

# Prompt-Guided Zero-Shot Voice Cloning for Music Generation

Anonymous ACL submission

## Abstract

Existing text-to-music (TTM) models rarely address voice cloning for singing, and most mainstream approaches do not support cross-domain voice cloning based on reference speech. However, speech-to-music voice transfer is highly valuable in practical applications, as speech data is easier to collect and provides more stable speaker representations. In this paper, we propose **S2M-Inject**, a cross-domain music generation framework that enables voice cloning based on reference speech. By injecting speaker representations extracted from speech into the music generation process, S2M-Inject produces music that preserves the target voice characteristics of the input speech. Experimental results demonstrate that S2M-Inject can effectively perform cross-domain voice cloning while maintaining reasonable music generation quality, and supports both Chinese and English music generation.

## 1 Introduction

In recent years, TTM models have achieved remarkable progress in controllable music generation (Copet et al., 2023; Liu et al., 2023; Gong et al., 2025; Chen et al., 2025a). However, voice cloning for singing remains a relatively underexplored problem, and most existing TTM approaches do not support cross-domain voice cloning based on reference speech (Copet et al., 2023; Liu et al., 2023; Gong et al., 2025; Chen et al., 2025a). Current models typically focus on controlling textual or musical attributes, such as style, emotion, or instrumentation, while exhibiting limited capability in modeling singer timbre, making precise and controllable timbre transfer difficult to achieve.

In contrast, speech-based voice cloning has been extensively studied (Wang et al., 2023; Du et al., 2024; Chen et al., 2025b), largely due to the large availability of speech data and the stability of speaker representations. From a practical perspec-

tive, speech-to-music voice transfer is highly desirable: speech data is easier to collect and provides reliable speaker identity information, making it a natural and scalable conditioning source for personalized music generation. Despite its potential, how to effectively incorporate speech-based voice characteristics into music generation remains insufficiently explored, as existing music generation models either rely on textual conditioning or music-intrinsic audio references, without explicitly modeling speech-derived timbre information (Copet et al., 2023; Liu et al., 2023; Agostinelli et al., 2023; Gong et al., 2025).

The lack of cross-domain voice cloning capability highlights a fundamental gap between speech generation and music generation. Although speech and music differ significantly in structure and expressive patterns, they share common acoustic factors related to speaker or singer timbre. Bridging this gap is crucial for enabling controllable singing generation driven by reference speech.

To address this challenge, we propose S2M-Inject, a cross-domain music generation framework that supports voice cloning based on reference speech. By injecting speaker representations extracted from speech into the music generation process, S2M-Inject enables music generation that preserves voice characteristics consistent with the reference speech. Experimental results demonstrate that S2M-Inject can effectively achieve cross-domain voice cloning while maintaining reasonable music generation quality, and supports both Chinese and English music generation.

Our main contributions are summarized as follows: (1) We study cross-domain voice cloning from speech to music, extending voice cloning beyond conventional speech synthesis; (2) We propose a fixed timbre injection strategy that integrates speech-based speaker representations into instruction-guided music generation; (3) Experiments show improved voice consistency with com-

petitive music generation quality.

## 2 Related Works

### 2.1 Zero-shot Voice Cloning

With the advancement of neural speech codecs and large-scale generative models, zero-shot voice cloning has achieved substantial progress in generation quality, stability, and generalization (Shen et al., 2023).

A prominent line of research is based on neural codec representations. VALL-E (Wang et al., 2023) represents speech using discrete codec tokens and adopts a cascaded autoregressive and non-autoregressive generation scheme, inheriting strong contextual modeling capability from language models while enabling high-quality zero-shot speech synthesis. NaturalSpeech 2 (Shen et al., 2023) instead employs continuous latent representations and introduces in-context learning into a diffusion-based acoustic modeling framework. Building upon this line of work, NaturalSpeech 3 (Shen et al., 2023) proposes a factorized diffusion-based zero-shot TTS system, improving naturalness through decoupled codec representations, although its codec prediction objectives remain constrained by identical textual content conditions. Another category of methods introduces intermediate semantic representations to bridge the modeling gap between text and acoustic features. SpearTTS (Kharitonov et al., 2023) and Make-a-Voice (Huang et al., 2023) leverage semantic tokens as an intermediate layer, enabling more robust zero-shot voice transfer across speakers and speaking styles. VoiceBox (Le et al., 2023) adopts a non-autoregressive flow-matching framework to perform speech infilling conditioned on audio context and text, demonstrating strong generation stability under partial conditioning. The Mega-TTS series (Jiang et al., 2023) employs mel-spectrograms as generation targets and explicitly disentangles timbre and prosody, modeling prosody in an autoregressive manner to improve timbre preservation. Similarly, P-Flow (Kim et al., 2023) applies flow-based models to zero-shot speech synthesis and achieves robust performance. In addition, some studies focus on improving inference efficiency and parallel generation capability. SoundStorm (Borsos et al., 2023b) and MobileSpeech (Ji et al., 2024) adopt masked non-autoregressive iterative generation strategies, achieving a favorable trade-off between synthesis speed and audio quality, which is

important for large-scale deployment scenarios.

Despite the significant progress of zero-shot voice cloning in the speech domain, existing approaches are largely restricted to speech-to-speech scenarios. Most methods assume that both the reference and target signals belong to the speech domain, and cross-domain voice cloning from speech to music or singing remains largely unexplored. To address this limitation, we propose S2M-Inject, which investigates speech-based cross-domain voice cloning by injecting speaker representations extracted from speech into the music generation process, enabling consistent voice modeling from speech to music.

### 2.2 Music Generation

Early automatic music generation methods primarily relied on rule-based systems and symbolic representations, such as MIDI, combined with sequence models including RNNs, LSTMs, and Transformers (Hadjeres et al., 2017; Huang et al., 2018). While these approaches are relatively mature in modeling musical structure, they heavily depend on large-scale, well-annotated MIDI datasets and exhibit clear limitations in audio quality and realism (Huang et al., 2018), making them difficult to scale to diverse and complex musical scenarios.

