

Beyond Acoustic Sparsity and Linguistic Bias: A Prompt-Free Paradigm for Mispronunciation Detection and Diagnosis

Anonymous submission to Interspeech 2026

Abstract

Mispronunciation Detection and Diagnosis (MDD) requires modeling fine-grained acoustic deviations. However, current ASR-derived MDD systems often face inherent limitations. In particular, CTC-based models favor sequence-level alignments that neglect transient mispronunciation cues, while explicit canonical priors bias predictions toward intended targets. To address these bottlenecks, we propose a prompt-free framework decoupling acoustic fidelity from canonical guidance. First, we introduce CROTTC, an acoustic model enforcing monotonic, frame-level alignment to accurately capture pronunciation deviations. Second, we implicitly inject mispronunciation information via the IF strategy under the knowledge transfer principle. Experiments show CROTTC-IF achieves a 71.77% F1-score on L2-ARCTIC and 71.70% F1 on the Iqra'Eval2 challenge. With empirical analysis, we demonstrate that decoupling acoustics from explicit priors provides highly robust MDD¹.

Index Terms: Mispronunciation detection and diagnosis Optimal transport, Consistency Regularization, Knowledge Transfer, Large Language Model.

1. Introduction

Mispronunciation Detection and Diagnosis (MDD) plays an indispensable role across a spectrum of applications, ranging from general Computer-Aided Pronunciation Training (CAPT) for L2 learners to religious domains such as Qur'anic recitation, as highlighted by the Iqra'Eval Challenges [1, 2]. The core challenge of MDD lies in acoustic fidelity: transcribing speech exactly as realized to pinpoint deviations from canonical norms. Historically, MDD was formulated as a downstream application of Automatic Speech Recognition (ASR). Conventional systems typically used ASR features to generate Goodness-of-Pronunciation (GOP) scores [3] or relied on Extended Recognition Networks (ERNs) to catch predefined error patterns [4, 5]. In this classical paradigm, MDD was essentially an ASR by-product.

With the rapid advancement of end-to-end (E2E) architectures and self-supervised learning (SSL) speech representations, the paradigm has shifted. MDD is now a standalone task, directly optimized for fine-grained phoneme diagnosis. However, despite this structural independence, MDD research inevitably remains under the heavy methodological shadow of ASR. While general ASR focuses on deducing semantic intent of a speaker despite acoustic imperfections, MDD demands objective, fine-grained acoustic fidelity. By uncritically adopting ASR's acoustic and linguistic optimization strategies, we believe that current

¹Anonymized for double-blind review. Source code, recipes, and checkpoints will be publicly available upon acceptance.

MDD studies frequently fall into two fatal traps.

The first is the *Acoustic Trap*, driven by the direct adoption of Connectionist Temporal Classification (CTC) [6]. As a de facto standard for E2E ASR, CTC usually serves as an acoustic model (AM) by maximizing the marginal probability of all valid alignment paths, forcing the model to prioritize global sequence correctness. However, when ported to the MDD task, this global optimization actively smooths over the subtle, transient acoustic variations (e.g., co-articulation, brief onset substitutions). In doing so, CTC might erase the exact phonetic evidence that is critical for mispronunciation diagnosis.

The second is the *Linguistic Trap*, driven by the reliance on explicit canonical prompts or strong language model (LM) post-processing. In typical ASR systems such as CTC/AED [7] or RNN-T [8], LMs are powerful components to correct acoustic errors. By adjusting the AM/LM weight, ASR models can shift their preference to favor either acoustic or linguistic features. In MDD, however, when models are guided by canonical priors, they naturally exhibit an over-correction tendency. Instead of objectively diagnosing the actual phonetic deviations from the acoustic signal, the system defaults to the linguistically most likely sequence, which is almost always the canonical text.

To break free from ASR's legacy and overcome these bottlenecks, we propose **CROTTC-IF**, a unified, canonical prompt-free paradigm, designed exclusively for the standalone MDD task. Our main contributions are summarized as follows:

- **Frame-wise Acoustic Modeling (CROTTC):** We introduce Consistency-Regularized Optimal Temporal Transport Classification. Unlike standard CTC, this acoustic model enforces strict monotonic, frame-wise alignments, effectively mitigating inherent sparsity and delayed emission issues to preserve fine-grained, transient mispronunciation cues.
- **Knowledge Transfer via Indirect Fusion (IF):** Grounded in the Learning Using Privileged Information (LUPI) paradigm, we introduce Indirect Fusion strategy to implicitly integrate linguistic priors. By treating canonical texts and mispronunciation patterns as privileged data during training, IF provides soft linguistic guidance to the language model, preventing the override of faithful acoustic evidence during inference.
- **Comprehensive Analysis of Canonical Bias:** With delicately designed Multi-modality Large Language Models (LLMs) and various prompt templates, we quantitatively investigate the impact of explicit canonical prior in MDD. Our findings demonstrated that over-reliance on canonical priors can degrade detection sensitivity, highlighting the necessity of balancing linguistic context with acoustic fidelity.
- **State-of-the-Art Performance:** Operating entirely without auxiliary data or explicit canonical prompts, the CROTTC-IF framework achieves highly competitive results. It demon-

46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94

95 strates strong generalization across diverse benchmarks,
96 ranging from general L2 English corpora (L2-ARCTIC, ERJ,
97 speechocean762) to the specialized Arabic Qur’anic recita-
98 tion task (Iqra’Eval2).

99 2. Related Works

100 Currently, modern MDD approaches broadly fall into two cate-
101 gories: dictation-style and text-prompting style.

102 2.1. Dictation-Style MDD and the Acoustic Trap

103 Dictation-style methods aim to recognize the uttered phoneme
104 sequence exclusively from acoustic-related features. For exam-
105 ple, Leung et al. introduced a CTC framework and demon-
106 strated the feasibility of capturing L2 phoneme information
107 from acoustic features alone [9]. More recent studies have lever-
108 aged self-supervised learning (SSL) speech representations,
109 such as wav2vec 2.0 [10], HuBERT [11, 12] and WavLM [13],
110 which provide robust phonemic context features. For in-
111 stance, [14] employed a fine-tuned wav2vec 2.0 model with
112 a CTC decoder, achieving promising diagnostic accuracy. To
113 strengthen free phoneme recognition, researchers have also
114 incorporated semi-supervised pseudo-labeling [15], dual-path
115 SSL for predicting manner-of-articulation [16], and acoustic-
116 to-articulatory inversion leveraging electromagnetic articulog-
117 raphy (EMA) data [17].

118 This dictation-style MDD also dominated the recent
119 Iqra’Eval 2025 challenge [1]. For instance, *Hafs2Vec* [18] uti-
120 lized extensive professional reciter data, while *Metapseud* [19]
121 applied curriculum learning and beam-search decoding. De-
122 spite these feature and data enhancements, the fundamental re-
123 liance on CTC decoders introduces a potential vulnerability: the
124 *Acoustic Trap*².

125 Conversely, frame-wise alignment methods have demon-
126 strated greater potential for MDD. For example, Feng et al.
127 pioneered the use of attention-based sequence labeling to map
128 individual acoustic frames directly to phone labels [20]. Simi-
129 larly, Lin et al. modeled MDD as a sequence labeling task via
130 a joint-alignment mechanism [21]. More recently, Tu et al. in-
131 troduced a retrieval-based framework for frame-wise classifica-
132 tion, achieving competitive performance without task-specific
133 training [22]. These advancements indicate that leveraging
134 **dense, frame-wise representations** provides superior perfor-
135 mance, outperforming sparse features for robust MDD.

136 2.2. Text-Prompting MDD and the Linguistic Trap

137 In reading-aloud scenarios, canonical phoneme sequences are
138 inherently available. Consequently, integrating this prior
139 knowledge into MDD models has become another standard
140 practice. For instance, Ye et al. proposed an approach to
141 leverage acoustic, phonetic and linguistic embeddings for MDD
142 [23]. Peng et al. employed a gating mechanism and contrastive
143 loss to dynamically regulate canonical embeddings [24]. Zheng
144 et al. utilized coupled-cross attention for explicit acoustic-
145 canonical fusion [25]. Similarly, Yan et al. injected implicit
146 linguistic priors using phoneme lookup tables and graph net-
147 works to model frequent mispronunciation patterns [26, 27].

148 Entering the era of Large Language Models (LLMs), re-
149 cent studies have shifted toward leveraging massive linguis-
150 tic priors for posterior correction in MDD. Current methods
151 range from adapting LLM-based ASR architectures [28] to fine-

152 tuning native SpeechLM [29]. For instance, Wu et al. adapted
153 LLMs for MDD by coupling a pretrained acoustic front-end
154 with a trainable projector and LLM backbone [30]. Other re-
155 searchers adapted multi-modality LLM to embed canonical se-
156 quences and predefined error patterns directly into the latent
157 space [31, 32]. Beyond simple transcription, LLM reasoning
158 is also deployed to generate educational feedback, offering nu-
159 anced diagnostics alongside canonical norms [33, 34].

