

Improving Mispronunciation Detection with Wav2vec2-based Momentum Pseudo-Labeling for Accentedness and Intelligibility Assessment

Mu Yang¹, Kevin Hirschi², Stephen D. Looney³, Okim Kang², John H. L. Hansen¹

¹Center for Robust Speech Systems (CRSS), University of Texas at Dallas, Richardson, TX, USA

²Northern Arizona University, Flagstaff, AZ, USA

³Pennsylvania State University, State College, PA, USA

¹{mu.yang, john.hansen}@utdallas.edu,

²{kevin.hirschi, okim.kang}@nau.edu, ³sd116@psu.edu

Abstract

Current leading mispronunciation detection and diagnosis (MDD) systems achieve promising performance via end-to-end phoneme recognition. One challenge of such end-to-end solutions is the scarcity of human-annotated phonemes on natural L2 speech. In this work, we leverage unlabeled L2 speech via a pseudo-labeling (PL) procedure and extend the fine-tuning approach based on pre-trained self-supervised learning (SSL) models. Specifically, we use Wav2vec 2.0 as our SSL model, and fine-tune it using original labeled L2 speech samples plus the created pseudo-labeled L2 speech samples. Our pseudo labels are dynamic and are produced by an ensemble of the on-line model on-the-fly, which ensures that our model is robust to pseudo label noise. We show that fine-tuning with pseudo labels achieves a 5.35% phoneme error rate reduction and 2.48% MDD F1 score improvement over a labeled-samples-only fine-tuning baseline. The proposed PL method is also shown to outperform conventional offline PL methods. Compared to the state-of-the-art MDD systems, our MDD solution produces a more accurate and consistent phonetic error diagnosis. In addition, we conduct an open test on a separate UTD-4Accents dataset, where our system recognition outputs show a strong correlation with human perception, based on accentedness and intelligibility.

Index Terms: Mispronunciation detection and diagnosis, wav2vec 2.0, pseudo-labeling, intelligibility assessment

1. Introduction

Second-language (L2) English learners typically present accents and mispronunciations, which highly impact their intelligibility in practical communication. In recent years, Computer Aided Pronunciation Training (CAPT) tools have been developed to provide diagnosis and feedback on phonetic-level errors (phoneme substitution, deletion, insertion [1, 2, 3, 4, 5]) and prosodic-level errors (e.g. lexical stress, intonation [6]). In this study, we focus on detecting phonetic-level pronunciation errors for L2 speech intelligibility and accentedness assessment.

Currently, most phonetic-level mispronunciation detection and diagnosis (MDD) systems perform end-to-end phoneme recognition on L2 speech, based on deep neural network (DNN) architectures [3, 4, 5, 7, 8]. One of the challenges of training such DNNs is data sparsity, due to the laborious process of annotating perceived phonemes on L2 speech. To address this issue, the “pre-training + fine-tuning” scheme has been shown to be effective [9]: in pre-training stage a model was trained on external large scale unlabeled data using self-supervised learning (SSL) objectives, and then fine-tuned on the data from

the downstream task with task-specific supervision. In speech realm, multiple SSL pre-trained models have been proposed [10, 11, 12, 13], and have shown promising results on many downstream tasks [14, 15, 16], including MDD task [4].

However, since data of the target task is limited, the vanilla fine-tuning approach may not be ideal due to domain mismatch between the pre-training and the target task [17, 18]. One solution is to use unlabeled data from the target domain [16, 17]. On MDD task, how to leverage unlabeled L2 speech remains unexplored. We approach this problem from a semi-supervised learning perspective based on pseudo-labeling (PL). We use Wav2vec 2.0 [10] as the SSL pre-trained model and extend the “pre-training + fine-tuning” scheme with one additional fine-tuning stage where pseudo-labeled L2 utterances are included in training. We propose to employ a recent momentum pseudo-labeling (MPL) method [19]. Unlike the conventional PL methods [20, 21, 22], MPL generates pseudo labels in a dynamic and online manner via teacher-student training: the online student model is trained using pseudo labels generated by an offline teacher model. The teacher model maintains a momentum-based moving average of the weights of the online model, which can be seen as an ensemble of the student model. This makes the online model robust to pseudo label noise and stabilizes training on unlabeled samples. We show that fine-tuning with MPL improves a vanilla fine-tuning baseline by 5.35% in phoneme error rate (PER), and 2.48% in MDD F1 score.

