

# TED-TTS: Training-Free Intra-Utterance Emotion and Duration Control for Text-to-Speech Synthesis

Anonymous ACL submission

## Abstract

While controllable Text-to-Speech (TTS) has achieved notable progress, most existing methods remain limited to inter-utterance-level control, making fine-grained intra-utterance expression challenging due to their reliance on non-public datasets or complex multi-stage training. In this paper, we propose TED-TTS, a training-free controllable framework for pretrained zero-shot TTS to enable intra-utterance emotion and duration expression. Specifically, we propose a segment-aware emotion conditioning strategy that combines causal masking with monotonic stream alignment filtering to isolate emotion conditioning and schedule mask transitions, enabling smooth intra-utterance emotion shifts while preserving global semantic coherence. Based on this, we further propose a segment-aware duration steering strategy to combine local duration embedding steering with global EOS logit modulation, allowing local duration adjustment while ensuring globally consistent termination. To eliminate the need for segment-level manual prompt engineering, we construct a 30,000-sample multi-emotion and duration-annotated text dataset to enable LLM-based automatic prompt construction. Extensive experiments demonstrate that our training-free method not only achieves state-of-the-art intra-utterance consistency in multi-emotion and duration control, but also maintains baseline-level speech quality of the underlying TTS model. Audio samples are available at <https://aclanonymou111.github.io/TED-TTS-DemoPage/>.

## 1 Introduction

Humans naturally regulate emotional expression and speaking pace during speech in a dynamic and flexible manner, reflecting changes in semantics, emphasis, and discourse intent. How to replicate such intra-utterance expressiveness remains a central challenge in building human-like Text-to-Speech (TTS) synthesis systems.

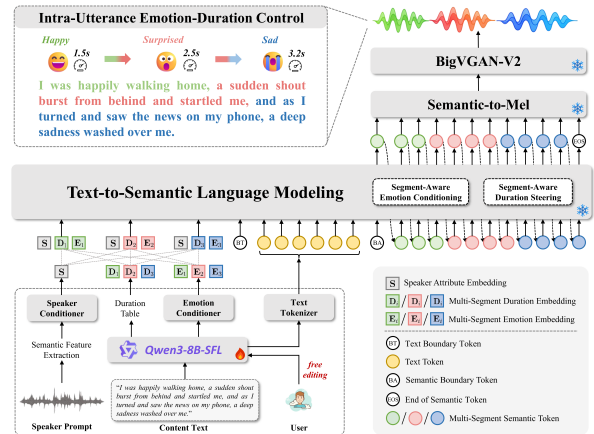


Figure 1: Overview of our training-free framework for intra-utterance emotion and duration control, where the green, red, and blue regions denote three segments with different emotion and duration settings within the same utterance.

Recent advances in controllable TTS have enabled zero-shot synthesis conditioned on attributes such as speaker identity, emotion, and speaking rate (Du et al., 2024b; Wang et al., 2025c,b; Chen et al., 2025; Gao et al., 2025; Yang et al., 2025; Zhou et al., 2025). Despite these advances, controllability in most existing methods remains confined to the utterance level, where a single emotional or prosodic condition is uniformly applied to an entire utterance, deviating from the dynamic expression naturally observed in human speech. To address this limitation, some methods (Luo et al., 2021; Tan et al., 2024) predict phoneme- or frame-level affective attributes directly from text, while others (Kanda et al., 2024; Wu et al., 2024) rely on emotional reference speech to guide localized expressive patterns, such as brief laughter or crying. Most recently, WeSCon (Wang et al., 2025a) proposes a self-training framework with transition smoothing and emotional-bias mechanisms, enabling the TTS model to render multiple emotions within an utterance through distillation. While these meaningful

066	progress, they typically rely on large-scale time-	specific emotion conditioning and perform on-	117
067	aligned annotated speech datasets or involve multi-	line text-semantic alignment, enabling stable	118
068	stage training pipelines, which substantially limit	multi-emotion transitions within a single ut-	119
069	their cross-model transferability and real-world de-	terance.	120
070	ployment.		
071	These challenges naturally raise an important	• We propose a segment-aware duration steering	121
072	question: <i>Is it possible to achieve stable segment-</i>	strategy that achieves local segment duration	122
073	<i>level emotion transitions and duration control</i>	control while preserving global consistent se-	123
074	<i>without retraining the model?</i> In this paper,	quence termination.	124
075	as shown in Fig. 1, we revisit controllable TTS		
076	from an inference-time perspective and propose	• Extensive experiments demonstrate that our	125
077	the first Training-free Intra-Utterance Emotion and	training-free method not only achieves state-	126
078	Duration control framework (TED-TTS). Rather	of-the-art intra-utterance consistency for	127
079	than introducing additional predictors or retrain-	multi-emotion and duration control, but also	128
080	ing the acoustic model, our approach focuses on	maintains the baseline-level speech quality of	129
081	restructuring how conditioning information is ac-	the underlying TTS model.	130
082	cessed and updated during autoregressive decoding.		
083	Specifically, for multi-emotion control, we propose	<b>2 Related Work</b>	131
084	a segment-aware emotion conditioning strategy that		
085	combines causal masking with monotonic stream	<b>2.1 Emotionally Controllable TTS</b>	132
086	alignment filtering algorithm, which jointly isolates	Emotion-controllable TTS methods can be broadly	133
087	segment-specific emotion conditioning and per-	categorized by the modality of emotion prompts.	134
088	forms online text-semantic alignment to schedule	<b>Speech-prompt-based methods</b> condition synthe-	135
089	mask transitions, enabling smooth intra-utterance	sis on reference emotional utterances and can trans-	136
090	emotion shifts while preserving global semantic	fer fine-grained affective cues such as intensity and	137
091	coherence. To enable multi-duration control, we	prosody (Eskimez et al., 2024; Du et al., 2024a;	138
092	further propose a segment-aware duration steering	Wang et al., 2025c,b; Chen et al., 2025), but their	139
093	strategy to incorporate local duration embedding	reliance on reference speech limits practical flex-	140
094	steering with global EOS logit modulation, allow-	ibility. In contrast, <b>Text-prompt-based methods</b>	141
095	ing segment-level pacing adjustment while ensur-	offer more flexible control, where early approaches	142
096	ing globally consistent sequence termination. Be-	rely on discrete emotion labels (Guo et al., 2023a;	143
097	sides, we construct a multi-emotion and duration-	Kang et al., 2023; Diatlova and Shutov, 2023; Tang	144
098	annotated text dataset with 30,000 samples and fine-	et al., 2024; Gao et al., 2025), while recent meth-	145
099	tune Qwen3-8B to enable LLM-based automatic	ods adopt natural language emotion descriptions	146
100	prompt construction, thereby eliminating the need	for richer and more continuous conditioning (Guo	147
101	for segment-level manual segmentation and prompt	et al., 2023b; Liu et al., 2023; Yang et al., 2024;	148
102	engineering. Extensive experiments demonstrate	Du et al., 2024b; Yang et al., 2025; Zhou et al.,	149
103	that our method achieves state-of-the-art perfor-	2025). However, these methods typically operate	150
104	mance in stable intra-utterance multi-emotion tran-	at the utterance level, assigning a single global	151
105	sitions and duration control, while preserving the	emotion to an entire utterance and thus failing to	152
106	strong zero-shot synthesis capability of the under-	capture intra-utterance emotional dynamics. To	153
107	lying TTS model without any additional training.	address this limitation, several <b>Intra-utterance</b>	154
108	Our contributions are summarized as follows:	<b>control methods</b> predict fine-grained affective at-	155
109		tributes directly from text (Im et al., 2022; Luo	156
110	• We propose a training-free controllable frame-	et al., 2021; Tan et al., 2024), or incorporate emo-	157
111	work for intra-utterance-level TTS, and elim-	tional reference speech to enable localized expres-	158
112	inate manual prompt engineering by construc-	sions such as laughter or crying (Kanda et al., 2024;	159
113	ing a 30,000-sample multi-emotion and	Wu et al., 2024). More recently, WeSCon (Wang	160
114	duration-annotated text dataset for LLM-	et al., 2025a) introduces a self-training framework	161
115	based automatic prompt construction.	to support multi-emotion rendering within a single	162
116		utterance. Despite these advances, existing meth-	163
	• We propose a segment-aware emotion con-	ods often rely on large-scale non-public emotional	164

165 datasets or multi-stage training pipelines that hinder  
166 their scalability and cross-model transferability,  
167 leaving training-free intra-utterance emotion control  
168 as an open and practically valuable challenge.

## 169 2.2 Duration Controllable TTS

170 Current exploration on duration control has advanced  
171 along both non-autoregressive and autoregressive  
172 approaches. **Non-autoregressive methods** achieve  
173 duration control via explicit duration predictors based  
174 on diffusion-transformers (Lee et al., 2025), flows  
175 (Kim et al., 2023), or language models (Du et al.,  
176 2025), but these predictors are trained separately  
177 and often struggle with temporal accuracy under  
178 prosodic variability. In contrast, **autoregressive  
179 methods** lack inherent duration control and typically  
180 rely on auxiliary cues, such as natural-language  
181 timing prompts (Zhou et al., 2024) or specialized  
182 attributes and labels (Li et al., 2025; Sahipjohn  
183 et al., 2024; Wang et al., 2025b). More recently,  
184 IndexTTS2 (Zhou et al., 2025) improves controllability  
185 by conditioning semantic token generation on duration  
186 positional embeddings, enabling more stable alignment  
187 between desired and produced token lengths than  
188 earlier autoregressive methods. However, existing  
189 approaches still struggle to decouple local pacing  
190 from global generation, failing to provide a unified  
191 framework that ensures stable intra-utterance  
192 duration control without compromising overall  
193 alignment.

## 194 2.3 Inference-Time Controllable TTS

195 Several approaches have explored inference-time  
196 controllable TTS, enabling flexible manipulation of  
197 speech attributes. EmoKnob (Chen et al., 2024) injects  
198 scaled emotion difference vectors into speaker  
199 embeddings for emotion control, while PRESENT  
200 (Lam et al., 2025) performs rule-based prosody  
201 shaping by adjusting pitch, duration, and energy  
202 predictions from text prompts. SPTTS (Sun et al.,  
203 2025) further operates in the latent embedding  
204 space, manipulating prosody and style directions  
205 derived via linear regression and vector arithmetic.  
206 More recently, EmoSteer-TTS (Xie et al., 2025)  
207 directly steers token-level activations in pretrained  
208 diffusion-based TTS models, enabling training-free  
209 emotion control with improved interpretability over  
210 global embedding methods. However, these methods  
211 predominantly focus on implicit latent manipulation  
212 or isolated feature editing, and lack a unified  
213 framework for jointly controlling segment-level  
214 emotion and pacing transitions.

## 215 3 Method

216 In this section, we introduce a training-free control-  
217 able framework for intra-utterance emotion and  
218 duration transitions, with details provided in the  
219 following subsections.

### 220 3.1 Automatic Prompt Construction

221 Existing intra-utterance controllable TTS systems  
222 require manual text segmentation and segment-  
223 level emotion and duration specification, which  
224 is labor-intensive in real-world scenarios. To  
225 eliminate manual prompt engineering, we fine-  
226 tune the Qwen3-8B LLM to automatically transform  
227 raw user text into structured multi-segment  
228 prompts. As a prerequisite, we construct a dedicated  
229 **Multi-Emotion and Duration-annotated text dataset**  
230 (MED-TTS) with 30,000 samples, which is used to  
231 supervise emotion-aware text segmentation,  
232 natural language emotion description generation,  
233 and segment-level speech duration estimation. As  
234 illustrated below, MED-TTS is synthesized using  
235 LLM through a structured pipeline consisting of  
236 generation, annotation, and verification.

**Step 1: Content text generation.** GPT-4o is  
237 prompted to generate emotion-rich English and  
238 Chinese texts with continuous intra-utterance  
239 emotional transitions spanning multiple phases from  
240 seven core emotions (*happy, sad, angry, surprised,  
241 fearful, disgusted, and neutral*), covering descriptive,  
242 dialog-style, and observational content in three  
243 text categories.

**Step 2: Multi-segment prompt annotation.** To  
244 enable precise segment-level control, DeepSeek-  
245 Chat is prompted to decompose each text into  
246 emotion-specific segments and assigns each segment  
247 a concise natural language emotion description  
248 and a realistic duration estimate, yielding structured  
249 emotion-duration sequences compatible with  
250 controllable TTS inputs.

**Step 3: Post-processing and manual verification.**  
251 Automatic checks are finally applied to filter  
252 samples with formatting errors, missing fields, or  
253 invalid segment boundaries, followed by systematic  
254 manual verification of outputs from both Step 1  
255 and Step 2. Based on this dataset, we perform  
256 supervised fine-tuning with LoRA on the Qwen3-8B  
257 large language model, enabling automatic construction  
258 of segment-level TTS prompts without manual  
259 prompt engineering. Detailed prompting strategies,  
260 step-wise checklists, dataset statistics, and fine-  
261 tuning details are provided in the Appendix A.