160 Despite reported successes, the practical superiority of text-
161 prompted and LLM-assisted methods remains debatable. Many
162 systems are evaluated on private datasets, precluding repro-
163 ducible comparisons. And on open benchmarks, they frequently
164 underperform compared to conventional models that rely solely
165 on acoustic modeling.

166 Beyond evaluation discrepancies, we believe that some of
167 these canonical-dependent methods inherently fall into the *Lin-*
168 *guistic Trap* due to a shared flaw: **canonical information leak-**
169 **age**. Whether by explicitly including canonical information
170 as part of the input, or by using a strong language model as
171 the model’s prior knowledge, these approaches can all be cat-
172 egorized as suffering from canonical information leakage. Al-
173 though many models carefully design their canonical informa-
174 tion usage, canonical information can easily override the sub-
175 tle acoustic information. Consequently, the model becomes
176 overly tolerant of errors, masking critical acoustic deviations
177 and sacrificing diagnostic objectivity. Furthermore, requiring
178 explicit canonical texts during inference renders these methods
179 entirely impractical for spontaneous speech or shadowing sce-
180 narios [35–37], restricting the real-world applicability of MDD.

181 3. Consistency Regularization on Optimal 182 Temporal Transport Classification

183 In this section, we analyze the inherent limitations of CTC and
184 introduce the architecture of our proposed frame-wise dense
185 acoustic model: **CROTTTC**. Our approach, as illustrated in the
186 left panel of Figure 1, is built upon two core pillars: Consistency
187 Regularization (**CR**) and Optimal Temporal Transport Classifi-
188 cation (**OTTC**).

189 3.1. Limitations of Connectionist Temporal Classification

190 Connectionist Temporal Classification (CTC) [6] has es-
191 tablished itself as a fundamental criterion for sequence-to-
192 sequence modeling, particularly in the ASR domain. By in-
193 troducing a blank token and optimizing the marginal proba-
194 bility over all valid alignments, CTC enables end-to-end train-
195 ing without requiring explicit frame-level annotations for target
196 units (e.g., phonemes or subwords). Consequently, the model
197 tends to produce a highly sparse posterior distribution, concen-
198 trating probability mass on a few highly discriminative frames
199 [38]. Moreover, these non-blank emission peaks are frequently
200 delayed, shifting towards the tail end of their corresponding
201 acoustic segments—a phenomenon widely recognized as CTC’s
202 delayed behavior [39, 40].

203 This sparsity and delayed behavior pose a significant chal-
204 lenge for precise segment detection tasks. In practice, a mis-
205 pronounced phoneme in MDD does not necessarily manifest
206 as a complete substitution starting from the very first frame;
207 rather, it often involves co-articulation errors or a partial blend
208 of canonical and abnormal acoustic traits. To illustrate, con-
209 sider a speaker attempting to produce the diphthong /aɪ/ (as in
210 “buy”) but erroneously substituting the initial vowel element to
211 produce /ɔɪ/ (as in “boy”). Acoustically, both diphthongs share

²Further detailed in Sec. 3.1

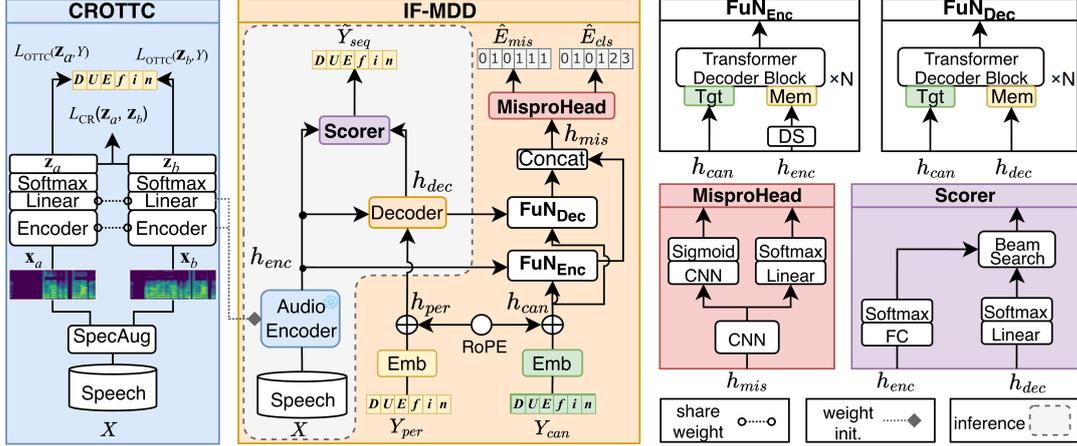


Figure 1: Overview of the CROTTIC-IF architecture. From left to right, it comprises the CROTTIC AM (blue), the lightweight IF-MDD LM (orange), and the detailed IF components.

212 an identical trailing off-glide (/1/). Driven by CTC’s inherent
 213 delayed emission tendency, the model will typically align the
 214 probability peak for the target label (/aɪ/) exclusively with this
 215 shared off-glide. Consequently, the preceding frames contain-
 216 ing the crucial evidence of the mispronunciation (/ɔ/) can be
 217 masked by blank tokens and entirely ignored. While this highly
 218 discriminative frame selection is beneficial for standard ASR,
 219 it discards the fine-grained acoustic cues that are essential for
 220 sensitive diagnostic tasks like MDD.

221 3.2. Optimal Temporal Transport Classification

222 To address the inherent shortcomings of CTC, we leverage Op-
 223 timal Temporal Transport Classification (OTTC) [41], a differ-
 224 entiable framework rooted in one-dimensional optimal transport
 225 that enables dense, frame-level alignments.

226 3.2.1. Sequence Alignment via Optimal Transport

227 Let $\mathbf{X} = \{x_i\}_{i=1}^n$ denote the acoustic frame sequence and
 228 $\mathbf{Y} = \{y_j\}_{j=1}^m$ the perceived phoneme sequence. We define
 229 that the alignment between \mathbf{X} and \mathbf{Y} can be solved with an
 230 optimal transport plan $\gamma \in \mathbb{R}_+^{n \times m}$, which serves as a mono-
 231 tonic mapping governing the frame-to-label correspondence. In
 232 optimal transport theory, γ is parameterized by two discrete
 233 distributions: the frame weights $\alpha \in \Delta^n$ derived from \mathbf{X} ,
 234 and the label weights $\beta \in \Delta^m$ associated with \mathbf{Y} , where
 235 $\Delta^k = \{v \in \mathbb{R}_+^k \mid \sum_{i=1}^k v_i = 1\}$ is the probability simplex.
 236 Here, the frame weights α can be predicted from the input via a
 237 neural network W :

$$\alpha[\mathbf{X}, W] = \text{softmax}(W(x_1), \dots, W(x_n))^\top, \quad (1)$$

238 and β is typically fixed as a uniform distribution to evenly as-
 239 sign each target token’s contribution. Subsequently, the align-
 240 ment plan $\gamma(\alpha, \beta)$ is computed as the solution to the 1D opti-
 241 mal transport problem:

$$\gamma(\alpha, \beta) = \arg \min_{\gamma \in \Gamma_{\alpha, \beta}} \sum_{i=1}^n \sum_{j=1}^m \gamma_{i,j} \|i - j\|_2^2, \quad (2)$$

242 where $\Gamma_{\alpha, \beta} = \{\gamma \in \mathbb{R}_+^{n \times m} \mid \gamma \mathbf{1}_m = \alpha, \gamma^\top \mathbf{1}_n = \beta\}$ denotes
 243 the set of valid couplings preserving the marginal distributions.
 244 Based on the *Sequence Optimal Transport Distance (SOTD)*
 245 framework [41], which defines a pseudo-metric over sequence
 246 space, we derive the OTTC loss as the expected transport cost
 247 under the cross-entropy metric:

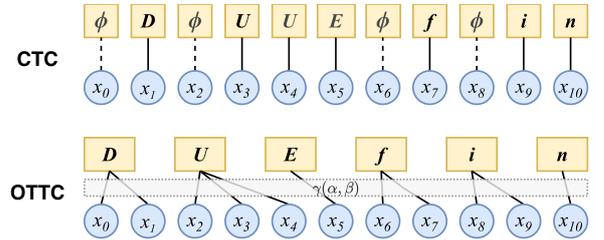


Figure 2: Comparison between CTC and OTTC, where ϕ is the blank token and $\gamma(\alpha, \beta)$ is the optimal transport plan.

$$\mathcal{L}_{\text{OTTC}} = - \sum_{i=1}^n \sum_{j=1}^m \gamma_{i,j}(\alpha, \beta) \cdot \log p(y_j | x_i), \quad (3)$$

248 where $p(y_j | x_i)$ is the posterior probability. Since α is pa-
 249 rameterized via a neural network W (Eq. 1) and β is fixed, the
 250 model parameters can be optimized in an end-to-end manner by
 251 minimizing $\mathcal{L}_{\text{OTTC}}$.

