

# Speaker-based Accented English Clustering Using a World English Archive

H.-P. Shen<sup>1,2</sup>, N. Minematsu<sup>2</sup>, T. Makino<sup>3</sup>, S. H. Weinberger<sup>4</sup>, T. Pongkittiphan<sup>2</sup>, C.-H. Wu<sup>1</sup>

<sup>1</sup> National Cheng Kung University, Tainan, Taiwan

<sup>2</sup> The University of Tokyo, Tokyo, Japan

<sup>3</sup> Chuo University, Tokyo, Japan

<sup>4</sup> George Mason University, Virginia, USA

<sup>2</sup>{happy,mine,bank}@gavo.t.u-tokyo.ac.jp, <sup>3</sup>mackinaw@tamacc.chuo-u.ac.jp,

<sup>4</sup>weinberg@gmu.edu, <sup>1</sup>chwu@csie.ncku.edu.tw

## Abstract

English is the only language available for global communication. Due to the influence of speakers' mother tongue, however, those from different regions often have different accents in their pronunciation of English. The ultimate goal of our project is automatic creation of a global pronunciation map of World Englishes on an individual basis, for speakers to use to locate similar English pronunciations. Creating the map mathematically requires a matrix of pronunciation distances among all the speakers considered. Our previous study proposed a good algorithm for that purpose [1], where, using phonetic reference pronunciation distances calculated from labeled data, a pronunciation distance predictor was trained and built for unlabeled data. Due to space limit in [1], the procedure for calculating the reference distances was not described in detail. Then in this paper, detailed descriptions are given and 498 world-wide native and non-native speakers in the Speech Accent Archive [2] are clustered using the phonetic reference distances. Results show high validity of using the calculated distances as reference distances for training a distance predictor.

**Index Terms:** World Englishes, IPA transcription, DTW, Speech Accent Archive, phonetic pronunciation clustering

## 1. Introduction

English is the only language available for global communication and it is true that English communication is done quite often between non-native speakers in international occasions. Due to the influence of the speakers' mother tongue, those from different regions inevitably have different accents in their pronunciation. Recently, more and more users of English accept the concept of World Englishes [3,4,5,6] and they regard US and UK pronunciations as just two major examples of accented English. Diversity of World Englishes is found in various aspects of speech acts such as dialogue, syntax, pragmatics, lexical choice, pronunciation etc. Among these kinds of diversity, this paper focuses on pronunciation. If one takes the philosophy of World Englishes as it is, he can claim that every kind of accented English is equally correct and incorrect. In this situation, there will be a great interest in how one type of pronunciation is *different* from another, not in how that type of pronunciation is *incorrect* compared to US or UK pronunciation. As shown in [7], the intelligibility of spoken English depends on the nature of the listeners as well as that of the speaker and the spoken content, and foreign accented English can indeed be more intelligible than native English. Generally speaking, speech intelligibility tends to be enhanced among speakers of similarly accented pronunciation.

The ultimate goal of our project is automatic creation of a global map of World Englishes on an individual basis, for a

speaker to use to locate similar Englishes and to find where his pronunciation is located in the diversity of English pronunciations. If the speaker is a learner, he can then find the best and easiest-to-communicate English conversation partner. A learner can also know how his pronunciation compares to other varieties. If he is too distant from these other varieties, he may have to correct his pronunciation for the first time to achieve smoother communication with these others. In real-world application, the global but individual pronunciation map may be popularized to the world of international business. Here, people often encounter new types of accented English pronunciation, some of which may be very problematic and cause some miscommunication. With this map, however, one can know in advance how his pronunciation is different from his new business partner's. He may find his colleague whose pronunciation is similar to that partner's and ask the colleague for help.

