

EFFECT OF SPEECH ENHANCEMENT AT PHONETIC LEVEL PERCEPTION OF ENGLISH SPEECH BETWEEN NATIVE ENGLISH AND MANDARIN LISTENERS

Yunqi C. Zhang¹, Catherine I. Watson², C.T Justine Hui¹, Yusuke Hioka¹

¹Acoustics Research Centre, Department of Mechanical and Mechatronics Engineering, University of Auckland, New Zealand

²Department of Electrical, Computer and Software Engineering, University of Auckland, New Zealand

yzhb694@aucklanduni.ac.nz; c.watson@auckland.ac.nz; justine.hui@auckland.ac.nz; yusuke.hioka@ieee.org

ABSTRACT

This study investigates the effect of commonly used speech enhancement algorithms in improving the intelligibility of noisy New Zealand English speech on native New Zealand English and native Mandarin listeners at the phonetic level. A phonetic error analysis was carried out to analyse the errors made by the listeners from a subjective listening test to find systematic errors. Results show that existing speech enhancement algorithms performed similarly by causing more onset deletion and coda deletion or insertion than when listening to clean speech. Overall, the algorithms did not improve the listeners' speech intelligibility on the phonetic level. Although native New Zealand English listeners always made fewer errors than native Mandarin listeners, they all made similar onset and coda errors. The nucleus errors made by the native Mandarin listeners were mainly caused by unfamiliarity with the accent.

Keywords: phonetic analysis, Mandarin, non-native, speech intelligibility, speech enhancement

1. INTRODUCTION

Speech intelligibility indicates how well speech can be comprehended by listeners. Non-native (L2) listeners have been known to experience more severe speech intelligibility degradation than native (L1) listeners under noise, such as listening or talking in a noisy environment [1, 2, 3]. To improve the speech intelligibility of listeners when listening to noisy speech, assistive listening devices implementing speech enhancement (SE) algorithms could be used [4]. While numerous numbers of SE algorithms have been developed and their performance has been thoroughly evaluated [5], to date their effect on L2 listeners is not well understood.

Listeners' speech perception was found to be influenced by their L1 language due to various factors such as the contrasts in phonetic familiarity, lexical structure, and acoustic structure between the L1 and L2 languages [3, 6]. As reported by the Ministry of Education of the People's Republic of China in 2019 [7], there were 703,500 Chinese students studying overseas and mainly entering English-speaking countries [8]. Hence, this study focuses on the perception of English speech by native Mandarin listeners.

Our previous study [9] investigated the effect of different SE algorithms on two groups of listeners: native English listeners residing in New Zealand as the control group and native Mandarin listeners residing in Mainland China who have never been immersed in an English-speaking country as the experimental group. It was found that all SE algorithms tested showed little improvement or even degradation in speech intelligibility compared to the original noisy speech for both groups, and the perception gap between the native English and native Mandarin listeners was not narrowed. So far, most studies on enhanced speech evaluated speech intelligibility by the correctness of each word, but none has investigated the response error pattern at the phonetic level. To ensure the corpus is in the L1 language of the control group, the previous study used a corpus in New Zealand English (NZE), despite most students in Mainland China learnt General American (GenAm) or British English at school [10]. Therefore, NZE is considered less familiar to the Mandarin participants and the effect of using NZE corpus in such listening tests remains unknown. NZE shares the same vowel system as other non-rhotic standard varieties of English such as Received Pronunciation, which consists of eleven monophthongs plus a neutral schwa [11], but differs from other non-rhotic varieties for the /ɪ, e, i:, æ/

vowels [12]. Despite NZE being less commonly used, 35% of the international students in New Zealand were imported from Mainland China [8], showing that NZE is used and can be understood by a part of the local Mandarin speakers in noisy environments.

The current study analyses the errors made by the two groups of listeners outlined in [9] for each phonetic component. The analysis explores the possible common or unique errors made by each group and gives a glimpse of the perception of native Mandarin listeners on NZE.

2. METHODOLOGY

Phonetic level analysis was applied to the results of two subjective listening tests conducted in [9], which were 1) the *baseline test* that evaluated the clean speech from a NZE corpus by native Mandarin listeners and 2) the *speech enhancement test* that investigated the intelligibility of noisy speech processed by selected SE algorithms on native English and native Mandarin listeners. The response keywords were divided into phonetic components for marking and the systematic errors made were analysed. The study was approved by the University of Auckland Human Participants Ethics Committee (UAHPEC24202).

