# EMOTION-ALIGNED GENERATION IN DIFFUSION TEXT TO SPEECH MODELS VIA PREFERENCE-GUIDED OPTIMIZATION

*Jiacheng Shi*[⋆]    *Hongfei Du*[⋆]    *Yangfan He*[†]    *Y. Alicia Hong*[‡]    *Ye Gao*[⋆]

[⋆] College of William & Mary, [†] University of Minnesota - Twin Cities, [‡] George Mason University

{jshi12, hdu02, ygao18}@wm.edu, he000577@umn.edu, yhong22@gmu.edu

## ABSTRACT

Emotional text-to-speech seeks to convey affect while preserving intelligibility and prosody, yet existing methods rely on coarse labels or proxy classifiers and receive only utterance-level feedback. We introduce Emotion-Aware Stepwise Preference Optimization (EASPO), a post-training framework that aligns diffusion TTS with fine-grained emotional preferences at intermediate denoising steps. Central to our approach is EASPM, a time-conditioned model that scores noisy intermediate speech states and enables automatic preference pair construction. EASPO optimizes generation to match these stepwise preferences, enabling controllable emotional shaping. Experiments show superior performance over existing methods in both expressiveness and naturalness.

*Index Terms*— Speech Synthesis, Diffusion Model, Emotional Text-To-Speech (Emo-TTS).

## 1. INTRODUCTION

Emotional text-to-speech (TTS) [1–7] aims to generate speech that remains intelligible while conveying affect and prosodic nuance. Effective emotional control supports conversational agents, accessibility, and content creation. Yet fine-grained control without sacrificing naturalness remains challenging, as emotional cues unfold over time and interact with linguistic content and speaker variation.

Early emotional TTS systems fused local style tokens with phoneme content via cross-attention for fine-grained control [1], guided synthesis with semantic prompts [2], and used variational decoders to model nuanced emotion–prosody relations [3]. Diffusion-based methods improved fidelity and control through emotion embeddings and noise-conditioned prosody [4]. Recent work adopts supervised speech tokens for semantic grounding and chunk-aware flow-based decoding [5, 6], and introduces parallel phoneme–audio branches in LLM-based generation to support fine-grained, freestyle emotional control [7]. Across these variants, two important gaps remain: supervision still targets proxy labels over preference-driven prompts, and feedback remains temporally sparse, limiting constraints on prosody–emotion dynamics.

Direct Preference Optimization (DPO) [8] aligns generative models to human choices via paired comparisons against a frozen reference policy. In diffusion settings [9], a common practice is to assign preference at the final step and propagate it to intermediate latents, then optimize a DPO-style log-likelihood ratio at each step. This removes the need for an explicit reward model and is effective for global preferences, but offers limited guidance for temporally evolving signals. DPO has also been used to align large language model–based TTS. EmoDPO [10], for instance, pairs utterances with identical text and marks the target emotion as preferred, achieving alignment without explicit reward modeling. However, attaching preference only at the endpoint yields sparse supervision for gradually varying cues, the assumption that all intermediate states on a preferred path are themselves preferred is often invalid, and curating domain-appropriate preference pairs is expensive. This leads to our central research question: *How can we provide dense, fine-grained supervision for emotionally expressive TTS generation, without relying on endpoint preferences or categorical labels?*

To better align diffusion-based TTS systems with fine-grained emotional preferences, we propose *Emotion-Aware Stepwise Preference Optimization* (EASPO). At each denoising step starting from a latent representation $x_t$, the model samples a small candidate set of $x_{t-1}$ mel-spectrograms. An *Emotion-aware Stepwise Preference Model* (EASPM) scores their emotional expressiveness and selects a win–lose pair that differs subtly in prosody while preserving linguistic content. One candidate is then randomly chosen to continue generation. Since these candidates originate from the same latent and differ by only a single denoising step, their variations are localized and emotion-focused. EASPM captures these nuanced differences and steers the model toward producing more emotionally consistent speech.