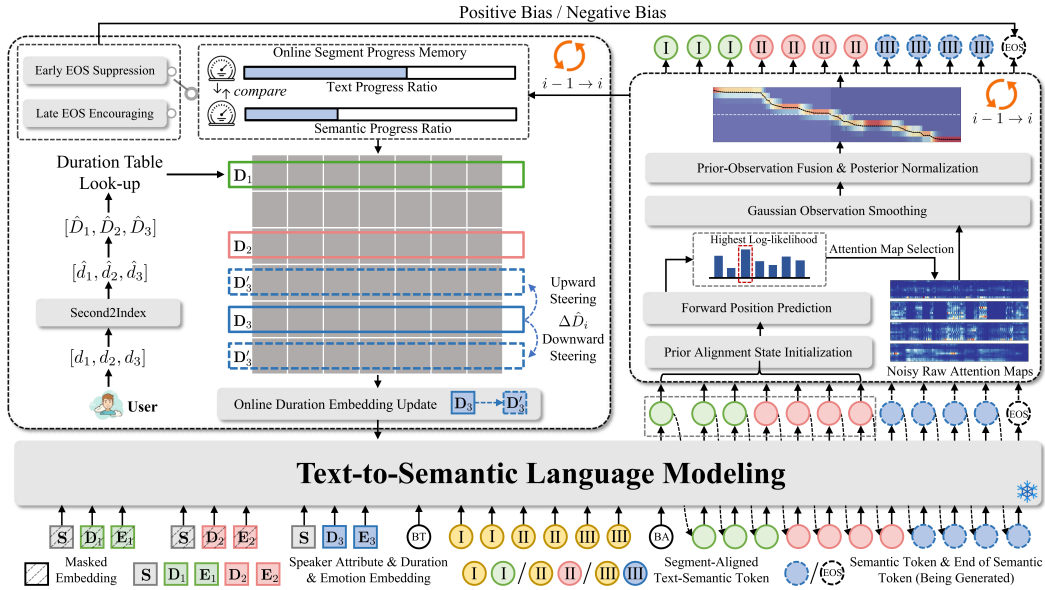


Figure 2: Overview of our training-free framework for fine-grained intra-utterance emotion and duration control, illustrating the transition from the second (red) segment to the third (blue) segment via segment-aware duration steering (left) and segment-aware emotion conditioning (right) strategy.

### 3.2 Segment-Aware Emotion Conditioning

Our TTS architecture follows the same configuration as the IndexTTS2 (Zhou et al., 2025) baseline, and we focus our design on its text-to-semantic (T2S) module to enable training-free intra-utterance emotion control. Specifically, T2S is formulated as an autoregressive semantic token prediction task conditioned on text and a set of control embeddings. Given an input text, represented by the yellow tokens in Fig. 2, we decompose it into  $M$  user-defined segments  $\mathbf{X} = \{X_1, X_2, \dots, X_M\}$ , and each segment  $X_m$  is assigned a condition embedding  $\mathbf{C}_m = \{\mathbf{I}, \mathbf{E}_m\}$ , where  $\mathbf{I}$  denotes a fixed speaker identity embedding shared across the segment, and  $\mathbf{E}_m$  represents segment-specific emotion conditions. However, in autoregressive T2S formulations, semantic tokens are generated as a continuous stream without explicit segment boundaries, making it non-trivial to apply segment-level conditions to their corresponding text segments while preserving semantic continuity. To address this challenge, we propose a 2D causal attention mask combined with a monotonic stream alignment algorithm to enable smooth intra-utterance emotion transitions.

**2D Causal Attention Mask.** To resolve the misalignment between continuous generation and segment-level conditions, we design a 2D causal attention mask that disentangles condition visibility from semantic context. The mask preserves

standard causal attention among text and semantic tokens across segment boundaries, ensuring globally coherent semantic generation while strictly restricting access to condition embeddings to be segment-local. Specifically, for any token that belongs to the  $m$ -th segment (either a text token in  $X_m$  or a generated semantic token that currently aligns to  $X_m$ ), attention is allowed to attend only to its corresponding condition embedding  $\mathbf{C}_m$ , while all other condition embeddings  $\{\mathbf{C}_j \mid j \neq m\}$  are masked out. Meanwhile, each condition embedding  $\mathbf{C}_m$  is prevented from attending to other condition embeddings, avoiding cross-condition information leakage. After that, as shown in the bottom of Fig. 2, emotional style is governed exclusively by the locally active condition, whereas semantic content remains globally visible through standard causal context.

However, applying 2D causal attention masks requires real-time knowledge of the alignment between generated semantic tokens and source text tokens. While transformer attention can provide alignment cues, raw attention maps are often noisy, head-dependent, and non-monotonic, making them unreliable for driving mask transitions. To address this, we propose an online Monotonic Stream Alignment (MSA) algorithm that performs Bayesian-style alignment tracking using attention as observation.

**Monotonic Stream Alignment (MSA).** As shown in Fig. 3, we use  $\mathbf{A}_i \in \mathbb{R}^{L \times H \times T}$  denote

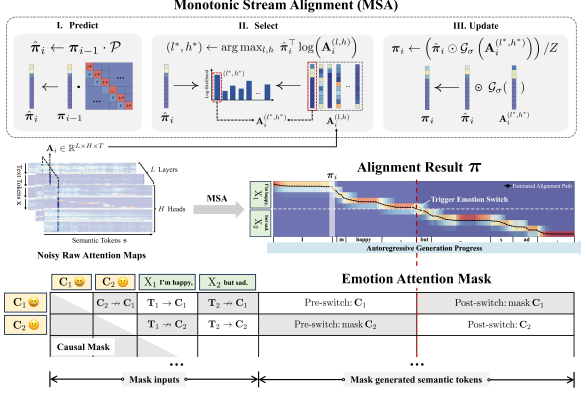


Figure 3: Detailed illustration of Monotonic Stream Alignment (MSA) in segment-aware emotion conditioning, where from top to bottom are MSA algorithm, MSA alignment result, and the visualization of 2D causal attention mask, respectively.

the raw attention maps from the current semantic token  $s_i$  to the  $T$  text tokens across  $L$  layers and  $H$  heads, where  $\mathbf{A}_i^{(l,h)}$  is the attention vector of head  $(l, h)$ . During online autoregressive decoding, MSA maintains a belief distribution over text positions to track the alignment of  $s_i$ , represented by a prior distribution  $\hat{\pi}_i$  and a posterior distribution  $\pi_i$ , both defined over the  $T$  text tokens. At each decoding step  $i$ , MSA first performs the *Predict* step by propagating the posterior  $\pi_{i-1}$  from the previous step forward along the text sequence using a monotonic transition operator  $\mathcal{P}$ . This propagation yields a prior distribution  $\hat{\pi}_i$  that encodes strong temporal monotonicity, encouraging gradual forward movement while suppressing backward alignment. After obtaining the monotonic prior  $\hat{\pi}_i$ , MSA enters the *Select* step to select the most reliable attention head by measuring how well each head’s attention distribution agrees with the predicted alignment:

$$(l^*, h^*) = \arg \max_{l,h} \hat{\pi}_i^\top \log \mathbf{A}_i^{(l,h)}, \quad (1)$$

where  $\mathbf{A}_i^{(l,h)}$  denotes the attention vector of head  $(l, h)$ . The resulting head  $(l^*, h^*)$  provides the most reliable attention observation used in the subsequent update. In the final *Update* step, MSA combines the selected attention observation  $\mathbf{A}_i^{(l^*, h^*)}$  with the monotonic prior  $\hat{\pi}_i$  to compute the posterior alignment belief as:

$$\pi_i = \frac{\hat{\pi}_i \odot \mathcal{G}_\sigma(\mathbf{A}_i^{(l^*, h^*)})}{Z}, \quad (2)$$

where  $\odot$  denotes element-wise multiplication,  $\mathcal{G}_\sigma(\cdot)$  is a Gaussian smoothing operator, and  $Z$  is a

normalization factor. This update incorporates real-time attention evidence while enforcing monotonicity, resulting in a stable alignment trajectory  $\pi_i$ . Benefiting from this alignment trajectory, segment-level causal mask switching is triggered by tracking the expected aligned text position, enabling subsequent semantic tokens to attend to the new segment condition. More detailed mathematical derivations are provided in Appendix B.

### 3.3 Segment-Aware Duration Steering

Beyond segment-level emotional expressiveness, we further extend our emotion control framework to enable multi-segment duration control in a fully training-free autoregressive setting.

**Local Duration Embedding Steering.** Inspired by IndexTTS2 (Zhou et al., 2025), we condition duration control on a dedicated duration embedding indexed by the semantic token length, and tie its embedding table  $\mathbf{W}_{dur}$  with the semantic positional embedding table  $\mathbf{W}_{sem}$  to align autoregressive positional progression with target duration. As shown in Fig.2, given an utterance with  $M$  segments and desired durations  $\mathbf{d} = \{d_1, d_2, \dots, d_M\}$ , each segment duration is converted into the corresponding number of semantic tokens according to the codec token rate (Wang et al., 2025c), yielding  $\hat{\mathbf{d}} = \{\hat{d}_1, \hat{d}_2, \dots, \hat{d}_M\}$ . We accumulate segment-level targets into cumulative token lengths  $\hat{D}_i = \sum_{k=1}^i \hat{d}_k$  and retrieve segment-wise initial duration embeddings as  $\mathbf{D}_i = \mathbf{W}_{dur}[\hat{D}_i]$ , which are concatenated into the segment-level conditioning inputs  $\mathbf{C}_m$  to guide subsequent generation.

During autoregressive decoding, the actual semantic token generation speed may deviate from the user-specified target due to alignment uncertainty and model stochasticity. To correct such deviations online, we introduce a local duration embedding steering mechanism that dynamically updates the duration embedding via adaptive duration table lookup. At each decoding step  $i$ , we leverage MSA (Section 3.2) to estimate the current aligned text position and compute two normalized progress indicators within the active segment: text progress  $r_{text}$  and semantic progress  $r_{sem}$ . Their discrepancy is defined as  $\Delta r = r_{text} - r_{sem}$ , where a positive value indicates lagging semantic generation, which is then used to adjust the effective semantic token length via a proportional controller:

$$\Delta \hat{D}_i = \text{clip}(\lfloor k \cdot \Delta r \rfloor, -\Delta_{\max}, \Delta_{\max}), \quad (3)$$

where  $k$  controls the correction strength,  $\lfloor \cdot \rfloor$  denotes rounding to the nearest integer, and  $\Delta_{\max}$  bounds the maximum adjustment. The effective segment-wise target is updated as  $\hat{D}_i + \Delta \hat{D}_i \rightarrow \hat{D}'_i$ , and the duration table  $\mathbf{W}_{dur}$  is re-queried only for the active segment to obtain the updated duration embedding  $\mathbf{D}'_i$ , while duration embeddings of other segments remain unchanged. For stability, updates are applied at a low temporal frequency, allowing multiple consecutive semantic tokens to share the same duration embedding.

**Global EOS Steering.** In autoregressive decoding, the End-Of-Semantic (EOS) token determines sequence termination and overall duration. While local duration embedding steering regulates local generation pace, it does not explicitly control when decoding ends. To address this, we introduce a global EOS steering strategy that modulates sequence termination by applying adaptive biases to the EOS logit. Specifically, EOS generation is suppressed in all non-final segments to prevent premature termination, and in the final segment, the EOS logit is progressively adjusted based on the remaining semantic budget, discouraging early termination while smoothly encouraging EOS emission as the target budget is approached. Detailed parameter settings are provided in Appendix C.

## 4 Experiments

**Datasets and Comparison Models.** We use the MED-TTS dataset for content text, text-based emotion prompts, and duration annotations, which contains 15,000 English and 15,000 Chinese pair samples. For each language, 500 samples are randomly held out for evaluation, while the remaining samples are used for Qwen3-8B fine-tuning. For identity and emotion speech prompts, we adopt the Emotional Speech Dataset (ESD) (Zhou et al., 2022), where a same-language speaker is fixed per test utterance to ensure timbre consistency, and the speaker’s emotional speech is used as segment-level emotion references.

We compare our method with representative controllable TTS methods spanning both non-autoregressive and autoregressive frameworks. The non-autoregressive models include MaskGCT (Wang et al., 2025c) and F5-TTS (Chen et al., 2025). The autoregressive models include CosyVoice2 (Du et al., 2024b), Spark-TTS (Wang et al., 2025b) and IndexTTS2 (Zhou et al., 2025).

