

ViT-TTS: Visual Text-to-Speech with Scalable Diffusion Transformer

Anonymous ACL submission

Abstract

Text-to-speech(TTS) has undergone remarkable improvements in performance, particularly with the advent of Denoising Diffusion Probabilistic Models (DDPMs). However, the perceived quality of audio depends not solely on its content, pitch, rhythm, and energy, but also on the physical environment. In this work, we propose ViT-TTS, the first visual TTS model with scalable diffusion transformers. ViT-TTS complement the phoneme sequence with the visual information to generate high-perceived audio, opening up new avenues for practical applications of AR and VR to allow a more immersive and realistic audio experience. To mitigate the data scarcity in learning visual acoustic information, we 1) introduce a self-supervised learning framework to enhance both the visual-text encoder and denoiser decoder; 2) leverage the diffusion transformer scalable in terms of parameters and capacity to learn visual scene information. Experimental results demonstrate that ViT-TTS achieves new state-of-the-art results, outperforming cascaded systems and other baselines regardless of the visibility of the scene. With low-resource data (1h, 2h, 5h), ViT-TTS achieves comparative results with rich-resource baselines.¹

1 Introduction

Text-to-speech(TTS) (Ren et al., 2019; Huang et al., 2022b; Huang et al.) has undergone remarkable improvements in performance, particularly with the advent of Denoising Diffusion Probabilistic Models (DDPMs). However, the perceived quality of audio depends not solely on its content, pitch, rhythm, and energy, but also on the physical environment. For instance, a room with hard surfaces like concrete or glass reflects sound waves, whereas a room with soft surfaces such as carpets or curtains absorbs them. This variance can drastically impact the clarity and quality of the sound we hear.

¹Audio samples are available at <https://ViT-TTS.github.io/>.

To ensure an authentic and captivating experience, it is imperative to accurately model the acoustics of a room, particularly in virtual reality (VR) and augmented reality (AR) applications. Recent years have seen a surge in significant research (Li et al., 2022; Radford et al., 2021; Li et al., 2023; Huang et al., 2023b) addressing the language-visual modeling problem. For instance, Li et al. (2022) have proposed a unified video-language pre-training framework for learning robust representation, while Radford et al. (2021) have focused on large-scale image-text pairs pre-training via contrastive learning. Visual TTS open-ups numerous practical applications, including dubbing archival films, providing a more immersive and realistic experience in virtual and augmented reality, or adding appropriate sound effects to games.

Despite the benefits of language-visual approaches, training visual TTS models typically requires a large amount of training data, while there are very few resources providing parallel text-visual-audio data due to the heavy workload. Besides, creating a sound experience that matches the visual content remains challenging when developing AR/VR applications, as it is still unclear how various regions of the image contribute to reverberation and how to incorporate the visual modality as auxiliary information in TTS.

In this work, we formulate the task of visual TTS to generate audio with reverberation effects in target scenarios given a text and environmental image, introducing ViT-TTS to address the issues of data scarcity and room acoustic modeling. To enhance visual-acoustic matching, we 1) propose the visual-text fusion to integrate visual and textual information, which provides fine-grained language-visual reasoning by attending to regions of the image; 2) leverage transformer architecture to promote the scalability of the diffusion model. Regarding the data shortage challenge, we pre-train the encoder and decoder in a self-supervised manner, showing

082 that large-scale pre-training reduces data require- 130
083 ments for training visual TTS models. 131

084 Experiments results demonstrate that ViT-TTS 132
085 generates speech samples with accurate reverbera- 133
086 tion effects in target scenarios, achieving new state- 134
087 of-the-art results in terms of perceptual quality. In 135
088 addition, we investigate the scalability of ViT-TTS 136
089 and its performance under low-resource conditions 137
090 (1h/2h/5h). The main contributions of this work 138
091 are summarized as follows: 139

- 092 • We propose the first visual Text-to-Speech 140
093 model ViT-TTS with vision-text fusion, which 141
094 enables the generation of high-perceived au- 142
095 dio that matches the physical environment. 143
- 096 • We show that large-scale pre-training allevi- 144
097 ates the data scarcity in training visual TTS 145
098 models. 146
- 099 • We introduce the diffusion transformer scal- 147
100 able in terms of parameters and capacity to 148
101 learn visual scene information. 149
- 102 • Experimental results on subjective and ob- 150
103 jective evaluation demonstrate the state-of- 151
104 the-art results in terms of perceptual qual- 152
105 ity. With low-resource data (1h, 2h, 5h), ViT- 153
106 TTS achieves comparative results with rich- 154
107 resource baselines. 155

108 2 Related Work 156

109 2.1 Text-To-Speech 157

110 Text-to-Speech(TTS) tasks are divided into two cat- 158
111 egories: (1) generating a mel-spectrogram from 159
112 text or phoneme sequence first (Wang et al., 2017; 160
113 Ren et al., 2019), and then converting the gener- 161
114 ated spectrum into a waveform via vocoder (Kong 162
115 et al., 2020; Lee et al., 2022; Huang et al., 2022b, 163
116 2021, 2022a); (2) generating audio directly from 164
117 text (Donahue et al., 2020; Kim et al., 2021). 165
118 The earlier TTS (Li et al., 2019; Wang et al., 166
119 2017) models adopt an autoregressive manner, 167
120 which suffers from the problem of slow infer- 168
121 ence speed. As a solution, non-autoregressive 169
122 models have been proposed to enable fast infer- 170
123 ence by generating mel-spectrograms in parallel. 171
124 More recently, Grad-TTS (Popov et al., 2021), 172
125 DiffSpeech (MoonInTheRiver, 2021), and ProD- 173
126 iff (Huang et al., 2022c) have employed diffu- 174
127 sion generative models to generate high-quality 175
128 audio, but they all rely on the convolutional archi- 176
129 tecture such as WaveNet (Oord et al., 2016) and 177

U-Net (Ronneberger et al., 2015) as the backbone. 130
In contrast, some studies (Peebles and Xie, 2023; 131
Bao et al., 2023) in image generation tasks have 132
explored transformers (Vaswani et al., 2017) as an 133
alternative to convolutional architectures, achieving 134
competitive results with U-Net. In this paper, we 135
present the first transformer-based diffusion model 136
as an alternative of convolutional architecture. By 137
harnessing the scalable properties of transformers, 138
we enhance the model capacity to more effectively 139
capture visual scene information and promote the 140
model performance. 141