With the advancement of deep learning and large-scale generative models, research has gradually shifted toward direct audio-level music generation. AudioLM (Borsos et al., 2023a) introduced the language modeling paradigm to audio generation, and MusicLM (Agostinelli et al., 2023) further incorporated pretrained models such as MuLan and w2v-BERT to achieve high-quality TTM synthesis. MusicGen (Copet et al., 2023) adopts a single-stage Transformer architecture with an efficient token interleaving strategy, striking a favorable balance among generation quality, efficiency, and textual controllability.

In parallel, diffusion-based approaches have demonstrated strong performance in music generation. Methods such as AudioLDM (Liu et al., 2023), Noise2Music (Huang et al., 2023), and MusicLDM (Liu et al., 2023) employ latent diffusion frameworks to enable text-guided high-fidelity music generation, and progressively extend to multi-track music and arrangement modeling.

Subsequent works, including JEN-1 (Yao et al., 2025b) and Mustango (Melechovsky et al., 2024), introduce two-stage diffusion architectures or explicit musical priors (e.g., rhythm, harmony, and

tonality) to further enhance structural modeling. In recent studies, MusicGen, AudioLDM, and MusicLDM represent mainstream text-guided music generation approaches, while DiffRhythm+ (Chen et al., 2025a) and ACE-Step (Gong et al., 2025) further target long-form song generation and achieve notable progress in musical coherence and structure modeling. However, the conditioning mechanisms of these methods are still largely restricted to textual descriptions or music-intrinsic attributes, and they generally do not support cross-domain voice cloning based on reference speech (Copet et al., 2023; Liu et al., 2023; Chen et al., 2024; Gong et al., 2025).

Overall, music generation research has evolved from symbolic modeling to end-to-end audio generation frameworks that integrate language models and diffusion models. Nevertheless, existing approaches predominantly focus on single-modality conditions, such as text or audio, and remain limited in cross-modal or cross-domain conditioning. To address this gap, we propose S2M-Inject, which explores speech-conditioned cross-domain music generation by injecting voice representations extracted from speech into the music generation process, enabling consistent voice modeling from speech to music and providing a new perspective for multi-source conditioned music generation.

## 3 Method

### 3.1 Overview

As illustrated in Fig. 1, S2M-Inject is built upon an MM-DiT architecture, drawing inspiration from the design of MM-Audio (Cheng et al., 2025) and targeting the modeling of Voice Cloning in cross-domain scenarios. Specifically, the proposed framework enables singer timbre modeling through cross-domain injection of *speech-derived timbre representations* for music generation. The model consists of two core components: joint diffusion transformer layers and single diffusion transformer layers (Esser et al., 2024).

During training, the instruction encoder, mel encoder, and mel decoder are employed as pre-trained modules and remain frozen, without participating in the optimization of the diffusion backbone. In the first joint diffusion transformer layer, the diffusion timestep is introduced, and the instruction embedding, text or lyric embedding, and speech-derived timbre embedding (Desplanques et al., 2020) are temporally concatenated as con-

ditioning signals. Meanwhile, a noise-perturbed Mel-VAE latent is used as the audio modality input (Liu et al., 2023).

Unlike existing non-autoregressive text-to-audio architectures (Chen et al., 2025b; Zhu et al., 2025), S2M-Inject does not require explicit upsampling of textual representations to align with audio features, thereby simplifying the overall modeling pipeline.

For the TTM task, S2M-Inject leverages natural language instructions together with cross-domain injected speech timbre conditions to achieve multi-attribute controllable music generation, including singer timbre injection, melodic and emotional expression, music genre, and instrumentation control, while supporting both Chinese and English singing generation.

### 3.2 Model Architecture

S2M-Inject is built upon a MM-DiT backbone and integrates specially designed encoders to process heterogeneous input streams from different modalities.

**Conditional Input Representation** To enable conditional modeling for singing generation, we propose a unified conditional input representation that supports target audio generation under multiple heterogeneous constraints. The model inputs include natural language audio descriptions, lyric content, speaker timbre conditions, latent audio representations, and stochastic Gaussian noise. Lyric text is first converted into phoneme sequences using a Grapheme-to-Phoneme (G2P) model (Qiang et al., 2022), providing a unified content representation. The target audio waveform is encoded into a latent representation via a Mel-VAE, which is then linearly interpolated with Gaussian noise under a diffusion timestep sampled from  $t \sim \mathcal{U}(0, 1)$  (Lipman et al., 2022), forming the noisy acoustic state required for diffusion modeling. This unified conditional representation provides a stable conditioning foundation for the subsequent generation process.

**Semantic and Textual Encoders** For semantic control, we feed structured audio descriptions generated by Gemini (Team et al., 2023) into Qwen3-Omni (Yang et al., 2025) and extract its dense hidden states as semantic embeddings  $E_{inst}$ , which characterize the global semantic attributes and expressive properties of the target audio. For content modeling, following the design principles of ZipVoice (Zhu et al., 2025) and M3-TTS (Wang et al., 2025), we adopt Zipformer (Yao et al., 2023)

