

ACOUSTIC AND PERCEPTUAL EFFECTS OF LISTENER ADAPTIVE TEMPORAL ADJUSTMENTS IN DIALOGUE

S. Imaizumi¹⁾, A. Hayashi²⁾, and T. Deguchi²⁾

¹⁾ RILP, University of Tokyo, ²⁾ Dept. of Education, Tokyo Gakugei University, Japan

ABSTRACT

Effects of listener adaptive temporal adjustments in dialogue were investigated by analyzing vowel devoicing in the speech of teachers directed to hearing-impaired (HI) or normal-hearing (NH) children, and read speech (RD). The teachers did reduce the devoicing rate more in the HI vs. NH and RD samples in such a manner that contrasts between the highly devoicable mora groups versus the others are enhanced within phonological and phonetic constraints of Japanese.

INTRODUCTION

Listener-oriented adaptation of speaking style appears to affect various stages in speech production processes. Our previous analyses [1, 2] of dialogues between professional teachers and normal-hearing (NH) or hearing-impaired (HI) children found that the teachers tended to use simpler and shorter sentences for the HI children than for the NH ones. They also reduced their speaking rates by inserting longer pauses at phonological phrase boundaries and producing longer syllable durations. The teachers also reduce their vowel devoicing, probably to improve the listener's comprehension.

The purpose of the present paper is to elucidate acoustical and perceptual effects of listener dependent adjustments of speaking style by analyzing vowel devoicing in dialogue between teachers and NH or HI children. The main focus was put on the relations between listener dependent adjustments of speaking style and a phonological constraint.

METHOD

Recording of Dialogues

Dialogues were recorded during a simple picture-searching game through which a teacher attempted to assess the

speech communication ability of a HI or NH child.

Two different panels were prepared, A or B, with each displaying 11 pictures (illustrations) of boys/girls labeled with their names. The A panel was set in front of the child and a copy of it in front of the teacher. The teacher instructed the child to point to a picture as fast as possible after a name was called out. The teacher randomly called out all the names one by one.

The question was fixed as "Donokoga /CVCVCV/ desuka?" (Who is /CVCVCV/ ?), where /CVCVCV/ represents the name of a picture. If the teacher mistakenly used a different form of question, the sample was not used.

Recording of Read Speech

To clarify the differences between dialogue and read speech, a read list was also recorded and analyzed. Six teachers read the target sentences "Donokoga /CVCVCV/ desuka?" five times at three tempo of fast (RDF), normal (RDN) and slow (RDS). The abbreviation RD is used to represent the read speech tokens.

Analyzed Samples

The names of the pictures consisted of target moras used to analyze the structure of dialogue. Each name consisted of three moras, that is, /C₁V₁C₂V₂C₃V₃/, where one mora is formed by one consonant C and one vowel V.

Six types of moras were analyzed, i.e., AffIni, FriIni, StoMed₁, AffMed, FriMed, and StoMed₂, which represent the manner of articulation of the component consonant (fricatives, affricates, or stops) and mora position (initial or medial). Accent was placed on the initial mora. AffIni and FriIni were the initial moras followed by /ki/, while StoMed₁ was /ki/ following the narrow

vowels /i/ or /u/. Both the initial and medial moras can be devoiced for AffIni, FriIni, and StoMed₁. All the moras in AffMed, FriMed, and StoMed₂ were preceded by an open vowel /a/. StoMed₁ and StoMed₂ were treated separately because devoicing can be different depending on whether the preceding mora can be frequently devoiced (StoMed₁) or not (StoMed₂).

Subjects

Six professional teachers and seven corresponding HI or NH children participated in the test. All were speakers of the Tokyo dialect of Japanese.

Measurements

All the target moras, totally 2740, were examined using an acoustic analysis system[1]. For each target mora sample, M_m , the length of the unvoiced segment, U_m , and that of the voiced segment, V_m , and their sum, $L_m = U_m + V_m$, were measured. M_m was classified as "Voiced" unless $V_m = 0$. For each sample, the classification variable *Voice* was as either "Devoiced" or "Voiced."

Modeling

Four classification variables (*Mode*, *Mora*, *Teacher* and *Voice*) were defined as follows: *Mode* with five levels (HI, NH, RDF, RDN, RDS), *Mora* with six levels (AffIni, FriIni, StoMed₁, AffMed, FriMed, StoMed₂), *Teacher* with six levels (TE1 - TE6), and *Voice* with two levels (Devoiced, Voiced). Each mora sample was characterized by the continuous variables U_m , V_m , and L_m , and by the classification variables of *Mode*, *Mora*, *Teacher*, and *Voice*.

