# PRONUNCIATION CHARACTERISTICS OF JAPANESE SPEAKERS' ENGLISH: A PRELIMINARY CORPUS-BASED STUDY

Takehiko Makino, Chuo University, Tokyo

The development of the English Read by Japanese (ERJ) Phonetic Corpus consists of computer-readable narrow phonetic transcriptions and their corresponding target phonemes of selected 800 utterances from ERJ speech database. In describing the pronunciation characteristics of English spoken by Japanese speakers (or speaker of any language), we have been relying on the "rules of thumb" based on informal observations or theoretical predictions from the L1-L2 phonological differences, such as L/R confusion or conflation of English vowels into a five-vowel system in the case of Japanese speakers. While such rules of thumb have had roles to play, corpus-based studies of other areas of linguistic research have proved that they cannot give us the total picture of what are being studied, and L2 pronunciation should not be an exception. Indeed, a preliminary survey of the ERJ Phonetic Corpus has revealed some rather unexpected findings. The most notable of such findings is the spirantization (fricative realization) of voiceless plosives. Such a process is not part of standard Japanese phonology and cannot be the case of a negative L2 transfer. We can expect that the Corpus will help make a more systematic description of Japanese speakers' pronunciation of English.

## Background

The purpose of this paper is to make an interim report on the findings from a corpus-based descriptive study of the characteristics of Japanese speakers' pronunciation of English. I am currently developing English Read by Japanese (ERJ) Phonetic Corpus (Makino, 2007, 2009; Makino & Aoki 2012). The Corpus consists of computer-readable narrow phonetic transcriptions and their corresponding target phonemes and words of selected 800 utterances from ERJ speech database, a large collection (more than 70,000 speech files) of read-aloud sentences by Japanese university students.

The rationale for the corpus development was the lack of a systematic survey of the characteristics of Japanese speakers' pronunciation of English. Corpus studies on L2

Makino Pronunciation: A Corpus-Based Study

pronunciation have been very rare (Gut, 2009; Li, Zhang, Li, Harrison, Lo, & Meng, 2011; Meng, Tseng, Kondo, Harrison, & Viscelgia, 2009). The Corpus is intended to help fill this gap.

**Introduction to ERJ speech database**

ERJ stands for "English Read by Japanese" and the database was collected mainly in order to help CALL system development (Minematsu, et al. 2002). 807 different sentences and 1,009 different words (or word sets) were read aloud by 200 (100 male and 100 female) speakers in 20 different recording sites in Japan. All of the sites were universities and all the speakers were students there. Probably because it was deemed unpractical to ask all the 200 speakers to record all the sentences and words, they were divided into several sets, and individual speakers were asked to record only one sentence set and one word set. As a result, each sentence was read aloud by about 24 speakers (12 males and 12 females) and each word by about 40 speakers (20 males and 20 females). In total, the ERJ speech database consists of more than 70,000 speech files: 24,744 sentence files and 45,495 word files.

**Training sheets**

Before recording, the speakers were asked to practice pronouncing the words and sentences with training sheets which presented phonemic notations as well as orthographic words and sentences. The phonemic symbols used in the training sheets are based on ARPAbet used in TIMIT database (Garofolo, et al. 1993) and the CMU Pronouncing Dictionary (Carnegie Mellon University 2008), listed below:

Consonants: P, T, K, B, D, G, CH, JH, F, TH, S, SH, HH, V, DH, Z, ZH, M, N, NG, W, Y, L, R

Vowels: IY, IH, EY, EH, AE, AA, AO, OW, UH, UW, AH, AX, AW, AY, OY, ER, AXR

The model of the pronunciation is therefore mainstream American English. Each vowel was specified for degrees of stress: "1" for primary, "2" for secondary and "0" for unstressed.

In order to ensure that the speakers understood these symbols correctly, a website was prepared where they could listen to word examples of each phonemic symbol.

