

LEVERAGING PHONE-LEVEL LINGUISTIC-ACOUSTIC SIMILARITY FOR UTTERANCE-LEVEL PRONUNCIATION SCORING

Wei Liu^{1,*}, Kaiqi Fu², Xiaohai Tian², Shuju Shi², Wei Li², Zejun Ma² and Tan Lee¹

¹Department of Electronic Engineering, The Chinese University of Hong Kong

²ByteDance

ABSTRACT

Recent studies on pronunciation scoring have explored the effect of introducing phone embeddings as reference pronunciation, but mostly in an implicit manner, i.e., addition or concatenation of reference phone embedding and actual pronunciation of the target phone as the phone-level pronunciation quality representation. In this paper, we propose to use linguistic-acoustic similarity to explicitly measure the deviation of non-native production from its native reference for pronunciation assessment. Specifically, the deviation is first estimated by the cosine similarity between reference phone embedding and corresponding acoustic embedding. Next, a phone-level Goodness of pronunciation (GOP) pre-training stage is introduced to guide this similarity-based learning for better initialization of the aforementioned two embeddings. Finally, a transformer-based hierarchical pronunciation scorer is used to map a sequence of phone embeddings, acoustic embeddings along with their similarity measures to predict the final utterance-level score. Experimental results on the non-native databases suggest that the proposed system significantly outperforms the baselines, where the acoustic and phone embeddings are simply added or concatenated. A further examination shows that the phone embeddings learned in the proposed approach are able to capture linguistic-acoustic attributes of native pronunciation as references.

Index Terms— Linguistic-Acoustic Similarity, Phone Embedding, Goodness of Pronunciation, Pronunciation Scoring.

1. INTRODUCTION

Pronunciation scoring is an essential component of Computer Assisted Pronunciation Training (CAPT) [1–5]. It is designed to automatically assess second language (L2) learners’ speech pronunciations [6, 7]. In general, the degree of proficiency/pronunciation level is measured as the amount of deviation of the L2 production from the reference native production. A typical assessment scenario is as follows: given a text prompt, the L2 learner is asked to read the text, and a scoring system is used to give a score based on the learner’s speech production.

Goodness of Pronunciation (GOP) [7] was a commonly used feature in automatic pronunciation assessment, mispronunciation detection, and related tasks [8–11]. In a deep neural network (DNN) based system, GOP is computed as the ratio of log phone posterior probability between the canonical reference phone and the hypothesized phone with the highest posterior probability [12]. It gives a general-sense measurement of the pronunciation quality, i.e., a lower value of GOP indicates a higher possibility of mispronunciation. To

improve the mispronunciation detection performance of GOP, various methods have been proposed. Transition probability between Hidden Markov Model (HMM) states was considered in [13] and a context-aware GOP score was investigated in [14].

A comparison-based framework was investigated in [15–17], where an utterance spoken by native speakers was adopted as a reference, and divergence-related features computed by dynamic time warping (DTW) between speech representations of native speakers and L2 learners were used to quantify the pronunciation deviation. However, parallel reference speech may not be available in real-world applications. Thus, it was considered to use phone embedding as the reference for pronunciation assessment [18–24]. One-hot representations of phoneme labels are fed into a trainable embedding layer to generate phone embedding vectors. The phone embeddings were used along with the corresponding phone-level acoustic embeddings for pronunciation score prediction. Addition [18–22] and concatenation [23, 24] of reference phone embedding and phone-level acoustic embedding are widely used methods to calculate phone-level pronunciation quality representation. The resultant representation is assumed to capture the deviation of non-native pronunciation from the reference production. However, either addition (add_phone) or concatenation (concat_phone) operation does not explicitly measure the degree of mismatch between what one native speaker pronounces (i.e., phone embedding) and how they actually pronounce (i.e., phone-level acoustic embedding). We hypothesize that explicit measurement of the degree of phone-level pronunciation deviation would better reflect L2 learners’ pronunciation quality.