For our project, however, we have two major problems. One is collecting data and labeling a part of them, and the other is creating a good algorithm of automatically drawing the global map for a huge amount of unlabeled data. Luckily enough, for the first problem, the fourth author has made a good effort in systematically collecting World Englishes from more than a thousand speakers from all over the world. This corpus is called the Speech Accent Archive (SAA) [2], which provides speech samples of a common elicitation paragraph with their narrow IPA transcriptions. The technical challenge in the second problem is that we need an algorithm that can focus exclusively on pronunciation differences between speakers by ignoring irrelevant differences such as those in age, gender, vocal tract length, etc. In our previous study [1], by using reference pronunciation distances calculated based on the IPA transcriptions, we built a pronunciation distance predictor using invariant pronunciation structure analysis. The invariant structure analysis was proposed in [8][9] inspired by Jakobson's structural phonology [10] and it can extract very robust features. The structural features were already introduced to various tasks such as pronunciation scoring [11], pronunciation error detection [12], language learners clustering [13], dialect analysis [14], automatic speech recognition [15,16], and speech synthesis [17]. In our previous study [1], our pronunciation distance predictor outperformed by far a baseline system that was built with a conventional HMM-based phoneme recognizer. Due to space limit in [1], however, the procedure for calculating reference distances was not described in detail. In this paper, detailed descriptions are given and 498 world-wide speakers in the Speech Accent Archive are clustered using the phonetic reference distances. For comparison between two IPA transcriptions, we adopt the DTW algorithm and the obtained alignment gives us a phonetic distance between them. For DTW, a phone-to-phone distance matrix is required and this is obtained through acoustic analysis of an expert phoneti-

cian's productions of all the IPA phones with/without a diacritic mark. It should be noted that pronunciation diversity of World Englishes is found in both segmental and prosodic aspects. In our previous study [1], reference distance was obtained by calculating distance between a pair of IPA transcriptions. This means that the reference distance in [1] ignored the prosodic diversity because IPA transcription gives us only phonetic information of a given utterance. We do not claim that the prosodic diversity is minor but, as will be shown in the current paper, it seems that the clustering of English users only based on the segmental aspect can still present visually and validly how World Englishes are diverse in terms of pronunciation.

This paper is organized as follows. In the following two sections, we describe the SAA corpus and how to estimate phone-to-phone distance information by acoustic analysis. In section 4, we explain how to estimate inter-speaker distances by using the DTW algorithm. Some results of speaker-based pronunciation clustering are presented in section 5. In section 6, a distance predictor constructed using the above inter-speaker distances is briefly introduced and its performance of inter-speaker distance prediction is shown. In section 7, this paper is concluded and future directions are also presented.

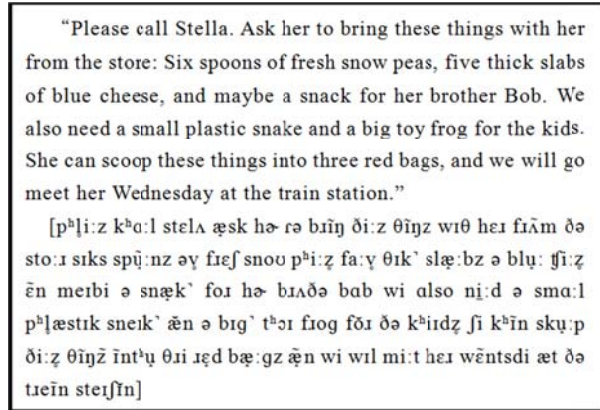


Fig. 1 The elicitation paragraph used in the SAA and an example of detailed IPA transcription with diacritic marks

## 2. Speech Accent Archive

The corpus is composed of read speech samples of more than 1,700 speakers and their corresponding narrow IPA transcriptions. The speakers are from different countries around the world and they read a common elicitation paragraph, shown in Fig. 1. It contains 69 words and can be divided into 221 phonemes by referring to the CMU dictionary [18]. Each sample has its narrow IPA transcription, which was provided by trained phoneticians, and an example is also shown in Fig. 1. The transcriptions will be used to calculate reference inter-speaker phonetic distances. Use of read speech for clustering is considered to reduce the pronunciation diversity because read speech will show us only “controlled” diversity. In [19] however, English sentences read by 200 Japanese university students showed a very large diversity in terms of pronunciation and [20] showed that the intelligibility of the individual utterances to American listeners covered a very wide range. Considering these facts, we considered that clustering of read speech samples can still capture well how diverse World Englishes are in their pronunciation. In the current study, only the data with no word-level insertion or deletion were used. The

speakers' files that had exactly 69 words were automatically selected as candidate files and then, 515 files were obtained. However, the word order in some files were found to be wrong and we manually removed them. At the end of the day, 498 speakers' data were obtained and used in our study.