Our contributions are: 1) We propose EASPO, a stepwise alignment framework that reformulates preference optimization as a local, time-conditioned task, replacing the flawed assumption that all intermediate states on a preferred trajectory are equally valid. By aligning win/lose candidates at each denoising step from a shared latent, it enables stepwise-controllable emotion shaping throughout the generation process. 2) We introduce EASPM, a time-aware reward model

that directly scores emotional expressiveness and prosody on noisy intermediate states, enabling dense, temporally grounded preference learning and on-the-fly scoring.

## 2. METHOD

We present a reinforcement learning framework for emotional TTS that fine-tunes diffusion models via stepwise preference supervision (Fig.1). Starting from a pretrained Grad-TTS [11], our method introduces dense emotion-aligned rewards at each denoising step. At every latent state, multiple mel-spectrograms are sampled and ranked by a frozen emotion preference model(Sec. 2.1). A preference pair is selected, and the model is optimized to favor the emotionally preferred sample by minimizing the gap between its advantage and the log-likelihood ratio(Sec. 2.2). Our EASPO objective extends DPO with stepwise preference signals, enabling fine-grained emotional control during generation.

### 2.1. Emotion-Aware Stepwise Preference Model

**Overall.** At reverse step $t-1$, given the current latent $x_t$, the generator draws $k$ candidates $\{x_{t-1}^1, \ldots, x_{t-1}^k\} \sim p_\theta(x_{t-1} \mid x_t)$. EASPM assigns each candidate a *timestep-aware* emotion–prompt consistency score and induces a ranking; the highest and lowest scored items are taken as the win/lose pair.
**Scoring.** EASPM is built on CLEP [12], a CLAP-based [13] contrastive audio–language encoder fine-tuned on large-scale emotional speech data. Let $f_{\text{CLEP-A}}(\cdot, t)$ and $f_{\text{CLEP-T}}(\cdot)$ denote the audio and text branches. Using $\ell_2$-normalized embeddings, the score for candidate $x_{t-1}^i$ is

$$s_i = \left\langle \frac{f_{\text{CLEP-A}}(x_{t-1}^i, t)}{\left\| f_{\text{CLEP-A}}(x_{t-1}^i, t) \right\|_2}, \frac{f_{\text{CLEP-T}}(c)}{\left\| f_{\text{CLEP-T}}(c) \right\|_2} \right\rangle. \quad (1)$$

For a pair $(x_{t-1}^w, x_{t-1}^l)$ with scores $s_w, s_l$, define $\Delta_t = s_w - s_l$. The probability that the win item is preferred is modeled by a pairwise logistic with temperature $\tau > 0$:

$$\hat{p}_w = \sigma(\tau \Delta_t) = \frac{1}{1 + \exp(-\tau \Delta_t)}. \quad (2)$$

Candidate selection is performed by $x_{t-1}^{\text{win}} = \arg\max_i s_i$ and $x_{t-1}^{\text{lose}} = \arg\min_i s_i$, and the stepwise preference loss is:

$$\mathcal{L}_{\text{pref}} = -\log \hat{p}_w = \log\left(1 + \exp(-\tau \Delta_t)\right). \quad (3)$$

**Training.** EASPM is adapted from CLEP to handle noisy intermediate representations. Given a win–lose pair $(x_0^w, x_0^l)$, we sample a timestep $t$ and apply the same forward diffusion to produce $(x_t^w, x_t^l)$. This tuple $(x_t^w, x_t^l, t, c)$ is used to optimize $\mathcal{L}_{\text{pref}}$, encouraging the model to recover the correct preference at step $t$. A time-aware normalization layer is added to CLEP's audio branch for timestep conditioning. To reduce mismatch with CLEP's pretraining domain, we optionally estimate a pseudo-clean $\hat{x}_0$ from $x_t$ via deterministic inversion

followed by [14]. After training, EASPM is frozen and used solely as a stepwise scorer for the RL objective.
**Random selection of the next state.** After EASPM ranks the candidate set $\{x_{t-1}^1, \ldots, x_{t-1}^k\}$ and selects a win–lose pair, we *do not* continue with the top sample as shown in Fig. 1. To avoid biased rollouts and degenerate paths, we uniformly sample the next state $\tilde{x}_{t-1}$ from the pool and proceed with $x_{t-2} \sim p_\theta(\cdot \mid \tilde{x}_{t-1}, c)$. This ensures all preference pairs originate from the same latent $x_t$ while decoupling supervision from sampling. Candidate pooling is applied only when $t \leq \kappa$; standard transitions are used for $t > \kappa$.

