SPEAKING GUIDED BY LISTENING: UNSUPERVISED TEXT-TO-SPEECH GENERATIVE MODEL GUIDED BY END-TO-END SPEECH RECOGNITION

Anonymous authors

Paper under double-blind review

Abstract

We propose to utilize end-to-end automatic speech recognition (E2EASR) as a guidance model to realize unsupervised text-to-speech (TTS). An unconditional score-based generative model (SGM) is trained with untranscribed speech data. In the sampling stage, the unconditional score estimated by the SGM is combined with the gradients from ASR models by the Bayes rule to get the conditional score. We use a set of small ASR models trained only on 80-hour labeled ASR data to guide the unconditional SGM and generate speech with high-quality scores in both objective and subjective evaluation. Similarly, we can also use additional speaker verification models to control speaker identity for the synthesized speech. That allows us to do the zero-shot TTS for the target speaker with a few seconds of enrollment speech. Our best unsupervised synthesized speech gets $\sim 8\%$ word error rate in testing, and the best speaker-controlled TTS gets 3.3 mean opinion score (MOS) in the speaker similarly testing.

025 026

006

008 009 010

011 012 013

014

015

016

017

018

019

021

1 INTRODUCTION

027 028 029

Text-to-speech (TTS) systems have made significant progress due to the need for natural, expressive speech in various applications such as virtual assistants, audiobooks, and automated customer ser-031 vice. With the rise of neural network-based models, TTS has transformed from traditional concatenative and parametric methods to more advanced deep learning approaches, leading to significant 032 improvements in quality. (Ren et al., 2019a; Li et al., 2019; Shen et al., 2018). The adoption of 033 autoregressive models like Tacotron and Transformer TTS (Wang et al., 2017; Li et al., 2019) has 034 allowed for significant improvements in the synthesis of natural-sounding speech. The emergence 035 of non-autoregressive models, such as FastSpeech (Ren et al., 2019a; 2020), further expanded the capabilities and inference speed of TTS by introducing parallel speech generation techniques. More-037 over, diffusion models Sohl-Dickstein et al. (2015); Ho et al. (2020); Song & Ermon (2019), which have shown high performance in image generation tasks, have recently been adapted for speech synthesis with promising results (Jeong et al., 2021; Popov et al., 2021; Huang et al., 2022; Tae et al., 040 2022). Despite these advancements, mainstream TTS models rely heavily on high-quality paired 041 text-speech data. This dependence remains a critical barrier to developing robust TTS systems for 042 languages and speakers with limited available data.

043 In this paper, we propose an E2EASR-guided method for unsupervised TTS¹. We trained an uncon-044 ditional score-based generative model (SGM) on the 522-hour LibriTTS-R (Koizumi et al., 2023) dataset without using any text or speaker label. Only 80-hour WSJ (Garofolo, John S. et al., 2007) 046 labeled ASR data is used for training the guidance ASR models. ASR data is usually much more 047 accessible to collect than TTS data. In fact, the ASR system can be directly trained on noisy or low-048 quality speech, while the TTS usually requires studio-quality recording. In addition, many of the 049 recent TTS works still need phoneme-level supervision to learn the phoneme duration (Ammar Abbas et al., 2022; Effendi et al., 2022; Kim et al., 2020; 2022b), while the majority of ASR systems, 050 in recent years, are simply trained in end-to-end style with sentence-level annotations (Graves et al., 051 2013; Bahdanau et al., 2016; Kim et al., 2017). We propose to use multiple independent ASR sys-052

¹Demos are available at https://asr-guided-tts.github.io/ (Anonymous version)

tems trained on the same dataset to provide more robust joint guidance. This allows us to generate
 higher-quality speech than what is possible by using guidance from only a single-ASR system.

We also propose using a speaker verification model for zero-shot target speaker TTS, in which only a few seconds of enrollment speech to generate speech for unseen speakers. By assuming the two statistically independent and using the Bayes rule we can combine the guidance from ASR and speaker models without conflicts or scale issues. The two can thus jointly guide the unconditional SGM to realize target speaker- and text-conditioned diffusion generation.

