INFINITEAUDIO: INFINITE-LENGTH AUDIO GENERA TION WITH CONSISTENT ACOUSTIC ATTRIBUTES

Anonymous authors

Paper under double-blind review

Abstract

This work aims to generate long-duration audio while preserving acoustic coherence, utilizing existing text-conditional audio generation models through diffusion-based approaches. Current diffusion models, however, encounter significant challenges in generating long audio sequences due to memory constraints, as output size scales with input length. While one possible solution is to concatenate short clips, this often leads to inconsistencies due to a lack of shared temporal information across segments. To address these challenges, we propose InfiniteAudio, a novel inference technique designed to generate long audio with consistent acoustic attributes. Our method is based on three key components. First, we implement a curved denoising approach with a fixed-size input, enabling theoretically infinite audio generation while maintaining a constant memory footprint. Second, we introduce conditional guidance alternation, a mechanism that enhances intelligibility in long speech generation. Finally, initial self-attention features are shared across future frames to maintain temporal coherence. The effectiveness of InfiniteAudio is demonstrated through comprehensive comparisons with existing text-to-audio generation baselines. Generated audio samples are available on our anonymous project page¹.

026 027

029

025

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Diffusion models (Ho et al., 2020; Song et al., 2020b) have received considerable attention across 031 various domains due to their ability to generate high-quality, diverse outputs. They have demon-032 strated impressive results in tasks including image generation (Dhariwal & Nichol, 2021; Rombach 033 et al., 2022), video generation (Ho et al., 2022; Singer et al., 2022; Wang et al., 2023), and text-to-034 audio (TTA) generation (Liu et al., 2023; Huang et al., 2023; Lee et al., 2024; Liu et al., 2024). TTA 035 models generate audio from text description prompts and typically utilize generative frameworks 036 such as latent diffusion models (Rombach et al., 2022) as illustrated in Fig. 1(a) or flow matching 037 models (Vyas et al., 2023). Recently, VoiceLDM (Lee et al., 2024) has advanced this capability by 038 generating both speech and background audio simultaneously, as shown in Fig. 1(b). The generated 039 speech not only reflects the background description prompt but also adapts to the content prompt. 040 For example, when prompted with "Hello" in a cathedral setting, the speech will naturally include reverberation to match the environment. 041

Despite these advancements, existing TTA generation models based on diffusion approaches face significant challenges when generating longer audio sequences. To extend the output size during inference, the input size must also be increased, given that diffusion models require the input and output dimensions to remain unchanged. Moreover, these models struggle to manage long text conditions when producing extended speech. While long audio can be generated by concatenating short clips created by existing TTA models, ensuring a smooth and continuous audio stream remains challenging due to the lack of temporal consistency between inter-clip segments.

To address these challenges, we introduce InfiniteAudio, a novel inference method for generating
long and consistent audio. InfiniteAudio generates extended audio by utilizing a fixed input size
with progressively increasing noise levels over time. As shown in Fig. 2, at each inference step, the
fully denoised audio segment at the beginning of the input is removed, while a new random noise

⁰⁵³

¹https://anonymousforcf.github.io/InfiniteAudio/

068

069

071

073 074

100

101

102 103



064 Figure 1: Overview of tasks. (a) InfiniteAudio enables the generation of longer audio using a pre-065 trained text-to-audio model, overcoming the memory limitations faced by existing models. (b) For 066 simultaneous audio and speech generation, InfiniteAudio can generate long speech that accurately reflects the audio description prompt.

Table 1: Comparison of existing diffusion inference methods with our approach. Our method generates longer audio with a fixed memory size.

Methods	Memory requirements	Long generation	Varying timesteps
Diffusion FIFO-Diffusion (Kim et al., 2024) InfiniteAudio	Various Small Very small	Limited capable capable	× ✓

075 latent is added at the end. In this manner, InfiniteAudio can theoretically generate infinite audio 076 frames using a fixed input size, effectively mitigating memory constraints.

077 While FIFO-Diffusion (Kim et al., 2024), which is designed for text-to-video (TTV) generation, 078 also employs a fixed input size, it utilizes all diffusion sampling steps. In contrast, as illustrated 079 in Fig. 2, our method chooses the more important steps rather than using the entire steps. This 080 selective approach, which we refer to as curved denoising, reduces the number of required sampling 081 steps while attaining high-quality generation, resulting in more efficient inference. Tab. 1 presents a 082 comparison of traditional diffusion inference methods, FIFO-Diffusion, and our proposed approach.

Additionally, InfiniteAudio addresses the challenge of generating long audio sequences from ex-084 tended text inputs by segmenting the text and applying a guidance alternation technique. By dividing 085 long text prompts into smaller segments, we reduce memory overhead. However, when processing consecutive prompts, the generated audio can be affected by preceding segments, leading to reduced 087 intelligibility. To resolve this, we propose a guidance alternation strategy that switches between con-880 ditional and unconditional guidance when processing following text inputs. This approach preserves the clarity of long speech while minimizing interference between segments.

While this method effectively handles extended speech, generating audio from multiple distinct 091 prompts within a single clip can disrupt coherence, as the prompts often lack the specificity required 092 to retain consistent speaker characteristics. To mitigate this issue, we share a query, key, and value (QKV) features within the self-attention layers of the diffusion model. Propagating the initial QKV 094 features across successive segments ensures uniform speaker attributes and continuity, preserving 095 vocal consistency and maintaining intelligibility of both background audio and speech, even with varying text inputs. Our experiments demonstrate that InfiniteAudio can generate extremely long 096 and coherent audio without any degradation in quality over time. Further details are provided in Sec 4. 098

- 099 Our contributions can be summarized as follows.
 - We propose InfiniteAudio, a method for generating long-duration audio without additional training, addressing memory limitations in existing TTA models using diffusion techniques.
 - We introduce curved denoising, which selectively applies key diffusion steps, improving efficiency.
- 105 • We suggest a conditional guidance alternation mechanism to support multiple speech conditions within a single audio stream, maintaining intelligibility. 107
 - We implement QKV sharing in self-attention, ensuring consistent speech generation.



Figure 2: Overall pipeline for the existing method and our method. Traditional diffusion models apply the same diffusion timestep across inputs during inference. Our method starts with a latent containing varying timesteps and skips unimportant timesteps for P multiple big steps. For every inference step, an audio frame reaching $\tau = 1$ is popped out and an audio frame with noise is inserted to maintain a same input size. This method theoretically allows infinite audio generation with constant memory usage, producing one audio frame per step.

131

132

123

124

125

126

127

2 RELATED WORKS

2.1 TEXT TO AUDIO AND SPEECH GENERATION

133 TTA generation (Liu et al., 2023; Kreuk et al., 2022; Yang et al., 2023) has attracted considerable 134 attention in recent years, driven by advancements in generative modeling techniques (Ho et al., 135 2020; Song et al., 2020b). Several works (Liu et al., 2023; Ghosal et al., 2023; Yang et al., 2023) 136 use the latent diffusion model (LDM) (Rombach et al., 2022) to generate audio, mitigating the 137 large computational costs of the original diffusion process. In the diffusion-based TTA models, 138 contrastive language audio pretraining (CLAP) (Wu et al., 2023) is utilized in many models (Liu 139 et al., 2023; Huang et al., 2023; Yuan et al., 2024), in order to align language and audio embeddings. 140 Additionally, large language models (LLMs) are exploited due to their strong text understanding 141 capabilities (Ghosal et al., 2023; Liu et al., 2024).

142 Besides TTA, text-to-speech (TTS) generation is also an active area of research, with early models 143 using autoregressive (AR) models (Wang et al., 2017; Oord et al., 2016). To address the issue of 144 slow inference speed that arises in AR models, researchers have introduced non-AR models (Ren 145 et al., 2019; Kim et al., 2020) that offer improved performance compared to AR models. Using the 146 diffusion model, Grad-TTS (Popov et al., 2021) produces high-quality speech with a score-based 147 decoder. Furthermore, several works have addressed environment-related speech generation. For example, VoiceLDM (Lee et al., 2024) introduces an efficient model that generates audio closely 148 aligned with both descriptive and content prompts. Audiobox (Vyas et al., 2023), a unified model 149 based on flow matching, can produce audio that contains various audio conditions, such as non-150 verbal sounds (e.g., coughing, screaming) or acoustic environments (e.g., rural, stadium, indoor). 151 While these environment-related speech generation models produce high-quality results, they strug-152 gle to generate long outputs that containing multiple sentences simultaneously. 153

154

2.2 LONG GENERATION USING DIFFUSION MODELS

155 156

Producing large-scale output with diffusion models is challenging due to the high computational costs and memory footprints. For image generation, Multi-Diffusion (Bar-Tal et al., 2023) and SyncDiffusion (Lee et al., 2023) use several windows to generate images with arbitrary aspect ratios but focus on smoothing the overlap regions of windows, falling short of solving the repetition problem. Scalecrafter (He et al., 2023) dynamically increases the receptive field and succeeds in generating ultra-high-resolution images up to 4096 × 4096.