In addition, we take one step forward towards using the MDD model for automatic L2 speech intelligibility and accentedness assessment. We conduct an open test of our MDD model on a separate Indian-accented L2 English corpus. Through a human listening test, we show that the phoneme recognition performance of the MDD model has strong correlations with human ratings of L2 speech intelligibility and accentedness. This finding reveals the alignment between the MDD model prediction and human perception.¹

2. Related work

MDD. Goodness-of-pronunciation (GOP) is among the first DNN-based methods to MDD, which relies on phone posterior outputs from an automatic speech recognizer (ASR) [1, 2, 23] to evaluate phonetic errors. More recently, end-to-end phoneme recognition has been studied [3, 4, 5, 8, 24], among which [4] and [24] also explored fine-tuning Wav2vec 2.0. Our proposed method differs from them in that we investigate the usage of

¹We provide an audio demo at https://mu-y.github.io/speech_samples/mdd_IS22/. Code will be available at <https://github.com/Mu-Y/mpl-mdd>.

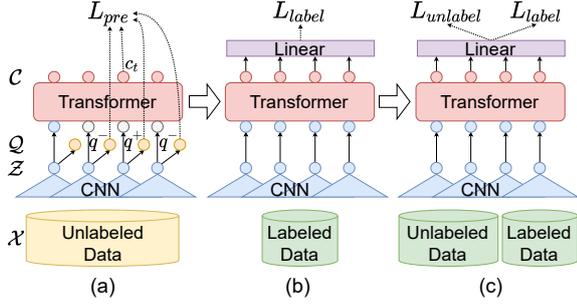


Figure 1: *Three-stage training: we extend the conventional pre-training (a) + fine-tuning (b) pipeline with an additional fine-tuning step with MPL (c). Training data from different domains are marked with different colors. Better viewed in color.*

unlabeled target domain speech to enhance MDD performance. **Pseudo-labeling in speech domain.** PL has been a widely used approach for semi-supervised ASR. In general, these methods can be divided into offline and online PL by the generation scheme of pseudo labels. Offline PL methods use a separately-trained teacher model to assign pseudo labels for unlabeled samples. A student model is then trained on labeled plus pseudo-labeled samples [25]. Filtering heuristics [20, 22] and iterative training [21] were shown to be useful to improve PL quality. On the other hand, in online PL methods, pseudo labels are generated *on-the-fly* by the online model itself [19, 26, 27]. We adopt the PL method in [19], but unlike [19], we combine PL with Wav2vec 2.0 fine-tuning, with phonemes as targets.

3. Method

We first review the pre-training and fine-tuning of Wav2vec 2.0 (Section. 3.1 and 3.2), and then describe our proposed fine-tuning method with MPL in Section. 3.3.

3.1. Wav2vec 2.0

Wav2vec 2.0 consists of Convolutional Neural Network (CNN) and Transformer layers (Figure. 1 (a)). The CNNs work as a feature extractor which converts input audio waveform \mathcal{X} into latent representation \mathcal{Z} . Before fed to Transformer layers, \mathcal{Z} is randomly masked by a certain portion (marked by grey in Figure. 1 (a)). The Transformer layers then contextualize \mathcal{Z} into \mathcal{C} . The unmasked latent representation \mathcal{Z} is further discretized to \mathcal{Q} via a learnable codebook. Given the contextualized representation c_t at masked time step t , we denote the discretized representation for time step t as q^+ , and q^- for other masked steps. During pre-training, Wav2vec 2.0 is trained by Contrastive Loss which aims to distinguish the true underlying discretized representation q^+ for each masked step t from those at other masked positions (q^-), based on the contextualized representation c_t . The full SSL loss, denoted as \mathcal{L}_{pre} , is a weighted sum of the Contrastive Loss and a codebook diversity loss [10]. Note that the model trained with such unsupervised SSL objective can be further fine-tuned by an ASR task using the text transcriptions of the pre-training audio, if available, to gain audio-text supervision [10]. We discuss the impact of this additional supervision on the downstream MDD task in Section. 5.1.