In [25], a linguistic-acoustic similarity based accent shift (LASAS) model was proposed for accent recognition. The accent shift is intended to capture the pronunciation variants of the same word in different accents. It is explicitly modeled by the similarity of acoustic embedding and aligned text anchor vectors. In the present study, we propose to use cosine similarity between a reference phone embedding and the corresponding acoustic embedding to explicitly measure the mismatch between standard and non-native pronunciation. A phone-level GOP pre-training process is developed to guide similarity-based learning for better initialization of the two embeddings. Lastly, a bottom-up hierarchical pronunciation scorer [19] is used to map a sequence of phone embeddings, acoustic embeddings along with the proposed similarity measures to predict the final utterance-level score. Experimental results show that the proposed system significantly improves the score prediction performance in terms of Pearson correlation coefficients (PCC) compared to its counterpart, where phone and acoustic embeddings are simply added or concatenated. In addition, it is shown that the learned phone embedding can capture linguistic-acoustic characteristics of native pronunciation as references.

* This work was done during an internship at ByteDance.

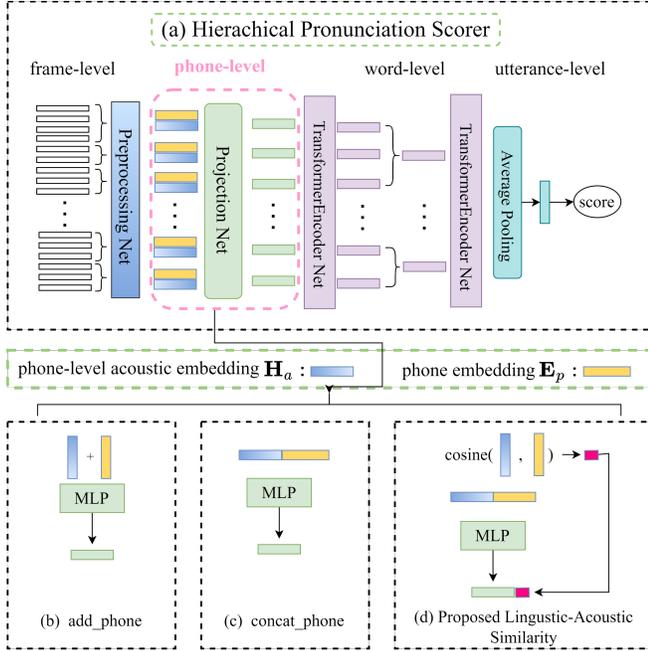


Fig. 1: (a) illustrates the hierarchical architecture of the pronunciation scoring network. (b) and (c) show the two conventional methods of combining phone-level acoustic embeddings and phone embeddings. (d) shows the proposed linguistic-acoustic similarity method.

2. METHOD

In this section, we first give an overview of the hierarchical pronunciation scoring network, followed by an introduction to the proposed methods, including linguistic-acoustic similarity measure and phone-level GOP pre-training.

2.1. Hierarchical Pronunciation Scorer

As shown in Figure 1 (a), the hierarchical pronunciation scoring network takes frame-level features as input, and aggregate and transform them into phone-level, word-level and utterance-level features layer by layer. The final output is a pronunciation score [19].

Previous research [18–24] was focused on different implementations of representation learning for phone-level pronunciation quality aiming to model the deviation of L2 pronunciation from reference native pronunciation. Figures 1 (b) and (c) depicts two such attempts. The vector in blue represents acoustic embeddings at phone-level after preprocessing network and the vector in yellow, i.e., the reference phone embedding, represents the native pronunciation of the current phone.

Given a pair of read speech utterance and text prompt, an acoustic model is used to extract frame-level features (e.g., deep feature as in [19]) and phone-level alignment. Phone-level acoustic features are obtained by averaging the aligned feature frames of each phone segment, denoted as $\mathbf{X} \in \mathbf{R}^{D_1 \times N}$. N is the number of phones, and D_1 is the feature dimension. Then, the phone-level acoustic feature and the corresponding phone ID are used as the input of the preprocessing network as shown in Eq. (1). A phone embedding layer encodes the phone ID \mathbf{e} into phone embedding vectors, $\mathbf{E}_p \in \mathbf{R}^{D_2 \times N}$. A fully connected (FC) layer projects the phone-level acoustic features \mathbf{X} into the same dimension as phone embedding, denoted as

$\mathbf{H}_a \in \mathbf{R}^{D_2 \times N}$. D_2 denotes the embedding dimension. \mathbf{H}_a is termed as the phone-level acoustic embedding in this paper.

$$\mathbf{H}_a, \mathbf{E}_p = \mathcal{F}(\mathbf{X}, \mathbf{e}), \quad (1)$$

where $\mathcal{F}(\cdot)$ denotes the preprocessing network.

