistake".

Paraphrasing has two merits. One is that it leads the solution of the biggest obstacle in natural language processing, synonymy in a wider sense (one to many correspondence from meaning to expression). Synonymy decreases the recall of information retrieval systems. The inference based on text is receiving much atten-tion for information exploitation these days, but the inference stops immediately without proper handling of synonymy. If the dictionary-based paraphrasing can transform any texts into expressions of a basic vocabulary, a part of the synonymy problem is can be handled successfully. Note that, for example, definitions in LDOCE only uses 2,000 words of the Longman defining vocabulary.

The other merit is related to a user interface. Flexible presentation of information according to the user is a challenging topic, and paraphrasing into simpler expressions is one of the key technologies. For example, paraphrasing of news articles into simpler expressions is useful for children or non-native speakers, and paraphrasing of technical terms into ordinal expressions is required in several situations.

The biggest problem in paraphrasing among synonym is, inversely, the polysemy (one to many correspondence from expression to meaning). Suppose a word, $A$, has two meanings, $A_1$ and $A_2$, and equivalent of $A$ in the sense $A_1$ is $B$; equivalent of $A$ in the sense $A_2$ is $C$. An ordinary dictionary exactly provides us with such information. In such a case, the problem is to disambiguate the sense of $A$ in a given context, $A_1$ or $A_2$. This situation is exactly the same as machine translation in which word selection is very problematic (this is natural, since machine translation is a kind of paraphrase).

Even if the disambiguation of $A$ is successful, there are several remaining problems. Since the equivalent of $A$ ($B$ and $C$) is not necessarily a word but a phrase, the proper equivalent should be extracted from the definitions of $A$ in the dictionary. Another problem is that the sentential

pattern of $B$ or $C$ might be different from that of $A$ (in case of Japanese, case markers might change).

This paper proposes a method to solve these problems simultaneously by aligning case frames of A and those of B and C. We utilize wide coverage and specific enough case frames which are automatically constructed from a raw corpus.

## 2 Dictionary based Paraphrase

### 2.1 Basic Idea

An ordinary dictionary provides us with definitions of headwords in simpler words and expressions. In case of verbs, the head of the definition (we call it a *def-head* hereafter) and the adverb modifying the def-head, if any, are an equivalent of the headword and can be used as a paraphrase of it. For example, the definition of *chiratsuku* (shimmer) is as follows:

**chiratsuku (shimmer)** yowaku hikaru (shine faintly).

The def-head, "shine" and its modifier "faintly" can be seen as an equivalent of *chiratsuku* (shimmer), and can be used as a paraphrase as follows:

- The lamp shimmers → The lamp shines faintly

### 2.2 Difficulties

Replacing a headword with the def-head and the adverb, however, does not always work because of the following problems.

#### Word sense ambiguity

When a headword has two or more meanings, the meaning of the headword in a given context must be chosen. For example, *keitou* (devote) has the following two meanings.

**keitou (devote)** **1** necchuu (be enthusiastic). **2** shitau (admire).

If an input is "literature ni keitou-suru (devote oneself to literature)", *keitou* has the first meaning and the input should be paraphrased into "literature ni nechuu-suru"[1][2]. If an input is

---

[1] In Japanese, postpositions function as case markers and a verb is final in a sentence

[2] Nouns in definitions and case frames are given in English translation in this paper.

"Lincoln ni keitou-suru (devote oneself to Lincoln)", *keitou* has the second meaning and the input should be paraphrased into "Lincoln wo shitau (admire Lincoln) ".

#### Size of the equivalent

Sometimes, the equivalent of a headword is larger than the def-head and its modifier. For example, as shown below, the equivalent of "taitoku (acquire)" is not "tsukeru (get)", but "mi ni tsukeru (get for oneself)".

**taitoku (acquire)** knowledge ya skill wo mi ni tsukeru (get knowledge or skill for oneself).

#### Transformation of case marker

When the verb is paraphrased to its equivalent, the sentential pattern (case markers) might have to be modified. For example, when "mistake wo miotosu (overlook a mistake)" is paraphrased by the following definition of "miotosu (overlook)", "mistake" shoule be transformed from wo case to ni case.