### 2.2. Objective Function of EASPO

We formulate denoising as a $T$-step MDP with state $s_t = (c, x_t)$, action $a_t = x_{t-1}$, and policy $\pi_\theta(a_t \mid s_t) = p_\theta(x_{t-1} \mid x_t, c)$. At each step $t$, we sample $k$ candidates $\{x_{t-1}^i\} \sim p_\theta(\cdot \mid x_t, c)$ and rank them via EASPM (Sec. 2.1). Let $x_{t-1}^w$ and $x_{t-1}^l$ denote the top and bottom-ranked samples from the *same* latent $x_t$.
**Dense stepwise reward.** EASPM supplies a dense emotional reward at step $t$:

$$\widehat{R}_t^j = s\left(x_{t-1}^j, c, t\right), \qquad \Delta \widehat{R}_t = \widehat{R}_t^w - \widehat{R}_t^l, \quad (4)$$

**Log-likelihood ratio against a reference policy.** Let the frozen reference be $\pi_{\text{ref}}(a_t \mid s_t) = p_{\theta_{\text{ref}}}(x_{t-1} \mid x_t, c)$. For $j \in \{w, l\}$ we define

$$\rho_t^j(\theta) = \log \pi_\theta\left(x_{t-1}^j \mid s_t\right) - \log \pi_{\text{ref}}\left(x_{t-1}^j \mid s_t\right), \quad (5)$$

and use the win–lose difference $\Delta \rho_t = \rho_t^w(\theta) - \rho_t^l(\theta)$ to measure the policy's preference change relative to the reference.
**Stepwise alignment objective.** Inspired by [15], we align the log-ratio difference with the dense reward difference via a mean-squared error with a time weight $\beta_t = \lambda^{T-t-1}/\eta$ ($\lambda \in (0, 1]$, $\eta > 0$):

$$\mathcal{L}_t(\theta) = \left(\beta_t \Delta \rho_t - \Delta \widehat{R}_t\right)^2. \quad (6)$$

**Final EASPO loss.** To improve sample efficiency, we optimize at a randomly shuffled step $\tau$ and skip the first $\kappa$ high-noise steps. Averaging over prompts $c$, initial noises $x_T$, and win/lose pairs from the policy gives

$$\mathcal{L}(\theta) = \mathbb{E}_{c \sim p(c), x_T \sim \mathcal{N}(0, I), \tau \sim \mathcal{U}[1, T-\kappa], x_{\tau-1}^w, x_{\tau-1}^l \sim p_\theta(\cdot \mid x_\tau, c)}$$
$$\left[\left(\left(\beta_\tau \left[\rho_\tau^w(\theta) - \rho_\tau^l(\theta)\right] - \left[s(x_{\tau-1}^w, c, \tau) - s(x_{\tau-1}^l, c, \tau)\right]\right)^2\right], \right.$$
$$(7)$$