**Evaluation Metrics.** We adopt both objective and

subjective metrics to comprehensively evaluate system performance. Intelligibility is measured by WER for English using Whisper-Large (Radford et al., 2023) and CER for Chinese using Paraformer (Gao et al., 2022). Speaker similarity (S-SIM) is computed as the cosine similarity between WavLM-Large speaker embeddings (Chen et al., 2022). Transition smoothness is assessed using DNSMOS-Pro (DNSM) (Cumlin et al., 2024) over sliding speech segments, while perceptual quality is evaluated with NISQA (Mittag et al., 2021) and OVRL (Reddy et al., 2022). Emotional accuracy is evaluated based on the emotion prompt type, using emotion2vec-Large embeddings (Ma et al., 2024) for speech prompts and a fine-tuned emotion2vec classifier for text prompts. Subjective evaluation is conducted using four MOS criteria: SMOS for speaker similarity, NMOS for the naturalness of emotion transitions, EMOS for emotion alignment, and SPMOS for speaking rate accuracy. All scores are collected on a 5-point scale and reported with mean values and 95% confidence intervals.

## 5 Results and Evaluation

### 5.1 Comparison with Reference Models

**Objective Evaluation.** Since comparative methods lack intra-utterance controllability, all segments are synthesized independently and concatenated for evaluation. Under this setting, we conduct objective evaluations for both emotion and duration control. For emotion control, results are reported under two prompting settings: speech emotion prompts and text emotion prompts. For duration control, emotion is fixed to neutral, and segment-level speech synthesis is evaluated under five duration scaling factors (0.75, 0.875, 1.0, 1.125, and 1.25). As shown in Tab. 1, our method achieves the best overall performance on most objective metrics across both languages and prompting settings, with consistent gains on DNSM and SSIM indicating smoother emotion transitions and improved speaker consistency. Although WER/CER and emotion recognition scores are not always optimal, they remain comparable to the IndexTTS2 baseline, which is expected for a training-free framework. For duration control, as shown in Tab. 2, our method attains the best DNSM, NISQA, and OVRL scores in both languages, reflecting more stable temporal pacing and improved perceptual quality. While some methods achieve higher SSIM, this advantage largely stems from segment-independent synthesis under

	Model	WER/CER↓	DNSM↑	SSIM↑	NISQA↑	OVRL↑	Emo2v↑	SMOS↑	NMOS↑	EMOS↑
<i>Speech Emotion Prompt</i>										
English	MaskGCT	3.520	3.829	0.347	4.475	3.275	<u>0.854</u>	2.96±0.34	2.77±0.28	<u>3.64</u> ±0.24
	F5TTS	2.632	3.674	0.353	4.427	3.330	0.832	3.33±0.36	<u>3.40</u> ±0.32	3.56±0.28
	SparkTTS	<u>2.433</u>	3.456	0.358	4.494	<b>3.404</b>	0.849	<u>3.49</u> ±0.29	<u>3.27</u> ±0.31	3.44±0.29
	CosyVoice2	<b>1.411</b>	3.605	0.402	<u>4.535</u>	3.316	0.831	3.33±0.31	<u>2.87</u> ±0.32	3.31±0.28
	IndexTTS2	2.454	<u>3.871</u>	<u>0.457</u>	4.465	3.304	<b>0.861</b>	3.20±0.36	2.98±0.30	<b>4.07</b> ±0.26
	<b>Ours</b>	2.519	<b>3.925</b>	<b>0.485</b>	<b>4.706</b>	<u>3.395</u>	0.837	<b>4.00</b> ±0.24	<b>4.20</b> ±0.23	3.42±0.30
Chinese	MaskGCT	7.221	3.693	0.350	4.309	3.278	<u>0.814</u>	2.80±0.34	2.33±0.31	3.64±0.27
	F5TTS	10.317	3.314	0.324	3.718	3.228	0.734	3.22±0.36	2.49±0.38	3.13±0.28
	SparkTTS	<b>3.107</b>	3.466	0.382	<u>4.338</u>	<u>3.345</u>	0.807	3.42±0.33	2.87±0.31	<u>3.80</u> ±0.24
	CosyVoice2	<u>3.375</u>	3.306	<u>0.423</u>	4.147	3.313	0.766	3.04±0.35	2.71±0.37	3.29±0.29
	IndexTTS2	4.015	<u>3.694</u>	0.401	4.146	3.289	<b>0.869</b>	<u>3.67</u> ±0.33	<u>3.02</u> ±0.30	<b>3.87</b> ±0.24
	<b>Ours</b>	3.792	<b>3.752</b>	<b>0.470</b>	<b>4.509</b>	<b>3.370</b>	0.724	<b>4.13</b> ±0.23	<b>4.07</b> ±0.30	3.62±0.32
<i>Text Emotion Prompt</i>										
English	CosyVoice2	<b>1.522</b>	3.465	<u>0.453</u>	<u>4.330</u>	<u>3.271</u>	0.303	3.33±0.37	<u>3.73</u> ±0.31	2.53±0.31
	IndexTTS2	2.246	<u>3.543</u>	0.424	4.299	3.216	<b>0.525</b>	<u>3.76</u> ±0.39	3.44±0.35	<u>3.42</u> ±0.29
	<b>Ours</b>	3.038	<b>3.694</b>	<b>0.462</b>	<b>4.569</b>	<b>3.335</b>	<u>0.433</u>	<b>4.04</b> ±0.29	<b>4.22</b> ±0.23	<b>3.64</b> ±0.31
Chinese	CosyVoice2	<b>4.488</b>	3.105	<b>0.477</b>	<u>4.346</u>	<u>3.206</u>	0.222	2.18±0.33	<u>3.56</u> ±0.31	2.84±0.35
	IndexTTS2	6.962	<u>3.212</u>	0.369	4.179	3.169	<b>0.702</b>	<u>3.29</u> ±0.36	2.71±0.33	<u>3.56</u> ±0.31
	<b>Ours</b>	<u>5.893</u>	<b>3.357</b>	<u>0.421</u>	<b>4.407</b>	<b>3.295</b>	<u>0.531</u>	<b>4.07</b> ±0.25	<b>4.04</b> ±0.22	<b>3.84</b> ±0.25

Table 1: Objective and subjective evaluation across different emotion prompt settings. ↓ indicates that lower values are better, while ↑ indicates that higher values are better. Subjective results are evaluated by 15 listeners, with 95% confidence intervals computed using a t-test. The best results are highlighted in **bold**, and the second-best results are underlined.

neutral emotion. In contrast, our method performs multi-segment duration control in a single generation, making SSIM preservation more challenging but better reflecting realistic controllable synthesis scenarios. Overall, these results demonstrate that our training-free framework supports effective intra-utterance emotion and duration control under more challenging settings, while consistently outperforming the baseline and comparative methods on most objective metrics and achieving state-of-the-art transition smoothness.

**Subjective Evaluation.** We report subjective results on SMOS, NMOS, EMOS, and SPMOS in Tab. 1 and 2. Unlike comparative methods that synthesize segments independently, our approach performs one-shot generation with all intra-utterance emotion and duration variations. Despite being training-free and inherently bounded by the baseline model, our framework achieves state-of-the-art or highly competitive performance across most MOS metrics in both emotion and duration control evaluations.

## 5.2 Ablation Study

**Emotion and Duration Control Evaluation.** We evaluate emotion conditioning and duration steering as two segment-level components of our frame-

work. For emotion control, restricting segments to local text only (w/o full-text access) or removing MSA alignment (w/o alignment) degrades expressive quality and cross-segment speaker consistency, as evidenced by reduced DNSM and SSIM in Tab. 3, indicating that full-text access and monotonic alignment primarily contribute to smooth emotional transitions rather than token-level accuracy. For duration control, disabling local steering (w/o local steering) leads to the largest performance drop, while removing global EOS control (w/o global EOS) causes a smaller but consistent degradation, suggesting that local pacing dominates segment-level naturalness and global EOS provides additional stabilization.

**Monotonic Stream Alignment Evaluation.** To evaluate the effectiveness of MSA, we visualize alignment results under different settings in Fig. 4 and report the mean absolute error (MAE) of segment boundary positions. Raw attention maps exhibit diffuse and locally non-monotonic patterns, making greedy alignment highly sensitive to noise and leading to unstable trajectories and frequent segment switching failures. Introducing the monotonic stream constraint alleviates this issue and reduces MAE to 0.216, but residual attention uncertainty still causes instability. By further incorporat-

	Model	WER/CER↓	DNSM↑	SSIM↑	NISQA↑	OVRL↑	SMOS↑	NMOS↑	SPMOS↑
<i>Speech Emotion Prompt</i>									
English	MaskGCT	<u>2.482</u>	<u>3.964</u>	0.539	4.536	3.301	4.00±0.32	3.42±0.38	3.47±0.31
	F5TTS	<b>1.941</b>	3.683	<u>0.543</u>	4.454	<u>3.307</u>	3.76±0.32	3.02±0.42	3.24±0.32
	IndexTTS2	2.597	3.899	<b>0.575</b>	<u>4.604</u>	3.273	3.89±0.33	<u>3.87±0.32</u>	<b>3.67</b> ±0.37
	<b>Ours</b>	3.227	<b>3.988</b>	0.532	<b>4.766</b>	<b>3.336</b>	<b>4.22</b> ±0.22	<b>4.20</b> ±0.25	<u>3.62</u> ±0.34
Chinese	MaskGCT	8.140	3.711	<b>0.614</b>	<u>4.366</u>	3.167	3.31±0.40	2.60±0.39	<u>3.02</u> ±0.31
	F5TTS	9.004	3.386	<u>0.598</u>	4.286	3.204	<u>3.82</u> ±0.35	2.59±0.40	2.89±0.36
	IndexTTS2	<b>1.623</b>	<u>3.715</u>	0.597	4.345	<u>3.248</u>	3.76±0.34	<u>3.27</u> ±0.34	2.84±0.38
	<b>Ours</b>	<u>2.732</u>	<b>3.803</b>	0.578	<b>4.536</b>	<b>3.291</b>	<b>3.98</b> ±0.28	<b>4.16</b> ±0.27	<b>3.62</b> ±0.30

Table 2: Objective and subjective evaluation on different duration scaling settings. ↓ indicates that lower values are better, while ↑ indicates that higher values are better. Subjective results are evaluated by 15 listeners, with 95% confidence intervals computed using a t-test. The best results are highlighted in **bold**, and the second-best results are underlined.

Method	WER/CER↓	DNSM↑	SSIM↑	NISQA↑
<i>Segment-aware Emotion Conditioning</i>				
<b>Ours</b>	2.519	<b>3.925</b>	<b>0.485</b>	<b>4.706</b>
w/o full-text access	2.409	3.855	0.449	4.578
w/o alignment	<b>2.043</b>	3.831	0.442	4.639
<i>Segment-aware Duration Steering</i>				
<b>Ours</b>	<b>3.227</b>	<b>3.988</b>	<b>0.460</b>	<b>4.766</b>
w/o local steering	3.861	3.032	0.437	4.750
w/o global EOS	3.513	3.885	0.451	4.717

Table 3: Ablation study of our segment-aware emotion conditioning and duration steering modules.

Method	*0.75	*0.875	*1	*1.125	*1.25
<b>Ours</b>	3.387	<b>1.704</b>	<b>3.218</b>	<b>3.210</b>	<b>3.211</b>
w/o local steering	3.728	3.203	5.670	8.179	11.594
w/o global EOS	<b>1.941</b>	2.404	5.638	7.650	9.158
Baseline	5.778	6.912	7.100	8.232	12.032

Table 4: Average semantic token number error rate (%) across segments for duration control under different settings. Lower indicates better duration accuracy.

ing the observation component, MSA effectively suppresses alignment uncertainty, enforces smooth monotonic trajectories, and reduces MAE to 0.157, yielding precise emotion transitions closely aligned with the ground-truth boundaries.

**Duration-Specified Evaluation.** We evaluate duration-specified speech synthesis under five segment-level scaling factors (\*0.75, \*0.875, \*1.0, \*1.125, and \*1.25), comparing our full system with ablated variants and an IndexTTS2 baseline. As shown in Tab. 4, our method consistently achieves the lowest semantic token number error across all settings, reducing the error by 3.53% and 2.41% on average compared to variants without local steering and global EOS control, respectively. Relative to the baseline without explicit duration control, our approach further yields a 5.07% average error

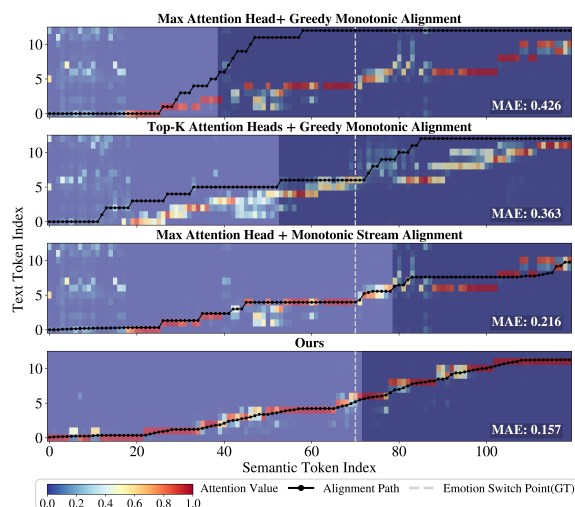


Figure 4: Visualization of alignment paths and emotion switching, with background shading denoting emotion segments and lower MAE indicating better alignment.

reduction, demonstrating accurate and robust duration control across diverse segment-level targets.