## 3. Phone-to-Phone Distance Estimation using Acoustic Analysis

In this study, the DTW algorithm is applied to compare two speakers' IPA transcriptions. Since the algorithm needs a distance matrix among all the existing IPA phones in the archive, we prepared the distance matrix firstly. In this paper, phone-to-phone distance was calculated through comparing acoustic characteristics of the two phones, which were produced by an expert phonetician. Before recording, we calculated frequency of each of the IPA phones, many of which were with a diacritical mark, and extracted the kinds of IPA phones that covered 95% of all the phones found in the archive. The resulting number of the kinds of the phones with/without a diacritical mark was 153. Table 1 shows the 153 phones. One expert phonetician, the third author, was asked to pronounce each of these phones twenty times. Here, he was asked to pay good attention to diacritical difference within the same kind of IPA phone. In the recording, the phonetician pronounced each vowel twenty times. For consonants, a consonant was succeeded and preceded at the same time by vowel [a]. For example, in order to collect data of phone [p], the phonetician spoke [apa] twenty times. In this way, each consonant was recorded.

Using the wav files and its phonetic transcription, a three-state HMM was built for each phone, where each state  $s_i$  ( $i \in 1, 2$ , and 3) contained a single Gaussian distribution with mean vector  $M_{s_i}$  and covariance matrix  $P_{s_i}$ . Here, MFCC(1-12) and its derivatives were used as acoustic features.

After training an HMM for each kind of the phones, the Bhattacharyya distance (BD) was calculated between two corresponding states of every phone pair. The equation of the BD between  $s_i$  of phone  $x$  and  $s_i$  of phone  $y$  is denoted below.

$$D_B(P_{s_i}^x, P_{s_i}^y) = \frac{1}{8} (M_{s_i}^x - M_{s_i}^y)^T P^{-1} (M_{s_i}^x - M_{s_i}^y) + \frac{1}{2} \ln \left( \frac{\det P}{\det P_{s_i}^x \det P_{s_i}^y} \right) \quad (1)$$

where  $M_{s_i}^x$  and  $M_{s_i}^y$  are mean vectors and  $P_{s_i}^x$  and  $P_{s_i}^y$  are covariance matrices of state  $s_i$  of  $x$  and state  $s_i$  of  $y$ , respectively. Note that  $P = (P_{s_i}^1 + P_{s_i}^2)/2$ .

For each phone pair, three Bhattacharyya distances were calculated, each corresponding to a state-to-state distance. By accumulating the distances and averaging them, we defined the acoustic distance between the phone pair. Equation 2 shows distance definition between two phones  $x$  and  $y$ .

$$d_{p_x, p_y} = \sqrt{\frac{D_B(P_{s_1}^x, P_{s_1}^y) + D_B(P_{s_2}^x, P_{s_2}^y) + D_B(P_{s_3}^x, P_{s_3}^y)}{3}} \quad (2)$$

We note here that, since the HMMs were trained in a speaker-dependent way, all the distances were calculated in the same and matched condition.

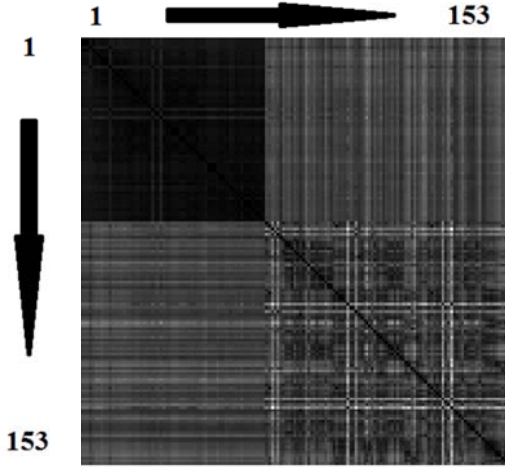


Fig. 2 The phone-to-phone distance matrix in gray-level image

The other 5% phones, which were not pronounced by the phonetician, were all with a diacritical mark and they were of very low frequency. For these phones, we substituted the HMMs of the same phones with no diacritical mark. With this substitution policy, the inter-phone distances among all the existing kinds of phones in the archive can be finally estimated. We converted the 153x153 phone-based distance matrix to its tree diagram. Although we do not show the diagram in this paper due to limit of space, the diagram confirmed us that we obtained a phonetically valid distance matrix. Fig. 2 shows a gray-level image of the 153x153 distance matrix instead.

