IMPACT-TTS: A Multimodal Prompt and Control Approach for Overcoming Low-Resource Constraints in Emotional TTS

Anonymous ACL submission

Abstract

Advancing emotional expressiveness in Text-to-Speech (TTS) systems remains a pivotal challenge for achieving natural and adaptive voice synthesis. Existing emotion-aware TTS models often struggle with limited emotional diversity, lack of fine-grained control, and reliance on small, labeled emotional speech-text datasets, making them less scalable and adaptable. To address these limitations, we propose IMPACT-TTS, an Integrated Multimodal Prompting and Adaptive Control for Text-to-Speech system 011 that effectively leverages a disentangled emo-012 tion module and a novel emotion modulation function. By incorporating large-scale pre-015 trained multimodal models, IMPACT-TTS mitigates dataset constraints while enabling flexible 017 emotional adjustments via prompt-based control. Our approach allows seamless blending of emotional intensities, significantly enhancing 019 expressiveness even in low-resource labeled datasets. Experimental results demonstrate that IMPACT-TTS outperforms existing models in emotional naturalness and adaptability, offering a scalable solution for emotion-aware TTS.

1 Introduction

027

033

037

041

Text-to-Speech (TTS) technology has made significant strides in recent years, particularly in improving intelligibility, naturalness, and speaker adaptation. However, achieving high levels of emotional expressiveness remains a critical challenge. Despite progress in prosody modeling and expressiveness enhancement, many existing TTS models struggle with capturing and generating nuanced emotional variations due to limited labeled emotional speech-text datasets and constrained emotion control mechanisms.

Existing approaches to emotional TTS primarily rely on textual cues or predefined emotion labels, which restrict the system's ability to generate diverse and contextually appropriate emotional speech. Moreover, models such as Hierspeech (Lee and et al, 2022) and Hierspeech++ (Lee and et al, 2023) highlight the *one-to-many mapping problem*, where a single input text can correspond to multiple valid emotional outputs, leading to ambiguity and reduced expressiveness. While approaches like Speech Slytherin (Jiang and et al, 2024) and limited-data two-stage models (Zhou et al., 2021) aim to address these issues through disentangled representations, they often require extensive labeled training data and computationally expensive architectures, making them less adaptable for low-resource scenarios.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Recent advancements in multimodal learning have introduced new possibilities for enhancing emotional expressiveness in TTS. By integrating diverse modalities such as text, audio, and vision, these methods aim to improve emotional fidelity and naturalness. For instance, MM-TTS (Li et al., 2024) leverages multimodal alignment techniques to enhance emotion modeling. However, it remains highly dependent on explicitly labeled datasets and does not provide fine-grained emotion control. Similarly, models such as VoiceLDM (Lee and et al, 2024) and ImaginaryVoice (Lee et al., 2023) incorporate multimodal integration but rely predominantly on textual or visual inputs, limiting the scope of emotional cues that can be effectively utilized.

In addition, prompt-based control mechanisms have gained traction for enabling user-guided emotional synthesis. Models like PromptStyle (Liu and et al, 2023), InstructTTS (Yang and et al, 2024), PromptTTS (Guo and et al, 2023), and PromptTTS2 (Leng and et al, 2023) explore the use of descriptive text prompts to facilitate controllable emotional modulation. While these methods provide greater flexibility, they are still constrained by their dependence on textual descriptions alone, which limits their ability to capture the full spectrum of emotional variation present in human speech.

To address these challenges, we introduce



Figure 1: Overview of the IMPACT-TTS architecture. The yellow box represents the emotion embedding, the red box represents the speaker embedding, the orange box represents the textual embedding, and the blue box represents the ground truth audio waveform embedding.

IMPACT-TTS. a novel framework that enhances emotional expressiveness by incorporating multimodal inputs and prompt-based emotion control while overcoming the limitations of small labeled emotional datasets. Unlike prior works that rely heavily on constrained emotion labels, IMPACT-TTS integrates large-scale pretrained multimodal models to extract meaningful emotion representations. Notably, we leverage ONE-PEACE (Wang and et al., 2023), a highly extensible multimodal model with 4B parameters, featuring separate adapters for audio, language, and vision. For audio representation, our system utilizes the feature extractor weights of WavLM (Sanyuan et al., 2022) for initialization, enhancing its ability to generalize from smaller datasets effectively.