Another benefit of the E2EASR guidance is that the duration model in conventional TTS, which estimates the phoneme duration, can be omitted. Our experiments show that we can generate highquality speech with the input of target text and a total desired speech length. The end-to-end training of ASR models enables them to implicitly handle phoneme duration, helping the SGM generate speech at a natural pace.

067 068

069

070

079

2 BACKGROUND

2.1 UNCONDITIONAL DIFFUSION FOR SPEECH GENERATION

071 We follow the unconditional generation model (SGM) introduced by Song et al. (2021). It unifies the 072 denoising score matching with Langevin dynamics (SMLD) (Song & Ermon, 2019) and denoising 073 diffusion probabilistic model (DDPM) (Ho et al., 2020) frameworks by using stochastic differential 074 equations (SDEs). For speech generation, we model the speech data in the mel-spectrogram domain 075 $\mathbf{X} \in \mathbb{R}^{L \times F}$, where L is the length of the spectrum and F is the number of frequency banks. Let 076 $p_0(\mathbf{X})$ be the probability density function (p.d.f.) of the clean speech data. The SDE describes 077 the forward diffusion process, which converts the distribution $p_0(\mathbf{X})$ to a simple prior distribution $p_T(\mathbf{X})$ by a continuous time variable $t \in [0, T]$: 078

$$d\mathbf{X} = \mathbf{F}(\mathbf{X}, t)dt + g(t)d\mathbf{W},\tag{1}$$

where dt is an infinitesimal timestep, $\mathbf{W} \in \mathbb{R}^{L \times F}$ is a Brownian motion, $\mathbf{F}(\cdot)$ is a matrix-valued *drift* function, and g(t) is the scalar value *diffusion* coefficient determined by t. $\mathbf{F}(\cdot)$ is the deterministic part of the SDE, while g(t) controls the scale of the noise-adding process. Based on previous research from Anderson (1982), there is a reverse SDE that describes the reverse diffusion process corresponding to the above forward process:

$$d\mathbf{X} = [-\mathbf{F}(\mathbf{X}, t) + g(t)^2 \nabla_{\mathbf{X}} \log p_t(\mathbf{X})] dt + g(t) d\bar{\mathbf{W}},$$
(2)

where $\bar{\mathbf{W}}$ is the Brownian motion in the reverse process, dt is a negative infinitesimal timestep, $\nabla_{\mathbf{X}} \log p_t(\mathbf{X})$ is the gradient of the logarithm data distribution $p_t(\mathbf{X})$ at timestep t, i.e., the *score* of the distribution. We can train a model to approximate such a score function, i.e., a score model $s_{\theta}(\mathbf{X}, t)$ parameterized by θ to estimate $\nabla_{\mathbf{X}} \log p_t(\mathbf{X})$ by using the score matching (Hyvärinen & Dayan, 2005; Song & Ermon, 2019) method:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\arg\min} \mathbb{E}_{t, \mathbf{X}_0, \mathbf{X}_t | \mathbf{X}_0} \left[\left\| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{X}_t, t) - \nabla_{\mathbf{X}_t} \log p_{0t}(\mathbf{X}_t | \mathbf{X}_0) \right\|^2 \right],$$
(3)

where $\mathbf{X}_0 \sim p_0(\mathbf{X})$ is the clean training data and $\mathbf{X}_t \sim p_{0t}(\mathbf{X}_t | \mathbf{X}_0)$ is the perturbed data sampled from the conditional distribution $p_{0t}(\mathbf{X}_t | \mathbf{X}_0)$ at timestep t.

After the score model has been trained, $\nabla_{\mathbf{X}} \log p_t(\mathbf{X})$ in Eq.2 can be replaced with $s_{\theta}(\mathbf{X}, t)$ for inference. We can start from the prior distribution $\mathbf{X}_T \sim p_T(\mathbf{X})$, and generate samples from the target distribution by using a numerical SDE solver to solve the reverse SDE. Commonly used numerical solvers include the Euler-Maruyama method and predictor-corrector sampler (Song et al., 2021).