Examples from training sheets:

S1_001          This was easy for us.
          [DH IH1 S] [W AA1 Z] [IY1 Z IY0] [F AO1 R] [AH1 S]
S1_002          Is this seesaw safe?
          [IH1 Z] [DH IH1 S] [S IY1 S AO2] [S EY1 F]
S1_003          Those thieves stole thirty jewels.
          [DH OW1 Z] [TH IY1 V Z] [S T OW1 L] [TH ER1 T IY0] [JH UW1 AX0 L Z]

**The recordings**

In the recording sessions, the scripts only presented orthographic words and sentences, without phonemic notations. The reason for the change was that reading phonemic notations can induce unnatural pronunciation.

**Corpus building procedure**
**Limitations of the ERJ data**

It follows from the nature of the recording procedure above that there are some limitations in the ERJ speech database:

a) The speech is not spontaneous but read-aloud. It does not represent what is happening in natural settings.
b) The sentences are all isolated and out-of-context. This could have led to improper prosody (accenting, intonation or rhythm).
c) The words in the TIMIT phonologically-balanced sentences were chosen for their sounds. As a result, many of them were unfamiliar to the subjects. Even though phonemic notations were presented at the training stage, this could have led to mispronunciations or awkward pronunciations.

**Choice of the materials used**

Obviously, it is absolutely unpractical to use the whole database for the corpus building because of its sheer size. I have chosen to transcribe the 800 sentence files used in another study (Minematsu, Okabe, Ogaki, & Hirose, 2011). In that study, the recordings were played over the telephone to Americans who were not familiar with Japanese speakers' English, who then repeated what they (thought they) heard, and those repetitions were transcribed orthographically.

These sentences are all from the phonologically-balanced sets of sentences. The exclusion of word sets is justified because we are not interested in the pronunciation of individual words.

**Transcription procedure**

The transcription procedure for ERJ Phonetic Corpus is listed below:

1) To reduce the effort of manual transcription, the files were pre-processed by the Penn Phonetics Lab Forced Aligner (Yuan & Liberman, 2008; http://www.ling.upenn.edu/phonetics/p2fa/), which produced forced aligned transcriptions of English words and phonemes for each file in the Praat (Boersma & Weenik, 2013) TextGrid format.
2) Then, using Praat, the TextGrids were re-formatted into three tiers (target words, target phones, actual phones). The actual phones were manually transcribed, and boundaries of target phones and target words were manually aligned with those of the actual phones.
3) The corrected TextGrids were then imported into ELAN (Wittenburg, Brugman, Russel, Klassmann, & Sloetjes, 2006; http://tla.mpi.nl/tools/tla-tools/elan/), which has much better searching functionality than Praat, allowing searching in multiple files and output in concordance format.

The resulting ELAN .eaf files and the original .wav files are the complete individual data of the Corpus.

**The data set for this study**

The ERJ Phonetic Corpus consists of 419 different sentences, read by 200 different speakers (100 males and 100 females). Most of the sentences were read by two (one

male and one female) speakers. The total number of words is 5,959, with 1,599 different words. The total number of target phonemes is 24,873. The number of different target phonemes is, of course, 41 (the number of phonemes in General American English. cf. §2.2). The total number of actual phones is 25,460. The total number of different actual phones is 481.

Table 1

*The tokens of each phoneme in ERJ Phonetic Corpus*

| Vowels | | | | Consonants | | |
|---|---|---|---|---|---|---|
| ARPAbet | IPA | Count | | ARPAbet | IPA | Count |
| IY | i | 900 | | P | p | 630 |
| IH | ɪ | 943 | | B | b | 540 |
| EY | eɪ | 448 | | T | t | 1510 |
| EH | ɛ | 605 | | D | d | 959 |
| AE | æ | 640 | | K | k | 901 |
| AA | ɑ | 605 | | G | g | 339 |
| AO | ɔ | 453 | | F | f | 480 |
| OW | oʊ | 345 | | V | v | 387 |
| UH | ʊ | 131 | | TH | θ | 144 |
| UW | u | 512 | | DH | ð | 565 |
| AH | ʌ | 372 | | S | s | 1223 |
| AX | ə | 2283 | | Z | z | 868 |
| AY | aɪ | 416 | | HH | h | 266 |
| OY | ɔɪ | 85 | | SH | ʃ | 247 |
| AW | aʊ | 153 | | ZH | ʒ | 27 |
| ER | ɝ | 200 | | CH | tʃ | 209 |
| AXR | ɚ | 546 | | JH | dʒ | 252 |
| | | | | M | m | 718 |
| | | | | N | n | 1496 |
| | | | | NG | ŋ | 258 |
| | | | | L | l | 1203 |
| | | | | R | r | 1257 |
| | | | | W | w | 453 |
| | | | | Y | j | 294 |