**Discussion.** Eq. (7) integrates stepwise preference optimization with reward-difference learning: EASPM provides a *dense emotional reward* (Eq. (4)), and the diffusion policy is updated so that its *log-likelihood ratio difference* between win/lose transitions (Eq. (5)) matches that reward difference (Eq. (6)). Empirical results from [15] show that using reward differences yields more stable reward optimization in diffusion models than standard policy gradient methods.
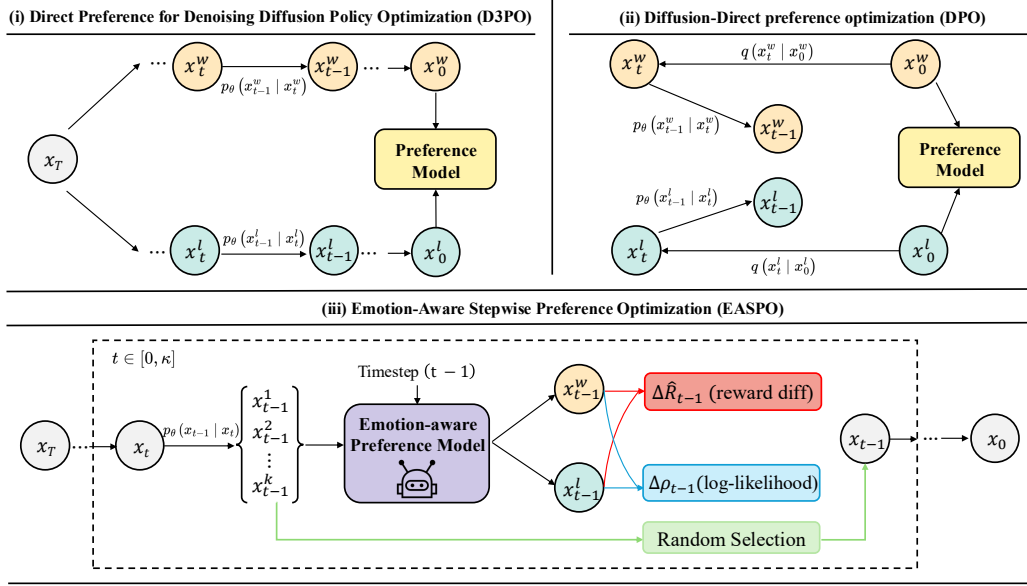
**Fig. 1**. Unlike prior DPO-based methods, EASPO avoids direct propagation of preferences across diffusion steps. At each step in EASPO, a set of candidate samples is produced, from which a suitable win–lose pair is chosen to update the diffusion model. Afterward, one sample is randomly picked to serve as the starting point for the next iteration.

**Table 1**. Objective comparison of our approach with other emotion-controllable TTS models in terms of emotion similarity, prosody similarity, WER, and UTMOS on ESD dataset.

| TTS Model | Emo_SIM↑ | Prosody_SIM↑ | WER↓ | UTMOS↑ |
|---|---|---|---|---|
| FG-TT [1] | 93.91 | 3.28 | 9.38 | 3.81 |
| PromptTTS [2] | 95.70 | 3.41 | **3.25** | 4.33 |
| Emospeech [3] | 96.35 | 3.39 | 7.13 | 4.24 |
| EmoDiff [4] | 96.62 | 3.55 | 5.62 | 4.35 |
| CosyVoice [5] | 97.07 | 3.64 | 4.32 | 4.41 |
| CosyVoice2 [6] | 98.47 | <u>3.78</u> | 3.83 | <u>4.43</u> |
| EmoVoice [7] | <u>98.59</u> | 3.67 | 4.16 | 4.39 |
| **Ours** | **99.15** | **3.89** | <u>3.74</u> | **4.47** |

**Table 2**. Subjective evaluation of naturalness, emotional expressiveness , emotion consistency, and emotion recall, evaluated by human raters.

| TTS Model | MOS↑ | Emo_MOS↑ | MOS_EC↑ | Recall↑ |
|---|---|---|---|---|
| PromptTTS [2] | 2.95 | 2.88 | 2.72 | 74.12 |
| EmoDiff [4] | 3.28 | 3.36 | 3.40 | 78.59 |
| CosyVoice2 [6] | <u>3.63</u> | 3.71 | <u>3.83</u> | <u>82.10</u> |
| EmoVoice [7] | 3.56 | <u>3.79</u> | 3.64 | 80.36 |
| **Ours** | **3.94** | **4.28** | **4.04** | **85.84** |

## 3. EXPERIMENTS

### 3.1. Datasets and Experimental Setup

We fine-tune EASPM on the English MSP-Podcast corpus [16] (∼55k utterances, >1,200 speakers), using textual prompts from emotion labels and acoustic descriptors (e.g., pitch, loudness, jitter, shimmer). Preference pairs are created by labeling target emotions (e.g., happy) as preferred over distractors (e.g., neutral) with the same text. For reinforcement learning, we use the English split of ESD (5 emotions × 10 speakers, 350 utterances/emotion), with an 8:1:1 train/val/test split per speaker–emotion. We evaluate against seven emotion-controllable TTS baselines: FG-TTS [1], PromptTTS [2], Emospeech [3], EmoDiff [4], CosyVoice [5], CosyVoice2 [6], and EmoVoice [7], using authors' code and checkpoints with default inference settings.

