

LINGUISTICS-ORIENTED LANGUAGE RESOURCE DEVELOPMENT AT THE NATIONAL INSTITUTE FOR JAPANESE LANGUAGE AND LINGUISTICS

Kikuo Maekawa

kikuo@ninjal.ac.jp

Dept. Corpus Studies, National Institute for Japanese Language and Linguistics

ABSTRACT

The aim of this talk consists in the introduction to the language-resource-related activities of the National Institute for the Japanese Language and Linguistics (NINJAL). Since the last half of the 1990s, the former National Language Research Institute (NLRI) played a central role in the development of Japanese language resources by constructing corpora like *Corpus of Spontaneous Japanese* (CSJ) and *Taiyo Corpus*. In 2006, the language resource group of NLRI started a Japanese corpus compilation initiative named KOTONOHA, and set about the construction of a 100 million words *Balanced Corpus of Contemporary Written Japanese* (BCCWJ). The activity of NLRI was inherited by the NINJAL Center for Corpus Development reestablished in 2009. Now that the construction of the BCCWJ was completed successfully in August 2011, the NINJAL center set about two new projects of exploratory nature: a historical corpus project and a 10-billion-word ultra-large-scale Web-based corpus project. In addition to the presentation of the NLRI-NINJAL activities, language resource development in Japanese institutions other than NINJAL will be introduced briefly in the beginning. Also, application of the CSJ to the study of phonetics will also be demonstrated at the end.

Index Terms— NINJAL, Corpus, KOTONOHA, BCCWJ

1. OVERVIEW OF LANGUAGE RESOURCE DEVELOPMENT IN JAPAN

The aim of the current paper consists in the introduction to the corpus development activities of the National Institute for the Japanese Language and Linguistics (NINJAL). Before going directly into this theme, however, it will be appropriate to present an overview of corpus development and circulation in Japan.

In Japan, currently, there are two national institutions that played important roles in language resource development, namely NICT (National Institute for Communications and Technology) and NINJAL. These institutions were formerly known as CRL (Communications Research Laboratory) and NLRI (National Language Research Institute) respectively up to the administrative reform in 2001.

Language resource developments by NICT and NINJAL show clear contrast in two respects. First, NICT resources contain a lot of NLP tools, while most of NINJAL resources are language data (corpus and dictionary). Second, the language data of NICT are mostly machine-generated, while the NINJAL resources are human-assisted to a large extent, hence highly accurate. Needless to say, these differences reflect the different installation purposes of the two institutions, i.e., communications research and linguistics.

In addition to these national laboratories, there are language resources, tools and corpora, developed by university laboratories. The largest contributions were made by Kyoto University and NAIST (NARA Institute of Science and Technology). Probably, the most well-known contributions of these laboratories are POS taggers (JUMAN developed in Kyoto Univ., on the one hand, and, Chasen and McCab developed in NAIST) and dependency structure taggers (KNP of Kyoto, and, Cabocha of NAIST). The laboratory of Kyoto University also made important contributions to the development of language data. Their recent contribution included the release of the “Kyoto University's case frame data”, which was automatically obtained from 1.6 billion Web pages. More detailed information of these and other language resource can be obtained in [1].

A peculiarity of Japanese language resources is the lack of national center for the circulation. There is no institution like LDC in the States. Instead, multiple organizations are established by different institutions with slightly different purposes. The following three are the most active.

The GSK, or *Gengo Shigen Kyokai* (language resource association), is a non-profit corporation established in 1999. As of August 2011, there are 12 resources distributed by GSK. They include dictionaries (including above-mentioned case frame data of Kyoto University and annotated version of Iwanami Japanese dictionary 5th edition), language data (like Google's web-based Japanese N-gram data developed), and an annotation tool.

The NII-SRC (National Institute of Informatics Speech Resources Consortium) was established in 2006. As its name suggests, this institution is specified for the distribution of speech-related resources. Currently, there are 27 corpora

including RWCP speech database, Chiba University Map Task corpus, Japanese Newspaper Article Sentences (JNAS) corpus and so forth.