252 **Comparison to CTC.** Fig. 2 illustrates the fundamental dis-
 253 tinction in alignment strategies. Specifically, CTC maximizes
 254 the marginal probability over *all* valid monotonic paths:

$$\mathcal{L}_{\text{CTC}} = - \log \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{Y})} \prod_{i=1}^n p(\pi_i | x_i), \quad (4)$$

255 where $\mathcal{B}^{-1}(\mathbf{Y})$ denotes the set of all paths collapsed to \mathbf{Y} . In
 256 contrast, OTTC maximizes the probability of a *single* opti-
 257 mal monotonic path $\gamma(\alpha, \beta)$ via Eq. (2) and (3). This formu-
 258 lation yields dense, frame-level alignments without blank-token do-
 259 minance, as γ explicitly models the frame-to-label correspondence
 260 rather than marginalizing over sparse paths.

261 3.3. Consistency Regularization

262 Since OTTC takes all frames into consideration, acoustic mod-
 263 els can be over-sensitive to local variations, leading to spurious
 264 insertions. To mitigate this, we also introduce *Consistency Reg-
 265 ularization (CR)* into our training pipeline [42]. Specifically, as
 266 shown in Fig 1, for each input utterance \mathbf{X} , we generate two
 267 augmented views \mathbf{X}_a and \mathbf{X}_b using stochastic perturbations³.
 268 Let \mathbf{Z}_a and \mathbf{Z}_b denote the posterior probability distributions
 269 output by the model for each view, where $\mathbf{Z}^{(i)} \in \Delta^K$ repre-
 270 sents the probability vector at frame i (equivalent to $p(y | x_i)$)

³See Sec. 6.1 for spectrogram augmentation details.

271 in Eq. 3). The consistency regularization loss is defined as
 272 the symmetric Kullback-Leibler (KL) divergence between these
 273 distributions:

$$\mathcal{L}_{CR} = \frac{1}{2n} \sum_{i=1}^n \left(\text{KL}(\mathbf{Z}_a^{(i)} \parallel \mathbf{Z}_b^{(i)}) + \text{KL}(\mathbf{Z}_b^{(i)} \parallel \mathbf{Z}_a^{(i)}) \right), \quad (5)$$

274 where n is the number of frames. By minimizing the distribu-
 275 tion distance and maximizing the mutual information between
 276 the two branches, CR enables each branch to receive supervi-
 277 sion from the other, which encourages the shared encoder to
 278 perform self-distillation on frame-level posterior distributions.
 279 Additionally, relatively long-term time masking implicitly en-
 280 courages *masked language modeling*-like behavior: the model
 281 learns to reconstruct occluded acoustic frames by leveraging
 282 contextual information from unmasked regions. Together, these
 283 mechanisms stabilize frame-level predictions and reduce sensi-
 284 tivity to local acoustic noise. Consequently, the training criteria
 285 for CROTTTC is defined as:

$$\mathcal{L}_{AM} = \mathcal{L}_{CR} + \eta (\mathcal{L}_{\text{OTTC}}(\mathbf{Z}_a, \mathbf{Y}) + \mathcal{L}_{\text{OTTC}}(\mathbf{Z}_b, \mathbf{Y})), \quad (6)$$

286 where hyperparameter η is set to 1.0.

287 4. Indirect Fusion of Mispronunciation 288 Information via Knowledge Transfer

289 In this section, we detail our Indirect Fusion (IF) strategy for in-
 290 corporating mispronunciation information into LM, named **IF-**
 291 **MDD**. As illustrated in the middle and right panels of Fig. 1,
 292 IF-MDD treats canonical phonemes and mispronunciation cues
 293 as privileged information available exclusively during the train-
 294 ing phase. These cues are leveraged to guide the model’s latent
 295 representations via backpropagation.

296 4.1. Learning using Privilege Information

297 Learning using Privileged Information (LUPI) is a training
 298 paradigm that fundamentally operates as a knowledge transfer
 299 mechanism, leveraging additional information available exclu-
 300 sively during the training phase to implicitly guide the model
 301 [43, 44]. While this paradigm has already been successfully
 302 applied to various speech-related tasks [45–47], its most com-
 303 pelling advantage lies in its effectiveness when training data is
 304 limited or labels are scarce [48–50]. Given the inherent data
 305 scarcity of expert-annotated MDD corpora, adopting the LUPI
 306 framework is a natural and highly principled choice.

307 4.2. Indirect Fusion of Mispronunciation Information via 308 Knowledge Transfer

309 As illustrated in Fig. 1, IF-MDD consists of two main com-
 310 ponents: 1) an encoder–decoder backbone; 2) an auxiliary
 311 mispronunciation-aware teacher network.

312 4.2.1. Encoder-decoder backbone

313 We employ an encoder–decoder architecture serving as the
 314 primary backbone of IF-MDD. Let the input waveform be
 315 $\mathbf{X} = \{x_i\}_{i=1}^n$ and the perceived phoneme sequence be $\mathbf{Y}_{\text{per}} =$
 316 $\{y_j^{\text{per}}\}_{j=1}^m$, where n and m denote the lengths of the speech sig-
 317 nal and phoneme sequence, respectively. The backbone is for-
 318 mulated as:

$$\mathbf{h}_{\text{enc}} = \text{Enc}(\text{RoPE}(\mathbf{X})), \quad (7)$$

$$\mathbf{h}_{\text{per}} = \text{RoPE}(\mathbf{Y}_{\text{per}}), \quad (8)$$

$$\mathbf{h}_{\text{dec}} = \text{Dec}(\mathbf{h}_{\text{per}}, \mathbf{h}_{\text{enc}}), \quad (9)$$

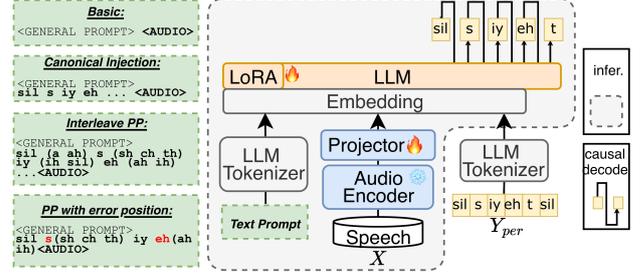


Figure 3: Illustration of LLM-based MDD, with alternative prompts.

where $\text{Enc}(\cdot)$ and $\text{Dec}(\cdot)$ represent the encoder (AM) and decoder (LM), respectively. Rotary position embeddings $\text{RoPE}(\cdot)$ are applied to both $\mathbf{h}_{\text{enc}} \in \mathbb{R}^{n \times D}$ and $\mathbf{h}_{\text{dec}} \in \mathbb{R}^{m \times D}$ to enhance positional modeling [51].

323 4.2.2. Auxiliary mispronunciation-detection teacher network

324 To explicitly model the mispronunciation information indicated
 325 by the misalignment between acoustic and canonical features,
 326 we introduce an auxiliary teacher network. This module lever-
 327 ages the canonical phoneme sequence $\mathbf{Y}_{\text{can}} = \{y_j^{\text{can}}\}_{j=1}^m$ and
 328 the corresponding phoneme-level ground-truth error sequence
 329 $\mathbf{E} = \{e_j\}_{j=1}^m$ as privileged information. It consists of two se-
 330 quential stages: feature fusion and comprehensive error detec-
 331 tion.

332 First, two parallel fusion networks (FuN) are employed to
 333 model multi-aspect feature interactions:

$$\mathbf{h}_{\text{can}} = \text{RoPE}(\mathbf{Y}_{\text{can}}), \quad (10)$$

$$\mathbf{h}_{\text{mis}}^{\text{enc}} = \text{FuN}_{\text{enc}}(\mathbf{h}_{\text{can}}, \text{DS}(\mathbf{h}_{\text{enc}})), \quad (11)$$

$$\mathbf{h}_{\text{mis}}^{\text{dec}} = \text{FuN}_{\text{dec}}(\mathbf{h}_{\text{can}}, \mathbf{h}_{\text{dec}}), \quad (12)$$

$$\mathbf{h}_{\text{mis}} = \text{Concat}(\mathbf{h}_{\text{mis}}^{\text{enc}}, \mathbf{h}_{\text{mis}}^{\text{dec}}), \quad (13)$$

334 where $\text{FuN}(\cdot)$ is composed of a transformer decoder structure
 335 and $\mathbf{h}_{\text{mis}} \in \mathbb{R}^{m \times D}$ denotes the fused representation. In both
 336 branches, $\text{FuN}(\cdot)$ uses \mathbf{h}_{can} as the query and \mathbf{h}_{enc} or \mathbf{h}_{dec}
 337 as the memory. Notably, a downsampler $\text{DS}(\cdot)$ is applied within
 338 $\text{FuN}_{\text{enc}}(\cdot)$ to improve alignment.