A four-dimensional contingency table, F_{ijkl} , was constructed first. F_{ijkl} represents the frequency of samples classified at the cell, C_{ijkl} , of the i -th *Mode*, j -th *Mora*, k -th *Teacher*, and l -th *Voice*. The devoicing rate, $Pr(C_{ijkl})$, was calculated as the ratio of the number of devoiced moras versus the total number of moras for each cell, i.e., $Pr(C_{ijkl}) = F_{ijk1} / (F_{ijk1} + F_{ijk2})$.

Two statistical models were constructed, i.e., a generalized linear model (GLM) describing the

relationship between the mora length and classification variables, and a logistic regression model predicting the devoicing rate using the mora length and classification variables [3].

Perceptual Analyses

The perceptual characteristics of the tokens were analyzed using the semantic differential method. As previously reported, 24 pairs of adjectives were used as 9-point dipole rating scales. The listening subjects were 8 normal hearing students. The tokens used were 30 samples of /Donokoga hikita desuka/? spoken by the 6 teachers in the five modes. Obtained rating scores were analyzed by a principal factor analysis, and then a regression analysis was carried out to extract any significant correlations with the temporal structure of the speech.

RESULTS AND DISCUSSION

Mora Length Adjustments

The ANOVA obtained by the GLM procedure showed that L_m was significantly affected by the four classification variables (*Mode*, *Mora*, *Teacher* and *Voice*).

Figure 1 shows bar plots for the total mora length, L_m , with respect to the *Mode* vs. *Mora*. As shown in Fig. 1, the teachers significantly lengthened the moras in speech directed to the HI children vs. the NH children and read speech. FriIni had the longest mora length which was significantly longer than the other mora groups. There was no significant difference in L_m between StoMed₁ vs. StoMed₂, AffIni vs. AffMed.

Devoicing Rate Variations

Figure 2 shows the predicted logistic regression curves for the devoicing rate of the HI, NH, and RD tokens. As L_m increases, the predicted devoicing rate clearly decreases more for HI than for NH, and even more than for RD, which confirms our previous report [1].

Some common tendencies, however, were observed regardless the modes. The accented initial mora groups, FriIni and AffIni, tended to have a lower devoicing rate than the medial

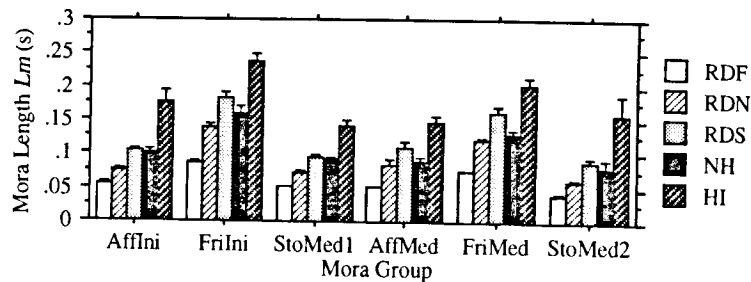


Figure 1. Interaction bar plot for mora length L_m between Mode vs. Mora.