Lastly, ALAGIN (Advanced LAnGuage INformation Forum) was established in 2008 under the joint support of Ministry of Internal Affairs and Communications, NICT, and some universities including Kyoto University. Most of the ALAGIN resources are the ones developed in the MASTAR project of NICT [2] that inherits speech databases developed by former ATR spoken language translation laboratories.

There are also resources that are distributed by the creators of the resources. For example, most of corpora developed in NINJAL are distributed by NINJAL.

To sum up, the distribution of language resources in Japan is in the state of “balkanization”, which is not a favorable situation for the users of language resources. In my opinion, there are two closely interrelated reasons for the balkanization. To begin with, as long as the Japanese language is concerned, the market of language resources is not large enough to be maintained on commercial basis, unlike the case of English. The development and distribution of Japanese language resources need continual support of governmental fund. As the result of this, the institutions of language resources are strongly influenced by so-called ‘vertically-segmented’ administrative system. In the present case, NICT and ALAGIN are under the control of the Ministry of Internal Affairs and Communications, and, NINJAL and NII are supervised by the Ministry of Education, Culture, Sports, Science and Technology. It is very difficult, if not impossible, to conceive of a unified governmental institution that crosses the border of the administrative segment.

GSK, on the other hand, is a non-profit corporation, and therefore is neutral to the administrative segmentation. It seems, however, that this favorable position is maintained only at the cost of financial instability.

Establishment of non-governmental institution having a stable financial basis is badly needed for the healthy development of Japanese language resources.

2. TAIYO CORPUS AND CSJ

2.1. Taiyo corpus

Although NINJAL (including NLRI) has a long tradition of statistical lexical survey that started in the beginning of 1950s, it was in the middle of 1990s that NLRI researchers began thinking seriously about the development of language corpus. The crucial difference between the traditional lexical survey and the recent corpus is in the open-accessibility to the sampled data. In the lexical surveys, it was only the results of the survey as summarized in the form of word frequency tables that were publicly accessible. Researchers outside the NLRI could not make access to the sampled data.

The first NLRI language corpus that was designed for public release was the *Taiyo Corpus*, which was a full-text corpus of a general-interest magazine named *Taiyo* (the sun) published in the years 1895-1925. This is the period of time when modern Japanese has changed quickly its writing style, from classical to colloquial. The compilation of *Taiyo Corpus* started in the middle of 1990s and ended in 2005. About 7 million words texts were annotated with respect to text structure, bibliographical information, and various features of writing such as errata in orthography using the tags in XML format. Note that the corpus was not POS analyzed because of extremely complex and changing writing system.

2.2. Corpus of Spontaneous Japanese

Subsequently, compilation of the CSJ, or *Corpus of Spontaneous Japanese*, started in 1999 as a part of a national project ‘Spontaneous Speech: Data and Processing Technology’. This project was a joint project of NLRI and CRL aiming at the development of new-generation automatic speech recognition (ASR) system for spontaneous speech, and was supervised by Professor Sadaoki Furui of Tokyo Institute of Technology.

CSJ is a richly annotated corpus of 7.5 million words or 652 hours [3]. In addition to the dual POS analysis (see 4.2 below as for dual POS analysis) which was applied to the whole CSJ, detailed phonetic and prosodic annotation was applied to the subset known as the CSJ-Core (about a half million words or 44 hours) using the X-JToBI annotation scheme developed for the CSJ [4].

The acoustic- and language-models developed using the CSJ showed remarkable success in the ASR of spontaneous speech. It turned out that the use of CSJ for the machine learning of acoustic- and language-models improved the mean word recognition rate from 50% to 70%; an extraordinary leap in the field of ASR [5]. The success of CSJ gave strong impetus for the establishment of new research fields like spoken document retrieval or spoken document summarization.

The CSJ was publicly released in June 2004, and as of February 2011, more than 450 copies of the corpus are used in academia and industry, covering both engineering and humanities, home and abroad. As far as I know, there are more than 800 papers that made reference to the CSJ.

3. KOTONOHA INITIATIVE

The remarkable success of CSJ in speech engineering and the high reputation of *Taiyo Corpus* among the specialists in Japanese linguistics gave strong impetus for NLRI to recognize language resource development as one of the most important missions of the institute.

