

Phonological Level wav2vec2-based Mispronunciation Detection and Diagnosis Method

Mostafa Shahin^a, Julien Epps^a, Beena Ahmed^a

^a*School of Electrical and Computer Engineering, University of New South Wales, Sydney, 2052, NSW, Australia*

Abstract

The automatic identification and analysis of pronunciation errors, known as Mispronunciation Detection and Diagnosis (MDD) plays a crucial role in Computer Aided Pronunciation Learning (CAPL) tools such as Second-Language (L2) learning or speech therapy applications. Existing MDD methods relying on analysing phonemes can only detect categorical errors of phonemes that have an adequate amount of training data to be modelled. With the unpredictable nature of the pronunciation errors of non-native or disordered speakers and the scarcity of training datasets, it is unfeasible to model all types of mispronunciations. Moreover, phoneme-level MDD approaches have a limited ability to provide detailed diagnostic information about the error made. In this paper, we propose a low-level MDD approach based on the detection of speech attribute features. Speech attribute features break down phoneme production into elementary components that are directly related to the articulatory system leading to more formative feedback to the learner. We further propose a multi-label variant of the Connectionist Temporal Classification (CTC) approach to jointly model the non-mutually exclusive speech attributes using a single model. The pre-trained wav2vec2 model was employed as a core model for the speech attribute detector. The proposed method was applied to L2 speech corpora collected from English learners from different native languages. The proposed speech attribute MDD method was further compared to the traditional phoneme-level MDD and achieved a significantly lower False Acceptance Rate (FAR), False Rejection Rate (FRR), and Diagnostic Error Rate (DER) over all speech attributes compared to the phoneme-level equivalent.

Keywords: Speech attributes, Self-supervised learning, wav2vec2.0, Mispronunciation detection and diagnosis

1. Introduction

Unlike acoustic modeling for Automatic Speech Recognition (ASR) applications where an abstract model that can handle all variations of the same word is desirable, pronunciation assessment applications such as second language (L2) learning, speech therapy, and language proficiency tests, require accurate detection of any deviation from standard pronunciation.

Mispronunciation detection and diagnosis (MDD) aims to automatically detect pronunciation errors in speech productions and diagnose them by providing error details such as the error location and type as well as the pronunciation quality score. MDD is a key feature of Computer-Aided Pronunciation Learning (CAPL) systems.

The level of diagnostic information provided by the MDD model depends on the level of evaluation the trained model can perform which in turn is restricted by the level of annotation of the training data. Only if the annotation is at the phoneme-level, i.e., speech is annotated with its pronounced phoneme sequence, can a phoneme-level model be trained and used to provide phoneme-level diagnostic information, e.g., the location of the mispronounced phoneme and the type of the error (substitution, deletion, or insertion).

Phoneme-level assessment is the most common assessment level used in MDD [1–8]. Different approaches have been proposed to achieve phoneme-level MDD, including the scoring approach [6, 9–15], the rule-based approach [3, 5, 16, 17], the classification approach [18–26], and the free-phoneme recognition approach [27–30].

One of the limitations of phoneme-level MDD is that it can only diagnose categorical pronunciation errors or mispronounced phonemes that exist in the acoustic model. Uncategorical errors, where the pronounced phoneme is a distorted version of the modelled phoneme or a new phoneme borrowed from other languages, are difficult to detect with phoneme-level MDD. To handle uncategorical errors, the acoustic model needs to include phonemes from multiple languages as well as all possible pronunciation variations of each phoneme, which is infeasible. Moreover, diagnostic information obtained from the phoneme-level MDD can't be used to construct formative feedback with corrective instruction to the learner.

In this work, we propose a low-level MDD system that detects and diagnoses pronunciation errors at the speech attribute level. Speech attributes,

such as manners and places of articulations, provide a low-level description of sound production in terms of which articulators are involved and how these articulators move to produce a specific sound. Any alteration in these attributes causes a pronunciation error. Therefore, accurate modeling of these attributes instead of phonemes can pave the way to a fully automated and interactive pronunciation assessment application where the learner receives informative and diagnostic automatic feedback not only about the existence of incorrect pronunciation, but also how the error is made. Furthermore, modeling speech attributes can be performed solely using typically pronounced datasets which are abundantly available, unlike atypical datasets such as disordered or non-native speech. Additionally, speech attributes are common across most spoken languages enabling modeling with speech corpora from multiple languages [31].

Our proposed speech attribute detection model is a wav2vec2-based Sequence-to-Sequence (Seq2Seq) classification model trained using a multi-label variant of the Connectionist Temporal Classification (CTC) criteria to handle the non-mutually exclusive nature of speech attributes, i.e, the same phoneme can be characterized by multiple speech attributes. To validate our model, it was applied to English as L2 speech corpora collected from different L1 speakers. Additionally, we compared the detection and diagnostic accuracy of a state-of-the-art phoneme-level MDD method with our proposed speech attribute-based approach. We also conducted experiments to demonstrate the low-level diagnosis capability of the model and the potential formative feedback message it can provide.

This paper presents three key contributions. Firstly, we established a new benchmark for speech attribute detection by utilizing the wav2vec2 speech representation upstream model. Secondly, we introduced a novel multi-label CTC approach to simultaneously learn 35 non-mutually exclusive speech attributes. Lastly, we proposed a low-level MDD method based on our speech attribute detection model.

The paper is structured as follows: In Section 2, the limitations of current phoneme-level MDD approaches are discussed. Section 3 presents our proposed method, followed by a description of the speech corpora used in Section 4. The experimental setup is outlined in Section 5, while the results are demonstrated and discussed in Section 6. Finally, Section 7 concludes the paper.

2. Related Work

2.1. Mispronunciation Detection and Diagnosis (MDD)

The Goodness Of Pronunciation (GOP) algorithm [6] is the earliest and most successful phoneme-level speech assessment method utilized in several applications to measure phoneme-level pronunciation quality [32–38]. The GOP approximates the posterior probability of each phoneme by taking the ratio between the forced alignment likelihood and the maximum likelihood of the free-phone loop decoding using a Gaussian Mixture Model-Hidden-Markov Model (GMM-HMM) acoustic model. This score is then used as a confidence score representing how close the pronunciation is to the target phoneme. Subsequent methods have leveraged DNN acoustic modelling and proposed a DNN-based method to estimate the GOP [14, 15, 39].

Despite its success, the GOP algorithm is very sensitive to the quality of the acoustic model used. The acoustic model affects not only the estimate of the posterior probability but also the accuracy of time boundaries obtained from forced alignment. In addition, as the decision threshold is determined using a mispronunciation dataset, it can be error specific and thus very hard to generalise to different types of pronunciation errors.

Phoneme-level error detection has also been treated as a binary classification problem, with each phoneme classified as “correct” or “mispronounced” using conventional classification methods such as SVM [20], decision tree [18], Linear Discriminant Analysis (LDA) [19], DNN [24], etc. Mel Frequency Cepstral Coefficients (MFCCs) are the most common features used with the classification methods [18]. Formant frequencies and GOP scores were also utilised to classify between correct and incorrect phoneme pronunciation [21].

Although all these classifiers led to a significant improvement compared with confidence score methods such as GOP, they still need large amounts of accurately annotated non-native data to model the mispronounced phonemes. Moreover, the mispronounced data needs to include all possible pronunciation errors, which is usually not feasible to collect.

In [7], an anomaly detection-based model was trained solely on the standard English pronunciation, namely native English speakers, and tested on foreign-accented speech and disordered speech. The method treated the mispronunciation as a deviation (anomaly) from the standard pronunciation. To detect the anomalies a One-Class SVM (OCSVM) model was trained for each phoneme using speech attribute features, namely manners and places of articulation. The method was shown to outperform the DNN-based GOP

method in both disordered and foreign-accented speech. Recently, [26] investigated fine-tuning a self-supervised speech representation pre-trained model, namely wav2vec2 [25] to perform phoneme level binary classification of L2 speech pronunciation as correct/mispronounced. Both the scoring and classification approaches can only provide either a soft (score) or hard (binary) decision on the phoneme pronunciation without any detailed description of the type of error.

The Extended Recognition Network (ERN) method was proposed to provide more detailed descriptions of the pronunciation error. It can identify the location of the mispronounced phoneme, and the type of the error, and recognise the erroneous phoneme. Unlike forced alignment used in the scoring method which contains only the canonical phonetic transcription of the prompt word, the ERN extends the single path alignment to a network that contains the correct (canonical) path and the expected mispronounced paths on phoneme level based on phonological rules designed for a specific learning domain [3, 5, 16].

The design of the search network is crucial for the reliability of the ERN-based MDD systems. Hand-crafted phonological rules are the most common designing criterion for the ERN [3, 5, 8, 16, 17]. Data-driven phonological approaches have also been proposed to automatically create the ERN. [40] first performed an automatic phonetic alignment step between the canonical phonetic transcription and the corresponding L2 transcription then all pairs of mismatched phonemes along with their contextual phones were grouped to form the initial set of rules. Finally, a rule selection criterion was performed to select the most significant rules. In [41], the authors proposed a grapheme-to-phoneme-like model trained on phonetic transcriptions of L2 learners.

The decoding of the ERN is performed by an acoustic model that is trained either on standard pronunciation speech corpora, such as L1 native speakers [5, 40], or non-standard pronunciation speech corpora, such as L2 non-native speakers [42, 43], or child disordered speech [3]. As the ERN method was proposed around 20 years ago, the most common acoustic model used was the GMM-HMM acoustic model [5, 16, 40, 42, 44]. Later the DNN-HMM acoustic model was utilised [3, 8, 17, 43] and shown to outperform the GMM-HMM in the ERN methods specifically when trained on a small non-standard dataset [3]. Moreover, the standard acoustic model has been adapted to non-standard domains using techniques such as Maximum Likelihood Linear Regression (MLLR) for GMM-HMM-based models [16, 44] or transfer learning for DNN-HMM-based acoustic models [8]. [42] proposed a

discriminative training criterion to directly optimise the acoustic model for MDD. The authors incorporated False Acceptance Rate (FAR), False Rejection Rate (FRR), and Diagnostic Error Rate (DER) in the acoustic model objective function and trained it using L2 non-native speech corpora.

With significant improvements in End-to-End (End2End) deep learning-based acoustic models, most of the current SOTA MDD approaches have adopted free-phone recognition criteria. Unlike the ERN, free-phone recognition has no restricted search space and therefore any pronounced phoneme sequence can be captured. That is, free-phone recognition MDD systems work by estimating the pronounced phoneme sequence and the evaluation of the system is achieved by aligning the recognised phoneme sequence with the annotated one. The free-phoneme recognition model is adapted to the pronunciation learning domain by incorporating linguistic information from the prompts used in pronunciation learning which, in most cases, are pre-designed and available. This linguistic information is extracted on the character level [2], phonemic level [27–30], graphemic level [28], or a combination of different levels [28].

In [1] the CNN-RNN-CTC system was proposed as an early attempt of an End2End MDD method. The system has a straightforward phoneme recognition architecture trained using a combination of native and L2 speech corpora and makes use of the Connectionist Temporal Classification (CTC) loss to avoid direct alignment between the phoneme sequence and the speech signal.

Recent methods adopted an encoder-decoder mechanism to achieve MDD by using multiple encoders to encode the audio signal [45] along with the prompt sentence phonemes [27, 29], characters [2], or words [46]. The encoder architecture is commonly a stacked CNN-RNN [2, 27] while the decoder is mostly an attention-based network.