**miotosu (overlook)** kizukanai (not notice).

### 2.3 Verb Paraphrase based on Case Frame Alignment

The problems mentioned above can be solved by finding a correspondence between an input and definitions. However, such a correspondence is very hard to find since predicate-argument structure of the def-head is not fully given in the definition, and the input also contains some omissions. That is, the information is too small to find an appropriate correspondence.

For example, as shown above, *keitou* has two meanings, but there is no information about its arguments in the dictionary. It is almost impossible to find which definition should be used in what way to paraphrase "literature ni keitou-suru (devote oneself to literature)".

Then, we have developed a method which grasps several usages of a headword and those of the def-heads as a form of their case frames and aligns those case frames, which means the acquisition of word sense disambiguation rules and the detection of the appropriate equivalent and case marker transformation. Case frames
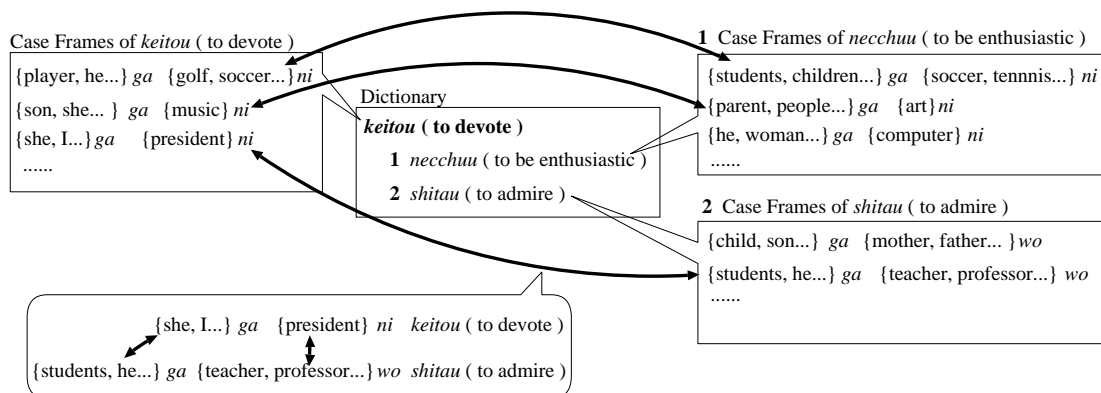
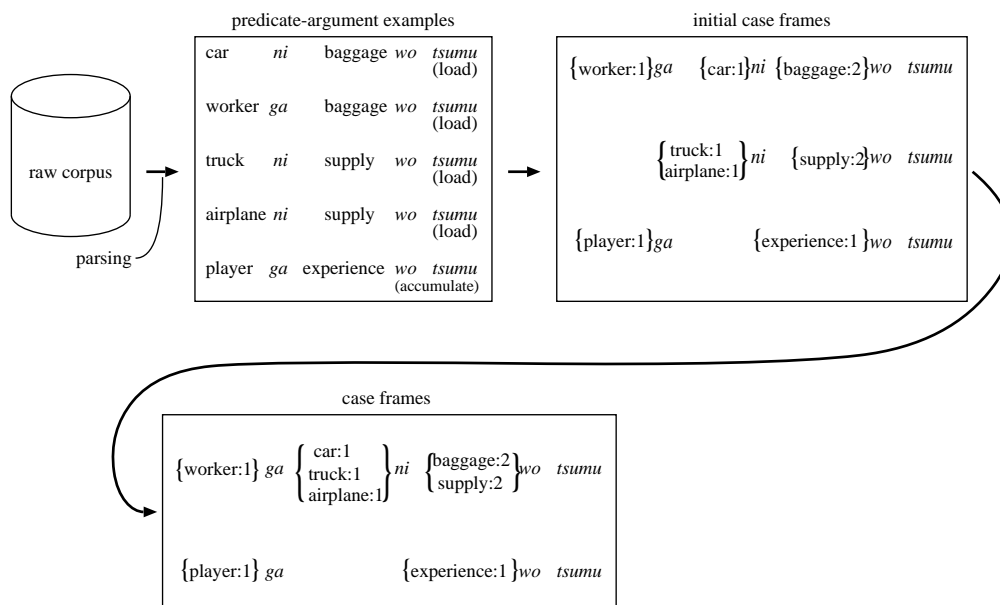Figure 1: Verb paraphrase based on case frame alignment.