162 In addition to image generation, research into long video generation has become increasingly active. 163 Many AR models (He et al., 2022; Voleti et al., 2022; Harvey et al., 2022; Chen et al., 2023) can 164 generate long videos, but due to error accumulation and a lack of temporal consistency between the 165 frames, there are quality issues. FreeNoise (Qiu et al., 2023) addresses this issue with a window-166 based function but cannot generate infinitely long videos as it requires memory proportional to the output length. FIFO-Diffusion (Kim et al., 2024) can produce infinitely long videos with a fixed 167 amount of memory by conducting diagonal diffusion across different timesteps. 168

169 In the audio domain, residual vector quantization (RVQ) is widely used to generate audio faster and 170 more efficiently (Défossez et al., 2022; Zeghidour et al., 2021). SoundStorm (Borsos et al., 2023b), 171 which combines RVQ with AudioLM (Borsos et al., 2023a), efficiently generates audio sequences 172 up to 30 seconds long, a relatively extended length. To generate longer audio, (Evans et al., 2024a) tackles this issue by leveraging LDMs, resulting in output up to 95 seconds long. Moreover, (Evans 173 et al., 2024b) leverage the diffusion transformer(DiT) to generate even longer audio, stretching up 174 to 4m 45s. However, these approaches (Evans et al., 2024a;b) require additional training on datasets 175 that match the desired output length. Our method does not require any additional training to generate 176 audio of theoretically infinite length. 177

178 179

211 212

214

3 **TEXT-TO-AUDIO DIFFUSION MODELS**

We briefly summarize the outline of existing text-to-audio generation models, focusing on two rep-181 resentative models: AudioLDM (Liu et al., 2023) and VoiceLDM (Lee et al., 2024). TTA models 182 produce audio based on given text prompts and deal with audio as an image since audio can be 183 represented as a 2D mel-spectrogram consisting of time and frequency axes. 184

185 Many TTA models commonly include the following modules: audio and text encoders, an audio decoder and vocoder, and a latent diffusion model. Through these modules, TTA models can learn the distribution of mel-spectrograms corresponding to a text prompt y. For the audio $f_{audio}(\cdot)$ and 187 text encoder $f_{text}(\cdot)$, many models (Liu et al., 2023; Lee et al., 2024; Huang et al., 2023) leverage 188 a contrastive language-audio pretraining (CLAP) encoder, which is trained to align text and audio 189 modalities (Wu et al., 2023). These encoders encode a start latent and conditions that are exploited 190 in the latent space, while the decoder reconstructs the mel-spectrogram denoted by $a \in \mathbb{R}^{T \times F}$ 191 from the latent $\mathbf{z}_1 \in \mathbb{R}^{C \times \frac{T}{r} \times \frac{F}{r}}$, where T represents the time dimension, F denotes the frequency 192 dimension, C refers to the channel dimension, $\tau \sim \mathcal{U}([1,...,M])$ is the diffusion timestep, and r 193 is the compression factor. The vocoder produces a waveform from the predicted mel-spectrogram. 194 The LDM is trained to denoise a perturbed version of the latent from z_{τ} to z_1 . 195

For noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the text condition $\mathbf{c} = f_{text}(y)$, AudioLDM is trained to minimize the 196 following loss: 197

$$\mathcal{L}_{AudioLDM} = \mathbb{E}_{\mathbf{z}_0, \boldsymbol{\epsilon}, \tau} \left[\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_{\tau}, \tau, \mathbf{c}) \|_2^2 \right].$$
(1)

, where ϵ_{θ} is the predicted diffusion score. Unlike AudioLDM, VoiceLDM generates not only 199 general audio but also produces clean speech and speech that reflects background sounds. Therefore, 200 the model uses two text prompts: a description prompt c_{desc} and a content prompt c_{cont} . Similar 201 to AudioLDM, VoiceLDM uses a CLAP encoder and latent diffusion model architectures. The 202 objective for VoiceLDM is as follows: 203

$$\mathcal{L}_{VoiceLDM} = \mathbb{E}_{\mathbf{z}_{0},\boldsymbol{\epsilon},\tau} \left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_{\tau},\tau,\mathbf{c}_{desc},\mathbf{c}_{cont})\|_{2}^{2} \right].$$
(2)

204 The model uses dual classifier-free guidance (Ho & Salimans, 2022) to control audio description 205 prompt and a text content prompt. Therefore, the diffusion score $\tilde{\epsilon}$ is formulated as follows: 206

$$\tilde{\epsilon}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) = \epsilon_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) + w_{desc} \left(\epsilon_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \emptyset_{cont} \right) - \epsilon_{\theta} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) + w_{cont} \left(\epsilon_{\theta} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \mathbf{c}_{cont} \right) - \epsilon_{\theta} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right)$$
(3)

where w is the guidance weight and \emptyset indicates the null condition.

INFINITEAUDIO 4 213

In this section, we describe how to generate infinite-length audio with a fixed memory footprint 215 using pretrained TTA models. Additionally, we introduce a method for generating longer, consistent speech. We explore two representative TTA models, AudioLDM and VoiceLDM, and suggest suitable inference techniques respectively.

219 220

221

4.1 CURVED DENOISING WITH REDUCED SAMPLING STEPS

Because diffusion models are trained to estimate the noise in the input, both the input and output must have the same size. This inherent property increases memory demands, as generating longer audio requires longer inputs. For instance, on a single GPU with 12GB of memory, AudioLDM is limited to generating audio no longer than 22.5 seconds.

To overcome memory limitations, we propose InfiniteAudio, which operates with a fixed input size but can theoretically generate infinite-length audio. Inspired by FIFO-Diffusion (Kim et al., 2024), which addresses text-to-video generation, we initiate the diffusion inference process with a fixed-length audio segment using a small portion of the output predicted by existing TTA models. Although mel-spectrograms, which consist of time and frequency axes, can be treated as images, audio generation tasks must handle temporal information, similar to video generation. Therefore, we treat the input latent $\mathbf{z}_{\tau} \in \mathbb{R}^{C \times \frac{T}{r} \times \frac{F}{r}}$ as containing $\frac{T}{r}$ audio frames, analogous to video frames. Each encoded mel-spectrogram frame corresponds to $\mathbf{z}_{1}^{i} \in \mathbb{R}^{C \times 1 \times \frac{F}{r}}$, where $i \in [1, \frac{T}{r}]$.

For infinite audio generation, noise is progressively added to the input audio frames over time, ex-235 cept for the initial frames, which act as a 'buffer zone' with no added noise. Since no further training 236 occurs in InfiniteAudio, using different diffusion timesteps during inference can produce a perfor-237 mance gap. The buffer zone mitigates this by applying the same timesteps as during training, helping 238 to reduce the performance gap. Beyond the buffer zone, the latent frames gradually transition: the 239 earlier frames are almost fully predicted, whereas the $\frac{T}{r}$ -th frame is treated as Gaussian noise. The 240 input to the inference stage consists of the initial buffer frames and $\frac{T}{r}$ frames with varying noise 241 levels. As represented in Fig. 2 (b), after each inference step, the first frame following the buffer 242 zone reaches diffusion timestep $\tau = 1$ and is then removed. To maintain a total of $\frac{T}{r}$ frames, we 243 insert a new noisy frame at the $\frac{T}{r}$ -th position. By repeating this process iteratively, we can generate 244 N frames in N inference steps. To effectively tackle memory limitation issues, InfiniteAudio keeps 245 the input size constant during inference, regardless of the desired output length. However, employ-246 ing the full set of diffusion timesteps still necessitates long input sequences. To mitigate this issue, 247 InfiniteAudio reduces input size by selecting only the most critical diffusion steps. By leveraging 248 deterministic denoising (Song et al., 2020a), existing models perform inference without requiring all 249 timesteps. Additionally, we found that we can further reduce the number of timesteps by skipping 250 unimportant steps while still preserving sample quality.