All the previous methods use hand-crafted features extracted from a short-time speech window, however, several recent SOTA speech recognition systems leverage the raw speech signal and use a learnable feature extraction network that can be integrated and trained within the whole network [47–49].

One of the shortcomings of the free-phoneme recognition method is that the phoneme sequence output is limited to the modelled phoneme set. Therefore, it is not able to detect distortion errors where the learner pronounces a distorted version of a phoneme that cannot be explicitly replaced with another phoneme within the target language phoneme set. This commonly occurs when the L2 learner is influenced by their L1 language.

In pronunciation learning, pronunciation errors made by the learner can be unpredictable. For instance, in L2 learning pronunciation errors are influenced by the proficiency of the learner and the degree of discrepancy between L1 and L2. Therefore, there are high variations in the way that the learner will adapt their pronunciation to try and match the target pronunciation. Given the limited amount of annotated mispronounced speech data available, it is infeasible to model all these variations. Existing MDD methods can only diagnose categorical pronunciations that can be modelled by the acoustic model, and struggle otherwise. Scoring approaches such as the GOP can help in detecting mispronunciations if the deviation from the correct pronunciation is high enough, however, the actual pronounced error and its type cannot be determined.

Here we have thus proposed a low-level MDD based on the speech attribute features. Speech attribute features, also known as phonological features, can break down the phoneme into elementary components that form the phoneme. These features include manners and places of articulations that describe the articulators' positions and movements during sound production. There are several advantages of using speech attributes in MDD: 1) the speech attributes are only limited by the possible positions of the articulators and thus shared among most spoken languages, 2) their models can be solely trained using correctly pronounced speech which alleviates the need for annotated mispronounced data, and 3) they can provide lower-level diagnostic information describing how the error is made and which attribute is missing allowing more formative feedback to be provided.

2.2. *Speech Attribute Modelling*

Speech attributes have been successfully utilized in various domains such as improving ASR performance [50, 51], language identification [31], speaker verification [52] and Mispronunciation Detection and Diagnosis (MDD) [53]. However, the speech attribute detection models are commonly trained at the frame level and hence frame-level labelling of the speech signal is required for all training samples. The attributes are obtained by first performing a forced alignment on a phoneme level and mapping phonemes to their corresponding attributes. The performance of the resultant speech attribute model is thus highly dependent on the quality of the acoustic model used in the forced alignment process. To overcome this constraint, some studies have explored using the CTC learning criteria, which eliminates the need for prior alignment between the input signal and output label [54–57].

The speech attributes are non-mutually exclusive; each phoneme can appear in multiple attributes. For example, the phoneme /z/ is fricative, voiced and alveolar. Consequently, modelling all attributes with a single model becomes a multi-label classification problem. Existing approaches address this by constructing multiple models, with each model dedicated to an individual attribute [58, 59], or a set of mutually exclusive attributes [60, 61]. Unfortunately, using multiple models is impractical specifically in real-time applications. In the inference time, multiple models need to be decoded which increases the application latency and consumes large memory making using speech attributes unfeasible for applications that need instant feedback to be provided to the user.

In this work, we proposed a multi-label CTC method to train a single model that can detect multiple non-mutually exclusive speech attributes. We also investigate using wav2vec2 speech representations to perform speech attribute detection as a downstream task.

3. Method

Two MDD systems were implemented, a conventional phoneme-based MDD similar to the wav2vec2-based model introduced in [25] and the proposed speech attribute-based, with the former used as a baseline to compare the performance of our proposed MDD. A block diagram describing the two methods is depicted in Figure 1. In the phoneme-based method, the raw speech signal was processed by the phonetic acoustic model to generate a sequence of recognised phonemes. The resultant phoneme sequence was then aligned with the reference (canonical) phoneme sequence to identify and diagnose the pronounced errors at the phoneme level. In the second speech attribute-based method, the raw speech signal was processed by the speech attributes model to generate multiple sequences of $+att/ - att$ representing the existence or absence of each attribute. The reference phoneme sequence was obtained by first mapping to multiple reference speech attribute sequences of $+att/ - att$. For each attribute, the reference and recognised sequences were then aligned to identify the missing, inserted, or replaced attributes at each position.

As shown in Figure 2, the wav2vec2 architecture was used as the backbone of both the speech attribute and the phonetic acoustic models while the Connectionist Temporal Classification (CTC) criterion was adopted to compute the training loss.

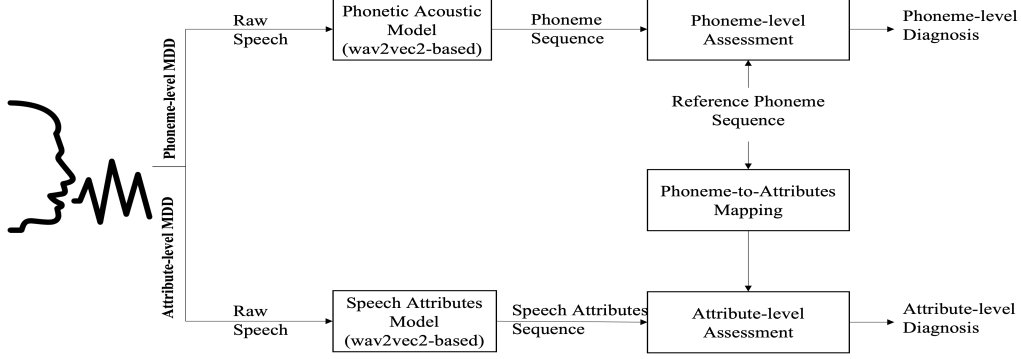


Figure 1: A block diagram of the phoneme and attribute level Mispronunciation Detection and Diagnosis (MDD) approaches implemented. The wav2vec2-based phonetic acoustic model processes the raw speech signal and outputs a sequence of recognised phonemes. The wav2vec-based speech attributes model processes the raw speech signal and produces multiple speech attribute sequences, one for each targeted attribute. The assessment was performed at the phoneme level by aligning the recognised phoneme sequence with the reference one, and at the speech attributes level by aligning each recognised attribute sequence of $+att/-att$ with the corresponding reference attribute sequence.

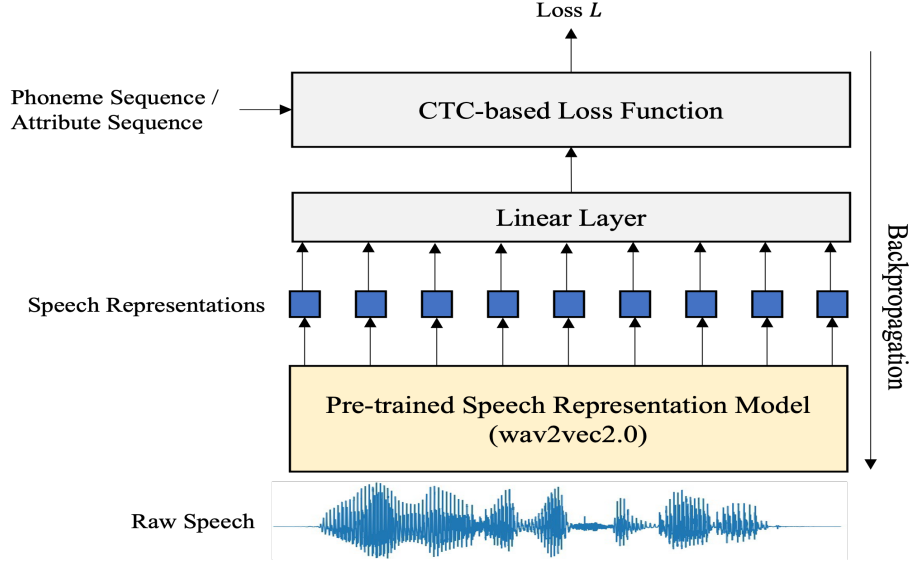
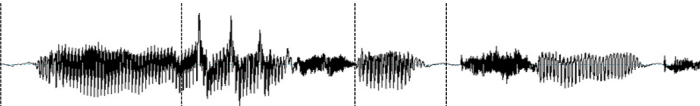


Figure 2: The architecture of the phoneme/attributes model. The speech representations generated by the pre-trained wav2vec2 model were passed to a linear layer with number of nodes equals to the number of target phonemes/attributes. The model was fed by the raw speech signal and the CTC-based loss function was calculated by considering all possible alignments between the input signal and the target phoneme/attribute sequence.

Figure 3 demonstrates an example of MDD of three substitution errors detected in the production of the sentence “There was a change” by an L2 speaker. Here the consonant sound /r/ was pronounced as /ah/ vowel while the voiced /z/ and /jh/ sounds were replaced with the voiceless /s/ and /ch/. As seen, the speech attribute detection model provides a detailed description of the pronunciation errors in terms of the manner and places of articulation.



		There			was			a	change			
Phoneme-level	Expected	dh	eh	r	w	ah	z	ah	ch	ey	n	jh
	Pronounced	dh	eh	ah	w	ah	s	ah	ch	ey	n	ch
Speech Attributes-level	Vowel	Expected	-	+	-	+	-	+	-	+	-	-
		Recognized	-	+	+	-	+	-	-	+	-	-
	Voiced	Expected	+	+	+	+	+	+	-	+	+	+
		Recognized	+	+	+	+	-	+	-	+	+	-
	Liquid	Expected	-	-	+	-	-	-	-	-	-	-
		Recognized	-	-	-	-	-	-	-	-	-	-

Figure 3: Mispronunciation detection and diagnosis example of sentence “*There was a change*”. /r/, /z/, and /jh/ phonemes are pronounced as /ah/, /s/, and /ch/ respectively. /r/ is *-vowel* and *+liquid* while /ah/ is *+vowel* and *-liquid*. Both /z/ and /jh/ are *+voiced* while their associated erroneous phonemes /s/ and /ch/ are *-voiced*. The proposed MDD breaks down the pronunciation error into elementary components (attributes) and identifies the incorrect attribute.

3.1. Connectionist Temporal Classification (CTC)

The Connectionist Temporal Classification (CTC) loss function, initially proposed for speech recognition [62], enables learning of a Seq2Seq model without explicit alignment between the input and target sequences. Therefore, it has become the state-of-the-art learning criterion for several Seq2Seq models in various domains, including computer vision [63], handwritten text recognition [64], and speech recognition [65]. The goal of the CTC algorithm is to compute $p(l|x)$ where $x = (x^1, x^2, x^3, \dots, x^T)$ is an input sequence

of length T and $l = (l^1, l^2, l^3, \dots, l^U)$ of length U its corresponding target sequence, with $U \leq T$.

L defines a finite alphabet containing all possible labels where $l^t \in L$. CTC also defines a *blank* label which is used when no output is assigned to a specific time slot and to differentiate time slots that belong to the same label from time slots that belong to a repeated label. Hence, let $L' = L \cup \text{blank}$. The output tensor y of the network is therefore of size $T \times |L'|$ where $|L'|$ is the total number of possible labels in addition to the *blank* output. The probability of each element $i \in L'$ at time t is denoted as y_i^t , with the softmax operation performed over y^t to give $\sum_{i=1}^{|L'|} y_i^t = 1$.