Figure 2: Automatic construction of case frame.

are automatically constructed beforehand from a raw corpus.

The outline of the method is depicted in Figure 1. The headword *keitou* and the def-head *necchuu* and *shitau* have several case frames depending on their several usages. For each case frame of *keitou*, we find the most similar case frame among the def-heads' case frames. At this matching stage, case components' correspondences are also found. For example, the case frame

{she, I ...} ga {president} ni

of *keitou* is matched to the case frame

{student, he ...} ga {teacher, professor ...} wo

of *shitau*, and their case components are aligned as shown in the lower part of Figure 1. From this correspondence, we found that this usage of *keitou* has the *shitau* meaning, ni case should be transformed into wo case, and the equivalent is *shitau*.

When an input is paraphrased, first its proper case frame is chosen by the matching of the input and case frames of the verb, then just by tracing the correspondence of the headword case frame and the def-head case frame, the predicate-argument structure of the input can be transformed into that of the def-head.
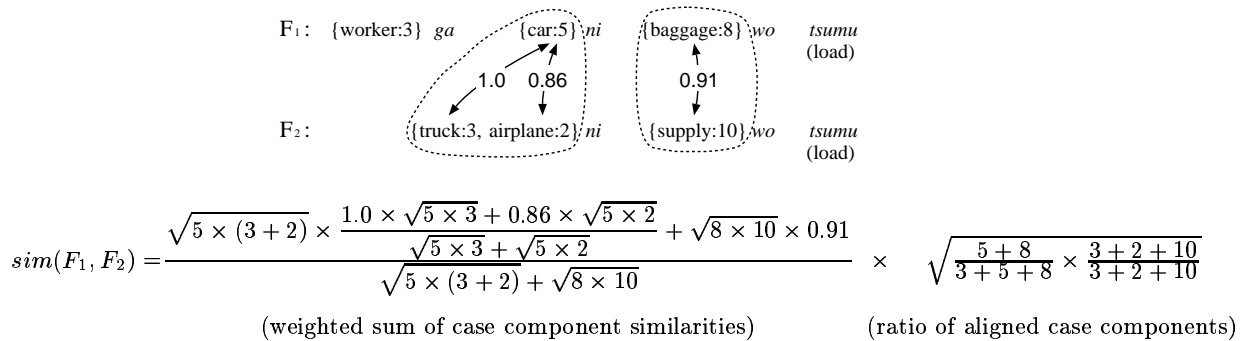
$F_1$: {worker:3} *ga*　{car:5} *ni*　{baggage:8} *wo*　*tsumu* (load)

1.0　0.86　　0.91

$F_2$:　{truck:3, airplane:2} *ni*　{supply:10} *wo*　*tsumu* (load)

$$sim(F_1, F_2) = \frac{\sqrt{5 \times (3+2)} \times \dfrac{1.0 \times \sqrt{5 \times 3} + 0.86 \times \sqrt{5 \times 2}}{\sqrt{5 \times 3} + \sqrt{5 \times 2}} + \sqrt{8 \times 10} \times 0.91}{\sqrt{5 \times (3+2)} + \sqrt{8 \times 10}} \times \sqrt{\frac{5+8}{3+5+8} \times \frac{3+2+10}{3+2+10}}$$

(weighted sum of case component similarities)　　　　(ratio of aligned case components)

Figure 3: Similarity between case frames.

# 3 Automatic Case Frame Construction

## 3.1 Outline

Case frames which are specific enough to the diversity of verb usages and have a wide coverage can be automatically constructed from a large raw corpus (Kawahara and Kurohashi, 2001).

The biggest problem in automatic case frame construction is verb sense ambiguity. Verbs which have different meanings should have different case frames, but it is hard to disambiguate verb senses very precisely. To deal with this problem, we collect predicate-argument examples, which are distinguished by the verb and its closest case component, and cluster them. That is, examples are not distinguished by verbs such as *naru* 'make/become' and *tsumu* 'load/accumulate', but by couples such as *tomodachi ni naru* 'make a friend', *byouki ni naru* 'become sick', *nimotsu wo tsumu* 'load baggage', and *keiken wo tsumu* 'accumulate experience'.