251 We first identify the most important timesteps 252 of the three segments, initial, middle, and fi-253 nal, during inference for both AudioLDM and 254 VoiceLDM. Since the attention scores for both 255 models reflect the relevance of one frame (Key) to another frame (Query), we analyze the self-256 attention maps in the diffusion U-Net decoder 257 modules. 258

As shown in Fig. 3, in AudioLDM, the query sequences are primarily influenced by the initial key sequences, which correspond to earlier frames or cleaner inputs. In contrast, VoiceLDM behaves differently: the query sequences are more influenced by later key sequences, which correspond to noisier inputs.



Figure 3: Attention maps denoting the importance of timesteps in the input sequences. In AudioLDM, the query in the last sequence segment focuses on the initial portions of the audio. In contrast, VoiceLDM demonstrates a stronger correlation with the later segments in its final query.

Since some initial frames lie within a buffer zone, we focus on regions beyond this zone.

267 Consequently, we allocate more timesteps to critical regions with high attention scores and skip
 268 less crucial timesteps, significantly reducing the overall number of inference steps and input size.
 269 This strategy, dubbed as curved denosing, enables us to achieve similar output quality with fewer
 269 computations compared to the traditional method, which uses N timesteps for N frames.



Figure 4: Illustration of guidance alternation method. Long sentence token c_{cont} is divided into several sentence tokens such as \mathbf{c}_{cont1} , \mathbf{c}_{cont2} , and so on. We apply existing conditional guidance to odd-numbered sentence prompts and switch to unconditional guidance for even-numbered sentence prompts. This alternation helps reduce the influence of one segment on the generation of the next, improving overall coherence in the generated audio.

285 4.2 LONG SPEECH GENERATION

287 In addition to generating long sounds, it is essential to generate extended speech. However, gen-288 erating long speech with LDM by encoding a long content prompt at once with a text encoder is 289 challenging due to memory limitations. To overcome memory limitation, we first split the long 290 content prompt \mathbf{c}_{cont} into smaller, manageable sentence segments: $\mathbf{c}_{cont1}, \mathbf{c}_{cont2}, ..., \mathbf{c}_{contk}$. Each 291 segment is then applied to its corresponding audio section.

292 After generating the first sentence segment, the model faces challenges in immediately processing 293 the next sentence prompt due to differing diffusion timesteps in the input. The part of latent which is in its final timesteps becomes confused when it receives a new sentence prompt, as it has already 295 processed using the previous sentence prompt. Therefore, we need to eliminate the residual effects 296 from the previous sentence before generating the next one. As shown in Fig. 4, rather than simply 297 sequencing the sentence segments $\mathbf{c}_{cont2}, ..., \mathbf{c}_{contk}$, which causes interference between sentences, 298 we devise a novel guidance alternation method.

299 Considering the guidance scale in Eq. 3, we utilize existing guidance differently depending on the 300 sentence prompts are odd-numbered or even-numbered prompts denoted in Eq. 4 and 5, respectively. 301

$$\tilde{\boldsymbol{\epsilon}}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) = \boldsymbol{\epsilon}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) + w_{desc} \left(\boldsymbol{\epsilon}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \boldsymbol{\vartheta}_{cont} \right) - \boldsymbol{\epsilon}_{\theta}^{even} \left(\mathbf{z}_{\tau}, \boldsymbol{\vartheta}_{desc}, \boldsymbol{\vartheta}_{cont} \right) \right) + w_{cont} \left(\boldsymbol{\epsilon}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \boldsymbol{\vartheta}_{desc}, \mathbf{c}_{cont} \right) - \boldsymbol{\epsilon}_{\theta}^{even} \left(\mathbf{z}_{\tau}, \boldsymbol{\vartheta}_{desc}, \boldsymbol{\vartheta}_{cont} \right) \right)$$
(4)

 $\tilde{\boldsymbol{\epsilon}}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) = \boldsymbol{\epsilon}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right)$

$$\int \left(\frac{1}{2} \right) dt = \left(\frac{1}{2$$

$$\begin{array}{l} = \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) = \mathbf{c}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont} \right) \\ + w_{desc} \left(\mathbf{\epsilon}_{\theta} \left(\mathbf{z}_{\tau}, \mathbf{c}_{desc}, \emptyset_{cont} \right) - \mathbf{\epsilon}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc}, \emptyset_{cont} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{z}_{\tau}, \emptyset_{desc} \right) \\ = \left(\mathbf{c}_{\theta}^{odd} \left(\mathbf{c}_{\tau}, \emptyset_{desc} \right) \\ = \left(\mathbf{c}_$$

$$-w_{cont}\left(\boldsymbol{\epsilon}_{\theta}^{odd}\left(\mathbf{z}_{\tau}, \boldsymbol{\emptyset}_{desc}, \boldsymbol{\emptyset}_{cont}\right) - \boldsymbol{\epsilon}_{\theta}^{even}\left(\mathbf{z}_{\tau}, \boldsymbol{\emptyset}_{desc}, \mathbf{c}_{cont}\right)\right)$$
(5)

For even-numbered sentence prompts, such as \mathbf{c}_{cont2} , \mathbf{c}_{cont4} , \mathbf{c}_{cont6} , we switch from conditional 311 guidance to unconditional guidance to mitigate the influence of the previous sentence. Specifically, 312 we treat existing unconditional guidance as conditional guidance and vice versa for odd-numbered 313 prompts. The red-highlighted sections indicate the existing conditional guidance, while the green-314 highlighted sections represent the alternated conditional guidance. This approach helps ensure that 315 sentences are generated accurately and completely. Furthermore, to balance the guidance alternation 316 method, we apply a negative sign to the guidance weights, wcont. As a result, we no longer need to 317 input the entire sentence prompt, which significantly reduces memory and computational demands 318 while enhancing speech intelligibility.

319

302 303

305 306

307

308

310

280

281

282

283

284

320 4.3 CONSISTENT SPEECH GENERATION

321

While long speech generation can be efficiently achieved by separating sentences and employing 322 the guidance alternation method, it is essential to maintain consistent speaker attributes throughout 323 the generation process. Splitting long sentences into multiple tokens can lead to inconsistencies,

such as fluctuations between male and female voices. To address this issue, we propose sharing
query, key, and value (QKV) features within the self-attention layers of the U-Net architecture in
the diffusion model, as demonstrated in Fig. 5. In video editing and image translation research,
self-attention layers in the diffusion model are critical as they can determine the overall structure of
the image (Tumanyan et al., 2023; Ceylan et al., 2023). Query, key, or value pairs of the previous
image are used for the next image generation to maintain the overall image concept.

330 Unlike previous works that focus on maintain-331 ing consistency within a single image by uti-332 lizing a single image (Tumanyan et al., 2023; 333 Ceylan et al., 2023), we aim for long-term con-334 sistency by sharing the QKV pairs from the initial speech throughout the generation of subse-335 quent sentences in the U-Net upsampling lay-336 ers. We retain a series of QKV features for the 337 initial T audio frames and utilize them during 338 the generation of later sentences. We empir-339 ically demonstrate a method for selecting the 340 appropriate QKV pairs and determining the op-341 timal number of frames to share. 342



Figure 5: To maintain a speaker characteristic, QKV paris in initial T audio frames are utilized. For following sentence generation, the model loads the QKV pairs for speech consistency.

5 EXPERIMENT

We present generated long audio using InfiniteAudio, built on pretrained AudioLDM and VoiceLDM models, and evaluate them quantitatively and qualitatively. Furthermore, we perform ablation studies on Sec. 4. For more audio samples and additional ablation studies, see App. B and D.

349 350 351

343 344

345 346

347

348

5.1 EXPERIMENTAL SETTINGS

352 **Datasets.** To evaluate our method on TTA generation, we randomly selected 500 audio-text pairs 353 from the 975 test files in the Audiocaps dataset (Kim et al., 2019), which is commonly used for evaluating existing TTA models. For text-to-long-speech (TTLS) generation, we constructed a test set 354 from the English subset of the CommonVoice 13.0 corpus (Ardila et al., 2019), randomly selecting 355 60 text samples, each consisting of more than five sentences, to assess long-form speech generation. 356 For TTAS, which involves generating both audio and speech, we used 60 randomly selected text 357 sets from the CommonVoice 13.0 corpus (Ardila et al., 2019). The test set for audio description was 358 sourced from the Audiocaps test set (Kim et al., 2019), specifically focusing on samples from the 359 "speech" category, such as those involving "talking" or "speaking." 360

Baselines. For comparison, we evaluate the performance of InfiniteAudio against two publicly available TTA models: AudioLDM² and VoiceLDM³. Notably, VoiceLDM is currently the only available model capable of generating both audio and speech simultaneously, making it the sole candidate for our TTLS and TTAS experiments.