090

094

By leveraging large pretrained representations, IMPACT-TTS significantly reduces reliance on explicitly labeled emotional speech datasets while 101 maintaining fine-grained emotion modulation ca-102 pabilities. Moreover, our approach employs a dis-103 entangled architecture to separate the core speech 104 105 synthesis process from the emotion module, ensuring robust and stable synthesis while enabling 106 precise control over emotional variations. This 107 disentanglement strategy allows IMPACT-TTS to generate expressive speech with a higher degree of 109

controllability, even in low-resource settings.

2 Method

IMPACT-TTS employs a disentangled structure for independent control of linguistic and emotional features, ensuring dynamic emotion modulation and consistency. The following sections describe the TTS model, emotion embedding with cross attention, emotion modulation function, and emotion consistency loss. 110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

2.1 Backbone Arthitecture

IMPACT-TTS is built upon a modified VITS architecture, inspired by YourTTS (Casanova and et al, 2022) and VECL-TTS (Gudmalwar and et al, 2024), utilizing variational inference and adversarial learning to generate high-quality speech. The backbone is designed to produce natural and highquality speech from input text, ensuring both linguistic and phonetic accuracy. The model features a transformer-based text encoder with 10 blocks, a decoder with 4 affine coupling layers, each comprising 4 WaveNet residual blocks (Oord et al., 2016), as inspired by previous work. HiFi-GAN (Kong et al., 2020) V1 is employed as the neural vocoder to generate speech waveforms. A Variational Auto-Encoder (VAE) with a Posterior



Figure 2: Process of Emotion Module

2.2

135 Encoder comprising 16 WaveNet residual blocks transforms linear spectrograms into latent variables, 136 seamlessly integrating the vocoder and the flow-137 based decoder for end-to-end synthesis. To enable 138 diverse rhythm generation from text, a Stochastic 139 Duration Predictor (SDP) is used, further enriching 140 the synthesized speech's naturalness. As shown 141 in Figure 1, the architecture combines emotion, 142 speaker, text, and waveform embeddings to control 143 speech attributes. Emotion embeddings (yellow) 144 from the Emotion Module modulate tones, speaker 145 embeddings (red) from H/ASP (Heo et al., 2020) 146 ensure consistency, text embeddings (orange) en-147 code linguistic information, and waveform embed-148 dings (blue) provide training supervision. The or-149 ange box corresponds to the textual embedding, encoding linguistic and phonetic information. Finally, 151 the blue box represents the ground truth waveform 152 embedding, used during training for supervision. 153 These embeddings are seamlessly integrated within 154 the TTS backbone, allowing precise control over 155 emotion, speaker, and linguistic features. To en-156 sure robustness and expressiveness of the model, 157 we incorporate Emotion Consistency Loss (ECL) 158 and Speaker Consistency Loss (SCL) to maximize 159 the similarity of emotional and speaker attributes 160 between the generated audio and the ground truth. 161

$$L_{CL} = -\frac{\alpha}{n} \cdot \sum_{i}^{n} \cos_sim (\phi(g_i), \phi(h_i)), \quad (1)$$

163where $\phi(g)$ and $\phi(h)$ represent the functions used164to extract emotion embeddings or speaker embed-165dings from the generated and ground truth audio,166respectively. The parameter α is a tunable weight167that determines the contribution of L_{CL} to the over-168all loss. The cosine similarity function is denoted as169 cos_sim , which measures the similarity between170the emotion embeddings.

162

To enhance emotional expressiveness in speech synthesis, our model integrates multimodal representations from text, vision, and audio through a unified embedding space. This section details how these modalities contribute to emotion synthesis, complementing the emotion disentanglement framework.

