

An Example-Based Approach to Japanese-to-English Translation of Tense, Aspect, and Modality

Masaki Murata Qing Ma Kiyotaka Uchimoto Hitoshi Isahara
Communications Research Laboratory,
Ministry of Posts and Telecommunications
588-2, Iwaoka, Nishi-ku, Kobe, 651-2401, JAPAN
{murata, qma, uchimoto, isahara}@crl.go.jp

Abstract

We have developed a new method for Japanese-to-English translation of tense, aspect, and modality that uses an example-based method. In this method the similarity between input and example sentences is defined as the degree of semantic matching between the expressions at the ends of the sentences. Our method also uses the k-nearest neighbor method in order to exclude the effects of noise; for example, wrongly tagged data in the bilingual corpora. Experiments show that our method can translate tenses, aspects, and modalities more accurately than the top-level MT software currently available on the market can. Moreover, it does not require hand-craft rules.

1 Introduction

The translation of Japanese tenses, aspects, and modalities into English are some of the most difficult problems in machine translation. Conventional approaches to these problems translate Japanese tenses, aspects and modalities according to hand-craft rules that use tense and aspect information (Kume et al. 1990) (Shirai et al. 1990). However, the complexity of Japanese tense/aspect/modality expressions makes it very difficult to formulate detailed rules. We therefore tried to translate Japanese tense/aspect/modality expressions using the example-based method, which was developed by Nagao (Nagao 1984). We prepared bilingual corpora containing pairs of Japanese and English sentences and tried translating tense/aspect/modality expressions by using the tense/aspect/modality expression of the English sentence corresponding to the most similar Japanese sentence.

The example-based method developed by Nagao in 1984 is effective but has rarely been used since it was used by Sumita et al. (Sumita et al. 1990) in the translation of the Japanese particle *no*¹. The method we describe here is the first application of the example-based method to tense/aspect/modality translation. It is based on a very simple measurement of the similarity between an input sentence and an example sentence. Similarity is defined as the degree of matching between the strings (or the

¹ The Japanese particle *no* has many English translations: “of,” “in,” “at,” “for,” and so on. Their work showed that an appropriate preposition can be chosen using the example-based method.

degree of semantic matching including the category number of the thesaurus and the inflectional form) in the expressions at the end of two sentences. Our method can also be used to analyze monolingual tense/aspect/modality if we substitute the corpora tagged with the correct tense/aspect/modality for the bilingual corpora.

Machine translation is very difficult because it requires both semantic analysis and discourse analysis, neither of which can be done well by the language-processing technology available today. Since our knowledge of language analysis and generation is insufficient, we lack the fundamental knowledge needed for high-quality machine translation. But machine translation is sometimes accomplished well enough by simply replacing words, as in a puzzle game. We want to use even a simple technique if it is at all useful. We therefore developed a simple method that can translate tenses, aspects, and modalities better than the top-class MT software can but that does not do deep processing such as semantic analysis, and that does not require hand-craft rules.

2 Example-Based Translation of Tense, Aspect, and Modality

2.1 Using a string matching at the end of sentences as the definition of similarity

Murata and Nagao have already used the example-based method to resolve the verb-phrase ellipsis in Japanese sentences (Murata & Nagao 1997). In the following sentence, for example, the verb *arimasu* “I have” at the end of the sentence is omitted.

[Input sentence]
jitsu-wa chotto onegaiga (arimasu).
 (actually) (a little) (request) (I have)
 Actually, (I have) a little request. (1)

[Example sentence]
 the matching part **the latter part**
anou chotto onegaiga arimasu. (2)
 (er) (a little) (request) (I have)
 Er, I have a little request.

In their method, for resolving this elliptical sentence, they search a corpus for sentences containing the longest string of characters matching those at the end of the input sentence (“*jitsu-wa chotto onegaiga.*”), get example sentences such as “*anou chotto onegaiga arimasu,*” and judge that the verb *arimasu* “I have” in the latter part of the detected sentences is an omitted verb. To find an example sentence similar to the input sentence, we must of course first define similarity. How similarity is defined is critical because the result of using an example-based method depends on the definition of the similarity. Murata and Nagao defined the similarity as the number of characters in the matching part from the end of the sentence, a definition that is both simple and appropriate for this problem and that resolves the elliptical verb phrase at the end of the sentence.