## 6 Conclusion

In this paper, we propose the first training-free controllable framework to enable intra-utterance emotion and duration control in pretrained zero-shot TTS. By introducing segment-aware emotion conditioning and duration steering from an inference-time perspective, our method achieves smooth emotional transitions and temporal regulation within a single utterance. Extensive experiments demonstrate that our method not only delivers state-of-the-art intra-utterance controllability, but also preserves baseline-level speech quality of the underlying TTS model.

## Limitations

Despite its advantages, our proposed training-free framework also has several limitations. First, the framework does not explicitly model gradual emotion transitions between adjacent segments. While segment-aware masking and alignment ensure smooth signal-level continuity, emotional variation is controlled in a segment-wise manner rather than through a continuous emotion trajectory, which may limit the representation of intermediate emotional states. Second, the precision of duration control is influenced by the duration representation learned in the pretrained baseline TTS model. Since our approach operates without parameter updates, the duration embedding may not always support strictly linear or fine-grained timing control, particularly under highly expressive or out-of-domain conditions. Future work will investigate training-free or minimally adaptive strategies to better model continuous emotion evolution and duration precision, while preserving the simplicity and generality of the proposed framework.

## Ethical Considerations

This work involves the use of large language models to generate a synthetic text dataset for Qwen3 fine-tuning and model evaluation, and therefore shares some general characteristics of LLM-based generation, such as occasional variations in factual precision or stylistic expression. All models and datasets used are publicly available and employed under their respective licenses, and no private or personally identifiable speech data is involved. While intra-utterance-level controllable TTS can benefit expressive speech synthesis and human-computer interaction research, high-fidelity speech generation also entails potential risks if misused, such as speaker impersonation or spoofing of voice-based authentication systems. In practical applications, it is important to incorporate appropriate safeguards, including audio watermarking, output traceability, or dedicated detection models, to facilitate the identification of synthesized speech and discourage unintended or malicious misuse.

## References

Haozhe Chen, Run Chen, and Julia Hirschberg. 2024. [Emoknob: Enhance voice cloning with fine-grained emotion control](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language*

*Processing, EMNLP 2024*, pages 8170–8180. Association for Computational Linguistics.

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian Wu, Michael Zeng, Xiangzhan Yu, and Furu Wei. 2022. [Wavlm: Large-scale self-supervised pre-training for full stack speech processing](#). *IEEE J. Sel. Top. Signal Process.*, 16(6):1505–1518.

Yushen Chen, Zhikang Niu, Ziyang Ma, Keqi Deng, Chunhui Wang, Jian Zhao, Kai Yu, and Xie Chen. 2025. [F5-TTS: A fairytaler that fakes fluent and faithful speech with flow matching](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, *ACL 2025*, pages 6255–6271. Association for Computational Linguistics.

Fredrik Cumlin, Xinyu Liang, Victor Ungureanu, Chandan K. A. Reddy, Christian Schüldt, and Saikat Chatterjee. 2024. [DNSMOS pro: A reduced-size DNN for probabilistic MOS of speech](#). In *25th Annual Conference of the International Speech Communication Association, Interspeech 2024*. ISCA.

Daria Diatlova and Vitalii Shutov. 2023. [Emospeech: guiding fastspeech2 towards emotional text to speech](#). In *12th ISCA Speech Synthesis Workshop, SSW 2023*, pages 106–112. ISCA.

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

Zhihao Du, Changfeng Gao, Yuxuan Wang, Fan Yu, Tianyu Zhao, Hao Wang, Xiang Lv, Hui Wang, Chongjia Ni, Xian Shi, and et al. 2025. [Cosyvoice 3: Towards in-the-wild speech generation via scaling-up and post-training](#). *arXiv preprint arXiv:2505.17589*.

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

Sefik Emre Eskimez, Xiaofei Wang, Manthan Thakker, Canrun Li, Chung-Hsien Tsai, Zhen Xiao, Hemin Yang, Zirun Zhu, Min Tang, Xu Tan, Yanqing Liu, Sheng Zhao, and Naoyuki Kanda. 2024. [E2 TTS: embarrassingly easy fully non-autoregressive zero-shot TTS](#). In *IEEE Spoken Language Technology Workshop, SLT 2024*, pages 682–689. IEEE.

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

695	Zhifu Gao, Shiliang Zhang, Ian McLoughlin, and Zhijie Yan. 2022. <a href="#">Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition</a> . In <i>23rd Annual Conference of the International Speech Communication Association, Interspeech 2022</i> , pages 2063–2067. ISCA.	751
696		752
697		753
698		
699		
700		
701	Yiwei Guo, Chenpeng Du, Xie Chen, and Kai Yu. 2023a. <a href="#">Emodiff: Intensity controllable emotional text-to-speech with soft-label guidance</a> . In <i>IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023</i> , pages 1–5. IEEE.	
702		
703		
704		
705		
706	Zhifang Guo, Yichong Leng, Yihan Wu, Sheng Zhao, and Xu Tan. 2023b. <a href="#">Promptts: Controllable text-to-speech with text descriptions</a> . In <i>IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023</i> , pages 1–5. IEEE.	
707		
708		
709		
710		
711	Chae-Bin Im, Sang-Hoon Lee, Seung-Bin Kim, and Seong-Wan Lee. 2022. <a href="#">EMOQ-TTS: emotion intensity quantization for fine-grained controllable emotional text-to-speech</a> . In <i>IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022</i> , pages 6317–6321. IEEE.	
712		
713		
714		
715		
716		
717	Naoyuki Kanda, Xiaofei Wang, Sefik Emre Eskimez, Manthan Thakker, Hemin Yang, Zirun Zhu, Min Tang, Canrun Li, Chung-Hsien Tsai, and Zhen Xiao. 2024. Making flow-matching-based zero-shot text-to-speech laugh as you like. <i>arXiv preprint arXiv:2402.07383</i> .	
718		
719		
720		
721		
722		
723	Minki Kang, Wooseok Han, Sung Ju Hwang, and Eunho Yang. 2023. <a href="#">Zet-speech: Zero-shot adaptive emotion-controllable text-to-speech synthesis with diffusion and style-based models</a> . In <i>24th Annual Conference of the International Speech Communication Association, Interspeech 2023</i> , pages 4339–4343. ISCA.	
724		
725		
726		
727		
728		
729	Daegyom Kim, Seongho Hong, and Yong-Hoon Choi. 2023. <a href="#">Sc vall-e: Style-controllable zero-shot text to speech synthesizer</a> . <i>arXiv preprint arXiv:2307.10550</i> .	
730		
731		
732		
733	Perry Lam, Huayun Zhang, Nancy F. Chen, Berrak Sisman, and Dorien Herremans. 2025. <a href="#">PRESENT: zero-shot text-to-prosody control</a> . <i>IEEE Signal Process. Lett.</i> , 32:776–780.	
734		
735		
736		
737	Keon Lee, Dong Won Kim, Jaehyeon Kim, Seungjun Chung, and Jaewoong Cho. 2025. <a href="#">Ditto-tts: Diffusion transformers for scalable text-to-speech without domain-specific factors</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025</i> .	
738		
739		
740		
741		
742		
743	Hanzhao Li, Yuke Li, Xinsheng Wang, Jingbin Hu, Qicong Xie, Shan Yang, and Lei Xie. 2025. <a href="#">Flespeech: Flexibly controllable speech generation with various prompts</a> . <i>arXiv preprint arXiv:2501.04644</i> .	
744		
745		
746		
747	Guanghou Liu, Yongmao Zhang, Yi Lei, Yunlin Chen, Rui Wang, Lei Xie, and Zhifei Li. 2023. <a href="#">Promptstyle: Controllable style transfer for text-to-speech with natural language descriptions</a> . In <i>24th Annual Conference of the International Speech Communication Association, Interspeech 2023</i> , pages 4888–4892. ISCA.	751
748		752
749		753
750		
	Xuan Luo, Shinnosuke Takamichi, Tomoki Koriyama, Yuki Saito, and Hiroshi Saruwatari. 2021. Emotion-controllable speech synthesis using emotion soft labels and fine-grained prosody factors. In <i>Asia-Pacific Signal and Information Processing Association Annual Summit and Conference, APSIPA ASC 2021</i> , pages 794–799. IEEE.	754
		755
		756
		757
		758
		759
		760
	Ziyang Ma, Zhisheng Zheng, Jiaxin Ye, Jinchao Li, Zhifu Gao, Shiliang Zhang, and Xie Chen. 2024. <a href="#">emotion2vec: Self-supervised pre-training for speech emotion representation</a> . In <i>Findings of the Association for Computational Linguistics, ACL 2024</i> , pages 15747–15760. Association for Computational Linguistics.	761
		762
		763
		764
		765
		766
		767
	Gabriel Mittag, Babak Naderi, Assmaa Chehadi, and Sebastian Möller. 2021. <a href="#">NISQA: A deep cnn-self-attention model for multidimensional speech quality prediction with crowdsourced datasets</a> . In <i>22nd Annual Conference of the International Speech Communication Association, Interspeech 2021</i> , pages 2127–2131. ISCA.	768
		769
		770
		771
		772
		773
		774
	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In <i>International Conference on Machine Learning, ICML 2023</i> , volume 202, pages 28492–28518. PMLR.	775
		776
		777
		778
		779
		780
	Chandan K. A. Reddy, Vishak Gopal, and Ross Cutler. 2022. <a href="#">Dnsmos P.835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors</a> . In <i>IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022</i> , pages 886–890. IEEE.	781
		782
		783
		784
		785
		786
	Neha Sahipjohn, Ashishkumar Gudmalwar, Nirmesh Shah, Pankaj Wasnik, and Rajiv Ratn Shah. 2024. <a href="#">Dubwise: Video-guided speech duration control in multimodal llm-based text-to-speech for dubbing</a> . In <i>25th Annual Conference of the International Speech Communication Association, Interspeech 2024</i> . ISCA.	787
		788
		789
		790
		791
		792
		793
	Antti Suni, Sébastien Le Maguer, Sofoklis Kakouros, Tuukka Törö, and Juraj Šimko. 2025. <a href="#">Style and Prosody control for Zero-shot Speech Synthesis</a> . In <i>13th edition of the Speech Synthesis Workshop</i> , pages 28–34.	794
		795
		796
		797
		798
	Xu Tan, Jiawei Chen, Haohe Liu, Jian Cong, Chen Zhang, Yanqing Liu, Xi Wang, Yichong Leng, Yuanhao Yi, Lei He, Sheng Zhao, Tao Qin, Frank K. Soong, and Tie-Yan Liu. 2024. <a href="#">Naturalspeech: End-to-end text-to-speech synthesis with human-level quality</a> . <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> , 46(6):4234–4245.	799
		800
		801
		802
		803
		804
		805

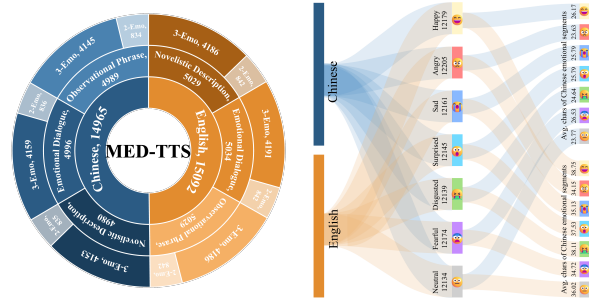


distributional constraints. For Step 1 text generation, each sample is verified for valid JSON formatting and required fields, including the text content, text category, and emotion sequence. Specifically, text length is constrained to 15-25 words for English and 15-25 characters for Chinese texts, and each emotion sequence is restricted to 2-3 segments drawn from a predefined set of 7 emotion categories. To reduce redundancy, we remove exact duplicates and filter near-duplicate texts using similarity-based criteria computed from normalized sequence-matching scores between token sequences, applying an overall similarity threshold of 0.85 and an opening similarity threshold of 0.5 over the first few tokens, with similarity comparisons primarily performed among samples sharing the same emotion sequence. For Step 2 multi-segment prompt annotation, automated checks enforce strict alignment between the text and its associated emotion annotations. Specifically, the number and order of segments are required to exactly match the predefined emotion sequence. Each segment must include a valid emotion label, a non-empty natural language emotion description with constrained length (5-25 words for English or 5-30 characters for Chinese), and an estimated speaking duration falling within a predefined range of 0.3-8.0 seconds.