If π is a label sequence of length equal to T and labels $\pi^t \in L'$, the probability of producing an output sequence π is computed as in (1)

$$p(\pi|x) = \prod_{t=1}^T y_{\pi^t}^t \quad (1)$$

CTC also defines a many-to-one mapping β that maps multiple label sequences of length T and labels $\pi^t \in L'$ to one label sequence l of length $U \leq T$ and labels $l^t \in L$. The mapping is performed by removing *blank* and repeated labels. For example, $\beta(a - ab -) = \beta(aa - ab) = \beta(-a - ab) = aab$ where $-$ denotes *blank*. Therefore, the probability of a label sequence l given an input sequence x can be computed as the sum of all paths mapped to l .

$$p(l|x) = \sum_{\pi \in \beta^{-1}(l)} p(\pi|x) \quad (2)$$

3.2. Connectionist Temporal Classification for Multi-label Sequences

As the speech attribute features are non-mutually exclusive, where each phoneme is described by more than one attribute, the speech attribute detection becomes a multi-label classification problem. Each speech utterance can thus be mapped to different attribute sequences. The traditional CTC method can only handle single-label inputs and therefore a multi-label variant of the CTC is needed to jointly model all speech attributes.

Two approaches were introduced to handle multi-label classification problem using CTC criteria, the Separable CTC (SCTC) and the Multi-label CTC (MCTC) [66]. The SCTC works by computing the CTC loss over each labelling category separately and then adding them together (i.e. multiplying the conditional probabilities) to get the target loss. In [67], the authors used

the SCTC criteria to train the CTC-based model to recognise both phones and tones in multilingual and cross-lingual scenarios.

On the other hand, the MCTC first computes the probability of each target element from its categorical components and then computes the CTC loss over the target sequence [66]. The MCTC approach has been successfully applied to polyphonic music audio to detect pitch classes [68] and to handwritten text recognition [66].

Wigington et. al [66] discuss the two approaches for the multi-label CTC. One major issue with SCTC is that each category is treated separately and therefore it is not guaranteed that at each frame the correct combination of components exists.

In this work, we adopted the SCTC approach as the objective is the classification of the separate speech attribute categories. However, to maintain the alignment between components, one blank node was shared among all categories. The reasoning behind using a shared blank token is that in speech attribute modelling, as the categories are binary representations of each attribute, each phoneme has one and only one representation in each category. The blank token over all categories thus represents the silence frames or the transition from one phoneme to another. Using a shared blank, increases the likelihood that all categories will produce a blank at the same time frame. We refer to this approach as Separable CTC with Shared Blank (SCTC-SB).

In the multi-label scenario, each input sequence has multiple target sequences with labels derived from different alphabet categories. N categories $C = \{C_1, C_2, C_3, \dots, C_N\}$ are then defined, where C_i represents the alphabet of category i . For any input x with target sequence l , each element in l is then decomposed into N components representing the N categories where:

$$l^t = (l_1^t, l_2^t, l_3^t, \dots, l_N^t) \quad (3)$$

Therefore, the input x of length T and target sequence l of length $U \leq T$ can be represented in N sequences of l_i all of length U .

The network output tensor y will be of size $T \times (\sum_{i=1}^N |C_i| + 1)$ where $|C_i|$ is the number of elements in category C_i plus the shared blank node. Let $C'_i = C_i \cup \text{blank}$, then the probability of output element j of category C'_i at time t is denoted as $y_{i,j}^t$. In this case, the softmax function is applied over the components of C'_i , hence $\sum_{j=1}^{|C'_i|} y_{i,j}^t = 1$.

The probability of each category is then computed separately as:

$$p(l_i|x) = \sum_{\pi_i \in \beta_i^{-1}(l_i)} p(\pi_i|x) \quad (4)$$

Where π_i is a label sequence of length T with $\pi_i^t \in C'_i$ while β_i is a many-to-one mapping defined for each category i that maps π_i to l_i , and the probability of the output path π_i is computed as:

$$p(\pi_i|x) = \prod_{t=1}^T y_{i,\pi_i^t}^t \quad (5)$$

The final objective function is then computed as the multiplication of the label probabilities of all categories:

$$p(l|x) = \prod_{i=1}^N p(l_i|x) \quad (6)$$

3.3. Speech Attribute Modeling

$N = 35$ speech attributes were adopted representing the manners and places of articulation along with other phonological features as listed in Table 1. These attributes were selected so that each phoneme has a unique binary representation in terms of the 35 attributes. The 35 speech attributes were jointly learnt using the proposed SCTC-SB criterion. A category for each attribute (*att*) was defined with items $C_i = +att, -att$. The category that represents the nasal attribute, for example, has possible outputs of $+nasal, -nasal$. Therefore, the number of network outputs is equal to 71, where 35 nodes represent the existence of each attribute ($+att$), 35 nodes represent the absence of each attribute ($-att$), and one node represents the blank output that is shared among all categories.

In the training phase, the phoneme sequence l of each utterance was mapped to 35 target sequences l_i of $+att, -att$ symbols, one for each attribute as shown in Figure 4 for an example phrase "How old are you?". The CTC loss of each category was then computed using (5) and the final loss function for the utterance was calculated according to (6) by multiplying the 35 CTC losses of all categories producing the phoneme sequence loss.

In the inference phase, for any speech input x of length T , the most probable labelling of each attribute i was obtained by applying an arg max

Table 1: The 35 learned speech attributes including manners and places of articulation and other phonological features such as voiceness [69]

Manners	Places	Others
Consonant, sonorant, fricative, nasal, stop, approximant, affricate, liquid, vowel, semivowel, continuant	Alveolar, dental, velar, front, anterior, retroflex, coronal, palatal, glottal, labial, mid, high, low, back, central, posterior, bilabial, coronal, dorsal	Long, short, monophthong, diphthong, round, voiced

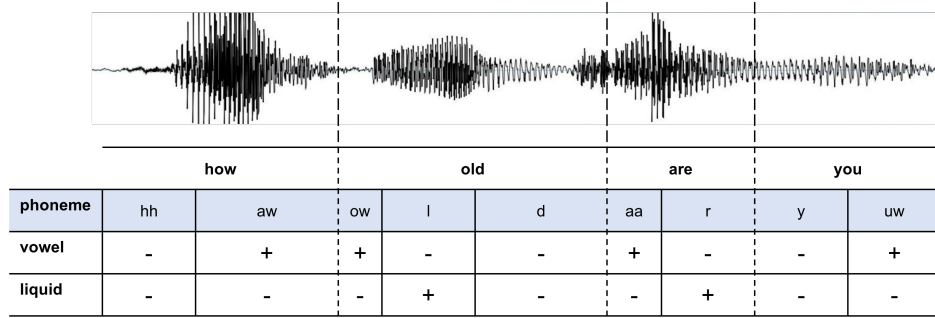


Figure 4: Example attribute mapping for the phrase ‘How old are you?’ showing the ground truth phoneme sequence and the corresponding vowel and liquid attributes labels. Here the vowel attribute is present for /aw/, /ow/, /aa/, and /uw/, while the liquid is present for /l/ and /r/ consonants.

function over the output nodes representing items in each category C'_i as follows

$$h_i(x) = \operatorname{argmax}_j y_{i,j}^t, \quad j = 1, 2, \dots, |C'_i| \quad (7)$$

The output is a sequence of $+att$, $-att$ and blank tokens of length T . Finally, the repeated tokens were merged, and all blank tokens were removed to get the final output sequence of length $U \leq T$ using predefined category mapping β_i .

4. Speech Corpora

Three English speech corpora were employed in this paper. The native Librispeech (LS) [70] corpus, which consists of audiobooks from the Lib-

rivox project [71], the native TIMIT dataset [72], and the non-native second language learning L2-ARCTIC corpus [73]. The LS corpus is a frequently utilized open-access database for evaluating diverse speech-processing tasks. On the other hand, the TIMIT dataset is a well-designed, small dataset that includes phonetically balanced recordings from various dialects. Additionally, the non-native L2-ARCTIC corpus is a relatively recent speech corpus that has become the standard dataset for evaluating MDD systems.

The LS dataset has “clean” and “other” versions, referred to as LS-clean and LS-other, respectively. The “other” portion of the LS dataset, LS-other, contains challenging speech data that has lower recognition accuracy compared to the “clean” portion [70].

On the other hand, L2-ARCTIC was collected from 24 non-native English speakers equally distributed over 6 native languages, namely Arabic, Hindi, Korean, Mandarin, Spanish, and Vietnamese. The L2-ARCTIC contains two subsets, scripted and spontaneous, referred to as L2-Scripted and L2-Suitcase respectively. The scripted one consists of 27 hours of speech, of which only 3.5 hours were manually annotated at the phoneme level. In contrast, the spontaneous speech subset contains 26 minutes of speech where each speaker recorded around one minute. Unlike the scripted speech, the whole spontaneous speech utterances were manually annotated at the phoneme level. The same speakers recorded both scripted and spontaneous subsets. The L2-ARCTIC was further split into training and testing subsets following the same split as in [27], where six speakers, one from each L1 language, formed the test set while the other 18 speakers were used for training.

The TIMIT corpus, 100 hours of the LS-clean, referred to as LS-clean-100, and a combination of TIMIT and L2-ARCTIC corpora, referred to as TIMIT+L2, were used separately for training, developing and testing of the speech attribute detection and the phoneme recognition acoustic models. While the L2-ARCTIC was for the evaluation of the MDD system.

Around 15% of the L2-ARCTIC has manual annotations at the phoneme level describing each pronunciation error by indicating the pronounced phoneme, if pronounced, and if the error is an insertion (I), deletion (D), or substitution (S). Table 3 shows the number of each type of error, along with the number of correctly pronounced phonemes (C), in the training and testing subsets of the L2-Scripted and L2-Suitcase datasets.

The CMU dictionary [74] was used to obtain the phonetic transcription of the LS orthographic transcription. To match the phoneme set of LS and L2-ARCTIC, the TIMIT phonemes were converted from 61 to 39 following

Table 2: Distribution of the utilised datasets.

	Name	Type	Hours	Speakers
Training	TIMIT	Native	3.9	462
	LS-clean-100	Native	100	251
	TIMIT+L2	Native+Non Native	6.5	480
Testing	LS-clean-test	Native	4.1	40
	LS-other-test	Native	4.4	33
	TIMIT-test	Native	1.4	168
	L2-Scripted	Non-Native	0.11	6
	L2-Suitcase	Non-Native	0.87	6

Table 3: The number of Correct (C), Substituted (S), Inserted (I), and Deleted (D) phonemes in the L2ARCTIC speech corpus.

	Subset	C	S	I	D
L2-Scripted	Training	79864	10474	772	2437
	Testing	28331	3198	214	939
L2-Suitcase	Training	5736	1032	55	290
	Testing	2333	438	28	137

the mapping proposed in [66]. However, /zh/ was not mapped to /sh/ as the former is *+voiced* while the latter is *−voiced* and merging them can confuse the *voiced* attribute model. All silence labels were further removed leaving the silence frames to be handled by the blank label.

5. Experimental Settings

5.1. Experiments

First, we performed a number of experiments to assess the performance of our proposed speech attribute detection method followed by a comparison

between the phonological (speech attribute)-level MDD and the phoneme-based MDD. We conducted the following experiments:

1. To investigate the influence of the pre-trained model size and the domain of the speech corpora used in the pre-training process, we compared the performance of the three wav2vec2-based pre-trained models with respect to speech attribute recognition: wav2vec2-base, wav2vec2-large, and wav2vec2-large-robust.
2. To explore the robustness of the the proposed speech attribute recognition approach, we compared its performance when applied to different domains (LS and TIMIT) as well as when used with out-of-domain data (native/non-native).
3. To demonstrate the effectiveness of the proposed wav2vec2-based speech attribute detection model, we compared its performance to an existing baseline based on the Deep Speech 2.0 [75] model.
4. We finally assessed the effectiveness of our speech attribute recognition approach when used for Mispronunciation Detection and Diagnosis (MDD) by applying the method to L2 speech and comparing it to phoneme-level MDD.