142 2.2 Self-supervised Pre-training 142

143 There are two main criteria for optimizing speech 144
145 pre-training: contrastive loss (Oord et al., 2018; 146
147 Chung and Glass, 2020; Baevski et al., 2020) and 148
149 masked prediction loss (Devlin et al., 2018). Con- 149
150 trastive loss is used to distinguish between positive 150
151 and negative samples with respect to a reference 151
152 sample, while masked prediction loss is originally 152
153 proposed for natural language processing (Devlin 153
154 et al., 2018; Lewis et al., 2019) and later applied to 154
155 speech processing (Baevski et al., 2020; Hsu et al., 155
156 2021). Some recent work (Chung et al., 2021) has 156
157 combined the two approaches, achieving good per- 157
158 formance for downstream automatic speech recog- 158
159 nition (ASR) tasks. In this work, we leverage the 159
160 success of self-supervised to enhance both the en- 160
161 coder and decoder to alleviate the data scarcity 161
162 issue. 162

163 2.3 Acoustic Matching 163

164 The primary objective of acoustic matching is to 164
165 convert audio from a source environment into au- 165
166 dio that resembles the target environment. In the 166
167 field of blind estimation (Mack et al., 2020; Xiong 167
168 et al., 2018; Murgai et al., 2017; Mezghani and 168
169 Swindlehurst, 2018), acoustic matching is applied 169
170 to generate a simple room impulse response (RIR) 170
171 that can be used to synthesize the corresponding 171
172 target audio using two critical acoustic metrics - 172
173 the direct-to-reverberant ratio (DRR) (Zahorik, 173
174 2002) and the reverberation time 60 (RT60) (Rat- 174
175 nam et al., 2003). The DRR is used to describe 175
176 the energy ratio between the direct-to-reverberant 176
177 sound and the reflected sound, while the RT60 is 177
178 used to measure the time taken for the sound to 178
179 decay by 60 dB. The music production community 179
also implements acoustic matching to modify the 179
reverberation, thus simulating the reverberation of 179
the target space or processing algorithm (Koo et al., 179

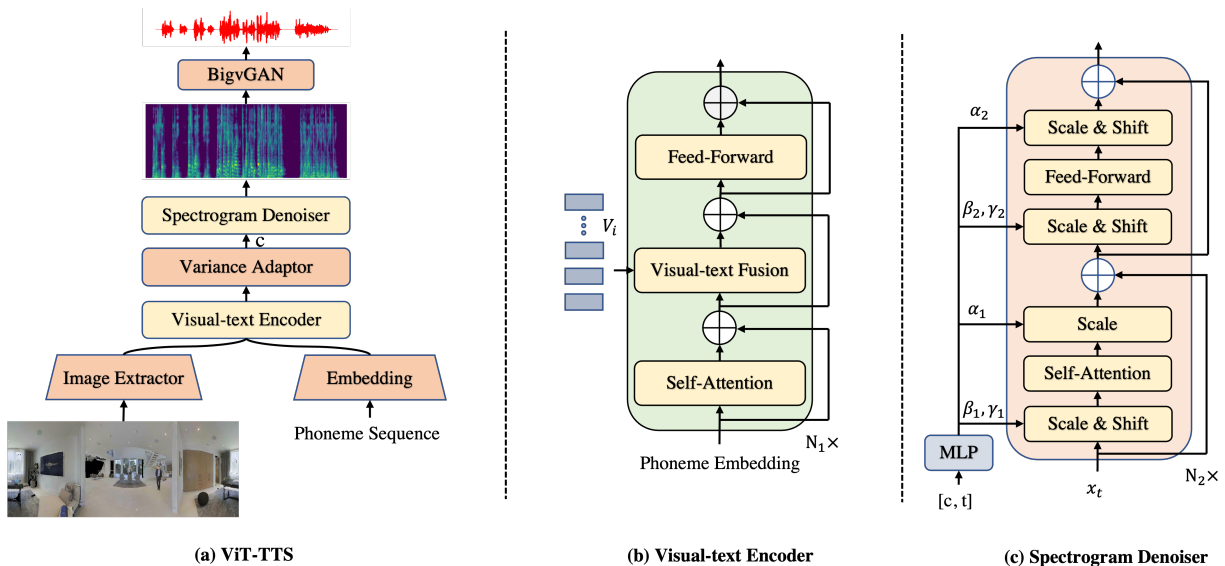


Figure 1: The overall architecture for ViT-TTS. In subfigure (b), V_i denotes the visual sequence and N_1 denotes the layers of Encoder. In subfigure (c), N_2 is the number of transformer layers. α and β are the dimension-wise scale parameters, while γ is the dimension-wise shift parameters. c is the variance adaptor’s output and t is the diffusion step.

2021; Sarroff and Michaels, 2020). Recently, there is research on visual acoustic matching (Chen et al., 2022), which involves generating audio recorded in the target environment based on the input source audio clip and an image of the target environment. However, our proposed visual TTS is distinct from those mentioned above as we aim to generate audio that captures the room acoustics in the target environment based on the written text and the target environment image.

3 Method

3.1 Overview

The overall architecture has been presented as Figure 1. To alleviate the issue of data scarcity, we leverage unlabeled data to pre-train the visual-text encoder and denoiser decoder with scalable transformers in a self-supervised manner. To capture the visual scene information, we employ the visual-text fusion module to reason about how different image patches contribute to texts. BigvGAN (Lee et al., 2022) converts the mel-spectrograms into audio that matches the target scene as a neural vocoder.

3.2 Enhanced visual-text Encoder

Self-supervised Pre-training The advent of the masked language model (Devlin et al., 2018; Clark et al., 2020) has marked a significant milestone in the field of natural language processing. To alleviate the data scarcity issue (Huang et al., 2022d; MoonInTheRiver, 2021) and learn robust contextual encoder, we are encouraged to adopt the mask-

ing strategy like BERT in the pre-training stage. Specifically, we randomly mask the 15% of each phoneme sequence and predict those masked tokens rather than reconstructing the entire input. The masked phoneme sequence is then input into the text encoder to obtain hidden states. The final hidden states are fed into a linear projection layer over the vocabulary to obtain the predicted tokens. Finally, we calculate the cross entropy loss between the predicted tokens and target tokens.

The masked token during the pre-training phase will not be used in the fine-tuning phase. To mitigate this mismatch between the pre-training and fine-tuning, we randomly choose the phonemes to be masked: 1) 80% probability to add masks; 2) 10% probability to keep phoneme unchanged, and 3) 10% probability to replace with a random token in the dictionary.