**Predictions about vowels**

Japanese has a five-vowel system: /i, e, a, o, u/, whose typical realizations are [i, e, ɐ, o, ɯ]. Based on this, the following predictions can be made about the realizations of English vowel phonemes by Japanese speakers:

| | | |
|---|---|---|
| IH, IY | | [i(ː)] |
| UH, UW | [ɯ(ː)] | |
| AA, AE, AH, ER, AXR, AX | | [ɐ(ː)] |
| EH, EY | | [e(ː)] |
| AO, OW | [oː] | |
| AY, OY, AW | | combinations of /a, i, o/: [ɐi], [oi], [ɐɯ] |

**Findings about vowels - IY /i/ and IH /ɪ/**

Before looking at the findings, I have to note that the data below represent only the one-to-one relationship between target phonemes and actual phones, due to the limitation in the searching function of ELAN. As a result, they cannot show a complete picture, which will eventually have to be addressed.

The realizations of English IY /i/ and IH /ɪ/ are shown in the following table, arranged by frequency:

Table 2

*The frequency of different phones for target /i, ɪ/ in ERJ Phonetic Corpus*

| IY /i/ | Count (N=874) | | IH /ɪ/ | Count (N=905) |
|---|---|---|---|---|
| i | 613 | | i | 445 |
| ï | 98 | | ɪ | 287 |
| ɪ | 66 | | ï | 48 |
| ĩ | 19 | | ĩ | 30 |
| ɨ | 16 | | ɨ | 24 |
| e | 15 | | ɐ | 13 |
| others | 47 | | ə | 10 |
| | | | others | 48 |

More than half of the realizations of IH are [i]-like, which is also the dominant realization of IY. This is probably due to negative L1 transfer. On the other hand, we

found a substantial number of IH's realized as [ɪ]. So the conflation of these phonemes is not complete.

The [ɪ] in IY is interesting because the Japanese language does not have this phone as an allophone of any of its phonemes. It may be a case of hypercorrection.

**UW /u/ and UH /ʊ/**

The realizations of English UW /u/ and UH /ʊ/ are shown in the following table, arranged by frequency:

Table 3

*The frequency of different phones for target /u, ʊ/ in ERJ Phonetic Corpus*

| UW /u/ | Count (N=483) | | UH /ʊ/ | Count (N=90) |
|---|---|---|---|---|
| ʉ | 163 | | ʉ | 23 |
| ʉ̈ | 102 | | ̈ɯ̈ | 19 |
| ɨ | 86 | | ʊ | 18 |
| ü | 75 | | ü | 8 |
| ɤ | 9 | | ɨ | 7 |
| ʊ | 9 | | other | 15 |
| other | 39 | | | |

Although the tokens of UH may not be sufficient, we do not see any difference in the distributions of different realizations of UW and UH. This is probably a case of negative L1 transfer.