Beyond automated filtering, we conduct manual verification through stratified random sampling of the validated outputs. In total, 1,000 samples are randomly selected for human review, comprising 500 English and 500 Chinese samples. The sampling process is stratified by language, text category, and the number of emotion segments to ensure broad coverage across diverse data conditions. Human reviewers then examine segmentation boundaries, emotion-text alignment, vocal-affect descriptions, and the plausibility of estimated speaking durations to identify subtle issues that may escape rule-based automatic checks. Insights obtained from this process are used to iteratively refine prompting strategies and validation thresholds. The manual verification checklist applied in Step 1 and Step 2 is provided in List 3.

### A.3 Dataset Statistics and Distribution

We summarize the statistics and distribution of the MED-TTS dataset across languages, text categories, and emotion segments. As illustrated in Fig. 5a, the dataset is well balanced across languages, comprising 14,965 Chinese and 15,092



(a) Distribution of Content Segments (b) Distribution of Emotion Segments

Figure 5: Statistics of the MED-TTS dataset. (a) Distribution of Chinese and English content segments. (b) Distribution of emotion segments with average character counts.

English samples. Within each language, samples are further evenly distributed across three text categories (vivid descriptions, emotional dialogues, and observational phrases), each contributing approximately 5,000 utterances. A finer-grained breakdown reveals that, within every text category, utterances containing three emotion segments consistently outnumber those with two emotion segments (e.g., roughly 4,100 vs. 800 per category), reflecting a deliberate emphasis on richer intra-utterance emotional transitions. Fig. 5b presents the segment-level emotion statistics. Across both Chinese and English, the 7 emotion types are uniformly represented, with each emotion accounting for approximately 1,200 segments. The average segment length remains stable within each language but differs across languages, with Chinese emotional segments typically spanning about 24-26 characters, while English segments are longer on average, ranging from roughly 34-39 words depending on emotion. Overall, MED-TTS achieves structured balance across languages, text categories, and emotion types, while maintaining sufficient emphasis on multi-emotion utterances to support modeling of continuous intra-utterance transitions.

As illustrated in Tab. 5, we further provide representative examples from the MED-TTS dataset across the three text categories for both Chinese and English. For each language-category pair, the table reports the sample count and total duration, along with illustrative text examples, corresponding emotion sequences, and natural language emotion descriptions. These examples demonstrate that each category consistently includes high-quality samples with different numbers of emotion segments,

Language	Category	Count	Duration(h)	Text Example	Emotion Sequence Example	Emotion Description Example
Chinese	Emotional Dialogue	4,996	9.89	失去你的日子里，心中满是空虚。⇒ 但与你重逢的那一刻，我的笑容重新绽放。	Sad ⇒ Happy	语速缓慢，语调低沉，带有失落和空虚感。⇒ 语速轻快，语调上扬，充满喜悦和温暖。
	Observational Phrase	4,989	10.00	茶杯中水波平静，⇒ 内心却如火山爆发。	Neutral ⇒ Angry	语调平稳，语速适中，声音自然放松。⇒ 语速加快，音调升高，声音紧张有力。
	Vivid Description	4,980	10.05	她无意中推开暗门，⇒ 霉味扑鼻，⇒ 脚步却不敢移动。	Surprised ⇒ Disgusted ⇒ Fearful	语调突然上扬，语速稍快，带有意外感。⇒ 声音压低，语速放缓，带有明显的嫌恶和停顿。⇒ 语调紧张、迟疑，语速缓慢，伴有轻微颤抖。
English	Emotional Dialogue	5,034	9.79	What in the world is that? ⇒ Ugh, it's revolting. ⇒ Well, I suppose it's just another part of life.	Surprised ⇒ Disgusted ⇒ Neutral	Voice rises sharply in pitch, with a quick, breathy delivery. ⇒ Tone is low, guttural, and drawn out with a visceral recoil. ⇒ Pace evens out to a calm, steady, and slightly resigned rhythm.
	Observational Phrase	5,029	9.56	A heated debate burned fiercely, each word adding fuel, ⇒ until playful banter extinguished the flames with lighthearted ease.	Angry ⇒ Happy	Voice is sharp, intense, and rapid, with a clipped, aggressive edge. ⇒ Tone becomes warm, relaxed, and lilting, with a cheerful, flowing cadence.
	Vivid Description	5,029	9.85	A high-pitched scream pierced his thoughts, ⇒ unraveling into a soft sigh, weighted with heartache and longing.	Fearful ⇒ Sad	Voice is sharp, tense, and sudden, with a quick, breathy delivery. ⇒ Tone is slow, breathy, and heavy, with a drawn-out, mournful quality.
<b>Total</b>		<b>30,057</b>	<b>59.14</b>			

Table 5: Dataset statistics and representative examples across languages and text categories.

highlighting the dataset’s coverage of diverse content types and intra-utterance emotional structures across different languages.

#### A.4 Fine-tuning Details

To enable automatic construction of segment-level TTS prompts, we fine-tune the Qwen3-8B large language model via supervised instruction tuning with parameter-efficient adaptation. We adopt LoRA to update only low-rank adapters while keeping the backbone frozen, thereby preserving general linguistic capabilities. Specifically, fine-tuning is carried out using the SFT-Trainer framework, with LoRA adapters applied to the attention and feed-forward projection layers, using a rank of 32, a scaling factor of 64, and a dropout rate of 0.1. Training is performed for 4 epochs with a per-device batch size of 2 and gradient accumulation over 4 steps, yielding an effective batch size of 8. We use a learning rate of  $1 \times 10^{-4}$  with a linear warmup of 100 steps and enable mixed-precision FP16 training for efficiency.

## B Segment-Aware Emotion Conditioning

This section provides a brief explanation of the symbols and steps used in Alg. 1.

**Inputs and outputs.**  $\mathbf{x} = \{x_1, x_2, \dots, x_T\}$  denotes the source text token sequence of length  $T$ .  $\mathbf{b} = \{b_1, b_2, \dots, b_M\}$  are segment boundaries on the text timeline, where  $M$  is the number of segments.  $\mathbf{C} = \{\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_M\}$  are the segment-wise condition embeddings (e.g., emotion/style prompts), and each condition may correspond to a short token span of length  $L_C$  in the decoder input. The algorithm autoregressively generates a semantic token stream  $\mathbf{s} = \{s_1, s_2, \dots, s_N\}$ .

**Segment index arrays.**  $\text{seg}_x[t]$  is the segment id assigned to the  $t$ -th text token  $x_t$  according to the boundaries.  $\text{seg}_s[i]$  stores the segment id used when generating the  $i$ -th semantic token  $s_i$ . The scalar  $m$  denotes the index of the *currently active* segment during decoding and determines which condition embedding is visible to the current generation step.

**Step 1: Direct construction of the 2D additive-bias mask  $\mathcal{M}_i$ .** At decoding step  $i$ , the decoder input is organized as a single concatenated token list: first all segment conditions, then the full text tokens, and finally the already-generated semantic tokens. Accordingly, the total query/key length is  $q = M \times L_C + T + i$ . We directly build an additive-bias mask  $\mathcal{M}_i \in \mathbb{R}^{q \times q}$ , where each entry

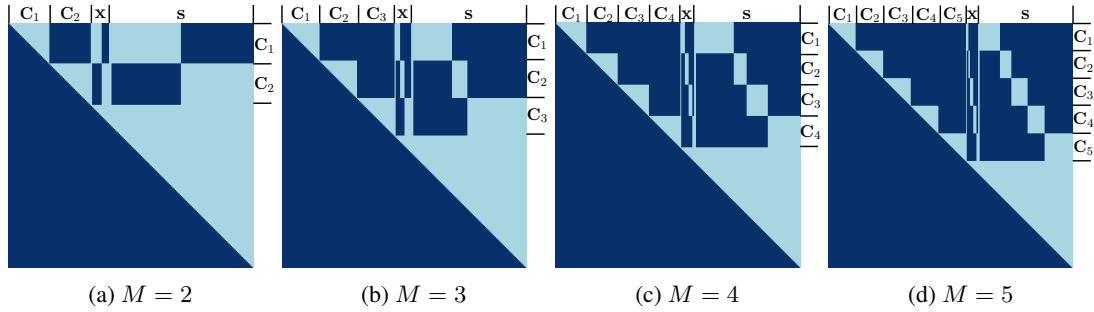


Figure 6: Visualization of the final attention mask under varying numbers of segment conditions ( $M$ ).

is either 0 (visible) or  $-\infty$  (masked). Compared to the previous block-matrix presentation, this version writes all constraints as in-place updates on  $\mathcal{M}_i$  with explicit offsets for the condition/text/semantic regions.

- **Standard causal visibility.** We first initialize  $\mathcal{M}_i$  and apply a standard causal mask so that each query token can only attend to itself and earlier tokens in the concatenated sequence. This ensures autoregressive consistency for semantic generation.
- **Text-to-condition visibility (segment-local control).** For each text token  $x_t$ , we overwrite its attention row to the *condition region* so that  $x_t$  can only see the condition tokens belonging to its own segment. Concretely, the entire condition region is masked out for that row, then only the span corresponding to  $C_{\text{seg}_x[t]}$  is unmasked. This prevents text tokens in one segment from reading condition prompts from other segments.
- **Semantic-to-condition visibility (segment-local control).** Similarly, for each previously generated semantic token  $s_r$ , we restrict its visibility to the condition region to be segment-local. The semantic token can only attend to the condition tokens of the segment recorded in  $\text{seg}_s[r]$ . This enforces that past semantic tokens do not leak information from conditions of unrelated segments.
- **Condition-to-condition isolation (no cross-condition leakage).** Condition tokens are not allowed to exchange information across different segments. We therefore mask each condition block’s attention to all other condition blocks, keeping only the within-block (diagonal) visibility. This makes each segment condition self-contained while still allowing the

overall model to read text/semantic context under the global causal structure.

**Step 2: One-step decoding and attention observation.** Given  $\mathcal{M}_i$ , the decoder performs one autoregressive step to produce the next semantic token  $s_i$ . During the same forward pass, it also returns the raw attention maps  $\mathbf{A}_i$  used as an on-line alignment observation. After generation, we append  $s_i$  to  $s$  and record  $\text{seg}_s[i] \leftarrow m$ .

**Step 3: Monotonic Stream Alignment (MSA).** MSA tracks where the semantic stream is aligned on the text in an online manner. It maintains a posterior belief over text positions and advances it with a monotonic prior transition to encourage forward progression. From the returned attentions  $\mathbf{A}_i$ , the algorithm selects a single layer/head whose attention pattern best matches the prior, optionally smooths it to reduce noise, and fuses it with the prior to obtain a stable posterior belief for the current step. In our implementation, the transition factor in  $\mathcal{P}$  is set to  $p = 0.1$ , and a Gaussian smoothing function  $\mathcal{G}_\sigma$  with  $\sigma = 1.2$  is applied.

**Step 4: Segment switching.** The active segment index  $m$  is updated by monitoring the expected aligned text position under the current posterior belief. Once this expected position passes the boundary of the current segment, we increment  $m$  to trigger an emotion/style switch for subsequent semantic tokens. Overall, this mechanism enforces segment-local control through restricted condition visibility, while preserving global coherence via standard causal decoding and online monotonic alignment.

**Mask visualization.** Fig. 6 illustrates the resulting mask pattern produced by the direct in-place construction in Step 1 when the utterance is partitioned into different numbers of segments

( $M=2, 3, 4, 5$ ). Each panel shows how the condition token blocks, text tokens, and semantic tokens are jointly constrained by (i) the global causal structure and (ii) the segment-local condition visibility. As  $M$  increases, the condition region is divided into more isolated blocks, and each text/semantic token is restricted to attend only to the condition block of its assigned (or currently active) segment. This visualization helps verify that the mask enforces local emotion/style control without allowing cross-condition leakage across segments.

