

HMamba: Towards Multifaceted Computer-assisted Pronunciation Training Leveraging Hierarchical Selective State Space Model and Decoupled Cross-entropy Loss

Anonymous ACL submission

Abstract

Prior efforts in building computer-assisted pronunciation training (CAPT) systems often treat automatic pronunciation assessment (APA) and mispronunciation detection and diagnosis (MDD) as separate fronts. APA aims to provide multiple pronunciation aspect scores across diverse linguistic levels, while MDD focuses instead on pinpointing the precise phonetic errors made by non-native language learners. However, a full-fledged CAPT system should integrate both features simultaneously. To address this pressing need, we in this work first propose **HMamba**, a novel hierarchical selective state space method that jointly tackles APA and MDD tasks. In addition, to enhance model performance, we introduce a novel loss function, decoupled cross-entropy loss (**deXent**), specifically tailored for the MDD task to facilitate better supervised label learning. A comprehensive set of empirical results carried out on the speechocean762 benchmark dataset demonstrate the effectiveness of our approach in multi-aspect multi-granular assessments. Furthermore, our proposed approach also yields considerable improvement in MDD performance over a competitive baseline, achieving an F1-score of 63.32%.

1 Introduction

In this era of globalization and technologization, computer-assisted pronunciation training (CAPT) systems have emerged as an appealing alternative to meet the surging demand for second language (L2) learning. In comparison with traditional curriculum learning, CAPT offers advantages in

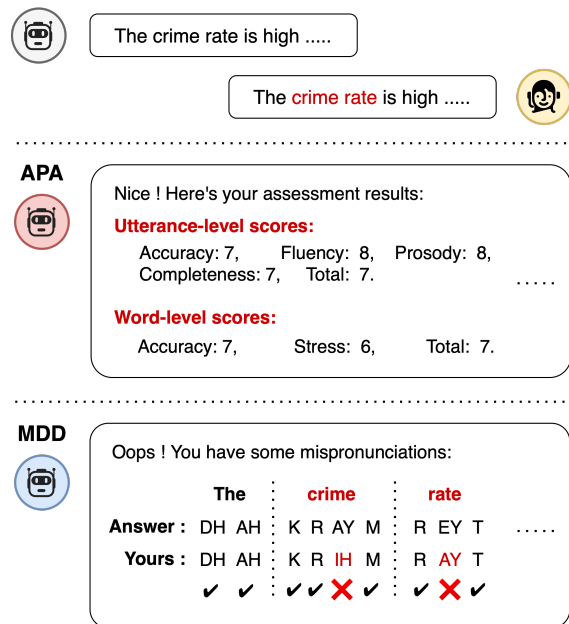


Figure 1: A running example depicts the evaluation differences between APA and MDD systems in the reading-aloud scenario.

terms of time-efficiency and cost-effectiveness. More importantly, it redefines the conventional pedagogical method from teacher-directed to self-directed learning, thereby providing a stress-free environment for L2 learners (Eskenazi et al., 2009). In addition, CAPT applications have achieved significant success in various commercial sectors or testing services, such as Duolingo (McCarthy et al., 2021) and the SpeechRater engine (Zechner et al., 2009) developed by Educational Testing Service (ETS). Typically, a de-facto archetype system for CAPT encompasses a “reading-aloud” scenario, where a non-native speaker is given a text prompt and instructed to pronounce it correctly. In this context, previous literature roughly divides applications of CAPT into two categories:

56 automatic pronunciation assessment (APA) and 108
57 mispronunciation detection and diagnosis (MDD), 109
58 with each category dedicated to specific facets of 110
59 pronunciation training. APA aims to evaluate the 111
60 L2 learners’ spoken proficiency by providing fine- 112
61 grained feedback on various aspect assessments 113
62 (e.g., accuracy, fluency) across multiple linguistic 114
63 levels (e.g., utterance level, word level) (Kheir et 115
64 al., 2023). To evaluate the extent of L2 learners’ 116
65 spoken proficiency, APA systems typically employ 117
66 scoring models that are either jointly trained (Gong 118
67 et al., 2022; Chao et al., 2022) or jointly exploit 119
68 multiple regressors (Bannò et al., 2022a; Bannò 120
69 and Matassoni, 2022b) to generate scores for each 121
70 aspect. As such, users can receive multi-aspect 122
71 assessment scores predicted by an APA system, as 123
72 illustrated in the example shown in Figure 1. 124
73 Compared with APA, MDD focuses more on non- 125
74 native speakers’ phonetic errors (Chen and Li, 126
75 2016). These errors usually have clear-cut 127
76 distinctions between correct and incorrect 128
77 pronunciations, and can be easily quantified 129
78 through deletions, substitutions, and insertions. 130
79 Therefore, MDD is often more deterministic than 131
80 APA. For instance, a number of MDD models are 132
81 capitalized on classifier-based approaches (Truong 133
82 et al., 2004; Strik et al., 2009; Harrison et al. 2009), 134
83 enabling precise identification of the exact 135
84 positions where pronunciation errors occur within 136
85 an utterance. This capability provides L2 learners 137
86 with specific feedback on discrepancies between 138
87 intended pronunciation and actual pronunciation. 139
88 Albeit the phonetic (segmental) errors are 140
89 crucial in the initial stages of non-native language 141
90 learning, prosodic (suprasegmental) errors may 142
91 often cause detrimental impact on the perception of 143
92 fluency and lead to poor intelligibility (Chen and 144
93 Li, 2016). This effect may be more pronounced in 145
94 learning stress-timed languages like English 146
95 compared with syllable-timed languages such as 147
96 Chinese (Ding and Xu, 2016). To tackle this 148
97 problem, APA can play a pivotal role by offering 149
98 prosodic assessment or intonation assessment for 150
99 L2 learners. For example, Lin et al. (2021a) 151
100 introduced rhythm rubrics to predict sentence-level 152
101 stress in L2 English, demonstrating a strong 153
102 correlation with the prosody scores assessed by the 154
103 human experts. In addition, Arias et al. (2010) 155
104 proposed text-independent systems for assessing 156
105 intonation and stress, focusing on measuring the 157
106 similarity between a student’s intonation or stress 158
107 curve and that of a reference response.