We think that because tense/aspect/modality expressions are at the ends of Japanese sentences, this definition of similarity can also be used in tense/aspect/modality translation. Our method searches a bilingual corpus for the Japanese example sentence containing the longest string matching that at the end of the input Japanese sentence, and it selects the tense/aspect/modality expression of the corresponding English sentence in the corpus as the tense/aspect/modality expression used for the English translation of the Japanese input sentence. Suppose that we translate the tense/aspect/modality expression of the following Japanese input sentence.

[Input sentence]
kare-wa yuumei-ni naritai-toiu yashin wo idai-teiru.
 (He) (famous) (to become) (an ambition) obj (have) (3)
 He has an ambition to become famous.

[Example sentence]
kare-wa hurusato-eno hageshii bojoh wo idai-teiru.
 (He) (home) (great) (a longing) obj (have) (4)
 He has a great longing for home.

The corresponding part

At first, we detect the example sentence containing the longest expression at the end of this input sentence. We then find that the above example sentence is the one containing the longest expression *wo idai-teiru*, “have.” We look at the verb of the English translation of the example, find that the tense/aspect/modality expression is the present tense form, and translate the tense/aspect/modality expression of the input sentence into the present tense. A rule-based method, in contrast, would likely determine the tense/aspect/modality expression to be the progressive form, since this sentence has a Japanese tense/aspect/modality expression *teiru*, which often means progression². Our example-based method, however, can correctly judge that the tense/aspect/modality expression of the input sentence (3) is the present tense form.

This similarity based on matching the strings at the end of sentences is simpler and more tractable than the similarity used in the translation of **Noun X no Noun Y**, to which the example-based method was first applied. In the problem of **Noun X no Noun Y**, there are some cases when **Noun X** is more important than **Noun Y** and some cases when **Noun Y** is more important than **Noun X**. We therefore need to appropriately weight **Noun X** and **Noun Y**, so the similarity is very complicated. But if we measure similarity by matching the strings at the end of sentences, we have only to check the string in order from the end of the sentence.

² The Japanese tense/aspect/modality expression *teiru* often means progression as in the following sentence.

kare-wa sentou-no sousha-ni pittari-kuttsuite hashit -teiru
 (He) (front) (runner) (at the heels of) (run) (-ing) (5)
 He is running at the heels of the front runner.

Table 1: Information obtained from language analysis

Morphology			Category number	Inflectional form
彼	<i>kare</i>	(He)	1200003012	
は	<i>wa</i>	topic	1195038023	
野望	<i>yabou</i>	(ambition)	1304207024	
を	<i>wo</i>	obj		
いだいて	<i>idaite</i>	(have)	2153417012	タ系連用テ形 (<i>ta-series predicative te-form</i>)
いる	<i>iru</i>	(be)	2120002012	基本形 (<i>the normal form</i>)

2.2 Two measures for matching strings at the end of sentences

In recent years, the technologies on natural language processing have developed and various morphological analyzers are open to the public. In our analysis, we check the degree of matching of strings at the end of a sentence in order to detect an example similar to an input sentence. At this time, we check the matching after recognizing words by morphological analysis. And we also check the similarity between words by using the semantic distance between words in the thesaurus rather than by matching strings of words. For checking the match at the end of a sentence, we therefore use the method using the result of language analysis in addition to the method using only strings. These two methods are explained below:

- **Method 1** Using simple strings

This method is the one mentioned in the previous section. It checks the degree of string matching from the end of a sentence and uses the length of the matching string as the similarity.

- **Method 2** Use of the result of language analysis

This method performs high-quality matching by using a morphological analyzer and a thesaurus. At first, we detect morphologies by using a morphological analyzer (Kurohashi & Nagao 1998). Next, we give each morphology a category number representing that morphology in a Japanese word thesaurus (NLRI 1964). When the morphology is an inflectional word, we also give it the inflectional form (e.g., the past tense) that is obtained from the output of the morphological analyzer.