364 365 366

361

362

363

Evaluation Metrics. We employ several quantitative metrics to evaluate the audio quality and the 367 alignment between the input text prompt and the generated audio. These metrics include Frechet 368 Distance (FD), Kullback-Leibler (KL) divergence, and the CLAP score, which are standard in text-369 to-audio generation evaluations (Liu et al., 2023; Lee et al., 2024; Vyas et al., 2023). FD and KL 370 divergence quantify how closely the generated audio matches the ground truth, with lower values 371 indicating better performance. The CLAP score, in contrast, assesses the relevance between the text 372 prompts and the generated audio, where higher values are preferable. For subjective evaluation of 373 the audio produced by TTA models, we use two metrics: (i) overall quality (OVL) and (ii) relevance to the input text description (REL). Both were rated on a scale of 1 to 5 by 20 domain experts, based 374 on 30 speech samples. Further details on the human evaluation process are available in App. C.3. 375

376 377

³https://github.com/glory20h/VoiceLDM

²https://github.com/haoheliu/AudioLDM

Method	CLAP↑	FD↓	$KL\downarrow$	OVL↑	REL↑
Ground Truth	0.5276	NA	NA	4.11±0.22	4.03±0
AudioLDM (Liu et al., 2023)	0.4908	44.6689	2.0805	3.03±0.23	3.06±0
InfiniteAudio w/ Equally spaced timesteps	0.3832	54.7479	2.4013	2.19 ± 0.21	2.33 ± 0
InfiniteAudio w/ Middle focused timesteps	0.3979	56.7792	2.6077	2.06 ± 0.19	2.18 ± 0
InfiniteAudio w/ Last focused timesteps	0.4559	43.3788	1.9650	2.63 ± 0.18	2.80 ± 0
InfiniteAudio w/ Initial focused timesteps	0.3110	67.0704	2.9838	2.13 ± 0.19	2.07 ± 0
VoiceLDM (Lee et al., 2024)	0.4199	51.4019	2.2749	2.53±0.24	2.41±0
InfiniteAudio w/ Equally spaced timesteps	0.3729	59.1521	2.4477	2.20 ± 0.21	2.33 ± 0
InfiniteAudio w/ Middle focused timesteps	0.3779	56.7321	2.4622	2.10 ± 0.20	2.41 ± 0
InfiniteAudio w/ Last focused timesteps	0.3542	64.8813	2.6227	2.38 ± 0.23	2.24 ± 0
InfiniteAudio w/ Initial focused timesteps	0.4107	51.5047	2.3498	2.38 ± 0.23	2.48±0

Table 2: Quantitative evaluations on TTA. Our method for both models achieves comparable results, 378 even surpassing original inference results. 379

393

394

395

396

397

398

399

380 381 382

> To assess speech intelligibility, we measure word error rate (WER) and character error rate (CER) using the Whisper automatic speech recognition (ASR) model (Radford et al., 2023), where lower scores indicate better intelligibility.

To evaluate voice consistency within a single audio sample, we utilize the Resemblyzer Python package⁴, which is commonly employed for extracting speaker embeddings, alongside the VoxCelebdisentangler model (Nam et al., 2024), which offers a high-level representation of speakers. We calculate the cosine similarity score C_{sim} based on the first 10 seconds of the speaker embedding $G(a_{:10})$, where G represents the speaker verification model used in the aforementioned methods. 400 Additionally, we assess 5-second segments of the embedding, $G(a_{5+5h:10+5h})$, where h = 1, 2, 3, 3, 3and 4. This score C_{sim} quantifies the similarity between two vectors in an inner product space.

401 402 403

404

417

5.2 QUANTITATIVE RESULTS

405 Memory Consumptions. We compare the 406 memory consumption of the existing TTA mod-407 els with our method. Since VoiceLDM re-408 stricts generation to 10 seconds, we conduct our 409 experiments using AudioLDM. AudioLDM's 410 memory usage increases as the length of the 411 generated audio grows. In contrast, our method 412 maintains consistent memory usage, regardless of the desired audio length, as demonstrated in 413 Fig. 6. Consequently, our method allows for 414 generating longer audio content without signif-415 icant performance degradation. 416



Figure 6: Comparisons on memory consumption between AudioLDM (Liu et al., 2023) and our method.

Generation Evaluations. We demonstrate the effectiveness of curved denoising strategy in Tab. 2. 418 While our framework is designed to generate long audio with fixed memory based on a pre-trained 419 model, it not only matches the performance of existing models but also achieves higher scores. We 420 demonstrate that our approach considering input attention relations represented in Fig. 3 achieves su-421 perior performance compared to the equally spaced timesteps that are used in FIFO-Diffusion (Kim 422 et al., 2024) or other strategies with the same steps. 423

Since VoiceLDM generates long speech in 10-second segments, it exhibits a significantly higher 424 WER than InfiniteAudio, as demonstrated in Tab. 3. In contrast, InfiniteAudio's alternating guid-425 ance strategy further reduces both WER and CER, enhancing sentence intelligibility by mitigating 426 interference between sentences in text-to-long speech (TTLS) generation. For TTAS evaluation, 427 our method delivers superior performance over existing approaches, particularly in WER and CER 428 scores, with only a slight reduction in the CLAP score. As there is no ground truth for generating 429 both audio and speech simultaneously, we focus our evaluation on WER, CER, and CLAP scores. 430

⁴https://github.com/resemble-ai/Resemblyzer

Task	Method	WER↓	CER↓	CLAP↑
TTLS	VoiceLDM w/ InfiniteAudio w/ InfiniteAudio and Guidance alternation w/ InfiniteAudio, Guidance alternation and QKV sharing	0.5363 0.5810 <u>0.3376</u> 0.3038	0.4595 0.5119 <u>0.2635</u> 0.2368	NA
TTAS	VoiceLDM w/ InfiniteAudio w/ InfiniteAudio and Guidance alternation w/ InfiniteAudio, Guidance alternation and QKV sharing	0.8038 0.4604 0.3863 0.3824	0.6070 0.3492 0.2825 0.2888	0.1252 0.0988 <u>0.1200</u> 0.0877

Table 3: Evaluation on TTS and TTAS. Guidance alternation and QKV sharing method can further decreases WER, which contribute to speech intelligibility.

Table 4: Speaker consistency evaluation across different time regions. C_{sim}^{h} represents the cosine similarity score between the first 10 seconds of speech and the subsequent 5 seconds. With QKV sharing, the speaker embeddings remain consistent throughout the entire duration.

Task	Method	$C^1_{sim}\uparrow$	$\underset{C_{sim}^{2}}{\operatorname{Resem}}\uparrow$	blyzer $C_{sim}^3 \uparrow$	$C^4_{sim}\uparrow$	VoxCele $C_{sim}^1 \uparrow$	eb-disentangl $C_{sim}^2 \uparrow$	er (Nam et al $C_{sim}^3 \uparrow$., 2024) $C_{sim}^4 \uparrow$
TTLS	InfiniteAudio	0.7810	0.7564	0.7479	0.8218	0.4802	0.4336	0.4524	0.5481
	+ QKV sharing	0.7900	0.7680	0.7658	0.8415	0.5310	0.5207	0.5400	0.6154
TTAS	InfiniteAudio	0.8101	0.8082	0.8019	0.8718	0.5016	0.4983	0.4997	0.6096
	+ QKV sharing	0.8406	0.8183	0.8254	0.8810	0.5699	0.5280	0.5463	0.6236

Voice Consistency. The test set comprises 60 long text samples, identical to those utilized in the TTLS and TTAS evaluations. As shown in Tab. 4, sharing QKV yields speaker embedding features that are more closely aligned across the entire speech segments while preserving speech intelligibility.

5.3 QUALITATIVE RESULTS

Sampling Strategies. We propose curved de-noising, where the sampling strategy is de-termined by considering attention scores. As shown in Fig. 7, in contrast to other strate-gies, which show interruptions in the gener-ated audio as observed in the spectrograms, our method using initial-focused timesteps ensures continuous audio generation, as evidenced by both the spectrogram and the CLAP score.