In Fig. 2, X-axis and Y-axis denote the ID of phones (See table 1 in Appendix). Each pixel represents the distance between two phones. The first 66 phones are vowels and the others are consonants. The darker a pixel is, the more similar the phone pair are. From this figure, it can be seen that the distances between vowels are smaller than those between consonants or between a vowel and a consonant. We can also know that the distances between consonants have higher variance than those between vowels and those between a vowel and a consonant. Some small squares aligned in the diagonal line can be found in the figure because phones of the same kind with different diacritical marks are aligned together. Elements on this phone-to-phone distance matrix will be used as local distance or penalty to calculate the inter-speaker distance through the DTW alignment of two IPA transcriptions.

#### 4. Dynamic Time Warping using Phone-to-Phone Distance Information

In this section, the DTW is done to compare every two IPA transcriptions in a word-by-word manner by using the distance matrix obtained above. The obtained DTW alignment gives us an accumulated distortion score, which will be used as reference pronunciation distance between the two speakers. This speaker-to-speaker phonetic distance can be used in automatically clustering speakers in terms of pronunciation. Since all the transcriptions contain exactly 69 words, word-level alignment is easy and we only have to deal with phone-level insertions, deletions, and substitutions between a word and its counterpart in the two transcriptions. The local and allowable path of the DTW used in this section is shown as Fig. 3.

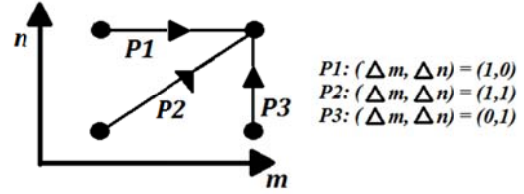


Fig. 3 Allowable paths of the DTW

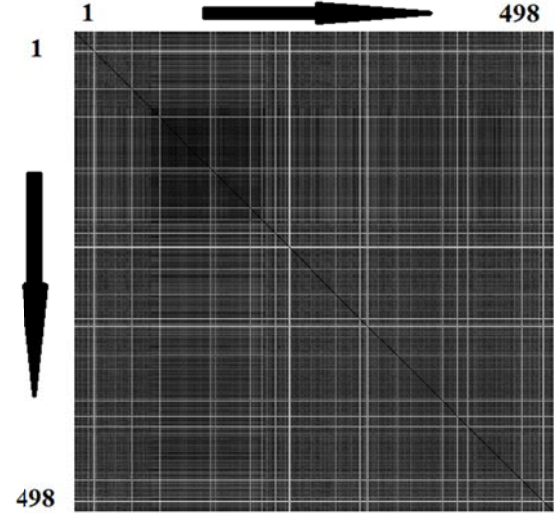


Fig. 4 The inter-speaker distance information matrix in gray-level image

P1, P2 and P3 are allowable paths of insertion, match and deletion. Path selection is done based on equation 3.

$$DTW[m, n] := \text{minimum} \left( \begin{aligned} &DTW[m-1, n] + \text{phone\_dist}[m, n], \\ &DTW[m-1, n-1] + 2 * \text{phone\_dist}[m, n], \\ &DTW[m, n-1] + \text{phone\_dist}[m, n] \end{aligned} \right) \quad (3)$$

$DTW[m, n]$  is the current accumulated cost at position  $(m, n)$  and  $\text{phone\_dist}[m, n]$  is a distance between the phone of time  $m$  and the phone of time  $n$ . Out of P1, P2, and P3, the path of which the accumulated cost at  $(m, n)$  is the minimum is selected. After normalizing this score by the total number of times of distortion accumulation, we can get a word-based distortion score. The 69 word-based scores are summed to be the final score for two given IPA transcriptions (speakers).