2.1. Subjective listening tests

The subjective listening tests were conducted online where the participants were asked to use their own headphones and transcribe the stimuli. The speech intelligibility was evaluated by counting the number of words correctly answered, which was normalised by the total number of words. Details of the test design can be found in [9].

2.1.1. Baseline Test

Nine native Mandarin listeners based in Mainland China (C_{clean} group) were recruited. All participants reported they had learnt or had been learning American ($n = 7$) or British ($n = 2$) English. Each participant transcribed 48 clean (without any noise added) English sentences.

2.1.2. Speech Enhancement (SE) Test

Twenty native NZE listeners residing in New Zealand (NZE_{SE} group) and 19 native Mandarin listeners residing in Mainland China (C_{SE} group, referred as CC group in [9]) were recruited. All participants in the C_{SE} group reported have learnt

either British ($n = 6$) and/or American English ($n = 13$). Each participant transcribed 108 noisy speech sentences, which included six types of speech: the original noisy speech (before SE was applied) and the speech enhanced by five widely used SE algorithms, namely Wiener filter (WF) [13], subspace (SS) [14], non-negative matrix factorisation (NMF) [15], Conv-TasNet (Conv) [16], and complex U-Net (Unet) [17].

2.1.3. Stimuli

To generate the stimuli used in the test, the Bamford-Kowal-Bench (BKB) [18] sentences from the Speech Perception Assessment New Zealand (SPANZ) corpus [19] were used. The corpus is designed for participants speaking NZE, i.e. the sentences were recorded in NZE accent and revised to accommodate commonly used expressions in New Zealand. According to the description in the original BKB sentences, it was assumed that the occurrence of the vowels and phonemes in the keywords was balanced. Each sentence contains 3 - 4 keywords to be marked; any words other than the keywords were not marked.

2.2. Phonetic Analysis

The correctness of participants' responses from the subjective listening tests was marked manually by the first author according to the suggestions in the BKB corpus [18]. All keywords that received incorrect answers were collected for phonetic analysis. This is because even though the response words differed for each participant, the phoneme they tend to make mistakes on may be the same.

A phoneme that occurs at different locations can have various pronunciations according to the surrounding phoneme environments, which is known as co-articulation [20]. Therefore, recording the location of the mistaken phonemes can be crucial for analysing the errors specifically. This was realised by dividing every syllable into its phonetic components: onset, nucleus, and coda. For multi-syllable words, the phonetic components were numbered to the syllable they were in.

Overall, each word was divided into their syllables and further divided by phonetic components. Absence of a component was recorded as *NA*. An example is the two syllable word, "happy" [hæ-pi:], being separated into onset 1 /h/, nucleus 1 /æ/, coda 1 *NA*, onset 2 /p/, nucleus 2 /i:/, and coda 2 *NA*. Responses with inconsistent number of syllables compared to that of the keywords were recorded with the difference

in the number of syllables, where a positive number indicates more syllables were answered than that in the keyword, and vice versa. The phonetic analysis was applied to the baseline and SE test results separately. The analysis for the SE test also recorded different SE algorithms and groups of participants separately.

This study only investigated the phonemes in the first two syllables in each word as all of the baseline test keywords (83.3% mono-, 16.7% double-syllable) and 98.4% of the SE test keywords were either mono- (69.4%) or double-syllable (28.8%). The responses of three-syllable keywords were excluded as the errors made by the participants were too few to provide any systematic error, i.e. all combinations of keyword and error occurred less than twice. Errors of the same phonetic component in different syllables were not combined as the second syllable had a different phoneme error pattern compared to the first syllable.

3. RESULTS AND DISCUSSION

According to the phonetic analysis, common phonetic errors can be classified into deletion, missing the whole phonetic component; insertion, making up non-existing phonetic component; and substitution, substituting certain phonemes or the whole phonetic component with other phonemes. The errors were analysed by each phonetic component separately and listed in tables. Since the main purpose of this study is to investigate the performance of SE algorithms on native Mandarin listeners, the tables in this paper will focus on the error occurrence of the C_{SE} group. Table 1 shows the three most common nucleus errors for the C_{SE} group for each SE algorithm from the first syllable, where the corresponding error occurrence of the NZE_{SE} group is also reported. Tables 2, 3, and 4 show the most common onset, nucleus, and coda errors of all three groups, respectively, where for the C_{SE} group only the errors occurred over 15 times are presented due to space constraint. The onset errors of the C_{clean} group are excluded due to the absence of systematic onset errors, i.e. all combinations of keyword and error appeared only once.