Subsequently, the phone-level projection network takes \mathbf{H}_a and \mathbf{E}_p as input and generates a pronunciation quality representation, denoted as \mathbf{P}_Q :

$$\mathbf{P}_Q = \mathcal{P}(\mathbf{H}_a, \mathbf{E}_p), \quad (2)$$

where $\mathcal{P}(\cdot)$ represents the phone-level projection network. Specifically, it operates in two different ways:

$$\mathbf{P}_Q = \begin{cases} \text{MLP}(\mathbf{H}_a + \mathbf{E}_p), & \text{add_phone,} \\ \text{MLP}([\mathbf{H}_a; \mathbf{E}_p]), & \text{concat_phone.} \end{cases} \quad (3)$$

where MLP refers to the multilayer perceptron function and $[\cdot; \cdot]$ denotes the concatenation of two vectors.

Finally, the utterance-level pronunciation score \hat{y} is predicted as,

$$\hat{y} = \mathcal{U}(\mathcal{T}_w(\mathcal{A}_w(\mathcal{T}_p(\mathbf{P}_Q)))), \quad (4)$$

where $\mathcal{T}_p(\cdot)$ and $\mathcal{T}_w(\cdot)$ are the phone-level and word-level TransformerEncoder networks [26], respectively. The \mathcal{A}_w denotes the operation of averaging phone-level features to be word-level and \mathcal{U} refers to the utterance-level output processing network.

2.2. Proposed Method

2.2.1. Linguistic-Acoustic Similarity

Eq. (3) gives the two frequently used means of producing the phone-level pronunciation quality representation. Neither add_phone nor concat_phone explicitly measures the degree of mismatch between what the native speaker should pronounce (i.e., phone embedding) and how the L2 speaker actually pronounces (i.e., phone-level acoustic embedding). To investigate the effect of modeling this pronunciation deviation in a more explicit manner, this study proposes a novel linguistic-acoustic similarity based learning method as illustrated in Figure 1 (d). The phone-level pronunciation quality representation \mathbf{P}_Q is calculated as in Eq. (2) but in a slightly different way:

$$\mathbf{P}_Q = [\text{MLP}([\mathbf{H}_a; \mathbf{E}_p]); s], \quad (5)$$

where s denotes a linguistic-acoustic similarity measure, which is given by the cosine similarity between \mathbf{H}_a and \mathbf{E}_p ,

$$s = \text{cosine}(\mathbf{H}_a, \mathbf{E}_p) \quad (6)$$

The computation of utterance-level predicted score \hat{y} remains unchanged.

2.2.2. GOP Pre-training

GOP is widely used for measuring phone-level pronunciation quality of non-native speech, i.e., how close the pronunciation is to that of a native speaker [7]. To enable the phone-level linguistic-acoustic similarity to reflect pronunciation quality more accurately (e.g., a lower similarity indicates a higher possibility of mispronunciation), we use GOP score to guide the proposed similarity based learning. Figure 2 shows the boxplot of the GOP pre-training of the phone-level preprocessing network. Here, the used GOP score [12] for phone p is calculated as follows:

$$\text{GOP}(p) = \mathcal{LPP}(p) - \max_{q \in Q} \mathcal{LPP}(q), \quad (7)$$

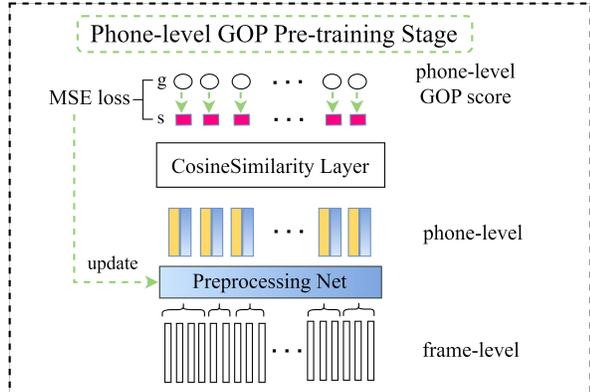


Fig. 2: The diagram of the phone-level GOP pre-training stage. Note that only the preprocessing network is updated at this stage.

where p is the phone in consideration and Q is the whole phone set. $\mathcal{L}^{PP}(p)$ is the log phone posterior and is computed as $\log p(p|\mathbf{o}; t_s, t_e)$, where t_s and t_e are the start and end frame indexes of phone p , and \mathbf{o} are the corresponding acoustic observations. Note that the phone-level time stamps and log posteriors are obtained using the extra acoustic model described in Section 3.2.1.