Visual-Text Fusion In the fine-tuning stage, we integrate the visual modal and module into the encoder to integrate visual and textual information. Before feeding into the visual-text encoder, we first extract image features of panoramic images through ResNet18 (Oord et al., 2018) and obtain phoneme embedding. Both the image features and phoneme embedding are fed into one of the variants of the transformer to get the hidden sequences. Specifically, we first pass the phoneme through relative self-attention, which is defined as follows.

$$\alpha(i, j) = \text{Softmax}\left(\frac{(Q_i W^Q)(K_j W^K + R_{ij}^K)^T}{\sqrt{d_k}}\right) \quad (1)$$

where n is the length of phoneme embedding, R_{ij}^K and R_{ij}^V are the relative position embedding of key and value, and Q, K, V are all the phoneme embedding. We use relative self-attention to model how much phoneme p_i attends to phoneme p_j . After that, we choose to use cross-attention instead of a simplistic concatenation approach as we can reason about how different image patches contribute to the text after feature extraction. The equation is defined as follows:

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

where Q is the phoneme embedding. K and V are the visual features. Finally, the feed-forward layer is applied to output the hidden sequence.

3.3 Enhanced Diffusion Transformer

Scalable Transformer As a rapidly growing category of generative models, denoising diffusion models (DDPMs) have demonstrated their exceptional ability to deliver top-notch results in both image (Zhang and Agrawala, 2023; Ho and Salimans, 2022) and audio synthesis (Huang et al., 2022c, 2023a; Lam et al., 2021). However, the most dominant diffusion TTS models adopt a convolutional architecture like WaveNet or U-Net as the de-factor choice of backbone, where it lacks the scalable ability to model additional visual information. This prevents the model’s incorporation of visual information as they lack scalability. Recent research (Peebles and Xie, 2023; Bao et al., 2023) in the image synthesis field has revealed that the inductive bias of convolutional structures is not a critical determinant of DDPMs’ performance. Instead, transformers have emerged as a viable alternative.

For this reason, we propose a diffusion transformer that leverages the scalability of transformers to expand model capacity and incorporate room acoustic information. Moreover, we leverage the adaptive normalization layers in GANs and initialize the full transformer block as the identity function to further improve the transformer performance.

Unconditional Pre-training In this part, we investigate self-supervised learning from orders of magnitude mel-spectrograms data to alleviate data scarcity. Specifically, assuming the target mel-spectrogram is x_0 , we first random select 0.065% of x_0 as starting indices and apply a mask that spans 10 steps following the Wav2vec2.0 (Baevski

et al., 2020). Then, we obtain x_t through a diffusion process, which is defined by a fixed Markov chain from data x_0 to the latent variable x_t .

$$q(x_1, \dots, x_T | x_0) = \prod_{t=1}^T q(x_t | x_{t-1}), \quad (3)$$

At each diffusion step $t \in [1, T]$, a tiny Gaussian noise is added to x_{t-1} to obtain x_t , according to a small positive constant β_t :

$$q(x_t | x_{t-1}) := \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}) \quad (4)$$

x_t obtained from the diffusion process is passed through the transformer to predict Gaussian noise ϵ_θ . Loss is defined as mean squared error in the ϵ space, and efficient training is optimizing a random term of t with stochastic gradient descent:

$$\mathcal{L}_\theta^{\text{Grad}} = \left\| \epsilon_\theta \left(\alpha_t x_0 + \sqrt{1 - \alpha_t^2} \epsilon \right) - \epsilon \right\|_2^2, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (5)$$

To this end, ViT-TTS takes advantage of the reconstruction loss to predict the self-supervised representations which largely alleviates the challenges of data scarcity. Detailed formulation of DDPM has been attached in Appendix C.

Controllable Fine-tuning During the fine-tuning stage, we will face the following challenges: (1) there is a data scarcity issue with the available panoramic images and target environmental audio for training; (2) a fast training method is equally crucial for optimizing the diffusion model, as it can save a significant amount of time and storage space. To address these challenges, we draw inspiration from Zhang and Agrawala (2023) and implement a swift fine-tuning technique. Specifically, we create two copies of the pre-trained diffusion model weights, namely a "trainable copy" and a "locked copy," to learn the input conditions. We fix all parameters of the pre-trained transformer, designated as Θ , and duplicate them into a trainable parameter Θ_t . We train these trainable parameters and connect them with the "locked copy" via zero convolution layers. These convolution layers are unique as they have a kernel size of one by one and weights and biases set to zero, progressively growing from zeros to optimized parameters in a learned fashion.

3.4 Architecture

As illustrated in Figure 1, our model comprises a visual-text encoder, variance adaptor, and spectrogram denoiser. The visual-text encoder converts

333 phoneme embeddings and visual features into hid- 382
334 den sequences, while the variance adaptor predicts 383
335 the duration of each hidden sequence to regulate 384
336 the length of the hidden sequences to match that 385
337 of speech frames. Furthermore, different variances 386
338 like pitch and speaker embedding are incorporated 387
339 with hidden sequences. Finally, the spectrogram 388
340 denoiser iteratively refines the length-regulated hid- 389
341 den states into mel-spectrograms. We put more 390
342 details in Appendix B. 391

343 **Visual-Text Encoder** The visual-text encoder 392
344 consists of relative position transformer blocks 393
345 based on the transformer architecture. Specifically, 394
346 it convolves a pre-net for phoneme embedding, a vi- 395
347 sual feature extractor for image, and a transformer 396
348 encoder which includes multi-head self-attention, 397
349 multi-head cross-attention, and feed-forward layer.