339 Subsequently, our teacher network employs a dual-path
 340 detection sub-network to provide a comprehensive diagnosis.
 341 Specifically, we apply a shared CNN trunk on the fused repre-
 342 sentation \mathbf{h}_{mis} , followed by two task-specific projection heads:

$$\mathbf{u} = \text{CNN}(\mathbf{h}_{\text{mis}}), \quad (14)$$

$$\hat{E}_{\text{mis}} = \sigma(\text{CNN}_{\text{bin}}(\mathbf{u})), \quad (15)$$

$$\hat{E}_{\text{cls}} = \text{Softmax}(\text{Linear}_{\text{cls}}(\mathbf{u})), \quad (16)$$

343 where \hat{E}_{mis} and \hat{E}_{cls} are the predicted counterparts for the posi-
 344 tion and classification facets of \mathbf{E} , respectively. Here, $\hat{E}_{\text{mis}} \in$
 345 $[0, 1]^m$ indicates the probability of mispronunciation, while \hat{E}_{cls}
 346 denotes the distribution over $\{\text{correct}, \text{substitution}, \text{deletion}, \text{in-}$
 347 $\text{sertion}\}$.

348 Crucially, this architectural design naturally facilitates
 349 knowledge transfer. Since the auxiliary mispronunciation heads
 350 are relatively lightweight and receive direct supervision from \mathbf{E} ,
 351 they converge significantly faster than the primary sequence-
 352 to-sequence backbone. During training, this fast-converging
 353 teacher provides a strong, discriminative gradient signal that
 354 backpropagates through the fusion networks into both \mathbf{h}_{enc} and
 355 \mathbf{h}_{dec} . Consequently, the diagnostic features from the teacher

network act as a catalyst, helping the deeper backbone to converge effectively and internalize the mispronunciation-aware representations required for prompt-free inference.

4.3. Inference and Training Criteria

During inference, the auxiliary teacher is discarded. A scorer searches for the optimal hypothesis $\hat{\mathbf{Y}}$ by maximizing the interpolated log-probabilities between the AM and LM:

$$\hat{\mathbf{Y}} = \arg \max_{\mathbf{Y}} \left\{ \lambda \log P_{\text{AM}}(\mathbf{Y}|\mathbf{h}_{\text{enc}}) + (1 - \lambda) \log P_{\text{LM}}(\mathbf{Y}|\mathbf{h}_{\text{enc}}) \right\}, \quad (17)$$

where $\lambda \in [0, 1]$ is a AM decoding weight. Here, P_{AM} represents the acoustic probability predicted by the CROTTTC branch, and P_{LM} is the contextual probability evaluated autoregressively by Dec(\cdot). The IF-MDD model is optimized via a multi-task objective:

$$\mathcal{L}_{\text{total}} = \omega_1 \mathcal{L}_{\text{AM}} + (1 - \omega_1) \mathcal{L}_{\text{LM}} + \omega_2 (\mathcal{L}_{\text{pos}} + \mathcal{L}_{\text{type}}) + \omega_3 \mathcal{L}_{\text{ga}}, \quad (18)$$

where \mathcal{L}_{ga} is the guided-attention loss [52] ensuring monotonic fusion in FuN. Hyperparameters $\omega_1, \omega_2, \omega_3$ are set to 0.3, 1.0, and 10.0, respectively.

5. Leveraging LLMs to Investigate Canonical Information Effect on MDD

To empirically investigate the impact of explicit canonical information and strong linguistic priors, this section introduces an LLM-based architecture, denoted as **LLM-MDD**. As illustrated in Fig. 3, we replace the conventional lightweight LM with an open-source LLM, which serves as a high-capacity linguistic processor for the acoustic embeddings. We adopt an LLM for this canonical information analysis for two primary reasons. First, LLMs inherently possess massive language modeling capabilities; therefore, leveraging an LLM to replace a lightweight LM itself represents a significant amplification of linguistic priors. Second, multimodal LLMs natively support the direct integration of text (canonical prompts) and speech (acoustic embeddings). This allows us to isolate and strictly verify the exact impact of various canonical injection strategies, while entirely bypassing complex and task-specific fusion mechanisms.

5.1. LLM-MDD Framework

Inspired by the SLAM-LLM architecture [28], our LLM-MDD system formulates the MDD task as mapping acoustic features and canonical cues to a perceived phoneme sequence. The architecture comprises three primary components: a pretrained AM, a trainable projector, and a backbone LLM.

During training, the AM remains frozen. We optimize the LLM via Low-Rank Adaptation (LoRA) [53] and fully fine-tune the projector to align acoustic features to the LLM’s embedding space. The system is optimized using standard cross-entropy loss against ground-truth perceived phoneme sequences. By injecting these massive linguistic priors, we aim to evaluate whether the emergent capabilities of LLMs can genuinely outperform traditional lightweight decoders in diagnosing subtle mispronunciations.

5.2. Canonical Information Injection via Prompting

To explore how explicit canonical information influences the model’s diagnostic accuracy, we design four distinct prompting strategies for LLM-MDD:

Basic: A standard SLAM-style sequence where the speech features are followed by the perceived phonemes: [Prompt]<speech><BOS><perc_phns><EOS>.

Canonical Injection: Following the implementation in [30], the canonical phoneme sequence is appended to a general prompt: [Prompt]<cano_phns><speech><BOS><perc_phns><EOS>.

Interleaved Potential Pronunciation: Following the implementation in [31, 32], we interleave canonical phonemes with potential pronunciation (PP) alternatives: [Prompt]<cano_phns_with_PP><speech><BOS><perc_phns><EOS>. In both the training and inference stages, these PP candidates are uniformly assigned to each canonical phoneme. This design ensures that the model receives no prior clues regarding the specific *location* or *type* of error.

Potential Pronunciation with Error Position: Unlike the uniform assignment, we only provide PP candidates for phonemes explicitly annotated as mispronounced: [Prompt]<cano_phns_with_PP_err_pos><speech><BOS><perc_phns><EOS>. Since the mispronunciation position is pinpointed, the task reduces to a forced-choice mispronunciation diagnosis. Although this non-causal prompting style is impractical for real-world applications, this configuration is designed to test the performance upper bound of LLMs—specifically, whether they prioritize the canonical or the acoustic features in MDD.

6. Experimental Evaluation

This section comprehensively evaluates our CROTTTC-IF framework against conventional dictation- and text-prompting-style baselines, alongside an empirical analysis of canonical priors using LLM-MDD.

6.1. Datasets

Table 1 summarizes the L2 English datasets utilized in our experiments. The most widely recognized among these is L2-ARCTIC, which standardized test set with six speakers (NJS, TLV, TNI, TXHC, YKWK, and ZHAA). Besides, we incorporated two supplementary datasets: Speechocean762 (SO762) and ERJ. SO762 was originally designed for pronunciation assessment, assigning each phoneme a score ranging from 0.0 to 2.0. Phonemes scoring below 0.5 are identified as mispronunciations, and their corresponding expert-annotated phonetic realizations are provided as ground truth. ERJ is a Japanese-accented L2-English corpus that serves as an out-of-domain (OOD) dataset to evaluate generalization ability. For vocabulary preparation, we adopted a 44-unit ARPAbet-style inventory (39 phonemes plus 5 special tokens <bos>, <eos>, <sil>, <error>, and <blank>).

6.2. Implementation Details

Acoustic Model (AM). For the AM backbone described in Sec. 3, we utilize WavLM Large [13] followed by a 2-layer Conformer [54] with a kernel and stride size of 3. The hidden dimension is set to 384. For the Consistency Regularization, we apply time warping with a factor of 80 alongside time and frequency masking following [42]. Specifically, we use a maximum of 3 mask blocks with a masking ratio $\rho \in (0.1, 0.3]$. In the time domain, the minimum masking length is enforced to ensure the CR mechanism leverages neighboring context rather than acting as simple data augmentation.

Language Model (LM). As introduced in Sec. 4, our LM comprises a 2-layer Transformer decoder with a hidden size

Table 1: Overview of the L2-ARCTIC, ERJ, and speechocean762 datasets.

Dataset	L1	Volume	Train Utt.	Val Utt.	Test Utt.	Test Spk.
L2-ARCTIC [57]	Mixed ⁴	3.66h	2429	268	900	6
SO762 [58]	Mandarin	5.58h	2250	250	2500	125
ERJ [59]	Japanese ⁵	0.99h	674	74	144	8

of 384. We employ a two-stage training strategy: (1) initial joint CTC/Transformer training, followed by (2) substituting the CTC AM with the pretrained CROTTTC AM and fine-tuning the LM until convergence. This approach resolves convergence issues arising from the absence of blank tokens in CROTTTC, which otherwise provide crucial implicit segmental boundaries. The AM/LM training weight ω_1 is set to 0.5, and the effect of the decoding weight is analyzed in Sec. 6.4.4.

Auxiliary Teacher Network. The teacher network in Sec. 4.2.2 comprises two parallel branches: an encoder-side $\text{FuN}_{\text{enc}}(\cdot)$ and a decoder-side $\text{FuN}_{\text{dec}}(\cdot)$, both implemented with 2-layer Transformer decoders with a hidden size of 384. Notably, the 1D CNN down-sampler (factor=4) is integrated into $\text{FuN}_{\text{enc}}(\cdot)$ for alignment improvement. The shared CNN trunk from the mispronunciation head has a dimension of 128, followed by a 64-dim CNN branch for error-position detection and a linear branch for error type detection, respectively.