For example, the sentence “*kare wa yabou wo idaite iru.*” (“He has an ambition.”), is represented by the information in Table 1. In the table, the input sentences are divided into morphologies such as *kare* “he” and *wa* topic, and each of them is given a category number and the inflectional form. In the thesaurus, each word has a 10-digit category number. This 10-digit category number indicates seven levels of an is-a hierarchy. The top five levels are expressed by the first five digits of the category number. The sixth level is expressed by the following two digits of the category number. And the last level is expressed by the last three digits of the category number.

After assigning the category numbers, we check the degree of matching at the end of a sentence by using the information in Table 1. At this time, to check the string match from the end of the sentence in the same as in Method 1, we use the following string combining all the information in Table 1. (We do not use the last three digits of the category number.)

彼 03 0 0 0 2 1 :
は 38 0 5 9 1 1 :
大望 07 2 4 0 3 1 :
を :
いだいて 17 4 3 5 1 2 タ系連用テ形 :
いる 02 0 0 2 1 2 基本形

In this information, we reverse the category number. This means that when we check the matching from the end of a sentence, we check the matching from the top of the category number and obtain the same result we would obtain if we used the normal way to check semantic similarity in the thesaurus.

In Method 2, we transform an input sentence into the above information and check the length of the matching characters from the end of the sentence. The length of the matching characters is treated as the similarity used in the example-based method. We can check, in order, the inflectional form, the similarity in the thesaurus, and the similarity of the strings of each morphology by checking the string match from the end in the above information.

2.3 Using the k-nearest neighbor method for preventing the problem of noise

The k-nearest neighbor method contains the example-based method (Fukunaga 1972). Instead of using the one-nearest example, this method uses the result obtained from the “voting” of the k nearest examples. The decision obtained by using only one example is unreliable since that example may be a noise. The decision using k examples makes a stable analysis possible even when the data include a little noise.

In the work reported here, we used 1, 3, 5, 7, and 9 as k. When one of the k-nearest examples has the same highest similarity as other examples, we should use all of them regardless of the value of k. In this work, however, we limited the number of examples to 10 in order to simplify the processing. When different tense/aspect/modality expressions had the same number of votes, the expression selected was that of the example obtained first.

Next we examine the k-nearest method by using the example of tense/aspect/modality translation in Table 2. Table 2 shows the analysis of the tense/aspect/modality expression of the input sentence “*kare wa watashi no shiriai da.*” (I am acquainted with him.) by using Method 2. The calculation of the similarity by using Method 2 is illustrated by the data listed in Table 3, where the bold-faced part matches the input sentence. One Japanese character consists of two bytes. So in this work, the

Table 2: Example of tense/aspect/modality translation

Input		Japanese	Category	English
		彼は私の知り合いだ	Present	I am acquainted with him.
No.	Sim.	Example sentence		
1	25	彼とは長年の知り合いだ	Present perfect	I have known him for a long time.
2	24	ふたりは長い間の知り合いだ	Present	The two are acquaintances of long standing.
3	11	彼とは10年余の顔見知りだ	Present perfect	I have known him for over ten years.
4	11	彼らは多年の知己だ	Present	They are friends of many years' standing.
5	10	彼はこのクラブの恩人だ	Present	He is a benefactor of this club.
6	10	彼は私の命の恩人だ	Present	I owe him my life.
7	10	彼はかたい人だ	Present	He is reliable.
8	10	彼はだれにも人当たりのいい人だ	Present	He is affable to everybody.
9	10	彼は胸触りのいい人だ	Present	He is a mild-mannered person.
10	10	なんと男振りもいい人だ	Present	What a handsome-looking man he is!

Sim. = Similarity

Category = Category of Tense/Aspect/Modality

Table 3: Calculation of similarity by string matching from the end of the sentence

Input	彼0300021:は3805911:私0100021:の0700011:知り合い0301221:だ基本形	
No.	Sim.	25242322212019181716151413121110987654321
1	25	3805911:は3805911:長年0724611:の0700011:知り合い0301221:だ基本形
2	24	ふた0305911:り:は3805911:長い間:の0700011:知り合い0301221:だ基本形
..
4	11	彼ら0300021:は3805911:多年0701611:の0700011:知己0301221:だ基本形
..

number of two-byte sequences in the matching part represents the similarity. Example 1 in Table 3, for example, has a similarity of 25 since the length of the matching part is 25 two-byte sequences. The results obtained from the 10 most-similar example sentences are listed in Table 2, where “Tense/Aspect/Modality” is that obtained from the tense/aspect/modality expression of the English sentence corresponding to the Japanese example sentence.