## C Segment-Aware Duration Steering

In practice, segment-aware duration steering is implemented as two lightweight inference-time controllers that operate entirely on semantic token counts and alignment signals. For local duration steering, a proportional controller described in Section 3.3 performs online correction by comparing the normalized semantic generation progress within the active segment to the normalized text progress obtained from the MSA algorithm. The correction is applied with gain  $k_p = 25.0$ , which determines the sensitivity of duration adjustment to progress mismatch, and is triggered only when the absolute progress error exceeds  $\varepsilon = 0.01$ , thereby preventing unnecessary updates caused by minor alignment fluctuations. To ensure stability, updates are performed at a fixed low frequency of one update every five decoding steps, and the per-update adjustment is clamped to a maximum magnitude of  $\Delta_{\max} = 10$  semantic tokens to avoid abrupt changes in generation pace. In addition, the effective target is constrained by adaptive lower bounds tied to the current global decoding cursor, with conservative and emergency regimes activated when the generated length exceeds  $1.2\times$  and  $1.5\times$  the planned segment budget, respectively, serving as a safeguard against uncontrolled over-generation.

For global EOS steering, an EOS controller is added to the logits processor list, where EOS logits are fully suppressed for all non-final segments, while being dynamically adjusted in the final segment based on the ratio between generated semantic tokens and the target semantic budget. Specifically, EOS is strongly suppressed when the ratio is below 0.5, gradually transitions to a neutral region over the interval  $[0.8, 1.1]$ , and is increasingly encouraged as the ratio approaches 1.2, with the applied bias bounded between  $-5.0$  and  $+15.0$ . These fixed hyperparameters were selected empiri-

cally and remain constant across all experiments, enabling robust intra-utterance duration control without modifying or retraining the underlying TTS model.

## D Ablation Models Implementation

### D.1 Emotion and Duration Control

In the segment-aware emotion conditioning part of Tab. 3, we compare our method with the following two ablated variants:

- **w/o full-text access:** In this variant, each segment condition can only attend to the local text tokens within its own segment, rather than the full text.
- **w/o alignment:** In this variant, we remove any alignment module and generate semantic tokens by randomly switching phases through a fixed probability at each step.

For the segment-aware duration steering part of Tab. 3, we further evaluate two ablated variants to analyze the contributions of local and global steering mechanisms:

- **w/o local steering:** In this variant, the local duration steering module is disabled, and segment-level pacing relies solely on the baseline duration embedding, while the global EOS control mechanism is retained.
- **w/o global EOS:** In this variant, the global EOS logit modulation is disabled, while the local duration steering module remains active.

### D.2 Monotonic Stream Alignment Evaluation

In Fig. 4, we compare our MSA method with the following ablated variants:

- **Max Attention Head + Greedy Monotonic Alignment:** In this variant, we replace our MSA with a deterministic heuristic. We firstly compute a score for each raw attention maps across all layers and heads through  $F^{(l,h)} = \frac{1}{T} \sum_{t=1}^T \mathbf{A}_{i,t}^{(l,h)}$ , where  $T$  is the length of text tokens and  $t$  is the text position. The optimal attention map ( $l^*, h^*$ ) is selected as the observation by the maximum score. For the update step, we restrict a monotonic constraint and simplify the posterior  $\pi_i$  to a one-hot vector, representing a hard alignment state. Let  $k$  be the active index at the previous step, i.e.,

$\pi_{i-1}(k) = 1$ , and the update rule follows a greedy local comparison between the current position  $k$  and the next position  $k + 1$ . The new belief is determined as:

$$\pi_i(t) = \mathbb{1} \left[ t = \arg \max_{m \in \{k, k+1\}} \mathbf{A}_{i,m}^{(l,h)} \right], \quad (4)$$

where  $\mathbb{1}[\cdot]$  is the indicator function.

- **Top- $k$  Attention Heads + Greedy Monotonic Alignment:** In this variant, we extend the previous method by selecting the top- $k$  attention heads as observations. Specifically, we first compute the scores  $F^{(l,h)}$  for all attention maps and select the top- $k$  heads with the highest scores. The observation is then derived as a weighted average of these selected attention maps based on their scores. The greedy monotonic alignment update remains the same as above.
- **Max Attention Head + Monotonic Stream Alignment:** In this variant, we retain alignment updates using our MSA algorithm. We replace the observation component by selecting a single attention head with the maximum score as described above, and get rid of the smoothing operation.

## E Evaluation Protocol

### E.1 Baseline and Comparative models

**Baseline.** IndexTTS<sup>3</sup> (Zhou et al., 2025) is an autoregressive zero-shot TTS model that supports utterance-level control of emotion and speech duration while maintaining high speech naturalness. It disentangles speaker identity from emotional expression, enabling faithful reconstruction of target timbre and accurate reproduction of the specified emotional style. By incorporating GPT-based latent representations, the model further improves semantic consistency and stability under expressive conditions.

We also adopt several strong zero-shot TTS as our comparative methods:

- **MaskGCT<sup>4</sup>** (Wang et al., 2025c) is a non-autoregressive TTS model that a masked generative transformer to predict semantic and

acoustic tokens, functioned with duration control. By leveraging two-stage mask prediction mechanism, it achieves high fidelity and robust voice synthesis.

- **F5TTS<sup>5</sup>** (Chen et al., 2025) is a non-autoregressive TTS system based on Diffusion Transformer (DiT). It eliminates explicit alignment by padding text to speech length. Trained on 100k hours of data, it employs Sway Sampling to achieve efficient, high-quality zero-shot multilingual synthesis.
- **SparkTTS<sup>6</sup>** (Wang et al., 2025b) is a powerful TTS system built upon Qwen2.5, which directly reconstructs audio from LLM-predicted codes and eliminates the need for complex intermediate models like flow matching. It excels in high-fidelity zero-shot voice cloning for bilingual scenarios while maintaining high efficiency.
- **CosyVoice2<sup>7</sup>** (Du et al., 2024b) is an autoregressive TTS model that combines a language model for semantic and prosodic modeling with flow matching for speaker identity reconstruction, utilizing a supervised speech tokenizer to achieve disentangled generation. Notably, it demonstrates superior performance in Chinese compared to English due to its training data distribution.

Our baseline and comparative models adopt a consistent segment-wise inference strategy. Each sentence is partitioned into multiple segments based on target emotions and speaking rates generated by our fine-tuned LLM. These segments are generated individually among these models and sequentially assembled to reconstruct the complete utterance for evaluation. All baseline and comparative models are implemented using their official open-source codebases and pretrained weights.

### E.2 Subjective Evaluation

We conduct a subjective Mean Opinion Score (MOS) evaluation focusing on four key dimensions: emotion consistency, speaking rate consistency, speaker similarity, and emotional transition smoothness. Participants were provided with explicit scoring criteria, and we report the mean scores along

<sup>3</sup><https://github.com/index-tts/index-tts>

<sup>4</sup><https://github.com/open-mmlab/Amphion/tree/main/models/tts/maskgct>

<sup>5</sup><https://github.com/SWivid/F5-TTS>

<sup>6</sup><https://github.com/SparkAudio/Spark-TTS>

<sup>7</sup><https://github.com/FunAudioLLM/CosyVoice?tab=readme-ov-file>

with 95% confidence intervals (CI) in Tab.1 and Tab.2. The evaluation involved 15 graduate students with relevant research backgrounds. Prior to the evaluation, participants were provided with detailed task protocols and informed of the specific usage of the data. Each participant evaluated 18 test samples (9 Chinese and 9 English) under different settings, with the entire session lasting approximately 40 minutes. Scores ranged from 1 to 5 with 1-point intervals. Each participant received compensation of 15 SGD for their participation. The user interface for MOS evaluation is illustrated in Fig. 8.

### E.3 Objective Evaluation

Our objective evaluation encompasses several metrics to assess various aspects of speech synthesis quality and controllability. Character accuracy is measured using an automatic speech recognition (ASR) model through comparison with ground-truth transcriptions. For English audio evaluation, we employ a Whisper Large V3 (Radford et al., 2023) ASR model to calculate Word Error Rate (WER)<sup>8</sup>, while for Chinese audio, we utilize a Paraformer (Gao et al., 2022) ASR model to calculate Character Error Rate (CER) for Chinese to quantify transcription accuracy<sup>9</sup>.

To evaluate the smoothness of transitions in both emotion and speaking rate, we adopt the DNSMOS Pro<sup>10</sup> (Cumlin et al., 2024), referred as DNSM. It is calculated by averaging the predicted MOS values obtained from a sliding window (2-second duration, 1-second stride) applied across the full utterance. Speaker similarity is assessed using fine-tuned WavLM-Large (Chen et al., 2022) for speaker verification<sup>11</sup> to extract speaker embeddings from synthesized and reference audios, followed by computing the cosine similarity, denoted as SSIM. We report the average scores of the two metrics: the similarity between the synthesized and reference audios, and the intra-utterance consistency measured across all segment pairs obtained via ASR-based segmentation within the generated speech.

For speech naturalness evaluation, we utilize NISQA<sup>12</sup> (Mittag et al., 2021) and OVRL from

<sup>8</sup><https://huggingface.co/openai/whisper-large-v3>

<sup>9</sup><https://huggingface.co/funasd/paraformer-zh>

<sup>10</sup><https://github.com/fcumlin/DNSMOSPro>

<sup>11</sup>[https://github.com/microsoft/UniSpeech/tree/main/downstreams/speaker\\_verification](https://github.com/microsoft/UniSpeech/tree/main/downstreams/speaker_verification)

<sup>12</sup><https://github.com/gabrielmittag/NISQA>

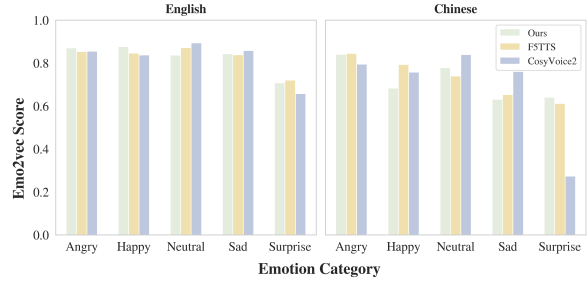


Figure 7: Comparison of Emo2Vec similarity scores across languages for five emotion categories. Our method is compared with F5TTS(Chen et al., 2025) and CosyVoice2(Du et al., 2024b).

DNSMOS<sup>13</sup> (Reddy et al., 2022) for overall quality of a synthesized sequence. Both of them are evaluated through the entire utterance without segmentation. The emotional expression accuracy is measured through extracting segment-level emotional embeddings from ASR-segmented audio clips using a pre-trained speech emotion recognition model emotion2vec-large<sup>14</sup> (Ma et al., 2024). We calculate the cosine similarity between synthesized and reference audios for in speech prompt settings, and utilize classification accuracy over 5 discrete emotional labels for text prompt settings.

### E.4 Experimental Result Supplements

**Emotion-Specific Control Evaluation.** We provide detailed experimental results for the emotional similarity across five discrete categories (Angry, Happy, Neutral, Sad, Surprise). Our results are shown in Fig. 7, highlighting distinct performance patterns across different input modalities. Our method demonstrates robust emotional fidelity, achieving emotional similarity scores that closely approach those of comparative methods. It’s worth emphasizing that while other methods could only generate single emotion clips where the global style is constant, our method generates continuous and multi-segment sequences with transitioning emotions. Despite the difficulty of modeling such dynamic emotional control, our model still maintains high emotional similarity, and even superior performance in certain categories such as Angry and Sad, proving its effectiveness in generating complex and varying prosody.

#### Category-Specific Emotion Control Evaluation.

<sup>13</sup><https://github.com/microsoft/DNS-Challenge/tree/master/DNSMOS>

<sup>14</sup>[https://huggingface.co/emotion2vec/emotion2vec\\_plus\\_large](https://huggingface.co/emotion2vec/emotion2vec_plus_large)

Prompt	Category	Method	WER↓	DNSM↑	SSIM↑	NISQA↑	OVRL↑	Emo2vec↑
Speech	Emotional Dialogue	MaskGCT	1.978	3.757	0.331	4.375	3.266	<b>0.862</b>
		CosyVoice	<b>0.459</b>	3.570	0.394	4.403	3.321	0.822
		Ours	5.462	<b>3.839</b>	<b>0.442</b>	<b>4.670</b>	<b>3.372</b>	0.811
	Observational Phrase	MaskGCT	4.468	3.849	0.352	4.518	3.277	<b>0.854</b>
		CosyVoice	1.021	3.567	0.410	4.515	3.333	0.844
		Ours	<b>0.834</b>	<b>3.903</b>	<b>0.491</b>	<b>4.710</b>	<b>3.418</b>	0.848
	Vivid Description	MaskGCT	3.612	3.866	0.377	4.463	3.267	0.831
		CosyVoice	<b>0.834</b>	3.621	0.443	4.535	3.313	0.826
		Ours	1.085	<b>3.940</b>	<b>0.508</b>	<b>4.636</b>	<b>3.391</b>	<b>0.836</b>
Text	Emotional Dialogue	CosyVoice	<b>0.956</b>	3.450	<b>0.473</b>	4.336	3.294	0.369
		IndexTTS2	4.059	3.506	0.399	4.201	3.194	<b>0.596</b>
		Ours	5.996	<b>3.634</b>	0.387	<b>4.515</b>	<b>3.324</b>	0.468
	Observational Phrase	CosyVoice	2.582	3.464	0.435	4.385	3.294	0.270
		IndexTTS2	<b>0.747</b>	3.531	0.436	4.316	3.220	<b>0.416</b>
		Ours	1.104	<b>3.657</b>	<b>0.477</b>	<b>4.606</b>	<b>3.327</b>	0.381
	Vivid Description	CosyVoice	1.631	3.483	0.437	4.373	3.284	0.289
		IndexTTS2	1.442	3.511	0.450	4.353	3.235	<b>0.489</b>
		Ours	<b>0.756</b>	<b>3.762</b>	<b>0.486</b>	<b>4.644</b>	<b>3.372</b>	0.393

Table 6: Objective evaluation results on English Speech and Text inputs across different text categories. ↓ indicates that lower values are better, while ↑ indicates that higher values are better. Best results are **bolded**.