5.2. Training Procedures

As aforementioned, the core model of the architecture is the self-supervised pre-trained wav2vec2 model. Figure 5 depicts the block diagram of the training procedure of our proposed model. The wav2vec2 model consists of a multi-layer CNN encoder that generates the latent representations followed by a transformer module with multiple blocks and multiple attention heads. A linear layer was added on top of the transformer module with the number of nodes representing the number of classes. The number of classes used to train the SCTC-SB-based speech attribute model was 71: 35 for the existence and 35 for the absence of each attribute plus one node for the *blank* output. For the CTC-based phoneme recognition model, 40 output nodes were used representing the 39 phonemes in addition to a *blank* node.

Except for the CNN encoder layer, the whole network was then fine-tuned to minimise either the SCTC-SB or CTC loss using backpropagation. As the feature extraction layer was already well-trained during pre-training, its parameters were fixed during the fine-tuning process. Furthermore, SpecAugment [76] was applied to the output of the CNN encoder to add more variations to the training data. The performance of three pre-trained models,

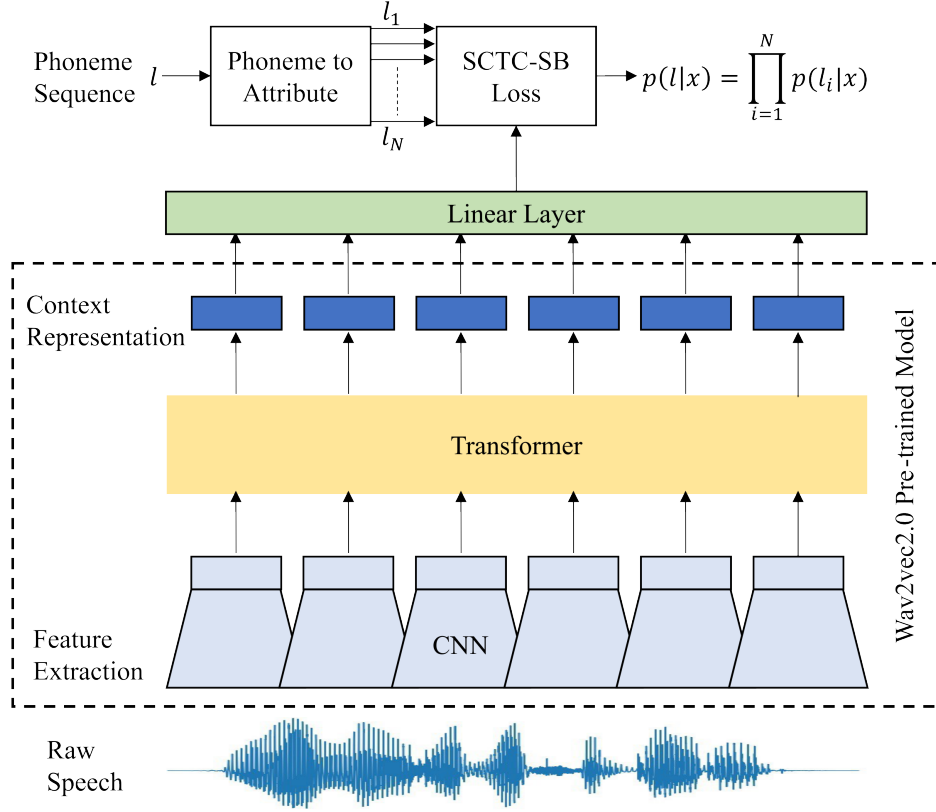


Figure 5: The training procedure for the proposed speech attribute recognition model. A linear layer was added on top of a wav2vec2-based pre-trained model. For each speech utterance, the linear layer converts the corresponding phoneme sequence to N binary sequences of speech attributes. The CTC loss was then computed for each attribute sequence. Finally, the SCTC-SB loss was computed by multiplying all speech attributes’ CTC losses.

namely wav2vec-base, wav2vec-large [25], and wav2vec-large-robust [77] were compared. The first two models were pre-trained on 960 hours of read-out books dataset, Librispeech, while the latter one was pre-trained on multiple datasets from different domains, including read-out books and telephone conversations. The wav2vec-base model has 95M parameters while both wav2vec-large and wav2vec-large-robust have 317M parameters.

AdamW optimization [78] was utilised for all experiments with a 0.005 weight decay and 0.0001 initial learning rate. The batch size was fixed to 32, and the fine-tuning ran for 30 epochs. 10% of the iteration steps were

consumed in the warmup phase to reach the initial learning rate.

5.3. Evaluation Metrics

In this work, we utilized several performance metrics to evaluate and compare the developed models based on their different tasks

5.3.1. Attribute Recognition Performance

For the attribute recognition model, as the output is a binary sequence of $+att/ -att$ symbols, we used the traditional error rate derived from the Levenshtein distance metric [79]. The Levenshtein distance metric works by measuring the difference between two sequences in terms of the number of Insertion (I), Deletion (D), and Substitution (S) edits. Therefore, the Attribute Error Rate (AER) was computed as follows:

$$AER = \frac{S + D + I}{N} \quad (8)$$

Where I , D , and S were estimated by comparing the recognized sequence to the reference sequence and N is the total number of reference symbols. In some experiments, we used Accuracy (ACC) instead of error rate which is simply computed as $1 - ErrorRate$.

Furthermore, for the speech attribute recognition task, we calculated the Precision (PRE), Recall (REC), and $F1$ score of each attribute as follows:

$$\begin{aligned} REC &= TP / (TP + FN) \\ PRE &= TP / (TP + FP) \\ F1 &= 2 \times \frac{PRE \times REC}{PRE + REC} \end{aligned} \quad (9)$$

Where TP , FP , and FN are the True-Positive, False-Positive, and False-Negative, respectively. The positive and negative here refer to the generation of the $+att$ or $-att$ symbols.

5.3.2. Mispronunciation Detection and Diagnosis (MDD) Performance

For the evaluation of the MDD task, we used metrics as proposed in [63] at the phoneme-level. Firstly, we counted the occurrences of the following:

1. True-Acceptances (TA), which represent the number of times a recognized phoneme matches a correctly pronounced phoneme.

2. True-Rejections (TR), which represent the number of times a recognized phoneme is different from a mispronounced phoneme.
3. False-Acceptances (FA), which represent the number of times a recognized phoneme matches a mispronounced phoneme.
4. False-Rejections (FR), which represent the number of times a recognized phoneme is different from a correctly pronounced phoneme.

The TR s were further split to Correctly Diagnosed (CD) when the recognized phoneme matches the phoneme that was pronounced, and Diagnosis Error (DE) otherwise. We then used these counts to estimate the False Acceptance Rate (FAR), the False Rejection Rate (FRR), and the Diagnostic Error Rate (DER) as follows:

$$\begin{aligned} FAR &= FA / (FA + TR) \\ FRR &= FR / (FR + TA) \\ DER &= DE / (CD + DE) \end{aligned} \tag{10}$$

Similar metrics were used for the attribute-level MDD. For instance, if the phoneme $/s/$ was mispronounced as $/z/$, this was considered a mispronunciation of the *voiced* attribute only and correct pronunciation for all other attributes.

6. Results and Discussion

6.1. Speech Attribute Recognition

6.1.1. Comparison of Pre-Trained *wav2vec2* Models

The goal of this experiment was to explore how the size of the model parameters and the nature of the training data impact the performance of speech attribute detection. Here, multi-label SCTC-SB speech attribute recognition models were trained by minimizing the SCTC-SB loss function in (7) for all binary groups of the 35 speech attributes listed in Table 1. At inference time, a sequence of 35 *+att/ - att* tokens was produced for each speech file representing a sequence of the existence/absence of each speech attribute. The attribute recognition accuracy (ACC), precision (PRE), recall (REC) and F1 metrics were used for performance evaluation and model comparison.

Figure 6 shows the $F1$ and ACC metrics of the 35 speech attributes as recognised by the speech attribute recognition model when fine-tuned with

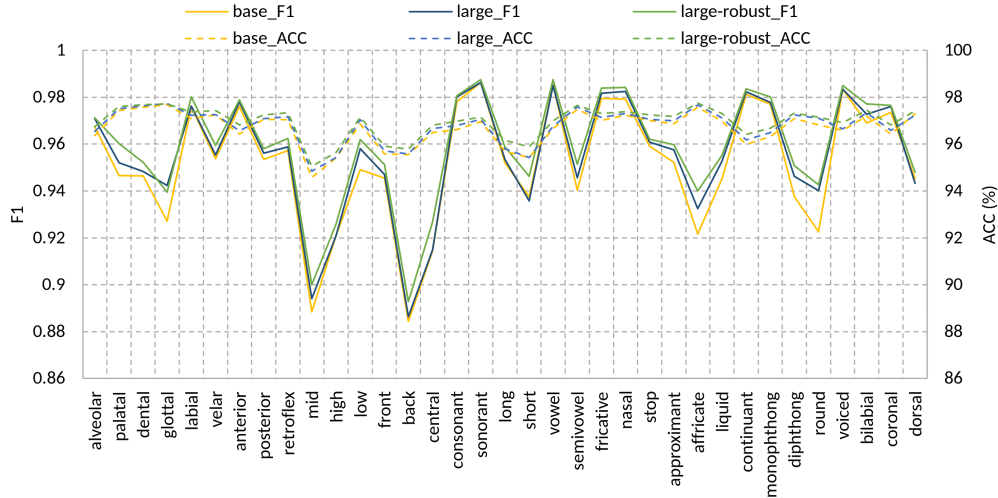


Figure 6: The performance of three wav2vec2-based pre-trained models for the speech attribute recognition task. The plots include the F1 (solid) and ACC (dashed) metrics (%) of the 35 speech attributes for the TIMIT complete test set obtained from a model fine-tuned with the TIMIT training set. The base model has the lowest performance while the large-robust performed consistently higher for all attributes.

the TIMIT training set and tested against the TIMIT test set. Here the performance of three wav2vec2-based pre-trained models, namely wav2vec2-base, wav2vec2-large, and wav2vec2-large-robust were compared.

Although both base and large models were pre-trained on the same dataset, the large model performed better than the base one. This demonstrates that increasing the model capacity improves the speech attributes recognition accuracy. On the other hand, the wav2vec2-large-robust model consistently slightly outperformed both base and large models for nearly all speech attributes. The average accuracies of the three models were 96.6 ± 0.73 , 96.8 ± 0.71 and 97 ± 0.65 for the base, large and large-robust models, respectively. The large-robust model was pre-trained on data from different domains to be more robust against out-of-domain data [77]. Similar behaviour was obtained when fine-tuning the three pre-trained models for the phoneme recognition task using the CTC criterion. The PERs of the base, large, and large-robust models were 8.1, 7.8, and 7.3 respectively.

In the subsequent experiments, the pre-trained wav2vec2-large-robust model was used in the phoneme and speech attribute recognition downstream tasks.