When $k = 1$, the tense/aspect/modality expression was analyzed by using only the example most similar to the input sentence, Example 1, which has the tense/aspect/modality expression “present perfect.” So our system judged that the target tense/aspect/modality expression was “present perfect,” even though the correct one was “present.” When $k = 3$, we tried to select the three most-similar example sentences but found that Examples 3 and 4 had the same similarity. So we used four examples, two of which voted “present perfect” and two of which voted “present.” The incorrect tense/aspect/modality expression “present perfect” was again selected because it was obtained earlier in the processing. When $k = 5$, we tried to select the five most-similar example sentences but found that Examples 5 through 10 had the same similarity. So we used all ten, two of which voted for “present perfect” and eight of which voted “present.” The correct tense/aspect/modality expression, “present,” was thus selected. When $k = 7$ or 9, we used all ten and got the correct tense/aspect/modality expression, “the present,” as when $k = 5$. The system output an incorrect answer when k is 1 or 3, and output a correct answer when k is 5, 7 or 9.

Table 4: Result

	All	Present	Past	Other
Software	80.6% (233/289)	91.1% (112/123)	96.3% (105/109)	28.1% (16/ 57)
Method 1(k=1)	76.8% (222/289)	89.4% (110/123)	89.9% (98/109)	24.6% (14/ 57)
Method 1(k=3)	82.0% (237/289)	93.5% (115/123)	97.2% (106/109)	28.1% (16/ 57)
Method 1(k=5)	83.0% (240/289)	95.1% (117/123)	97.2% (106/109)	29.8% (17/ 57)
Method 1(k=7)	82.4% (238/289)	94.3% (116/123)	96.3% (105/109)	29.8% (17/ 57)
Method 1(k=9)	82.4% (238/289)	94.3% (116/123)	96.3% (105/109)	29.8% (17/ 57)
Method 2(k=1)	78.5% (227/289)	88.6% (109/123)	89.9% (98/109)	35.1% (20/ 57)
Method 2(k=3)	81.7% (236/289)	91.1% (112/123)	95.4% (104/109)	35.1% (20/ 57)
Method 2(k=5)	81.3% (235/289)	92.7% (114/123)	93.6% (102/109)	33.3% (19/ 57)
Method 2(k=7)	81.7% (236/289)	92.7% (114/123)	93.6% (102/109)	35.1% (20/ 57)
Method 2(k=9)	81.7% (236/289)	92.7% (114/123)	93.6% (102/109)	35.1% (20/ 57)

3 Experiment and Discussion

3.1 Experiment

We carried out the experiments on tense/aspect/modality translation in order to verify the method described in Section 2. We used the bilingual corpus (36,617 sentences) in the Kodansha Japanese-English dictionary (Shimizu & Narita 1976) as the database of examples. From this corpus, we randomly selected 300 sentences as input sentences and compared the results obtained by using our method with those obtained by using the top-level software currently available on the market. When we ran the software on the 300 input sentences, the verb parts of 11 of them could not be translated and the tense/aspect/modality expressions could not be obtained from them. We therefore eliminated these 11 sentences from our experiments.

We classified the tense/aspect/modality into the following 27 categories:

1. all the combinations of {Present, Past}, {Progressive, Not-progressive}, and {Perfect, Not-perfect} (8 categories),
2. imperative mood (1 category),
3. auxiliary verbs ({Present, Past} of “be able to”, {Present, Past} of “be going to”, can, could, have to, had to, let, may, might, must, need, ought, shall, should, will, would) (18 categories).

“Must” and “have to” or “can” and “be able to” should really be grouped together, but since they may have different meanings, we defined the tense/aspect/modality according to the English surface expression strictly and handled these cases as different tenses/aspects/modalities. We used the tense/aspect/modality expression of the corresponding verb in the English sentence as the correct tense/aspect/modality³.

³ In the experiment the criterion for judging whether the result was correct was very strict: the output tense/aspect/modality must be the same as the tense/aspect/modality of the English translation of the input sentence in our bilingual database. As in 2(b) in Section 3.2, there are some cases when English tense/aspect/modality expressions that express the same tense/aspect/modality are different. The real accuracy rates may be much higher than listed those in Table 4.