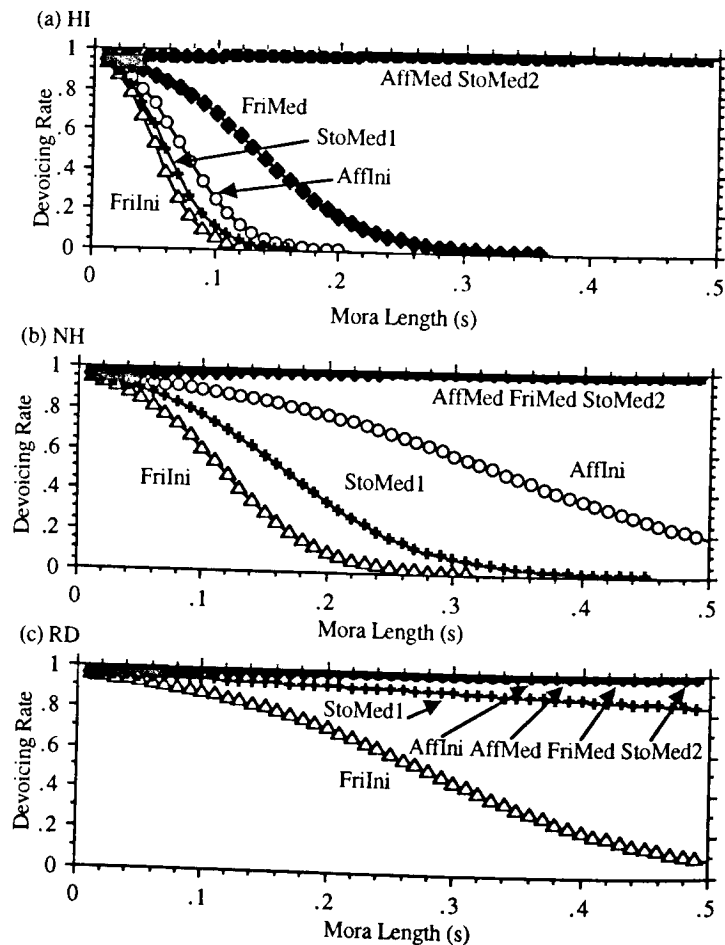


Figure 2. The predicted devoicing rate by the CATMOD logistic regression analysis for the modes of (a) HI, (b) NH, and (c) RD. The mode RD includes RDF, RDN and RDS.

mora groups, FriMed, AffMed and StoMed2. FriIni, the accented initial moras followed by /ki/, had the lowest devoicing rate, while AffMed and StoMed2, the medial unaccented moras preceded by an open vowel /a/, had the highest devoicing rate.

Furthermore, the devoicing rate was significantly different between StoMed1 and StoMed2. Both were the unaccented medial /ki/. StoMed1 was preceded by a high vowel which was frequently devoiced, while StoMed2 by an open vowel /a/ which was seldom devoiced. This significant difference in the devoicing rate cannot be accounted for by the mora length variation because there was no significant difference in L_m . This result suggests that the devoicing rate depends on the devoicing probability of the preceding vowel.

These common tendencies observed regardless the modes may be explained by a phonological rule proposed by Kondo [4]. She showed that vowels tended not to be devoiced consecutively over two moras to avoid creating series of consonant clusters on the surface level, which is not a favored structure in Japanese.

From our point of view, these common tendencies further suggest that the teachers reduced the devoicing rate more in the HI vs. NH samples, and even more in the HI vs. RD samples, within the phonological constrain. The mora groups, AffMed and StoMed2, which were highly devoicable from the phonological point of view were kept devoiced even when the teachers tried to talk carefully to HI children, while the others were highly voiced. The teachers did reduce the devoicing rate so as to enhance the contrasts between the highly devoicable mora groups versus the others within the phonological constraint of Japanese. Connecting the results of our previous report [1], listener-oriented adaptation of devoicing occurred within phonological and phonetic [2, 5] constraints of Japanese.

Perceptual Characteristics

The perceptual profiles of tokens could be represented by four factors F1, F2, F3 and F4. F1 represents the contrast between discomfort ("Rough, Uneasy,

Busy,") and pleasant ("Easy, Kind, Friendly, Restful, Polite"), corresponding to the perceptual difference between the RD and the other modes (NH and HI). F3 represents the contrast between "Slow, Stiff, Unnatural, Intelligible, Strong" and "Busy, Lifeless, Tense, Rough, Dull," corresponding to the differences between HI and the other modes (RD and NH). F2 and F4 could be interpreted as representing differences among the teachers. These results suggest that listener-oriented adaptation of speaking style produced significant perceptual effects.

CONCLUSION

The teachers did reduce the devoicing rate more in the HI than NH and RD samples in such a manner that contrasts between the highly devoicable mora groups versus the others are enhanced within phonological and phonetical constraints of Japanese. Listener-oriented adaptation of speaking style created significant acoustical and perceptual effects.

ACKNOWLEDGMENTS

Sincere gratitude is extended to all the teachers and students at Ohji Elementary School. This research was supported by a Grant-in-Aid for Scientific Research on Priority Areas of "Spoken Dialogue," Ministry of Education, Science and Culture, Japan.

References

- [1] Imaizumi, S., et al. (1993). "Listener adaptive characteristics in dialogue speech," *Proc. of ISSD*, (Waseda University Printing, Tokyo), 279-282.
- [2] Imaizumi, S., et al. (1995). "Listener adaptive characteristics of vowel devoicing in Japanese dialogue," *J. Acoust. Soc. Am.*, in press.
- [3] SAS Institute Inc. (1989), *SAS/STAT User's Guide*, Ver. 6, Fourth Edition (Cary, NC, USA).
- [4] Kondo, M. (1994). "Mechanisms of vowel devoicing in Japanese," *Proc. of ICSLP 94* (Yokohama, 1994), 1, 61-64.
- [5] Jun, S-A, and Beckman, M. (1993). "A gestural-overlap analysis of vowel devoicing in Japanese and Korean," *Annual Meeting of the Linguistic Society of America*, (Los Angeles, USA, 1993).