**AA /ɑ/, AE /æ/, AH /ʌ/, ER /ɝ/, AXR /ɚ/ and AX /ə/**

The realizations of English AA /ɑ/, AE /æ/, AH /ʌ/, ER /ɝ/, and AX /ə/ are shown in the following table, arranged by frequency:

Table 4

*The frequency of different phones for target /a, æ, ʌ, ɝ, ɚ/ in ERJ Phonetic Corpus*

| AA | Count | AE | Count | AH | Count | ER | Count | AXR | Count |
|---|---|---|---|---|---|---|---|---|---|

| /ɑ/ | (N=436) | /æ/ | (N=623) | /ʌ/ | (N=362) | /ɝ/ | (N=194) | /ɚ/ | (N=517) |
|---|---|---|---|---|---|---|---|---|---|
| ɐ | 94 | ɐ | 329 | ɐ | 169 | ɐ | 105 | ə | 172 |
| ö | 58 | a | 68 | ə | 42 | ə | 29 | ɐ | 158 |
| ɔ | 53 | ə | 63 | ɐ̃ | 30 | ɚ | 27 | ɚ | 84 |
| ɔ̈ | 44 | æ | 60 | ö | 25 | etc | 33 | ɨ | 14 |
| o | 39 | ɐ̃ | 37 | a | 19 | | | etc | 89 |
| ə | 32 | etc | 66 | ɵ | 11 | | | | |
| ɵ | 32 | | | etc | 66 | | | | |
| etc | 84 | | | | | | | | |

It is evident that most of the phonemes here are realized as a phone typically representing the Japanese phoneme /a/.

The o-like realizations for AA have occurred because of the influence from spelling. AA is typically represented with <o>, which is used in Romanized representation of the Japanese /o/, e.g. *tori* 'bird' /tori/ [toɾi].

The realizations of English AX (schwa) are shown in Table 5, arranged by frequency. Although the phones [ə, ɐ] rank the highest, the realizations of schwa is much more diverse than other a-like phonemes. This again may reflect its spelling, because schwa is represented with various spellings in English.

Table 5

*The frequency of different phones for target /ə/ in ERJ Phonetic Corpus*

| AX /ə/ | Count (N=2081) | | | | | | | |
|--------|----------------|---|----|----|----|---|------|----|
| ə | 495 | ĩ | 58 | ö | 18 | | others (ɔ̃, eɪ, ɪ̥, ··ʏ, ɑ̃, ɛ̃, ɜ, ʌ, y, ə̥, əɚ, əɪ, ɛɪ, ɨ̥, ʔɐ, ä, ɒ, ɐ̥, ɐ̥, ɑˤ, ɐðɨ, æ̃, ɐi, ɐɪ, ɐ̃ɪ, ẽ, əd̄, əɪ̥, əɯ, ən, ət, ɦɐ, ɦɐ, ɨ̥, ɪ̥, ĩɪ, ɪɥ, l, ɯ, ɯ̃, ɯ̈, ø, öʊ, ɸ, ɹɪ, s, ü, ũ, ʉ̃, ʊ, ÿ, z̃, ʔə) | 94 |
| | | ɵ | 46 | o | 16 | | | |
| ɐ | 347 | ẽ | 34 | ĩ | 14 | | | |
| | | ɛ | 33 | ɯ̈ | 14 | | | |
| i | 225 | a | 27 | ɨ̃ | 12 | | | |
| | | ə̃ | 23 | ɔ̈ | 11 | | | |
| ɪ | 192 | ʉ | 22 | ɛ̈ | 11 | | | |
| | | õ | 20 | æ | 10 | | | |
| ɨ | 136 | ɵ̃ | 20 | ɔ | 10 | | | |
| | | ɐ̃ | 19 | ɤ | 10 | | | |
| e | 108 | ë | 19 | ɚ | 9 | | | |
| | | ï | 19 | ɨ̥ | 9 | | | |

**Predictions about consonants**

Japanese has the following consonant phonemes: /p, t, k, b, d, g, s, z, h, m, n, r, j/. It also has special moraic phonemes which are usually realized as consonants: a moraic nasal /N/ and a moraic obstruent /Q/. Some of the realizational rules of the consonantal phonemes are:

a) Plosives /p, t, k, b, d, g/: voiced set tends to be realized as fricatives (spirantized) between vowels; /g/ can also be realized as [ŋ] in the same environment.

b) Fricatives /h, s, z/: /h/ is realized as [ç] before /i/ and [ɸ] before /u/; /z/ is often realized as an affricate [dz] word-initially.

c) Dental/alveolar sounds /t, d, s, z, n/ are heavily palatalized before /i/ and /j/ and realized as [tɕ, (d)ʑ, ɕ, (d)ʑ, ɲ] respectively, although non-palatalized pronunciation is possible for loanwords.

d) A liquid /r/: typically a tap [ɾ], but [l, ɖ] are also possible.

e) A moraic nasal /N/ appearing syllable-finally; a nasal stop homorganic to the following consonant; a nasalized vowel before a vowel; when absolute final a uvular nasal [ɴ].

f) No phonemic consonant cluster; no syllable final consonants, except for the moraic nasal /N/ and obstruent /Q/.