101 102 103

107

092

2.2 CLASSIFIER GUIDANCE

Classifier guidance (Song et al., 2021; Ho et al., 2020) can generate desired data by using the un conditional diffusion model and an external classifier. To generate the class conditioned data, Eq.2
 can be updated into the following conditioned form:

$$d\mathbf{X} = [-\mathbf{F}(\mathbf{X}, t) + g(t)^2 \nabla_{\mathbf{X}} \log p_t(\mathbf{X}|y)] dt + g(t) d\mathbf{\bar{W}},$$
(4)

where $p_t(\mathbf{X}|y)$ is the data p.d.f conditioned by discrete class y. According to the Bayes rule, $\nabla_{\mathbf{X}} \log p_t(\mathbf{X}|y)$ can be rewritten as:

$$\nabla_{\mathbf{X}} \log p_t(\mathbf{X}|y) = \nabla_{\mathbf{X}} \log p_t(\mathbf{X}) + \nabla_{\mathbf{X}} \log P_t(y|\mathbf{X}), \tag{5}$$

The first term in Eq.5 can be estimated by the unconditional SGM $s_{\theta}(\mathbf{X}, t)$, and the second term in Eq.5 can be estimated by the external classifier $P_{\phi}(y|\mathbf{X})$ parameterized by ϕ . The labeled data used to train the classifier can also be different from the data used for the unconditional SGM. This allows a TTS system to be built without transcribed TTS data but only with ASR data. The latter is more accessible to collect since ASR does not require high-quality recordings as TTS does and thus can be scaled more easily.

117 118 119

3 RELATED WORKS

120 Recently, many diffusion-based TTS models have proven to be remarkably successful in the text-to-121 speech task (TTS) (Jeong et al., 2021; Popov et al., 2021; Huang et al., 2022; Tae et al., 2022). This 122 emerging class of generative models adopts an iterative generative approach, where, during training, 123 a complex data distribution is gradually corrupted by Gaussian noise (Song et al., 2021). These 124 models are trained to estimate the gradient fields that reverse this process from a noisy prior, guiding 125 the data sample back to its original distribution. And most of those diffusion models are conditioned 126 on the text or semantic tokens. To train such models conditioned by the input text, large scale of 127 high-quality paired text-speech data are required, posing a major practical problem for TTS (Ren et al., 2019b). 128

A recent work named Guided-TTS (Kim et al., 2022a) uses the classifier guidance method for TTS by using a phoneme classifier as the guidance model. Guided-TTS is the most related work to this paper. The main difference between Guided-TTS and this paper includes:

1) Guided-TTS needs to train a frame-level phoneme classifier, an additional duration predictor is required to estimate the duration of each phoneme, and the diffusion model is trained on single-speaker datasets. They use large-scale 960-hour labeled ASR data to train a phone-level classifier. while we only use 80-hour labeled ASR data to train end-to-end ASR guidance models.

2) Guided-TTS trains unconditional diffusion models on single-speaker training datasets. It can not synthesize the target speaker differently from the training speaker. A follow-up work named Guided-TTS2 (Kim et al., 2022b) replaced the unconditional diffusion model with a speaker-conditioned diffusion model to make it able for speaker-conditional TTS. That requires the diffusion training data to be labeled for the speaker. In this work, we proposed to use speaker verification models for TTS guidance. That allows us to do the zero-shot target speaker TTS with a few seconds of enrollment speech. We do not need any speaker labels for the diffusion model training.

3) In the Guided-TTS, the authors use norm-based gradient scaling methods to combine the gradients from the unconditional diffusion model and the guidance classifier. In this paper, we follow Eq.5 to combine $\nabla_{\mathbf{X}} \log p_t(\mathbf{X})$ and $\nabla_{\mathbf{X}} \log P_t(y|\mathbf{X})$ without any scaling weight that hurt the Bayes rule. The impact of guidance gradients is instead tuned by the temperature in the Softmax function when estimating $P_t(y|\mathbf{X})$. Following the Bayes rule allows us to combine multiple guidance models more easily and safely by making an independent assumption between the guidance models.

149 150

151

4 TEXT-TO-SPEECH GUIDED BY END-TO-END ASR

In this section, we propose using E2EASR as guidance for unsupervised TTS. Our method involves three different modules: 1) an unconditional score-based generative model, which models the distribution of clean speech from the speech data without transcription; 2) E2EASR systems that provide the gradient guidance conditioned by the target text; 3) optional speaker verification models for controlling the target speaker identity, enabling zero-shot target speaker TTS with only a few seconds of reference speech.