In 2005, KOTONOHA Initiative was publicly announced by NLRI (Fig.1). This was a grand design of the corpus development by the NLRI in the coming 10-15 years.

This initiative covered both written (the upper half of the figure) and spoken (lower half) varieties and the period of time between 1867 and the present, i.e., after the beginning of Meiji era (or the end of Tokugawa feudal government).

At this time, the most important corpus to be compiled was a balanced corpus of contemporary Japanese; a Japanese counterpart of *British National Corpus*. Lack of such a corpus caused several serious problems in the corpus-based analysis of Japanese.

At the time we designed BCCWJ, most of corpus-based linguistic analyses were based upon electronic archives of newspaper articles, but in Japanese, as well as in many other languages, newspaper texts are very specific from a point of view of stylistics. For example, use of kanji (Chinese logograph) in newspaper articles is too systematic and lack the variation observed in ordinary texts, hence inappropriate for the study of variation in Japanese writings.

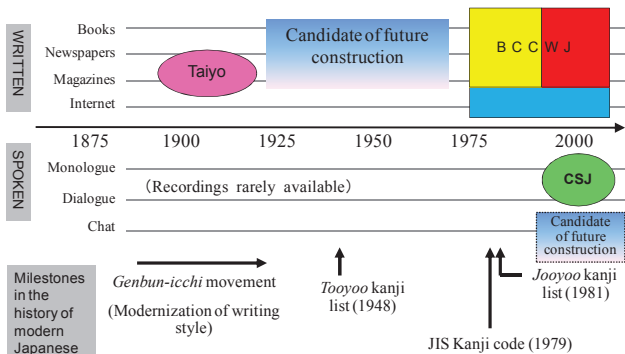


Fig.1: The KOTONOHA Initiative

There was one more text archive that was used frequently by linguists. *Aozorabunko* was a collection of copyright-expired literary texts. This is a very important contribution for the study of modern Japanese literature, but copyright-expired texts are too old to be the main resource for the study of contemporary Japanese.

Moreover, lack of balanced corpus caused serious problems in the fields of science like psycholinguistics and brain science of language, where the occurrence frequencies of the words and phrases used as the stimulus of experiment are needed to be controlled.

4. BCCWJ

4.1. Structure and size

BCCWJ, or *Balanced Corpus of Contemporary Written Japanese*, was designed in 2005 to be Japan's first balanced corpus [6]. It consists of 3 main subcorpora. Publication Subcorpus (or Production subcorpus) consists of texts randomly sampled from the populations of books, magazines, and newspapers published during the years 2001-2005, whose total size is about 35 million words.

Library subcorpus (or circulation subcorpus) contains samples of books found in public libraries. The statistical population consists of the totality of books that are registered in more than 13 public libraries of Tokyo metropolis; the size of this population is almost equal to that of the books in the publication sub-corpus. This subcorpus contains about 30 million words.

Lastly, special-purpose subcorpus is a collection of various mini corpora that were indispensable for the purposes of linguistic and language planning studies of the Japanese language. This subcorpus contains the samples of texts in the Web (bulletin board *Yahoo! Chiebukuro* and *Yahoo! blog*), law, governmental white papers, textbooks in elementary, junior high- and high schools, minutes of the National Diet, verse, and various reports issued by local governments for public relations purposes. See Table 1 for more details.

Table 1: Size of the BCCWJ subcorpora

SUBCORPUS	REGISTER	N SAMPLE	N WORD (K word)
PUBLICATION SUBCORPUS	Books	10,105	28,490
	Magazines	1,989	4,430
	Newspapers	1,479	1,380
LIBRARY SUBCORPUS	Books	10,461	30,110
SPECIAL PURPOSE SUBCORPUS	White papers	1,500	4,880
	Textbooks	412	930
	Local government reports	354	4,000
	Bestselling books	1,377	3,710
	Bulletin board texts	91,445	10,280
	Blog texts	52,680	10,270
	Verse	253	230
	Law	346	1,000
Minute of Diet	159	5,100	
TOTAL		172,560	104,810

4.2. Annotation

BCCWJ is annotated with respect to POS information, character and text structure information, and, bibliographical information. Like CSJ, the whole texts of BCCWJ are dually POS analyzed by using the SUW-LUW dual POS systems. SUW, or short-unit-word, approximates the size of entry words found in ordinary Japanese dictionaries, while LUW, or long-unit-word, contains compounds of various sizes. This dual analysis is necessary because Japanese is a so-called agglutinative language, and its writing system has not established the social convention of word segmentation.