6.1.2. Performance of Speech Attribute Model Over Different Domains

In this experiment, we investigated the robustness of the proposed model against challenging audio, like the "other" portion of the librispeech corpus (LS-other), and in the case of domain shift from native to non-native speech. Figure 7 shows the 35 speech attribute recognition accuracies of the proposed multi-label SCTC-SB model trained on the LS-clean-100 datasets and tested against the clean and other test and development sets, LS-clean-test, LS-clean-dev, LS-other-test, and LS-other-dev. The figure shows that the accuracies of all speech attributes of the clean test and development sets are above 99%. However, applying the model to the more challenging LS-other speech datasets causes degradation across all speech attributes ranging from 31% up to 45%. The highest drop in the performance occurred in the vowel-related attributes such as *high*, *front*, *mid* and *back* while the *fricative*, *stop*, *alveolar* and *dental* were less affected attributes.

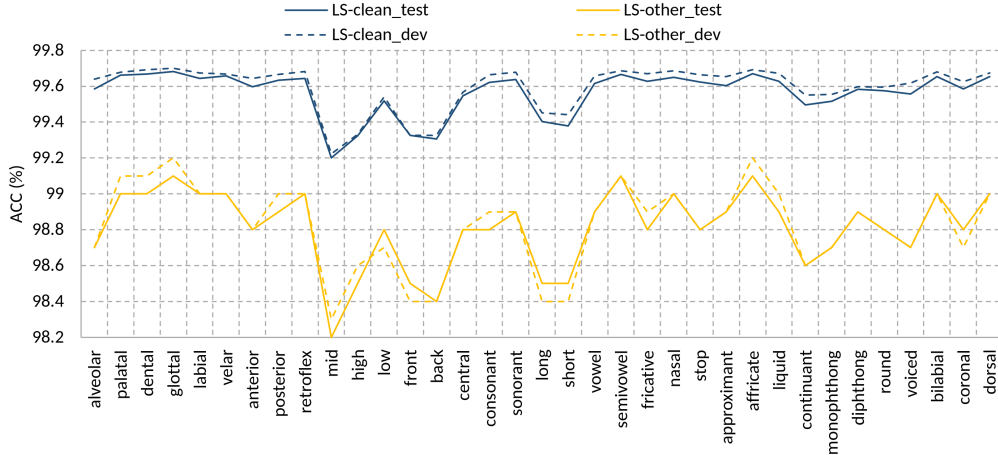


Figure 7: The 35 speech attribute recognition accuracies (%) of the multi-label SCTC-SB model trained on the LS-clean-100 datasets and tested against the “clean” and “other” test and development sets. The “other” refers to a more challenging portion of the LS corpus. When applied to LS-other data, the model performance dropped consistently across all attributes with ratios ranging from 31% - 45%. Model performance declined the most for attributes associated with vowels, including *high*, *front*, *mid*, and *back*. On the other hand, attributes like *fricative*, *stop*, *alveolar*, and *dental* were less impacted.

To demonstrate the effect of the domain-mismatch between training and testing data, two multi-label SCTC-SB models were trained on LS-clean-100 and the TIMIT training set and tested against the TIMIT test set. As

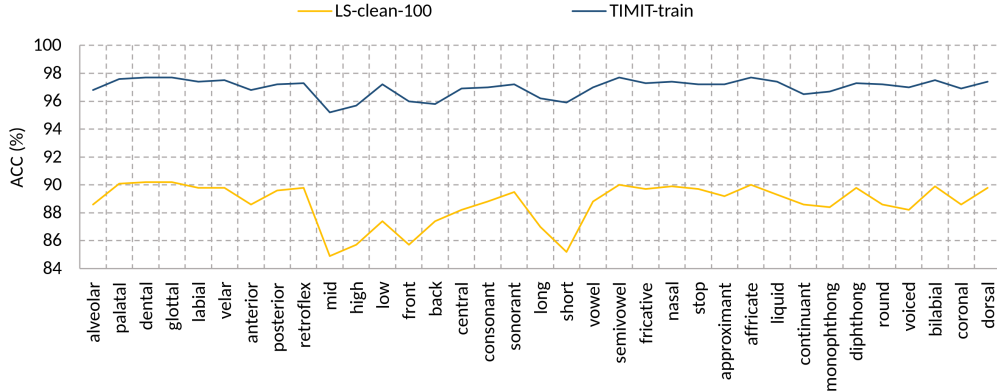


Figure 8: The effect of domain-mismatch between training and testing data. The plots show the ACC (%) of the 35 speech attributes for multi-label SCTC-SB models trained on the LS-clean-100 and TIMIT training set and tested with the TIMIT test set. Testing the model with the out-of-domain data causes a relative degradation in the ACC ranging from 22% to 33%.

shown in Figure 8, the in-domain experiment, where the model was trained and tested on TIMIT speech corpus, achieved the lowest AER of around 2.3% for dental, glottal, semivowel, and affricate attributes while the highest AER of 4.8% was obtained by the mid attribute followed by high, back and short attributes. On the other hand, the model accuracy degradation ranged from 22% to 33% when trained on out-of-domain data, LS-clean-100.

Figure 9 shows how the speech attributes recognition accuracy of models fine-tuned on native (TIMIT) and native+non-native (TIMIT-L2) datasets were impacted when tested on non-native sets, L2-Scripted and L2-Suitcase. For the L2- Scripted test set, it is evident that adding non-native data to the training set significantly improves the recognition accuracy of all speech attributes. However, for L2-Suitcase, not all attributes show the same improvement when using non-native data in the model’s fine-tuning. The glottal, retroflex, nasal and diphthong attributes achieved almost the same accuracy with and without non-native data. Moreover, accuracies of labial and semivowel degraded slightly with non-native data.

Figure 10 summarises the precision, recall, and F1 quartiles of different pairs of training and testing sets. It is noticeable that the performances of both TIMIT and L2-Scripted with a multi-label SCTC-SB model trained on LS-clean-100 were comparable despite both TIMIT and LS-clean-100 being

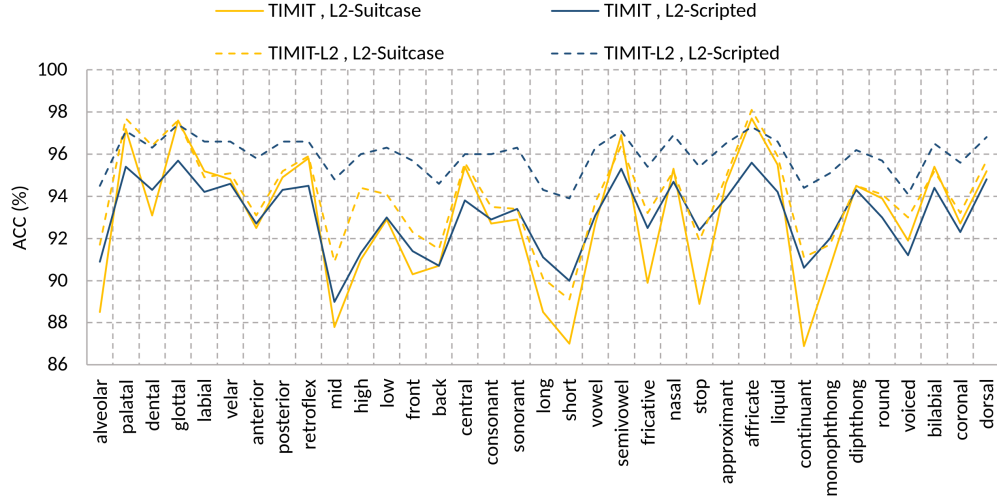


Figure 9: The accuracy of speech attributes recognition when used with non-native test sets. Here multi-label SCTC-SB models are fine-tuned with native (TIMIT) and native+non-native (TIMIT-L2) datasets and tested with non-native L2-Scripted and L2-Suitcase datasets. The L2-Scripted data achieves significant improvements when adding L2 data to the training set over all attributes.

native English data. Moreover, although both L2-suitcase and L2-Scripted were recorded by the same speakers, the L2-Suitcase seems more challenging than the L2-Scripted. This is explained by the L2-Suitcase recordings’ longer duration and spontaneous nature compared to the short and scripted L2-Scripted dataset.

6.1.3. Comparison With Existing Work

Most of the existing work on speech attribute detection reports frame-level performance. However, the proposed model is based on CTC, which outputs a sequence of symbols. As the CTC blank node allows the model to output no label in some frames when uncertain, obtaining frame-level alignment from the CTC model is not accurate. Consequently, it is hard to make a fair comparison with a frame-level attribute detection system. Therefore, the proposed SCTC-SB model was compared with two CTC-based systems. Both systems were based on the Deep Speech 2 model [75] and trained using CTC. The first system was proposed in [57] for nasality detection where the output is a sequence of +nasal,-nasal tokens. The model achieved an error rate of 4.4%. The second system, introduced in [80], discriminates be-

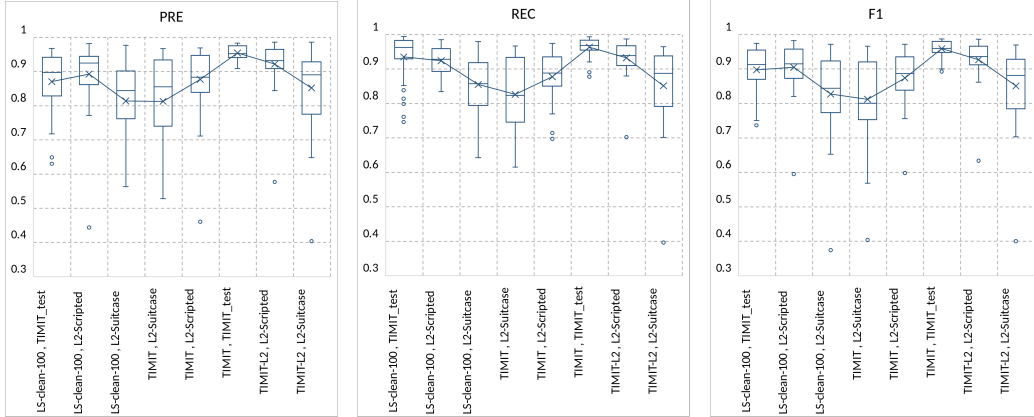


Figure 10: A comparison of the speech attribute recognition performance metrics across different pairs of training and testing datasets. Both TIMIT and LS-clean-100 are native English data, while L2-Scripted and L2-Suitcase are non-native English scripted and spontaneous speech, respectively. The boxplots show the precision (PRE), recall (REC), and F1 quartiles of the different (training, testing) data pairs. The spontaneous L2-Suitcase test set achieves lower PRE, REC, and F1 compared with the L2-Scripted test set over all training sets.

tween five manners of articulations, namely vowel, semi-vowel, nasal, stop, and fricative. The system achieved an error rate of 2.7%. Both systems were trained and tested on the LS-clean-100 and LS-clean datasets, respectively. As shown in Figure 8, when the proposed SCTC-SB system was trained and tested on the same dataset, the error rate for all attributes was less than 1%, significantly outperforming both baseline systems.

6.2. Mispronunciation Detection and Diagnosis (MDD)

The next experiments aim to demonstrate the effectiveness of the proposed speech attribute-based MDD in terms of detection and diagnosis accuracies. For comparison, we implemented a SOTA phoneme-level MDD and applied it to the same native and non-native speech corpora.