We initialize EASPM from CLEP. Audio is resampled to 16 kHz and cropped/padded to 5 s. The text encoder is frozen, while the audio encoder and projection head are trained using Adam (batch size 64, 80 epochs), with learning rates $1 \times 10^{-5}$ and $1 \times 10^{-3}$, respectively. Preference supervision is derived from emotion-labeled recordings, marking the target emotion as preferred and another as dis-preferred. To make scoring step-aware, both waveforms in a pair are perturbed by identical diffusion noise at a sampled denoising step. Our base TTS model is Grad-TTS with 80-dim mel-spectrograms. We freeze the encoder and duration predictor and fine-tune only the decoder (score network), initialized from Grad-TTS pre-training settings (Adam, $1 \times 10^{-4}$ LR, batch size 16, random 2 s mel segments). During EASPO, we apply step shuffling and candidate pooling. A denoising step $\tau \sim \mathcal{U}(1, T - \kappa)$ is chosen (skipping noisy early steps), $k=4$ candidates are sampled from $p_\theta(\cdot \mid x_\tau, c)$, ranked by EASPM, and a win–lose pair is used to minimize the stepwise difference–matching loss between policy log-ratio difference and reward difference. Unless specified, we set $\kappa=0.25T$, batch size 32, and decoder learning rate $1 \times 10^{-5}$. Waveforms are synthesized using a pretrained HiFi-GAN vocoder [17].

### 3.2. Evaluation Metrics

**Objective metrics.** *Emo_SIM* quantifies emotional alignment as the average cosine similarity between emotion2vec-base embeddings of generated and reference utterances. *Prosod_SIM* is computed via AutoPCP [18], which com-

**Table 3**. Comparing EASPM with variants: no time condition.

| Prefer. model | E-S | P-S | WER | UTMOS |
|---|---|---|---|---|
| EASPM | **99.15** | **3.89** | **3.74** | **4.47** |
| w/o step con. | 98.79 | 3.81 | 3.83 | 4.36 |
| CLAP | 95.84 | 3.36 | 3.96 | 4.05 |

**Table 4**. Comparing random sampling with other sampling strategies.

| Initial. | E-S | P-S | WER | UTMOS |
|---|---|---|---|---|
| $x_{t-1}^w$ | 97.78 | 3.63 | 3.81 | 4.20 |
| $x_{t-1}^l$ | 98.39 | 3.75 | 3.79 | 4.33 |
| random | **99.15** | **3.89** | **3.74** | **4.47** |

**Table 5**. Impact of number of sampled images $k$ at each step. We use $k = 4$.

| #samples $k$ | E-S | P-S | WER | UTMOS |
|---|---|---|---|---|
| 2 | 98.31 | 3.76 | 3.78 | 4.23 |
| 4 | **99.15** | 3.89 | 3.74 | **4.47** |
| 8 | 98.84 | **3.93** | **3.71** | 4.27 |

**Table 6**. Impact of timestep range.

| Timestep Range | E-S | P-S | WER | UTMOS |
|---|---|---|---|---|
| [0–250] | 98.16 | 3.35 | 3.81 | 4.14 |
| [0–500] | 98.79 | 3.62 | 3.76 | 4.39 |
| [0–750] | **99.15** | **3.89** | 3.74 | **4.47** |
| [0–1000] | 97.92 | 3.57 | **3.69** | 4.26 |
| [250–750] | 98.85 | 3.81 | 3.71 | 4.42 |
| [500–750] | 98.34 | 3.69 | 3.75 | 4.37 |
| [250–500] | 98.68 | 3.77 | 3.73 | 4.28 |