For example, "kokuritsukokugokenkyuujo" (NINJAL) is a single LUW consisting of {kokuritsu} (national), {kokugo} (national language), kenkyuu {research}, and {jo} (institution), where {jo} is a suffix and all others are nouns. Similarly, "arukidasu" (start walking) is a compound verb consisting of two SUW, {aruku} (walk) and {dasu} (start an action). Moreover, Japanese has a lot of compound particles. For example, a LUW {nioite} (at) is a compound particle consisting of three SUW, namely {ni} (case particle), {oku}

(verb, to place), and {te} (conjunction particle). Notice that all numbers of words are counted by means of SUW in this paper. Main XML tags used in the BCCWJ annotation are summarized in Table 2.

Table 2: Main tags used in the annotation of BCCWJ

TYPE OF TAG	TAG	USAGE
Sample	sample	Covers the totality of a sample
Hierarchical structure	artice	Text written by a single author
	cluster	Extent covered by a title tag
	title	
	list	Extent of enlisted elements
	paragraph	
Tables & Figs	figure	
	caption	Caption of a figure
Citations	citation	
	speech	Extent of a direct speech
	quote	
Notes	noteBody	Text of a note
Miscs	abstract	
	authorsData	
	verse	Haiku, waka, etc.
Characters	ruby	
	correction	Errata in the original text
	missingCharacters	Characters outside JIS X02132004

4.3. UniDic

A very important by-product of BCCWJ is UniDic, a SUW-based machine-readable dictionary for natural language and speech processing [7]. Important characteristics of UniDic include coherence in SUW analysis and adaptation to the writing variation.

As for the former problem, the the Chasen POS tagger analyzes {kokuritsukokkaitoshokan} (National Diet Library) into two words like {kokuritsu} (national) and {kokkaitoshokan} (Diet library), while the same tagger analyzes {kokuritsukoobunshokan} (National Archives of Japan) into 4 words as {kokuritsu} (national), {koobunsho} (official records), and {kan} (house), reflecting the inconsistent word boundary definition in the IPA dictionary, which is the default dictionary of Chasen.

If UniDic is used instead of IPA dictionary, the same tagger yields the following results: {kokuritsu}, {kokkai}, {toshokan}, {kan}, and, {kokuritsu}, {koo} (official) {bunsho} (documents), and {kan}. In UniDic analyses, affix {kan} and {koo} are coherently extracted, while IPA dictionary treated them inconsistently depending, perhaps, on the frequencies and/or familiarities of the words in question.

As for the second problem, UniDic systematically registers the writing variations of SUWs as observed in the BCCWJ data. The writing variants are classified with respect to the variation of word-form and that of writing forms. For example, adverb lexeme {yahari} (after all) has two variants /yahari/ and /yappari/ at the level of word-form, and, there are 3 variants at the level of writing form, because the word-form {yahari} can be written either in hiragana (や

はり) or combination of kanji and hiragana (矢張り). Currently, UniDic registers about 220,000 lexemes (SUW).