6.2.1. Phoneme-level MDD

In this experiment, MDD at the phoneme level was performed using three phonetic acoustic models trained on the LS-clean-100, TIMIT, and TIMIT-L2 datasets. Table 4 summarises the evaluation results of the different models tested with L2-Scripted and L2-Suitcase test sets. The table shows that the model fine-tuned with the LS-clean-100 dataset has the highest $FARs$ of 63%

and 57% for L2-Scripted and L2-Suitcase respectively, and an FRR of 8% for both datasets. In contrast, the model fine-tuned with the TIMIT dataset achieved a better balance between FRR and FAR . Adding L2 data to the TIMIT training dataset significantly reduced the FRR from 18% and 21%, when using only the TIMIT training set, to 7% and 12% for the L2-Scripted and L2-Suitcase test sets, respectively. Furthermore, the model fine-tuned with the TIMIT-L2 dataset achieved the lowest DER s of 15% and 17% for L2-Scripted and L2-Suitcase test sets, respectively. These results suggest that by increasing the amount of native data in the training of the phonetic acoustic model, the model’s tendency to accept pronunciation errors (FAR) is increased significantly. By comparing the TIMIT and TIMIT-L2 models, it is noticeable that although both models can detect the existence of error with a comparable FAR , adding in-domain L2 data improved the ability of the model to Correctly Diagnose (CD) the error and significantly reduce the DER .

Table 4: Evaluation results for the non-native test data using phonetic acoustic models trained on native and non-native datasets for a phoneme-level MDD task. Using a combination of native and non-native training data achieved the lowest DER over the two test sets.

Train Data	Test Data	FA	FR	TA	TR		FRR (%)	FAR (%)	DER (%)
					CD	DE			
LS-clean-100	L2-Scripted	2686	2261	23920	1077	497	8.64	63.05	31.58
	L2-Suitcase	345	256	1884	191	66	11.96	57.31	25.68
TIMIT	L2-Scripted	1649	4808	21300	1896	715	18.42	38.71	27.38
	L2-Suitcase	209	455	1649	264	129	21.63	34.72	32.82
TIMIT-L2	L2-Scripted	1683	1899	24079	2170	407	7.31	39.51	15.79
	L2-Suitcase	155	268	1829	370	77	12.78	25.75	17.23

6.2.2. Speech Attribute-level MDD

In this experiment, the three datasets, LS-clean-100, TIMIT, and TIMIT-L2 were used to train three SCTC-SB speech attributes models to detect the existence or absence of 35 speech attributes listed in Table 1. The output of each model was a sequence of $+att/-att$ for each input speech signal. The evaluation was performed for each attribute separately by aligning each output sequence with the corresponding reference sequence obtained by mapping the target phoneme sequence to the target attribute sequence. The FAR , FRR , and DER were computed as explained in section 5.3.

Figure 11 depicts boxplots summarizing the FAR , FRR , and DER of the 35 speech attributes MDD obtained using models trained on the three mentioned training sets and tested using the L2-Scripted and L2-Suitcase test sets. The results show that the FAR of most speech attributes was less than 30% when using the purely native model trained on LS-clean-100 which is significantly lower than the phoneme-level counterpart FAR s of 57% and 63% for L2-Suitcase and L2-Scripted respectively (see Table 4). The DER s of the three models tested with L2-Scripted for all speech attributes lies below 10% which is also significantly lower than the phoneme-level DER s of 31%, 27%, and 15% of the LS-clean-100, TIMIT, and TIMIT-L2 respectively.

Detailed results for each attribute obtained from the TIMIT-L2 model are shown in Figure 12 for both L2-Scripted and L2-Suitcase test sets. The speech attributes are sorted from lowest to highest FAR values. It is obvious that there are high variations in the FAR among speech attributes while the DER and FRR are more consistent. The results demonstrate also that in both L2 test sets, all the speech attributes have achieved FAR , FRR , and DER lower than the phoneme level equivalents (shown as straight lines).

To better understand the benefits and limitations of the speech attribute model in MDD, the phoneme-based and speech attribute-based models were used to detect 30 common substitution errors obtained from the manual annotation of the L2-Scripted dataset. A complete pronunciation error matrix for L2 speakers is depicted in Figure 13.

Figure 14 shows the FAR of each substitution error computed from phoneme-level assessment and speech attribute-level assessment using attributes that can discriminate between the expected and mistaken phonemes. For instance, the top-left bar graph shows the FAR of pairs of confused phonemes using the phoneme-based model and the voiced attribute as the only attribute that discriminates between each pair.

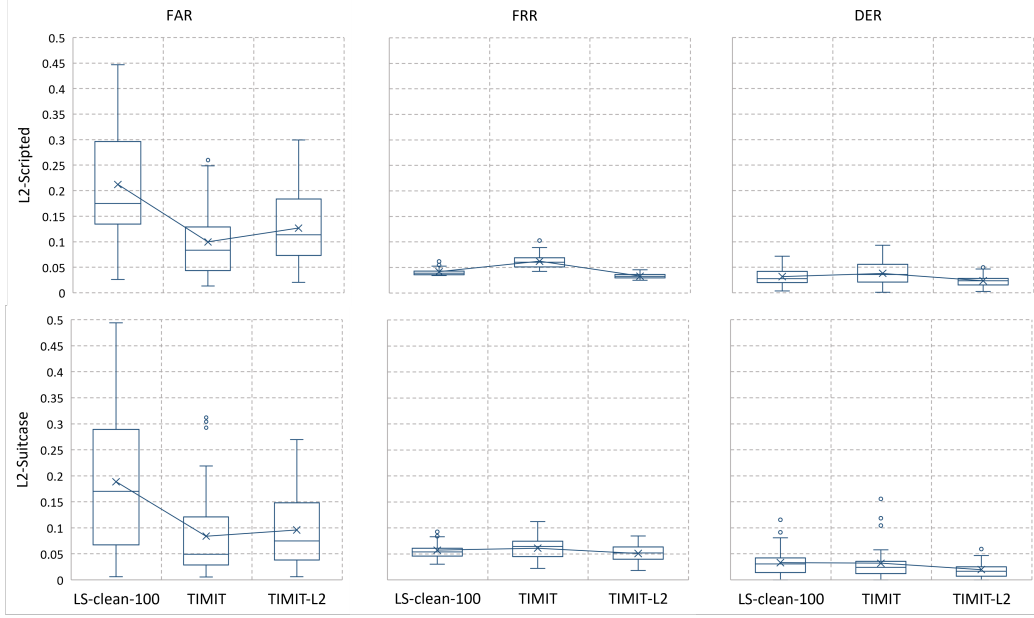


Figure 11: Evaluation of three speech attributes models in performing speech attribute level MDD of two test sets, L2-Scripted and L2-Suitcase. Two models were trained on native speech corpora, namely LS-clean-100, and TIMIT, while the third model was trained on a combination of native and non-native data TIMIT-L2.

The results show that the voiced attribute reduced the FAR of “/d/, /t/”, “/p/, /b/”, “/v/, /f/” and “/g/, /k/” compared to the phoneme-based model, however, it failed to effectively discriminate between /jh/ and /ch/ and /z/ and /s/. On the other hand, the dental attribute significantly reduced the FAR of /th/ that was mispronounced as /s/ from 72% when using the phoneme-based model to 37%. Also, the fricative attribute reduced the FAR of “/v/, /b/” substitution error by 54%. The results shown here are promising. They indicate that the speech attributes model, in addition to providing low-level diagnostic information, can also improve phoneme-level mispronunciation detection.

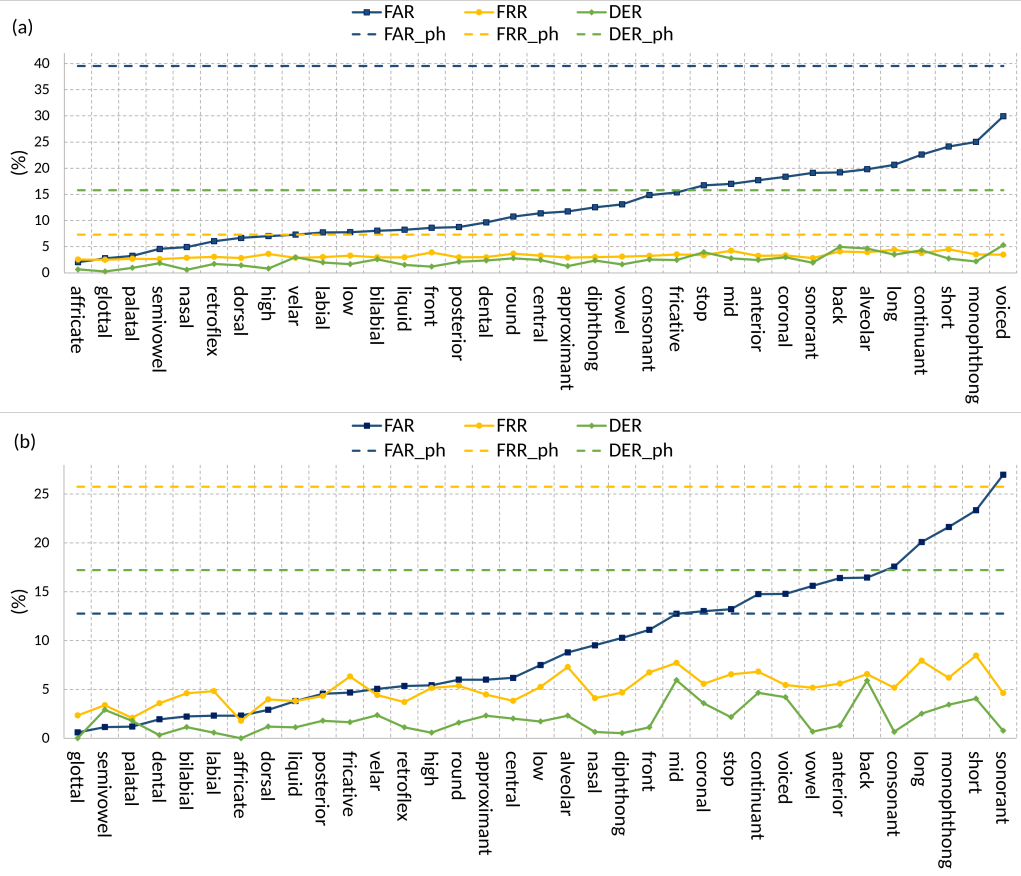


Figure 12: The speech attribute-level MDD performance of (a) L2-Scripted and (b) L2-Suitcase. The attributes are sorted from lowest to highest FAR. The straight lines FAR_ph, FRR_ph and DER_ph are the phoneme-level MDD metrics. For L2-scripted the FAR for all attributes detected by the attribute level MDD was lower than the phoneme level MDD. The DER of the scripted test set (L2-Scripted) was consistently lower than that of the spontaneous test set (L2-Suitcase).