Method	DNSM↑	SSIM↑	NISQA↑	OVRL↑	Emo2vec↑
<b>Ours</b>	<b>3.925</b>	<b>0.485</b>	<b>4.706</b>	<b>3.395</b>	<b>0.837</b>
Max Head + Greedy	3.901	0.443	4.683	3.372	0.803
Top- <i>k</i> + Greedy	3.878	0.442	4.702	3.383	0.815
Max Head + MSA	3.907	0.462	4.697	3.393	0.828

Table 7: Objective Comparison of Different Alignment Strategies. ↑ indicates that higher values are better. Best results are **bolded**.

Tab. 6 extends our evaluation to three distinct synthesis scenarios in our dataset: Emotional Dialogue, Observational Phrase, and Vivid Description. Our method demonstrates a consistent advantage in audio quality, achieving the highest NISQA and OVRL scores in almost all settings. It also excels in naturalness of emotional transitions, as reflected by the DNSM metric, where our method consistently outperforms all methods across different text categories and input modalities. This further confirms that our segment-aware generation effectively maintains naturalistic acoustic synthesis even when handling complex emotional transitions.

Nevertheless, we observe that our method still faces challenges in speaker similarity in certain scenarios. In the text prompt setting, our method lags behind CosyVoice2 in Emotional Dialogue, where the generated sentences often contain extensive emotionally charged and oral conversational elements, such as modal particles and emphatic punctuation. In this scenario, our method and

baseline prioritize the expressive prosody, which may lead to deviations from the target speaker’s timbre. While CosyVoice2 keeps a more stable and consistent prosody during generation, and preserves speaker identity, it fails to convey the intended emotional expressions. This highlights the inherent trade-off between emotional expressiveness and speaker fidelity in zero-shot TTS, especially when generating highly dynamic prosody from text alone.

**MSA Ablation Studies.** We further validate the effectiveness of our proposed Monotonic Stream Alignment (MSA) in Tab. 7. We compare our MSA-based alignment strategy against several variants illustrated in Appendix D.2. These protocols are the same as those in Fig. 4. The results show that replacing our MSA with greedy monotonic alignment leads to noticeable performance drops across all metrics, indicating that the MSA mechanism is crucial not just for text-audio synchronization, but also for stabilizing the emotional contents. Through maintaining a robust posterior belief of the current position, MAS prevents the model from drifting off the complex segment boundaries, thereby ensuring the naturalness and coherence of emotional transitions. Notably, even without the full MSA mechanism, these ablation variants still maintain relatively high performance levels. This suggests that synthesizing speech containing multiple com-

1429 plex emotions in a single continuous streaming pro-  
1430 cess, rather than generating each segment indepen-  
1431 dently, inherently preserves semantic and acoustic  
1432 coherence, which benefits the overall quality of the  
1433 generated speech.

Role:  
You are an expert creative screenwriter and emotional expression specialist.  
Your task is to generate high-quality text utterances for text-to-speech synthesis evaluation.

Task:  
Given an ordered emotion sequence, generate a single-sentence text utterance that reflects the emotional journey described by the sequence.  
The text should naturally transition through these emotions in order.

Emotion Sequence:  
1. \${Emotion\_1}\$  
2. \${Emotion\_2}\$  
3. \${Emotion\_3}\$

Requirements:

- Text Utterance:
  - Length: 15-25 words (corresponding to 5-10 seconds of speech).
  - The text MUST contain all emotions in the given sequence, each clearly identifiable.
  - Emotional transitions MUST be conveyed through changes in language tone, imagery, internal reactions, or perspective.
  - CRITICAL: Do NOT use explicit temporal markers such as "then", "now", "afterward", "at first", "later", "next", "suddenly", or "finally".
  - The sentence must be semantically coherent and flow naturally as a single utterance.
  - Avoid clichéd or overused expressions, especially as opening phrases.
  - The opening MUST be unique and creative; avoid common narrative patterns.
- Text Category Constraint:

```

${
- vivid_descriptive: Vivid descriptive sentences (novel prose style). Example: "Wind whispered through the parched cornstalks, its voice fraying like worn silk." |
- emotional_dialogue: Emotionally charged dialogue excerpts (natural spoken lines). Example: "I've asked you three times! Why is the door still locked?" |
- observational_phrase: Observational phrases (subtle situational commentary). Example: "Rain taps the window like it's bruising the glass-rhythmic, insistent, all night."
}$

```
- Output Format:  
Provide your response in the following JSON structure ONLY:

```

{
  "text": "<generated single-sentence utterance>",
  "text_category": "${text_category: vivid_descriptive | emotional_dialogue | observational_phrase}$"
}

```

Examples:

Example 1  
\${Example:  
Vivid Descriptive  
Input Emotion Sequence:  
1. Happy  
2. Surprised  
3. Sad  
Output:  

```

{
  "text": "Warm light drifts around me, a sudden sharp gust jolts the calm, and a muted heaviness settles quietly over my thoughts.",
  "text_category": "vivid_descriptive"
}
}$
...

```

Now generate a text utterance for the given emotion sequence.

Listing 1: Example prompt for generating content text with emotion shifts using GPT-4o.

Role:  
You are an expert linguistic annotator specialized in emotional prosody for TTS datasets.  
Your task is to segment the given sentence into emotion-aligned segments while preserving the exact original wording.

Task:  
Segment the following text into contiguous spans that correspond to the emotions in the sequence.  
Each segment must represent a natural linguistic unit and reflect its assigned emotion through tone, sensory cues, or attitude-NOT through explicit time markers.

Input Text:  
\${Original text generated in Step-1}\$

Emotion Sequence:  
1. \${Emotion\_1}\$  
2. \${Emotion\_2}\$  
3. \${Emotion\_3}\$

Requirements:

1. Segmentation Rules:
  - Produce EXACTLY the same number of segments as emotions in the sequence.
  - CRITICAL: Segments MUST correspond to the emotion sequence IN ORDER.  
The first segment maps to the first emotion, the second to the second emotion, etc.
  - Each segment MUST be a continuous span from the original text.  
Do NOT rewrite, reorder, omit, or add any words.
  - All punctuation marks from the original text MUST be preserved in their exact positions.
  - The concatenation of all segments MUST reconstruct the original text exactly.
  - Segment boundaries should align with natural linguistic or prosodic boundaries (e.g., phrase or clause boundaries). Do NOT split inside tight phrases.
2. Emotion Description (for TTS prosody reference):
  - Provide a short vocal-affect description (5-15 words) focusing on auditory qualities.
  - The description should focus on auditory characteristics (e.g., pitch, intensity, pacing), not on events or semantics.
  - The description MUST align with the assigned emotion.
3. Speaking Time Estimation:
  - Estimate speaking duration in seconds using the guideline:  
0.18-0.30 seconds per word as a baseline.
  - The estimated duration should also reflect the emotional tone of the segment, as different emotions naturally influence speaking pace (e.g., excited or tense delivery tends to be quicker, while somber or reflective delivery tends to slow down).
  - The final time MUST be a realistic approximation of how the segment would be delivered aloud.
  - IMPORTANT: The sum of all segment durations MUST fall within 5-13 seconds.
  - Output time values as decimal strings (e.g., "2.4").

Output Format (JSON ONLY):

```
{
  "original_text": "${original input text}$",
  "segments": [
    {
      "lines_seg": "<text segment>",
      "emotion": "<emotion label from the sequence>",
      "emotion_description": "<vocal-affect description>",
      "time": "<estimated speaking time in seconds>"
    },
    ...
  ]
}
```

Example:  
...

Now generate the segmentation for the given input text and emotion sequence.

Listing 2: Example prompt for emotion-aligned segmentation and duration annotation using DeepSeek-Chat.

Manual Review Checklist (total 1,000 samples: 500 EN / 500 ZH)

-----  
[Step 1] Content Text Generation

- Text validity:  
the text is complete, fluent, and natural, without obvious truncation, repetition, or unfinished clauses (typically a single well-formed sentence).
- Length appropriateness:  
text length falls within the intended range (EN: 15-25 words; ZH: 15-30 characters), and does not appear unnaturally compressed or padded to meet length requirements.
- Semantic coherence:  
the text conveys a single coherent idea or situation, rather than a loose collection of phrases or unrelated clauses.
- Category consistency:  
the assigned text category matches the content style (vivid descriptive / emotional dialogue / observational phrase), with category cues clearly identifiable within the text.
- Emotion sequence correctness:  
the emotion sequence contains 2-3 valid emotions drawn from the predefined set, and all emotions are meaningfully reflected somewhere in the text.
- Emotion progression naturalness:  
emotional transitions implied by the text occur in a plausible order, without abrupt or logically unsupported emotion jumps.
- Language quality:  
the text does not contain obvious grammatical errors, unnatural phrasing, or machine-like constructions that would hinder natural speech rendering.

-----  
[Step 2] Multi-segment Prompt Annotation

- Segmentation boundaries:  
segment splits occur at natural linguistic or prosodic boundaries, such as phrase or clause breaks, and avoid splitting fixed expressions or tight collocations.
- Emotion-text alignment:  
the semantic content of each segment clearly supports its assigned emotion, and the intended emotion is perceivable without relying on the description.
- Vocal-affect specificity:  
emotion description includes concrete auditory cues (e.g., energy level, pitch tendency, speaking rate, intensity) rather than abstract emotion names.
- Description naturalness:  
emotion description reads as a natural speaking instruction and typically spans one short phrase or sentence, rather than a list of keywords.
- Duration plausibility:  
estimated speaking durations are reasonable given segment length and linguistic complexity, and fall within the expected range of 0.3-8.0 seconds per segment.
- Duration consistency:  
duration differences across segments reflect intuitive pacing differences, such as faster delivery for excited emotions and slower delivery for calm or reflective ones.
- Coverage consistency:  
concatenated segment texts fully reconstruct the original text, with no missing, duplicated, or reordered content.

-----  
Manual Review Protocol

- Reviewers:  
all sampled items are independently inspected by at least two reviewers, covering both English and Chinese samples.
- Disagreement handling:  
cases with inconsistent judgments are discussed and resolved through consensus review, and recurring issues are recorded for prompt or rule refinement.

Listing 3: Manual verification checklist used in our human review process for Step 1 and Step 2 outputs.