#### 5. Speaker-based Pronunciation Clustering

After obtaining the inter-speaker distances, all the speakers can be clustered using Ward's method, one of the hierarchical clustering methods. Since the clustering result of the 498 speakers is too complicated, we firstly show the gray-level image of the distance matrix of the 498 speakers in Fig. 4.

In the gray-level image, X-axis and Y-axis denote the ID of speakers. The IDs of speaker are assigned based on the alphabetical order of their countries' name. Each pixel represents the distance between two speakers. We can find a darker square in the top left. The distances in this region are from between native speakers of American English and this means that they have similar and stable English pronunciations. For non-native speakers, larger distances tend to be found to native

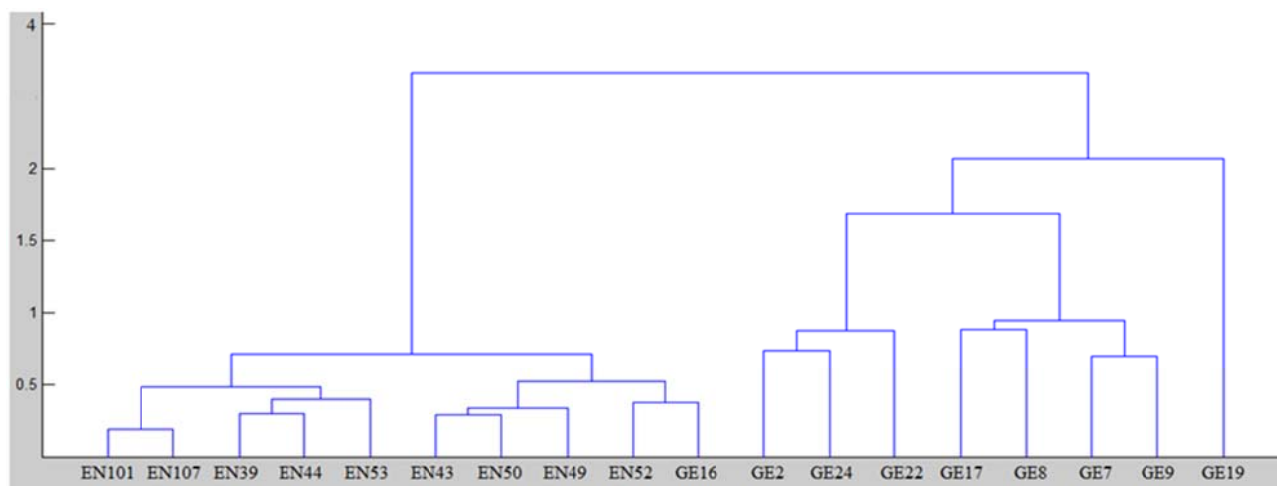


Fig.5 The clustering result of 18 selected speakers

speakers and to other non-native speakers. Non-native pronunciations can be affected by their mother tongue in different ways and to different degrees. In Fig.5, the clustering result of 18 selected speakers is shown. We picked up German speakers who were born in Germany in the archive, the number of whom was 9, and 9 native American English speakers were randomly selected. “EN” and “GE” denote American and German, respectively. The numbers succeeding “EN” or “GE” in the figure are speaker IDs. From Fig. 5, it can be seen that the all American speakers are clustered into one sub-tree and eight German speakers are clustered into the other sub-tree. Although GE16 is clustered into the same sub-tree with American speakers, by inspecting his biography included in the SAA, it is found that he had lived in USA for 4 years. It seems that his pronunciations has been reasonably affected by and adapted to American accent. On the other hand, most of the other German speakers live in America within or less than 1 year and this is supposed to be the reason why they are clustered into the other sub-tree. The 9 American and the 9 German samples will be included in the CD-ROM as media files. Interested readers should listen to those samples.