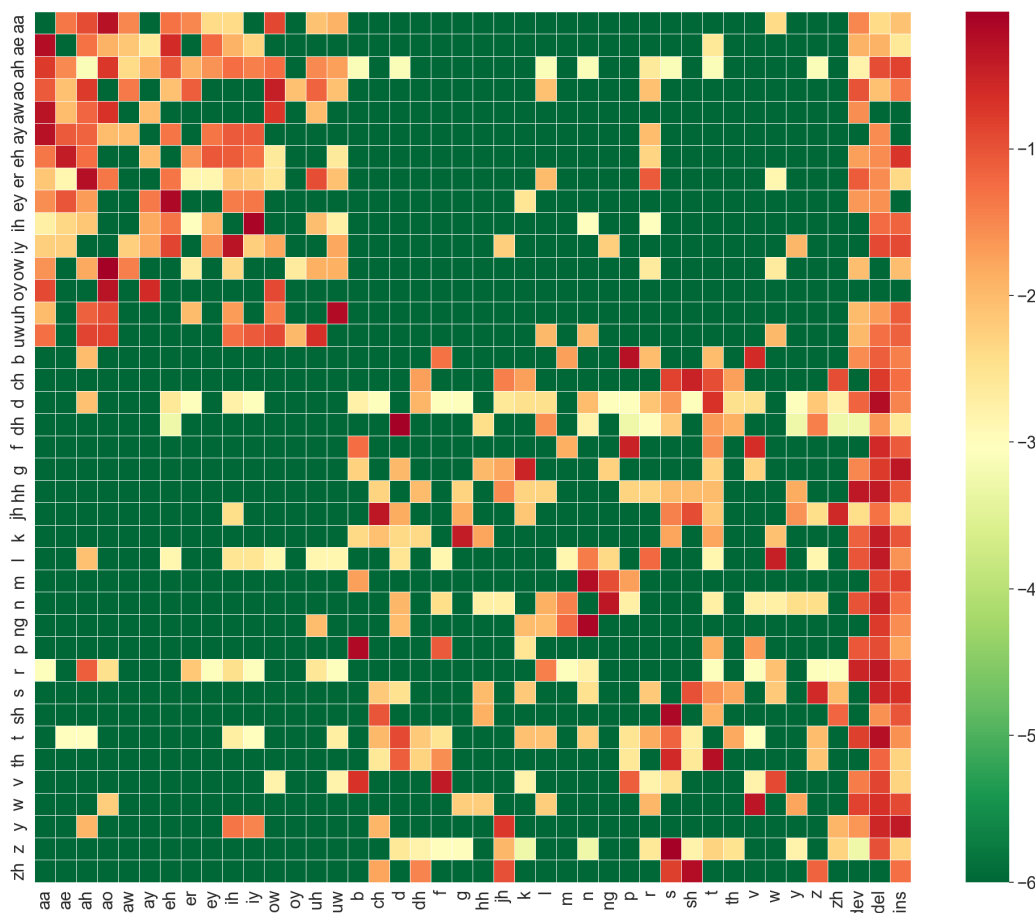


Figure 13: Pronunciation errors made by L2 speakers of L2-ARCTIC speech corpora. The vertical axis labels show the canonical phoneme, while the horizontal axis labels show the replaced phoneme, in case of substitution error, the deviation (dev), insertion (ins), or deletion (del). The colour represents the percentage of each error relative to the total number of errors of each phoneme on a log scale with red indicating the highest confusion rate and green indicating no confusion. There is a notable occurrence of both deletion and insertion errors across all sounds, particularly consonants. Additionally, within the category of consonants, there is a significant level of confusion between voiced and unvoiced sounds, such as /s/ and /z/, /p/ and /b/, as well as /d/ and /t/.



Figure 14: The FAR of substitution errors computed from phoneme-level assessment and speech attribute-level assessment using attributes that can discriminate between the expected and mistaken phonemes. The voiced attribute improves the discrimination between (/d/,/t/), (/p/,/b/), (/v/,/f/) and (/k/,/g/), while the diphthong improves the discrimination between (/ey/,/eh/), (/ow/,/ao/), and (/aw/,/aa/).

7. Conclusion

This paper introduced a novel MDD method to detect the existence of pronunciation errors and provide comprehensive diagnostic information. Unlike the current SOTA phoneme-level MDD which can only recognise the incorrect phoneme, our method, in addition to locating the mistaken phoneme, gives a detailed description of the pronunciation error in terms of which articulators are involved in the error’s production.

This was achieved by first modelling the speech attribute (phonological) features, which include the manners and places of articulation. Given the recent superiority of the pre-trained speech representation model in different downstream tasks, we adopted the wav2vec2 pre-trained model as the core architecture of our speech attribute detection model. The results show that the large wav2vec2 model, which was pre-trained using data from different domains (clean, noisy, telephone, crowd-sourcing), outperformed single-domain models (see Figure 6).

We further proposed a novel multi-label variant of the CTC loss function (SCTC-SB) to handle the non-mutually exclusive nature of speech attributes. This enables the use of a single network for the joint modelling of all speech attributes, making the model efficient in speed and memory. This is in contrast to current methods which require separate models for each attribute [58, 59] or a group of mutually exclusive attributes [60, 61]. Furthermore, the results demonstrate the superiority of the proposed speech attribute detection model over End2End Deep Speech 2-based models [75] for the detection of nasality [57] and manner of articulations [80]. (See Section 6.1.3)

The resultant speech attribute detection model was then used in performing MDD on non-native English speech corpus collected from 24 speakers of 6 different native languages. The experimental results show that our speech attribute-based MDD can achieve a reasonable performance (with *FAR* below 30% and *DER* below 10%) when trained solely on native speech corpora and applied to non-native speech (see Figure 11). On the other hand, phoneme-level MDD which was trained and tested on the same native and non-native datasets achieved *FAR* and *DER* of 57% and 31% respectively (see Table 4). This is a significant advantage of the system as one of the major impediments in developing an accurate MDD is the scarcity of the non-native training dataset.

Furthermore, adding 2 hours of annotated non-native data with 4 hours of native data achieved an average correct detection rate of 88% and an average

diagnosis accuracy of 97% in the speech attribute MDD. (see Figure 11)

Moreover, for certain frequently occurring substitution error pairs, employing a single speech attribute yields higher discrimination accuracy compared with phoneme models. For instance, utilizing the voicing attribute proves more effective in distinguishing between the sounds /d/ and /t/ than their phonetic acoustic models. In addition to a higher accuracy compared to phoneme-based MDD, speech attribute-based MDD provides a low-level description of the pronunciation errors that is directly related to the articulatory system allowing formative feedback to be constructed.

Our future work includes leveraging the universal nature of the phonological features and incorporating speech corpora from multiple languages to train the speech attribute detection model. This will add more variations to the model leading to more robust MDD system. Furthermore, we are planning to extend the MDD system to more challenging domains such as adult and child disordered speech.

References

- [1] W.-K. Leung, X. Liu, H. Meng, Cnn-rnn-ctc based end-to-end mispronunciation detection and diagnosis, in: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp. 8132–8136.
- [2] Y. Feng, G. Fu, Q. Chen, K. Chen, Sed-mdd: Towards sentence dependent end-to-end mispronunciation detection and diagnosis, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp. 3492–3496.
- [3] M. Shahin, B. Ahmed, J. McKechnie, K. Ballard, R. Gutierrez-Osuna, A comparison of gmm-hmm and dnn-hmm based pronunciation verification techniques for use in the assessment of childhood apraxia of speech, in: Fifteenth Annual Conference of the International Speech Communication Association.
- [4] B.-C. Yan, M.-C. Wu, H.-T. Hung, B. Chen, An end-to-end mispronunciation detection system for l2 english speech leveraging novel anti-phone modeling, arXiv preprint arXiv:2005.11950 (2020).
- [5] A. M. Harrison, W.-K. Lo, X.-j. Qian, H. Meng, Implementation of an extended recognition network for mispronunciation detection and diagnosis in computer-assisted pronunciation training, in: International Workshop on Speech and Language Technology in Education.
- [6] S. M. Witt, S. J. Young, Phone-level pronunciation scoring and assessment for interactive language learning, *Speech communication* 30 (2-3) (2000) 95–108.
- [7] M. Shahin, B. Ahmed, Anomaly detection based pronunciation verification approach using speech attribute features, *Speech Communication* 111 (2019) 29–43.
- [8] D. V. Smith, A. Sneddon, L. Ward, A. Duenser, J. Freyne, D. Silvera-Tawil, A. Morgan, Improving child speech disorder assessment by incorporating out-of-domain adult speech, in: INTERSPEECH, pp. 2690–2694.

- [9] F. K. Soong, W.-K. Lo, S. Nakamura, Generalized word posterior probability (gwpp) for measuring reliability of recognized words, Proc. SWIM2004 5 (2004).
- [10] J. Zheng, C. Huang, M. Chu, F. K. Soong, W.-p. Ye, Generalized segment posterior probability for automatic mandarin pronunciation evaluation, in: 2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07, Vol. 4, IEEE, pp. IV-201-IV-204.
- [11] F. Zhang, C. Huang, F. K. Soong, M. Chu, R. Wang, Automatic mispronunciation detection for mandarin, in: 2008 IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, pp. 5077-5080.
- [12] K. Yan, S. Gong, Pronunciation proficiency evaluation based on discriminatively refined acoustic models, International Journal of Information Technology and Computer Science 3 (2) (2011) 17-23.
- [13] K. C. Sim, A phone verification approach to pronunciation quality assessment for spoken language learning, in: Proceedings: APSIPA ASC 2009: Asia-Pacific Signal and Information Processing Association, 2009 Annual Summit and Conference, Asia-Pacific Signal and Information Processing Association, 2009 Annual Summit and Conference, International Organizing Committee, pp. 619-622.
- [14] [link].
URL <https://www.clapa.com/treatment/early-years-1-4/speech/>
- [15] W. Hu, Y. Qian, F. K. Soong, Y. Wang, Improved mispronunciation detection with deep neural network trained acoustic models and transfer learning based logistic regression classifiers, Speech Communication 67 (2015) 154-166.
- [16] S. M. Abdou, S. E. Hamid, M. Rashwan, A. Samir, O. Abdel-Hamid, M. Shahin, W. Nazih, Computer aided pronunciation learning system using speech recognition techniques, in: Ninth International Conference on Spoken Language Processing.
- [17] M. S. Elaraby, M. Abdallah, S. Abdou, M. Rashwan, A deep neural networks (dnn) based models for a computer aided pronunciation learning

- system, in: International Conference on Speech and Computer, Springer, pp. 51–58.
- [18] K. Truong, A. Neri, C. Cucchiarini, H. Strik, Automatic pronunciation error detection: an acoustic-phonetic approach (2004).
 - [19] H. Strik, K. P. Truong, F. d. Wet, C. Cucchiarini, Comparing classifiers for pronunciation error detection, in: Eighth Annual Conference of the International Speech Communication Association.
 - [20] H. Franco, L. Ferrer, H. Bratt, Adaptive and discriminative modeling for improved mispronunciation detection, in: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp. 7709–7713.
 - [21] J. v. Doremalen, C. Cucchiarini, H. Strik, Automatic detection of vowel pronunciation errors using multiple information sources, in: 2009 IEEE Workshop on Automatic Speech Recognition & Understanding, pp. 580–585. doi:10.1109/ASRU.2009.5373335.
 - [22] T. Zhao, A. Hoshino, M. Suzuki, N. Minematsu, K. Hirose, Automatic chinese pronunciation error detection using svm trained with structural features, in: 2012 IEEE Spoken Language Technology Workshop (SLT), IEEE, pp. 473–478.
 - [23] S. Wei, G. Hu, Y. Hu, R.-H. Wang, A new method for mispronunciation detection using support vector machine based on pronunciation space models, *Speech Communication* 51 (10) (2009) 896–905.
 - [24] W. Hu, Y. Qian, F. K. Soong, A new neural network based logistic regression classifier for improving mispronunciation detection of 12 language learners, in: The 9th International Symposium on Chinese Spoken Language Processing, IEEE, pp. 245–249.
 - [25] A. Baevski, Y. Zhou, A. Mohamed, M. Auli, wav2vec 2.0: A framework for self-supervised learning of speech representations, *Advances in Neural Information Processing Systems* 33 (2020) 12449–12460.
 - [26] X. Xu, Y. Kang, S. Cao, B. Lin, L. Ma, Explore wav2vec 2.0 for mispronunciation detection, in: Interspeech, pp. 4428–4432.