We employ the following procedure for automatic case frame construction (Figure 2):

1. A large raw corpus is parsed by KNP (Kurohashi and Nagao, 1994), and reliable predicate-argument examples are extracted from the parse results.

2. The extracted examples are bundled according to the verb and its closest case component, making initial case frames.

3. The initial case frames are clustered using a similarity measure, resulting in the final case frames.

## 3.2 Similarity between Case Frames

The clustering of initial case frames is performed using a similarity measure. Two case frames, $F_1$ and $F_2$, are first aligned according to the agreement of case markers. Suppose the result of the case component alignment of $F_1$ and $F_2$ is as follows:

$$
\begin{array}{cccccc}
F_1: & C_{11}, & C_{12}, & \cdots & C_{1l} & \cdots & C_{1m} \\
& \updownarrow & \updownarrow & & \updownarrow & & \\
F_2: & C_{21}, & C_{22}, & \cdots & C_{2l} & \cdots & C_{2n}
\end{array}
$$

where $C_{xx}$ denotes a case component which contains several case examples. This result means that $l$ case components are aligned between $F_1$ and $F_2$ and $(m-l)$ and $(n-l)$ case components remained in $F_1$ and $F_2$ respectively.

The similarity between two case frames is based on the weighted sum of aligned case components' similarities and the ratio of aligned case components. We explain how to calculate the similarity in detail (Figure 3 shows an example).

**Similarity between two words**

The similarity between two words, $e_1$ and $e_2$, is calculated using the NTT thesaurus (Ikehara et al., 1997) as follows:

$$sim(e_1, e_2) = \frac{2d_c}{d_{e_1} + d_{e_2}} \tag{1}$$

where $d_{e_1}$ and $d_{e_2}$ are the depths of $e_1$ and $e_2$ in the thesaurus, and $d_c$ is the depth of their lowest (most specific) common node. If $e_1$ and $e_2$ are in the same node of the thesaurus, the similarity is 1.0, the maximum score based on this criterion.

**Similarity between two case components**

The similarity between two case components, $C_{1i}$ and $C_{2i}$, is the weighted sum of the similar-

ities of case examples as follows:

$$sim(C_{1i}, C_{2i}) =$$
$$\frac{\sum_{e_1 \in C_{1i}} \sum_{e_2 \in C_{2i}} \sqrt{|e_1||e_2|} \cdot sim(e_1, e_2)}{\sum_{e_1 \in C_{1i}} \sum_{e_2 \in C_{2i}} \sqrt{|e_1||e_2|}} \quad (2)$$

where $|e_1|$ and $|e_2|$ represent the frequencies of $e_1$ and $e_2$ respectively.

## Weighted sum of case component similarities

The case components' similarities calculated by the above measure are summed up with the weight of case components' frequencies as follows:

$$WSofCCS =$$
$$\frac{\sum_{i=1}^{l} \sqrt{|C_{1i}||C_{2i}|} \cdot sim(C_{1i}, C_{2i})}{\sum_{i=1}^{l} \sqrt{|C_{1i}||C_{2i}|}} \quad (3)$$

where

$$|C_{1i}| = \sum_{e_1 \in C_{1i}} |e_1|, \qquad |C_{2i}| = \sum_{e_2 \in C_{2i}} |e_2|$$

## Ratio of aligned case components

The ratio of aligned case components is calculated as follows:

$$RofACC = \sqrt{\frac{\sum_{i=1}^{l} |C_{1i}|}{\sum_{i=1}^{m} |C_{1i}|} \times \frac{\sum_{i=1}^{l} |C_{2i}|}{\sum_{i=1}^{n} |C_{2i}|}} \quad (4)$$

## Similarity between two case frames

Finally, similarity between two case frames is calculated as follows:

$$sim(F_1, F_2) = WSofCCS \times RofACC \quad (5)$$

Based on this similarity measure, case frames whose similarity is more than the threshold, 0.9, are merged.