Mean squared error (MSE) between the cosine similarity s and phone-level GOP score g are used as the target loss function. We normalize both s and g into a range of $[0,1]$. Note that only the parameters of preprocessing network $\mathcal{F}(\cdot)$, as shown in Eq. (1), are optimized at this stage, which results in an updated version of the phone-level acoustic embedding \mathbf{H}_a and phone embedding \mathbf{E}_p .

3. EXPERIMENTAL SETUP

3.1. Datasets

Two L2 speech datasets are used in this study, namely Speechocean762 and ByteRate. Speechocean762 [27] is an open-sourced corpus designed for pronunciation assessment, in which 5,000 English utterances are collected from 250 learners. The corpus is split into train/test sets of equal size, each with 2,500 utterances from 125 English learners. Each utterance is rated by five experts in a range of 0 to 10, and the median value of the five scores is selected as the final score. ByteRate is an internal dataset at ByteDance, including a total of 10k utterances from 4k English learners. The train/dev/test sets are split as 3k/5k/2k, respectively. Each utterance is rated by three experts in a range of 0 to 4, and the final score is the average of the scores by all three experts. For both datasets, a higher rating indicates more native-like pronunciation and vice versa, and the scores are normalized into a range of 0 to 1. The first language (L1) of all L2 speakers is Mandarin.

3.2. Model Configurations

3.2.1. Acoustic Model

The deep feedforward sequential memory network and HMM, i.e., DFSMN-HMM, is adopted as the acoustic model [28]. DFSMN consists of 2 convolution layers and 24 FSMN layers followed by two FC layers. The input features are 39-dimension Mel-frequency cepstral coefficients (MFCCs). The acoustic model is trained on about 970-hour English speech, including an internal corpus of 10 hours of non-native English speech by L1 Mandarin learners and 960 hours

Table 1: The detailed structure of the proposed pronunciation scorer. LN refers to LayerNorm. Concat. is short for concatenation. Note that the time sequence information is omitted here.

Network	Structure	in \times out size	
Preprocessing	[FC, LN, Tanh]	512×32	
	[Embedding, LN, Tanh]	1×32	
Projection	Concat. of \mathbf{H}_a and \mathbf{E}_p	$[32, 32] \times 64$	
	MLP	[FC, ReLU]	64×32
		FC	32×32
	Concat. with s	$[32, 1] \times 33$	
Phone-level TransformerEncoder	[LN, FC, Tanh]	33×32	
	att_dim: 32, nhead: 4, ff_dim: 32, nlayer: 1	32×32	
Word-level TransformerEncoder	[FC, Tanh]	32×32	
	att_dim: 32, nhead: 4, ff_dim: 32, nlayer: 1	32×32	
Output	[FC, Sigmoid]	32×1	

of native English speech from the Librispeech corpus (Libri) [29]. 512-dimensional deep feature is extracted from the penultimate layer of the acoustic model. The same acoustic model is used to force-align speech with the corresponding text prompt to obtain phone-level time stamps and compute GOP scores as shown in Eq. (7).

3.2.2. Pronunciation Scorer

Table 1 presents the detailed network configuration of the pronunciation scorer. The output dimension of preprocessing, projection, and TransformerEncoder networks is equal to 32 [19]. For training the pronunciation scoring network, MSE between the predicted scores and the true scores is used as the loss function to be minimized. The training setups for pronunciation scorer training and GOP pre-training (Section 2.2.2) are the same. The Adam optimizer is utilized with a learning rate of 0.002 [30]. The maximum number of epochs is set as 50, and early stopping is activated if the loss stops decreasing for seven consecutive epochs. It should be noted that the GOP pre-training stage does not involve any additional speech data.

4. RESULTS AND ANALYSIS

In this section, we present the experimental results of the proposed and baseline systems, analyze the effect of the GOP pre-training stage by comparing system performance using acoustic models trained on different amounts of non-native data, and further examine the linguistic-acoustic characteristics captured by the learned phone embeddings. For performance evaluation, PCC between machine-predicted scores and human-predicted scores is calculated.