3.2. Fine-tuning

After pre-training, we add linear layers on top of the Transformer layers, and remove the discretization module. The en-

tire model (with CNN layers frozen) is fine-tuned end-to-end on the downstream L2 speech phoneme recognition task using connectionist temporal classification (CTC) loss [28] (Figure. 1 (b)). Let $X = (x_1, x_2, \dots, x_T)$ denote an input audio waveform, $Y = (y_1, y_2, \dots, y_L)$ denote training targets, which in our MDD task are the human-labeled phoneme sequences (i.e. what the L2 speaker actually pronounced). The CTC loss L_{label} on labeled samples can be expressed as

$$L_{label}(\theta) = -\log \sum_{A \in \beta(Y)} \prod_i P_{\theta}(a_i | X) \quad (1)$$

where $A = (a_1, a_2, \dots, a_T)$ denotes a compatible latent alignment between X and Y , $\beta(Y)$ denotes the set of all such compatible alignments, and θ denotes model parameters. We apply speed perturbation [29] on X , plus a modified SpecAugment on latent representations \mathcal{Z} as in [10] for data augmentation.

3.3. Fine-tuning with momentum pseudo-labeling

In addition to the fine-tuning stage where only labeled samples are used, we consider including unlabeled samples from the target domain into fine-tuning (Figure. 1 (c)). Formally, given the available labeled L2 speech samples \mathcal{D}^L (same samples as in Figure. 1 (b)) and additional unlabeled speech samples \mathcal{D}^U , our goal is to continually learn a new model ξ based on a base model θ learned in the previous fine-tuning stage, using both \mathcal{D}^L and \mathcal{D}^U . We propose to leverage the unlabeled samples via momentum pseudo-labeling (MPL) [19]. In MPL, an offline teacher model ϕ is used to assign pseudo labels for the unlabeled samples, which guide the learning of an online student model ξ . Pseudo labels \hat{Y} are inferred by the teacher model ϕ :

$$\hat{Y} = \underset{Y}{\operatorname{argmax}} P_{\phi}(Y | X), X \in \mathcal{D}^U \quad (2)$$

where we use argmax to represent greedy CTC decoding. Then, similar to Equation. 1, \hat{Y} can be used as the training targets of the unlabeled samples for the student model ξ :

$$L_{unlabel}(\xi) = -\log \sum_{A \in \beta(\hat{Y})} \prod_i P_{\xi}(a_i | X), X \in \mathcal{D}^U \quad (3)$$

The online model can then be trained on both labeled and unlabeled samples using a unified loss L :

$$L(\xi) = \begin{cases} L_{label}(\xi) & (X, Y) \in \mathcal{D}^L \\ L_{unlabel}(\xi) & X \in \mathcal{D}^U \end{cases} \quad (4)$$

To create dynamic pseudo labels, the teacher model ϕ is also updated during training by a moving average of itself and the online student model ξ , controlled by a momentum factor α :

$$\phi \leftarrow \alpha \phi + (1 - \alpha) \xi \quad (5)$$

The momentum factor here controls the update magnitude of ϕ and makes ϕ evolve more smoothly than ξ , which prevents the pseudo labels from changing drastically and helps stabilize training. We follow the heuristic in [19] to determine α based on total training steps. Both ξ and ϕ are initialized with θ before entering MPL fine-tuning, and are updated according to Equation. 4 and 5 respectively at each training step. The data augmentation in Section. 3.2 is also applied in this stage. Compared to static PL methods [20, 22], since the teacher model can be seen as an ensemble of the online model at different training steps, MPL is expected to be more robust to pseudo label noise. We compare MPL with static PL in Section. 5.2.