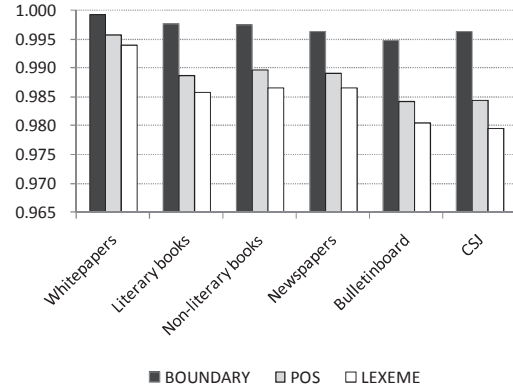


Fig.2: Performance of POS analysis by MeCab + UniDic

Fig.2 compares the performance of SUW analysis by means of MeCab tagger and UniDic as a function of text registers. The performance was shown with respect to 3 criteria: the correctness of SUW boundary detection, the correctness of POS information assignment, and, the correctness of lexeme assignment (in the case of homographs). Note that correct lexeme assignment implies correct boundary detection and POS assignments. The performance degenerates gradually from left to right, but as long as the written registers are concerned, performance of lexeme assignment is higher than 98%.

4.4 Web interface

POS analyzed data of the whole BCCWJ can be retrieved via a Web interface named Chuunagon [8]. Users can make choice of 3 ways of data retrieval, i.e., full-text search, SUW search, and LUW search. In the SUW or LUW search, users can specify the N-gram of SUW or LUW up to N=10 using the POS information of the corpus. So, for example, it is possible to make query of sentences whose predicates consist of an adjective followed immediately by an auxiliary verb which is immediately followed by a phrase-final particle. In this particular case, 18,763 hits are obtained in about 90 seconds. All retrieved samples can be downloaded to the user's local PC.

Chuunagon can be used free of charge, but users are requested to make contract with the NINJAL to obtain user licenses. There is another Web interface to BCCWJ, known as Shoonagon that provides only the function of full-text search [9]. *Shoonagon* is free of charge and no contract is necessary, but users cannot download the retrieved results.

5. REVISION OF KOTONOHA

Now that the public release of the BCCWJ was almost done (the DVD version of the corpus will be released in October

2011), it is time to start new corpus development projects. Originally, as shown in Fig.1, compilation of a balanced corpus filling the temporal gap between the BCCWJ and Taiyo, or a corpus of highly spontaneous conversation were expected to get started upon the completion of BCCWJ.

Administrative reform of NINAJL that took place in 2009, however, forced us to revise KOTONOHA initiative and start two new projects.

One of them is the development of historical corpora that, in the long run, will cover the whole history of the Japanese language. Historical study has been outside the scope of research by NLRI=NINJAL since its establishment in 1948, but the 2009 reform changed drastically the purpose of establishment of the institute so that historical study of Japanese was regarded to be one of the main research themes of the institute.

The historical corpus project currently conducts a feasibility study using the electronic texts of 18 literary works provided by the courtesy of *Shogakkan* publishing company including *Tosa nikki*, *Ise monogatari*, *Genji monogatari*, *Makuransoushi* etc. The objectives of the feasibility study includes, among other things, developments of POS taggers for classical and middle Japanese, and definition of XML format for annotated texts.

Another new project is the development of ultra-large-scale corpus of the present-day Japanese. By ultra-large-scale is meant the corpus size of the order of 10 giga words, (i.e., hundred times larger than the BCCWJ). Needless to say, a corpus of this size can't be developed by the same compilation techniques used for BCCWJ. The corpus can be developed only by using the Web as the source of texts.

This corpus size is needed for several reasons. For example, the number of lexical items covered by the BCCWJ is not sufficient for corpus-based dictionary compilation. Although more than 220,000 lexemes can be found in the BCCWJ, 75,000 occurred only once in the corpus. The numbers of lexemes for which BCCWJ can provide more than 100 and 50 examples are 31,000 and 44,000 respectively; while most of today's middle-sized Japanese dictionaries cover more than 70,000 lexical items.

Ultra-large-scale corpus is also needed for systematic examination of the possibility of word combinations. For example, Japanese adverb {sura} and {sae} (both mean 'even' or 'just') occurred 54,000 and 4,300 times in the BCCWJ, but their combination "sura sae" did not occur in the corpus. It doesn't mean, however, the multiword expression doesn't exist. The phrase like "sore ga nani o imi shite iru no ka sura sae wakara nai" (I don't understand even what is meant by that) is natural and can be found in the Web.