470
471
471
472
472
473
474
475
476
476
477
477
478
478
479
479
479
470
470
470
470
470
470
471
471
472
473
473
473
474
474
475
475
476
476
477
477
478
478
478
479
479
479
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470
470



Figure 7: Analysis on various diffusion sampling strategies on VoiceLDM (Lee et al., 2024).

ing audio that adheres to the description prompts for durations exceeding 10 seconds. For evaluating TTLS generation, we utilize WER as the primary metric. Our method consistently produces intelligible speech, whereas VoiceLDM often struggles, frequently distorting speech segments to fit within a 10-second constraint. Generating coherent audio and speech simultaneously, especially for extended durations, is challenging due to the need to satisfy both the content prompt c_{cont} and the description prompt c_{desc} . In contrast, InfiniteAudio effectively generates speech that aligns with both prompts.

482 5.4 ABLATION STUDY ON QKV SHARING

Figure 9 demonstrates our method for selecting QKV features, showing that sharing QKV features
 consistently outperforms other approaches across all metrics. Notably, it achieves higher cosine similarity scores, indicating better voice consistency, compared to the non-sharing method.



Figure 8: Qualitative results for TTA, TTLS, and TTAS. InfiniteAudio generates high-quality long audio that accurately follows both the audio description prompt and the speech content prompt.





5.5 ANALYSIS ON SAMPLING STEPS AND AUDIO LENGTH

510 InfiniteAudio is designed to minimize the number of sampling steps while preserving high au-511 dio quality. As shown in Tab. 2, our method 512 outperforms other strategies with the same 513 number of steps. Furthermore, even when com-514 pared to methods that increase sampling steps 515 to 200 or 250 using equally spaced timesteps, 516 our approach consistently achieves excellent 517 scores across all metrics, despite utilizing fewer 518 than 150 steps, as denoted in Tab. 5.

519 InfiniteAudio aims to generate longer audio se-520 quences while maintaining high quality. As 521 shown in Tab. 6, compared to the fixed 10-522 second generation, it produces comparable re-523 sults across a range of lengths, from 10 to 20 524 seconds. Notably, the CLAP score for this ex-525 periment is measured using a different check-526 point from the one used in other tables, as ex-

Table 5: Comparison of sampling steps between VoiceLDM (Lee et al., 2024) and InfiniteAudio. InfiniteAudio requires fewer than 150 steps to achieve superior results.

Sampling steps	CLAP↑	FD↓	$KL {\downarrow}$
w/ 200 equally spaced steps w/ 250 equally spaced steps	0.3923 <u>0.3941</u>	53.0555 50.5447	2.3334 2.3937
InfiniteAudio	0.4107	<u>51.5047</u>	<u>2.3498</u>

Table 6: Comparison of generated audio lengths between a fixed duration of 10 seconds and variable-length generation approaches.

Generated audio length	CLAP↑	FD↓	KL↓
Fix	0.3207	43.3788	1.9650
Various	0.3257	48.3701	1.9058

527 periments involving varying audio lengths require a distinct CLAP model⁵.

6 CONCLUSION

530 531

528 529

498

499

500

501

502

504

505

506

507

509

0 CONCLUSION

We introduce InfiniteAudio, a novel inference method designed to generate infinitely long, consistent audio using pretrained text-to-audio models. InfiniteAudio effectively maintains a fixed memory footprint, addressing the memory limitations of existing models. Additionally, we propose a new guidance alternation method that can produce long speech with high intelligibility. By sharing QKV pairs in the self-attention layers, InfiniteAudio ensures consistent speech generation and mitigates issues such as voice variations. These contributions open up possibilities for long text-to-audio generation and pave the way for continuous, coherent long audio content.

⁵³⁹

⁵https://github.com/LAION-AI/CLAP

540 REFERENCES

547

554

561

562

565

570

576

- Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer,
 Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. Common voice: A
 massively-multilingual speech corpus. *arXiv:1912.06670*, 2019.
- Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. *Proc. ICML*, 2023.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. Audiolm: a language modeling approach to audio generation. *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, 31:2523–2533, 2023a.
- Zalán Borsos, Matt Sharifi, Damien Vincent, Eugene Kharitonov, Neil Zeghidour, and Marco
 Tagliasacchi. Soundstorm: Efficient parallel audio generation. *arXiv:2305.09636*, 2023b.
- Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *Proc. ICCV*, 2023.
- Xinyuan Chen, Yaohui Wang, Lingjun Zhang, Shaobin Zhuang, Xin Ma, Jiashuo Yu, Yali Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. Seine: Short-to-long video diffusion model for generative transition and prediction. In *Proc. ICLR*, 2023.
 - Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio compression. *arXiv:2210.13438*, 2022.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 2021.
- Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. Fast timing-conditioned latent audio diffusion. *arXiv:2402.04825*, 2024a.
- Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Long-form music generation with latent diffusion. *arXiv:2404.10301*, 2024b.
- Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing
 Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for
 audio events. In *Proc. ICASSP*, 2017.
- Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio gener ation using instruction-tuned llm and latent diffusion model. *arXiv:2304.13731*, 2023.
- William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible
 diffusion modeling of long videos. *NeurIPS*, 2022.
- 579 Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv:2211.13221*, 2022.
 581
- Yingqing He, Shaoshu Yang, Haoxin Chen, Xiaodong Cun, Menghan Xia, Yong Zhang, Xintao
 Wang, Ran He, Qifeng Chen, and Ying Shan. Scalecrafter: Tuning-free higher-resolution visual
 generation with diffusion models. In *Proc. ICLR*, 2023.
- ⁵⁸⁵ Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv*:2207.12598, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
 Fleet. Video diffusion models. *NeurIPS*, 2022.
- Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin
 Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. In *Proc. ICML*, 2023.

594 Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. AudioCaps: Generating 595 Captions for Audios in The Wild. In NAACL-HLT, 2019. 596 Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-tts: A generative flow for 597 text-to-speech via monotonic alignment search. NeurIPS, 2020. 598 Jihwan Kim, Junoh Kang, Jinyoung Choi, and Bohyung Han. Fifo-diffusion: Generating infinite 600 videos from text without training. arXiv:2405.11473, 2024. 601 Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, 602 Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. 603 arXiv:2209.15352, 2022. 604 605 Gihyun Kwon, Jangho Park, and Jong Chul Ye. Unified editing of panorama, 3d scenes, and videos 606 through disentangled self-attention injectionf. arXiv:2405.16823, 2024. 607 608 Yeonghyeon Lee, Inmo Yeon, Juhan Nam, and Joon Son Chung. Voiceldm: Text-to-speech with environmental context. In Proc. ICASSP, 2024. 609 610 Yuseung Lee, Kunho Kim, Hyunjin Kim, and Minhyuk Sung. Syncdiffusion: Coherent montage via 611 synchronized joint diffusions. *NeurIPS*, 2023. 612 613 Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, 614 Audioldm: Text-to-audio generation with latent diffusion models. and Mark D Plumbley. 615 arXiv:2301.12503, 2023. 616 Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu 617 Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation 618 with self-supervised pretraining. IEEE/ACM Trans. on Audio, Speech, and Language Processing, 619 2024. 620 621 KiHyun Nam, Hee-Soo Heo, Jee-weon Jung, and Joon Son Chung. Disentangled representation 622 learning for environment-agnostic speaker recognition. Proc. Interspeech, 2024. 623 Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, 624 Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for 625 raw audio. arXiv:1609.03499, 2016. 626 627 Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-tts: A diffusion probabilistic model for text-to-speech. In Proc. ICML, 2021. 628 629 Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu. 630 Freenoise: Tuning-free longer video diffusion via noise rescheduling. arXiv:2310.15169, 2023. 631 632 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 633 Robust speech recognition via large-scale weak supervision. In Proc. ICML, 2023. 634 Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech: 635 Fast, robust and controllable text to speech. NeurIPS, 2019. 636 637 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-638 resolution image synthesis with latent diffusion models. In Proc. CVPR, 2022. 639 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry 640 Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video 641 data. arXiv:2209.14792, 2022. 642 643 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. 644 arXiv:2010.02502, 2020a. 645 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, 646 and Ben Poole. Score-based generative modeling through stochastic differential equations. 647

arXiv:2011.13456, 2020b.

