

Data Collection and Transliteration of Japanese Spontaneous Database in the Travel Arrangement Task Domain

Akira Kurematsu¹, Youichi Akegami², Tanja Schultz³, Susanne Burger⁴

¹ Univ. of Electro-Communications, JAPAN, e-mail: kure@apple.ee.uec.ac.jp

² ATR International, JAPAN, e-mail: akegami@ctr.atr.co.jp

³ Univ. of Karlsruhe, Germany, e-mail: tanja@ira.uka.de

⁴ Univ. of Munchen, German, e-mail: burger@phonetik.uni-muenchen.de

ABSTRACT

This paper describes the method to construct and transcribe Japanese spontaneous speech data for VERBMOBIL, the German research project of speech translation. Spontaneous spoken dialogue database is the basis for developing speech and language processing for dialogue systems such as speech translation system. The extended data of human-to-human spoken dialogue in the scenario of travel arrangement has been initiated to be collected in German, English and Japanese in the travel arrangement task. Romanized transcription is used to develop acoustic model and language model in speech recognition system, and natural language translation system. In this paper, issues of transliteration method and several rules and conventions to transcribe Japanese spoken dialogue will be described.

1. Introduction

The corpora of person to person dialogue are needed especially for a spoken language system and a speech translation system. The speech situation with two people is more natural than a guided speech by the machine. The study of simultaneous speech from two speakers is important for research on dialogue or dialogue analysis, intention analysis, and spoken language understanding. In addition, speech recognition system, which are now only able to deal with one speaker at a time, would also have to be able to deal with different speaker talking simultaneously. Speech corpora containing dialogues will be used for the training and testing data for advanced recognizers. The issues of spoken language corpus design is found in detail in [1].

The spontaneous speech corpora has been collected under VERBMOBIL project, the German project of VERBMOBIL of face-to-face speech translation system. The project have supported to collect spontaneous spoken dialogue data in German, English and Japanese in the same scenario. The corpora of VERBMOBIL2 in succession to VERBMOBIL1 is gathered in the travel arrangement task with the role playing manner. The purpose of the VERBMOBIL2 speech data collection is to record situational dialogues as close as possible to actual dialogues between businessman. In the following chapters, issues of data collection of the travel arrangement task, and the transliteration method and several rules and conventions to transcribe Japanese spoken dialogue will be described.

2. Data Collection

2.1. Scenario

In the VERBMOBIL data collection, there are several assumptions for conversation. Two speakers are businessmen from the same company. As for the scenario of the Japanese data collection, two speakers is assumed to trip together and the trip is planned to visit the

office in Hannover in Germany from Japan and starting point of the trip is the same place of the recording place. The information on the schedule of the trip and flight schedule and hotel name change according to the selected trip document. Speakers look at separate calendars and flight schedule and hotel information. The schedules are conceived that the date of the dialogue will be set properly in reference to the calendar of the speakers. The calendar of three months are given to each speaker. We assume that the place where the conversation is to be held is the same as the starting point for the trip, i.e., the place of recording. The information on the flight schedule changes according to the place of recording.

The whole business trip is to take five days. The business in Hannover will require one and a half days, and the remaining time is for sight-seeing.

Although the entire recording is to be performed over several days, a schedule is to be prepared every year. The speech data is made to cover all twelve months. Therefore, the day of conversation will differ from the day of recording.

Three set schedule sheets of the same month are prepared. The same speaker presupposes that a different schedule is to be used. The total number of the dialogues is assumed to be approximately more than 200. Conversations by a calendar of the same month are taken as 14 conversations.

As information on the stock of a hotel, information on three hotels, a place, some equipment, and a price list of 10 hotels are used. The combination of hotel selections changes each time. The visit of sights and events is only broadly conversed. No special documents are prepared on this business trip. This is due to reduce the complexity of conversations.

2.2. Data Collection

Data sheet for recording As the data sheet for recording, following documents are prepared for the recordings.

- For experimental persons, a name-tag on which the name is written is placed. With this the speakers can easily address themselves with names.
- Calendar sheet containing three months. The schedule of the calendar is designed in rather general terms such as meeting, holiday, seminar etc.
- For speaker A, information on the timetable of flights from Kansai or Narita airport to Hannover. The schedule of the plane is designed in rather general terms of the plane company name and flight number and departure time and arrival time. The airport of the transit is also shown.
- For speaker B, information on hotels in Hannover.
- For each experimental person, consent documents on using the speech data for the research.