LLM-MDD. For the modules mentioned in Sec. 5, we adapted AMs fine-tuned with both CTC and CROTTTC. The projector is a 1D-CNN with a down-sampling factor of 4. We employed LLaMA-3.2-1B-Instruct [55] and Qwen3-4B-Thinking-2507 [56] as backbone LLMs, with rank=16 for LoRA fine-tuning. For the potential pronunciation prompting, we conducted a corpus-level statistical analysis, where the top 5 alternatives were selected as the PP candidates.

Miscellaneous. The AMs were trained for 300 epochs with a batch size of 32; learning rates were set to 10^{-5} for the SSL backbone and 3×10^{-4} for the following modules. The LM and Teacher Network were trained for 200 epochs (batch size 32, learning rate 2×10^{-4}). LLM training was limited to 20 epochs (batch size 4, learning rate 2×10^{-4}). During inference, both the LM and LLMs utilized beam search decoding with a beam size of 10 and a temperature of 1.1. All experiments were conducted on a single NVIDIA GH200 GPU. Model selection was determined by the optimal F1-score on the validation set.

6.3. Evaluation metrics

We followed the hierarchical evaluation protocol of [61]. Detection was measured by true acceptance (TA), true rejection (TR), false acceptance (FA), and false rejection (FR), while diagnosis was evaluated with correct diagnosis (CD) and error diagnosis (ED) as illustrated in Table 2. The main performance is then summarized with the F1 score, treating TR as the positive class:

$$P = \frac{TR}{TR + FR}, \quad R = \frac{TR}{TR + FA}, \quad F1 = \frac{2PR}{P + R}.$$

We also report phoneme error rate (PER) and Correction Rate (COR) as a measure of recognition accuracy:

$$\text{PER} = \frac{S + D + I}{N}, \quad \text{COR} = 1 - \frac{S + D}{N},$$

where S , D , and I denote the numbers of substitutions, deletions, and insertions, and N is the total number of **perceived phonemes**.

⁴Hindi, Korean, Mandarin, Spanish, Arabic, and Vietnamese.

⁵Only expert-annotated speech were utilized as described in [60].

Table 2: An illustration of MDD metrics.

Canonical phoneme	s	p	iy	k	t
Perceived phoneme	<u>s</u>	<u>b</u>	<u>iy</u>	<u>g</u>	<u>d</u>
Predicted phoneme	s	p	ih	g	th
Evaluation result	TA	FA	FR	TR CD	TR ED

Table 3: Comparison of the proposed methods against various baseline models on L2-ARCTIC. **Bold** and underlined values indicate the best and second-best results among fair-comparison models, respectively. \uparrow/\downarrow : higher/lower is better.

Model	Detection			Diagnosis			Recognition	
	F1 \uparrow	P \uparrow	R \uparrow	FRR \downarrow	FAR \downarrow	EDR \downarrow	PER \downarrow	COR \uparrow
<i>Dictation-Style Baselines</i>								
MPL-MDD [15]	55.42	60.39	51.20	5.60	48.80	22.71	14.36	-
RNN-T [62]	59.10	63.40	55.30	5.30	44.70	-	15.47	-
MV-w2v2 [16]	60.31	59.23	61.43	-	-	-	14.13	-
w2v2-CTC [14]	60.44	62.86	58.57	5.70	41.80	29.28	16.20	-
Meta-Learn [63]	61.45	91.60	46.24	29.75	8.40	-	42.25	-
<i>Text-prompting-Style Baselines</i>								
Qwen2 [34]	50.60	71.62	39.12	-	-	-	-	-
Qwen2-sub [34]	55.00	63.27	48.64	-	-	-	-	-
AEL w/o Pos. [64]	56.33	58.36	55.00	6.55	45.00	25.72	14.81	-
MDDGCN [26]	56.49	51.90	61.97	9.18	38.03	25.27	-	-
Peppanet [27]	56.81	64.53	51.38	9.61	36.47	25.88	-	-
PG-MDD [17]	60.10	60.06	60.15	6.65	39.94	23.13	<u>13.92</u>	-
TG+Contrast. [24]	61.75	62.12	61.38	6.19	38.62	28.92	-	-
<i>Frame-wise MDD Baselines</i>								
Joint-align [21]	63.04	<u>77.12</u>	53.31	-	-	-	-	-
PER-MDD [22]	<u>69.60</u>	71.78	67.56	4.43	<u>32.44</u>	37.77	104.08	<u>90.42</u>
<i>Proposed Methods (AM)</i>								
OTTC	63.18	66.36	60.29	5.14	39.71	22.12	18.07	89.96
CROTTTC	62.39	69.70	56.47	4.13	43.53	22.06	17.48	90.29
<i>Proposed Methods (LM & LLM)</i>								
CTC-IF	58.37	61.81	55.29	5.75	44.71	19.98	13.72	88.34
CROTTTC-IF	71.77	76.94	<u>67.24</u>	3.39	32.76	27.47	46.52	92.42
CROTTTC-LLaMA	56.87	54.81	59.08	8.20	40.92	21.98	15.85	86.55
CROTTTC-Qwen	55.19	58.00	52.64	6.42	47.36	23.80	15.42	86.78

6.4. Experimental evaluation on L2-ARCTIC

Table 3 compares our proposed framework against recent representative MDD methods on the L2-ARCTIC benchmark. Notably, CROTTTC-IF achieves a peak F1 score of 71.77%, yielding the lowest FRR and highest COR among all baselines. In the following subsections, we systematically analyze how the AM, LM, and canonical injection strategies collectively drive this superior performance.

6.4.1. Escaping the Acoustic Trap: Ablation on CROTTTC

Table 4(a) details the performance differences between CTC- and OTTC-based AMs. Compared to standard CTC, CROTTTC yields improvements in both F1 score and PER. For OT-based methods, the F1 score exhibits a significant improvement, increasing from 57.89% to 63.18%. Concurrently, the PER also rises from 14.45% to 18.35%, with the most notable change being a surge in insertion errors from 2.21% to 8.31%. The integration of CR effectively mitigates these insertion errors down to 7.76% and reduces the FRR to 4.13%. In MDD, FRR is a critical usability metric; incorrectly flagging a student’s correct pronunciation as an error can lead to learner frustration and undermine the pedagogical value of the system. Therefore, we select CROTTTC as our primary AM, albeit with a slight trade-off in the overall F1 score.

To understand these acoustic mechanisms, we illustrate the frame-wise probability distributions of argmax predictions

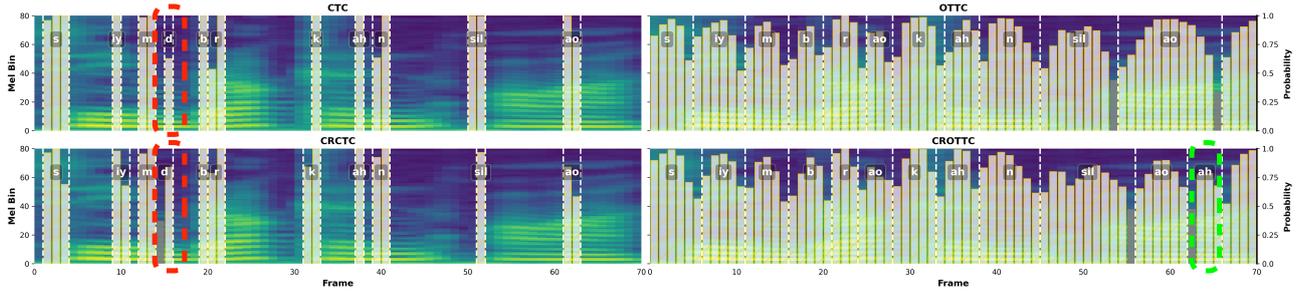


Figure 4: Comparison of frame-level probability distributions across different AMs (blank tokens omitted). The perceived phoneme sequence is */s iy m (d) b r aa k ah n sil ao ah/*, corresponding to the utterance “seemed broken or” → “seem broken or-ah”. While the CTC-based method hallucinated the */d/* and omitted the */ah/*, our CROTTTC faithfully captures the actual phonetic realization.

Table 4: Comprehensive ablation studies on L2-ARCTIC. (a) Acoustic model variations. (b) Indirect Fusion components. (c) Canonical prompting strategies on LLM-MDD.