Errors of the phonetic components were first analysed for each SE algorithm, but little difference was found between algorithms. For example, as shown in Table 1, most of the common nucleus errors reappear in every algorithm. The unique [ɜ:/-ou/] confusion in noisy speech also occurred in most SE algorithms for the C_{SE} group (twice for SS and Unet, three times for WF and NMF). The onset

Table 1: The three most common occurred nucleus errors for C_{SE} group in the first syllable for each SE algorithm.

SE Algorithm	Nucleus error (C_{SE} count, NZE_{SE} count)		
Noisy	eɪ → aɪ (7, 3)	e → i: (4, 1)	ɜ: → ou (4, 0)
WF	eɪ → aɪ (7, 7)	ə → NA (5, 5)	ɒ → ou (5, 1)
SS	eɪ → aɪ (9, 9)	e → eɪ (8, 0)	ɒ → ʌ (5, 0)
NMF	eɪ → aɪ (9, 10)	ə → NA (8, 5)	ɒ → o: (7, 1)
Conv	eɪ → aɪ (11, 9)	ə → NA (6, 4)	e → eɪ (6, 0)
Unet	ɒ → ou (7, 0)	ə → NA (6, 2)	ɒ → o: (5, 0)

Table 2: Most frequent onset errors (C_{clean} group is excluded as no systematic onset errors).

C_{SE} group	Count	NZE_{SE} group	Count
Onset 1 (first syllable)			
ð → NA	38	ð → NA	39
ð → h	17	p → NA	15
k → h	16	f → NA	13
p → NA	15	h → NA	13
f → NA	15	p → p l	11
Onset 2 (second syllable)			
d → NA	17	l → NA	6
s → NA	17	r → NA	6
l → NA	10	b → NA	5
d → t	9	d → NA	5
r → NA	9	s → NA	5

Table 3: Most frequent nucleus errors.

C_{clean} group	Count	C_{SE} group	Count	NZE_{SE} group	Count
Nucleus 1 (first syllable)					
e → ɪ	9	eɪ → aɪ	46	eɪ → aɪ	41
ʌ → æ	7	ə → NA	31	ə → NA	25
e → i:	6	e → eɪ	26	æ → a:	14
æ → e	3	ɒ → ou	20	æ → e	11
ɪə → i:	3	ɒ → o:	19	eɪ → ə	11
ou → o:	2	e → ɪ	17	au → e	11
æ → ʌ	2	ɒ → ʌ	17	ou → ʌ	7
eɪ → aɪ	2	æ → e	16	eɪ → æ	7
ʌ → a:	2	e → i:	16	u: → i:	7
ɪ → i:	2	ɜ: → u:	16	i: → ə	7
ɪ → ʌ	2	eɪ → i:	15	æ → eɪ	6
Nucleus 2 (second syllable)					
i: → ə	2	ə → NA	65	ə → NA	32
		ɪ → NA	21	ɪ → NA	8
		i: → ɪ	21	i: → ɪ	7

and coda errors caused by each SE algorithm were mostly deletion errors, which are not shown in this paper. The only exception was the /j/-/h/ confusion in speech enhanced by NMF for the C_{SE} group ($n = 5$) in onset. As a result, analysis of aggregated results from all SE algorithms was conducted.

For common onset errors, as reported earlier, no systematic errors ($n = 1$) were found from the C_{clean} group. However, for the C_{SE} group, the occurrence of the most frequent systematic onset errors increased ($n = 39$) similar to that of the nucleus ($n = 46$) and coda errors ($n = 32$). This increase is expected as the consonants at the beginning of an English word are often voiceless [21]. Therefore, the onset, especially the onset of the first syllable, can be masked easily by noise and/or removed by the SE algorithm. The most frequent

Table 4: Most frequent coda errors.

<i>C_{clean}</i> group	Count	<i>C_{SE}</i> group	Count	<i>NZE_{SE}</i> group	Count
Coda 1 (first syllable)					
l → t	3	NA → n	32	NA → n	12
NA → d	3	NA → d	20	n t → t	11
t → tʃ	3	n → NA	15	k → t	8
d → t	2	NA → l	15	NA → d	8
Coda 2 (second syllable)					
		ŋ → NA	27	ŋ → NA	10
		NA → ŋ	18	n → NA	8
		n → NA	16	NA → ŋ	8

type of onset error in the *C_{SE}* group was deletion, which was the same for the *NZE_{SE}* group, indicating that the deletion of phonemes is not due to their language nativity but a consequence of the noise and SE algorithms.