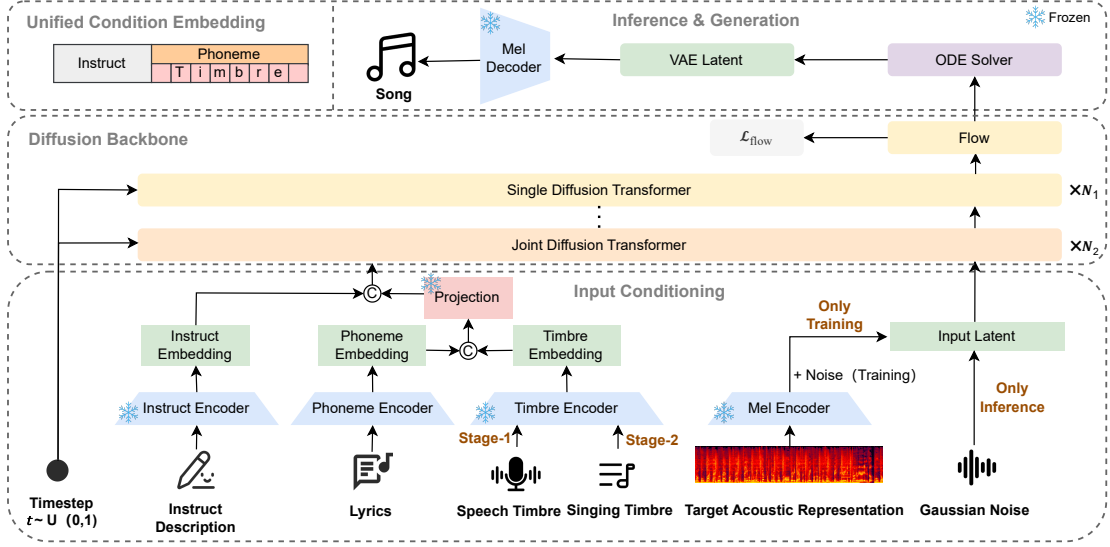


Figure 1: S2M-Inject is built upon a multimodal diffusion transformer (MM-DiT) for music generation, supporting multi-attribute control and voice cloning conditioned on natural language instructions and reference speech timbre. Audio is represented using continuous latent variables extracted from a pretrained Mel-VAE. During inference, the target music latent is obtained by solving an ordinary differential equation (ODE) and decoded into waveform audio.

as the textual encoder. The input text is tokenized using Byte Pair Encoding (BPE) (Sennrich et al., 2016) and then processed by the Zipformer to efficiently extract temporally aligned feature representations, enabling accurate modeling of singing content.

### Acoustic Conditioning and Timbre Injection

The proposed framework adopts a densely connected time-delay neural network (D-TDNN) (Yu and Li, 2020) to extract highly discriminative speaker timbre representations from short speech segments. To mitigate interference from expressive singing variations, we adopt a two-stage timbre conditioning strategy that uses singing-derived timbre information in Stage 1 and speech-derived timbre information from constructed speech–singing pairs in Stage 2, with data construction detailed in Section 3.3.

To enable cross-domain voice cloning, speaker timbre information is derived from the original speech, emotional attributes are neutralized using a predefined neutral embedding, and neutral speech for the target text is synthesized with IndexTTS2 as the timbre reference.

$$\mathbf{E}_{\text{spk}} \in \mathbb{R}^{B \times D}. \quad (1)$$

In the joint diffusion layers, the instruction embeddings and phoneme embeddings are used as

textual conditioning inputs. Specifically, the instruction embeddings are denoted as

$$\mathbf{E}_{\text{inst}} \in \mathbb{R}^{B \times L_1 \times D}, \quad (2)$$

and the phoneme embeddings are represented as

$$\mathbf{E}_{\text{phn}} \in \mathbb{R}^{B \times L_2 \times D}. \quad (3)$$

Here,  $B$  denotes the batch size,  $D$  is the shared embedding dimension, and  $L_1$  and  $L_2$  represent the sequence lengths of instruction tokens and phoneme tokens, respectively.

The speaker timbre embedding  $\mathbf{E}_{\text{spk}}$ , derived from the stage-specific timbre condition  $e_{\text{spk}}$  (i.e.,  $e_{\text{spk}}^{\text{sing}}$  in Stage 1 and  $e_{\text{spk}}^{\text{cross}}$  in Stage 2), is first projected and expanded along the temporal dimension to align with the phoneme sequence, and then concatenated with the phoneme embeddings along the feature dimension:

$$\tilde{\mathbf{E}}_{\text{phn}} = [\mathbf{E}_{\text{phn}} \parallel \mathbf{E}_{\text{spk}}] \in \mathbb{R}^{B \times L_2 \times 2D}. \quad (4)$$

A linear projection is subsequently applied to restore the unified feature dimension:

$$\hat{\mathbf{E}}_{\text{phn}} \in \mathbb{R}^{B \times L_2 \times D}. \quad (5)$$

Finally, the instruction embeddings and the fused phoneme–timbre representations are concatenated

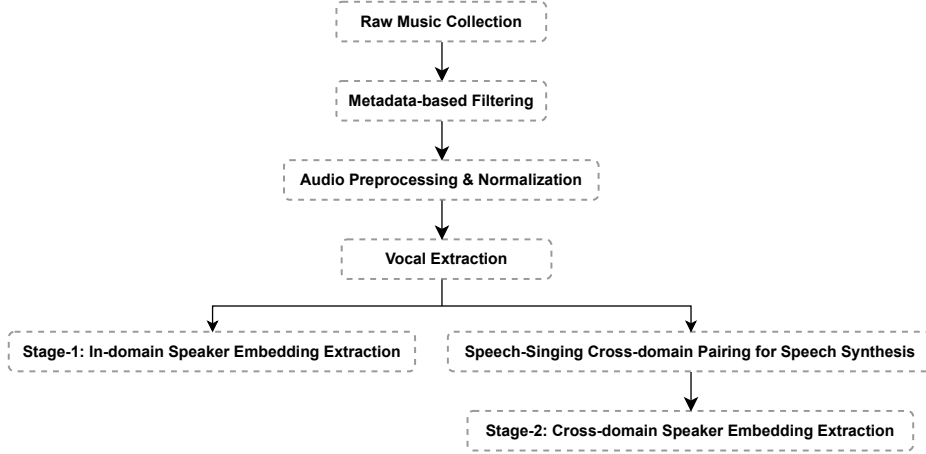


Figure 2: Overall Two-stage Training Data Construction Pipeline.

along the token (temporal) dimension to form the unified textual conditioning input:

$$\mathbf{C}_{\text{text}} = \left[ \mathbf{E}_{\text{inst}} \parallel_T \hat{\mathbf{E}}_{\text{phn}} \right] \in \mathbb{R}^{B \times (L_1 + L_2) \times D}. \quad (6)$$

Audio generation is performed in the latent space of a Mel-VAE, which compresses 44.1 kHz waveforms into low-dimensional continuous representations, effectively reducing computational cost while preserving essential acoustic characteristics.