On these grounds, it is evident that both APA and MDD are indispensable ingredients of CAPT, playing complementary roles in its success. However, previous studies on APA and MDD appear to have developed independently, with limited research exploring their integration or combined use. Ryu et al. (2023) proposed a joint model for APA and MDD, leveraging knowledge transfer and multi-task learning. Their findings revealed high negative correlations between several assessment scores and mispronunciations, suggesting that the human assessors may be influenced by phonetic errors when evaluating overall proficiency scores for various aspects. While jointly modeling both tasks can achieve better performance than a single task, only utterance-level holistic scores are considered in their experiments. In order to provide more comprehensive and fine-grained feedback for L2 learners, other granularities, such as the phone or word level, should also be aptly modeled. In this paper, we propose a novel hierarchical selective state space model, dubbed HMamba, for multifaceted CAPT. Unlike previous studies that used Transformer-based structures (Gong et al., 2022; Chao et al., 2022; Do et al., 2023a), HMamba leverages Mamba (Gu and Dao, 2023), a selective state space model (SSM) approach, is capable of addressing both APA and MDD tasks simultaneously. Being aware of linguistic hierarchy, HMamba can render the intrinsic multi-layer speech structure and provide more detailed, multi-granular pronunciation assessments while offering accurate mispronunciation feedback.

The main contributions of this paper can be summarized as follows:

1. We introduce HMamba, a unified and linguistically hierarchy-aware model that jointly tackles APA and MDD tasks, achieving superior overall performance compared to prior arts that are either single-task or multi-task models.
2. We propose a novel loss function, decoupled cross-entropy loss (termed deXent), which effectively addresses the inherent issue of text prompt-aware MDD methods. Additionally, deXent is feasible and well-suited for optimizing the MDD performance, particularly in striking the balance between precision and recall.
3. To the best of our knowledge, this is the first work to adopt and extend Mamba in the APA and MDD tasks for comprehensive CAPT.

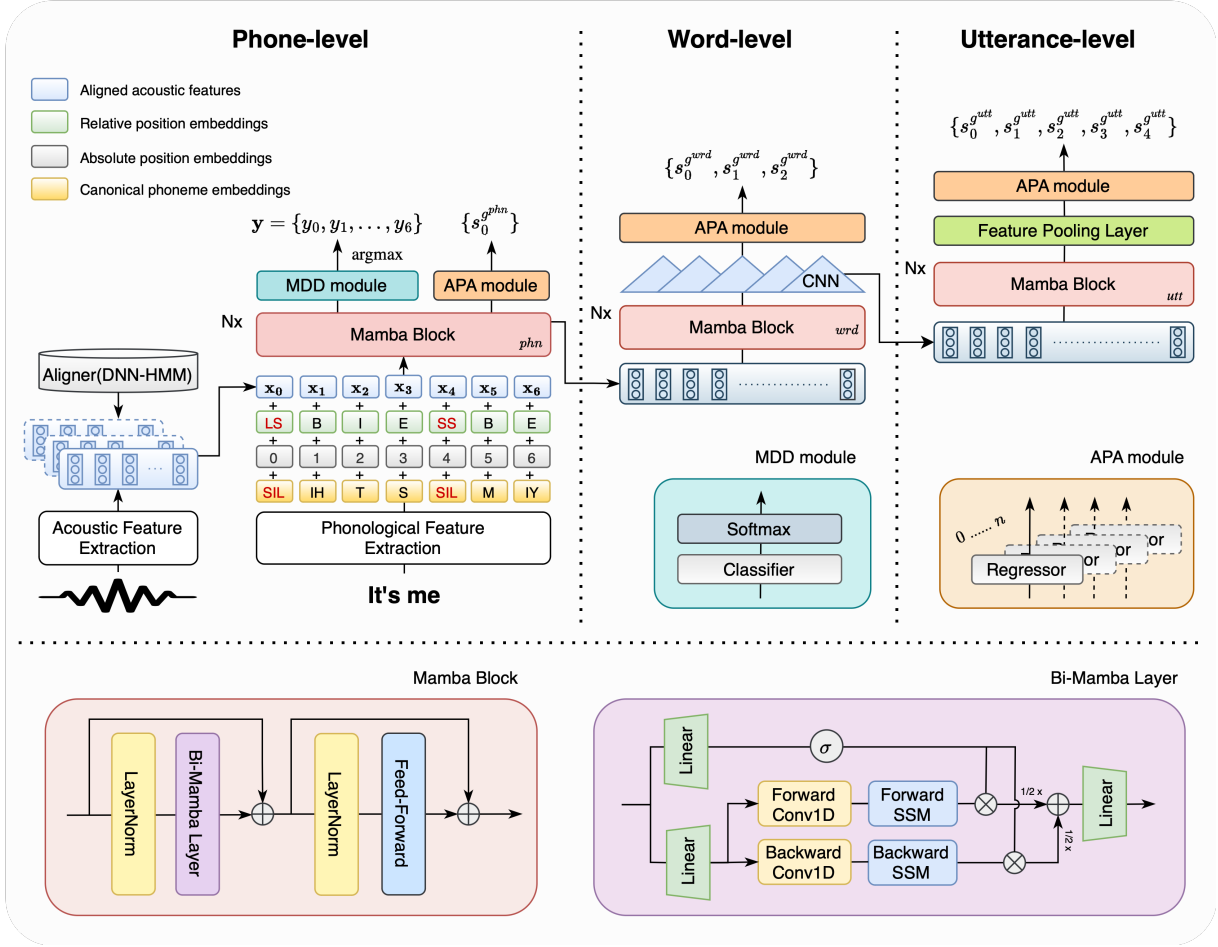


Figure 2: An overall architectural overview of HMamba, which consists of a bottom-up hierarchical modeling structure with several Mamba blocks across three levels (viz. phone, word, and utterance levels) that can perform multi-granular APA and MDD in parallel.

159 2 Methodology

160 2.1 Problem Definition

161 Considering an input time sequence of speech
 162 signal \mathbf{u} uttered by an L2 learner and a reference
 163 text prompt \mathbf{p} that contains N -length canonical
 164 phone sequence $\mathbf{p} = \{p_0, p_1, \dots, p_{N-1}\}$, we adopt
 165 a set of feature extractors along with an aligner to
 166 extract an acoustic feature sequence $\mathbf{X} =$
 167 $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{N-1}\}$ that aligned with \mathbf{p} from \mathbf{u} .
 168 Our model aims to address APA and MDD tasks
 169 simultaneously but with separate processing flows:
 170 First, we define G as a set of linguistic
 171 granularities, and for each granularity $g \in G$ the
 172 model manages to predict a set of aspect scores
 173 $\mathbf{s}^g = \{s_0^g, s_1^g, \dots, s_{M_g-1}^g\}$, where M_g refers to the
 174 number of aspect scores of target granularity g . In
 175 this work, $G = \{g^{phn}, g^{wrd}, g^{utt}\}$, where we have
 176 granularities of g^{phn} (phone level), g^{wrd} (word
 177 level), and g^{utt} (utterance level) for the APA task.