- [27] W. Ye, S. Mao, F. Soong, W. Wu, Y. Xia, J. Tien, Z. Wu, An approach to mispronunciation detection and diagnosis with acoustic, phonetic and linguistic (apl) embeddings, in: ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp. 6827–6831.
- [28] K. Li, X. Qian, H. Meng, Mispronunciation detection and diagnosis in l2 english speech using multidistribution deep neural networks, IEEE/ACM Transactions on Audio, Speech, and Language Processing 25 (1) (2016) 193–207.
- [29] K. Fu, J. Lin, D. Ke, Y. Xie, J. Zhang, B. Lin, A full text-dependent end to end mispronunciation detection and diagnosis with easy data augmentation techniques, arXiv preprint arXiv:2104.08428 (2021).
- [30] K. Li, X. Qian, S. Kang, P. Liu, H. Meng, Integrating acoustic and state-transition models for free phone recognition in l2 english speech using multi-distribution deep neural networks, in: SLaTE, pp. 119–124.
- [31] S. M. Siniscalchi, J. Reed, T. Svendsen, C.-H. Lee, Universal attribute characterization of spoken languages for automatic spoken language recognition, Computer Speech & Language 27 (1) (2013) 209–227.
- [32] O. Saz, S.-C. Yin, E. Lleida, R. Rose, C. Vaquero, W. R. Rodríguez, Tools and technologies for computer-aided speech and language therapy, Speech Communication 51 (10) (2009) 948–967.
- [33] L. He, J. Zhang, Q. Liu, H. Yin, M. Lech, Y. Huang, Automatic evaluation of hypernasality based on a cleft palate speech database, Journal of medical systems 39 (5) (2015) 1–7.
- [34] T. Pellegrini, L. Fontan, J. Maclair, J. Farinas, M. Robert, The goodness of pronunciation algorithm applied to disordered speech, in: Fifteenth Annual Conference of the International Speech Communication Association.
- [35] B. Mak, M. Siu, M. Ng, Y.-C. Tam, Y.-C. Chan, K.-W. Chan, K.-Y. Leung, S. Ho, J. Wong, J. Lo, Plaser: Pronunciation learning via automatic speech recognition, in: Proceedings of the HLT-NAACL 03 workshop on Building educational applications using natural language processing, pp. 23–29.

- [36] A. A. Hindi, M. Alsulaiman, G. Muhammad, S. Al-Kahtani, Automatic pronunciation error detection of nonnative arabic speech, in: 2014 IEEE/ACS 11th International Conference on Computer Systems and Applications (AICCSA), pp. 190–197. doi:10.1109/AICCSA.2014.7073198.
- [37] F. de Wet, C. Van der Walt, T. R. Niesler, Automatic assessment of oral language proficiency and listening comprehension, *Speech Communication* 51 (10) (2009) 864–874. doi:<https://doi.org/10.1016/j.specom.2009.03.002>.
URL <https://www.sciencedirect.com/science/article/pii/S016763930900034X>
- [38] D. Luo, N. Minematsu, Y. Yamauchi, K. Hirose, Analysis and comparison of automatic language proficiency assessment between shadowed sentences and read sentences, in: *International Workshop on Speech and Language Technology in Education*.
- [39] S. Sudhakara, M. K. Ramanathi, C. Yarra, P. K. Ghosh, An improved goodness of pronunciation (gop) measure for pronunciation evaluation with dnn-hmm system considering hmm transition probabilities.
- [40] W.-K. Lo, S. Zhang, H. Meng, Automatic derivation of phonological rules for mispronunciation detection in a computer-assisted pronunciation training system, in: *Eleventh Annual Conference of the International Speech Communication Association*.
- [41] X. Qian, H. Meng, F. Soong, Capturing l2 segmental mispronunciations with joint-sequence models in computer-aided pronunciation training (capt), in: *2010 7th International Symposium on Chinese Spoken Language Processing, IEEE*, pp. 84–88.
- [42] X. Qian, H. Meng, F. Soong, Discriminatively trained acoustic model for improving mispronunciation detection and diagnosis in computer-aided pronunciation training (capt), in: *Proc. Interspeech, Citeseer*.
- [43] K. Kyriakopoulos, K. Knill, M. Gales, Automatic detection of accent and lexical pronunciation errors in spontaneous non-native english speech, *ISCA*.

- [44] Y.-B. Wang, L.-S. Lee, Improved approaches of modeling and detecting error patterns with empirical analysis for computer-aided pronunciation training, in: 2012 IEEE international conference on acoustics, speech and signal processing (ICASSP), IEEE, pp. 5049–5052.
- [45] S.-W. F. Jiang, B.-C. Yan, T.-H. Lo, F.-A. Chao, B. Chen, Towards robust mispronunciation detection and diagnosis for l2 english learners with accent-modulating methods, in: 2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), IEEE, pp. 1065–1070.
- [46] T.-H. Lo, S.-Y. Weng, H.-J. Chang, B. Chen, An effective end-to-end modeling approach for mispronunciation detection, arXiv preprint arXiv:2005.08440 (2020).
- [47] B.-C. Yan, B. Chen, End-to-end mispronunciation detection and diagnosis from raw waveforms, in: 2021 29th European Signal Processing Conference (EUSIPCO), IEEE, pp. 61–65.
- [48] M. Ravanelli, Y. Bengio, Interpretable convolutional filters with sincnet, arXiv preprint arXiv:1811.09725 (2018).
- [49] L. Peng, K. Fu, B. Lin, D. Ke, J. Zhang, A study on fine-tuning wav2vec2. 0 model for the task of mispronunciation detection and diagnosis, in: Interspeech, pp. 4448–4452.
- [50] I.-F. Chen, S. M. Siniscalchi, C.-H. Lee, Attribute based lattice rescoring in spontaneous speech recognition, in: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2014, pp. 3325–3329.
- [51] C.-H. Lee, S. M. Siniscalchi, An information-extraction approach to speech processing: Analysis, detection, verification, and recognition, *Proceedings of the IEEE* 101 (5) (2013) 1089–1115.
- [52] J. Qi, W. Guo, J. Shi, Y. Chen, T. Liu, Exploring universal speech attributes for speaker verification with an improved cross-stitch network, arXiv preprint arXiv:2010.06248 (2020).
- [53] W. Li, K. Li, S. M. Siniscalchi, N. F. Chen, C.-H. Lee, Detecting mispronunciations of l2 learners and providing corrective feedback using

- knowledge-guided and data-driven decision trees, in: *Interspeech*, pp. 3127–3131.
- [54] R. Pradeep, K. S. Rao, Incorporation of manner of articulation constraint in lstm for speech recognition, *Circuits, Systems, and Signal Processing* 38 (2019) 3482–3500.
 - [55] I. Karaulov, D. Tkanov, Attention model for articulatory features detection, *arXiv preprint arXiv:1907.01914* (2019).
 - [56] L. Qu, C. Weber, E. Lakomkin, J. Twiefel, S. Wermter, Combining articulatory features with end-to-end learning in speech recognition, in: *Artificial Neural Networks and Machine Learning–ICANN 2018: 27th International Conference on Artificial Neural Networks, Rhodes, Greece, October 4–7, 2018, Proceedings, Part III* 27, Springer, 2018, pp. 500–510.
 - [57] M. Cernak, S. Tong, Nasal speech sounds detection using connectionist temporal classification, in: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, pp. 5574–5578.
 - [58] C.-H. Lee, M. A. Clements, S. Dusan, E. Fosler-Lussier, K. Johnson, B.-H. Juang, L. R. Rabiner, An overview on automatic speech attribute transcription (asat), in: *Eighth annual conference of the international speech communication association*.
 - [59] S. King, P. Taylor, Detection of phonological features in continuous speech using neural networks, *Computer Speech & Language* 14 (4) (2000) 333–353.
 - [60] B. Abraham, S. Umesh, N. M. Joy, Articulatory feature extraction using etc to build articulatory classifiers without forced frame alignments for speech recognition, in: *INTERSPEECH*, pp. 798–802.
 - [61] D. Merks, O. Scharenborg, Articulatory feature classification using convolutional neural networks, in: *INTERSPEECH*, pp. 2142–2146.
 - [62] A. Graves, S. Fernández, F. Gomez, J. Schmidhuber, Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks, in: *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376.

- [63] B. Shi, X. Bai, C. Yao, An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition, *IEEE transactions on pattern analysis and machine intelligence* 39 (11) (2016) 2298–2304.
- [64] D. Coquenat, C. Chatelain, T. Paquet, End-to-end handwritten paragraph text recognition using a vertical attention network, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
- [65] W. Chan, N. Jaitly, Q. Le, O. Vinyals, Listen, attend and spell: A neural network for large vocabulary conversational speech recognition, in: *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, IEEE, pp. 4960–4964.
- [66] C. Wigington, B. Price, S. Cohen, Multi-label connectionist temporal classification, *Proceedings of the International Conference on Document Analysis and Recognition, ICDAR* (2019) 979–986doi:10.1109/ICDAR.2019.00161.
- [67] J. Li, M. Hasegawa-Johnson, Autosegmental neural nets: Should phones and tones be synchronous or asynchronous?, *arXiv preprint arXiv:2007.14351* (2020).
- [68] C. Weiss, G. Peeters, Training deep pitch-class representations with a multi-label ctc loss, in: *International Society for Music Information Retrieval Conference (ISMIR)*.
- [69] N. Chomsky, M. Halle, *The sound pattern of english*. (1968).
- [70] V. Panayotov, G. Chen, D. Povey, S. Khudanpur, Librispeech: an asr corpus based on public domain audio books, in: *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, IEEE, pp. 5206–5210.
- [71] J. Kearns, Librivox: Free public domain audiobooks, *Reference Reviews* 28 (1) (2014) 7–8.
- [72] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, *Darpa timit acoustic-phonetic continuous speech corpus cd-rom*. nist speech disc 1-1.1, NASA STI/Recon technical report n 93 (1993) 27403.

- [73] G. Zhao, S. Sonsaat, A. Silpachai, I. Lucic, E. Chukharev-Hudilainen, J. Levis, R. Gutierrez-Osuna, L2-arctic: A non-native english speech corpus, in: INTERSPEECH, pp. 2783–2787.
- [74] R. L. Weide, The cmu pronouncing dictionary (1998).
URL <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>
- [75] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen, Deep speech 2: End-to-end speech recognition in english and mandarin, in: International conference on machine learning, PMLR, pp. 173–182.
- [76] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, Q. V. Le, SpecAugment: A simple data augmentation method for automatic speech recognition, arXiv preprint arXiv:1904.08779 (2019).
- [77] W.-N. Hsu, A. Sriram, A. Baevski, T. Likhomanenko, Q. Xu, V. Pratap, J. Kahn, A. Lee, R. Collobert, G. Synnaeve, Robust wav2vec 2.0: Analyzing domain shift in self-supervised pre-training, arXiv preprint arXiv:2104.01027 (2021).
- [78] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, arXiv preprint arXiv:1711.05101 (2017).
- [79] V. I. Levenshtein, et al., Binary codes capable of correcting deletions, insertions, and reversals, in: Soviet physics doklady, Vol. 10, Soviet Union, 1966, pp. 707–710.
- [80] R. Pradeep, Manner of articulation detection using connectionist temporal classification to improve automatic speech recognition performance (2018).