Table 1: *MDD evaluation metrics and PER on L2-ARCTIC test set. Numbers before and after “/” represent the percentage and the absolute number of occurrences of a particular case. P and R stand for Precision and Recall, respectively. PER numbers with a star (*) are not directly comparable to ours due to different settings (see Section. 5.3).*

Models		Correct Pronunciations		Mispronunciations			P (%)	R (%)	F1 (%)	PER (%)
		True Accept (%/#)	False Reject (%/#)	False Accept (%/#)	True Reject (%/#)					
					Corr. Diag.	Err. Diag.				
Prior	CTC-Attn + Anti-Phone [8]	–	–	–	–	–	46.57	70.28	56.02	–
	wav2vec2-large-lv60 [4]	94.01/24198	5.99/1542	43.37/1850	68.08/1645	31.91/771	61.04	56.63	58.75	16.01*
	wav2vec2-large-XLSR [4]	94.57/24343	5.43/1397	43.95/1875	65.75/1572	34.25/819	63.12	56.05	59.37	15.43*
Ours	wav2vec2-base	92.84/23873	7.16/1841	46.07/1977	75.84/1755	24.16/559	55.69	53.93	54.80	15.52
	+ one-shot PL (scratch)	93.26/23982	6.74/1732	46.05/1976	76.11/1762	23.89/553	57.20	53.95	55.53	14.85
	+ one-shot PL (continual)	93.16/23955	6.84/1759	46.17/1981	76.19/1760	23.81/550	56.77	53.83	55.26	15.04
	+ MPL	93.54/24052	6.46/1662	45.84/1967	77.24/1795	22.76/529	58.30	54.16	56.16	14.69
	wav2vec2-base-960h	93.83/24128	6.17/1586	47.91/2056	76.64/1713	23.36/522	58.49	52.09	55.10	14.87
	+ MPL	94.40/24273	5.60/1441	48.80/2094	77.29/1698	22.71/499	60.39	51.20	55.42	14.36

4. Experimental setup

4.1. Datasets

L2-ARCTIC [30] is used to train our MDD model. L2-ARCTIC includes L2 speech from 24 non-native English speakers with different L1 backgrounds (Indian, Mandarin, Vietnamese, Korean, Arabic, Spanish). It provides human-labeled perceived phonemes for around 15% of the utterances per speaker. Following [4, 5], we set labeled samples from 6 speakers as test set and labeled samples for the remaining speakers as labeled training set. Unlabeled samples of the remaining 18 speakers are used as the unlabeled training set. We randomly split 10% of the labeled training set as development set. Statistics of our data splits are shown in Table. 2. Since L2-ARCTIC uses artificial *sil* tokens to represent phoneme deletions and insertions, to construct target training phoneme sequences, we remove the artificial *sil* tokens, while preserving the *sil* tokens that correspond to true pauses and silences.

UTD-4Accents [31] is an in-house dataset that consists of 4 English accents: US (native), Australian, Spanish and Indian. We use the Indian-accent part for the open test (see Section. 5.4), which includes 112 speakers (balanced for gender and age). The utterances are read speech from diverse domains, including general vocabulary, voice search, etc.

4.2. Evaluation

We evaluate the PER between recognized phonemes and training targets, as well as the MDD metrics following [4, 32]. For cases of correct pronunciations where human-perceived phonemes agree with canonical phonemes (i.e. the phonemes that a L2 speaker was supposed to pronounce), we have True Accept (TA) and False Reject (FR) cases based on whether model predictions match both canonical and perceived phonemes. In the same spirit, for mispronunciation cases where human perceived phonemes are inconsistent with canonical phonemes, we could have False Accept (FA) and True Reject (TR) cases. TR can be further divided into Correct Diagnosis and Erroneous Diagnosis, based on whether model predictions match human labels. Precision (P), Recall (R) and F1 can then

be computed from TR, FR, FA: $P = TR / (FR + TR)$; $R = TR / (FA + TR)$; $F1 = 2PR / (P + R)$.