Conditions for speakers As for the conditions for speakers, two speakers sit along a desk and have a face to face conversation. Although the dialog of the two may be superimposed, it is made to be as short as possible. Both speakers' native language is Japanese. Speakers are allowed to speak standard spoken Japanese, but they are not required to have had such special training in speaking. No strict rules are set on pronunciation.

Conditions for recording The recording system uses SennHeiser HMD414 microphones and records with a DAT deck. The DAT's sampling rate is 48KHz and the signal data is resampled at 16KHz by using workstations. Speech data are recorded to one channel per speaker.

Instructions for testconductor Data is collected so that the data are as close as possible to realistic dialogues of business men. The purpose of the data collection of the dialogue is explained that the data will be used for the development of speech recognition and language translation of the speech translation system. Corresponding to the pattern of the dialogue, the instructor explains the situation of the dialogue and the stream of the topics of dialogue.

3. Transliteration

3.1. Transliteration of Spontaneous Speech

The standard method to transliterate spontaneous speech is defined in VERBMOBIL [2]. Japanese spontaneous speech data for VERBMOBIL is transcribed in Kanji-kana and Roman script with the segmentation into words. By using an orthographic transliteration, the data contained in the dialogues are made accessible in written symbols to a wider range of research and usage. The user of the corpora listen to the recorded dialogues and transliterate these at the level of lexical elements.

Transliteration conventions are used for various reasons:

- phenomena occur in spontaneous such as disruptions of sentences, corrections and repetitions of utterances, reductions and hesitations,
- dialogues technical artifacts may occur, such as technically caused disruptions of recordings or noises resulting from the microphone,
- situation in dialogues may lead to interference of speakers by the partner.

The outline of conventions to transliterate various phenomena of spontaneous is described in [3].

3.2. Basic Requirements and Limitations

The basic requirements of the transliteration is given in two issues; a computational processibility and requirements regarding contents. As for the computational processibility, a unified file structure and consistent transliteration is maintained. Transliteration conventions which can be parsed by a parser which will be a tool for filtering with various filter options. As for the requirements regarding contents, all audible dialogues are to be transliterated. Interference by noise or speaker are to be indicated. Certain word categories (names, numbers, foreign words) are to be indicated. There remains some limitations of transliteration. Audible events will not be described exactly. Orthography can represent dialogues in a word level and not a phonetic description of oral utterances. For the recording of

noise and nonverbal oral sound production only limited categories may be used for identification.

3.3. Transliteration Conventions

The format of the transliteration file is defined. A file format consists of a header and turns. Each turn begins with an identification of the turn or a turn name. The turn body contains all audible events, syntactical and semantic markers and comments. At the end of a line within a line there is a carriage return and the newline is indented by one white space. Transliteration rules of turn body for Japanese text are defined. Punctuation marks separate sentences, phrases and words. For Roman script sentences, a period is set after a sentence.

3.4. Transliteration of Turn Elements

Lexical elements are defined as words, interjections, regular reduced forms of words, classified words, words with articulatory irregularities, and words with comments regarding their pronunciation. Words of a dictionary are words which may be listed in their transliterated form in the VERBMOBIL word list and can be translated to different languages. They are also well intelligible and not distorted. Interjections are short sounds inserted for the purpose of investigating surprise, affirmation, or doubt. Interjections a, ee, oo, uu, ara, ma, eeto, etto, etc., are handled as general words. Like e and ee, long duration interjections determined to be short duration interjections are regarded as "other" interjections.

Compound words are multiple words in a series. In Kanji-kana sentences, compound words are simply "proper nouns" + "common nouns" and are connected with a hyphen. Before each compound word, a marker () is added.

There are several classified lexical elements; names, numbers and foreign words. They are attached with a specific symbols. Proper nouns have a marker before them. For numbers, marker # is used before them. (Days are included.) Even when a number is used as an adverb, marker # is used. For compound words that include numbers, marker # is added in front. For foreign loan words, marker <*ENG> is used in front of them. Foreign words are non-Japanese words which do not exist in the Japanese dictionary.

Lexical elements with irregularities in pronunciation includes lengthening and poor comprehensibility of words and interruption of lexical elements. They are marked with symbols.