**MM-DiT Backbone** Our generative model adopts a MM-DiT architecture, whose overall design is based on Stable Diffusion 3 (Esser et al., 2024) and trained with Conditional Flow Matching (CFM) (Lipman et al., 2022). The backbone consists of  $N_2$  Joint Diffusion Transformer layers and  $N_1$  Single Diffusion Transformer layers.

In the joint diffusion layers, a unified textual conditioning input and a latent acoustic state are used as the two interacting modalities. Following the timbre injection formulation in Sec. 3.2, the instruction embedding  $\mathbf{E}_{\text{inst}}$  and the fused phoneme-timbre representation  $\hat{\mathbf{E}}_{\text{phn}}$  are concatenated along the token (temporal) dimension to form  $\mathbf{C}_{\text{text}}$ .

For the audio modality, the latent acoustic state  $\mathbf{x}_t$  is constructed via linear interpolation between Gaussian noise and the Mel-VAE latent under a diffusion timestep sampled from  $t \sim \mathcal{U}(0, 1)$ . The two modalities interact through joint attention, where queries, keys, and values from different modalities are concatenated and processed by scaled dot-product attention to enable cross-modal information fusion. The output preserves the input dimensionality and is subsequently split back into the corresponding modalities.

In the single diffusion transformer layers, the model operates solely on the latent acoustic state to further refine audio generation quality, where the joint attention mechanism degenerates into standard self-attention over the audio modality. During training, We define the CFM loss as

$$\mathcal{L}_{\text{flow}}(\theta) = \mathbb{E} \|v_\theta(t, \mathbf{C}_{\text{text}}, \mathbf{x}_t) - u(t, \mathbf{x}_t)\|^2. \quad (7)$$

where  $v_\theta$  denotes the learned conditional velocity field,  $u$  is the target conditional vector field, and  $t$  is the randomly sampled timestep. During inference, the target Mel-VAE latent is obtained by solving the corresponding ODE and finally decoded into waveform audio.

### 3.3 Two-Stage Data Construction and Training Strategy

We construct a two-stage dataset and adopt a corresponding training strategy to enable stable speech-to-singing timbre injection. The data construction pipeline is illustrated in Fig. 2. We collect approximately 5,000 hours of raw music data and apply metadata-based filtering (e.g., duration, sampling rate, compression), resulting in about 4,000 hours of curated data (around 1M clips). All samples are processed with unified preprocessing and normalization, including ASR-based lyric transcription, natural language audio captions, and vocal extraction using Demucs (Défossez et al., 2019).

Timbre supervision is constructed in two stages. In Stage 1, the original singing vocals from 500k clips provide in-domain timbre information that is acoustically consistent with the target singing audio. In Stage 2, original speech recordings provide speaker timbre information, while emotional

---

**Algorithm 1: Two-Stage Cross-Style Training Strategy**

---

**Input:**Stage-1 condition  $C_{\text{text}}^{(1)} = [C_{\text{text}}, e_{\text{spk}}^{\text{sing}}]$ ;Stage-2 condition  $C_{\text{text}}^{(2)} = [C_{\text{text}}, e_{\text{spk}}^{\text{cross}}]$ ;Stage transition step  $s_{\text{stage1}}$ **for each training step  $s$  do****if  $s \leq s_{\text{stage1}}$  then**

Stage 1: In-domain timbre injection;

Use condition  $C_{\text{text}}^{(1)}$ ;

Optimize model with

 $\mathcal{L}_{\text{flow}}(x_t, t, C_{\text{text}}^{(1)})$ ;**else**

Stage 2: Cross-style timbre injection;

Use condition  $C_{\text{text}}^{(2)}$ ;

Optimize model with

 $\mathcal{L}_{\text{flow}}(x_t, t, C_{\text{text}}^{(2)})$ ;**end if****end for**

---

factors are removed by using a predefined neutral emotion embedding. Neutral speech corresponding to the target texts is synthesized using IndexTTS2 (Zhou et al., 2025), and the synthesized speech signals serve as cross-domain timbre conditions aligned with the target singing voice.

Based on the constructed data, we adopt a two-stage training strategy summarized in Algorithm 1. The model is first trained with in-domain singing timbre conditions and then fine-tuned with speech-derived timbre conditions, using the same CFM objective in both stages.

## 4 Experiments

### 4.1 Model Details and Datasets

Following the overall architectural design described in Sec. 3.2, the proposed S2M-Inject model contains 1.34 billion parameters, with a feed-forward network dimension of 1024. The model consists of 14 Joint Diffusion Transformer layers and 6 Single Diffusion Transformer layers, and employs Rotary Positional Embedding (RoPE) (Su et al., 2024) to enhance sequence modeling capability.

For textual and instruction modeling, we adopt a Zipformer-based phoneme encoder (Yao et al., 2023) with a feed-forward dimension of 512. The instruction encoder is instantiated using Qwen2.5-7B (Yang et al., 2025), which maps natural-

language instructions into high-dimensional semantic representations. The mel encoder takes raw waveforms sampled at 44.1 kHz as input and produces latent embeddings at approximately 43 Hz, corresponding to a temporal downsampling factor of about  $1024\times$  relative to the input sampling rate. This design significantly reduces computational cost while preserving sufficient representational capacity.

Model training is conducted on 32 NVIDIA Tesla A800 GPUs with 80 GB memory each, using a per-GPU batch size of 16. We employ the Adam optimizer (Kingma, 2014) with an initial learning rate of  $1e-4$ , following common practice in diffusion-based audio generation models.