**Consonants: Some possible predictions**

Based on the above, the following predictions, among others, can be made about the realization of English consonant phonemes:

a) Lack of distinction between /l/ and /r/ before vowels

b) Lack of distinctions between:
   /s/ and /θ/
   /z/ and /ð/ and /dz/
   /dʒ/ and /ʒ/
   /b/ and /v/
   /h/ and /f/ before [u]-like phonemes

c) Addition of a weak vowel ([ɯ, ɨ, ə]) after a word-final consonant or in consonant clusters

**Unexpected findings: Spirantization of a voiceless plosive /p/**

The realizations of English /p/ are shown in the table below. The list is categorized according to the phonological contexts and arranged by frequency:

Table 6

*The frequency of different phones for target /p/ in different phonological contexts in ERJ Phonetic Corpus*

| realization | Phonological contexts (sil < silence) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
|  | V_C | C_V | V_V | C_C | sil_C | sil_V | V_sil | C_sil | Total |
| p | 25 | 94 | 74 | 39 | 12 | 15 |  |  | 259 |
| ɸ | 49 | 12 | 18 | 22 | 8 | 5 | 6 | 2 | 122 |
| pʰ | 19 | 12 | 21 | 10 | 7 | 4 | 10 | 1 | 84 |
| pɸ | 8 | 1 | 1 | 5 | 1 |  | 3 | 1 | 20 |
| pɨ | 7 |  | 1 | 5 | 2 |  |  |  | 15 |
| p̚ | 12 |  |  |  |  |  | 2 |  | 14 |
| other | 10 | 2 | 4 | 13 | 3 | 1 | 1 | 0 | 34 |
| **Total** | **130** | **121** | **119** | **94** | **33** | **25** | **22** | **4** | **548** |

Vowel insertion before another consonant is what theories predict, but it is hardly the most frequent case. We find numerous cases of spirantized phones. They cannot be predicted from Japanese phonology, where voiceless plosives are not spirantized.

Spirantization is most likely to occur syllable-finally, especially before consonants. But note that it occurs prevocalically as well, which is not the most likely position for consonantal weakening. This fact can have repercussions to phonological theories.

**Unexpected findings: Spirantization of a voiceless plosive /k/**

The realizations of English /k/ are shown in the table below. The list is categorized according to the phonological contexts and arranged by frequency:

Table 7

*The frequency of different phones for target /k/ in different phonological contexts in ERJ Phonetic Corpus*

| | Phonological contexts (sil < silence) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **realization** | **V_C** | **V_V** | **C_V** | **C_C** | **sil_V** | **V_sil** | **sil_C** | **C_sil** | **Total** |
| k | 80 | 100 | 87 | 23 | 30 | 1 | 8 | | 329 |
| kʰ | 76 | 25 | 35 | 26 | 17 | 20 | 4 | 1 | 204 |
| k̟ | 1 | 17 | 20 | 4 | 11 | | | | 53 |
| x | 36 | 4 | | 1 | | 2 | 1 | | 44 |
| kɨ | 9 | 4 | | 6 | | | 4 | | 23 |
| k̟ʰ | 1 | 11 | 5 | | 4 | | | | 21 |
| k̚ | 16 | | | | | | | | 16 |
| k' | 6 | 3 | | 1 | | 4 | | | 14 |
| kx | 6 | | | 2 | | 3 | | 1 | 12 |
| others | 23 | 7 | | 16 | | 2 | 5 | | 53 |
| **Total** | **254** | **171** | **147** | **79** | **62** | **32** | **22** | **2** | **769** |

Here the spirantization is less frequent than in the case of /p/, but the most likely contexts are the same. Spirantization is not at all frequent in /t/, which is not addressed here.