Multimodal Representation and Prompts

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

196

197

198

199

200

201

202

203

204

205

2.2.1 Multimodal Representation Model

Traditional emotional TTS systems rely heavily on text-based emotional conditioning, limiting expressiveness. Our model overcomes this limitation by incorporating vision and audio cues, ensuring richer emotion transfer. The extracted representations are aligned into a shared multimodal space to balance textual, auditory, and visual emotion information.

Vision and Audio Embedding Fusion

- Vision Representation: Facial expressions offer valuable nonverbal emotional cues. We extract visual embeddings from a pretrained vision-audio-language model (ONE-PEACE) and map them into our emotion embedding space.
- Audio Emotion Features: Instead of relying only on text-based emotion modeling, we extract pitch contours, energy variation, and prosodic rhythm features from speech waveforms, which provide additional emotion grounding.

Multimodal Cross-Attention for Emotion Integration The extracted text, vision, and audio embeddings are integrated using a cross-attention mechanism. Unlike direct concatenation, this allows for dynamic weighting of each modality based

242

on its relevance to the input context. This mechanism enhances context-aware emotion modulation and provides finer control over expressiveness.

2.2.2 Prompt-Based Emotion Control

207

209

224

231

237

Rather than relying on categorical emotion labels, 210 we employ prompt-based modulation, where textual prompts dynamically influence synthesized 212 speech expressions. These text prompts, described 213 in the Datasets section, are mapped into semantic 214 embedding spaces and fused with vision and audio 215 representations. 216

LLM-based Emotion Prompt Expansion In-217 stead of fixed emotion labels, we generate emotion-218 ally diverse paraphrases of input sentences using a structured synonym mapping algorithm. Given a neutral sentence, our system selects semantically appropriate modifications to control emotion inten-222 sity and expressiveness.

Multimodal Prompt Adaptation The generated text prompts are conditioned alongside vision and audio cues, ensuring that emotional modulation is not solely text-driven. Vision embeddings provide global emotional state awareness, while audio features introduce prosodic characteristics, leading to a richer emotional representation.

2.2.3 **Ablation Study: Impact of Different Modalities**

To quantify the contribution of each modality (text, vision, and audio) to emotional speech synthesis, we conduct an ablation study where we systematically remove individual modalities and analyze their impact on expressiveness.

Table 1: Ablation Study: Impact of Different Modalities on Emotion Expression. Higher values for std-F0 (Pitch Variability), std-RMS (Energy Variability), and WavLM Score (Expressiveness) indicate greater expressiveness.

Configuration	std-F0 \uparrow	std-RMS ↑	WavLM Score ↑
Audio Cues Only	0.21	0.18	3.5
Vision Cues Only	0.15	0.14	3.2
Text Prompts Only	0.12	0.10	2.8
Audio + Vision	0.25	0.22	3.9
Audio + Text	0.28	0.24	4.1
Text + Vision	0.20	0.17	3.4
Text + Vision + Audio	0.32	0.28	4.5

The results demonstrate that integrating all three modalities (Text, Vision, and Audio) significantly enhances emotional speech synthesis, as indicated by higher pitch and energy variability, as well as improved expressiveness scores.

2.3 **Disentangled Emotion Module**

IMPACT-TTS adopts a multi-modal emotion encoder to get emotion embeddings from various modal, text and vision, and combine these embeddings with cross attention. The emotion module operates independently of the TTS backbone, modulating the emotional tone of the synthesized speech using the emotion embedding E_{emotion} with the process of Figure 2. By disentangling the emotion module from the TTS backbone, IMPACT-TTS achieves precise control over emotional attributes. This disentangled structure ensures that the linguistic content remains unaffected by emotional variability.