178 Meanwhile, the model also requires to detect error
 179 states $\mathbf{e} = \{e_0, e_1, \dots, e_{N-1}\}$ with respect to \mathbf{p}
 180 and in turn generate the correct diagnostic output
 181 $\mathbf{y} = \{y_0, y_1, \dots, y_{N-1}\}$, where y_n denotes the
 182 realized phone of the learner corresponds to p_n .

183 2.2 Hierarchical Selective State Space Model

184 In this subsection, we elucidate the details of the
 185 proposed model, HMamba, which is devised as a
 186 hierarchical structure built upon the paradigm of
 187 selective SSM. An overview of the complete
 188 architecture is depicted in Figure 2. Specifically,
 189 HMamba leverages the APA and MDD modules,
 190 which contain multiple regressors and a classifier,
 191 respectively. These modules collectively generate
 192 the corresponding aspect score sequence \mathbf{s}^g for
 193 each linguistic granularity g , as well as the
 194 phonetic error states \mathbf{e} and diagnosis \mathbf{y} .
 195 Furthermore, each classifier and regressor is
 196 implemented with a simple feed-forward network
 197 (FFN).

198 **Acoustic Feature Extraction:** In order to portray 241
 199 the non-native speaker’s pronunciation quality, 242
 200 previous studies on either APA or MDD generally 243
 201 adopt a pre-trained acoustic model to extract 244
 202 goodness of pronunciation (GOP)-based features 245
 203 (Witt and Young, 2000; Hu et al., 2015; Shi et al., 246
 204 2020). However, these GOP-based features merely 247
 205 offer the segmental-level information that may not 248
 206 be amenable for capturing the prosodic errors of an 249
 207 L2 learner. Given this limitation, apart from GOP, 250
 208 we utilize a pre-trained acoustic model as an 251
 209 aligner to identify phone boundaries (including 252
 210 silence), facilitating the extraction of other 253
 211 prosodic features such as the phone duration and 254
 212 statistics of root mean squared energy (Dong et al.,
 213 2024). To alleviate the low-resourced data problem
 214 (Chao et al., 2022), we also consider other self-
 215 supervised learning (SSL) features including
 216 wav2vec 2.0¹ (Baevski et al., 2020), HuBERT²
 217 (Hsu et al., 2021), and WavLM³ (Chen et al., 2022).
 218 All these features are then concatenated and
 219 subsequently projected through a linear layer to
 220 form a sequence of acoustic features \mathbf{X} . The
 221 transformation of each time step t is given by:

$$\mathbf{a}_t = [\mathbf{a}_t^{gop}; \mathbf{a}_t^{dur}; \mathbf{a}_t^{eng}; \mathbf{a}_t^{w2v}; \mathbf{a}_t^{hbt}; \mathbf{a}_t^{wlm}] \quad (1)$$

$$\mathbf{x}_t = \mathbf{W}\mathbf{a}_t + \mathbf{b} \quad (2)$$

222 where \mathbf{W} and \mathbf{b} are trainable parameters. Notably,
 223 a dropout rate of 10% is applied to all SSL features
 224 prior to the concatenation due to the discrepancy in
 225 dimensionality between these and other features.

226 **Phonological Feature Extraction:** In addition to
 227 acoustic cues, a common practice in CAPT is to
 228 inject the phonological information by introducing
 229 the reference text prompt features such as
 230 canonical phoneme embeddings (Gong et al.,
 231 2022), context-aware sup-phoneme embeddings
 232 (Chao et al., 2023), and vowel/consonant
 233 embeddings (Fu et al., 2021). In contrast to
 234 previous studies (Gong et al., 2022; Chao et al.,
 235 2022; Do et al., 2023a), we extract the canonical
 236 phoneme embeddings \mathbf{E}^{phn} from \mathbf{p} using a phone
 237 embedding layer that includes the silence (SIL)
 238 information, which has been shown to be crucial
 239 when evaluating a learner’s spoken proficiency. In
 240 addition, an absolute positional embedding \mathbf{E}^{abs}

and a relative position embedding \mathbf{E}^{rel} are
 extracted. Distinct from \mathbf{E}^{abs} , \mathbf{E}^{rel} denotes
 relative positions of phones in a word using tokens
 such as begin [B], internal [I], end [E], and
 single-phone word [S] tokens. For special cases
 of silence positions, we explicitly categorize them
 as either long silence [LS] or short silence [SS]
 based on their duration. According to the guideline
 suggested by ETS (Evanini et al., 2015), positions
 with a silence duration exceeding 0.495 seconds
 are assigned to [LS]; otherwise, they are assigned
 to [SS]. Finally, all these embedding features are
 point-wise added to \mathbf{X} to obtain phone-level input
 features for subsequent modeling:

$$\mathbf{H}_0^{phn} = \mathbf{X} + \mathbf{E}^{phn} + \mathbf{E}^{abs} + \mathbf{E}^{rel} \quad (3)$$

255 **Mamba Blocks:** Recently, the state space model
 256 (SSM) and its variants have gained widespread
 257 adoption for sequence modeling. Among them,
 258 Mamba (Gu and Dao, 2023) has shown outstanding
 259 performance over Transformer (Vaswani et al.,
 260 2017) across various domains and tasks, including
 261 natural language processing (NLP) (Gu and Dao,
 262 2023), computer vision (CV) (Zhu et al., 2024), and
 263 also speech processing (Zhang et al., 2024).
 264 Different from previous SSM instantiations,
 265 Mamba features an input-dependent selection
 266 mechanism and a hardware-aware algorithm,
 267 allowing for efficient input information filtering by
 268 dynamically adjusting the SSM parameters based
 269 on the input data. This also facilitate faster
 270 recurrent computation of the model using scan.

Nevertheless, the original Mamba conducts causal
 computations in a unidirectional manner, relying
 solely on historical information, which prevents it
 from capturing global dependencies as effectively
 as the multi-head self-attention (MHSA) module
 involved in Transformer. To address this, we
 explore bidirectional variant of Mamba as the basic
 modeling block. In this approach, we replace the
 MHSA module in the Transformer encoder with a
 bidirectional Mamba layer, as depicted in Figure 2.
 Specifically, for input \mathbf{H}^{g_i} to the Mamba block at
 granularity level g , the output $\mathbf{H}^{g_{i+1}}$ of the block is:

$$\mathbf{H}'^{g_i} = \text{BiMamba}(\text{LayerNorm}(\mathbf{H}^{g_i})) + \mathbf{H}^{g_i} \quad (4)$$

$$\mathbf{H}^{g_{i+1}} = \text{FFN}(\text{LayerNorm}(\mathbf{H}'^{g_i})) + \mathbf{H}'^{g_i} \quad (5)$$