350 **Variance Adaptor** In variance adaptor, the du- 398
351 ration and pitch predictors share a similar model 399
352 structure consisting of a 2-layer 1D-convolutional 400
353 network with ReLU activation, each followed by 401
354 the layer normalization and the dropout layer, and 402
355 an extra linear layer to project the hidden states 403
356 into the output sequence. 404

357 **Spectrogram Denoiser** Spectrogram denoiser 405
358 takes in x_t as input to predict ϵ added in diffusion 406
359 process conditioned on the step embedding E_t and 407
360 encoder output. We adopt a variant of the trans- 408
361 former as our backbone and make some improve- 409
362 ments upon the standard transformer motivated 410
363 by Peebles and Xie (2023), mainly includes:(1) 411
364 we explore replacing standard layer norm layers 412
365 in transformer blocks with adaptive layer norm 413
366 (adaLN) to regress scale and shift parameters from 414
367 the sum of the embedding vector of t and hidden 415
368 sequence. (2) Inspired by ResNets (Oord et al., 416
369 2018), we initialize the transformer block as the 417
370 identity function and initialize the MLP to output 418
371 the zero-vector. 419

372 3.5 Pre-training, Fine-tuning, and Inference 420 373 Procedures 421

374 **Pre-training** The pre-training has two stages: 1) 422
375 encoder stage: pre-train the visual-text encoder 423
376 via masked LM loss \mathcal{L}_{CE} (ie. cross-entropy loss) 424
377 to predict the masked tokens. 2) decoder stage: 425
378 the masked x_0 is puted into denoiser to predict 426
379 Gaussian noise ϵ_θ . Then, the MSE loss is applied 427
380 to the predicted Gaussian noise and target Gaussian 428
381 noise.

Fine-tuning We begin by loading model weights 382
from the pre-trained visual-text encoder and uncon- 383
ditional diffusion decoder, after which we finetune 384
both of them until the model converges. The fi- 385
nal loss term consists of the following parts: (1) 386
sample reconstruction loss \mathcal{L}_θ : MSE between the 387
predicted Gaussian noise and target Gaussian noise. 388
(2) variance reconstruction loss $\mathcal{L}_{dur}, \mathcal{L}_p$: MSE be- 389
tween the predicted and the target phoneme-level 390
duration, pitch. 391

Inference In inference, DDPM iteratively runs 392
the reverse process to obtain the data sample x_0 , 393
and then we use a pre-trained BigvGAN-16khz- 394
80band as the vocoder to transform the generated 395
mel-spectrograms into waveforms. 396

397 4 Experiment 397

398 4.1 Experimental Setup 398

Dataset We use the SoundSpaces-Speech 399
dataset (Chen et al., 2023), which is constructed 400
on the SoundSpaces platform based on real-world 401
3D scans to obtain environmental audio. The 402
dataset includes 28,853/1,441/1,489 samples for 403
training/validation/testing, each consisting of clean 404
text, reverberant audio, and panoramic camera 405
angle images. Following (Chen et al., 2022), we 406
remove out-of-view samples and divide the test set 407
into test-unseen and test-seen, where the unseen 408
set injects room acoustics depicted in novel images 409
while the seen set only contains the scenes we 410
have seen in the training stage. We convert the 411
text sequence into the phoneme sequence with an 412
open-source grapheme-to-phoneme conversion 413
tool (Sun et al., 2019) ². 414

Following the common practice (Ren et al., 2019; 415
MoonInTheRiver, 2021), we conduct preprocess- 416
ing on the speech and text data: 1) extract the spec- 417
trogram with the FFT size of 1024, hop size of 256, 418
and window size of 1024 samples; 2) convert it to 419
a mel-spectrogram with 80 frequency bins; and 3) 420
extract F0 (fundamental frequency) from the raw 421
waveform using Parselmouth tool ³. 422

Model Configurations We extract the mel- 423
spectrogram from the raw waveform and set the 424
hop size and frame size to 256 and 1024 in respect 425
of the sample rate 16kHz. The size of the phoneme 426
vocabulary is 73. The dimension of phoneme em- 427
beddings and the hidden size of the visual-text 428

²<https://github.com/Kyubyong/g2p>

³<https://github.com/YannickJadoul/Parselmouth>

Method	Test-Seen			Test-Unseen			Params
	MOS(\uparrow)	RTE (\downarrow)	MCD (\downarrow)	MOS(\uparrow)	RTE (\downarrow)	MCD (\downarrow)	
GT	4.34 \pm 0.07	/	/	4.24 \pm 0.07	/	/	/
GT (voc.)	4.18 \pm 0.05	0.006	1.46	4.19 \pm 0.07	0.008	1.50	/
WaveNet	3.85 \pm 0.09	0.091	4.61	3.78 \pm 0.12	0.110	4.69	42.3M
Transformer-S	3.92 \pm 0.07	0.068	4.57	3.80 \pm 0.06	0.077	4.68	32.38M
Transformer-B	3.98 \pm 0.06	0.061	4.53	3.90 \pm 0.07	0.066	4.62	41.36M
Transformer-L	4.02 \pm 0.08	0.056	4.37	3.95 \pm 0.07	0.061	4.50	56.96M
Transformer-XL	4.05\pm0.07	0.047	4.35	4.00\pm0.05	0.053	4.39	115.12M

Table 1: Comparison between the diffusion WaveNet and diffusion transformers sweeping over model config(S, B, L, XL). All models remove the pre-training stage and other conditions not related to backbone in training and inference remain the same.

transformer block are both 256. We use the pre-trained ResNet18 as an image feature extractor. As for the pitch encoder, the size of the lookup table and encoded pitch embedding are set to 300 and 256. In the denoiser, the number of transformer-B layers is 5 with the hidden size 384 and head 12. We initialize each transformer block as the identity function and set T to 100 and β to constants increasing linearly from $\beta_1 = 10^{-4}$ to $\beta_T = 0.06$. We have attached more detailed information on the model configuration in Appendix B

Pre-training, Fine-tuning, and Inference During the pre-training stage, we pre-train the encoder for 120k steps and the decoder for 160k until convergence. The diffusion probabilistic models have been trained using 1 NVIDIA A100 GPU with a batch size of 48 sentences. In the inference stage, we uniformly use a pre-trained BigvGAN-16khz-80band (Lee et al., 2022) as a vocoder to transform the generated mel-spectrograms into waveforms.

4.2 Scalable Diffusion Transformer

We compare and examine diffusion transformer sweeping over model config(S, B, L, XL), and conduct evaluations in terms of audio quality and parameters. Appendix A gives the details of the model configs. The results have been shown in Table 1. We have some observations from the results: (1) Increasing the depth and number of layers in the transformer can significantly enhance the performance of the diffusion model, resulting in an improvement in both objective metrics and subjective metrics, which demonstrates that expanding the model size enables finer-grained room acoustic modeling. (2) Our proposed diffusion transformer outperforms WaveNet backbone under similar parameters across both test-unseen and test-seen sets,

significantly in the rt60 metric. We attribute this to the fact that instead of directly concatenating the condition input like WaveNet, we replace standard layer norm layers in transformer blocks with adaptive layer norm to regress dimension-wise scale and shift parameters from the sum of the embedding vectors of diffusion step and encoder output, which can better incorporate the conditional information, as proven in GANs (Brock et al., 2018; Karras et al., 2019).