243

245

246

247

248

249

250

251

253

254

255

257

258

259

260

261

262

263

265

266

267

270

271

272

273

274

275

276

277

278

279

284

287

2.3.1 **Emotion Encoder**

The emotion embedding extraction module utilizes finetuned representation model with retrieval function as there are some important strengths. Since retrieval models focus on semantic similarity rather than discrete classification, they offer rich, crossmodal feature representations suitable for TTS task compared to classification models. Through this, they can capture fine-grained relationships between the input and output, which is crucial for generating high-quality, natural speech. We finetuned ONE-PEACE (Wang and et al., 2023) with MEAD-TTS dataset(Guan et al., 2024) and ESD(Zhou and et al, 2022b), and use it as the Emotion Representation Model in Figure 2.

The cross-attention mechanism computes attention between different modalities. We integrate cross-attention module into a emotion encoder to fuse modalities by allowing one modality to attend to the other modal's features. Given two modalities: modality A(e.g., text embeddings) and modality **B**(e.g., visual features or audio features): Ouerv (Q) from modality A, and Keys (K) and Values (V) from modality B. The cross-attention can be formulated as:

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V,$$
 (2)

Embeddings from text-speech modalities and text-image modalities are concatenated. At the end, we make it transform the multimodal representation into a continuous emotion embedding, which is used to guide the synthesis of emotionally expressive speech.

Model Variant	Prosody	Emotion Metrics		Speech Quality	
	DTW (F0) \downarrow	Pitch Var \uparrow	ECA \uparrow	$\mathrm{KLD}\downarrow$	$MCD\downarrow$
Full Model (Ours)	15.3	0.32	85.2%	0.12	4.25
w/o Emotion Modulation	20.8	0.21	78.9%	0.19	4.67
w/o Emotion Classification	18.1	0.27	80.3%	0.17	4.50
w/o Both	24.5	0.15	73.4%	0.24	4.80

Table 2: Ablation Study on Emotion Modulation and Classification

292

294

295

305

307

311

Emotion Classifier 2.3.2

The Emotion Classifier Module serves to categorize the high-dimensional emotion embeddings generated by the Emotion Encoder into predefined emotion classes. This module operates in two stages: 1. Linear Transformation The emotion embedding is first projected into a lower-dimensional space that corresponds to the number of predefined using a linear layer, enabling the embedding to be directly mapped to emotion class probabilities; 2. Softmax **Normalization** The output of the linear layer is passed through a softmax activation function. The final output of the Emotion Classifier Module is a set of probabilities, where each value corresponds to the likelihood of the input belonging to a specific emotion class. For example, in Figure 2, an output of [0.3, 0.6, 0.1] for three classes (e.g., neutral, happy, sad) indicates a higher confidence for the "happy" emotion. These probabilities serve as critical inputs for the subsequent Emotion Modulation, where they are utilized to determine the blending weights (α_i) in fine-grained emotion interpolation. 310

Emotion Modulation Function 2.3.3

312 To achieve fine-grained control over blended emotions, IMPACT-TTS incorporates a spherical in-313 terpolation method inspired by (Cho and et al., 314 2024), (Zhou and et al, 2022a) and (Im and et al., 315 2022). This method allows for smooth transitions 316 and precise mixing of multiple emotional attributes by treating emotions as points on a hypersphere. The blending weights for interpolation, denoted as 319 α_i , are derived directly from the probability distribution generated by the Emotion Classifier Module. 322 These weights ensure that the contribution of each emotion vector to the final blended representation aligns with the classifier's confidence levels. Given two emotions represented as vectors e_1 and e_2 , the blended emotion e_{blend} is computed as: 326

$$e_{\text{blend}} = \frac{\sin((1-\alpha)\theta)}{\sin(\theta)}e_1 + \frac{\sin(\alpha\theta)}{\sin(\theta)}e_2$$
, (3)

328

329

331

333

334

335

336

337

338

339

340

341

343

344

345

346

347

349

350

351

352

354

356

357

360

where $\alpha \in [0, 1]$ is the interpolation parameter that controls the influence of each emotion, and θ = $\arccos\left(\frac{e_1 \cdot e_2}{\|e_1\|\|e_2\|}\right)$ is the angle between e_1 and e_2 on the hypersphere. This formulation ensures that the resulting vector remains normalized, preserving the geometric properties of the emotion space.