Table 5: Accuracy when determining each tense/aspect/modality

	All	Pr.	Past	Pr.-ing	P.-ing	Perf.	Imp.	can	could	let	may	must	will	would
No.	289	123	109	7	1	15	12	3	2	1	2	4	9	1
Software														
	81%	91%	96%	29%	0%	0%	67%	67%	100%	100%	0%	25%	0%	0%
Method 1														
k=1	77%	89%	90%	14%	0%	7%	58%	33%	50%	100%	0%	0%	22%	0%
k=3	82%	93%	97%	29%	0%	0%	75%	33%	0%	100%	0%	0%	33%	0%
k=5	83%	95%	97%	14%	0%	0%	83%	33%	0%	100%	0%	25%	33%	0%
k=7	82%	94%	96%	14%	0%	0%	83%	33%	0%	100%	50%	25%	22%	0%
k=9	82%	94%	96%	14%	0%	0%	83%	33%	0%	100%	50%	25%	22%	0%
Method 2														
k=1	79%	89%	90%	29%	0%	13%	75%	33%	100%	100%	0%	0%	33%	0%
k=3	82%	91%	95%	29%	0%	13%	83%	0%	50%	100%	0%	25%	33%	0%
k=5	81%	93%	94%	29%	0%	7%	92%	0%	0%	100%	0%	25%	33%	0%
k=7	82%	93%	94%	14%	0%	7%	92%	0%	0%	100%	50%	50%	33%	0%
k=9	82%	93%	94%	14%	0%	7%	92%	0%	0%	100%	50%	50%	33%	0%

The accuracies obtained when determining each tense/aspect/modality are listed in Table 4. Only 13 of the 27 tense/aspect/modality expressions were found in the 289 sentences, and the accuracy rates for each of them are listed in Table 5. “Pr,” “P.,” “-ing,” “Perf.,” and “Imp.” respectively indicate “Present,” “Past,” “Progressive,” “Perfect,” and “Imperative.”

3.2 Discussion

1. Accuracy rates

- (a) Method 1 when k=5 is best (83%) (Table 4). This result shows that even simple string-matching can yield comparatively high accuracy rates.
- (b) All the overall “All” accuracy rates obtained using our methods when $k \neq 1$ are higher than those obtained using the software.

When $k = 1$, our method suffers from noise and the accuracy rate is low. And the results in Table 4 clearly show that the k-nearest neighbor method is effective.

- (c) As listed in Table 4, when determining “Other” tenses/aspects/modalities (those other than “Present” and “Past”), Method 2, using the result of language analysis, yields higher accuracy rates than those of Method 1 or the software.

It is important to examine the accuracy rate when determining difficult tenses/aspects/modalities if we want to implement high-quality machine translation. Even if the overall accuracy rate (“All”) is high, high-quality translations will not be produced if only “Present” or “Past” are selected correctly. Although the overall accuracy of Method 2 is a little low, the accuracy for determining the difficult tenses/aspects/modalities “Other” is high. We therefore think that Method 2 is more promising for high-quality machine translation than Method 1.

2. Problems of our method

- (a) In Japanese, two sentences that have the same surface expression for the verb phrase sometimes have different tenses/aspects/modalities. For example, in the two-verb phrases *tokeru* “can solve” or “thaw” in the following examples, the first one has the modality “Potential” and the second one has only the tense “Present.”

shougakusei-nara taiteiwa konomondai-wa tokeru.
(elementary schoolchildren) (most) (this problem) (can solve)
Most elementary schoolchildren can solve this problem.

ike-no koori-wa sangatsu-ni tokeru.
(pond) (ice) (in March) (thaw)
The pond thaws in March.

To handle these examples, we must disambiguate the word senses of *tokeru*: “can solve” or “thaw.”

- (b) Although our method uses only Japanese tense/aspect/modality expressions and does not consider the structure of the English translated sentence, the tense/aspect/modality that should be used sometimes depends on the structure of the translated English sentence. The first of the following sentences has the aspect “Non-Progressive,” and the second has the aspect “Progressive.”

kare-wa shitsujituna seikatsu-wo okut-teiru.
(He) (sober and simple) (life) (live)
He lives a sober and simple life.

kare-wa taidana seikatsu-wo okut-teiru.
(He) (lazy) (life) (be leading)
He is leading a lazy life.