4.1. Performance of the Proposed and Baseline Systems

Table 2 presents the results of different systems on ByteRate and Speechocean762 datasets, respectively. We first examine the effectiveness of using phone-level linguistic-acoustic similarity for pronunciation assessment. Compared to add_phone and concat_phone baselines, the proposed method improves the performance by a large margin on both ByteRate ($\uparrow 0.06$ PCC) and Speechocean762 ($\uparrow 0.04$ PCC) datasets. The results suggest that the proposed linguistic-acoustic similarity can capture pronunciation deviation more effectively for pronunciation scoring.

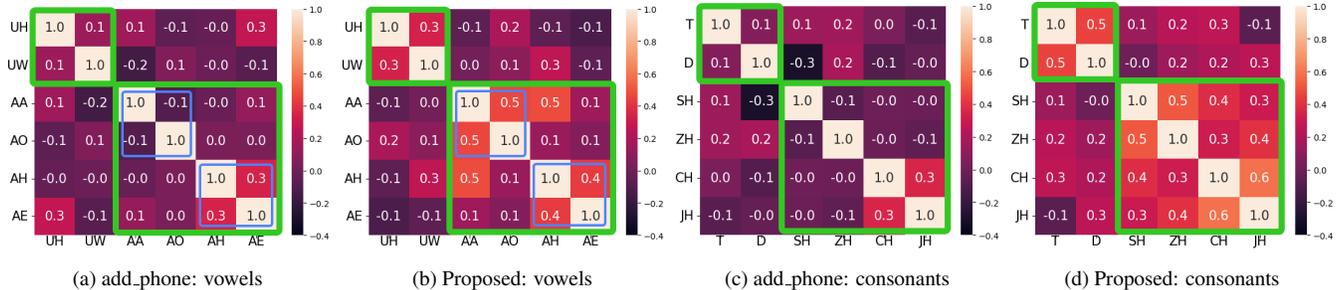


Fig. 3: Similarity heatmaps of phone embeddings between different phones (show-case examples).

Table 2: The PCC results of the proposed and baseline systems. Note that the proposed and baseline systems only differ in how to produce that phone-level pronunciation quality representation.

Datasets	ByteRate		Speechocean762	
	✗	✓	✗	✓
add_phone	0.764	0.781	0.598	0.610
concat_phone	0.763	0.823	0.618	0.652
Proposed	0.825	0.858	0.644	0.702

We then examine the effectiveness of GOP pre-training for pronunciation assessment. It is observed that, for all the three approaches, add_phone, concat_phone and the proposed system, the systems with GOP pre-training consistently outperform their counter-parts in terms of PCC. The proposed method achieves the best performance, with a PCC of 0.858 and 0.702 on ByteRate and Speechocean762 datasets, respectively. Hence, we conclude that the proposed phone-level linguistic-acoustic similarity framework with GOP pre-training has a clear advantage over the baselines.

Table 3: The PCC results of the proposed system using two different acoustic models (AM) trained with different amount of L2 speech.

Datasets	ByteRate		Speechocean762	
	✗	✓	✗	✓
AM: Libri + 10h	0.825	0.858	0.644	0.702
AM: Libri + 4000h	0.860	0.893	0.704	0.766

4.2. GOP Pre-training Stage: Less Can Be More

Previous research has shown the benefits of introducing more non-native data in acoustic model training for L2 pronunciation assessment [19, 21]. Acoustic models trained with both native and non-native data could provide more accurate phoneme segmentation of the L2 speech, hence better L2 phone representations for subsequent modeling processes. Unfortunately, non-native data of large size and high-quality annotation is not always available. In this study, we further examine how the proposed GOP pre-training process could help accommodate a lack of non-native speech data by conducting two more experiments which differ only in the amount of non-native data used during acoustic model training: 10 hours vs. 4,000 hours. The results are given in Table 3, which show that: (1) Including more non-native data in acoustic model training improves

the system performance which is consistent with previous findings; (2) GOP pre-training is beneficial for both the 10 hour and 4,000 hours of non-native data conditions; and (3) The results of including 10h non-native data with GOP pre-training are comparable with results including 4,000 hour non-native data without GOP pre-training, suggesting that the GOP pre-training process could serve as an alternative when the amount of non-native data is limited.