648	Narek Tumanyan, Michal Gever, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for
649	text-driven image-to-image translation. In <i>Proc. CVPR</i> , 2023.
650	

- Vikram Voleti, Alexia Jolicoeur-Martineau, and Chris Pal. Mcvd-masked conditional video diffusion
 for prediction, generation, and interpolation. *NeurIPS*, 2022.
- Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang,
 Xinyue Zhang, Robert Adkins, William Ngan, et al. Audiobox: Unified audio generation with
 natural language prompts. *arXiv:2312.15821*, 2023.
- Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv:2308.06571*, 2023.
- Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly,
 Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: Towards end-to-end
 speech synthesis. *arXiv*:1703.10135, 2017.
- Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov.
 Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *Proc. ICASSP*, 2023.
- Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu.
 Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, 31:1720–1733, 2023.
- Yi Yuan, Haohe Liu, Xubo Liu, Qiushi Huang, Mark D Plumbley, and Wenwu Wang. Retrieval-augmented text-to-audio generation. In *Proc. ICASSP*, 2024.
- Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. Sound stream: An end-to-end neural audio codec. *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, 30:495–507, 2021.

702 A ALGORITHMS OF INFINITEAUDIO

We present pseudo-code for InfiniteAudio for TTA, TTLS, and TTAS respectively.

A.1 TTA GENERATION WITH CURVED DENOISING

Inpu	t:	
	• N: number of frames	
	• <i>T</i> : audio frames	
	• <i>r</i> : compression factor	
	• $\frac{T}{T} = f$: total timesteps	
	• $\epsilon_{\alpha}(\cdot)$: noise prediction model	
	• Dec(·): decoder	
	• $\{\mathbf{z}^i\}^f$ · initial latent variables	
	$\{\mathbf{z}_{\tau_i}\}_{i=2}^{f}$. Initial factor variables	
	• $\{\tau_i\}_{i=1}$: timesteps	
_	• c _{desc} : description prompt	
Outp	but : <i>v</i> : generated audio sequence	
1: ı	$v \leftarrow []$	
2: 7	$\tau \leftarrow [\tau_1, \tau_{1+P}, \tau_{1+2P};; \tau_{f-2}, \tau_{f-1}, \tau_f]$	▷ Curved denoising with focused initial timeste
3: C	$\mathbf{Q} \leftarrow [\mathbf{z}_1, \mathbf{z}_{1+P}, \mathbf{z}_{1+2P};; \mathbf{z}_{f-2}, \mathbf{z}_{f-1}, \mathbf{z}_f]$	▷ Latent variables for diffusion ste
4: f	for l to N do	▷ Generate N frames of aud
5:	$oldsymbol{Q} \leftarrow oldsymbol{\epsilon}_{ heta}(oldsymbol{Q}, oldsymbol{ au}, \mathbf{c}_{desc}, \emptyset)$	▷ Update latent variables with noise prediction
6:	$\mathbf{z}_{ au_0}^{\iota} \leftarrow oldsymbol{Q}$. dequeue ()	\triangleright Pop out the clean audio fram
7:	$m{v}$.append(Dec $(\mathbf{z}_{ au_0^l}))$	▷ Decode the frame and add to outp
8:	$\mathbf{z}_{ au_f}^{l+ t Len(oldsymbol{Q})} \sim \mathcal{N}(0, \mathbf{I})$	Generate new random noi
9:	$oldsymbol{Q}$.enqueue $(\mathbf{z}_{ au_f}^{l+ t Len(oldsymbol{Q})})$	Insert new noise to latent sequen
10: e	end for	
11: r	eturn v	▷ Return the generated audio sequent

Algori	ithm 2 InfiniteAudio for TTLS (Text-to-Lon	g Speech)
Input:		
	• N: Number of frames	
	• <i>T</i> : audio frames	
	• r: compression factor	
	• f : total timesteps, $\frac{T}{r} = f$	
	• $\epsilon_{\theta_{odd}}(\cdot), \epsilon_{\theta_{even}}(\cdot)$: Conditional/uncondit	onal guidance models
	• $Dec(\cdot)$: Decoder function to generate aud	lio from latent states
	• $\{\mathbf{z}_{\tau}^i\}_{i=2}^f$: Initial latent variables for each	timestep
	• $\{\tau_i\}_{i=1}^{r_i=2}$: Timestep schedule	
	• c _{dase} : Description condition (high-level to	ext description)
	• c _{comt} : Content conditions (detailed text si	plit into segments)
	 sn: Sentence number (starts at 0) 	file into begineins)
Outru	t: Generated audio sequence a	
		∧ Initialize output a
$\frac{1}{2} \tau$		Curved denoising timesten sche
2. 1	$ \begin{bmatrix} 1, 1+P, 1+2P, \dots, 1 \\ f-2, 1 \\ f-1, 1 \end{bmatrix} $	▷ Initialize latent queue for diffu
3. Q	$\leftarrow [\mathbf{z}_1, \mathbf{z}_{1+P}, \mathbf{z}_{1+2P},, \mathbf{z}_{f-2}, \mathbf{z}_{f-1}, \mathbf{z}_f]$	Solit content promote into com
4: [C	$cont1, \mathbf{C}_{cont2}, \dots, \mathbf{C}_{contK}] \leftarrow \mathbf{C}_{cont}$	Spin content prompts into segment lon
5: 5	$\leftarrow [\text{Len}(\mathbf{c}_{cont1}), \text{Len}(\mathbf{c}_{cont2}),, \text{Len}(\mathbf{c}_{cont2}),,,,,,,, .$	htK)] \triangleright Store segment len
6: ϵ_{θ}	$(\cdot) \leftarrow \epsilon_{ heta_{odd}}(\cdot)$	▷ Initialize with conditional guidance (c
7: fo	$\mathbf{r} \ l \ \mathbf{to} \ N \ \mathbf{do}$	▷ Iterate over N fra
8:	while $l \leq S[sn]$ do	▷ Generate frames for the current sente
9:	$oldsymbol{Q} \leftarrow oldsymbol{\epsilon}_{ heta}(oldsymbol{Q},oldsymbol{ au}, \emptyset, \mathbf{c}_{cont})$	▷ Apply noise prediction to latent st
10:	$\mathbf{z}_{ au_0}^\iota \leftarrow oldsymbol{Q}$. dequeue ()	⊳ Pop out clean audio fr
11:	$oldsymbol{v}$. append(Dec $(\mathbf{z}_{ au_0^l}))$	\triangleright Decode and append to out
12:	$\mathbf{z}_{ au_f}^{l+ t Len(oldsymbol{Q})} \sim \mathcal{N}(0, \mathbf{I})$	⊳ Generate new random n
13:	Q .engueue $(\mathbf{z}_{ au_{\ell}}^{l+ ext{Len}(m{Q})})$	▷ Insert new noise into latent qu
14:	if $l = S[sn] - 1$ then	\triangleright Check if we reached the end of the sente
15:	$sn \leftarrow sn + 1$	▷ Move to the next sente
16:	if $\epsilon_{\theta}(\cdot) = \epsilon_{\theta}$, (\cdot) then	▷ Switch conditional guida
17:	$\boldsymbol{\epsilon}_{\boldsymbol{ heta}}(\cdot) \leftarrow \boldsymbol{\epsilon}_{\boldsymbol{ heta}}$ (·)	⊳ Switch to unconditional guida
18:	else	8
19:	$\epsilon_{a}(\cdot) \leftarrow \epsilon_{a}(\cdot)$	⊳ Switch back to conditional guida
20:	end if	8
21.	end if	
21. 22.	end while	
22. 23. en	d for	
23. Ch	turn »	► Return the final generated audio segue
2 4 . 10		