Extension to Multiple Emotions In scenarios where more than two emotions are blended, the spherical interpolation is generalized to handle nemotions. Using the probabilities from the classifier $(\{\alpha_1, \alpha_2, \dots, \alpha_n\})$, where $\sum_{i=1}^n \alpha_i = 1$), the blended emotion e_{blend} is computed as:

$$e_{\text{blend}} = \frac{\sum_{i=1}^{n} \alpha_i \sin(\theta_i) e_i}{\sum_{i=1}^{n} \sin(\theta_i)} , \qquad (4)$$

where θ_i is the angle between each emotion vector e_i and the reference point (e.g., a neutral emotion vector). This integration of classifier outputs into the modulation function ensures that the model can generate expressive, contextually accurate speech by leveraging both high-level classification results and continuous fine-grained emotion control.

2.3.4 **Ablation Study on Emotion Modulation** and Classification

To assess the contribution of each component in the disentangled emotion module, we conduct an ablation study by systematically removing (i) emotion modulation, (ii) emotion classification, and (iii) both components. Table 2 presents the results across prosody metrics, emotion classification accuracy, and speech quality.

The results indicate that both emotion modulation and classification significantly contribute to the overall performance of our model. Removing

Method	Intra-domain			Out-of-domain				
	MOS	STOI	ECS	MCD	MOS	STOI	ECS	MCD
GT	4.53 ± 0.03	0.87	-	-	4.73 ± 0.05	0.85	-	-
GT (Mel+Voc)	4.51 ± 0.01	0.74	-	-	4.71 ± 0.04	0.79	-	-
MM-TTS (Li et al., 2024)	4.32 ± 0.07	0.42	0.87	3.21	3.96 ± 0.09	0.33	0.74	6.68
Instruct-TTS	4.31 ± 0.05	0.40	0.89	3.10	3.83 ± 0.08	0.31	0.75	6.59
Prompt-TTS	4.23 ± 0.06	0.21	0.75	3.78	3.66 ± 0.08	0.27	0.73	6.73
IMPACT-TTS (proposed)	$\textbf{4.41} \pm \textbf{0.02}$	0.63	0.92	3.17	$\textbf{4.02} \pm \textbf{0.03}$	0.54	0.77	6.12

Table 3: The performance comparison of text prompt based TTS.

emotion modulation results in a higher DTW (Dynamic Time Warping; F0) score (20.8 vs. 15.3), suggesting a reduction in prosodic expressiveness. Additionally, pitch variability drops from 0.32 to 0.21, indicating that synthesized speech becomes less emotionally dynamic.

361

367

371

373

374

375

377

386

395

Similarly, excluding emotion classification lowers the Emotion Classification Accuracy (ECA) from 85.2% to 80.3%, showing that the model struggles to preserve intended emotional expressiveness without explicit classification. The Kullback-Leibler Divergence (KLD) increases from 0.12 to 0.17, further confirming reduced alignment with expected emotion distributions.

The most significant degradation occurs when both components are removed, with ECA dropping to 73.4% and MCD increasing to 4.80, indicating a loss in both emotional clarity and speech quality. These findings demonstrate that the disentangled emotion module plays a crucial role in achieving both expressive and high-quality synthesized speech.

3 Inference-Time Emotion Adjustment and Expressiveness Control

3.1 Real-Time Emotion Modulation

Traditional emotional TTS models often rely on pretrained emotion embeddings that remain fixed during inference, limiting their ability to adapt expressiveness dynamically. Our model introduces **Inference-Time Emotion Adjustment**, which allows for real-time modification of prosody and emotional intensity during speech synthesis.

After obtaining a pretrained emotion embedding, we apply an emotion modulation function that dynamically adjusts: • **Pitch (F0)**: Controls intonation variation, allowing for smooth shifts in emotion intensity.

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

• Energy: Determines speech loudness and emphasis, enabling finer expressiveness control.