### 3.3 Selection of Case Components

If a case component whose frequency is much lower than other case components in a case frame, it might be collected because of parsing errors, or has little relation to its verb. We set the threshold for the case component frequency experimentally as follows:

$$2 \times \sqrt{\left( \begin{array}{c} \text{frequency of the most} \\ \text{frequent case component} \end{array} \right)}$$

Table 1: Pruning of def-head case frames

| Obligatory case (ga, wo, ni) |
| --- |
| Case-1: book wo read ⇒ Case frames whose closest case component is "book". |
| Case-2: book or article wo read ⇒ Case frames whose closest case component is similar to "book" or "article". |
| Case-3: people ni ask ⇒ Case frames whose closest case component belongs the category < HUMAN >. |
| Optional case |
| Case-4: room de read ⇒ Case frames similar to the definition. |

If the frequency of a case component is less than the threshold, it is discarded. For example, suppose the most frequent case component in a case frame is *wo* case, 100 times, and the frequency of *ni* case is 16, *ni* case is discarded (since it is less than the threshold, 20).

## 4 Case Frame Alignment

This section describes how to align the case frames of a headword and those of the def-heads, resulting in word sense disambiguation and detection of equivalents and case marker transformation.

### 4.1 Pruning of Def-Head Case Frames

When a def-head has sense ambiguity, a part of its case frames corresponds to headword's case frames. Suppose that *dokusyo* has a definition "*hon wo yomu*(read a book)" and its def-head is "*yomu* (read)". In the case frame dictionary, there are many case frames of *yomu*, such as "{magazine,article}*wo yomu*" or "{mind}*wo yomu*". But the usage of *yomu* in "{mind}*wo yomu*" is obviously different from that in the definition. So, that case frame will never correspond to headword's case frame. If all case frames of the def-head are blindly used, the accuracy of case frame alignment becomes lower.

In order to mitigate such a problem, we utilize the definition sentence to prune def-head case frames. The point of this process is similar to the automatic case frame construction in section 3. That is, we exploit the information of the closest case component of the def-head.

Table 1 summarizes how to prune def-head case frames. If the closest case component of the def-head is an obligatory case (ga, wo, or ni

case), we prune case frames in three ways according to the case component word.

If the closest case component of the def-head is an ordinal noun such as "book wo read" (case-1 of Table 1), the definition is so specific and we need to consider only this usage of the def-head. Therefore, we only use the case frame whose closest case component contains the same word in the next alignment step.

If the closest case component has a coordinate structure such as "book or article" (case-2), the definition is not so specific as the previous case, but we probably can limit the usages of the def-head whose closest case component is similar to the conjuncts. Therefore, we use the case frames whose closest case component is very similar to one of conjucts (word similarity is larger than 0.9).

If the closest case component is a general term such as "people" (case-3), it specifies the usage of the def-head categorically. We use the case frames whose closest case component belongs to the specified category, that is, has appropriate semantic markers in the NTT thesaurus as follows:

> people/opponent : < HUMAN >
> things: < ABSTRACT >
> place: < PLACE >

If the closest case component is an optional case (case-4), it cannot be used to prune case frames by itself. In such a case, we calculate the similarity between the definition sentence and the def-head case frames by the same criteria in Section 3.2, and use case frames whose similarity is 0.8 or more.

If there is no case component in the definition, of course we cannot prune the def-head case frames.

## 4.2 Alignment of Headword Case Frames and Def-Head Case Frames

Headword case frames and def-head case frames which survived in the pruning process are aligned (see Figure 1). For each headword case frame, we try to find the most similar def-head case frame and the best correspondences of their case components. The difference between this process and the case frame clustering in Section 3.2 is that case components can correspond without case marker agreement in order to allow case marker transformation in paraphrasing.

The similarity measure between the two case frames is almost the same as the similarity measure in Section 3.2. The only difference is the calculation of the similarity between two case components. Instead of the formula (2), the following formula is used:

$$sim'(C_{1i}, C_{2i}) = \frac{\sum_{e_1 \in C_{1i}} |e_1| \cdot \max\{sim(e_1, e_2) | e_2 \in C_{2i}\}}{\sum_{e_1 \in C_{1i}} |e_1|} \quad (6)$$

Probably because a def-head is more general term than a headword, a case component (case examples) of a def-head often covers wider semantic range than that of the headword. Therefore, we do not consider the similarities of all possible pairs of case examples, but the best correspondence for each case example of the headword.