158 159

- 4.1 UNCONDITIONAL SCORE-BASED GENERATIVE MODEL GUIDED BY E2EASR
- Here, we introduce the E2EASR guidance mechanism that enables TTS with the unconditional SGM. We denote the one-hot text tokens as $\mathbf{Y} \in \{0, 1\}^{K \times V}$, where K is the length of the text and

162 V is the vocabulary size. The TTS task by score-based modeling is to model the distribution of 163 $p(\mathbf{X}|\mathbf{Y})$, where $\mathbf{X} \in \mathbb{R}^{L \times F}$ is the mel-spectrogram with length² L and the number of frequency 164 banks F. We train the unconditional SGM with the objective function in Eq.3. To generate the 165 speech conditioned by \mathbf{Y} , the reverse SDE can be written as:

- 166
- 167 168

177

180 181

$$d\mathbf{X} = \left[-\mathbf{F}(\mathbf{X}, t) + g(t)^2 \nabla_{\mathbf{X}} \log p_t(\mathbf{X}|\mathbf{Y})\right] dt + g(t) d\bar{\mathbf{W}},\tag{6}$$

similar to that discussed in Sec.2.2, $\nabla_{\mathbf{X}} \log p_t(\mathbf{X}|\mathbf{Y})$ can be decomposed into the sum of $\nabla_{\mathbf{X}} \log p_t(\mathbf{X})$ and $\nabla_{\mathbf{X}} \log P_t(\mathbf{Y}|\mathbf{X})$: The former gradient can be estimated by the unconditional SGM $s_{\theta}(\mathbf{X}, t)$, while the latter gradient can be calculated by differentiating the ASR models, which were trained to estimate $P_t(\mathbf{Y}|\mathbf{X})$.

173 We use the joint CTC-Attention (Kim et al., 2017; Watanabe et al., 2018) E2EASR, which combines 174 the connectionist temporal classification (CTC) (Graves et al., 2006; 2013) and Attention-based 175 Encoder-Decoder (AED) (Chan et al., 2016; Bahdanau et al., 2016). With the joint CTC-Attention 176 ASR, $P_t(\mathbf{Y}|\mathbf{X})$ can be jointly represented as:

$$P_t(\mathbf{Y}|\mathbf{X}) = P_t(Y_{CTC} = \mathbf{Y}, Y_{AED} = \mathbf{Y}|\mathbf{X}), \tag{7}$$

If we assume the CTC task and AED components estimate the target text Y independently in their
 different feature spaces, then:

$$\nabla_{\mathbf{X}} \log P_t(\mathbf{Y}|\mathbf{X}) = \nabla_{\mathbf{X}} \log P_t^{CTC}(\mathbf{Y}|\mathbf{X}) + \nabla_{\mathbf{X}} \log P_t^{AED}(\mathbf{Y}|\mathbf{X}).$$
(8)

We train the ASR model with perturbed data \mathbf{X}_t and text label \mathbf{Y} to make sure that the model can estimate a reliable $P_t(\mathbf{Y}|\mathbf{X})$ at each SDE timestamp t.

So far, with the above-proposed method it should be theoretically possible to control the generated speech from the unconditional SGM. However, the mel-spectrogram space $\mathbb{R}^{L \times F}$ is sparse, and in the sampling stage, the numerical SDE solver tends to sample a locally optimal $\hat{\mathbf{X}}_0$ leading to high log-likelihood with guidance ASR but poor quality. To alleviate this problem, we propose joint guidance by multiple ASR systems. We train N compact ASR systems with a small number of parameters (≈ 15 M) using the same perturbed training data and assume that each ASR system is independent of the others when estimating Y, then the 8 can be rewritten as:

191 192

193

194

195 196

197

199 200

204 205

209 210

211 212

$$\nabla_{\mathbf{X}} \log P_t(\mathbf{Y}|\mathbf{X}) = \sum_{n=1}^{N} (\nabla_{\mathbf{X}} \log P_t^{CTC_n}(\mathbf{Y}|\mathbf{X}) + \nabla_{\mathbf{X}} \log P_t^{AED_n}(\mathbf{Y}|\mathbf{X})),$$
(9)

where $P_t^{CTC_n}(\mathbf{Y}|\mathbf{X})$ and $P_t^{CTC_n}(\mathbf{Y}|\mathbf{X})$ are the CTC-ASR and AED-ASR the *n*-th joint E2EASR model, respectively.