We can consider that the verbs of these Japanese sentences have almost the same meaning, and the same tense/aspect/modality. But changing the verb used in English translation from “live” to “lead” makes the difference between “Present” and “Progressive.” If we want to use our method in high-quality translation, it is necessary to match not only Japanese sentences but also English sentences when detecting a similar example. In an overall machine translation system, the structure of the English translation of a Japanese input sentence is made up of results of the structure analysis. By using the results, we will only be able to detect examples whose structure is similar to the structure of the English translation.

- (c) In some cases it would be better to use not only expressions at the end of the sentence but also adverbs at the beginning (Kume et al. 1990). For example, the difference between *mou* “already” and *kinou* “yesterday” makes the difference between “Past perfect” and “Past.”

mou *touroku-shita*. I've already registered,
(already) (register)

kinou *touroku-shita*. I registeredyesterday.
(yesterday) (register)

Our method would have to be changed if it were to handle the above case.

3. Advantages of our method

- (a) It does not require hand-craft rules.
- (b) It is very easy to implement.

Our method determined tense/aspect/modality more accurately than the top-level MT software currently available on the market. This indicates that our method is useful.

4 Conclusion

To translate Japanese tense/aspect/modality expressions into English by using the example-based method, we defined the similarity between input and example sentences as the degree of semantic match between expressions at the end of sentences. We used the k-nearest neighbor method in order to exclude the effects of noise. In experiments, our method translated tense/aspect/modality expressions more accurately than the top-level MT software currently available on the market did. Another advantage of our method is that it does not require hand-craft rules.

We used two methods to evaluate the degree of similarity: one that simply matches character strings, and the other that uses the result of language analysis. The overall accuracies obtained by using the string-matching are only a little better than those obtained by using language analysis. However, the results of translating tenses/aspects/modalities other than “Present” and “Past” are quite a bit better when the language analysis was used. Because high-quality machine translation requires effective handling of difficult tenses/aspects/modalities, we think that the latter method will be more promising.

The tense/aspect/modality translation method we developed can also be applied to English-to-Japanese translation by eliminating the subject of the English input sentence and using string-matching from the beginning of the remainders; that is, from the beginning of a verb phrase. And because this method does not need hand-craft rules, it is very useful for many other languages where hand-craft rules have not been prepared well. We will also be able to use our method for monolingual tense/aspect/modality analysis. For example, if instead of the bilingual corpora we use the monolingual corpora tagged with the correct tense/aspect/modality, we will be able to identify the tense/aspect/modality immediately.

References

- Fukunaga, Keinosuke: 1972, *Introduction to Statistical Pattern Recognition*, Academic Press Inc.
- Kume, Masako, Takayuki Toyoshima & Masaaki Nagata: 1990, 'Japanese aspect processing for spoken language translation', in *Information Processing Society of Japan, the 40th National Convention, IF-7*, pp. 415-416, (in Japanese).
- Kurohashi, Sadao & Makoto Nagao: 1998, *Japanese Morphological Analysis System JUMAN version 3.5*, Department of Informatics, Kyoto University, (in Japanese).
- Murata, Masaki & Makoto Nagao: 1997, 'Resolution of verb ellipsis in Japanese sentence using surface expressions and examples', in *NLPRS'97*.
- Nagao, Makoto: 1984, 'A Framework of a Mechanical Translation between Japanese and English by Analogy Principle', *Artificial and Human Intelligence*, pp. 173-180.
- NLRI: 1964, (National Language Research Institute). *Word List by Semantic Principles*, Syuei Syuppan, (in Japanese).
- Shimizu, Mamoru & Narimasu Narita, eds.: 1976, *The KODANSHA Japanese-English Dictionary*, Kodansha.
- Shirai, Satoshi, Akio Yokoo & Francis Bond: 1990, 'Generation of tense in newspaper translation', in *The Institute of Electronics, Information and Communication Engineers, Autumn Convention*, pp. D-69, (in Japanese).
- Sumita, Eiichiro, Hitoshi Iida & Hideo Kohyama: 1990, 'Translating with examples : A new approach to machine translation', in *The Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Language*, TMI, no. 3, pp. 203-212.