4.3. Implementation details

We tune hyper-parameters on the development set. The best model (in terms of PER) on the development set is evaluated on the test set. We experimented with two pre-trained Hugging-Face [33] Wav2vec 2.0 models, *wav2vec2-base* and *wav2vec2-base-960h*. Both models have identical base-size network architecture and are pre-trained with the same SSL objective (Section. 3.1) on 960-hour LibriSpeech audio [34]. *wav2vec2-base-960h* has been additionally fine-tuned by an ASR task on LibriSpeech after pre-training, which learns explicit audio-text mapping compared to *wav2vec2-base*. We use separate Adam optimizers for linear layers and Wav2vec 2.0 layers, with fixed learning rates of $3e - 4$ and $1e - 5$, respectively. Gradient accumulation is used to obtain an effective batch size of 32. All experiments run for 50 epochs on a single NVIDIA RTX 2080 GPU. Our implementation is based on SpeechBrain toolkit [35].

5. Results

5.1. Effect of pseudo-labeling and pre-trained models

In Table. 1, we use *wav2vec2-base* and *wav2vec2-base-960h* to denote the vanilla fine-tuning baselines based on the corresponding pre-trained SSL models without using any unlabeled samples (i.e. stage (b) in Figure. 1). First, we can see that with MPL, both *wav2vec2-base* and *wav2vec2-base-960h* gain a significant improvement over the vanilla fine-tuning baseline, in terms of F1 score (*wav2vec2-base*: 56.16 vs. 54.80; *wav2vec2-base-960h*: 55.42 vs. 55.10) and PER (*wav2vec2-base*: 14.69 vs. 15.52; *wav2vec2-base-960h*: 14.36 vs. 14.87). This demonstrates the benefit of leveraging unlabeled L2 samples. Second, an interesting finding is that *wav2vec2-base-960h* produces more False Accepts, less False Rejects and less overall Reject cases than *wav2vec2-base*, leading to higher Precision but lower Recall. This means that it tends to be a more “tolerant” judge by rejecting less L2 pronunciations. One possible reason is that the extra audio-text supervision has biased it to the canonical pronunciations, which makes it “over-robust” to mispronunciations. However, for MDD task this may not be desired, because we expect the MDD model to faithfully reflect what a L2 speaker actually pronounced and the “over-robustness” may conceal the mispronunciations. In contrast, *wav2vec2-base* has a more balanced Precision and Recall.

Table 2: *L2-ARCTIC data splits and statistics.*

	Train		Development	Test
	labeled	unlabeled		
# utterances	2429	17381	268	900
# hours	2.51	17.73	0.27	0.88

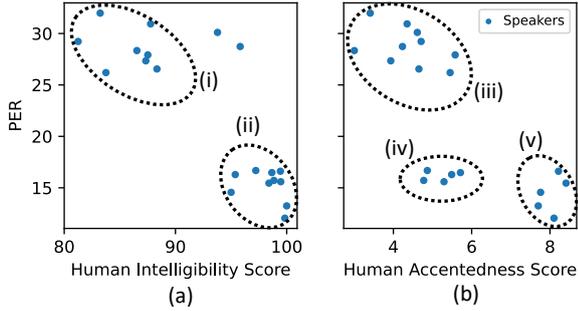


Figure 2: (a) *PER vs. Human Intelligibility Score*. Pearson Correlation = -0.84 ($p < 10^{-5}$); (b) *PER vs. Human Accentedness Score*. Pearson Correlation = -0.75 ($p < 10^{-3}$).

5.2. Effect of the momentum update mechanism

To study whether MPL is superior to static PL methods, we compare MPL with two PL baselines: **(1) one-shot PL (scratch)**: pseudo labels are generated offline by the model fine-tuned with labeled samples only (i.e. stage (b) in Figure. 1). A new model is then trained from scratch using the pseudo-labeled samples plus the original labeled samples; **(2) one-shot PL (continual)**: same as the above method, except that the new model is initialized with the weights of the fine-tuned model. Since in both baselines pseudo labels are fixed once generated, we refer to them as *one-shot*. Table. 1 compares MPL with the two baselines on *wav2vec2-base*. Although both baselines outperform the labeled-samples-only fine-tuning, a larger improvement comes from MPL, demonstrating the benefits of having dynamic pseudo labels. The performances of the two baselines are comparable, with *one-shot PL (continual)* being slightly worse. This may be caused by over-fitting.