For music generation experiments, we collect and curate approximately 5,000 hours of music data from the Internet, and automatically generate aligned natural-language instruction descriptions and lyric annotations through an internal data processing pipeline. The instruction descriptions cover a wide range of musical and vocal attributes, including music genre, instrumentation, singer gender and age, rhythmic characteristics, and overall atmosphere. Audio clips range from 2 to 20 seconds in duration, and the dataset maintains balanced distributions in both language (Chinese/English) and singer gender at a 1:1 ratio. The dataset primarily consists of single-singer segments to reduce the interference of multi-voice overlap on timbre modeling, and covers diverse musical styles, accompaniment patterns, tempo variations, and accompaniment complexity. All audio samples are uniformly resampled to 44.1 kHz for both training and evaluation.

### 4.2 Comparisons and Evaluation Metrics

Most existing TTM models do not support cross-domain Voice Cloning based on reference speech, such as speech-to-singing timbre transfer. Even in music generation scenarios, models explicitly designed for singer Voice Cloning remain scarce, with the majority of approaches providing only coarse-grained singer attribute control rather than strict Voice Cloning.

To comprehensively evaluate the performance of S2M-Inject in both Voice Cloning and music generation, we conduct comparative experiments under Voice Cloning and TTM settings against representative baseline methods. All baseline models are deployed using their local inference implementations and evaluated on the same held-out test set

Table 1: Cross-domain voice cloning results evaluated by speaker similarity, WER, and subjective evaluation.

Model	Task	SIM $\uparrow$	WER $\downarrow$	SIM-MOS $\uparrow$
Ground Truth	–	1.00	0.00	–
CosyVoice2 (Du et al., 2024)	Speech $\rightarrow$ Speech	<b>0.68</b>	2.01	<b>3.75<math>\pm</math>0.13</b>
Stage-1 S2M-Inject	Speech $\rightarrow$ Music	0.34	1.77	2.35 $\pm$ 0.22
<b>Stage-2 S2M-Inject</b>	Speech $\rightarrow$ Music	0.60	<b>1.56</b>	3.32 $\pm$ 0.34

Table 2: Music generation results evaluated by SongEval and subjective evaluation.

Model	Params	SongEval $\uparrow$					MOS $\uparrow$	
		Coh.	Mus.	Mem.	Clar.	Nat.	QMOS	MMOS
Ground Truth	–	3.60	3.52	3.56	3.40	3.34	–	–
DiffRhythm+ (Chen et al., 2025a)	1B	2.68	2.61	2.57	2.48	2.37	3.04 $\pm$ 0.46	2.79 $\pm$ 0.54
ACE-Step (Gong et al., 2025)	3B	2.89	2.87	2.83	2.77	2.71	<b>3.30 <math>\pm</math> 0.28</b>	2.88 $\pm$ 0.20
InstructAudio (Qiang et al., 2025)	1.3B	<b>3.08</b>	2.98	3.00	2.89	<b>2.82</b>	2.82 $\pm$ 0.26	<b>2.91 <math>\pm</math> 0.35</b>
Stage-1 S2M-Inject	1.3B	2.88	2.79	2.69	2.64	2.66	2.73 $\pm$ 0.36	2.78 $\pm$ 0.42
<b>Stage-2 S2M-Inject</b>	1.3B	3.06	<b>2.99</b>	<b>3.03</b>	<b>2.95</b>	2.77	2.92 $\pm$ 0.33	2.90 $\pm$ 0.37

consisting of 500 samples, ensuring fair and consistent comparison. Since DiffRhythm+ (Chen et al., 2025a) does not natively support generating music clips shorter than 90 seconds, we first generate full-length sequences and subsequently crop them to the target duration; this implementation difference may introduce minor bias in certain evaluation metrics.

For the Voice Cloning task, we compare S2M-Inject with CosyVoice2 (Du et al., 2024), and the results are summarized in Table 1. For music generation, we compare S2M-Inject with DiffRhythm+ (Chen et al., 2025a), ACE-Step (Gong et al., 2025), and InstructAudio (Qiang et al., 2025), with quantitative results summarized in Table 2.

We adopt a combination of objective and subjective evaluation metrics. For Voice Cloning, objective metrics include Speaker Similarity (SIM) and Word Error Rate (WER) (Anastassiou et al., 2024). For music generation, we employ the SongEval benchmark (Yao et al., 2025a), which evaluates Coherence, Musicality, Memorability, Clarity, and Naturalness. Subjective evaluations are conducted using QMOS, MMOS, and SIM-MOS.

### 4.3 Results and Analysis

**Voice Cloning Performance Analysis** Table 1 reports both objective and subjective evaluation results on the voice cloning task. We focus on

analyzing the performance variation of S2M-Inject across the two training stages (Stage 1 and Stage 2).

From the SIM metric, we observe that Stage 1 achieves a similarity score of 0.34, while Stage 2 significantly improves this score to 0.60 after introducing the speech–music cross-pair based training strategy. This improvement indicates that incorporating reference speech timbre constraints in the second stage effectively enhances the model’s ability to capture and transfer target speaker characteristics. In contrast, Stage 1, which relies solely on in-domain singing data, tends to depend more on music-domain statistics, resulting in limited performance for cross-domain voice alignment.

Compared with the speech-to-speech voice cloning model CosyVoice2, Stage 2 still exhibits a performance gap in terms of SIM. It is important to note that CosyVoice2 is specifically designed for in-domain speech voice cloning, whereas S2M-Inject targets the more challenging speech-to-music cross-domain voice cloning scenario. Given the fundamental differences in task difficulty and input conditions, this performance gap is within reasonable expectations.

In terms of intelligibility, Stage 2 achieves a lower WER of 1.56 compared to 1.77 for Stage 1, indicating that introducing cross-domain voice conditioning does not degrade linguistic content. Subjective evaluation results are consistent with these observations: Stage 2 attains a SIM-MOS score of

3.32, significantly higher than Stage 1, demonstrating improved voice similarity and overall perceptual quality.