¹<https://huggingface.co/facebook/wav2vec2-large-xlsr-53>

²<https://huggingface.co/facebook/hubert-large-1160k>

³<https://huggingface.co/microsoft/wavlm-large>

where BiMamba denotes the bidirectional Mamba layer and FFN refers to the feed-forward module, respectively. Notably, there are several studies investigating the bidirectional processing of Mamba (Liang et al., 2024; Zhang et al., 2024; Jiang et al., 2024). In this work, we use a similar structure as Jiang et al. (2024) to implement the bidirectional Mamba layer. For input \mathbf{N}^{g_i} from the output of layer normalization of \mathbf{H}^{g_i} to a bidirectional Mamba layer, the corresponding output \mathbf{M}^{g_i} is computed as follows:

$$\mathbf{Z}^{g_i} = \text{Linear}(\mathbf{N}^{g_i}) \quad (6)$$

$$\mathbf{S}^{g_i \rightarrow} = \text{Linear}(\mathbf{N}^{g_i}), \quad \mathbf{S}^{g_i \leftarrow} = \text{Flip}(\mathbf{S}^{g_i \rightarrow}) \quad (7)$$

$$\begin{cases} \mathbf{C}^{g_i \rightarrow} = \text{Conv1D}^{\rightarrow}(\mathbf{S}^{g_i \rightarrow}) \\ \mathbf{C}^{g_i \leftarrow} = \text{Conv1D}^{\leftarrow}(\mathbf{S}^{g_i \leftarrow}) \end{cases} \quad (8)$$

$$\begin{cases} \mathbf{O}^{g_i \rightarrow} = \sigma(\mathbf{Z}^{g_i}) \otimes \text{SSM}^{\rightarrow}(\mathbf{C}^{g_i \rightarrow}) \\ \mathbf{O}^{g_i \leftarrow} = \sigma(\mathbf{Z}^{g_i}) \otimes \text{SSM}^{\leftarrow}(\mathbf{C}^{g_i \leftarrow}) \end{cases} \quad (9)$$

$$\mathbf{M}^{g_i} = \text{Linear}\left(\frac{1}{2}\mathbf{O}^{g_i \rightarrow} + \frac{1}{2}\text{Flip}(\mathbf{O}^{g_i \leftarrow})\right) \quad (10)$$

where $\mathbf{S}^{g_i \rightarrow}$ and $\mathbf{S}^{g_i \leftarrow}$ denote the forward and backward sequence features, respectively. Specifically, $\mathbf{S}^{g_i \leftarrow}$ is derived from $\mathbf{S}^{g_i \rightarrow}$ by a flipping operation $\text{Flip}(\cdot)$. $\text{Conv1D}(\cdot)$, $\sigma(\cdot)$, and $\text{SSM}(\cdot)$ represents the 1-D convolution, activation function, and selective SSM algorithm described in Mamba (Gu and Dao, 2023), respectively.

Hierarchical Mamba: Since the speech signals are typically distinguished by the complex hierarchical composition, prior studies (Do et al., 2023a; Chao et al., 2023) have suggested that hierarchical modeling structures is more amenable than parallel modeling structures (Gong et al., 2022). To capture the linguistic hierarchy while retaining the cross-aspect relations within the same linguistic unit, we design and instantiate our model with a hierarchical structure and introduce Mamba blocks to model the global dependencies at each granularity level. More concretely, our approach generates finer granularity scores at the lower layers and coarser granularity scores at the higher layers, as exhibited in Figure 2. In phone-level modeling, we first use \mathbf{H}_0^{phn} as the input into L_p -layer Mamba blocks to obtain the phone-level contextualized representations $\mathbf{H}_{L_p}^{phn}$:

$$\mathbf{H}_{L_p}^{phn} = \text{MambaBlock}_{phn}(\mathbf{H}_0^{phn}) \quad (11)$$

Subsequently, $\mathbf{H}_{L_p}^{phn}$ are then propagated forward into the APA module and the MDD module for solving a regression and a sequence classification problem, respectively. The APA module contains one regressor that aims to predict the phone-level aspect score s_0^{phn} (accuracy). On the other hand, the MDD module comprises a classifier and a softmax function that cooperatively learn a distribution \hat{y}_t over the phoneme classes C for each time step t . The diagnosis y_t can then be identified by applying the argmax function to \hat{y}_t . In this work, we streamline the MDD task by treating it as a process of free phone recognition (Li et al., 2015). As a result, we can directly detect the corresponding error state e_t by comparing y_t with p_t , eliminating the need for a separate detection module. Meanwhile, the resulting $\mathbf{H}_{L_p}^{phn}$ is served as \mathbf{H}_0^{word} for subsequent modeling.

In word-level modeling, L_w -layer Mamba blocks are first adopted and followed by a 1-D convolution layer to capture the local dependencies (Lee, 2016). The reason for utilizing the convolution layer is that the convolution operation can accommodate different realizations of the same underlying phone from various L2 speakers, thereby mitigating the temporal variability. The word-level representations $\mathbf{H}_{L_w}^{word}$ can be derived as follows:

$$\mathbf{H}'_{L_w}^{word} = \text{MambaBlock}_{word}(\mathbf{H}_0^{word}) \quad (12)$$

$$\mathbf{H}_{L_w}^{word} = \text{Conv1D}_{word}(\mathbf{H}'_{L_w}^{word}) \quad (13)$$

To obtain word-level aspect scores, we put $\mathbf{H}_{L_w}^{word}$ into the word-level APA module which contains three regressors to predict the word-level aspect scores s_0^{word} , s_1^{word} , s_2^{word} (accuracy, stress, and total scores), respectively. To facilitate training efficiency, we propagate the word score to each of its phones during the training stage. In the inference phase, we ensure consistency by averaging the outputs corresponding to each word. In addition, $\mathbf{H}_{L_w}^{word}$ is viewed as \mathbf{H}_0^{utt} for further modeling.

As for the utterance-level assessments, instead of prepending the [CLS] tokens to learn the utterance-level representation (Gong et al., 2022), we explore pooling-based approaches to aggregate the hidden information. To this end, we utilize an attention pooling layer similar to Peng et al. (2022). Specifically, assuming that a d -dimensional input

364 sequence to the attention pooling layer is
 365 $\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_{T-1}$, the pooling output is $\mathbf{h} =$
 366 $\sum_{i=0}^{T-1} \alpha_i \mathbf{h}_i$, where α_i is calculated as follows:

$$\alpha_i = \frac{\exp(\mathbf{w}^T \mathbf{q}_i / \tau)}{\sum_{j=0}^{T-1} \exp(\mathbf{w}^T \mathbf{q}_j / \tau)} \quad (14)$$