5.3. Comparison with prior works

Finally, we compare our MDD models with the current leading MDD methods. [8] uses a CTC-Attention model with Anti-Phone augmentation. [4] is also based on fine-tuning Wav2vec 2.0 models and achieves state-of-the-art MDD performance. They used large-size Wav2vec 2.0 models which are pre-trained on larger-scale audio corpora. Models in [4] have 300M+ parameters, while ours have around 90M parameters.² From Table. 1, we can see that [8] achieves higher Recall but much lower Precision. Our proposed MPL model outperforms [8] in terms of overall F1. Compared with [4], although our proposed method does not outperform theirs, we observe a higher Correct Diagnosis and lower Erroneous Diagnosis in the percentage and the absolute number of occurrences. This implies that our models are able to provide more accurate pronunciation diagnosis feedback which may help L2 learners correct their mispronunciations more effectively. Note that the PER reported in [4] is not directly comparable to ours, as they did not pre-process the target phoneme sequences as we do (Section. 4.1).

5.4. Open test: Indian accent and intelligibility assessment

We hypothesize that a proper MDD model should perceive the L2 pronunciations in a similar way as humans, and thus there should exist a correlation between its phoneme recognition performance and L2 speech accentedness and comprehensibility,

²Due to computation resources limit, we are not able to run experiments on large models, but our proposed methods are generic and we leave applications on larger-size models for future investigation.

i.e. a higher PER (more mispronunciations) corresponds to a heavier accent and lower comprehensibility. Our goal is to investigate the existence of such relations, which shed some light on the applicability of the MDD model towards automatic intelligibility assessment for L2 speech.

For this purpose, we run an open test of our best MDD model (*wav2vec2-base* + MPL) on Indian-accented L2 speech from the UTD-4Accents dataset. We compute per-speaker PERs between the recognized phonemes and the canonical phonemes given by a grapheme-to-phoneme model.³ Then we select 10 highest-PER speakers and 10 lowest-PER speakers, and randomly sample 10 utterances for each of the 20 speakers. 17 human listeners score the accentedness (scale: 1-9 where 1 means heavy accent) and intelligibility (scale: 0-100 where 0 means not intelligible at all) of the sampled utterances. Each speaker receives 30+ ratings from different raters. We then aggregate per-speaker accentedness and intelligibility scores. All raters are graduate students at Northern Arizona University with 2+ years of experience in L2 pronunciation teaching or research. To create unbiased ratings, raters are presented with the L2 audio only, without any text transcription or other information.

We plot the per-speaker PER-Intelligibility and PER-Accentedness relations in Figure. 2 (a) and (b), respectively. Figure. 2 (a) shows that the speakers are roughly grouped into 2 clusters ((i) and (ii)). This is consistent with our expectation: since we selected 20 speakers with highest and lowest PER, we expect those speakers are also clustered into more-intelligible and less-intelligible groups. Further, the PERs and human intelligibility scores present a strong negative correlation, which means that higher PERs align with less intelligible speech. Such a negative correlation can also be observed in the PER-Accentedness plot (Figure. 2 (b)). Interestingly, besides the two aforementioned clusters, we observe an additional cluster (iv). Speakers in this cluster reside in cluster (ii) in Figure. 2 (a). This indicates that for these speakers, humans perceive relatively heavy accents, while they are still highly intelligible. This is possibly because apart from phonetic errors, accentedness is also highly impacted by prosodic factors, such as intonation, lexical stress, etc, which may not be as important for intelligibility. Since our MDD model only detects phonetic-level errors, assessing the accentedness of these speaker is beyond its capacity. In summary, the alignment between phoneme recognition of the MDD model and human perception has validated our motivation of using the MDD model as a component towards automatic L2 speech intelligibility assessment.