**Music Generation Quality Analysis** Table 2 reports both objective and subjective evaluation results for music generation. We focus on whether S2M-Inject can maintain reasonable music generation quality while introducing cross-domain speech-to-singing timbre conditioning.

Comparing the two training stages, Stage 2 consistently outperforms Stage 1 across all SongEval metrics, with clear improvements in coherence, musicality, memorability, clarity, and naturalness. This indicates that the second-stage training successfully incorporates speech-derived timbre information without degrading musical structure or perceptual consistency, and in fact improves overall music generation quality under short-segment settings. These gains suggest that the introduced cross-domain timbre conditioning is well aligned with the underlying music generation objectives.

When compared with existing TTM models, S2M-Inject (Stage 2) achieves competitive SongEval performance. It outperforms DiffRhythm+ across all SongEval dimensions and is comparable to ACE-Step, while remaining slightly below InstructAudio on several metrics. We note that these comparisons are conducted on short audio clips ranging from 5 to 20 seconds, whereas models such as ACE-Step and DiffRhythm+ are primarily optimized for long-form music generation. Therefore, the reported results mainly reflect performance under short-segment generation scenarios.

In terms of subjective evaluation, S2M-Inject (Stage 2) shows consistent improvements over Stage 1 on both QMOS and MMOS, indicating that introducing cross-domain timbre conditioning does not negatively affect perceived audio quality. Although its QMOS remains lower than ACE-Step and InstructAudio, S2M-Inject achieves comparable MMOS, suggesting that the generated music remains musically acceptable while prioritizing cross-domain voice cloning capability. Overall, these results demonstrate that S2M-Inject preserves stable and competitive music generation quality while enabling speech-driven singing voice cloning.

## 5 Discussion

This work investigates speech-conditioned cross-domain voice cloning for music generation under

a unified instruction-guided framework. Using speech as the timbre source provides stable and scalable identity cues, which contributes to the consistent improvement from Stage 1 to Stage 2; however, the intrinsic acoustic and expressive mismatch between speech and singing still limits the achievable upper bound of timbre similarity compared to in-domain speech-to-speech cloning. Moreover, while instruction-based conditioning enables unified control across speech and music generation, it inevitably compresses acoustic information and introduces one-to-many mappings, leading to averaged realizations and a remaining gap in perceptual quality relative to reference-audio-based methods. These results indicate that cross-domain timbre transfer remains challenging, especially when balancing controllability and perceptual quality. To ensure stable joint modeling, music generation is restricted to short segments, which facilitates fair cross-domain evaluation aligned with speech durations but does not fully reflect long-form music generation capability. Overall, S2M-Inject demonstrates the feasibility of speech-driven cross-domain voice cloning for music generation and serves as an initial step toward more flexible conditioning mechanisms across audio domains.

## 6 Conclusions and Future Work

In this paper, we propose S2M-Inject, a cross-domain voice cloning framework based on the MM-DiT architecture, which enables timbre injection and cloning from speech to music generation. Experimental results demonstrate that the proposed method can effectively achieve cross-domain timbre transfer, validating both the feasibility and effectiveness of the framework. The main contributions of this work can be summarized as follows. (1) We explore the problem of cross-domain voice cloning from speech to music, extending voice cloning beyond conventional speech synthesis scenarios to music generation. (2) We propose a fixed timbre injection strategy that integrates speech-based speaker representations into instruction-guided music generation. (3) Extensive experiments show that the proposed approach improves timbre consistency while maintaining reasonable music generation quality. For future work, we plan to further investigate long-form music generation, enhance cross-domain voice cloning capability, and extend voice cloning applications to a broader range of audio domains.

## 638 Limitations

639 Despite the promising results in unified speech and  
640 music generation, this work has several limitations.

641 First, our timbre control relies on reference  
642 speech signals, and cross-domain transfer from  
643 speech to singing/music may introduce a domain  
644 mismatch, which can affect timbre faithfulness and  
645 perceptual consistency in some cases.

646 Second, to accommodate the joint modeling of  
647 speech and music, music generation is restricted  
648 to short clips of 5–20 seconds, which limits the  
649 model’s ability to capture long-term musical struc-  
650 ture.

651 Finally, unified speech and music generation  
652 still lacks mature and comprehensive evaluation  
653 benchmarks, and existing metrics may not fully  
654 reflect long-term musical coherence or creative di-  
655 versity. We also note that AI-based tools were  
656 used for proofreading and language polishing of  
657 the manuscript, while all technical content and con-  
658 clusions were determined by the authors.

## 659 Ethical Statement

660 This work investigates instruction-guided speech  
661 and music generation with timbre control, which  
662 may introduce potential ethical risks. We outline  
663 the following considerations:

- 664 • **Voice misuse and impersonation.** The pro-  
665 posed method enables timbre-controlled gen-  
666 eration, which could be misused for imper-  
667 sonation. This work is intended solely for  
668 research purposes and does not aim to repli-  
669 cate or impersonate any specific individual  
670 without authorization.
- 671 • **Data usage and privacy.** All training data are  
672 collected from publicly available or properly  
673 licensed sources and processed in accordance  
674 with applicable data usage and privacy poli-  
675 cies.
- 676 • **Responsible deployment.** We encourage re-  
677 sponsible use of the proposed method and em-  
678 phasize that appropriate safeguards, consent  
679 mechanisms, and usage restrictions should be  
680 considered in real-world applications to pre-  
681 vent malicious or deceptive use.