## 5 Experiments and Discussion

### 5.1 Data

We applied the automatic case frame construction method to Mainichi Newspaper Corpus and Nikkei Newspaper Corpus (20 years in total, 15,000,000 sentences). From these corpora, case frames of about 20,000 verbs are constructed; the average number of example case frames of a verb is 16.6; the average number of case slots of a verb is 2.2; the average number of example nouns in a case slot is 4.4.

We used the dictionary, *Reikai Shougaku* dictionary for paraphrasing (Tadika, 1997). The top 2,000 frequent words in definitions of the dictionary were considered as a basic vocabulary of Japanese. Other than these basic words, 220 verbs were randomly chosen, and for each verb a test sentence including the verb was collected from another dictionary, *Shinmeikai* dictionary. We used those 220 sentences as a test set.

### 5.2 Experimental Results

Table 2 shows the result of word sense disambiguation. Out of 220 verbs, 115 verbs have

Table 2: Result of word sense disambiguation.

|            | Correct | Incorrect | Accuracy |
|------------|---------|-----------|----------|
| Baseline   | 60      | 55        | 52.1 %   |
| Our method | 82      | 33        | 71.3 %   |

Table 3: Result of verb paraphrase.

Test sentences whose sense disambiguation was successful and those without sense ambiguity.

|            | Correct | Incorrect | Accuracy |
|------------|---------|-----------|----------|
| Baseline   | 163     | 24        | 87.1 %   |
| Our method | 170     | 17        | 90.9 %   |

Total

|            | Correct | Incorrect | Accuracy |
|------------|---------|-----------|----------|
| Baseline   | 147     | 73        | 66.8 %   |
| Our method | 170     | 50        | 77.2 %   |

word sense ambiguity, that is, they have two or more definitions. On average one verb has 2.5 definitions.

The accuracy of our method is 71.3% (the detection of equivalent and case marker transformation are not checked here). The accuracy of the baseline method which just chooses the first definition is 52.1%.

In SENSEVAL-2 Japanese dictionary task, and the accuracy of the participant systems were 75% to 78%, but they were all supervised systems which used large training data (Shirai, 2001). In SENSEVAL-2 English lexical sample task, the best accuracy of supervised systems was 64%, and that of unsupervised systems was 40% (Kilgarriff, 2001). Considering these scores, we can say the accuracy of our unsupervised WSD method seems reasonably good. Table 3 shows the accuracy of paraphrasing, indicating the effectiveness of our method compared to the baselines. Here, when the equivalent is properly detected and case markers are properly changed, if necessary, the analysis was regarded as correct. Our method could paraphrase 13 sentences correctly in 24 sentences which the baseline method failed, but failed to paraphrase 6 sentences in the 163 sentences which the baseline method could paraphrase correctly.

The upper part of Table 3 is the result for test sentences whose sense disambiguation was successful (82 sentences) and those without sense ambiguity (105 sentences). The baseline regarded a def-head alone as the equivalent and did not change case markers. The lower part

of Table 3 shows the total performance of paraphrasing of 220 sentences. The baseline means the same as the above methods.

### 5.3 Discussion

Table 4 shows examples of successful paraphrasing, including examples which properly detect larger equivalents and transform case markers.

The main cause of incorrect paraphrase is the data sparseness problem of automatically constructed case frames. Since they are constructed from newspaper corpus, they do not cover daily-life expressions well. Furthermore, since definition sentences are sometimes strange from the view point of standard usage of language, it was harder to collect def-head case frames than those of headwords.

### 5.4 Related Work

Kondou et al. used a dictionary as a knowledge base of paraphrase (Kondo et al., 1999), but did not handle word sense ambiguity and case marker transformation.

Takahashi et al. proposed a software environment for manual construction of sentential pattern transformation (Takahashi et al., 2001).