---

**Algorithm 1** Segment-Aware Emotion Conditioning with Monotonic Stream Alignment (MSA)

---

**Input** : Text tokens  $\mathbf{x} = \{x_1, x_2, \dots, x_T\}$ , segment boundaries  $\mathbf{b} = \{b_1, b_2, \dots, b_M\}$ , condition embeddings  $\mathbf{C} = \{\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_M\}$

**Output** : Generated semantic tokens  $\mathbf{s} = \{s_1, s_2, \dots, s_N\}$

// Compute segment id for text tokens  $\text{seg}_{\mathbf{x}}[1..T]$ :

```
1 for  $t \leftarrow 1$  to  $T$  do
2   |  $\text{seg}_{\mathbf{x}}[t] \leftarrow 1 + \sum_{r=1}^{M-1} \mathbb{I}[t > b_r]$ 
3 end
4 Initialize  $\text{seg}_{\mathbf{s}} \leftarrow []$  // store segment id for each generated semantic token
5 Initialize semantic index  $i \leftarrow 0$ , segment index  $m \leftarrow 1$ ,  $\mathbf{s} \leftarrow \emptyset$  Initialize posterior alignment belief  $\boldsymbol{\pi}_0 \in \mathbb{R}^T$  (one-hot at  $t = 1$ )
6 while not EndOfSentence do
7   |  $i \leftarrow i + 1$ 
8   | // 1) Build 2D additive-bias offset mask directly on  $\mathcal{M}_i$ 
9   |  $q \leftarrow M \times L_C + T + i$  // #Q tokens:  $[\mathbf{C}_{1:M}, x_{1:T}, s_{1:i}]$ 
10  |  $\mathcal{M}_i \leftarrow (-\infty) \cdot \mathbf{1}_{(q) \times (q)}$ 
11  | // 1.1) Standard causal mask
12  | for  $u \leftarrow 1$  to  $q$  do
13  |   | for  $v \leftarrow 1$  to  $u$  do
14  |     |  $\mathcal{M}_i[u, v] \leftarrow 0$ 
15  |   end
16  | end
17  | // 1.2)  $x \rightarrow C$ : Text tokens to condition embeddings (segment-local condition visibility)
18  |  $\text{off} \leftarrow M \times L_C$ 
19  | for  $t \leftarrow 1$  to  $T$  do
20  |   |  $\mathcal{M}_i[\text{off} + t, 0 : M \times L_C] \leftarrow -\infty$ 
21  |   |  $\mathcal{M}_i[\text{off} + t, L_C \times (\text{seg}_{\mathbf{x}}[t] - 1) : L_C \times \text{seg}_{\mathbf{x}}[t]] \leftarrow 0$ 
22  | end
23  | // 1.3)  $S \rightarrow C$ : Semantic tokens to condition embeddings (segment-local condition visibility)
24  |  $\text{off} \leftarrow M \times L_C + T$ 
25  | for  $r \leftarrow 1$  to  $i - 1$  do
26  |   |  $\mathcal{M}_i[\text{off} + r, 0 : M \times L_C] \leftarrow -\infty$ 
27  |   |  $\mathcal{M}_i[\text{off} + r, L_C \times (\text{seg}_{\mathbf{s}}[r] - 1) : L_C \times \text{seg}_{\mathbf{s}}[r]] \leftarrow 0$ 
28  | end
29  | // 1.4)  $C \rightarrow C$ : Condition embeddings to condition embeddings (no cross-condition leakage)
30  | for  $u \leftarrow 1$  to  $M$  do
31  |   |  $\mathcal{M}_i[(u-1) \times L_C : u \times L_C, 0 : M \times L_C] \leftarrow -\infty$ 
32  |   |  $\mathcal{M}_i[(u-1) \times L_C : u \times L_C, (u-1) \times L_C : u \times L_C] \leftarrow 0$ 
33  | end
34  | // 2) Decode one step with mask and get raw attentions as observation
35  |  $(s_i, \mathbf{A}_i) \leftarrow f_{\theta}^{\text{decode-step}}(\mathbf{x}, \mathbf{s}_{<i}, \{\mathbf{C}_j\}_{j=1}^M, \mathcal{M}_i, \text{return\_attn} = \text{True})$ 
36  |  $\mathbf{s} \leftarrow \{\mathbf{s}, s_i\}$ 
37  |  $\text{seg}_{\mathbf{s}}[i] \leftarrow m$ 
38  | // 3) MSA: Predict-Select-Update
39  |  $\hat{\boldsymbol{\pi}}_i \leftarrow \boldsymbol{\pi}_{i-1} \cdot \mathcal{P}$  // Predict (Prior)
40  |  $(l^*, h^*) \leftarrow \arg \max_{l, h} \hat{\boldsymbol{\pi}}_i^{\top} \log(\mathbf{A}_i^{(l, h)})$ 
41  |  $\mathbf{a}^* \leftarrow \mathcal{G}_{\sigma}(\mathbf{A}_i^{(l^*, h^*)})$  // Select (Observation)
42  |  $\boldsymbol{\pi}_i \leftarrow (\hat{\boldsymbol{\pi}}_i \odot \mathbf{a}^*) / Z$  // Update (Posterior)
43  | // 4) Segment switching via expected aligned position
44  | if  $m < M$  and  $\sum_{t=1}^T t \cdot \boldsymbol{\pi}_i[t] > b_m$  then
45  |   |  $m \leftarrow m + 1$  // Trigger emotion switch
46  | end
47 end
48 return  $\mathbf{s}$ 
```

---

**Section 1.1: Speech-Reference Assessment**  
(8 questions)

In this section, the synthesized audio is generated using reference speech clips as a guide for timbre and emotion. Please evaluate the generated audio by comparing it against the reference, strictly adhering to the scoring criteria below:

**1. SpeakerMOS (Speaker Similarity)**  
Assessment of how closely the voice resembles the target speaker's identity. **Timbre consistency between segments** should be considered in this metric.  
5 (Excellent): Nearly indistinguishable from the target speaker.  
4 (Good): Timbre is very close to the target.  
3 (Fair): Largely similar, with minor pronunciation differences.  
2 (Poor): Vague resemblance but clearly different.  
1 (Bad): Completely different timbre or identity.  
**Note:** All segments are spoken by the same target speaker. When evaluating speaker similarity, you should assess whether the synthesized audio consistently reflects the identity of the target speaker as a whole. **Do not judge similarity by matching each segment independently. The speaker timbre in the generated audio should remain unified and consistent across all segments, rather than varying between segments.**

**2. NaturalnessMOS (Smoothness & Transitions)**  
Assessment of the speech's natural flow and the absence of artifacts. **Both inter-segment and intra-segment naturalness** should be considered in this metric.  
5 (Excellent): Seamless transitions with perfectly natural flow.  
4 (Good): Smooth and fluent transitions.  
3 (Fair): Mostly natural, with slight discontinuities.  
2 (Poor): Unnatural and abrupt transitions.  
1 (Bad): Highly unnatural with noticeable breaks.

**3. EmoSOS (Emotional Similarity)**  
Assessment of how well the synthesized emotion aligns with the reference.  
5 (Excellent): Completely aligned, strong, and accurate.  
4 (Good): Highly consistent with only subtle differences.  
3 (Fair): Generally aligned but with minor deviations.  
2 (Poor): Directionally similar but clearly different.  
1 (Bad): Completely mismatched and inconsistent.

**Note:** Please familiarize yourself with the criteria above before proceeding. This page will not be shown again.

Next

**(Question 1/18) Please evaluate the audio samples according to the listed metrics.**

1. Generated Audio:  
▶ 0:00 / 0:10

[Emotion Sequence: Angry -> Happy -> Neutral]  
[Text: You missed the deadline again! I But, seeing your efforts brightens my day. I Let's discuss the next steps calmly.]  
[Note: The "I" symbol separates segments with different emotions.]

2. Reference Audio:  
Angry:  
▶ 0:00 / 0:02  
Happy:  
▶ 0:00 / 0:01  
Neutral:  
▶ 0:00 / 0:03

Emotional Similarity (Between the Generated & Reference Audio)	1 (Bad)	2 (Poor)	3 (Fair)	4 (Good)	5 (Excellent)
Speaker Similarity (Between the Generated & Reference Audio, include timbre consistency across segments)					
Naturalness (Smoothness & Transitions of the Generated Audio, including inter-segment and intra-segment naturalness)					

Next

### (a) Speech-Prompted Emotion Control Evaluation

**Section 1.2: Text-Reference Assessment**  
(9 questions)

In this section, the synthesized audio is generated using reference text as an emotion guide, and reference audio as a timbre guide. Please evaluate the generated audio by comparing it against the reference, strictly adhering to the scoring criteria below:

**1. SpeakerMOS (Speaker Similarity)**  
Assessment of how closely the voice resembles the target speaker's identity. **Timbre consistency between segments** should be considered in this metric.  
5 (Excellent): Nearly indistinguishable from the target speaker.  
4 (Good): Timbre is very close to the target.  
3 (Fair): Largely similar, with minor pronunciation differences.  
2 (Poor): Vague resemblance but clearly different.  
1 (Bad): Completely different timbre or identity.  
**Note:** All segments are spoken by the same target speaker. When evaluating speaker similarity, you should assess whether the synthesized audio consistently reflects the identity of the target speaker as a whole. **Do not judge similarity by matching each segment independently. The speaker timbre in the generated audio should remain unified and consistent across all segments, rather than varying between segments.**

**2. NaturalnessMOS (Smoothness & Transitions)**  
Assessment of the speech's natural flow and the absence of artifacts. **Both inter-segment and intra-segment naturalness** should be considered in this metric.  
5 (Excellent): Seamless transitions with perfectly natural flow.  
4 (Good): Smooth and fluent transitions.  
3 (Fair): Mostly natural, with slight discontinuities.  
2 (Poor): Unnatural and abrupt transitions.  
1 (Bad): Highly unnatural with noticeable breaks.

**3. EmoSOS (Emotional Similarity)**  
Assessment of how well the synthesized emotion aligns with the reference.  
5 (Excellent): Completely aligned, strong, and accurate.  
4 (Good): Highly consistent with only subtle differences.  
3 (Fair): Generally aligned but with minor deviations.  
2 (Poor): Directionally similar but clearly different.  
1 (Bad): Completely mismatched and inconsistent.

**Note:** Please familiarize yourself with the criteria above before proceeding. This page will not be shown again.

Next

**(Question 9/9) Please evaluate the audio samples according to the listed metrics.**

1. Generated Audio:  
▶ 0:00 / 0:08

[Emotion Sequence: Disgusted -> Surprised -> Happy]  
[Text: Ugh, what is that smell? Wait, is it chocolate cake baking? Oh wow, this reminds me of Grandma's kitchen!]  
2. Reference Text:  
(1)  
"segment": "Ugh, what is that smell? ",  
"emotion": "Disgusted",  
"emotion\_description": "The voice is tense and sharp, with a clear note of displeasure."  
(2)  
"segment": "Wait, is it chocolate cake baking? ",  
"emotion": "Surprised",  
"emotion\_description": "The tone is quick and lifted, capturing sudden curiosity."  
(3)  
"segment": "Oh wow, this reminds me of Grandma's kitchen!",  
"emotion": "Happy",  
"emotion\_description": "The voice is warm and nostalgic, with a gentle, joyful cadence."

3. Reference Audio:  
▶ 0:00 / 0:02

Emotional Similarity (Between the Generated Audio & Reference Text)	1 (Bad)	2 (Poor)	3 (Fair)	4 (Good)	5 (Excellent)
Speaker Similarity (Between the Generated & Reference Audio, include timbre consistency across segments)					
Naturalness (Smoothness & Transitions of the Generated Audio, including inter-segment and intra-segment naturalness)					

Next

### (b) Text-Prompted Emotion Control Evaluation

**Section 2: Duration Assessment**  
(12 questions)

In this section, the synthesized audio is generated using reference audio as a timbre guide, with the generation duration controlled for a specific segment. Please evaluate the generated audio by comparing it against the reference, strictly adhering to the scoring criteria below:

**1. SpeakerMOS (Speaker Similarity)**  
Assessment of how closely the voice resembles the target speaker's identity. **Timbre consistency between segments** should be considered in this metric.  
5 (Excellent): Nearly indistinguishable from the target speaker.  
4 (Good): Timbre is very close to the target.  
3 (Fair): Largely similar, with minor pronunciation differences.  
2 (Poor): Vague resemblance but clearly different.  
1 (Bad): Completely different timbre or identity.

**2. NaturalnessMOS (Smoothness & Transitions)**  
Assessment of the speech's natural flow and the absence of artifacts. **Both inter-segment and intra-segment naturalness** should be considered in this metric.  
5 (Excellent): Seamless transitions with perfectly natural flow.  
4 (Good): Smooth and fluent transitions.  
3 (Fair): Mostly natural, with slight discontinuities.  
2 (Poor): Unnatural and abrupt transitions.  
1 (Bad): Highly unnatural with noticeable breaks.

**3. SPMOS (Speaking Rate Matching)**  
Assessment of how accurate the speaking rate is compared to the reference or description.  
5 (Excellent): Perfect match with the reference or description.  
4 (Good): Close match with a natural rhythm.  
3 (Fair): Close to reference but with slight differences.  
2 (Poor): Significantly deviates from the reference.  
1 (Bad): Severely deviates (e.g., fast vs. slow).

**Note:** Please familiarize yourself with the criteria above before proceeding. This page will not be shown again. Speaking Rate Matching

Next

**(Question 1/12) Please evaluate the audio samples according to the listed metrics.**

1. Original Audio:  
▶ 0:00 / 0:06

[Text: Is it so hard to understand my point? **Your silence cuts deeper than any harsh words could.**]

2. Processed Audio: Apply a 1.125x duration scaling to the bold segment (slower)  
▶ 0:00 / 0:07

3. Reference Audio:  
▶ 0:00 / 0:02

Speaking Rate Matching (Between the Processed Audio and the Duration Target)	1 (Bad)	2 (Poor)	3 (Fair)	4 (Good)	5 (Excellent)
Speaker Similarity (Between the Generated & Reference Audio, include timbre consistency across segments)					
Naturalness (Smoothness & Transitions of the Generated Audio, including inter-segment and intra-segment naturalness)					

Next

### (c) Duration Control Evaluation

Figure 8: User interface for MOS evaluation across different evaluation tasks.