## References

- 682  
683 Andrea Agostinelli, Timo I Denk, Zalán Borsos,  
684 Jesse Engel, Mauro Verzetti, Antoine Caillon,  
685 Qingqing Huang, Aren Jansen, Adam Roberts,  
686 Marco Tagliasacchi, and 1 others. 2023. Musi-  
687 clm: Generating music from text. *arXiv preprint*  
688 *arXiv:2301.11325*.
- 689 Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe  
690 Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng,  
691 Chuang Ding, Lu Gao, and 1 others. 2024. Seed-tts:  
692 A family of high-quality versatile speech generation  
693 models. *arXiv preprint arXiv:2406.02430*.
- 694 Zalán Borsos, Raphaël Marinier, Damien Vincent, Eu-  
695 gene Kharitonov, Olivier Pietquin, Matt Sharifi,  
696 Dominik Roblek, Olivier Teboul, David Grangier,  
697 Marco Tagliasacchi, and 1 others. 2023a. Audioldm:  
698 a language modeling approach to audio generation.  
699 *IEEE/ACM transactions on audio, speech, and lan-  
700 guage processing*, 31:2523–2533.
- 701 Zalán Borsos, Matt Sharifi, Damien Vincent, Eugene  
702 Kharitonov, Neil Zeghidour, and Marco Tagliasacchi.  
703 2023b. Soundstorm: Efficient parallel audio genera-  
704 tion. *arXiv preprint arXiv:2305.09636*.
- 705 Huakang Chen, Yuepeng Jiang, Guobin Ma, Chunbo  
706 Hao, Shuai Wang, Jixun Yao, Ziqian Ning, Meng  
707 Meng, Jian Luan, and Lei Xie. 2025a. Diffrrhythm+:  
708 Controllable and flexible full-length song genera-  
709 tion with preference optimization. *arXiv preprint*  
710 *arXiv:2507.12890*.
- 711 Ke Chen, Yusong Wu, Haohe Liu, Marianna Nezhurina,  
712 Taylor Berg-Kirkpatrick, and Shlomo Dubnov.  
713 2024. Musicldm: Enhancing novelty in text-to-music  
714 generation using beat-synchronous mixup strategies.  
715 In *ICASSP 2024-2024 IEEE International Confer-  
716 ence on Acoustics, Speech and Signal Processing*  
717 (*ICASSP*), pages 1206–1210. IEEE.
- 718 Yushen Chen, Zhikang Niu, Ziyang Ma, Keqi Deng,  
719 Chunhui Wang, JianZhao JianZhao, Kai Yu, and Xie  
720 Chen. 2025b. F5-tts: A fairytaler that fakes fluent and  
721 faithful speech with flow matching. In *Proceedings*  
722 *of the 63rd Annual Meeting of the Association for*  
723 *Computational Linguistics (Volume 1: Long Papers)*,  
724 pages 6255–6271.
- 725 Ho Kei Cheng, Masato Ishii, Akio Hayakawa, Takashi  
726 Shibuya, Alexander Schwing, and Yuki Mitsufuji.  
727 2025. Mmaudio: Taming multimodal joint training  
728 for high-quality video-to-audio synthesis. In *Pro-  
729 ceedings of the Computer Vision and Pattern Recog-  
730 nition Conference*, pages 28901–28911.
- 731 Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David  
732 Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre  
733 Défossez. 2023. Simple and controllable music gen-  
734 eration. *Advances in Neural Information Processing*  
735 *Systems*, 36:47704–47720.

736	Alexandre Défossez, Nicolas Usunier, Léon Bottou, and Francis Bach. 2019. Demucs: Deep extractor for music sources with extra unlabeled data remixed. <i>arXiv preprint arXiv:1909.01174</i> .	Sungwon Kim, Kevin Shih, Joao Felipe Santos, Evelina Bakhturina, Mikyas Desta, Rafael Valle, Sungroh Yoon, Bryan Catanzaro, and 1 others. 2023. P-flow: A fast and data-efficient zero-shot tts through speech prompting. <i>Advances in Neural Information Processing Systems</i> , 36:74213–74228.	792
737			793
738			794
739			795
740	Brecht Desplanques, Jenthe Thienpondt, and Kris Demuynck. 2020. Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification. <i>arXiv preprint arXiv:2005.07143</i> .	Diederik P Kingma. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	796
741			797
742			
743			
744			
745	Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, and 1 others. 2024. Cosyvoice 2: Scalable streaming speech synthesis with large language models. <i>arXiv preprint arXiv:2412.10117</i> .	Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, Vimal Manohar, Yossi Adi, Jay Mahadeokar, and 1 others. 2023. Voicebox: Text-guided multilingual universal speech generation at scale. <i>Advances in neural information processing systems</i> , 36:14005–14034.	800
746			801
747			802
748			803
749			804
750	Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, and 1 others. 2024. Scaling rectified flow transformers for high-resolution image synthesis. In <i>Forty-first international conference on machine learning</i> .	Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. 2022. Flow matching for generative modeling. <i>arXiv preprint arXiv:2210.02747</i> .	806
751			807
752			808
753			809
754			
755			
756	Junmin Gong, Sean Zhao, Sen Wang, Shengyuan Xu, and Joe Guo. 2025. Ace-step: A step towards music generation foundation model. <i>arXiv preprint arXiv:2506.00045</i> .	Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. 2023. Audioldm: Text-to-audio generation with latent diffusion models. <i>arXiv preprint arXiv:2301.12503</i> .	810
757			811
758			812
759			813
760			814
761	Gaëtan Hadjeres, François Pachet, and Frank Nielsen. 2017. Deepbach: a steerable model for bach chorales generation. In <i>International conference on machine learning</i> , pages 1362–1371. PMLR.	Jan Melechovsky, Zixun Guo, Deepanway Ghosal, Navonil Majumder, Dorien Herremans, and Soujanya Poria. 2024. Mustango: Toward controllable text-to-music generation. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8293–8316.	815
762			816
763			817
764	Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, Monica Dinulescu, and Douglas Eck. 2018. Music transformer. <i>arXiv preprint arXiv:1809.04281</i> .		818
765			819
766			820
767			821
768			822
769	Qingqing Huang, Daniel S Park, Tao Wang, Timo I Denk, Andy Ly, Nanxin Chen, Zhengdong Zhang, Zhishuai Zhang, Jiahui Yu, Christian Frank, and 1 others. 2023. Noise2music: Text-conditioned music generation with diffusion models. <i>arXiv preprint arXiv:2302.03917</i> .	Chunyu Qiang, Peng Yang, Hao Che, Jinba Xiao, Xiaorui Wang, and Zhongyuan Wang. 2022. Back-translation-style data augmentation for mandarin chinese polyphone disambiguation. In <i>2022 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)</i> , pages 1915–1919. IEEE.	823
770			824
771			825
772			826
773			827
774			828
775	Shengpeng Ji, Ziyue Jiang, Hanting Wang, Jialong Zuo, and Zhou Zhao. 2024. Mobilespeech: A fast and high-fidelity framework for mobile zero-shot text-to-speech. <i>arXiv preprint arXiv:2402.09378</i> .	Chunyu Qiang, Kang Yin, Xiaopeng Wang, Yuzhe Liang, Jiahui Zhao, Ruibo Fu, Tianrui Wang, Cheng Gong, Chen Zhang, Longbiao Wang, and 1 others. 2025. Instructaudio: Unified speech and music generation with natural language instruction. <i>arXiv preprint arXiv:2511.18487</i> .	830
776			831
777			832
778			833
779	Ziyue Jiang, Yi Ren, Zhenhui Ye, Jinglin Liu, Chen Zhang, Qian Yang, Shengpeng Ji, Rongjie Huang, Chunfeng Wang, Xiang Yin, and 1 others. 2023. Mega-tts: Zero-shot text-to-speech at scale with intrinsic inductive bias. <i>arXiv preprint arXiv:2306.03509</i> .	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In <i>Proceedings of the 54th annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 1715–1725.	834
780			835
781			
782			
783			
784			
785	Eugene Kharitonov, Damien Vincent, Zalán Borsos, Raphaël Marinier, Sertan Girgin, Olivier Pietquin, Matt Sharifi, Marco Tagliasacchi, and Neil Zeghidour. 2023. Speak, read and prompt: High-fidelity text-to-speech with minimal supervision. <i>Transactions of the Association for Computational Linguistics</i> , 11:1703–1718.	Kai Shen, Zeqian Ju, Xu Tan, Yanqing Liu, Yichong Leng, Lei He, Tao Qin, Sheng Zhao, and Jiang Bian. 2023. Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers. <i>arXiv preprint arXiv:2304.09116</i> .	836
786			837
787			838
788			839
789			840
790			
791			