4.3. Linguistic-Acoustic Attributes of the Phone Embeddings

In this section, we examine how (well) the learned phone embedding could relate to linguistic-acoustic attributes of its corresponding phoneme. In particular, we plot the similarity heatmaps between the phonemes based on the cosine similarity of their respective phone embeddings. The results for a group of vowels and a group of consonants are given in Figure 3 as an example. In Figure 3, (a) and (c) show the results by the add_phone approach, (b) and (d) by the proposed approach. Figure 3 (b) clearly shows the pattern that the six vowels could be firstly divided into two clusters [UH, UW] and [AA, AO, AH, AE], and the second cluster could be further divided into two smaller clusters [AA, AO] and [AH, AE]. Specifically, the first two clusters, i.e., [UH, UW] and [AA, AO, AH, AE], differ in terms of vowel height, and the second in terms of tenseness. Similarly, in Figure 3 (d), the six consonants seem to form two clusters [T, D] and [SH, ZH, CH, JH], with the sounds in the first cluster being plosives and those in the second one being fricatives or affricates. In either the vowel or the consonant group, similar patterns could not be observed from the phone embeddings obtained by the add_phone approach. This shows that the phone embeddings learned by the proposed approach could reflect linguistic-acoustic attributes of their corresponding phonemes. Thus they are believed to provide more accurate reference representation of phone-level pronunciation.

5. CONCLUSION

In this paper, we proposed to use linguistic-acoustic similarity as additional feature to explicitly measure phone-level pronunciation deviation for pronunciation assessment. Moreover, a phone-level GOP pre-training stage was also proposed, which leads to better network initialization and more meaningful acoustic and phone embedding learning. The experiments conducted on both ByteRate and Speechocean762 datasets suggested that both linguistic-acoustic similarity and GOP pre-training contribute to the performance improvement in terms of PCC. It is also shown that the phone embeddings learned in the proposed approach can capture linguistic-acoustic attributes of standard pronunciation as reference. In the future, we plan to improve the system by using more contextualized acoustic features (e.g. wav2vec2.0) for the linguistic-acoustic similarity calculation.