Smooth Emotion Interpolation. Unlike models that use predefined emotion categories, our system enables continuous control over emotion intensity. By leveraging a spherical interpolation mechanism, the model allows smooth transitions between emotions:

- Neutral → Happy: Gradual increase in pitch and speech rate.
- Sad \rightarrow Angry: Lowered pitch at the start, followed by a rapid increase in energy.
- Excited \rightarrow Calm: Controlled pitch decay with reduced energy.

This capability ensures that synthesized speech adapts dynamically rather than being constrained by fixed, pre-trained emotion embeddings.

3.2 Comparison with Contrastive Learning-Based Emotion Alignment

Many multimodal emotional TTS models incorporate contrastive learning to pre-align emotion embeddings from text, vision, and audio into a shared representation space. This approach ensures consistency across modalities but relies on fixed emotion representations at inference time.

Our approach introduces an alternative strategy by integrating real-time emotion modulation, where emotion embeddings are adaptively adjusted at inference time based on synthesis needs. This enables:

474

475

- Fine-Grained Emotion Interpolation Blending between emotions smoothly rather than switching between predefined categories.
 - Intensity Scaling Adjusting the strength of an emotion dynamically (e.g., slightly sad vs. deeply sad).
 - Context-Aware Expressiveness Modulating prosody in real-time based on speaker intent.

4 Datasets

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

453

Pretraining The TTS backbone is initially trained on neutral speech data to establish a robust foundation for speech generation. We use VCTK(Yamagishi et al., 2019) dataset as the pretraining data, which comprises of 44 hours of speech from 109 different speakers.

Finetuning The emotion module is introduced, and 443 the entire system is finetuned using datasets en-444 riched with emotional speech samples to ensure 445 consistent and expressive output. Specifically, we 446 leverage two datasets for finetuning: ESD (Zhou 447 and et al, 2022b), which covers 5 emotion classes 448 (neutral, angry, happy, sad, and surprised), and 449 MEAD-TTS (Guan et al., 2024) dataset, which en-450 compasses 8 emotion classes (neutral, angry, con-451 452 tempt, disgusted, fear, happy, sad and surprised).

4.1 Generation of Emotional Text Prompt

Given a dataset annotated with emotional tags, we 454 generate multiple emotionally nuanced variations 455 456 of each sentence. This ensures that the text better reflects the speaker's emotional state, enhancing the 457 system's ability to produce more accurate prosodic 458 features during synthesis. The input data includes 459 columns indicating the speaker's sex (male, female) 460 and the associated emotion from the 8 classes. Us-461 ing this information, we generate five emotionally 462 enriched variations of the sentence as text prompts 463 using GPT-4.0 for each input sentences based on 464 (Primus et al., 2023). 465

Synonym Mapping for Diversity We employ a 466 synonym mapping technique to introduce lexical 467 variety in the emotion-contained sentences. Each 468 469 emotion is associated with a set of semantically similar words, such as "angry", "furious", "irate", 470 etc. This lexical variety helps diversify the emo-471 tional cues present in the input text, thereby im-472 proving the robustness of the generated speech. 473

4.2 **Comparison of Emotion-Labeled Dataset** Sizes

To analyze the effectiveness of IMPACT-TTS in low-resource scenarios, we compare the size of labeled emotional speech datasets used in our model with those of MM-TTS, InstructTTS, and PromptTTS. The dataset size significantly impacts the ability of a model to generalize emotional expressiveness effectively.

As shown in Table 4, IMPACT-TTS effectively reduces reliance on large labeled emotional speech datasets by leveraging pretrained multimodal models, whereas other models require substantially larger explicitly labeled datasets.