846 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan,  
847 Wen Bo, and Yunfeng Liu. 2024. Roformer: En-  
848 hanced transformer with rotary position embedding.  
849 *Neurocomputing*, 568:127063.

850 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-  
851 Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan  
852 Schalkwyk, Andrew M Dai, Anja Hauth, Katie Mil-  
853 lican, and 1 others. 2023. Gemini: a family of  
854 highly capable multimodal models. *arXiv preprint*  
855 *arXiv:2312.11805*.

856 Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang,  
857 Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu,  
858 Huaming Wang, Jinyu Li, and 1 others. 2023. Neural  
859 codec language models are zero-shot text to speech  
860 synthesizers. *arXiv preprint arXiv:2301.02111*.

861 Xiaopeng Wang, Chunyu Qiang, Ruibo Fu, Zhengqi  
862 Wen, Xuefei Liu, Yukun Liu, Yuzhe Liang, Kang  
863 Yin, Yuankun Xie, Heng Xie, and 1 others. 2025. M3-  
864 tts: Multi-modal dit alignment & mel-latent for zero-  
865 shot high-fidelity speech synthesis. *arXiv preprint*  
866 *arXiv:2512.04720*.

867 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
868 Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
869 Gao, Chengen Huang, Chenxu Lv, and 1 others.  
870 2025. Qwen3 technical report. *arXiv preprint*  
871 *arXiv:2505.09388*.

872 Jixun Yao, Guobin Ma, Huixin Xue, Huakang Chen,  
873 Chunbo Hao, Yuepeng Jiang, Haohe Liu, Ruibin  
874 Yuan, Jin Xu, Wei Xue, and 1 others. 2025a.  
875 Songeval: A benchmark dataset for song aesthetics  
876 evaluation. *arXiv preprint arXiv:2505.10793*.

877 Yao Yao, Peike Li, Boyu Chen, and Alex Wang. 2025b.  
878 Jen-1 composer: A unified framework for high-  
879 fidelity multi-track music generation. In *Proceedings*  
880 *of the AAAI Conference on Artificial Intelligence*,  
881 volume 39, pages 14459–14467.

882 Zengwei Yao, Liyong Guo, Xiaoyu Yang, Wei Kang,  
883 Fangjun Kuang, Yifan Yang, Zengrui Jin, Long Lin,  
884 and Daniel Povey. 2023. Zipformer: A faster and bet-  
885 ter encoder for automatic speech recognition. *arXiv*  
886 *preprint arXiv:2310.11230*.

887 Ya-Qi Yu and Wu-Jun Li. 2020. Densely connected  
888 time delay neural network for speaker verification.  
889 In *Interspeech*, pages 921–925.

890 Siyi Zhou, Yiquan Zhou, Yi He, Xun Zhou, Jinchao  
891 Wang, Wei Deng, and Jingchen Shu. 2025. In-  
892 dextts2: A breakthrough in emotionally expressive  
893 and duration-controlled auto-regressive zero-shot  
894 text-to-speech. *arXiv preprint arXiv:2506.21619*.

895 Han Zhu, Wei Kang, Zengwei Yao, Liyong Guo,  
896 Fangjun Kuang, Zhaoqing Li, Weiji Zhuang, Long  
897 Lin, and Daniel Povey. 2025. Zipvoice: Fast  
898 and high-quality zero-shot text-to-speech with flow  
899 matching. *arXiv preprint arXiv:2506.13053*.