6. Conclusions and acknowledgements

We have presented an approach to use unlabeled L2 speech to enhance MDD performance via pseudo-labeling. In addition, we take one step forward towards using the MDD model for automatic L2 speech intelligibility and accentedness assessment. Through a human listening test, we have shown that the MDD model recognition performance shows a strong correlation with human perception. In future, we plan to include more speech attributes, such as lexical stress, speech rate, into the L2 speech intelligibility assessment framework. This study is supported by NSF EAGER CISE Project 2140415, and partially by the University of Texas at Dallas from the Distinguished University Chair in Telecommunications Engineering held by J. H. L. Hansen. We would like to thank the raters at Northern Arizona University for their participation in our listening test.

³<https://github.com/Kyubyong/g2p>

7. References

- [1] S. Sudhakara, M. K. Ramanathi, C. Yarra, and P. K. Ghosh, "An improved goodness of pronunciation (gop) measure for pronunciation evaluation with dnn-hmm system considering hmm transition probabilities." in *Proc. Interspeech 2019*, 2019, pp. 954–958.
- [2] W. Li, S. M. Siniscalchi, N. F. Chen, and C.-H. Lee, "Improving non-native mispronunciation detection and enriching diagnostic feedback with dnn-based speech attribute modeling," in *ICASSP 2016*. IEEE, 2016, pp. 6135–6139.
- [3] M. Wu, K. Li, W.-K. Leung, and H. Meng, "Transformer Based End-to-End Mispronunciation Detection and Diagnosis," in *Proc. Interspeech 2021*, 2021, pp. 3954–3958.
- [4] L. Peng, K. Fu, B. Lin, D. Ke, and J. Zhan, "A Study on Fine-Tuning wav2vec2.0 Model for the Task of Mispronunciation Detection and Diagnosis," in *Proc. Interspeech 2021*, 2021, pp. 4448–4452.
- [5] Y. Feng, G. Fu, Q. Chen, and K. Chen, "Sed-mdd: Towards sentence dependent end-to-end mispronunciation detection and diagnosis," in *ICASSP 2020*. IEEE, 2020, pp. 3492–3496.
- [6] D. Korzekwa, R. Barra-Chicote, S. Zaporowski, G. Beringer, J. Lorenzo-Trueba, A. Serafinowicz, J. Droppo, T. Drugman, and B. Kostek, "Detection of Lexical Stress Errors in Non-Native (L2) English with Data Augmentation and Attention," in *Proc. Interspeech 2021*, 2021, pp. 3915–3919.
- [7] T.-H. Lo, S.-Y. Weng, H.-J. Chang, and B. Chen, "An Effective End-to-End Modeling Approach for Mispronunciation Detection," in *Proc. Interspeech 2020*, 2020, pp. 3027–3031.
- [8] B.-C. Yan, M.-C. Wu, H.-T. Hung, and B. Chen, "An End-to-End Mispronunciation Detection System for L2 English Speech Leveraging Novel Anti-Phone Modeling," in *Proc. Interspeech 2020*, 2020, pp. 3032–3036.
- [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 4171–4186.
- [10] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.
- [11] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [12] C. Wang, Y. Wu, Y. Qian, K. Kumatani, S. Liu, F. Wei, M. Zeng, and X. Huang, "Unispeech: Unified speech representation learning with labeled and unlabeled data," in *International Conference on Machine Learning*. PMLR, 2021, pp. 10937–10947.
- [13] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao *et al.*, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *arXiv preprint arXiv:2110.13900*, 2021.
- [14] Y. Y. Lin, C.-M. Chien, J.-H. Lin, H.-y. Lee, and L.-s. Lee, "Fragmentvc: Any-to-any voice conversion by end-to-end extracting and fusing fine-grained voice fragments with attention," in *ICASSP 2021*. IEEE, 2021, pp. 5939–5943.