5 **Experimental Results**

We evaluate the generated speech quality and similarity by objective metrics and subjective eval-Short-Time Objective Intelligibility uations. (STOI) is an objective measure of how understandable and clear speech is. The STOI score ranges from 0 to 1, where 1 indicates perfectly intelligible speech. Emotion Cosine Similarity (ECS) measures the similarity between the synthesized and target emotional embeddings, with higher ECS values indicating better alignment. Mel Cepstral **Distortion (MCD)** measures the spectral distance between the ground truth and synthesized Melspectrum features. As for subjective evaluations, we conduct 5-scale Mean Opinion Score (MOS) using MOSNet (Chen-Chou et al., 2019) for 5 times. We generate 50 speech samples for each model and select 50 samples from MEAD-TTS and LibriTTS testing sets for intra-domain and out-of-domain evaluation respectively. Table 3 presents the experimental results of text prompt-based speech synthesis, encompassing MOS, audio quality and classification accuracy for emotion and gender. For text prompt based TTS, we conduct experiments on the 511 following systems : 1) GT: This is the ground-truth 512 recording; 2) GT (Mel + Voc): This is the speech 513 synthesized using pretrained HiFi-GAN vocoder 514 for GT Mel-spectrogram; 3) MM-TTS (Guan et al., 515 2024): This is a model for multi-modal prompt 516 based TTS, which prompt encoder based on CLIP; 517 4) IMPACT-TTS: This is the proposed model for 518 multi-modal prompt based expressive TTS. In the 519 context of prompts, IMPACT-TTS model surpasses 520 the baseline model in terms of audio naturalness 521 and classification performance on the MEAD-TTS 522 dataset as well as out-of-domain datasets. These 523

Model	Emotion-Labeled Speech Data (Hours)
IMPACT-TTS (Ours)	65
MM-TTS (Li et al., 2024)	100
InstructTTS (Yang and et al, 2024)	300
PromptTTS (Guo and et al, 2023)	152

Table 4: Comparison of emotion-labeled dataset sizes across emotional TTS models. The dataset sizes (in hours) are estimated. MM-TTS also utilizes 15,433 emotion-labeled images from the RAF-DB dataset.

results highlight the effectiveness of the method in accurately capturing and extracting emotion attributes from various text prompts.

6 Limitation

524

525

526

544

545

546

547

549

550

551

558

559

562

One limitation of the current model is the lack of vision datasets for facial expressions, leading 529 to confusion with prompts like "sorrowful eyes," 530 where happy emotions may be generated due to 531 overlapping features like watery eyes. In the future, this limitation could be addressed by incorporating more diverse and abundant datasets and 534 employing fine-grained image detection techniques to distinguish specific facial components. Another challenge lies in the large size of the representation model, which requires substantial server capacity for effective processing. Future work will focus on exploring lightweight alternatives, such 540 541 as compressing the emotion encoder through pruning and quantization or using parameter-efficient 542 fine-tuning methods. 543

7 Conclusion and Future Work

In this paper, we introduced IMPACT-TTS, a multimodal prompt-based TTS system that enhances emotional expressiveness while overcoming the limitations of small labeled datasets. By integrating large-scale pretrained models, generating text prompts with diverse synonyms and employing spherical interpolation for emotion modulation, IMPACT-TTS offers fine-grained control over emotional expression with minimal reliance on explicit emotion labels.

Future work will focus on improving efficiency by developing lightweight versions of the multimodal components. Additionally, we aim to extend the framework to support a broader range of emotional variations, including nuanced styles and speaker-adaptive expressions. We will also explore integrating ethical safeguards, such as audio watermarking for authenticity verification and usage monitoring, to prevent misuse. Incorporating broader ethical considerations into the design and deployment process will strengthen the trustworthiness and social responsibility of emotional TTS systems like IMPACT-TTS. 563