4.2 SPEAKER CONDITIONAL GUIDANCE FOR ZERO-SHOT TTS

If speaker-id conditioning is added, the reverse SDE becomes:

$$d\mathbf{X} = [-\mathbf{F}(\mathbf{X}, t) + g(t)^2 \nabla_{\mathbf{X}} \log p_t(\mathbf{X} | \mathbf{Y}, c)] dt + g(t) d\bar{\mathbf{W}},$$
(10)

where c is the class of the target speaker, $\nabla_{\mathbf{X}} \log p_t(\mathbf{X}|\mathbf{Y}, c)$ is the score conditioned by target speaker and text. Given that c and **Y** are conditionally independent, the conditioned score can be rewritten as:

$$\nabla_{\mathbf{X}} \log p_t(\mathbf{X}|\mathbf{Y}, c) = \nabla_{\mathbf{X}} \log p_t(\mathbf{X}) + \nabla_{\mathbf{X}} \log P_t(\mathbf{Y}, c|\mathbf{X}), \tag{11}$$

$$= \nabla_{\mathbf{X}} \log p_t(\mathbf{X}) + \nabla_{\mathbf{X}} \log P_t(\mathbf{Y}|\mathbf{X}) + \nabla_{\mathbf{X}} \log P_t(c|\mathbf{X}), \quad (12)$$

where $P_t(c|\mathbf{X})$ can be estimated by a speaker verification model \mathbf{m}_{μ} , where μ is the pretained parameters. In the zero-shot target speaker TTS scenario, we can estimate $P_t(c|\mathbf{X})$ with the following method:

$$\mathbf{s}^c = \mathbf{m}_{\boldsymbol{\mu}}(\mathbf{X}^c),\tag{13}$$

 $\mathbf{s}_t = \mathbf{m}_{\boldsymbol{\mu}}(\hat{\mathbf{X}}_t),\tag{14}$

$$P_t(c|X = \hat{\mathbf{X}}_t) \approx \frac{e^{\beta d(\mathbf{s}^c, \mathbf{s}_t)}}{e^{\beta d(\mathbf{s}^c, \mathbf{s}_t)} + \sum_i^M e^{\beta d(\mathbf{s}^i, \mathbf{s}_t)}},\tag{15}$$

²The proposed method does not rely on a duration model. In the inference stage, L can be empirically set or estimated by some simple algorithm based on the text length K.

where \mathbf{X}^c is an enrollment utterance from target speaker c, \mathbf{s}^c is the enrollment speaker vector extracted from \mathbf{X}^c ; $\hat{\mathbf{X}}_t$ is the data sampled at current timestep t in the diffusion reverse process, \mathbf{s}_t is the current speaker vector extracted by the speaker model; scalar β is a manually set parameter³ which controls the sharpness of the distribution, whose inverse is also known as *temperature*. By tuning β , we can tune the impact of guidance gradients; function d is a metric that measures the similarity of two vector embeddings, which is cosine similarity in our implementation;

 s^i is the *i*-th speaker embedding in the pre-trained parameters μ , which learned from the training data, and M is the number of speaker embeddings in the pertained model. Similar to Eq.9 introduced in the Sec.4.1, we can also use multiple independent speaker models to get a more robust guidance.

225 226 227

228 229

230

222

223

224

- 5 EXPERIMENTS
- 5.1 Dataset

We use three datasets in this paper; one to train the unsupervised score generation model, one for the ASR guidance model, and another for the speaker guidance model.

The first one is the LibriTTS-R (Koizumi et al., 2023) dataset. It is a quality-improved dataset derived from LibriTTS (Zen et al., 2019). The sampling rate is 24 kHz. LibriTTS-R provides the text annotation for each sample, but we do not use it during training (*unsupervised*). We use the 522-hour speech partition to train the unconditional score-based generative model. As said, such an unconditional score-based diffusion process is performed on 100-dimensional mel-scaled spectrum features (i.e., F = 100) extracted from the speech signal.