## 6. Use of the Reference Distances to Build a Distance Predictor for New Data

Using the inter-speaker distances calculated in the previous section as reference distances, a distance predictor for new speakers was trained and built in [1]. Here, only speech data of the new speakers were used and their IPA transcriptions were not. Invariant pronunciation analysis was adopted for pronunciation representation and Support Vector Regression was used for prediction. The correlation between manually prepared IPA-based distances and automatically predicted distances was 0.77. For comparison, an HMM-based phoneme recognizer was tested with word-based network grammar to convert new speakers’ utterances into phoneme sequences, not phone sequences. Here the network grammar was built in order to cover word-based pronunciation variations found in the SAA. Then, two generated phoneme sequences were aligned through the DTW by using the HMM-based phoneme-to-phoneme distance matrix. Since almost all the data were non-native and the recording environment varied from sample to sample, the phoneme recognition performance was so low as 46 % and the resulting correlation between the IPA-based

reference inter-speaker distances and the HMM-based distances was 0.043. The proposed predictor outperformed by far the HMM-based baseline system. Interested readers should refer to [1].

## 7. Conclusions

With the ultimate goal of drawing a global map of World Englishes on an individual basis, we’re developing a method of predicting the pronunciation distance between any pair of speakers [1]. For this project, the reference pronunciation distances are required and this paper describes how to prepare these distances in detail. Since the SAA archive provides a narrow IPA transcription for each accented utterance of the fixed elicitation paragraph, the DTW was applied to those IPA transcriptions with a phone-to-phone distance matrix obtained from recordings by an expert phonetician. Using the obtained distances, speaker clustering was done. Results showed that speaker clustering was effectively and validly performed only in terms of pronunciation. Although we’re focusing on only the segmental aspect of pronunciation, the obtained clustering result indicates that clustering only based on the segmental aspect can still capture how diverse World Englishes are in their pronunciation rather well. In future work, we are planning to collect a more data using social network infrastructure and incorporate the prosodic diversity into pronunciation distance calculation. Pedagogical application of the World and individual English map will also be considered in collaboration with language teachers.

## 8. Acknowledgements

The authors would like to thank National Science Council of Taiwan for their financial support. This work was supported in part by the National Science Council of Taiwan under the Grants NSC101-2917-I-006-011.

## 9. References

- [1] H.-P. Shen, N. Minematsu, S. H. Weinberger, T. Makino, J. Novak, T. Pongkittiphan, C.-H. Wu, "Speaker-based pronunciation clustering of World Englishes based on pronunciation structure analysis," *IEICE Technical Report*, SP2012-116, pp.7-12 (2013-2)



- [2] S. H. Weinberger, Speech Accent Archive, George Mason University, <http://accent.gmu.edu>
- [3] D. Crystal, *English as a global language*, Cambridge University Press, New York, 1995.
- [4] J. Jenkins, *World Englishes: a resource book for students*, Routledge, 2009.
- [5] B. Kachru, Y. Kachru, and C. Nelson, *The handbook of World Englishes*, Wiley-Blackwell, 2009.
- [6] A. Kirkpatrick, *The Routledge handbook of World Englishes*, Routledge, 2012.
- [7] M. Pinet, Paul Iverson, Mark Huckvale, "Second language experience and speech-in-noise recognition: the role of L2 experience in the talker-listener accent interaction", *Proc. of SLaTE*, CD-ROM, 2010.
- [8] N. Minematsu, "Mathematical evidence of the acoustic universal structure in speech," in *Proc. ICASSP*, pp.889-892, 2005.
- [9] N. Minematsu, Y. Qiao, S. Asakawa, M. Suzuki, "Speech structure and its application to robust speech processing", *Journal of New Generation Computing*, 28, 3, pp. 299-319, 2010.
- [10] R. Jakobson and L. R. Waugh, *Sound shape of language*, Branch Line, 1979.
- [11] M. Suzuki, Y. Qiao, N. Minematsu, and K. Hirose, "Pronunciation proficiency estimation based on multilayer regression analysis using speaker-independent structural features," in *Proc. SLaTE*, CD-ROM, 2010.
- [12] T. Zhao, A. Hoshino, M. Suzuki, N. Minematsu, K. Hirose, "Automatic Chinese pronunciation error detection using SVM with structural features," in *Proc. Spoken Language Technology*, pp.473-476, 2012.
- [13] X. Ma, R. Xu, N. Minematsu, Y. Qiao, K. Hirose, A. Li, "Dialect-based speaker classification using speaker invariant dialect features", in *Proc. of Int. Symposium on Chinese Spoken Language Processing*, pp.171-176, 2010.
- [14] N. Minematsu, K. Kamata, S. Asakawa, T. Makino, and K. Hirose, "Structural representation of the pronunciation and its use for clustering Japanese learners of English," in *Proc. SLaTE*, CD-ROM, 2007.
- [15] Y. Qiao, N. Minematsu, "A study on invariance of f-divergence and its application to speech recognition," in *IEEE Trans. on Signal Processing*, vol.58, no.7, pp.3884-3890, 2010.
- [16] M. Suzuki, G. Kurata, M. Nishimura, N. Minematsu, "Discriminative reranking for LVCSR leveraging invariant structure," in *Proc. INTERSPEECH*, CD-ROM, 2012.
- [17] D. Saito, S. Asakawa, N. Minematsu, and K. Hirose, "Structure to speech -- speech generation based on infant-like vocal imitation --," in *Proc. INTERSPEECH*, pp.1837-1840, 2008.
- [18] The CMU pronunciation dictionary, <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>
- [19] N. Minematsu, Y. Tomiyama, K. Yoshimoto, K. Shimizu, S. Nakagawa, M. Dantsuji, and S. Makino, "Development of English speech database read by Japanese to support CALL research," in *Proc. ICA*, pp.557-560, 2004.
- [20] N. Minematsu, K. Okabe, K. Ogaki, K. Hirose, "Measurement of objective intelligibility of Japanese accented English using ERJ (English Read by Japanese) database," in *Proc. INTERSPEECH*, pp.1481-1484, 2011.