**L and R**

The realizations of English /l/ and /r/ are shown in the table below. The list is categorized according to the phonological contexts and arranged by frequency:

Table 8

*The frequency of different phones for target /l, r/ in different phonological contexts in ERJ Phonetic Corpus*

| /l/ | Phonological contexts (sil < silence) | | | | | | | /r/ | Phonological contexts (sil < silence) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | C_V | V_C | V_V | V_sil | sil_V | Total | | | C_V | V_V | V_C | sil_V | V_sil | Total |
| l | 134 | 102 | 119 | 20 | 9 | 384 | | ɹ | 303 | 61 | 3 | 22 | | 389 |
| ɫ | 32 | 109 | 32 | 30 | | 203 | | ɾ | 120 | 30 | | 12 | | 162 |
| ɾ | 79 | 2 | 50 | | 8 | 139 | | ə | 1 | 4 | 54 | | 10 | 69 |
| ɹ | 58 | | 41 | 2 | 1 | 102 | | l | 33 | 17 | | 7 | | 57 |
| lɨ | | 25 | 4 | 8 | | 37 | | ɚ | 1 | 7 | 26 | | 6 | 40 |
| rɨ | | 27 | 1 | 8 | | 36 | | ɐ | | 3 | 20 | | 13 | 36 |
| others | 43 | 47 | 23 | 20 | 2 | 135 | | ɯ | 28 | 1 | | | | 29 |
| | | | | | | | | ɾ | 20 | 5 | | 2 | | 27 |
| | | | | | | | | others | 32 | 5 | 20 | 7 | 5 | 69 |
| Total | 346 | 312 | 270 | 88 | 20 | 1036 | | Total | 538 | 133 | 123 | 50 | 34 | 878 |

Japanese speakers' English /r/ is notoriously mispronounced, but the data shows that it is not that bad. Nearly half the occurrence of this phoneme conforms to the native-speaker target. [ɾ] is what theories predict in both phonemes, but it is not actually in the majority.

/l/ is probably the more problematic phoneme for Japanese speakers of English. Note the case of [ɹ] for /l/. This could be the case of hypercorrection, probably because /r/ is one of the most emphasized sound in the teaching of the pronunciation of English in Japan, whereas /l/ is hardly ever emphasized.

**Implications for pronunciation teaching**

The study presented here is preliminary and incomplete, but I believe that I have shown that the reality is far more complicated than what phonological theory predicts.
In Japan, instruction of pronunciation of English at an introductory level is patchy at best. /r/ is emphasized but /l/ is not, /θ/ and /ð/ but not /s/ and /z/, vowels are not taught as an overall system but only some of the individual vowels deemed difficult such as /æ/

and /ɚ/ are treated, if at all. This may have resulted in the sort of disparate pronunciations presented above. We have to study what is happening more closely and help make a better pronunciation teaching syllabus.

**Prosodic notation**

So far only the segmental transcription has been completed for ERJ Phonetic Corpus. Since the Corpus is intended to be a source of all the phonetic characteristics of Japanese speakers' English speech, prosodic transcription is also necessary.

L2 prosody is very difficult to describe. Studies of non-native prosody such as Gut (2009) and Li, et al. (2011) use (modified) English ToBI, which I think is a wrong thing to do. L2 prosodic system is neither that of L1 nor of the target language, but something of the mixture of the two. I am now addressing this problem and devising a notational system of Japanese speakers' English prosody.

**Mispronunciation and misperception**

As noted above, the data set in this study was the same as that used in Minematsu, et al. (2011), where the recordings were played over the telephone to Americans who were not familiar with Japanese speakers' English and they were asked to repeat the sentences they heard. Their responses were orthographically transcribed.

With this data, we will be able to explore what sort of actual phone deviations are likely to lead to misunderstandings of what sort. This can be the basis for a study of intelligibility.