4.3 Model Performances

In this study, we conduct a comprehensive comparison of the generated audio quality with other systems, including 1) GT, the ground-truth audio; 2) GT(voc.), where we first convert the ground-truth audio into mel-spectrograms and then convert them to audio using BigvGAN; 3) DiffSpeech (MoonInTheRiver, 2021), one of the most popular DDPM based on WaveNet; 4) ProDiff (Huang et al., 2022c), a recent generator-based diffusion model proposed to reduce the sampling time; 5) Visual-DiffSpeech, incorporate visual-text fusion module into DiffSpeech; 6) Cascaded, the system composed of DiffSpeech and Visual Acoustic Matching(VAM) (Chen et al., 2022). The results, compiled and presented in Table 2, provide valuable insights into the effectiveness of our approach:

(1) As expected, the results in the test-unseen set do poorer than the test-seen part because there are invisible scenarios among the test-unseen set. However, our proposed model has achieved the best performance compared to baseline systems in both sets, indicating that our model generates the best-perceived audio that matches the target environment from written text. (2) Our model surpassed TTS diffusion models(i.e.DiffSpeech and ProDiff) across all metric scores, especially in terms of RTE

Method	Test-Seen			Test-Unseen		
	MOS (\uparrow)	RTE (\downarrow)	MCD (\downarrow)	MOS (\uparrow)	RTE (\downarrow)	MCD (\downarrow)
GT	4.34 \pm 0.07	/	/	4.24 \pm 0.07	/	/
GT(voc.)	4.18 \pm 0.05	0.006	1.46	4.19 \pm 0.07	0.008	1.50
DiffSpeech	3.79 \pm 0.08	0.104	4.65	3.67 \pm 0.05	0.120	4.71
ProDiff	3.76 \pm 0.13	0.121	4.67	3.65 \pm 0.06	0.137	4.72
Visual-DiffSpeech	3.85 \pm 0.09	0.091	4.61	3.78 \pm 0.12	0.110	4.69
Cascaded	3.61 \pm 0.08	0.071	5.13	3.59 \pm 0.08	0.082	5.25
ViT-TTS	3.95\pm0.06	0.066	4.52	3.86\pm0.05	0.076	4.59

Table 2: Comparison with baselines on the SoundSpaces-Speech for Seen and Unseen scenarios. The diffusion step of all diffusion models is set to 100. We use the pre-trained model provided by VAM for the evaluation of cascaded.

values. This suggests that conventional diffusion models in TTS do poorly in modeling room acoustic information, as they mainly focus on audio content, pitch, energy, etc. Our proposed visual-text fusion module addresses this challenge by injecting visual properties into the model, resulting in a more accurate prediction of the correct acoustics from images and high-perceived audio synthesis. (3) The results of comparison with Visual-DiffSpeech highlight the advantages of our choice of transformer and self-supervised pre-training. Although Visual-DiffSpeech adds the visual-text module, the choice of WaveNet and the lack of a self-supervised pre-training strategy make it perform worse in predicting the correct acoustics from images and synthesizing high-perceived audio. (4) The cascaded system composed of DiffSpeech and Visual Acoustic Matching model visual properties is better than other baselines. However, compared to our proposed model, it performed worse in both test-unseen and test-seen environments. This suggests that our direct visual text-to-speech system eliminates the influence of error propagation caused by the cascaded manner, resulting in high-perceived audio. In conclusion, our comprehensive evaluation results demonstrate the effectiveness of our proposed model in generating high-quality audio that matches the target environment.

4.4 Low Resource Evaluation

Training visual text-to-speech models typically requires a large amount of parallel target environment image and audio training data, while there may be very few resources due to the heavy workload. In this section, we prepare low-resource audio-visual data (1h/2h/5h) and leverage large-scale text-only and audio-only data to boost the performance of the visual TTS system, to investigate the effectiveness

of our self-supervised learning methods. The results are compiled and presented in Table 3, and we have the following observations: 1) As training data is reduced in the low-resource scenario, a distinct degradation in generated audio quality could be witnessed in both test sets (test-seen and test-unseen). 2) Leveraging orders of magnitude text-only and audio-only data with self-supervised learning, the ViT-TTS achieve RTE scores of 0.082 and 0.068 respectively in test-unseen and test-seen, showing a significant promotion regardless of the unseen scene. In this way, the dependence on a large number of parallel audio-visual data can be reduced for constructing visual text-to-speech systems.

Method	MOS (\uparrow)	RTE (\downarrow)	MCD (\downarrow)
Finetune with 1 hour data			
Test-Seen	3.72 \pm 0.05	0.092	5.04
Test-Unseen	3.67 \pm 0.06	0.101	5.11
Finetune with 2 hours data			
Test-Seen	3.75 \pm 0.06	0.089	4.85
Test-Unseen	3.70 \pm 0.07	0.097	4.89
Finetune with 5 hours data			
Test-Seen	3.83 \pm 0.05	0.068	4.65
Test-Unseen	3.73 \pm 0.09	0.082	4.72

Table 3: Low resource evaluation results.

4.5 Case Study

We provide two examples of generation sampled from a large empty room with significant reverberation in the Test-Seen environment depicted in Figure 2, and have the following observations: 1) Mel-spectrograms produced by ViT-TTS are noticeably more similar to the target counterpart. 2) Moreover in challenging scenarios with invisible

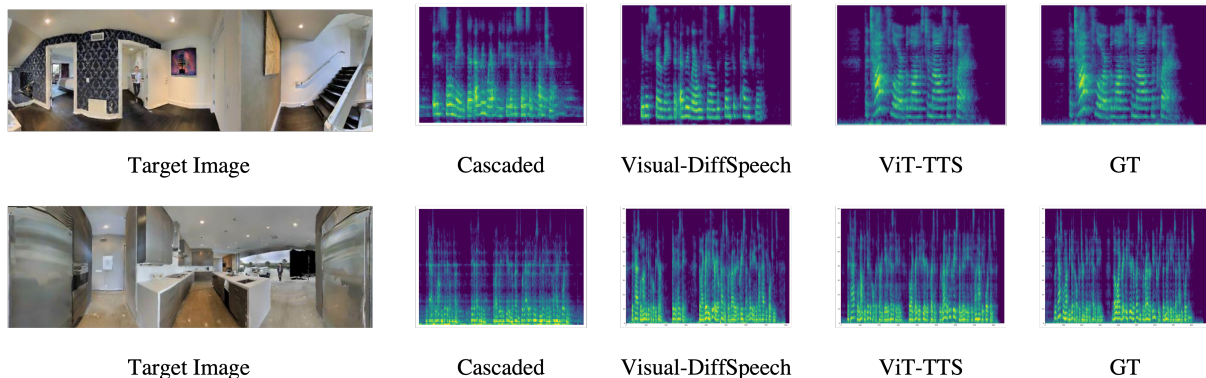


Figure 2: Visualizations of the ground truth and generated mel-spectrograms by different Visual TTS models. The text corresponding to the first line in test-seen is "it is so made that everywhere we feel the sense of punishment" while the second line in test-unseen is "the task will not be difficult returned david hesitating though i greatly fear your presence would rather increase than mitigate his unhappy fortunes".