367 where \mathbf{w} is a learnable vector, \mathbf{q} is the
 368 concatenated scores of $[s_0^{g^{phn}}, s_0^{g^{wrđ}}, s_1^{g^{wrđ}}, s_2^{g^{wrđ}}]$,
 369 and τ is a controllable temperature hyperparameter.
 370 The whole process of utterance-level modeling can
 371 then be formulated as follows:

$$\mathbf{H}^{g_{L_u}^{utt}} = \text{MambaBlock}_{utt}(\mathbf{H}^{g_0^{utt}}) \quad (15)$$

$$\mathbf{h}^{g_{L_u}^{utt}} = \text{AttentionPooling}_{utt}(\mathbf{H}^{g_{L_u}^{utt}}) \quad (16)$$

372 After obtaining $\mathbf{H}^{g_{L_u}^{utt}}$ from L_u -layer Mamba
 373 blocks, $\mathbf{h}^{g_{L_u}^{utt}}$ is derived through the attention
 374 pooling layer to predict the utterance-level aspect
 375 scores $s_0^{g_{L_u}^{utt}}, s_1^{g_{L_u}^{utt}}, s_2^{g_{L_u}^{utt}}, s_3^{g_{L_u}^{utt}}, s_4^{g_{L_u}^{utt}}$ (accuracy,
 376 completeness, fluency, prosody, and total scores)
 377 via an utterance-level APA module which contains
 378 five regressors corresponding to each score.

379 2.3 Optimization

380 **Automatic Pronunciation Assessment Loss:** In
 381 the proposed model, each APA module is
 382 optimized using Mean Square Error (MSE). The
 383 loss for multi-aspect multi-granular assessment,
 384 \mathcal{L}_{APA} , is calculated by assigning weights to each
 385 granularity level g :

$$\mathcal{L}_{APA} = \sum_{g \in G} \omega_g \cdot \frac{1}{N_g} \sum_{k=0}^{N_g-1} \mathcal{L}_{g_k} \quad (17)$$

386 where ω_g and N_g are the tunable parameter and
 387 number of aspect scores at granularity level g ,
 388 respectively. \mathcal{L}_{g_k} refers to the MSE loss computed
 389 for k -th aspect score at granularity level g .

390 **Mispronunciation Detection and Diagnosis Loss:**
 391 To be in line with previous MDD studies, our
 392 model incorporates canonical phoneme
 393 embeddings to enhance text prompt-awareness.
 394 Despite some performance improvements, the
 395 mismatch between the L2 learner’s realized phones
 396 and canonical phones can still cause some
 397 deteriorating effects. This discrepancy can
 398 introduce inaccurate predictions that may
 399 potentially affect the overall quality of phonetic
 400 analysis. To mitigate this negative impact, we
 401 devise a new loss function tailored for the MDD

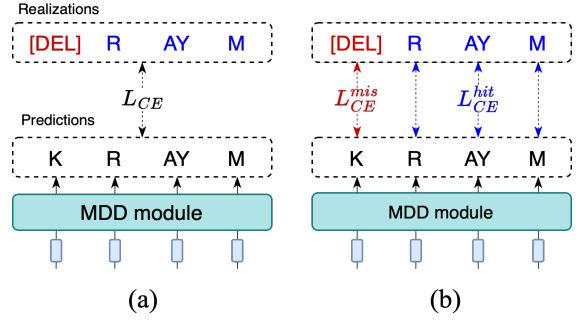


Figure 3: Difference between (a) the original cross-entropy loss and (b) the decoupled cross-entropy loss, given the text prompt “crime.”

402 task, as illustrated in Figure 3. Specifically, we first
 403 decouple the original cross-entropy loss into two
 404 separate losses, one for mispronunciations and the
 405 other for correct pronunciations:

$$\mathcal{L}_{Xent}^{mis} = - \sum_{t \in \mathcal{M}} \log(\hat{y}_t[y_t]) \quad (18)$$

$$\mathcal{L}_{Xent}^{hit} = - \sum_{t \in \mathcal{H}} \log(\hat{y}_t[y_t]) \quad (19)$$

406 where \mathcal{M} and \mathcal{H} are mispronunciation and
 407 correct pronunciation positions, respectively, and
 408 $\hat{y}_t[y_t]$ is the predicted probability of the true label
 409 y_t at time step t . After obtaining two decoupled
 410 losses, we re-weight them using the following
 411 formulation:

$$\mathcal{L}_{MDD} = \mathcal{L}_{Xent}^{hit} + \left(\frac{\mu^h}{\mu^m}\right)^\alpha \mathcal{L}_{Xent}^{mis} \quad (20)$$

412 where μ^m and μ^h denote the frequency of the
 413 mispronunciations and correct pronunciations in
 414 the training set, respectively, and α controls the
 415 weight magnitude. After that, we use \mathcal{L}_{MDD} to
 416 optimize the MDD module, and the overall loss
 417 thus can be expressed by:

$$\mathcal{L} = \mathcal{L}_{APA} + \beta \cdot \mathcal{L}_{MDD} \quad (21)$$

where β is a tunable parameter.

419 3 Experimental Setup

420 3.1 Dataset and Evaluation Metrics

421 We conducted experiments on speechocean762, a
 422 widely-used open-source dataset curated for APA
 423 and MDD research (Zhang et al., 2021). The
 424 dataset consists of 5,000 English-speaking
 425 recordings from 250 Mandarin L2 learners, divided
 426 equally into training and test sets. For the APA task,
 427 pronunciation proficiency scores were assessed at

Model	Year	Phone Score		Word Score (PCC)			Utterance Score (PCC)				
		MSE↓	PCC↑	Accuracy↑	Stress↑	Total↑	Accuracy↑	Completeness↑	Fluency↑	Prosody↑	Total↑
Deep Feature	2021	-	-	-	-	-	-	-	-	-	0.720
HuBERT Large	2022	-	-	-	-	-	-	-	0.780	0.770	-
Joint-CAPT-L1	2023	-	-	-	-	-	0.719	-	0.775	0.773	0.743
LSTM	2022	0.089	0.591	0.514	0.294	0.531	0.720	0.076	0.745	0.747	0.741
GOPT	2022	0.085	0.612	0.533	0.291	0.549	0.714	0.155	0.753	0.760	0.742
3M	2022	0.078	0.656	0.598	0.289	0.617	0.760	0.325	0.828	0.827	0.796
HiPAMA	2023	0.084	0.616	0.575	0.320	0.591	0.730	0.276	0.749	0.751	0.754
3MH	2023	0.071	0.693	0.682	0.361	0.694	0.782	0.374	0.843	0.836	0.811
HMamba	2024	0.063	0.732	0.701	0.309	0.710	0.802	0.210	0.846	0.841	0.825

Table 1: APA performance evaluations of our model and all strong baselines on the speechocean762 test set.

various linguistic granularities and across different pronunciation aspects. Each score is evaluated by five experienced experts using standardized rubrics. For the MDD task, the dataset provides an extra mispronunciation transcription annotated using a set of 46 phones. This set comprises 39 phones from the CMU dictionary⁴, 6 L2-specific phones, and a [unk] token for unknown phones. Notably, there are no insertion errors in the utterances, and a [DEL] token is introduced to mark deletion errors of L2 learners. Therefore, the realized phones can be aligned with canonical phones in this dataset. The evaluation metrics employed include the Pearson Correlation Coefficient (PCC) and Mean Square Error (MSE) for the APA task. On the other hand, we use precision, recall, F1-score, and phone error rate (PER) to evaluate the MDD performance, so as to be in accordance with prior studies.