The dataset we use to train the ASR models is the WSJ SI-284 dataset. It contains about 80 hours of training data. The original dataset is 16kHz, and we upsample the data to 24 kHz and extract mel-spectrum features with the same parameter as that for LibriTTS-R.

The third dataset we use is the Voxceleb2 (Chung et al., 2018) dataset, which is used to train the speaker verification model. It includes 5994 speakers for model training, and the total amount of training data is about 2000 hours.

- 246
- 247 5.2 MODEL CONFIGURATIONS248

249 We use the sub-variance preserving (sub-VP) introduced in Song et al. (2021) as the SDEs for un-250 conditional score model training and set hyperparameters $\beta_{max} = 16.0, \beta_{min} = 0.1$. The NCSN++ (Song et al., 2021; Richter et al., 2023) network is used for the score model. The diffusion model 251 is trained using Adam optimizer Kingma & Ba (2014) with an initial learning rate of 10^{-4} . The 252 learning rate decays by a factor of 0.97 in every epoch. In each epoch, we optimize the model by 253 2000 steps with a batch of 2. The predictor-corrector sampler Song et al. (2021) is used to solve 254 the reverse SDE with 100 discrete time step in the sampling stage. We apply the neural vocoder 255 BigVGAN Lee et al. (2022), to resynthesize speech from the generated mel-spectrogram feature. 256

We trained 12 ASR models with the data perturbed by SDE on WSJ SI-284 as the ASR guidance
models with the ESPNet toolkit (Watanabe et al., 2018). The 12 models include 4 different model
sizes, each of them with 3 different byte pair encoding (BPE) (Sennrich et al., 2016). The detailed
configurations are listed in the Appendix.A.1. All the ASR models have an AED structure and are
optimized with joint CTC+attention loss (Kim et al., 2017).

For speaker guidance, we trained the speaker guidance models using the wespeaker toolkit Wang et al. (2024). We chose the ResNet34-based r-vector Zeinali et al. (2019) as the speaker embedding extractor and trained the systems on the SDE perturbed Voxceleb2 dataset following the wespeaker recipe⁴. To enhance the speaker guidance in the diffusion inference process, we trained two speaker guidance models by applying different random seeds.

267 268

 $^{^{3}\}beta = 1000$ in our experiments.

⁴https://github.com/wenet-e2e/wespeaker/tree/master/examples/voxceleb/ v2

270 5.3 EVALUATION METRICS271

The generated speech is evaluated using both objective and subjective metrics. We generated 100 samples from the LibriTTS-R testing set for objective evaluation and 15 for subjective evaluation. In addition to the target text, we need to specify the total length that needs to be generated for each sample before sampling. We adopt the length of ground truth audio in the generation in most experiments unless otherwise stated. A detailed ablation study of target speech length will be conducted in Sec. 6.4.

278 **Objective evaluation**. We first report the word error rate (WER), which is tested on an ASR model trained on the Librispeech (Panayotov et al., 2015) dataset. The evaluation model is available online 279 ⁵. The second metric is UTMOS Saeki et al. (2022), a pseudo mean opinion score (MOS) predicted 280 by a neural network. The third is the SpeechBERT (Chuang et al., 2020; Saeki et al., 2024) score, 281 which measures the BERTScore (Zhang et al., 2019) for generated and reference speech with self-282 supervised dense speech features. The last metric is speaker similarity, which is used to evaluate the 283 effect of speaker guidance. We extract the speaker embeddings using a publicly available speaker-id 284 embedding model⁶, and calculate the cosine similarity between the enrolment speech and generated 285 speech. The UTMOS, SpeechBERT score, and speaker similarity are evaluated by the VERSA toolkit⁷. 286 287

Subjective evaluation. We perform a human evaluation on the generated examples based on three criteria: N-MOS (naturalness Mean Opinion Score) for fluidity and naturalness, M-MOS (meaning-fulness Mean Opinion Score) for meaningfulness and content quality, and S-MOS (Speaker Similarity Mean Opinion Score) for the speaker similarity between generated speech and the enrollment utterance. 20 listeners are asked to evaluate 15 utterances on a scale from 1 to 5. The instructions for subjective evaluations are provided in Appendix A.2.