Many of the frequent nucleus errors made by the Mandarin listeners (*C_{clean}* and *C_{SE}* groups) can be explained by the unfamiliarity with the NZE accent, which are shown in bold in Table 3, including the well-known centralised /ɪ/, the raised /e/ [11], and the more recently studied merging [ɜ:/-/u:/] [22]. Some vowel confusions may be caused by the listeners being more familiar with the rhotic GenAm accent, which is learnt by most of the Mandarin participants, than NZE accent. For example, the [ʌ/-/æ/] confusion can be explained by the fact that /æ/ in GenAm is usually used as /ʌ/ in NZE; the [ɒ/-/o:/] confusion may be due to the fact that NZE distinguishes /ɒ/ and /o:/ while /ɒ/ and /a/ merged in GenAm [23, 21]. The most common nucleus errors in the second syllable were identical between the *C_{SE}* and *NZE_{SE}* groups, indicating that under noisy conditions, even the native listeners have difficulty distinguishing the vowel lengths between /ɪ/ and /i:/. Moreover, the schwa in both syllables was ignored frequently, 96 times by the *C_{SE}* group, 57 times by the *NZE_{SE}* group in total, as it is an unstressed vowel with a very short duration and usually appears in the unstressed syllable in the corpus, e.g. “along” to “long”. Hence, it can be easily missed or masked by noise. All three groups made two common vowel errors, the [eɪ/-/aɪ/] confusion (e.g. “hay” to “hi”) and the [æ/-/e/] confusion (e.g. “bag” to “bed”), where the latter is regarded as a characteristic of NZE accent. Both confusions imply that under noisy conditions, even native listeners encounter difficulties in distinguishing certain vowel pairs.

As seen in Table 4, the common coda errors from the *C_{clean}* group were mainly substitution errors, which can be explained by the language structure of Mandarin in which no obstruent consonant in the word-final position is allowed [24]. Hence, the *C_{clean}* group may be less sensitive to the word-final sound even when listening to clean speech. For the SE test, most coda errors by both groups can be

concluded as deletion or insertion errors. This can be explained by the same reason for the missing onset, as the consonants in English at the end of a word are also mostly voiceless [21]. The *NZE_{SE}* group showed much fewer coda errors than the *C_{SE}* group, which may be because the *NZE_{SE}* group was able to notice small but significant vowel duration differences to distinguish voicing minimal pairs compared to the *C_{SE}* group [25]. The coda errors in the second syllable for both *C_{SE}* and *NZE_{SE}* groups were the same, related to the deletion or insertion of /n/ and /ŋ/. One possible explanation of this finding is that the SE algorithms tend to manipulate the noisy signal in a way that the nasal consonants sound similar to noise. The high occurrence of /ŋ/ may be due to the fact that many double-syllable keywords in the corpus are gerund, e.g. “playing”.

4. CONCLUSION

This study analysed the phonetic errors from three groups of results from listening tests that evaluated the intelligibility of speech: native Mandarin listeners’ responses to clean NZE speech (*C_{clean}* group), native Mandarin (*C_{SE}* group) and English listeners’ (*NZE_{SE}* group) responses to enhanced NZE speech. The errors between different algorithms showed little difference, indicating that the existing SE algorithms did not improve the speech intelligibility of either native English or native Mandarin listeners at phonetic level. The result showed that the *NZE_{SE}* group had much fewer errors than the *C_{SE}* group for all phonetic components. Compared to the *C_{clean}* group, the *C_{SE}* group tended to have much more deletion errors for onset, and deletion and insertion errors for coda. Such error patterns in the onset and coda were also similar for the *NZE_{SE}* group.

The main limitation of this study is that the corpus is not designed for phonetic analysis. Since meaningful sentences were used in the test, participants with higher language proficiency would have been able to use semantic cues to guess the correct keywords. This top-down effect should have been considerably reduced as most of the keywords in a sentence had multiple possible substitutions. The meaningful sentences also lead to the test containing unbalance phones, resulting in a less comprehensive analysis for each phoneme. Also, part of the manual marking involved guessing the response, hence, the result may not fully reflect the participants’ perception. Further experiments with a corpus that is designed for phonetic analysis would be recommended.

ACKNOWLEDGEMENT

We thank the participants for their participation, Dr. Suzanne Purdy for providing the SPANZ corpus, and Dr. Elaine Ballard for her insightful advice.