**ACKNOWLEDGEMENTS**

**ABOUT THE AUTHOR**

Takehiko Makino is an Associate Professor of English as a Foreign Language at Chuo University (Tokyo, Japan) and has taught English phonetics at a number of other universities including his alma mater Tokyo University of Foreign Studies, from which he received an MA in English linguistics in 1991. He also studied linguistics at the University of Kansas, USA in 1987-88. In 2004, he was awarded a Certificate of Proficiency in the Phonetics of English by the International Phonetic Association. From April 2012 till March 2014, he had an appointment as a Visiting Scholar at the University of Pennsylvania and the Linguistic Data Consortium under Chuo University Overseas Research Program.

Address:
  Chuo University
  742-1 Higashi Nakano
  Hachioji-shi, Tokyo 192-0393, Japan
Telephone: +81-(0)42-674-3401
Email: mackinaw@tamacc.chuo-u.ac.jp

**REFERENCES**

Boersma, P., & Weenink, D. (2013). Praat: doing phonetics by computer (version 5.3.55) [Computer software]. Available: http://www.praat.org/

Carnegie Mellon University (2008). CMU Pronouncing Dictionary (v. 0.7a). [Electronic database]. Available: http://www.speech.cs.cmu.edu/cgi-bin/cmudict

Garofolo, John, Lori Lamel, William Fisher, Jonathan Fiscus, David Pallett, Nancy Dahlgren, & Victor Zue. (1993). TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1. Web Download. Philadelphia: Linguistic Data Consortium. [Electronic database]. Available: https://catalog.ldc.upenn.edu/LDC93S1

Gut, U. (2009). *Non-native speech: A corpus-based analysis of phonological and phonetic properties of L2 English and German*. Frankfurt: Peter Lang

Li, M., Zhang, S., Li, K., Harrison, A., Lo, W.-K., & Meng, H. (2011). Design and collection of an L2 English corpus with a suprasegmental focus for Chinese learners of English. *Proceedings from the 17th International Congress of Phonetic Sciences*, Hong Kong, 1210-1213.

Makino, T. (2007). A corpus of Japanese speakers' pronunciation of American English: preliminary research. Paper presented at Phonetics Teaching and Learning Conference 2007. Abstract retrieved from http://www.ucl.ac.uk/psychlangsci/ptlc/proceedings_2007/ptlcpaper_02e

Makino, T. (2009). Vowel substitution patterns in Japanese speakers' English. In Biljana Čubrović & Tatjana Paunović (Eds.), *Ta(l)king English Phonetics Across Frontiers* (pp.19-31). Newcastle: Cambridge Scholars Publishing.

Makino, T., & Aoki, R. (2012). English Read by Japanese Phonetic Corpus: An interim report [Electronic version]. *Research in Language*, 10(1), 79–95.

Meng, H., Tseng, C., Kondo, M., Harrison A., & Viscelgia, T. (2009). Studying L2 suprasegmental features in Asian Englishes: a position paper. *Proceedings from Interspeech 2009*, Brighton, UK, 1683-1686.

Minematsu, N., Okabe, K., Ogaki, K. & Hirose, K. (2011). Measurement of objective intelligibility of Japanese accented English using ERJ (English Read by Japanese) database. *Proceedings from Interspeech 2011*, Florence, Italy, 1481-84.

Minematsu, N., Tomiyama, Y., Yoshimoto, K., Shimizu, K., Nakagawa, S., Dantsuji, M., & Makino, S. (2002). English speech database read by Japanese learners for CALL system development. *Proceedings from the 3rd International Conference of Language Resources and Evaluation*, Las Palmas, Spain, 896-903.

Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., & Sloetjes, H. (2006). ELAN: A professional framework for multimodality Research. *Proceedings from the 5th International Conference of Language Resources and Evaluation*, Genoa, Italy, 1556-59.

Yuan, J., & Liberman, M. (2008). Speaker identification on the SCOTUS corpus. *Proceedings from Acoustics '08*, Paris, France, 5687-5690.