scene images, cascaded systems suffer severely from the issue of noisy and reverb details missing, which is largely alleviated in ViT-TTS.

4.6 Ablation Studies

We conduct ablation studies to demonstrate the effectiveness of several key techniques on the Test-Unseen set in our model, including the encoder pre-training(EP), decoder pre-training(DP), visual input, and random image. The results of both subjective and objective evaluations have been presented in Table 4, and we have the following observations: 1) Removing the self-supervised encoder and decoder pre-training strategy results in a decline in all indicators, which demonstrates the effectiveness and efficiency of the proposed pre-training strategy in reducing data variance and promoting model convergence. 2) Without the input of RGB-D image and removing all of the modules related to the image causes a distinct degradation in RTE values, which demonstrates that our model successfully learns acoustics from the visual scene.

Method	MOS (\uparrow)	RTE (\downarrow)	MCD (\downarrow)
GT(voc.)	4.18 ± 0.07	0.008	1.50
ViT-TTS	3.86 ± 0.05	0.076	4.59
w/o EP	3.82 ± 0.07	0.078	4.63
w/o DP	3.83 ± 0.06	0.081	4.65
w/o Visual	3.78 ± 0.07	0.102	4.68
w/ RI	3.73 ± 0.08	0.103	4.75

Table 4: Ablation study results. EP, DP, and RI are encoder pre-training, decoder pre-training, and random images respectively.

Furthermore, we conducted a more detailed exploration of our model’s processing and reasoning

about different patches in the RGB-D images. To achieve this, we deliberately substituted the target image with random images, allowing us to determine whether the model can derive meaningful representations from visual inputs. Our findings show that after replacing the target image with a random image, the performance of our model significantly degraded, indicating that our model could model the room acoustic information of visual input.

5 Conclusion

In this paper, we proposed ViT-TTS, the first visual text-to-speech synthesis model that aimed to convert written text and target environmental images into audio that matches the target environment. ViT-TTS complemented the phoneme sequence with the visual information to generate high-perceived audio, opening up new avenues for practical applications of AR and VR, as it allows for a more immersive and realistic audio experience. To mitigate the data scarcity for training visual TTS tasks and model visual acoustic information, we 1) introduced a self-supervised learning framework to enhance both the visual-text encoder and denoiser decoder; 2) leveraged the diffusion transformer scalable in terms of parameters and capacity to learn visual scene information.

Experimental results demonstrate that ViT-TTS achieved new state-of-the-art results, outperforming cascaded systems and other baselines regardless of the visibility of the scene. With low-resource data (1h, 2h, 5h), ViT-TTS achieves comparative results with rich-resource baselines. To this end, ViT-TTS provided a solid foundation for future visual text-to-speech studies, and we envision that our approach will have far-reaching impacts on the fields of AR and VR.

6 Limitation and Potential Risks

As indicated in the experimental setup, we utilized ResNet-18 as our image feature extractor. While it is a classic extractor, there may be newer extractors that perform better. In future work, we will explore the use of superior extractors to enhance the quality of generated audio.

Moreover, our pre-trained encoder and decoder are based on the SoundSpace-Speech dataset, which, as described in the dataset section, is not sufficiently large. To address this limitation in future work, we will pre-train on a large-scale dataset to achieve better performance in low-resource scenarios.

ViT-TTS lowers the requirements for visual text-to-speech generation, which may cause fraud and scams by impersonating someone else’s voice. Furthermore, there is the potential for leading to the spread of false information and rumors.

References

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang Su, and Jun Zhu. 2023. One transformer fits all distributions in multi-modal diffusion at scale. *arXiv preprint arXiv:2303.06555*.

Andrew Brock, Jeff Donahue, and Karen Simonyan. 2018. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*.

Changan Chen, Ruohan Gao, Paul Calamia, and Kristen Grauman. 2022. Visual acoustic matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18858–18868.

Changan Chen, Wei Sun, David Harwath, and Kristen Grauman. 2023. [Learning audio-visual dereverberation](#).

Yu-An Chung and James Glass. 2020. Improved speech representations with multi-target autoregressive predictive coding. *arXiv preprint arXiv:2004.05274*.

Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng Chiu, James Qin, Ruoming Pang, and Yonghui Wu. 2021. w2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training. *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 244–250.

Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Jeff Donahue, Sander Dieleman, Mikołaj Bińkowski, Erich Elsen, and Karen Simonyan. 2020. End-to-end adversarial text-to-speech. *arXiv preprint arXiv:2006.03575*.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models.

Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460.

Rongjie Huang, Feiyang Chen, Yi Ren, Jinglin Liu, Chenye Cui, and Zhou Zhao. 2021. Multi-singer: Fast multi-singer singing voice vocoder with a large-scale corpus. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 3945–3954.

Rongjie Huang, Chenye Cui, Feiyang Chen, Yi Ren, Jinglin Liu, Zhou Zhao, Baoxing Huai, and Zhefeng Wang. 2022a. Singgan: Generative adversarial network for high-fidelity singing voice generation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 2525–2535.

Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. 2023a. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. *arXiv preprint arXiv:2301.12661*.

Rongjie Huang, Max WY Lam, Jun Wang, Dan Su, Dong Yu, Yi Ren, and Zhou Zhao. 2022b. Fastdiff: A fast conditional diffusion model for high-quality speech synthesis. *arXiv preprint arXiv:2204.09934*.

Rongjie Huang, Mingze Li, Dongchao Yang, Jia-tong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. 2023b. Audiogpt: Understanding and generating speech, music, sound, and talking head. *arXiv preprint arXiv:2304.12995*.