Input		
-	• N: Number of frames	
	• 1: audio frames	
	• r: compression factor	
	• f : total timesteps, $\frac{T}{r} = f$	
	• $\epsilon_{\theta_{odd}}(\cdot), \epsilon_{\theta_{even}}(\cdot)$: Conditional/uncondi	tional guidance models
	• $Dec(\cdot)$: Decoder function to generate as	udio from latent states
	• $\{\mathbf{z}_{\tau_i}^i\}_{i=2}^f$: Initial latent variables for each	h timestep
	• $\{\tau_i\}_{i=1}^{f}$: Timestep schedule	
	• c _{docc} : Description condition (high-level	text description)
	• c Content conditions (detailed text	split into segments)
	Cont. Content conditions (detailed text	spit into segments)
-	• <i>sn</i> : Sentence number (starts at 0)	
Outpu	ut : Generated audio sequence v	
1: v	← []	⊳ Initialize output au
2: $ au$	$\leftarrow [\tau_1, \tau_{1+P}, \tau_{1+2P};; \tau_{f-2}, \tau_{f-1}, \tau_f]$	Curved denoising timestep sched
3: Q	$P \leftarrow [\mathbf{z}_1, \mathbf{z}_{1+P}, \mathbf{z}_{1+2P};; \mathbf{z}_{f-2}, \mathbf{z}_{f-1}, \mathbf{z}_f]$	Initialize latent queue for diffus
4: [c	$[\mathbf{c}_{cont1}, \mathbf{c}_{cont2},, \mathbf{c}_{contK}] \leftarrow \mathbf{c}_{cont}$	Split content prompts into segme
5: S	$\leftarrow [\texttt{Len}(\mathbf{c}_{cont1}),\texttt{Len}(\mathbf{c}_{cont2}),,,\texttt{Len}(\mathbf{c}_{cont2}),,,,,,,,,,,,.$	$contK$)] \triangleright Store segment length
6: ϵ_{ϵ}	$\phi(\cdot) \leftarrow \boldsymbol{\epsilon}_{\boldsymbol{ heta}_{odd}}(\cdot)$	▷ Initialize with conditional guidance (o
7: fo	or <i>l</i> to <i>N</i> do	⊳ Iterate over N fra
8:	while $l < S[sn]$ do	▷ Generate frames for the current sente
9:	$\boldsymbol{Q} \leftarrow \boldsymbol{\epsilon}_{ heta}(\boldsymbol{Q}, \boldsymbol{\tau}, \mathbf{c}_{desc}, \mathbf{c}_{cont})$	> Apply noise prediction to latent sta
10:	$\mathbf{z}^l \leftarrow Q$, dequeue ()	▷ Pop out clean audio fra
11:	$v_{\text{append}}(\text{Dec}(\mathbf{z}_{l}))$	\triangleright Decode and append to out
12.	$r^{l+\text{Len}(\mathbf{Q})} \sim \Lambda(\mathbf{Q},\mathbf{I})$	S Generate new rendem n
12:	$\mathbf{z}_{ au_f} \sim \mathcal{N}(0,\mathbf{I})$	
13:	Q .enqueue $(\mathbf{z}_{\tau_f}, \cdot, \cdot)$	▷ Insert new noise into latent qu
14:	If $l = S[sn] - 1$ then	> Check if we reached the end of the sente
15:	$sn \leftarrow sn + 1$	▷ Move to the next sente
16:	if $\epsilon_{\theta}(\cdot) = \epsilon_{\theta_{odd}}(\cdot)$ then	▷ Switch conditional guida
17:	$oldsymbol{\epsilon}_{ heta}(\cdot) \leftarrow oldsymbol{\epsilon}_{ heta_{even}}(\cdot)$	▷ Switch to unconditional guida
18:	else	
19:	$oldsymbol{\epsilon}_{ heta}(\cdot) \leftarrow oldsymbol{\epsilon}_{ heta_{odd}}(\cdot)$	Switch back to conditional guida
20:	end if	
21:	end if	
22:	end while	
23: er	nd for	
24: re	eturn v	Return the final generated audio seque



We present additional audio samples that showcase the capabilities of InfiniteAudio, which builds upon AudioLDM. As illustrated in Fig. 10, InfiniteAudio can generate a variety of sounds exceeding 10 seconds in duration. Remarkably, the generated audio maintains consistency with the provided text prompts for up to 20 seconds. Furthermore, the system demonstrates the ability to produce coherent audio for durations of up to 60 seconds, as evidenced by CLAP scores.

918			I				
919	c _{cont}	C _{desc}	Methods	Spectrograms for generated Audio	CLAP	WER	
920	The closet door stuck a little as he pulled it open revealing a bare wooden floor						
921	covered on one side by a stack of old		VoiceLDM		0.1576	1.0	
922	For centuries after her death, Welshmen			This room			
923	when engaging in battle.	A person talking which later	网络新闻制制 网络新闻 化热气热 二乙酮 法主义法人				
924	The mixture will produce a chlorine solution of approximately five hundred	imitates a couple of meow sounds	AudioFIFO		0.2409	0.2386	
925	milligrams per liter. The boy wanted to believe that his			The closet door stuck a little as he pulled it open revealing a bare wooden fl	oor covered on a	one side by a	
926	friend had simply become separated from him by accident.			stack of old shoes boxes. For centuries after her death welshmen cried out revenge for gwenllian who	en engaging in l	battle.	
927	Sherman's mother-in-law, Margaret Odding, married secondly John Porter,			The mixture will produce a chlorine solution of approximately 500 milligrat The boy wanted to believe that his friend had simply become separated from	ns per liter. him by acciden	ıt.	
928	another signer of the Portsmouth Compact.			Sherman is mother in law margaret odding married secondly john porter an portsmouth compact.	other signer of	the	
929							
930	Nevertheless, Hussey's "Country Life" articles on contemporary houses are		VoiceLDM		0.1781	0.5195	
931	often overlooked. When she called her friend for help	A gun cocking then					
932	with her computer, she repeatedly emphasized, that she hadn't done	clanks on a hard		vevertnetess nussey is country life articles on contemporary houses are often overlooked. When she called her friend to help with her computer she repeatedly emphasized that she had not done			
933	anything. When I have read a million of these	surface followed by a man talking		anything. To goad a 1000000 of these scented tools as to invite as many as mainstream	ı americans as s	she could.	
934	sentences, I am going to apply as a professional narrator.	during an electronic laser	AudioFIFO		0 1865	0 1429	
935	The goal of the listed tools is to embed accessibility into various mainstream	effect as gunshots and explosions go			011000	011129	
936	technologies. These clubs organize inter-school	off in the distance		Hussie is country life articles on contemporary houses are often overlooked. For help with her computer she repeatedly emphasized that she had not don	e anything.		
937	workshops and symposia annually for students, teachers and parents.			When i have read a 1000000 of these sentences i am going to apply as a pro The goal of the listed tools is to embed accessibility into mainstream technol	fessional narrat logies.	tor.	
938				These clubs organize inter school workshops and symposia annually for stu	dents.		
939	The Queensland Greens support the						
940	reintroduction of an upper house elected by proportional representation.		VoiceLDM		0.2662	0.4857	
941	Three new low powered relay stations were built, allowing easier access to			Queensland greens support the reintroduction of an upper house elected by	proportional re	presentation.	
942	Anglia transmissions. The local economy of Ninomiya is based	A woman speaks		3 new low powered rail stations were built allowing easier access to anglia to The local economy and agriculture and commercial fishing the next month	ransmissions. will be a great y	vear.	
943	on primarily on agriculture and commercial fishing.	followed by a girl speaking faintly					
944	Hash trees allow efficient and secure verification of the contents of large data		AudioFIFO		0.2404	0.2143	
945	structures. The next month, "Ticonderoga"						
946	recovered her second set of space			The queensland greens support the reintroduction of an upper house by pro Relay stations were built allowing easier access to anglia transmissions.	portional zentai	tion greenloo.	
947				The local economy of nomea is based primarily on agriculture and commer- Hash trees allow and securefication of the contents of large data structures.	cial fishing.		
948				Teconderoga recovered her 2nd set of space voyagers near american samoa.			
949							

Figure 11: Generated audio samples based on VoiceLDM.

B.2 VOICELDM

In our exploration of VoiceLDM, we generate audio based on two distinct text prompts: c_{cont} and c_{desc} . As shown in Fig. 11, our approach produces significantly improved audio quality when provided with longer content prompts. However, it's worth noting that VoiceLDM is limited to generating audio of no more than 10 seconds, often resulting in truncated or incomplete sentences. For clarity, we have included the transcriptions of the sentences below each spectrogram, with any cut-off portions highlighted in red. In contrast, our method ensures the generation of entire sentences without omitting sections. The superiority of our approach is further supported by CLAP and WER scores, which validate the enhanced intelligibility and coherence of the generated audio.