Interruptions of lexical elements Articulatory abortion or termination is indicated with a mark to allocation within that word. For the restating of words after they had been completely uttered markers are placed so that there is no unnaturalness even if the enclosed part is omitted.

If a pause is inserted between words and then the remaining pronunciation continues, the marker is used.

When the pronunciation is colloquial or is not correct in the case of spoken language, the correct expression is added. Colloquial pronunciation is transcribed as it is heard in bracket and the normal expression is shown in front of the colloquial pronunciation. The dictionary of terms has a base of Japanese dictionary such as 'DAIJIRIN' dictionary [4]. When a word is not found in the dictionary, it is regarded as an "incorrect expression".

3.5. Syntactic Semantic Structure

Syntactic semantic structures are markers for the structuring of the sentence flow. The transliterations try to mark regular sentences by a subset of punctuation marks. Non grammatical phenom-

ena likecorrection of sentence abortions are marked so that they may beremoved by specific text filters and correct grammatical structures in a syntactic and semantic sense are left.

Nonverbal articulatory productions and sounds For noise not expressed with words uttered by a person, an appropriate mark is placed (depending on the type). A “< >” mark is used for noise heard but with completely no meaning (not words).

A respiration is a sound made when breathing. The symbol “” is used.

For completely difficult utterances to hear, a marker is added. Poor comprehensibility of words, symbol “<%>” is used.

Filled pause or hesitations are when long duration sounds are made midway through saying a word, stopping, or inserting a vowel. Hesitation i.e. human noises and articulating human noises i.e. human non-phonemic noises are enclosed by angle brackets. Typical examples of filled pause are as follows; <uh>, <ah>, <uhm>, <ahm>, <e>, <eeto>, <ma>, <ano>. All other articulation which may not be classified among the other listed above is indicated as <hes>.

Nonverbal Articulatory Sounds

For noise not expressed with words uttered by a person, an appropriate mark is placed depending on the type. Predetermined limited symbols are used for nonverbal articulatory sounds. Those symbols are: <Smack>, <Throat>, <Cough>, <Laugh>, <Swallow>, and <Noise>.

{subsubsectionNoise and technical artifact Nonverbal Articulatory Sounds

For machine noise from recording equipment or peripherals, several limited numbers of these sounds are marked. They are <#Click>, <#Ring>, <#Knock>, <#Mtouch>(A microphone that is touched), <#Mwind>(Blowing on a microphone), <#Rustle>, <#Squeak>, and <#>(Others). When a word is cut into parts by a noise expression, XX_<marking>_XX is used.

Acoustic interference Interference By Dialog Partners

Overlapping utterances with two speakers are written down with numbers attached in the order of the utterances. At the points of overlapping utterances, application is done with the utterances of the two. Markers are added to the descriptions of both speakers. The markers are n@ and @n. n is the serial number used for overlapping voices; it starts from 1. For passive interference, n@ is put at the starting position and @n is put at the end position; corresponding to the active interference. For active interference, @n is put at the starting position and n@ is put at the end position. It is possible for many small utterances to be made between multiple sentences over the topics of conversation.

Interference by noise There may be one or more lexical elements in an utterance, so when noise covers all of those parts, markers are added. Comprising the interruption noise are the noise voiced by the speakers, the background noise, and the surrounding noise. The symbol, “< >” denotes interruption noise. In a part of a language, when noise covers the whole body, the marker “< :>” is used.

4. Romanized Transcription and Segmentation

4.1. Romanization

The collected dialogues are transcribed in kanji-kana sentences and Roman script sentences. In Roman transcription, there are two types

of transliteration. One is segmented into words and the other is without segmentation.

1. Orthographic representation

In Kanji-kana sentence, the romanization depending on Japanese orthography [5] is used. Three kinds of particles, i.e. “ha”, “wo”, “he”, which are written in the orthogonal way in the kanji-kana text are transcribed in romanization in the way to consistent to the pronunciation as “wa”, “o”, “e”.

2. Long Vowels

In the Japanese colloquial speech, some words are pronounced as if they are long vowels. Long vowels are transcribed by duplicating vowels. For example, “iu” as “yuu”, etc. A long vowel is used in the interjections and the foreign word. It is indicated by duplicated vowel. (Example: aa, @sukejuuru)

3. Double Consonant

The double consonant is transcribed by duplicating consonant in Romanization. (Example: jikkeN, chotto).