**Table 7**. Comparison with other diffusion based RL alignment methods.

| Method | E-S | P-S | WER | UTMOS |
|---|---|---|---|---|
| Vanilla-DM | 96.62 | 3.55 | 5.62 | 4.35 |
| DDPO | 98.37 | 3.63 | 4.07 | 4.41 |
| D3PO | 97.51 | 3.59 | 4.41 | 4.40 |
| Diff.-DPO | 97.85 | 3.67 | 3.82 | 4.37 |
| EASPO | **99.15** | **3.89** | **3.74** | **4.47** |

**Table 8**. Comparison of win–lose pair selection strategies. Using candidates with highest and lowest emotional scores yields stronger contrast and better alignment than random sampling.

| win-lose sample | E-S | P-S | WER | UTMOS |
|---|---|---|---|---|
| best & worst SPO | **99.15** | **3.93** | **3.74** | **4.47** |
| random | 98.82 | 3.89 | 3.95 | 4.36 |

pares utterance-level prosody (rhythm, stress, intonation). *WER* (word-error-rate) is calculated using Whisper Large-v3 transcripts. UTMOS [19] is employed to evaluate speech naturalness and perceptual quality.

**Subjective metrics.** We conduct listening tests with 20 raters. Each system is evaluated on 30 clips per rater (six per emotion across five emotions). *MOS* assesses naturalness, and *Emotion MOS* evaluates how well the target emotion is conveyed given the prompt, both on a 1–5 scale (0.5 increments). *MOS_EC* measures consistency between the generated audio and the instruction (emotion + text). For *Emotion Recall*, raters identify the perceived emotion from five choices; accuracy is averaged over both utterances and emotion classes.

### 3.3. Main Results

To validate the effectiveness of our proposed emotionally preference-aligned method for TTS, we compare it with seven recent emotion-controllable baselines. As shown in Table 1, our method achieves superior performance across emotion similarity, prosody similarity, intelligibility (WER), and perceptual quality (UTMOS), with a notable 2.07% gain over CosyVoice in Emo_SIM. A similar trend is observed in Table 2, where our model consistently outperforms baselines in naturalness (MOS), emotional expressiveness (Emo_MOS), emotion consistency (MOS_EC), and emotion classification accuracy (Recall). These results demonstrate the effectiveness of our approach in generating emotionally aligned speech with enhanced coherence, while preserving natural prosody and speech quality. Demo page is available .

### 3.4. Ablation Study

**Effectiveness of the Emotion-aware Stepwise Preference Model.** We ablate timestep conditioning and CLEP initialization to assess their roles in Tab. 3. Removing either results in consistent performance drops, confirming that both components are essential for accurate step-aware emotional scoring.
**Random Selection for Next Iteration Initialization.** We compare random selection from the candidate pool with

reusing the previous lose sample $x_{t-1}^l$ to initialize the next denoising step in Tab. 4. Random selection consistently improves overall performance by avoiding bias toward dispreferred regions and encouraging trajectory diversity.
**Impact of Candidate Pool Size.** We vary the number of candidates k at each step and observe its effect in Tab. 5. A moderate k balances contrastive supervision and fidelity, improving the learning of emotional and prosodic preferences, while large k introduces artifacts that weaken supervision.
**Impact of Timestep Range.** We apply EASPO on a subset of denoising steps $[0, \kappa]$ in Tab. 6 and find that skipping noisy early steps improves alignment due to limited speech structure, while omitting fine-grained late steps weakens prosodic refinement. A mid-range window (e.g., $[0, 750]$) balances diversity and emotional clarity, yielding optimal performance.
**Comparison with other diffusion-based RL alignment methods.** We compare EASPO to prior diffusion-based RL methods (DDPO, D3PO, Diff-DPO) in Tab. 7, observing consistent metric gains that highlight EASPO's strength in aligning fine-grained emotional and prosodic preferences.
**Choice of Win/Lose Pairs.** We compare selecting the top–bottom EASPM-scored candidates versus randomly sampled pairs from the same latent state in Fig. 8. Using highest–lowest scoring samples ensures stronger emotional contrast while maintaining comparable content and noise levels, leading to more stable and informative supervision.