6 RESULTS AND ANALYSIS

Table 1: Objective evaluation for different ASR guidance. The number of parameters of the unconditional SGM is 147.4 M. We also list the additional parameters in guidance ASR models.

ASR Guidance ID	# total ASR param. (M)	WER(%) \downarrow	UTMOS \uparrow	SpeechBERT \uparrow
{1}	14.9	81.3	2.84	0.64
$\{1, 2\}$	29.8	39.9	3.18	0.69
$\{1, 2, 3\}$	44.8	24.4	3.28	0.71
$\{1, \cdots, 6\}$	89.7	14.9	3.42	0.72
$\{1, \cdots, 9\}$	147.2	11.1	3.44	0.73
$\{1, \cdots, 12\}$	229.9	10.1	3.47	0.74
Ground Truth	-	3.3	4.15	1.00

311

322

323

305 306

288

289

290

291

292

293

295 296 297

298

6.1 EXPERIMENTAL RESULTS ON ASR GUIDANCE

312 We compare the effect of ASR guidance in Table.1. We trained the 12 ASR guidance model on 313 80-hour WSJ data, and they are identified by ID 1-12. The detailed configurations of them can be found in the Appendix A.1. We first tried to guide the unconditional SGM trained on LibriTTS-R 314 with a single ASR, the results of which are listed in the first line of Table.1. We found that the 315 speech generated by the guided diffusion by just one ASR system performed poorly in the objective 316 evaluation. In contrast, if the WER is evaluated using the same ASR model used for guidance, 317 it is near zero. That is because the generated speech data is optimized directly for the guidance 318 ASR in the sampling process, similarly as it happens for white-box attacks (Wang et al., 2022). 319 In other words, the guidance model can be easily fooled by low-quality samples generated by the 320 guided sampling, especially as the data space $\mathbb{R}^{L \times F}$ is high-dimensional and sparse. A simple and 321

⁵https://huggingface.co/asapp/e_branchformer_librispeech

⁶https://huggingface.co/espnet/voxcelebs12_rawnet3

⁷https://github.com/shinjiwlab/versa/tree/main

straightforward approach that we adopt is to use multiple independent ASR models to provide more
 robust guidance, as described in Eq.9. The joint likelihood assessment of multiple ASR models
 during the sampling stage makes the proposed approach more robust and less prone to collapse to
 trivial solutions.

We identify the ASR guidance models trained on WSJ with ID $1 \sim 12$. Their details can be found in the Appendix.A.1. We gradually increased the number of guidance ASR models from 1 to 12, the WER can be significantly reduced from 81.3% to 10.1%. It is worth noting that all the guidance ASR models are trained on the same 80-hour WSJ training set. Although we used guidance models up to 12, each model's parameter scale is deliberately kept relatively small for efficiency. When using 6 models as guidance, their total parameters (89.7M) do not exceed those of the unconditional score model (147.4M), and the word error rates can be reduced from 81.3% to 14.9%.

6.2 EXPERIMENTAL RESULTS ON SPEAKER GUIDANCE

Table 2: Objective evaluation for speaker guidance. SIM is the speaker similarity between the generated speech and enrollment speech.

Guidance	# total guide param. (M)	WER \downarrow	SpeechBERT \uparrow	UTMOS \uparrow	SIM \uparrow
ASR {1}	14.9	81.3	0.64	2.84	0.16
+ 1 Spk. model	21.6	83.1	0.62	2.86	0.29
+ 2 Spk. model	28.3	80.0	0.63	2.87	0.37
ASR {1,2,3}	44.8	24.4	0.71	3.28	0.16
+ 1 Spk. model	51.5	23.7	0.71	3.41	0.31
+ 2 Spk. model	58.2	23.5	0.71	3.37	0.38
ASR {1,,6}	89.7	14.9	0.72	3.42	0.15
+ 1 Spk. model	96.4	14.8	0.73	3.45	0.29
+ 2 Spk. model	103.1	14.9	0.73	3.48	0.38
ASR {1,, 12}	229.9	10.1	0.74	3.47	0.12
+1 Spk. model	236.6	9.0	0.74	3.47	0.27
+2 Spk. model	243.3	10.3	0.74	3.46	0.34
Ground Truth	-	3.3	1.00	4.15	0.59