3.2 Implementation Details

For input feature extraction, we adopt a publicly available acoustic model⁵ to extract GOP features, which also serves as an aligner for force alignment. Subsequently, the phone-level duration, energy statistics, and SSL features are computed by a time aggregation method (Kim et al., 2022) according to the alignment. The resulting acoustic features X and all embeddings are 128 dimensions. For all Mamba blocks, we set the number of hidden units to 128 and use a kernel size of 4 for the 1-D convolution. The SSM modules follow the original configuration used in Mamba. L_p , L_w , L_u are set to 3, 1, 1, respectively. In addition, the word-level

⁴ <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

Model	Mispronunciations			PER↓
	Precision↑	Recall↑	F1↑	
Joint-CAPT-L1	26.70%	91.40%	41.50%	9.93%
HMamba	64.50%	62.34%	63.32%	2.78%

Table 2: MDD performance evaluations of our model, compared with a representative multi-task approach (Ryu et al., 2023) on the speechocean762 test set.

1-D convolution has 256 kernels, each with a size of 3. Regarding hyperparameters, τ in attention pooling layer is set to 1.0. The combining weights ω_g for APA loss are uniformly set to 1.0 for each granularity level g . Parameters α and β are tuned to be 0.7 and 0.003, respectively. To ensure the validity of our experimental results, we conducted 5 independent trials for each experiment, running 20 epochs with different seeds. The metrics for each task are reported as the average of these trials.

3.3 Compared Baselines

For the APA task, we compare our proposed approach, HMamba, with various cutting-edge baselines which can be categorized into two families: single-aspect (or partial-aspect) pronunciation assessment models or multi-granular multi-aspect pronunciation assessment models. The first group includes the Deep Feature (Lin et al., 2021b), HuBERT Large (Kim et al., 2022), and Joint-CAPT-L1 (Ryu et al., 2023). The second group encompasses LSTM, GOPT (Gong et al.,

⁵ <https://kaldi-asr.org/models/m13>

2022), 3M (Chao et al., 2022), HiPAMA (Do et al., 2023a), and 3MH (Chao et al., 2023). As for the MDD task, we compare HMamba with the Joint-CAPT-L1 model, as to our knowledge it is the only attempt that jointly addresses the APA and MDD tasks with the speechocean762 dataset.

4 Experimental Results and Discussion

4.1 APA Performance

In Table 1, we compare the APA performance of HMamba with other competitive baselines, leading to several key observations. Firstly, it is notable that our approach, HMamba, consistently outperforms all other methods in nearly all assessment tasks, particularly in terms of accuracy scores at phone, word, and utterance levels. This improvement stems from the joint modeling paradigm of APA and MDD, highlighting that pronunciation assessment can also benefit from phonetic error discovery, consistent with prior research findings. In addition, by adopting SSL features, HMamba along with other approaches like HuBERT Large, 3M, and 3MH, achieves significant improvements over the other APA methods in terms of utterance-level assessments. In comparison to other hierarchical models such as HiPAMA and 3MH, HMamba leverages an SSM structure instead of the Transformer structure, demonstrating superior performance on a variety of assessment tasks. Furthermore, due to severe imbalance issues in the aspects of utterance completeness and word stress, where over 90% of assessments consistently receive the highest score (Do et al., 2023b), our approach slightly falls behind the other approaches.

4.2 MDD Performance

In the second set of experiments, we evaluate the MDD performance of HMamba by comparing it with another advanced multi-task learning approach, Joint-CAPT-L1. As shown in Table 2, HMamba achieves a significant improvement in terms of F1-score over Joint-CAPT-L1, with a relative increase of 21.82%. Additionally, there is a marked reduction in PER by 7.15%. These substantial enhancements demonstrate that HMamba can produce more robust and reliable mispronunciation detection and diagnosis results.

4.3 Effects of Decoupled Cross-entropy Loss

On the grounds of the distinct improvements in the MDD performance, we further analyze the

Loss	α	Mispronunciations			PER ↓
		Precision ↑	Recall ↑	F1 ↑	
Xent	-	74.15%	40.21%	52.12%	2.58%
	0.3	68.60%	55.74%	61.49%	2.62%
	0.5	63.51%	61.43%	62.32%	2.83%
deXent	0.7	64.50%	62.34%	63.32%	2.78%
	0.9	57.73%	70.11%	63.19%	3.14%

Table 3: Comparison of MDD performance between the original cross-entropy loss (Xent) and proposed decoupled cross-entropy loss (deXent).

underlying effects of proposed decoupled cross-entropy loss on model performance. As illustrated in Table 3, training a text prompt-aware MDD model using the original cross-entropy often yields high precision but low recall. This is because the model primarily relies on input canonical phones, leading it to predict prior phones and overlook the actual mispronunciations of a learner. Such a model may not be suitable for educational settings where accurately detecting potential mispronunciations is critical. To remedy this, the proposed decoupled cross-entropy loss provides a feasible solution. By adjusting the weighting factor α , we can better strike the balance between precision and recall, thus optimizing the MDD performance. This flexibility is particularly prominent across different CAPT applications. For example, in a clinical setting such as speech therapy, prioritizing precision can help prevent incorrect diagnoses of speech disorders.

5 Conclusion

In this paper, we have presented a novel hierarchical selective state space model (dubbed HMamba) for multifaceted CAPT application. Extensive experimental results substantiate the viability and efficacy of the proposed method compared to several top-of-the-line approaches in terms of both the APA and MDD performance. In future work, we envisage mitigating the issue of data imbalance from an optimization perspective. In addition, another key area for future research involves tackling the assessment of open-response scenarios in CAPT.

562 Limitations

563 **Lack of Accent Diversity.** The dataset used in this
564 study comprises only Mandarin L2 learners,
565 limiting the generalizability of the proposed model.
566 As a result, it may be inapplicable when assessing
567 L2 learners with diverse accents. This lack of
568 accent diversity could lead to biases and
569 inaccuracies in pronunciation assessment for
570 learners from different linguistic backgrounds.