6. REFERENCES

- [1] Silke M Witt, “Automatic error detection in pronunciation training: Where we are and where we need to go,” in *International Symposium on automatic detection on errors in pronunciation training*, 2012, pp. 1–8.
- [2] Jonás Fouz-González, “Trends and directions in computer-assisted pronunciation training,” *Investigating English Pronunciation*, pp. 314–342, 2015.
- [3] Nancy F Chen and Haizhou Li, “Computer-assisted pronunciation training: From pronunciation scoring towards spoken language learning,” in *Proc. of APSIPA*. IEEE, 2016, pp. 1–7.
- [4] Wei Li, *Improving mispronunciation detection and enriching diagnostic feedback for non-native learners of Mandarin*, Ph.D. thesis, Georgia Institute of Technology, 2019.
- [5] Pamela M Rogerson-Revell, “Computer-assisted pronunciation training (CAPT): Current issues and future directions,” *RELC Journal*, vol. 52, no. 1, pp. 189–205, 2021.
- [6] Leonardo Neumeyer, Horacio Franco, Mitchel Weintraub, and Patti Price, “Automatic text-independent pronunciation scoring of foreign language student speech,” in *Proc. of ICSLP*. IEEE, 1996, vol. 3, pp. 1457–1460.
- [7] Silke M Witt and Steve J Young, “Phone-level pronunciation scoring and assessment for interactive language learning,” *Speech communication*, vol. 30, no. 2-3, pp. 95–108, 2000.
- [8] Sandra Kanters, Catia Cucchiari, and Helmer Strik, “The goodness of pronunciation algorithm: a detailed performance study,” 2009.
- [9] Helmer Strik, Khiet Truong, Febe De Wet, and Catia Cucchiari, “Comparing different approaches for automatic pronunciation error detection,” *Speech communication*, vol. 51, no. 10, pp. 845–852, 2009.
- [10] Hao Huang, Haihua Xu, Ying Hu, and Gang Zhou, “A transfer learning approach to goodness of pronunciation based automatic mispronunciation detection,” *The Journal of the Acoustical Society of America*, vol. 142, no. 5, pp. 3165–3177, 2017.
- [11] Hyuksu Ryu and Minhwa Chung, “Mispronunciation Diagnosis of L2 English at Articulatory Level Using Articulatory Goodness-Of-Pronunciation Features,” in *SLaTE*, 2017, pp. 65–70.
- [12] Wenping Hu, Yao Qian, Frank K Soong, and Yong Wang, “Improved mispronunciation detection with deep neural network trained acoustic models and transfer learning based logistic regression classifiers,” *Speech Communication*, vol. 67, pp. 154–166, 2015.
- [13] Sweekar Sudhakara, Manoj Kumar Ramanathi, Chiranjeevi Yarra, and Prasanta Kumar Ghosh, “An Improved Goodness of Pronunciation (GoP) Measure for Pronunciation Evaluation with DNN-HMM System Considering HMM Transition Probabilities,” in *Proc. of Interspeech*, 2019, pp. 954–958.
- [14] Jiatong Shi, Nan Huo, and Qin Jin, “Context-aware goodness of pronunciation for computer-assisted pronunciation training,” *arXiv preprint arXiv:2008.08647*, 2020.
- [15] Ann Lee and James Glass, “A comparison-based approach to mispronunciation detection,” in *Proc. of SLT*. IEEE, 2012, pp. 382–387.
- [16] Ann Lee and James Glass, “Pronunciation assessment via a comparison-based system,” in *Speech and Language Technology in Education*, 2013.
- [17] Junwei Yue, Fumiya Shiozawa, Shohei Toyama, Yutaka Yamauchi, Kayoko Ito, Daisuke Saito, and Nobuaki Minematsu, “Automatic Scoring of Shadowing Speech Based on DNN Posteriors and Their DTW,” in *Proc. of Interspeech*, 2017, pp. 1422–1426.
- [18] Binghuai Lin, Liyuan Wang, Xiaoli Feng, and Jinsong Zhang, “Automatic Scoring at Multi-Granularity for L2 Pronunciation,” in *Interspeech*, 2020, pp. 3022–3026.
- [19] Binghuai Lin and Liyuan Wang, “Deep feature transfer learning for automatic pronunciation assessment,” in *Proc. of Interspeech*, 2021, pp. 4438–4442.
- [20] Kaiqi Fu, Shaojun Gao, Kai Wang, Wei Li, Xiaohai Tian, and Zejun Ma, “Improving non-native word-level pronunciation scoring with phone-level mixup data augmentation and multi-source information,” *arXiv preprint arXiv:2203.01826*, 2022.
- [21] Yuan Gong, Ziyi Chen, Iek-Heng Chu, Peng Chang, and James Glass, “Transformer-based multi-aspect multi-granularity non-native english speaker pronunciation assessment,” in *Proc. of ICASSP*. IEEE, 2022, pp. 7262–7266.
- [22] Fu-An Chao, Tien-Hong Lo, Tzu-I Wu, Yao-Ting Sung, and Berlin Chen, “3M: An Effective Multi-view, Multi-granularity, and Multi-aspect Modeling Approach to English Pronunciation Assessment,” *arXiv preprint arXiv:2208.09110*, 2022.
- [23] Huayun Zhang, Ke Shi, and Nancy F Chen, “Multilingual Speech Evaluation: Case Studies on English, Malay and Tamil,” *arXiv preprint arXiv:2107.03675*, 2021.
- [24] Binghuai Lin and Liyuan Wang, “A neural network-based noise compensation method for pronunciation assessment,” in *Interspeech*, 2021, pp. 3939–3943.
- [25] Qijie Shao, Jinghao Yan, Jian Kang, Pengcheng Guo, Xian Shi, Pengfei Hu, and Lei Xie, “Linguistic-acoustic similarity based accent shift for accent recognition,” *arXiv preprint arXiv:2204.03398*, 2022.
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Proc. of NeurIPS*, vol. 30, 2017.
- [27] Junbo Zhang, Zhiwen Zhang, Yongqing Wang, Zhiyong Yan, Qiong Song, Yukai Huang, Ke Li, Daniel Povey, and Yujun Wang, “speechocean762: An open-source non-native English speech corpus for pronunciation assessment,” *arXiv preprint arXiv:2104.01378*, 2021.
- [28] Shiliang Zhang, Ming Lei, Zhijie Yan, and Lirong Dai, “Deep-FSMN for large vocabulary continuous speech recognition,” in *Proc. of ICASSP*. IEEE, 2018, pp. 5869–5873.
- [29] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *Proc. of ICASSP*. IEEE, 2015, pp. 5206–5210.
- [30] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.