- [15] Z. Fan, M. Li, S. Zhou, and B. Xu, "Exploring wav2vec 2.0 on speaker verification and language identification," *arXiv preprint arXiv:2012.06185*, 2020.
- [16] L.-W. Chen and A. Rudnicky, "Exploring wav2vec 2.0 fine-tuning for improved speech emotion recognition," *arXiv preprint arXiv:2110.06309*, 2021.
- [17] S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith, "Don't stop pretraining: Adapt language models to domains and tasks," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 8342–8360.
- [18] L. M. Dery, P. Michel, A. Talwalkar, and G. Neubig, "Should we be pre-training? an argument for end-task aware training as an alternative," *arXiv preprint arXiv:2109.07437*, 2021.
- [19] Y. Higuchi, N. Moritz, J. L. Roux, and T. Hori, "Momentum pseudo-labeling for semi-supervised speech recognition," *arXiv preprint arXiv:2106.08922*, 2021.
- [20] J. Kahn, A. Lee, and A. Hannun, "Self-training for end-to-end speech recognition," in *ICASSP 2020*. IEEE, 2020, pp. 7084–7088.
- [21] Q. Xu, T. Likhomanenko, J. Kahn, A. Hannun, G. Synnaeve, and R. Collobert, "Iterative Pseudo-Labeling for Speech Recognition," in *Proc. Interspeech 2020*, 2020, pp. 1006–1010.
- [22] D. S. Park, Y. Zhang, Y. Jia, W. Han, C.-C. Chiu, B. Li, Y. Wu, and Q. V. Le, "Improved Noisy Student Training for Automatic Speech Recognition," in *Proc. Interspeech 2020*, 2020, pp. 2817–2821.
- [23] W. Hu, Y. Qian, F. K. Soong, and Y. 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.
- [24] X. Xu, Y. Kang, S. Cao, B. Lin, and L. Ma, "Explore wav2vec 2.0 for Mispronunciation Detection," in *Proc. Interspeech 2021*, 2021, pp. 4428–4432.
- [25] Q. Xu, A. Baevski, T. Likhomanenko, P. Tomasello, A. Conneau, R. Collobert, G. Synnaeve, and M. Auli, "Self-training and pre-training are complementary for speech recognition," in *ICASSP 2021*. IEEE, 2021, pp. 3030–3034.
- [26] Y. Chen, W. Wang, and C. Wang, "Semi-Supervised ASR by End-to-End Self-Training," in *Proc. Interspeech 2020*, 2020, pp. 2787–2791.
- [27] H. Zhu, L. Wang, Y. Hou, J. Wang, G. Cheng, P. Zhang, and Y. Yan, "Wav2vec-s: Semi-supervised pre-training for speech recognition," *arXiv preprint arXiv:2110.04484*, 2021.
- [28] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [29] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, "Audio augmentation for speech recognition," in *Proc. Interspeech 2015*, 2015, pp. 3586–3589.
- [30] G. Zhao, S. Sonsaat, A. Silpachai, I. Lucic, E. Chukharev-Hudilainen, J. Levis, and R. Gutierrez-Osuna, "L2-ARCTIC: A Non-native English Speech Corpus," in *Proc. Interspeech 2018*, 2018, pp. 2783–2787.
- [31] S. Ghorbani and J. H. Hansen, "Leveraging Native Language Information for Improved Accented Speech Recognition," in *Proc. Interspeech 2018*, 2018, pp. 2449–2453.
- [32] K. Li, X. Qian, and H. Meng, "Mispronunciation detection and diagnosis in l2 english speech using multidistribution deep neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 1, pp. 193–207, 2016.
- [33] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz *et al.*, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, 2020, pp. 38–45.
- [34] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *ICASSP 2015*. IEEE, 2015, pp. 5206–5210.
- [35] M. Ravanelli, T. Parcollet, P. Plantinga, A. Rouhe, S. Cornell, L. Lugosch, C. Subakan, N. Dawalatabad, A. Heba, J. Zhong *et al.*, "Speechbrain: A general-purpose speech toolkit," *arXiv preprint arXiv:2106.04624*, 2021.