Rongjie Huang, Yi Ren, Jinglin Liu, Chenye Cui, and Zhou Zhao. Generspeech: Towards style transfer for generalizable out-of-domain text-to-speech. In *Advances in Neural Information Processing Systems*.

726	Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren. 2022c. Prodiff: Progressive fast diffusion model for high-quality text-to-speech. In <i>Proceedings of the 30th ACM International Conference on Multimedia</i> , pages 2595–2605.		
727			
728			
729			
730			
731	Rongjie Huang, Zhou Zhao, Jinglin Liu, Huadai Liu, Yi Ren, Lichao Zhang, and Jinzheng He. 2022d. Transpeech: Speech-to-speech translation with bilateral perturbation. <i>arXiv preprint arXiv:2205.12523</i> .		
732			
733			
734			
735	Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 4401–4410.		
736			
737			
738			
739			
740	Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In <i>International Conference on Machine Learning</i> , pages 5530–5540. PMLR.		
741			
742			
743			
744			
745	Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. <i>Proc. of NeurIPS</i> .		
746			
747			
748			
749	Junghyun Koo, Seungryeol Paik, and Kyogu Lee. 2021. Reverb conversion of mixed vocal tracks using an end-to-end convolutional deep neural network. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 81–85. IEEE.		
750			
751			
752			
753			
754			
755	Max WY Lam, Jun Wang, Rongjie Huang, Dan Su, and Dong Yu. 2021. Bilateral denoising diffusion models. <i>arXiv preprint arXiv:2108.11514</i> .		
756			
757			
758	Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catanzaro, and Sungroh Yoon. 2022. Bigvgan: A universal neural vocoder with large-scale training. <i>arXiv preprint arXiv:2206.04658</i> .		
759			
760			
761			
762	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Annual Meeting of the Association for Computational Linguistics</i> .		
763			
764			
765			
766			
767			
768			
769	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. <i>arXiv preprint arXiv:2301.12597</i> .		
770			
771			
772			
773	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>International Conference on Machine Learning</i> , pages 12888–12900. PMLR.		
774			
775			
776			
777			
	Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. 2019. Neural speech synthesis with transformer network. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 6706–6713.		778 779 780 781 782
	Wolfgang Mack, Shuwen Deng, and Emanuël Habets. 2020. Single-channel blind direct-to-reverberation ratio estimation using masking. In <i>Interspeech</i> .		783 784 785
	Amine Mezghani and A. Lee Swindlehurst. 2018. Blind estimation of sparse broadband massive MIMO channels with ideal and one-bit ADCs. <i>IEEE Transactions on Signal Processing</i> , 66(11):2972–2983.		786 787 788 789
	MoonInTheRiver. 2021. DiffSinger. https://github.com/MoonInTheRiver/DiffSinger .		790 791 792
	Prateek Murgai, Mark Rau, and Jean-Marc Jot. 2017. Blind estimation of the reverberation fingerprint of unknown acoustic environments. <i>Journal of The Audio Engineering Society</i> .		793 794 795 796
	Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. Wavenet: A generative model for raw audio. <i>arXiv preprint arXiv:1609.03499</i> .		797 798 799 800 801
	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. <i>arXiv preprint arXiv:1807.03748</i> .		802 803 804
	William Peebles and Saining Xie. 2023. Scalable diffusion models with transformers.		805 806
	Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. 2021. Grad-tts: A diffusion probabilistic model for text-to-speech. In <i>International Conference on Machine Learning</i> , pages 8599–8608. PMLR.		807 808 809 810 811
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.		812 813 814 815 816 817
	Rama Ratnam, Douglas L Jones, Bruce C Wheeler, William D O’Brien Jr, Charissa R Lansing, and Albert S Feng. 2003. Blind estimation of reverberation time. <i>The Journal of the Acoustical Society of America</i> , 114(5):2877–2892.		818 819 820 821 822
	Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. 2019. FastSpeech: Fast, robust and controllable text to speech. <i>Advances in Neural Information Processing Systems</i> , 32.		823 824 825 826
	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In <i>Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany</i> ,		827 828 829 830 831

- 832 *October 5-9, 2015, Proceedings, Part III 18*, pages
833 234–241. Springer.
- 834 Andy Sarroff and Roth Michaels. 2020. Blind arbitrary
835 reverb matching. In *Proceedings of the 23rd*
836 *International Conference on Digital Audio Effects*
837 *(DAFx-2020)*, volume 2.
- 838 Manfred R Schroeder. 1965. New method of measuring
839 reverberation time. *The Journal of the Acoustical*
840 *Society of America*, 37(6):1187–1188.
- 841 Hao Sun, Xu Tan, Jun-Wei Gan, Hongzhi Liu, Sheng
842 Zhao, Tao Qin, and Tie-Yan Liu. 2019. Token-level
843 ensemble distillation for grapheme-to-phoneme con-
844 version. *arXiv preprint arXiv:1904.03446*.
- 845 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
846 Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
847 Kaiser, and Illia Polosukhin. 2017. [Attention is all](#)
848 [you need](#).
- 849 Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui
850 Wu, Ron J Weiss, Navdeep Jaitly, Zongheng Yang,
851 Ying Xiao, Zhifeng Chen, Samy Bengio, et al.
852 2017. Tacotron: Towards end-to-end speech syn-
853 thesis. *arXiv preprint arXiv:1703.10135*.
- 854 Feifei Xiong, Stefan Goetze, Birger Kollmeier, and
855 Bernd T Meyer. 2018. Joint estimation of reverberation
856 time and early-to-late reverberation ratio from
857 single-channel speech signals. *IEEE/ACM Transactions on*
858 *Audio, Speech, and Language Processing*,
859 27(2):255–267.
- 860 Pavel Zahorik. 2002. Direct-to-reverberant energy ratio
861 sensitivity. *The Journal of the Acoustical Society of*
862 *America*, 112(5):2110–2117.
- 863 Lvmin Zhang and Maneesh Agrawala. 2023. [Adding](#)
864 [conditional control to text-to-image diffusion models](#).

865 A TRANSFORMER CONFIGURATION

866 The details of transformer denoisers are shown in
867 Table 5, while B, M, L, and XL means the base,
medium, large, extra large respectively.