571 **Limited Interpretability.** The proposed model is
572 designed to replicate expert annotations without
573 relying on manual assessment rubrics or external
574 knowledge databases, which makes it challenging
575 to provide clear and reasonable explanations for the
576 assessment results. This lack of interpretability
577 may hinder its acceptance and trustworthiness
578 among educators and learners who require
579 transparent and justifiable assessments.

580 **Limited Generalizability** This research is
581 centered on the “reading-aloud” pronunciation
582 training scenario, where it is assumed that the L2
583 learner accurately pronounces a predetermined text
584 prompt. This narrow focus limits the applicability
585 of our models to other learning contexts, such as
586 spontaneous speech or open-ended conversations.

587 Ethics Statement

588 We acknowledge that all co-authors of this work
589 comply with the ACL Code of Ethics and adhere to
590 the code of conduct. Our experimental corpus,
591 speechocean762, is widely used and publicly
592 available, and we believe there are no potential
593 risks associated with this work.

594 References

595 Juan Pablo Arias, Nestor Becerra Yoma, and Hiram
596 Vivanco. 2010. Automatic intonation assessment for
597 computer aided language learning. *Speech*
598 *Communication*, volume 52, pages 254–267.

599 Stefano Bannò, Bhanu Balusu, Mark Gales, Kate Knill,
600 and Konstantinos Kyriakopoulos. 2022a. View-
601 specific assessment of L2 spoken English. In
602 *Proceedings of the Annual Conference of the*
603 *International Speech Communication Association*
604 *(INTERSPEECH)*, pages 4471–4475.

605 Stefano Bannò and Marco Matassoni. 2022b.
606 Proficiency assessment of L2 spoken English using
607 wav2vec 2.0. In *Proceedings of IEEE Spoken*
608 *Language Technology Workshop (SLT)*, pages 1088-
609 1095.

610 Alexei Baeviski, Yuhao Zhou, Abdelrahman Mohamed,
611 and Michael Auli. 2020. wav2vec 2.0: A framework
612 for self-supervised learning of speech
613 representations. In *Proceedings of the Conference*
614 *on Neural Information Processing Systems*
615 *(NeurIPS)*, pages 12449–12460.

616 Fu An Chao, Tien Hong Lo, Tzu I. Wu, Yao Ting Sung,
617 Berlin Chen. 2022. 3M: An effective multi-view,
618 multigranularity, and multi-aspect modeling
619 approach to English pronunciation assessment. In
620 *Proceedings of the Asia-Pacific Signal and*
621 *Information Processing Association Annual Summit*
622 *and Conference (APSIPA ASC)*, pages 575–582.

623 Fu-An Chao, Tien-Hong Lo, Tzu-I Wu, Yao-Ting Sung,
624 Berlin Chen. 2023. A hierarchical context-aware
625 modeling approach for multi-aspect and multi-
626 granular pronunciation assessment. In *Proceedings*
627 *of the Annual Conference of the International*
628 *Speech Communication Association*
629 *(INTERSPEECH)*, pages 974–978.

630 Nancy F. Chen, and Haizhou Li. 2016. Computer-
631 assisted pronunciation training: From pronunciation
632 scoring towards spoken language learning. In
633 *Proceedings of the Asia-Pacific Signal and*
634 *Information Processing Association Annual Summit*
635 *and Conference (APSIPA ASC)*, pages 1–7.

636 Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu
637 Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki
638 Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu,
639 Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian
640 Wu, Michael Zeng, Xiangzhan Yu, and Furu Wei.
641 2022. Wavlm: Large-scale self-supervised pre-
642 training for full stack speech processing. *IEEE*
643 *Journal of Selected Topics in Signal Processing*,
644 volume 16, number 6, pages 1505-1518.

645 Heejin Do, Yunsu Kim, and Gary Geunbae Lee. 2023a.
646 Hierarchical pronunciation assessment with multi-
647 aspect attention. In *Proceedings of the IEEE*
648 *International Conference on Acoustics, Speech and*
649 *Signal Processing (ICASSP)*, pages 1–5.

650 Heejin Do, Yunsu Kim, and Gary Geunbae Lee. 2023b.
651 Score-balanced loss for multi-aspect pronunciation
652 assessment. In *Proceedings of the Annual*
653 *Conference of the International Speech*
654 *Communication Association (INTERSPEECH)*,
655 pages 4998–5002.

656 Hongwei Ding, Xinpeng Xu. 2016. L2 English rhythm
657 in read speech by Chinese students. In *Proceedings*
658 *of the Annual Conference of the International*
659 *Speech Communication Association*
660 *(INTERSPEECH)*, pages 2696-2700.

661 Bin Dong, Qingwei Zhao, Jianping Zhang, and
662 Yonghong Yan. 2004. Automatic assessment of
663 pronunciation quality. In *Proceedings of IEEE*