(a) Ablation of CROTTTC AM								(b) Ablation of IF LM					(c) Ablation of LLM Canonical Prompts								
Method	F1↑	P↑	R↑	FRR↓	FAR↓	PER↓	INS↓	Method	F1↑	P↑	R↑	FAR↓	EDR↓	Method	F1↑	P↑	R↑	FRR↓	FAR↓	EDR↓	PER↓
CTC	57.89	59.40	56.45	6.49	43.55	14.45	2.21	CTC-IF	58.37	61.81	55.29	44.71	19.98	CTC-LLaMA	55.16	54.76	55.57	7.73	44.43	24.17	16.23
CRCTC	58.45	61.40	55.78	5.90	44.22	14.01	2.05	w/o FuN _{Enc}	57.33	60.29	54.65	45.35	21.74	CROTTTC-LLaMA	56.87	54.81	59.08	8.20	40.92	21.98	15.85
OTTC	63.18	66.36	60.29	5.14	39.71	18.35	8.31	w/o FuN _{Dec}	56.51	59.42	53.86	46.14	22.10	w/ cano.	40.52	68.22	28.83	2.26	71.17	32.56	13.55
CROTTTC	62.39	69.70	56.47	4.13	43.53	17.48	7.76	w/o Error Head	57.56	59.54	55.71	44.29	22.28	w/ PP	42.63	54.11	35.18	<u>5.02</u>	64.82	35.15	<u>14.91</u>
								CTC-LM	54.95	59.23	51.25	48.75	22.60	w/ pos. (oracle)	91.78	95.02	88.74	0.78	11.16	24.72	5.04

539 across different loss functions in Fig. 4. Unlike CTC, OT-based
 540 methods assign non-blank tokens to every frame, which in-
 541 creases inference sensitivity but also leads to higher insertion
 542 errors during greedy decoding. By incorporating CR, the model
 543 effectively captures more subtle pronunciation deviations, suc-
 544 cessfully inheriting the alignment flexibility of OTTC while
 545 maintaining the recognition stability.

546 6.4.2. Escaping the Linguistic Trap: Implicit Knowledge 547 Transfer via IF

548 Table 4(b) presents the ablation of the auxiliary teacher net-
 549 work, where CTC served as the AM backbone. Removing either
 550 FuN_{Enc} or FuN_{Dec} leads to a clear drop in the F1 score by 1.04%
 551 to 1.86%, demonstrating the effectiveness of multi-aspect mon-
 552 itoring in our indirect fusion strategy. Ablating the error-type
 553 classification head yields a smaller but consistent decrease in F1
 554 alongside an increase in EDR, suggesting that this multi-view
 555 design helps to regularize the mispronunciation representations.
 556 Visualizations in Fig. 5 display the cross-attention heatmaps of
 557 the fusion network’s last layer: removing these components
 558 produces blurrier, less-peaky patterns with reduced diagonal
 559 dominance, reflecting the degraded representation quality of the
 560 mispronunciation cues. Furthermore, without the teacher net-
 561 work’s knowledge transfer, the model devolves into a standard
 562 CTC/Transformer ASR architecture (CTC-LM), resulting in a
 563 significant decrease in the F1 score. This proves that uncriti-
 564 cally adapting general LM-based ASR models to the standalone
 565 MDD task generates suboptimal diagnostic results.

566 6.4.3. The Danger of Explicit Priors: Canonical Injection with 567 LLMs

568 Table 4(c) presents our investigation into the effect of strong
 569 canonical priors on MDD performance. First, compared to
 570 CTC-LLaMA, CROTTTC-LLaMA yields better detection perfor-
 571 mance, boosting the F1 score from 55.16% to 56.87%. This
 572 supports our hypothesis that the sparsity inherent in CTC causes
 573 MDD models to overlook transient mispronunciation cues. Sec-
 574 ond, by varying the canonical information injected into the

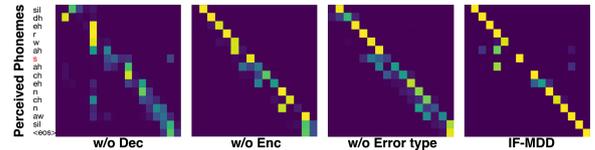


Figure 5: Attention heatmaps of the Fusion Network under different ablation conditions. The x-axis denotes the memory, and the y-axis denotes the canonical-phoneme embeddings.

575 prompt, we observe that explicitly providing canonical targets
 576 is highly detrimental to MDD, resulting in a drastic drop in the
 577 F1 score to 40.52%. While prompting with Potential Pronun-
 578 ciations (PP) offers the LLM slightly more flexibility, the perfor-
 579 mance recovery remains severely limited (F1 = 42.63%), con-
 580 sistent with the observation in [34]. These observations strongly
 581 demonstrate the “Linguistic Trap” in MDD: strong textual pri-
 582 ors bias the LLM toward over-correction, masking objective
 583 acoustic fidelity. Third, when prompting the LLM with both
 584 PP and explicit mispronunciation positions, the F1 score natu-
 585 rally sees a significant improvement to 91.78%. However, a gap
 586 from perfect MDD remains: while FRR drops to 0.78%, EDR
 587 and FAR remain high at 24.72% and 11.16%. This proves the
 588 fundamental bottleneck for multimodal LLMs in MDD is not
 589 merely detecting error positions, but also the inability to utilize
 590 the fine-grained acoustic resolution required for precise diag-
 591 nosis. Even with explicit positional hints, the LLM struggles
 592 to accurately identify specific substitutions and stubbornly de-
 593 faults to the textual prior. Therefore, we believe that if future
 594 research aims to utilize multimodal LLMs for MDD, develop-
 595 ing more refined acoustic processing stages is essential.

596 6.4.4. Balancing Acoustic Fidelity and Linguistic Context

597 As formulated in Eq. 17, CROTTTC-IF performs decoding un-
 598 der a joint AM/LM strategy. In standard CTC/LM architectures
 599 for general speech recognition, a sweet spot for the AM weight
 600 (typically $\lambda \in [0.2, 0.5]$) is usually found to achieve the optimal
 601 Word Error Rate (WER). However, for the standalone MDD
 602 task, we hypothesize that a significantly higher AM weight is
 603 crucial for preserving detection sensitivity. To validate this, we

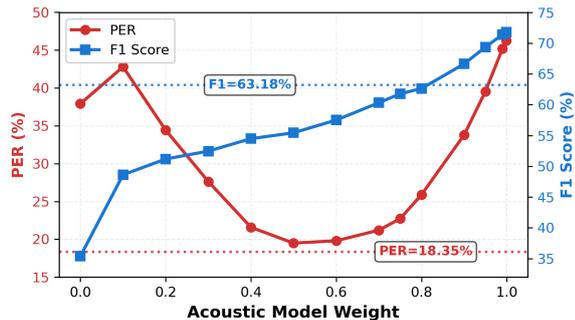


Figure 6: Trade-off between PER and F1 on the L2-arctic over varying w_{AM} .

Table 5: Fine-tuned and zero-shot results on SO762 and ERJ.

Dataset	Cond.	Model	F1 \uparrow	P \uparrow	R \uparrow	FRR \downarrow	FAR \downarrow	EDR \downarrow
SO762	Fine-tuned	Ryu2023 [67]	41.50	26.90	91.60	-	-	-
		JAM [65]	45.01	<u>61.10</u>	34.76	0.58	64.32	<u>45.23</u>
		MuFFIN [66]	67.98	67.60	68.37	<u>1.01</u>	31.63	58.82
	Fine-tuned	CTC-IF	46.68	31.86	87.23	6.32	12.77	50.30
		CROTTTC	51.86	36.81	87.67	5.09	<u>12.32</u>	57.40
		CROTTTC-IF	<u>57.16</u>	41.90	<u>89.92</u>	8.63	6.53	10.08
ERJ	Fine-tuned	CTC-IF	83.98	83.91	84.06	9.70	15.94	<u>26.60</u>
		CROTTTC	<u>85.79</u>	<u>85.79</u>	<u>85.79</u>	<u>8.63</u>	<u>14.21</u>	29.09
		CROTTTC-IF	89.27	89.12	89.43	6.63	10.57	25.78
	Zero-shot	CTC-IF	67.13	79.57	58.06	-	41.94	54.03
		CROTTTC	<u>69.17</u>	<u>83.37</u>	<u>59.10</u>	<u>7.16</u>	<u>40.90</u>	<u>46.92</u>
		CROTTTC-IF	78.44	87.26	71.23	6.32	28.77	44.77

conducted a spectrum-control study on the AM weight λ .

Figure 6 clearly demonstrates this divergence: while the PER reaches its optimum around $\lambda = 0.5$, the F1 score exhibits a completely different trajectory.⁶ When $\lambda = 0$, the F1 drops to 36.83%, proving that relying solely on the LM’s contextual priors leads to mispronunciation detection failure. As λ increases, F1 shows a monotonically increasing trend, confirming our hypothesis that acoustic fidelity is fundamentally more critical than linguistic priors in MDD. Notably, when $\lambda > 0.8$, the joint decoding achieves superior performance compared to CROTTTC AM alone, demonstrating that IF-MDD provides effective and complementary soft linguistic guidance.

6.5. Generalization Study on SO762 and ERJ

Table 5 presents our experimental results on two additional L2 datasets: SO762 and ERJ. For the zero-shot experiments, we directly evaluated the models trained on L2-ARCTIC (as described in Section 6.4), and further adapted to the mentioned datasets for the fine-tuned settings. Compared to conventional baselines on SO762, CROTTTC-IF achieves a highly competitive F1 score of 57.16% with significantly better FAR and EDR, without requiring the explicit fine-grained phonetic score labels as used in [65,66]. Furthermore, CROTTTC-IF delivers outstanding fine-tuned and remarkably strong zero-shot results on ERJ, clearly demonstrating the robust out-of-domain (OOD) generalization of our prompt-free architecture.