Model	layers	Hidden Size	Heads
Transformer-S	4	256	8
Transformer-B	5	384	12
Transformer-L	6	512	16
Transformer-XL	8	768	16

868 Table 5: Diffusion Transformer Configs.

869 B ARCHITECTURE

870 We list the model hyper-parameters of ViT-TTS in
871 Table 6.

872 C DIFFUSION POSTERIOR 873 DISTRIBUTION

874 Firstly we compute the corresponding constants
875 respective to diffusion and reverse process:

$$876 \alpha_t = \prod_{i=1}^t \sqrt{1 - \beta_i} \quad \sigma_t = \sqrt{1 - \alpha_t^2} \quad (6)$$

877 The Gaussian posterior in diffusion process is
878 defined through the Markov chain, where each iter-
879 ation adds Gaussian noise.

$$880 q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}),$$

$$881 q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (7)$$

882 We emphasize the property observed by (Ho
883 et al., 2020), the diffusion process can be computed
in a closed form:

$$884 q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{x}_0, \sigma_t \mathbf{I}) \quad (8)$$

885 Applying Bayes' rule, we can obtain the forward
886 process posterior when conditioned on \mathbf{x}_0

$$887 q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)}$$

$$888 = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\boldsymbol{\beta}}_t \mathbf{I}), \quad (9)$$

$$889 \text{where } \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\alpha_{t-1} \beta_t}{\sigma_t} \mathbf{x}_0 + \frac{\sqrt{1 - \beta_t} (\sigma_{t-1})}{\sigma_t} \mathbf{x}_t, \quad \tilde{\boldsymbol{\beta}}_t = \frac{\sigma_{t-1}}{\sigma_t} \beta_t$$

Algorithm 1 Training procedure

- 1: **Input:** The denoiser ϵ_θ , diffusion step T and variance condition c.
- 2: **repeat**
- 3: Sample $\mathbf{x}_0 \sim q_{data}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4: Take gradient descent steps on $\nabla_\theta \| \epsilon - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, c, t) \|^2$.
- 5: **until** convergence

Algorithm 2 Sampling

- 1: **Input:** The denoiser ϵ_θ , and variance condition c.
- 2: Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 3: **for** $t = T, \dots, 1$ **do**
- 4: **if** $t = 1$ **then**
- 5: $\mathbf{z} = 0$
- 6: **else**
- 7: Sample $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 8: **end if**
- 9: Sample $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(\mathbf{x}_t, c, t)) + \sigma_t \mathbf{z}$
- 10: **end for**

D DIFFUSION ALGORITHM

E EVALUATION MATRIX

E.1 Evaluation Metrics

890 We measure the sample quality of the generated
891 waveform using both objective metrics and subjective
892 indicators. The objective metrics we collected
893 are designed to measure varied aspects of wave-
894 form quality between the ground-truth audio and
895 the generated sample. Following the common prac-
896 tice of (Huang et al., 2022c; MoonInTheRiver,
897 2021; Popov et al., 2021), we randomly select a
898 part of the test set for objective evaluation, here is
899 50. We provide the following metrics: (1) **RT60**
900 **Error(RTE)**-the correctness of the room acoustics
901 between the predicted waveform and target wave-
902 form's RT60 values. RT60 indicates the reverbera-
903 tion time in seconds for the audio signal to decay
904 by 60 dB, a standard metric to characterize room
905 acoustics. We estimate the RT60 directly from
906 magnitude spectrograms of the output audio, using
907 a model trained with disjoint SoundSpaces data.
908 (2) **Mel Cepstral Distortion(MCD)**-measures the
909 spectral distance between the synthesized and refer-
910 ence mel-spectrum features. The utilization of RTE
911 is solely intended for evaluating the room acoustic
912 performance of the generated audio, and as an ad-
913

Hyperparameter		ViT-TTS
Visual-Text Encoder	Phoneme Embedding	256
	Pre-net Layers	3
	Pre-net Hidden	256
	Visual Conv2d Kernel	(7, 7)
	Visual Conv2d Stride	(2, 2)
	Encoder Layers	4
	Encoder Hidden	256
	Encoder Conv1d Kernel	9
	Encoder Conv1D Filter Size	1024
	Encoder Attention Heads	2
	Encoder Dropout	0.1
Variance Predictor	Variance Predictor Conv1D Kernel	3
	Variance Predictor Conv1D Filter Size	256
	Variance Predictor Dropout	0.5
Spectrogram Denoiser	Diffusion Embedding	256
	Transformer Layers	4
	Transformer Hidden	256
	Transformer Attention Heads	8
	Position Embedding	256
	Scale/Shift Size	256
Total Number of Parameters		32.38M

Table 6: Hyperparameters of ViT-TTS models.

ditional measure, we have incorporated the MCD metric to assess the quality of the mel-spectrogram.

For subjective metrics, we use crowd-sourced human evaluation via Amazon Mechanical Turk, where raters are asked to rate **Mean Opinion Score(MOS)** on a 1-5 Likert scale.

E.2 RT60 Estimator

Following (Chen et al., 2022), we first encode the 2.56s speech clips as spectrograms, process them with a ResNet18 (Oord et al., 2018) and predict the RT60 of the speech. The ground truth RT60 is calculated with the Schroeder (Schroeder, 1965). We optimize the MSE loss between the predicted RT60 and the ground truth RT60.

E.3 MOS Evaluation

To probe audio quality, we conduct the MOS (mean opinion score) tests and explicitly instruct the raters to “focus on examining the audio quality, naturalness and whether the audio matches with the given image.”. The testers present and rate the samples, and each tester is asked to evaluate the subjective naturalness on a 1-5 Likert scale.

Our subjective evaluation tests are crowd-sourced and conducted via Amazon Mechanical Turk. These ratings are obtained independently for model samples and reference audio, and both are reported. The screenshots of instructions for testers have been shown in Figure 3. A small subset


of speech samples used in the test is available at <https://ViT-TTS.github.io/>

944

945

Instructions Shortcuts How natural (i.e. human-sounding) and compatibility (i.e. matching degree with picture) is this recording? Please focus on examining the audio quality, naturalness and whether it matches the environment in the picture, and ignore the differences of sty...

Environment picture:



Transcripts: in truth she seemed absolutely hidden behind it

0:00 / 0:03

Select an option

Excellent - Completely natural and matched speech - 5	1
4.5	2
Good - Mostly natural and matched speech - 4	3
3.5	4
Fair - Equally natural and matched speech - 3	5
2.5	6
Poor - Mostly unnatural and unmatched speech - 2	7
1.5	8

Figure 3: Screenshots of subjective evaluations.