- 664 *International Symposium on Chinese Spoken*
665 *Language Processing (ISCSLP)*, pages 137-140.
- 666 Maxine Eskenazi. 2009. An overview of spoken
667 language technology for education. *Speech*
668 *Communication*, volume 51, pages 832–844.
- 669 Keelan Evanini, Michael Heilman, Xinhao Wang, and
670 Daniel Blanchard. 2015. Automated scoring for the
671 TOEFL Junior® comprehensive writing and
672 speaking test. *ETS Research Report Series*
673 2015(1):1–11.
- 674 Kaiqi Fu, Jones Lin, Dengfeng Ke, Yanlu Xie, Jinsong
675 Zhang, and Binghuai Lin. 2021. A full text-
676 dependent end to end mispronunciation detection
677 and diagnosis with easy data augmentation
678 techniques. *arXiv preprint arXiv:2104.08428*.
- 679 Yuan Gong, Ziyi Chen, Iek-Heng Chu, Peng Chang,
680 and James Glass. 2022. Transformer-based multi-
681 aspect multigranularity non-native English speaker
682 pronunciation assessment. In *Proceedings of the*
683 *IEEE International Conference on Acoustics,*
684 *Speech and Signal Processing (ICASSP)*, pages
685 7262–7266.
- 686 Albert Gu and Tri Dao. 2023. Mamba: Linear-time
687 sequence modeling with selective state spaces.
688 *arXiv preprint arXiv:2312.00752*.
- 689 Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,
690 Kushal Lakhota, Ruslan Salakhutdinov, and
691 Abdelrahman Mohamed. 2021. Hubert: Self-
692 supervised speech representation learning by
693 masked prediction of hidden units. *IEEE/ACM*
694 *Transactions on Audio, Speech, and Language*
695 *Processing*, volume 29, pages 3451–3460.
- 696 Alissa M. Harrison, Wai-Kit Lo, Xiao-Jun Qian, and
697 Helen Meng. 2009. Implementation of an extended
698 recognition network for mispronunciation detection
699 and diagnosis in computer-assisted pronunciation
700 training. In *Proceedings of the Workshop on Speech*
701 *and Language Technology in Education (SLaTE)*,
702 pages 45-48.
- 703 Wenping Hu, Yao Qian, Frank K. Soong, and Yong
704 Wang. 2015. Improved mispronunciation detection
705 with deep neural network trained acoustic models
706 and transfer learning based logistic regression
707 classifiers. *Speech Communication*, volume 67,
708 pages 154–166.
- 709 Xilin Jiang, Cong Han, and Nima Mesgarani. Dual-
710 path mamba: Short and long-term bidirectional
711 selective structured state space models for speech
712 separation. 2024. *arXiv preprint arXiv:2403.18257*.
- 713 Yassine Kheir, Ahmed Ali, and Shammur Chowdhury.
714 2023. Automatic pronunciation assessment - a
715 review. In *Findings of the Association for*
716 *Computational Linguistics: EMNLP*, pages 8304–
717 8324.
- 718 Eesung Kim, Jae-Jin Jeon, Hyeji Seo, Hoon Kim. 2022.
719 Automatic pronunciation assessment using self-
720 supervised speech representation learning. In
721 *Proceedings of the Annual Conference of the*
722 *International Speech Communication Association*
723 *(INTERSPEECH)*, pages 1411–1415.
- 724 Ann Lee. 2016. Language-independent methods for
725 computer-assisted pronunciation training, *Ph.D.*
726 *thesis, Massachusetts Institute of Technology*.
- 727 Aobo Liang, Xingguo Jiang, Yan Sun, Xiaohou Shi,
728 and Ke Li. 2024. Bi-Mamba4TS: Bidirectional
729 mamba for time series forecasting. *arXiv preprint*
730 *arXiv:2404.15772*.
- 731 Binghuai Lin, Liyuan Wang, Hongwei Ding, Xiaoli
732 Feng. 2021a. Improving L2 English rhythm
733 evaluation with automatic sentence stress detection.
734 In *Proceedings of IEEE Spoken Language*
735 *Technology Workshop (SLT)*, pages 713-719.
- 736 Binghuai Lin and Liyuan Wang. 2021b. Deep feature
737 transfer learning for automatic pronunciation
738 assessment. In *Proceedings of the Annual*
739 *Conference of the International Speech*
740 *Communication Association (INTERSPEECH)*,
741 pages 4438–4442.
- 742 Kun Li, Xiaojun Qian, Shiyong Kang, Pengfei Liu, and
743 Helen Meng. 2015. Integrating acoustic and state-
744 transition models for free phone recognition in L2
745 English speech using multi-distribution deep neural
746 networks. In *Proceedings of the Workshop on*
747 *Speech and Language Technology in Education*
748 *(SLaTE)*, pages. 119-124.
- 749 Arya D. McCarthy, Kevin P. Yancey, Geoffrey T.
750 LaFlair, Jesse Egbert, Manqian Liao, and Burr
751 Settles. 2021. Jump-starting item parameters for
752 adaptive language tests. In *Proceedings of the*
753 *Conference on Empirical Methods in Natural*
754 *Language Processing (EMNLP)*, pages 883–899.
- 755 Yifan Peng, Siddharth Dalmia, Ian Lane, and Shinji
756 Watanabe. 2022. Branchformer: Parallel mlp-
757 attention architectures to capture local and global
758 context for speech recognition and understanding.
759 In *Proceedings of the International Conference on*
760 *Machine Learning (PMLR)*, pages 17627–17643.
- 761 Hyungshin Ryu, Sunhee Kim, and Minhwa Chung.
762 2023. A joint model for pronunciation assessment
763 and mispronunciation detection and diagnosis with
764 multi-task learning. In *Proceedings of the Annual*
765 *Conference of the International Speech*
766 *Communication Association (INTERSPEECH)*,
767 pages 959-963.
- 768 Jiatong Shi, Nan Huo, and Qin Jin. 2020. Context-
769 aware goodness of pronunciation for computer-
770 assisted pronunciation training. In *Proceedings of*
771 *the Annual Conference of the International Speech*

772 *Communication Association (INTERSPEECH)*,
773 pages 3057-3061.

774 Helmer Strik, Khiet Truong, Febe De Wet, and Catia
775 Cucchiarini. 2009. Comparing different approaches
776 for automatic pronunciation error detection. *Speech*
777 *Communication*, volume 51, number 10, pages 845-
778 852.

779 Khiet Truong, Ambra Neri, Catia Cucchiarini, and
780 Helmer Strik. 2004. Automatic pronunciation error
781 detection: an acoustic-phonetic approach.
782 *InSTIL/ICALL Symposium*.

783 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
784 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
785 Kaiser, and Illia Polosukhin. 2017. Attention is all
786 you need. In *Proceedings of the Conference on*
787 *Neural Information Processing Systems (NeurIPS)*,
788 pages 5998–6008.

789 Anjana S. Vakil and Jürgen Trouvain. 2015.
790 “Automatic classification of lexical stress errors for
791 German CAPT,” in *Proceedings of the Workshop on*
792 *Speech and Language Technology in Education*
793 *(SLaTE)*, pages 47–52.

794 Silke M. Witt and Steve J. Young. 2000. Phone-level
795 pronunciation scoring and assessment for
796 interactive language learning. *Speech*
797 *Communication*, volume 30, pages 95–108.

798 Klaus Zechner, Derrick Higgins, Xiaoming Xi, and
799 David M. Williamson. 2009. Automatic scoring of
800 non-native spontaneous speech in tests of spoken
801 English. *Speech Communication*, volume 51,
802 number 10, pages 883-895.

803 Lianghai Zhu, Bencheng Liao, Qian Zhang, Xinlong
804 Wang, Wenyu Liu, and Xinggang Wang. 2024.
805 Vision mamba: Efficient visual representation
806 learning with bidirectional state space model. *arXiv*
807 *preprint arXiv:2401.09417*.

808 Xiangyu Zhang, Qiquan Zhang, Hexin Liu, Tianyi
809 Xiao, Xinyuan Qian, Beena Ahmed, Eliathamby
810 Ambikairajah, Haizhou Li, and Julien Epps. 2024.
811 Mamba in Speech: Towards an alternative to self-
812 attention. *arXiv preprint arXiv:2405.12609*.

813 Junbo Zhang, Zhiwen Zhang, Yongqing Wang,
814 Zhiyong Yan, Qiong Song, Yukai Huang, Ke Li,
815 Daniel Povey, and Yujun Wang. 2021.
816 Speechocean762: An open-source non-native
817 English speech corpus for pronunciation assessment.
818 In *Proceedings of the Annual Conference of the*
819 *International Speech Communication Association*
820 *(INTERSPEECH)*, pages 3710–3714.