7. Iqra’Eval2 Challenge

The Iqra’Eval2 Challenge [2] provides an ideal arena for our framework. Unlike standard L2 English benchmarks, Qur’anic

⁶In MDD, PER is considered secondary to F1 and FRR [9,22]. High PER often reflects redundant insertions that do not degrade diagnostic accuracy, as evidenced by our robust COR at 92.42%.

Table 6: Overview of the Iqra’Eval2 datasets.

Dataset	Volume	Error Type	Description
Iqra	~79h	None	Native speakers; golden reference
TTS	~80h	Synthetic	TTS-injected mispronunciations
Extra	~2h	Real errors	Real-world human mispronunciations

Table 7: Performance of CROTTTC-IF on the Iqra’Eval2 leaderboard.

Model	F1 \uparrow	Pre. \uparrow	Rec. \uparrow	PER \downarrow
3rd-team	71.57	67.69	75.93	4.05
CROTTTC	68.77	70.07	67.52	4.11
w/ IF ($w_{AM} = 0.3$)	70.72	72.67	68.89	3.82
w/ IF ($w_{AM} = 0.9$)	71.70	73.25	70.20	3.72
1st-team	72.01	74.16	69.98	3.65

recitation is governed by *Tajweed* rules⁷, demanding precise modeling of nuanced articulatory deviations. Based on the empirical findings from our L2 English experiments, we deployed our CROTTTC-IF architecture for this challenge.

Table 6 summarizes the challenge’s datasets. Among the provided corpora, we exclusively utilized the TTS and extra datasets. Specifically, we applied the TTS corpus for the AM’s baseline training, followed by fine-tuning on the extra corpus, which contains genuine human pronunciation errors. For the extra corpus, we randomly allocated 10% of the utterances as a local test set, with the remainder used for training and validation. For vocabulary preparation, we followed [68], utilizing an inventory of 67 Arabic phonemes⁸ plus 4 special tokens (<bos>, <eos>, <sil>, and <blank>), resulting in a vocabulary size of 71. Other training configurations remain consistent with those detailed in Sec 6.2. During inference, we applied a high acoustic weight ($w_{AM} = 0.9$) during joint decoding as discussed in Sec 6.4.4. As shown in Table 7, this principled approach yielded outstanding results. Operating entirely without explicit canonical prompts, our CROTTTC-IF system ranked 2nd on the official Iqra’Eval2 leaderboard, achieving a highly competitive F1-score of 71.70% and a remarkably low PER of 3.72%. These results further demonstrate our core theory: decoupling acoustic modeling from explicit canonical priors can provide a robust and objective paradigm for real-world MDD tasks.

8. Conclusion and Future Work

In this paper, we proposed **CROTTTC-IF**, a canonical prompt-free framework designed to overcome the inherent acoustic and linguistic bottlenecks in MDD. By integrating a dedicated acoustic frontend with an implicit knowledge transfer language model, our architecture achieves robust, state-of-the-art performance across diverse scenarios. Furthermore, through comprehensive LLM prompting experiments, we empirically demonstrated that uncritically injecting explicit canonical priors severely compromises acoustic fidelity. This highlights the critical need for careful and principled designs when incorporating canonical information. By successfully decoupling acoustic modeling from explicit textual priors, we believe this prompt-free paradigm establishes a highly objective and robust foundation for the MDD field. Future work will focus on further optimizing the detection-recognition trade-off and extending the framework to other CAPT applications.

⁷A strict set of phonetic rules governing Qur’anic recitation, distinct from Modern Standard Arabic (MSA).

⁸We omitted “<<”, the geminated glottal stop, as it is not represented in our selected training data.

9. Generative AI Use Disclosure

Generative AI technology was employed strictly for minor grammatical corrections and stylistic improvements during the drafting process. The authors have verified all content and retain full responsibility for the accuracy and originality of this work.

10. References

- [1] Y. El Kheir, A. Meghanani, H. O. Toyin *et al.*, “Iqra’eval: A shared task on qur’anic pronunciation assessment,” in *ArabicNLP*, 2025, pp. 443–452.
- [2] Y. E. Kheir, O. Ibrahim, A. Meghanani *et al.*, “Towards a unified benchmark for arabic pronunciation assessment: Quranic recitation as case study,” 2025.
- [3] S. M. Witt and S. J. Young, “Phone-level pronunciation scoring and assessment for interactive language learning,” *Speech communication*, vol. 30, no. 2-3, pp. 95–108, 2000.
- [4] A. M. Harrison, W.-K. Lo, X. Qian *et al.*, “Implementation of an extended recognition network for mispronunciation detection and diagnosis in computer-assisted pronunciation training,” in *SLaTE*, 2009, pp. 45–48.
- [5] K. Li, X. Qian, and H. Meng, “Mispronunciation detection and diagnosis in l2 english speech using multidistribution deep neural networks,” *IEEE/ACM TASLP*, vol. 25, no. 1, pp. 193–207, 2016.
- [6] A. Graves, S. Fernández, F. Gomez *et al.*, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *ICML*, 2006, pp. 369–376.
- [7] S. Watanabe, T. Hori, S. Kim *et al.*, “Hybrid CTC/attention architecture for end-to-end speech recognition,” *IEEE JSTSP*, vol. 11, no. 8, pp. 1240–1253, 2017.
- [8] A. Graves, “Sequence transduction with recurrent neural networks,” *ICML*, 2012.
- [9] W.-K. Leung, X. Liu, and H. Meng, “CNN-RNN-CTC based end-to-end mispronunciation detection and diagnosis,” in *ICASSP*, 2019, pp. 8132–8136.
- [10] A. Baevski, H. Zhou, A. Mohamed *et al.*, “wav2vec 2.0: a framework for self-supervised learning of speech representations,” in *NeurIPS*, 2020.
- [11] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai *et al.*, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM TASLP*, vol. 29, pp. 3451–3460, 2021.
- [12] M. Z. Boito, V. Iyer, N. Lagos *et al.*, “mHuBERT-147: A Compact Multilingual HuBERT Model,” in *Interspeech 2024*, 2024.
- [13] S. Chen, C. Wang, Z. Chen *et al.*, “WavLM: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE JSTSP*, vol. 16, no. 6, pp. 1505–1518, 2022.
- [14] L. Peng, K. Fu, B. Lin *et al.*, “A study on fine-tuning wav2vec2.0 model for the task of mispronunciation detection and diagnosis,” in *Interspeech*, 2021, pp. 4448–4452.
- [15] M. Yang, K. Hirschi, S. D. Looney *et al.*, “Improving mispronunciation detection with wav2vec2-based momentum pseudo-labeling for accentedness and intelligibility assessment,” in *Interspeech*, 2022, pp. 4481–4485.
- [16] Y. EL Kheir, S. Chowdhury, and A. Ali, “Multi-view multi-task representation learning for mispronunciation detection,” in *SLaTE 2023*, 2023, pp. 86–90.
- [17] M.-S. Lin, B.-C. Yan, T.-H. Lo *et al.*, “Pg-mdd: Prompt-guided mispronunciation detection and diagnosis leveraging articulatory features,” in *APSIPA ASC*. IEEE, 2024, pp. 1–6.
- [18] A. Ibrahim, “Hafs2Vec: A system for the IqraEval Arabic and qur’anic phoneme-level pronunciation assessment,” in *ArabicNLP*, Nov. 2025, pp. 453–456.
- [19] A. Mansour, “Metapseud on iqra’eval: Domain adaptation with multi-stage fine-tuning for phoneme-level qur’anic mispronunciation detection,” in *ArabicNLP*, Nov. 2025, pp. 475–479.
- [20] Y. Feng, G. Fu, Q. Chen *et al.*, “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)*, 2020, pp. 3492–3496.
- [21] B. Lin and L. Wang, “Phoneme mispronunciation detection by jointly learning to align,” in *ICASSP*, 2022, pp. 6822–6826.
- [22] H. T. Tu, H. V. Khanh, T. T. Dat *et al.*, “Mispronunciation detection and diagnosis without model training: A retrieval-based approach,” *arXiv preprint arXiv:2511.20107*, 2025.
- [23] W. Ye, S. Mao, F. Soong *et al.*, “An approach to mispronunciation detection and diagnosis with acoustic, phonetic and linguistic (apl) embeddings,” in *ICASSP*. IEEE, 2022, pp. 6827–6831.
- [24] L. Peng, Y. Gao, B. Lin *et al.*, “Text-aware end-to-end mispronunciation detection and diagnosis,” *arXiv preprint arXiv:2206.07289*, 2022.
- [25] N. Zheng, L. Deng, W. Huang *et al.*, “Coca-mdd: A coupled cross-attention based framework for streaming mispronunciation detection and diagnosis,” in *Interspeech 2022*, 2022, pp. 4352–4356.
- [26] B.-C. Yan, H.-W. Wang, Y.-C. Wang *et al.*, “Effective graph-based modeling of articulation traits for mispronunciation detection and diagnosis,” in *ICASSP*. IEEE, 2023, pp. 1–5.
- [27] B.-C. Yan, H.-W. Wang, and B. Chen, “Peppanet: Effective mispronunciation detection and diagnosis leveraging phonetic, phonological, and acoustic cues,” in *IEEE SLT*, 2022, pp. 1045–1051.
- [28] Z. Ma, G. Yang, W. Chen *et al.*, “Slam-llm: A modular, open-source multimodal large language model framework and best practice for speech, language, audio and music processing,” *Proc. IEEE Journal of Selected Topics in Signal Processing*, 2026.
- [29] Y. Chu, J. Xu, Q. Yang *et al.*, “Qwen2-audio technical report,” *arXiv preprint arXiv:2407.10759*, 2024.
- [30] M. Wu, J. Xu, X. Wu *et al.*, “Prompting large language models with mispronunciation detection and diagnosis abilities,” in *Interspeech*, 2024, pp. 2990–2994.
- [31] M. Wu, J. Xu, X. Chen *et al.*, “Integrating potential pronunciations for enhanced mispronunciation detection and diagnosis ability in LLMs,” in *ICASSP*. IEEE, 2025, pp. 1–5.
- [32] Z. Song, Z. Kadeer, M. Kahaer *et al.*, “Phoneme-controlled llm with self-supervised speech prompts for mispronunciation detection,” *ACMMM Asia*, 2025.
- [33] H. Zhong, Y. Xie, and Z. Yao, “Leveraging large language models to refine automatic feedback generation at articulatory level in computer aided pronunciation training,” in *INTERSPEECH*, 2024.
- [34] Y. Xie, H. Zhong, X. Lan *et al.*, “Mispronunciation detection and diagnosis based on large language models,” *Computer Speech & Language*, p. 101942, 2026.
- [35] N. Minematsu, C. Zhu, G. Dangtran *et al.*, “Development of shadowing speech corpora to measure instantaneous intelligibility as sequential annotation on L2 speech,” in *Tech. Rep. Speech, Acoust. Soc. Jpn*, 2022, pp. 7–12.
- [36] H. Geng, D. Saito, and N. Minematsu, “A pilot study of applying sequence-to-sequence voice conversion to evaluate the intelligibility of l2 speech using a native speaker’s shadowings,” in *Proc. APSIPA ASC*, 2024, pp. 1–6.
- [37] —, “A perception-based l2 speech intelligibility indicator: Leveraging a rater’s shadowing and sequence-to-sequence voice conversion,” in *Proc. Interspeech 2025*, 2025, pp. 2420–2424.
- [38] A. Zeyer, R. Schlüter, and H. Ney, “Why does CTC result in peaky behavior?” *arXiv preprint arXiv:2105.14849*, 2021.