## 4. CONCLUSION

We present EASPO, a diffusion-based speech synthesis framework that introduces stepwise preference optimization for fine-grained emotional alignment. By leveraging an emotion-aware scoring model to compare candidate samples at each denoising step, EASPO progressively guides generation toward emotionally expressive and prosodically natural speech. This step-conditioned training strategy enables the model to capture subtle affective cues through contrastive supervision. Extensive experimental demonstrate the effectiveness across both objective and subjective evaluations.

# 5. REFERENCES

[1] Li-Wei Chen and Alexander Rudnicky, "Fine-grained style control in transformer-based text-to-speech synthesis," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 7907–7911.

[2] Zhifang Guo, Yichong Leng, Yihan Wu, Sheng Zhao, and Xu Tan, "Promptts: Controllable text-to-speech with text descriptions," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[3] Daria Diatlova and Vitaly Shutov, "Emospeech: Guiding fastspeech2 towards emotional text to speech," *arXiv preprint arXiv:2307.00024*, 2023.

[4] Yiwei Guo, Chenpeng Du, Xie Chen, and Kai Yu, "Emodiff: Intensity controllable emotional text-to-speech with soft-label guidance," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[5] Zhihao Du, Qian Chen, Shiliang Zhang, Kai Hu, Heng Lu, Yexin Yang, Hangrui Hu, Siqi Zheng, Yue Gu, Ziyang Ma, et al., "Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised semantic tokens," *arXiv preprint arXiv:2407.05407*, 2024.

[6] Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, et al., "Cosyvoice 2: Scalable streaming speech synthesis with large language models," *arXiv preprint arXiv:2412.10117*, 2024.

[7] Guanrou Yang, Chen Yang, Qian Chen, Ziyang Ma, Wenxi Chen, Wen Wang, Tianrui Wang, Yifan Yang, Zhikang Niu, Wenrui Liu, et al., "Emovoice: Llm-based emotional text-to-speech model with freestyle text prompting," *arXiv preprint arXiv:2504.12867*, 2025.

[8] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn, "Direct preference optimization: Your language model is secretly a reward model," *Advances in neural information processing systems*, vol. 36, 2023.

[9] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik, "Diffusion model alignment using direct preference optimization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8228–8238.

[10] Xiaoxue Gao, Chen Zhang, Yiming Chen, Huayun Zhang, and Nancy F Chen, "Emo-dpo: Controllable emotional speech synthesis through direct preference optimization," in *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2025, pp. 1–5.

[11] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov, "Grad-tts: A diffusion probabilistic model for text-to-speech," in *International conference on machine learning*. PMLR, 2021.

[12] Jiacheng Shi, Yanfu Zhang, and Ye Gao, "Clep-dg: Contrastive learning for speech emotion domain generalization via soft prompt tuning," in *Proc. Interspeech 2025*, 2025, pp. 4498–4502.

[13] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang, "Clap learning audio concepts from natural language supervision," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[14] Jiaming Song, Chenlin Meng, and Stefano Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.

[15] Fei Deng, Qifei Wang, Wei Wei, Tingbo Hou, and Matthias Grundmann, "Prdp: Proximal reward difference prediction for large-scale reward finetuning of diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 7423–7433.

[16] Reza Lotfian and Carlos Busso, "Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings," *IEEE Transactions on Affective Computing*, vol. 10, no. 4, pp. 471–483, 2017.

[17] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," *Advances in neural information processing systems*, vol. 33, pp. 17022–17033, 2020.

[18] Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, et al., "Seamless: Multilingual expressive and streaming speech translation," *arXiv preprint arXiv:2312.05187*, 2023.

[19] Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi Saruwatari, "Utmos: Utokyo-sarulab system for voicemos challenge 2022," *arXiv preprint arXiv:2204.02152*, 2022.