5. REFERENCES

- [1] M. Cooke, M. L. G. Lecumberri, and J. Barker, "The foreign language cocktail party problem: Energetic and informational masking effects in non-native speech perception." *The Journal of the Acoustical Society of America* 123, no. 1, pp. 414–427, 2008.
- [2] S. L. L. M. C. Mattys, C. K. W. Li, and S. L. Y. Chan, "Effects of energetic and informational masking on speech segmentation by native and non-native speakers." *Speech Communication, Non-native Speech Perception in Adverse Conditions*, 52, no. 11, pp. 887–99, 2010.
- [3] M. L. G. Lecumberri, M. Cooke, and A. Cutle, "Non-native speech perception in adverse conditions: A review." *Speech Communication*, 52, no. 11, pp. 864–86, 2010.
- [4] P. Loizou, *Speech Enhancement: Theory and Practice*, 2007.
- [5] D. S., R. Deshmukh, and P. P., "A review of speech signal enhancement techniques." *International Journal of Computer Applications* 139, pp. 23–26, 2016.
- [6] L. Polka, "Characterizing the influence of native language experience on adult speech perception." *Perception & Psychophysics* 52, no. 1, pp. 37–52, 1992.
- [7] Statistics on chinese learners studying overseas in 2019 - ministry of education of the people's republic of china [www document]. [Online]. Available: http://en.moe.gov.cn/news/press_releases/202012/t20201224_507474.html
- [8] OECD, "What is the profile of internationally mobile students?" 2021.
- [9] Y. C. Zhang, C. J. Hui, Y. Hioka, and C. I. Watson, "Performance of speech enhancement algorithms for native Mandarin listeners on English perception," in *Proc. 15th ICA International congress on acoustics*, South Korea, 2022, pp. 111–116.
- [10] W. Wenfeng and G. Xuesong, "English language education in china: A review of selected research." *Journal of Multilingual and Multicultural Development* 29, no. 5, pp. 380–99, 2008.
- [11] C. I. Watson, M. Maclagan, and J. Harrington, "Acoustic evidence for vowel change in new zealand english." *Language Variation and Change* 12, no. 1, pp. 51–68, 2000.
- [12] M. Maclagan and J. Hay, "Getting fed up with our feet: Contrast maintenance and the new zealand english 'short' front vowel shift." *Language Variation and Change* 19, no. 1, pp. 1–25, 2007.
- [13] P. Scalart and J. Filho, "Speech enhancement based on a priori signal to noise estimation," in *1996 IEEE International Conference on Acoustics, Speech, and Signal Processing Conference Proceedings*, vol. 2, 1996, pp. 629–632 vol. 2.
- [14] Y. Hu and P. Loizou, "A generalized subspace approach for enhancing speech corrupted by colored noise," *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 4, pp. 334–341, 2003.
- [15] N. Mohammadiha, P. Smaragdis, and A. Leijon, "Supervised and unsupervised speech enhancement using nonnegative matrix factorization," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 10, pp. 2140–2151, 2013.
- [16] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time-frequency magnitude masking for speech separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [17] H.-S. Choi, J.-H. Kim, J. Huh, A. Kim, J.-W. Ha, and K. Lee, "Phase-aware speech enhancement with deep complex U-Net," *ICLR*, p. 20, 2019.
- [18] J. Bench, A. Kowal, and J. Bamford, "The BKB (bamford-kowal-bench) sentence lists for partially-hearing children," *British Journal of Audiology*, vol. 13, no. 3, pp. 108–112, 1979.
- [19] J.-H. Kim and S. Purdy, "Speech perception assessments new zealand (SPANZ)," *New Zealand Audiological Society Bulletin*, vol. 24, pp. 9–16, 2014.
- [20] R. Daniloff and R. Hammarberg, "On defining coarticulation." *Journal of Phonetics* 1, pp. 239–248, 1973.
- [21] H. Rogers, *The Sounds of Language: An Introduction to Phonetics*. Taylor & Francis Group, London, United Kingdom, 2000.
- [22] M. Maclagan, C. I. Watson, R. Harlow, J. King, and P. Keegan, "Investigating the sound change in the new zealand english nurse vowel /ɜ:/," *Australian Journal of Linguistics* 37, no. 4, pp. 465–85, 2017.
- [23] M. Hay, J. and Maclagan, E. Gordon, J. Beal, P. Honeybone, and A. McMahon, *New Zealand English*. Edinburgh University Press, Edinburgh, United Kingdom, 2008.
- [24] C.-C. Cheng and C.-C. Cheng, *A Synchronic Phonology of Mandarin Chinese*. De Gruyter, Inc., 1973.
- [25] J. E. Flege, M. J. Munro, and L. Skelton, "Production of the word-final english /t-/d/ contrast by native speakers of english, mandarin, and spanish." *The Journal of the Acoustical Society of America* 92, no. 1, pp. 128–43, 1992.