564

565

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

References

- Edresson Casanova and et al. 2022. Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone. *International Conference on Machine Learning. PMLR*.
- Lo Chen-Chou, Fu Szu-Wei, and Huang Wen-Chin. 2019. Mosnet: Deep learning-based objective assessment for voice conversion. *arXiv preprint arXiv:1904.08352*.
- Deok-Hyeon Cho and et al. 2024. Emosphere-tts: Emotional style and intensity modeling via spherical emotion vector for controllable emotional text-to-speech. *INTERSPEECH 2024*.
- Wenhao Guan, Yishuang Li, Tao Li, and Hukai Huang. 2024. Mm-tts: Multi-modal prompt based style transfer for expressive text-to-speech synthesis. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38.
- Ashishkumar Gudmalwar and et al. 2024. Vecl-tts: Voice identity and emotional style controllable crosslingual text-to-speech. *INTERSPEECH 2024*.
- Zhifang Guo and et al. 2023. Promptts: Controllable text-to-speech with text descriptions. *ICASSP 2023-*2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE.
- Hee Soo Heo, Bong-Jin Lee, and Jaesung Huh. 2020. Clova baseline system for the voxceleb speaker recognition challenge 2020. *arXiv preprint arXiv:2009.14153*.
- Chae-Bin Im and et al. 2022. Emoq-tts: Emotion intensity quantization for fine-grained controllable emotional text-to-speech. *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE.*
- Xilin Jiang and et al. 2024. Speech slytherin: Examining the performance and efficiency of mamba for 603

- 610 611 612 613 614 615 616 617 618 619 622 623 625 629 630 631

- 641 643

654

speech separation, recognition, and synthesis. arXiv preprint arXiv:2407.09732.

- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. NeurIPS 2020.
- Jiyoung Lee, Joon Son Chung, and Soo-Whan Chung. 2023. Imaginary voice: Face-styled diffusion model for text-to-speech. ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE.
- Sang-Hoon Lee and et al. 2022. Hierspeech: Bridging the gap between text and speech by hierarchical variational inference using self-supervised representations for speech synthesis. Advances in Neural Information Processing Systems 35, pages 16624–16636.
 - Sang-Hoon Lee and et al. 2023. Hierspeech++: Bridging the gap between semantic and acoustic representation of speech by hierarchical variational inference for zero-shot speech synthesis. arXiv preprint arXiv:2311.12454.
- Yeonghyeon Lee and et al. 2024. Voiceldm: Text-tospeech with environmental context. ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE.
- Yichong Leng and et al. 2023. Prompttts 2: Describing and generating voices with text prompt. arXiv preprint arXiv:2309.02285.
- Xiang Li, Zhi-Qi Cheng, Jun-Yan He, and XiaoJiang Peng. 2024. Mm-tts: A unified framework for multimodal, prompt-induced emotional text-to-speech synthesis. arXiv preprint arXiv:2404.18398.
- Guanghou Liu and et al. 2023. Promptstyle: Controllable style transfer for text-to-speech with natural language descriptions. arXiv preprint arXiv:2305.19522.
- Aaron van den Oord, Sander Dieleman, and Heiga Zen. 2016. Wavenet: A generative model for raw audio. arXiv preprint arXiv:1609.03499.
- Paul Primus, Khaled Koutini, and Gerhard Widmer. 2023. Advancing natural-language based audio retrieval with passt and large audio-caption data sets. arXiv preprint arXiv:2308.04258.
- Chen Sanyuan, Wang Chengwi, and Chen Zhengyang. 2022. Wavlm: Large-scale self-supervised pretraining for full stack speech processing. arXiv preprint arXiv:2110.13900.
- Peng Wang and et al. 2023. One-peace: Exploring one general representation model toward unlimited modalities. arXiv preprint arXiv:2305.11172.
- Junichi Yamagishi, Christophe Veaux, and Kirsten MacDonald. 2019. The vctk corpus: The centre for speech technology research (cstr) voice cloning toolkit. https://datashare.ed.ac.uk/handle/ 10283/3443. Version 0.92.

Dongchao Yang and et al. 2024. Instructtts: Modelling expressive tts in discrete latent space with natural language style prompt. IEEE/ACM Transactions on Audio, Speech, and Language Processing.

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

- Kun Zhou and et al. 2022a. Emotion intensity and its control for emotional voice conversion. IEEE Transactions on Affective Computing 14.1, pages 31– 48.
- Kun Zhou and et al. 2022b. Emotional voice conversion: Theory, databases and esd. Speech Communication, 137:1-18.
- Kun Zhou, Berrak Sisman, and Haizhou Li. 2021. Limited data emotional voice conversion leveraging textto-speech: Two-stage sequence-to-sequence training. INTERSPEECH 2021.