## 10. Appendix

Table 1: 153 phones used in acoustic analysis

Vowels and Consonants used in Acoustic Analysis					
1. i	2. ĭ	3. i:	4. ĵ	5. ĭ̄	6. ĭ̄
7. y	8. ĵ	9. ɪ	10. ɪ:	11. ĵ̄	12. ĵ̄
13. e	14. ě	15. ě̄	16. ɛ	17. ě̄	18. ě̄
19. æ	20. æ	21. æ:	22. ǣ	23. a	24. ā
25. i	26. ĵ	27. ʔ	28. u	29. ŭ	30. ʊ
31. ɜ	32. ɜ̄	33. ɐ	34. ẽ	35. ũ	36. o
37. ẽ	38. ɐ	39. ẽ̄	40. ɐ̄	41. ẽ̄	42. ɐ̄
43. u	44. ũ	45. ũ̄	46. u	47. ũ̄	48. u:
49. ũ	50. ũ̄	51. ũ̄:	52. ʊ	53. ʊ̄	54. o
55. ʊ̄	56. ʊ̄	57. ʌ	58. ʌ̄	59. ɔ	60. ɔ:
61. ʌ̄	62. ʌ̄	63. ɑ	64. ɑ:	65. ă	66. ă̄
67. p	68. p <sup>h</sup>	69. p̄	70. b	71. b̄	72. b̄
73. f	74. β	75. β̄	76. β̄	77. f	78. v
79. ɣ	80. v	81. m	82. m̄	83. m̄	84. n
85. n̄	86. n̄	87. n̄	88. n̄	89. ŋ	90. n̄
91. t	92. t <sup>h</sup>	93. t̄	94. t̄	95. t̄	96. t̄
97. d	98. d̄	99. d̄	100. d̄	101. s	102. s̄
103. s̄	104. z	105. z̄	106. ɹ	107. ɹ̄	108. ɹ̄
109. r	110. ɹ	111. ɹ̄	112. l	113. l̄	114. l̄
115. θ	116. ð	117. ɸ	118. z	119. z̄	120. f
121. ʒ	122. ʒ̄	123. j	124. j̄	125. k	126. k <sup>h</sup>
127. k̄	128. k'	129. k <sup>h</sup>	130. k̄	131. g	132. g
133. ġ̄	134. ġ̄	135. x	136. ɣ̄	137. ɣ̄	138. ɰ
139. ʔ̄	140. h	141. fi	142. w	143. ɰ̄	144. pɰ̄
145. tθ	146. dð	147. ts	148. dz	149. tɛ	150. dz̄
151. tɰ̄	152. dʒ̄	153. kx			