6. APPLICATION OF CSJ TO PHONETICS

Lastly, some results of the application of CSJ to the phonetic analysis of spontaneous speech will be presented. Phonetic

study of spontaneous speech was one of the main application domains conceived at the time of corpus design.

6.1 Articulatory variations of /z/ and /b, d, g/

Japanese /z/ phoneme is known to vary between affricate [dz] and fricative [z]. Following the influential book of Shiro Hattori published in 1951 [10], this variation has long been explained as a conditional variation in terms of the phoneme's location in a word; affricate in word-initial position and fricative elsewhere.

Conjoined analysis of the morphological, prosodic, and segmental annotations of the CSJ-Core revealed that this explanation was far from being perfect. It was rather the time that speaker can use for the articulation of the phoneme that plays fundamental role in the realization of /z/. Regardless of the phoneme's location in word, the phoneme tends to be realized as an affricate if there is enough time for articulation.

The time allotted for consonant articulation (TACA) increases naturally when speaking rate diminishes or when the phoneme in question is preceded by a pause, but these are not the only factors of increase in TACA. It increases considerably also when the phoneme in question is preceded by phonemes whose points of articulation are underspecified but assimilate to that of the following consonant (i.e., the /z/). In Japanese, there are two such phonemes, /N/ (morai nasal) and /Q/ (geminate). The rate of affricate articulation increases considerably, when /z/ is preceded by these consonants [11].

By the way, voiced plosives /b/, /d/, and /g/ are often realized as homorganic voiced fricative, viz., [β], [ð], and [ɣ] in Japanese. The mechanism of this weakening can be explained by the same device as in the case of the /z/ [12]. Fig. 3 shows the relationship between the TACA and the rate of stop (in the case of plosives) or affricate (/z/) realizations of consonants.

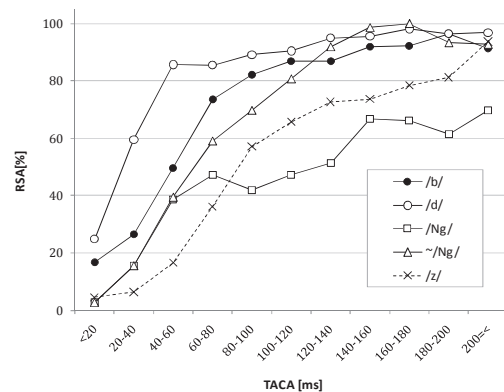


Fig. 3: TACA and rate of stop articulation (RSA).

In this figure, the data of /g/ was sub-classified into two cases where /g/ is preceded by a moraic nasal (/Ng/) and otherwise (~Ng/). This was because in Standard Japanese,

/g/ is mostly realized as velar nasal [ŋ] after a /N/. It is very interesting that the amount of TACA needed to achieve 50% RSA increases in the order of /d/ < /b/ < /g/. This order seems to reflect the differences of phonological complexity at the corresponding places of articulation. In Japanese, there is 3-way (/d/~z/~n/) contrast at alveolar and 2-way contrast at bilabial (/b/~m/), but there is no phonological contrast at velar position (/g/ only).

6.2 Linguistic function of a BPM

Standard Japanese has a rich inventory of phrase final intonation or BPM (boundary pitch movement). It was pointed out in 1959 by Ooishi that there is a variant of rising-falling (HL%) BPM with respect to the timing of F0 peak. While in ordinary HL%, the F0 peak is located in the final mora of the phrase, there is a variant whose F0 peak is located in the penultimate mora. The latter type is called PNLP, or penultimate non-lexical prominence. The linguistic function of PNLP has been an open question during the last 50 years.

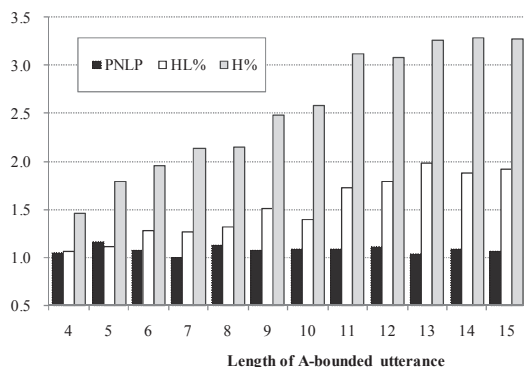


Fig.4: Frequencies of BPM and utterance length.