4.2. Segmentation

Word separation is undertaken based on the morphological unit of Japanese sentence. The separation of word unit is based on Japanese morphological program CHASEN [6]. In order to apply the segmented words for speech recognition and language analysis, some modification is undertaken to form compound words according to the rules on common word in the dictionary, compound words, and collapsed forms.

- Common word
Common word which is written in the dictionary is treated as one word.
- Compound word
If the compound word is used common as one word in the dictionary, it will be treated as a word. If the suffix or prefix to a common word is listed in the DAJJIRIN dictionary, it will be treated as a word.
- Collapsed form
Collapsed form of inflated verb and particle connection is treated as compound word of suffix.
- Particles
Almost all particles are segmented into words. When a particle such as “de” or “te” is connected to verb or adjective or auxiliary verb in a “Renyou” inflection, the connected words is treated as a compound word. (Example: susuN-de, itashi-te).
- Auxiliary verb
An auxiliary verb has an inflection forms. The following auxiliary verbs in the list will be listed as suffixes and treated as a compound form. “masu”, “nai”, “rareru”, “reru”, “saseru”, “seru”, “ta”, “tai”, “you” (Example: ari-mase-N, de-rareru) When an auxiliary verb such as “ta” or the inflated form of the auxiliary verb is connected to adjective in a “Mizen” inflection, the connected words is treated as a compound word. (Example: atsukat-ta). When a concatenation of inflated auxiliary verb of “Renyo” form such as “mashi”, and an inflated auxiliary verb such as “ta”, the connected words is treated as a compound word. (Example: ari-mase-N).

4.3. Example of transliteration

Example of transliteration is shown in Figure 1.

Figure 1: Example of transliteration

```
Kanji-Kana transcription j002ach1_007_BAB_150000: <> #płA <P> ZB Roman transcription with phrase segmentation j002ach1_007_BAB_150000: <ah> #ku-gatsuno kouhaNwa potsupotsuto youjiga hait-teirunode chotto <P> isogashiINdesukeredomo . Roman transcription with word segmentationj002ach1_007_BAB_150000: <ah> #kugatsu no kouhaN wa potsupotsu to youji ga hait-te iru node chotto <P> isogashiIN desu keredomo . Parsed text by trpparserj002ach1_007_BAB_150000: kugatsu no kouhaN wa potsupotsu to youji ga hait-te iru node chotto isogashiIN desu keredomo
```

5. Conclusion

We have described the method to construct and transcribe Japanese spontaneous speech data under VERBMOBIL project. The design method of spontaneous corpora in the travel arrangement task in a role playing situation has been explained in detail. Standard way of transliteration method and several rules and conventions to transcribe Japanese spoken dialogue have also described. Up to now, most of the spontaneous speech corpora have been constructed in the specific task domain in the feasible scale. When it is required to design the spoken language system in the different task domain, the problem of data collection arises and it takes much time to collect the data. The future issues will be how to make the spontaneous data corpora universal in the sense of dialog expressions and coverage of the usual usage of spoken language. It is also required to collect the unified data for the research of spoken language processing technology. There are many cases to be identified in the spoken dialogue in the pronunciation of variants association with the same spelling but different meaning. The tagging of other information than the word category, prosodic information and dialog tagging will also necessary to be used for the dialogue research.

6. Acknowledgement

The authors gratefully acknowledge the opportunity of German Ministry for Scientific Research and ATR International. I wish to thank Prof. Waibel@and colleagues of VERBMOBIL Japanese database and language processing group for their helpful suggestions and discussions on data collection of Japanese spontaneous dialogues.

References

1. D. Gibbon, et. al. ed. "Spoken Language System and Corpus Design", Mouton de Gruyter, (1998)
2. S. Burger, "Transliteration of spontaneous speech data: Manual of transliteration conventions VERBMOBIL 2 English version Release 2", (August 1997)
3. A. Kurematsu, "Roman Transcription Of Japanese Spoken Dialogue In VERBMOBIL Project", Proc. of First Int. Workshop on East-Asian Language resources and Evaluation, @pp.83-86, (May. 1998)
4. A. Matsumura, "Daijirin Dictionary", Sanseido Publishing, (1985)
5. K. Masuda, "New Japanese-English Dictionary", Kenkyuusha, (1988)
6. Y. Matsumoto, "Japanese Morphological Analysis System: CHASEN", Information Science Technical Report NAIST-IS-TR97007, Nara Institute of Science and Technology, (1997)