- 801 [39] J. Tian, B. Yan, J. Yu *et al.*, “Bayes risk ctc: Controllable
802 ctc alignment in sequence-to-sequence tasks,” *arXiv preprint*
803 *arXiv:2210.07499*, 2022.
- 804 [40] Z. Yao, W. Kang, F. Kuang *et al.*, “Delay-penalized ctc imple-
805 mented based on finite state transducer,” in *Proc. Interspeech*
806 *2023*, 2023, pp. 1329–1333.
- 807 [41] Y. Kaloga, S. Kumar, P. Motlicek *et al.*, “A differentiable align-
808 ment framework for sequence-to-sequence modeling via optimal
809 transport,” *arXiv preprint arXiv:2502.01588*, 2025.
- 810 [42] Z. Yao, W. Kang, X. Yang *et al.*, “CR-CTC: Consistency regular-
811 ization on ctc for improved speech recognition,” in *ICLR*, 2024.
- 812 [43] V. Vapnik and A. Vashist, “A new learning paradigm: Learning
813 using privileged information,” *Neural networks*, vol. 22, no. 5-6,
814 pp. 544–557, 2009.
- 815 [44] V. Vapnik and R. Izmailov, “Learning using privileged informa-
816 tion: similarity control and knowledge transfer,” *The Journal of*
817 *Machine Learning Research*, vol. 16, no. 1, pp. 2023–2049, 2015.
- 818 [45] T. Fukuda and S. Thomas, “Implicit transfer of privileged acoustic
819 information in a generalized knowledge distillation framework”
820 in *Interspeech*, 2020, pp. 41–45.
- 821 [46] T. Fukuda, M. Suzuki, G. Kurata *et al.*, “Efficient knowledge dis-
822 tillation from an ensemble of teachers,” in *Interspeech*, 2017, pp.
823 3697–3701.
- 824 [47] X. Lu, P. Shen, Y. Tsao *et al.*, “Temporal order preserved optimal
825 transport-based cross-modal knowledge transfer learning for asr,”
826 in *2024 SLT*, 2024, pp. 1–8.
- 827 [48] D. Lopez-Paz, L. Bottou, B. Schölkopf *et al.*, “Unifying distilla-
828 tion and privileged information,” *ICLR*, 2016.
- 829 [49] K. Wang, G. Ortiz-Jimenez, R. Jenatton *et al.*, “Pi-dual: using
830 privileged information to distinguish clean from noisy labels,” in
831 *ICML*, 2024.
- 832 [50] G. Ortiz-Jimenez, M. Collier, A. Nawalgaria *et al.*, “When does
833 privileged information explain away label noise?” in *ICML*, 2023,
834 pp. 26 646–26 669.
- 835 [51] J. Su, M. Ahmed, Y. Lu *et al.*, “Roformer: Enhanced transformer
836 with rotary position embedding,” *Neurocomputing*, vol. 568, p.
837 127063, 2024.
- 838 [52] H. Tachibana, K. Uenoyama, and S. Aihara, “Efficiently trainable
839 text-to-speech system based on deep convolutional networks with
840 guided attention,” in *ICASSP*. IEEE, 2018, pp. 4784–4788.
- 841 [53] E. J. Hu, Y. Shen, P. Wallis *et al.*, “Lora: Low-rank adaptation of
842 large language models.” *Iclr*, vol. 1, no. 2, p. 3, 2022.
- 843 [54] A. Gulati, J. Qin, C.-C. Chiu *et al.*, “Conformer: Convolution-
844 augmented transformer for speech recognition,” in *Interspeech*,
845 2020, pp. 5036–5040.
- 846 [55] A. Grattafiori, A. Dubey, A. Jauhri *et al.*, “The llama 3 herd of
847 models,” *arXiv preprint arXiv:2407.21783*, 2024.
- 848 [56] Q. Team, “Qwen3 technical report,” 2025. [Online]. Available:
849 <https://arxiv.org/abs/2505.09388>
- 850 [57] G. Zhao, S. Sonsaat, A. Silpachai *et al.*, “L2-arctic: A non-native
851 english speech corpus,” in *Interspeech*, 2018, p. 2783–2787.
- 852 [58] J. Zhang, Z. Zhang, Y. Wang *et al.*, “speechocean762: An open-
853 source non-native english speech corpus for pronunciation assess-
854 ment,” in *Proc. Interspeech 2021*, 2021.
- 855 [59] S. D. C. of the Priority Areas Project, “Ume english speech
856 database read by japanese students (ume-erj),” jun 2007.
- 857 [60] T. Makino and R. Aoki, “English read by japanese phonetic cor-
858 pus: An interim report,” *Research in Language*, vol. 9, no. 2, pp.
859 79–95, 2012.
- 860 [61] K. Li, X. Qian, and H. Meng, “Mispronunciation detection and
861 diagnosis in L2 english speech using multidistribution deep neural
862 networks,” *IEEE/ACM TASLP*, vol. 25, no. 1, pp. 193–207, 2017.
- 863 [62] D. Y. Zhang, S. Saha, and S. Campbell, “Phonetic rnn-transducer
864 for mispronunciation diagnosis,” in *ICASSP*, 2023, pp. 1–5.
- [63] Y. Wan, Y. Shi, B. Lin *et al.*, “A study of mispronunciation detec- 865
tion and diagnosis based on meta-learning,” in *ICASSP*, 2024. 866
- [64] C. Zhu, A. Wumaier, D. Wei *et al.*, “Pronunciation error detection 867
model based on feature fusion,” *Speech Communication*, vol. 156, 868
p. 103009, 2024. 869
- [65] Y.-Y. He, B.-C. Yan, T.-H. Lo *et al.*, “JAM: A unified neural archi- 870
tecture for joint multi-granularity pronunciation assessment and 871
phone-level mispronunciation detection and diagnosis towards a 872
comprehensive capt system,” in *2024 APSIPA ASC*, 2024, pp. 1– 873
6. 874
- [66] B.-C. Yan, M.-K. Tsai, and B. Chen, “Muffin: Multifaceted pro- 875
nunciation feedback model with interactive hierarchical neural 876
modeling,” *IEEE Transactions on Audio, Speech and Language* 877
Processing, vol. 33, pp. 4295–4310, 2025. 878
- [67] H. Ryu, S. Kim, and M. Chung, “A Joint Model for Pronunciation 879
Assessment and Mispronunciation Detection and Diagnosis with 880
Multi-task Learning,” in *Interspeech 2023*, 2023, pp. 959–963. 881
- [68] N. Halabi and M. Wald, “Phonetic inventory for an Arabic speech 882
corpus,” in *LREC*, May 2016, pp. 734–738. 883