Barzilay et al. pointed out that there are many paraphrases in a corpus of multiple English translations of the same source text, and proposed an unsupervised method of extracting paraphrases from such

a parallel corpus.(Barzilay and McKeown, 2001)

## 6 Conclusion

This paper proposed a method of translating a predicate-argument structure of a verb into that of an equivalent verb, which is a core component of the dictionary-based paraphrasing. Our method grasps several usages of a headword and those of the def-heads as a form of their case frames and aligns those case frames, which means the acquisition of word sense disambiguation rules and the detection of the appropriate equivalent and case marker transformation.

We are planning to extend our dictionary-based paraphrasing system to more complicated phrases and sentences. We also would like to apply the proposed method to a word selection task in machine translation.

Table 4: Examples of successful paraphrases.

| kakageru | 1 | takaku (highly) | ageru (raise) | | |
|---|---|---|---|---|---|
| | 2 | kangae ya (idea) | shutyou wo (request) | hiroku (publicly) | shirareru youni suru (make known) |
| | 3 | shinbun (newspaper) | nado ni | noseru (publish) | |

| kokki wo (national flag) | kakageru | ⇒ | kokki wo (national flag) | takaku (highly) | ageru (raise) |
|---|---|---|---|---|---|

| kouryaku | 1 | teki (enemy) | no | jinchi ya (encampment) | shiro wo (castle) | ubau (obtain) |
|---|---|---|---|---|---|---|
| | 2 | aite wo (opponent) | semmete (attack) | makasu (defeat) | | |

| yokozuna wo (grand champion) | kouryaku suru | ⇒ | yokozuna wo (grand champion) | makasu (defeat) |
|---|---|---|---|---|

| tozakeru | 1 | tooku (away) | he | hanare saseru (get) |
|---|---|---|---|---|
| | 2 | tsukiawa naku suru (not to get along) | | |

| akuyuu wo (bad friend) | tozakeru | ⇒ | akuyuu to (bad friend) | tsukiawa naku suru (not to get along) |
|---|---|---|---|---|

| narihibiku | 1 | naru (ring) | oto ga (sound) | kikoeru (hear) |
|---|---|---|---|---|
| | 2 | hyouban ga (reputation) | hiroku (widely) | sirewataru (have) |

| bell ga | narihibiku | ⇒ | bell no | oto ga (sound) | kikoeru (hear) |
|---|---|---|---|---|---|

| nagabiku | zikan (time) | ga | kakaru (take) |
|---|---|---|---|

| kousyou ga (negotiation) | nagabiku | ⇒ | kousyou ni (negotiation) | zikan (time) | ga | kakaru (take) |
|---|---|---|---|---|---|---|

## References

Regina Barzilay and Kathleen R. McKeown. 2001. Extracting paraphrases from a parallel corpus. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, pages 50–57.

Satoru Ikehara, Masahiro Miyazaki, Satoshi Shirai, Akio Yokoo, Hiromi Nakaiwa, Kentarou Ogura, and Yoshifumi Oyama Yoshihiko Hayashi, editors. 1997. *Japanese Lexicon*. Iwanami Publishing.

Daisuke Kawahara and Sadao Kurohashi. 2001. Japanese case frame construction by coupling the verb and its closest case component. In *Proceedings of the Human Language Technology Conference*, pages 204–210.

Adam Kilgarriff. 2001. English lexical sample task description. In *Proceedings of SENSEVAL-2*, pages 17–20.

Keiko Kondo, Satoshi Sato, and Manabu Okumura. 1999. Paraphrasing of "sahen-noun + suru" (in Japanese). *Journal of Information Processing Society of Japan*, 40(11):4064–4074.

Sadao Kurohashi and Makoto Nagao. 1994. A syntactic analysis method of long japanese sentences based on the detection of conjunctive structures. *Computational Linguistics*, 20(4).

Kiyoaki Shirai. 2001. Senseval-2 japanese dictionary task. In *Proceedings of SENSEVAL-2*, pages 33–36.

Jyunichi Tadika, editor. 1997. *Reikai Shougaku Kokugojiten (Japanese dictionary for children)*. Sanseido.

Tetsuro Takahashi, Tomoya Iwakura, Ryu Iida, Atsushi Fujita, and Kentaro Inui. 2001. Kura: A revision-based lexico-structural paraphrasing engine. In *Proceedings of the Natural Language Processing Pacific Rim Symposium Post-Conference Workshop*, pages 37–46.