Analysis of the CSJ-Core provides new insight into this long-standing problem [13]. Fig. 4 compares the occurrence frequencies of the rising (H%), ordinary rising-falling (HL%), and PNLP version of rising-falling BPMs in utterances punctuated by typical sentence ending forms. For each of the 3 BPMs, only the utterances that contained at least one occurrence of the BPM in question were analyzed. Mean occurrence frequencies (ordinate) of the H% and HL% increase nearly monotonically as the length of utterance (abscissa, as measured by the number of constituent accentual phrases) increases, while the frequency of PNLP stays unchanged regardless of the utterance length. This finding suggests strongly the interpretation that PNLP has a kind of “culminative” function. Moreover, analysis of the location of PNLP in an utterance revealed that the mean occurrence probability of PNLP became the highest in the penultimate accentual phrase of an utterance. This finding suggests the interpretation that PNLP has the function of foretelling the end of an utterance, most presumably before a

topic boundary. From a point of view of discourse prosody, the behavior of PNLP is quite interesting in that it provides the case where local tonal event, as opposed to global prosodic manipulation like pitch range control, plays crucial role in the transmission of discourse structure.

6.3 Prosodic registers

The last application is the automatic classification of speech register by use of the frequency information of prosodic events. It is widely believed by laymen and specialists that prosody differs systematically according to the purpose and/or situation of speech, but no quantitative evidence has been presented so far.

Discriminant analysis of the normalized relative occurrence frequencies of all X-JToBI labels and speaking rate revealed that it was possible to classify 4 speech registers (academic presentation, extemporaneous speech of daily topic, interview speech, and, reading aloud of transcribed monologue speech) with about 85% accuracy (closed data). Moreover, it turned out that it was possible to classify the speech register with about 75% accuracy by use of the subset data of only 60 seconds long extracted from each of the original speeches [14].

These results provide important suggestion for the design of a new corpus for the study of register diversity of spontaneous speech.

7. REFERENCES

- [1] http://nlp.kuee.kyoto-u.ac.jp/NLP_Portal/lr-cat-e.html
- [2] <http://mastar.jp/index-e.html>
- [3] Maekawa, K. "Corpus of Spontaneous Japanese: Its Design and Evaluation", *Proc. SSPR2003*, 7-12, 2003
- [4] Maekawa et al. "X-JToBI: An extended J_ToBI for spontaneous speech", *Proc. ICSLP2002*, 1545-1548, 2002.
- [5] Furui, S. "Recent progress in corpus-based spontaneous speech recognition." *IEICE Trans. Information and Systems*, E88-D (3), 366-375, 2005.
- [6] Maekawa et al. "Design, Compilation, and Preliminary Analyses of Balanced Corpus of Contemporary Written Japanese." *Proc. LREC 2010*, 1483-1486, 2010.
- [7] Den, Y. et al. "A proper approach to Japanese morphological analysis: Dictionary, model, and evaluation." *Proc. LREC2008*, 1019-1024, 2008. (See also <http://download.unidic.org/>)
- [8] <https://chunagon.ninjal.ac.jp/>
- [9] <http://www.kotonoha.gr.jp/shonagon/>
- [10] Hattori, S. *Onseigaku*. Tokyo; Iwanami, 1951.
- [11] Maekawa, K. "Coarticulatory reinterpretation of allophonic variation: Corpus-based analysis of /z/ in spontaneous Japanese." *J. of Phonetics*, 38(3), 360-374, 2010.
- [12] Maekawa, K. "Weakening of stop articulation in Japanese voiced plosives." *J. Phonetics Soc. of Jpn*, 14(2), 1-15, 2010.
- [13] Maekawa, K. "Phonetic shape and linguistic function of penultimate non-lexical prominence." *J. Phonetic Soc. of Jpn*, 15 (1), 1-13, 2011.
- [14] Maekawa, K. "Discrimination of speech registers by prosody." *Proc. 17th ICPHs*, 1302-1305, 2011.