972 C EXPERIMENT DETAILS 973

974 C.1 DATASET

Existing TTA models: AudioLDM is trained on a diverse combination of datasets, including AudioSet (Gemmeke et al., 2017), the largest audio dataset with over 5,000 hours of data, as well as AudioCaps Kim et al. (2019), Freesound (FS), and the BBC Sound Effect (SFX) library, covering a wide range of sounds. Similarly, VoiceLDM is trained on AudioSet for TTA, the English subset of the CommonVoice 13.0 corpus and VoxCeleb1 for speech generation, and the DEMAND dataset for non-speech segments. AudioLDM is evaluated on both the AudioSet and AudioCaps datasets, while VoiceLDM is tested exclusively on the AudioCaps dataset.

InfiniteAudio: We utilize the Audiocaps test set for text-to-audio generation, which comprises audio files paired with corresponding caption texts. Each audio file is accompanied by several text captions, from which we randomly select 860 audio-text pairs for evaluation.

For text-to-speech (TTS) evaluation, we employ the CommonVoice 13.0K test set. Unlike traditional TTS evaluations, our focus is on generating longer speech segments. Therefore, we specifically target sentences exceeding 90 characters in length. For each text input, we utilize more than four selected sentences, resulting in a total of 60 text pairs and approximately 300 sentences for testing.

Both audio and speech generation evaluations leverage the aforementioned datasets. We include
captions categorized as "speech" from the Audiocaps test set as prompts for audio descriptions and
randomly select 60 long speech pairs from the CommonVoice 13.0K test set for content prompts.

995 C.2 CONFIGURATION

We conducted experiments using InfiniteAudio alongside existing text-to-audio generation models,
AudioLDM and VoiceLDM. Both models are based on Latent Diffusion Models (LDM) utilizing a
U-Net architecture. We increased the number of inference steps to optimize performance, deviating
from the default settings of the original models, and employed DDIM sampling.

For AudioLDM, we set the inference steps to 300, aligning with the original model but omitting the initial and middle regions, except for multiples of 4, while retaining the final steps, which introduce slight noise into the spectrograms. In contrast, for VoiceLDM, we increased the inference steps from the original 50 to 200, skipping the middle and final regions except for multiples of 5, while including the initial timesteps to introduce Gaussian-like noise.

In text-to-long speech generation (TTLS), we segment long content prompts at the sentence level.
 For QKV sharing, we apply the sharing mechanism every 200 audio frames. However, at the start of each new sentence, the sharing process is reset, beginning again with the first 200 frames.

1009 1010 C.3 HUMAN EVALUATION

Subjective evaluation plays a vital role in the text-to-audio generation domain. For our assessment, we randomly selected 30 generated audio samples, which were rated by 20 domain experts on a scale from 1 to 5. The evaluation criteria focused on overall audio quality and the relevance of the generated audio to the corresponding descriptive text.

Our model aims to generate longer audio segments while preserving the performance of pretrained text-to-audio models. Despite the inherent challenges in evaluating these samples, our results indicate comparable performance to existing models, with ground truth scores averaging around 4.

1019

- 1020
- 1021
- 1022
- 1023 1024
- 1024

CLAP↑ WER↓ WER↓ **CLAP**[↑] 0.090 0.45 0.09 0.40 0.44 0.39 0.085 0.085 0.43 0.38 0.42 0.080 0.08 0.41 0.37 0.40 0.075 0.075 0.3 0.39 0.070 0.070 0.38 Cosine similarity1 Embedding Distance↓ 0.88 Inner Middle 1.06 200 0.86 1.04 0.84 1.02 0.82 1.00 0.80 S2 S4 C1 C^{2} C4 S³ C3 (a) Impacts on long QKV sharing (b) Analysis on Unet

1041 Figure 12: Analysis of QKV sharing: (a) Impact of extended QKV sharing and (b) Comparison of 1042 U-Net decoder modules.

1044 D ADDITIONAL EXPERIMENTS 1045

1046 D.1 IMPACT OF EXTENDED QKV SHARING 1047

1048 As illustrated in Fig. 12, sharing Query, Key, and Value (QKV) representations over 200 audio 1049 frames proves effective based on CLAP, Word Error Rate (WER), and speaker embedding distance 1050 scores. While sharing QKV for a single frame yields smaller speaker embedding distances, it can 1051 adversely affect CLAP and WER metrics. In contrast, sharing QKV across 200 frames reduces 1052 speaker embedding distances compared to the case with no sharing, while simultaneously enhancing both CLAP and WER scores. 1053

1054 1055 D.2 COMPARATIVE ANALYSIS OF U-NET DECODER MODULES

1056 The U-Net architecture comprises both encoder 1057 and decoder components, typically organized 1058 into multiple downsampling and upsampling 1059 modules. In this study, we categorize these 1060 modules into inner, middle, and outer groups. 1061 Previous works, such as (Ceylan et al., 2023; 1062 Kwon et al., 2024; Tumanyan et al., 2023), have 1063 explored sharing mechanisms across various U-1064 Net decoder modules. We empirically determined the most effective module for sharing. Our findings indicate that the outer module en-1066 hances speaker embedding similarity and im-1067 proves CLAP scores. 1068



Figure 13: Architecture of the U-Net decoder.

1069

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039 1040

- 1075
- 1076
- 1077
- 1078
- 1079



1134 D.4 EXISTING LONG GENERATION METHODS BASED ON DIFFUSION MODELS. 1135

1136 We briefly review key papers on long-form generation using pretrained diffusion models, focusing specifically on methods that do not involve additional training in other generative domains. 1137

1138

D.4.1 MULTIDIFFUSION FOR TTI 1139

1140 MultiDiffusion⁶ is a robust framework for text-to-image generation that leverages a pre-trained dif-1141 fusion model without requiring additional training or fine-tuning. At its core, MultiDiffusion merges 1142 multiple independent diffusion processes through an optimization algorithm, reconciling them into a 1143 coherent and high-quality image. This approach enables user-controllable image generation, making 1144 it highly adaptable for a wide range of tasks.

1145 The primary innovation of MultiDiffusion lies in its ability to simultaneously process multiple image 1146 regions, adhering to user-defined constraints such as text prompts, aspect ratios, and spatial layout 1147 signals (e.g., segmentation masks or bounding boxes). However, a key limitation of the framework 1148 is its inability to incorporate temporal information, restricting its utility to image generation. This 1149 makes it unsuitable for video generation, where maintaining temporal coherence across frames is 1150 essential.

1151 Despite this limitation, the underlying optimization process ensures that all image regions conform 1152 closely to the reference diffusion model, preserving both high image quality and visual consistency. 1153 While not applicable to temporal tasks, MultiDiffusion remains a flexible and efficient solution for 1154 generating complex images that meet diverse spatial constraints. 1155

1156 D.4.2 FREENOISE FOR TTV 1157

FreeNoise⁷ is a framework designed to extend the capabilities of text-to-video diffusion models 1158 for generating longer, temporally coherent videos. Traditional text-to-video models are typically 1159 trained on a limited number of frames, restricting their ability to generate high-fidelity long videos 1160 during inference. FreeNoise addresses this limitation by introducing a **tuning-free** paradigm that 1161 dynamically reschedules noise over time to maintain consistency across frames. 1162

1163 **Key Innovations** 1164

- 1165 • Noise Rescheduling: Unlike traditional methods that initialize noise uniformly for all frames, FreeNoise dynamically adjusts the noise distribution during video generation. This method captures long-range temporal correlations, ensuring that visual consistency is pre-1168 served across extended video sequences.
 - Temporal Attention Mechanism: FreeNoise incorporates a window-based temporal attention mechanism, which helps maintain coherence over longer time frames. By focusing attention over localized windows, the model can efficiently capture and retain relevant temporal dependencies.
 - Motion Injection for Multi-Prompt Videos: The framework supports multi-prompt video generation by enabling dynamic changes in video content based on evolving text prompts. This allows FreeNoise to generate videos where different segments adhere to different prompts, accommodating more complex narrative transitions over time.
- 1176 1177

1166

1167

1169

1170

1171

1172

1173

1174

1175

1178 **Limitations** Despite its strengths, FreeNoise has certain limitations. Since the framework does 1179 not involve any fine-tuning of the pre-trained models, it might not be optimally adapted to domainspecific datasets, potentially leading to suboptimal performance in specialized contexts. Moreover, 1180 while it addresses temporal consistency, the model's reliance on pre-trained diffusion models can 1181 limit its ability to handle diverse or complex motion dynamics inherent in specific generative tasks. 1182

- 1183
- 1184
- 1185

¹¹⁸⁶

⁶https://github.com/omerbt/MultiDiffusion 1187

⁷https://github.com/AILab-CVC/FreeNoise