We apply the target speaker guidance by following Eq. 12 and Eq. 13. The enrollment speech X^c is randomly picked from other speech samples from the same speaker of the ground truth speech in the LibriTTS-R testing set. The speaker similarity evaluated between the ground truth and their enrollment speech is 0.59. We compare the speaker guidance with one and two guidance speaker models. The results are listed in Table.2. In the speaker guidance experiments, we found that only using one speaker model as guidance can clearly improve speaker similarity by introducing speaker guidance, most of the systems slightly reduce the WER, while some of them get a slightly worse WER. No obvious conflicts between the speaker and ASR guidance are observed, which is consistent with the independent assumption used in Eq.11 and Eq.12. When more ASR guidance models were used (12), speaker similarity improved less, which can be observed in the results of 12 ASR guidance.

6.3 RESULTS ON MEAN OPINION SCORE EVALUATIONS

We report the MOS evaluation results in Table 3. We get M-MOS and N-MOS scores around 3.9,
which still has a gap between the TTS speech and the real speech (~ 4.9). Relative comparisons
between different guidance are consistent with those in Sec.6.1 and Sec.6.2. The 12-ASR guidance
setup shows better M-MOS and N-MOS than the 6-ASR guidance but gets worse S-MOS than
the latter one. However, both systems equipped with speaker guidance show significant S-MOS improvement when compared to the system without speaker guidance.



Table 3: Mean Opinion Score (MOS) for the Guided TTS.

Figure 1: WER(%) w.r.t. the target speech length. λ is the scaling factor to the ground truth length. Where $\lambda = 1$ is identical to the 12-ASR guidance in Table. 1.

6.4 Ablation Studies on Speech Duration

In our proposed approach, the E2EASRs are leveraged to control speech generation for target text. E2EASRs learn an implicit alignment between the speech and text labels in their training stage. When using E2EASRs as guidance to control the TTS generation, we do not explicitly control the speed of speech or duration of phonemes. Our method only needs the desired total speech length, target text, and optional enrollment speech as input. In our previous experiments, we did not investigate the total length of the target speech in detail but simiply adopted the ground truth length as the target length. In real applications that do not have a ground truth length, an algorithm is needed to predict the target speech length from the text.

To understand the effect that the defects of the length prediction algorithm may have on the quality of TTS, we conducted an ablation experiment on the total length of speech. We scale the length of ground truth speech by a factor λ , and use it as the target length in the sampling. The guidance models are all the 12 WSJ E2EASR models. The curves of WER with respect to λ are plotted in Fig 2.

As Fig.2 shows, the proposed methods are generally more sensitive to short speech length. If the total length of speech generation is shortened, the WER will increase. However, if the length of speech is longer than the ground truth in a reasonable range (from 1.0 to 1.5), the WER evaluation becomes on par or even better. A possible explanation for this phenomenon is that the LibriTT-R training data are well-segmented, and most of the onset and offset silence audio are clipped out. On the other hand, the WSJ ASR training data has more silence on both the onset and offset of the speech. So, the ground truth speech length in LibriTTS-R may be shorter for the distribution learned by the ASR guidance model. We give more detailed examples in the Appendix. A.3, the model is trying to generate onset and offset silence if the target length is set too long.

431 This finding can guide the design of length prediction algorithms in real applications. For example, the predictor can be biased to output a longer length than the ground truth.

⁴³² 7 DISCUSSIONS AND CONCLUSION

433

434 In this paper, we proposed to use E2EASR models to guide an unconditional score-based generative 435 model (SGM) and enable TTS. The unconditional SGM can be trained on large-scare unlabeled 436 speech data. We show that using ASR models trained only on 80-hour can guide the unconditional 437 SGM to generate high-quality speech. Meanwhile, using speaker verification models as guidance, we can also conduct zero-shot TTS with a few seconds of the target speaker's enrollment speech. 438 By utilizing the end-to-end training for the guidance ASR, We found that we can synthesize high-439 quality speech without using the phoneme duration model. An ablation study shows that our model 440 is robust against the mismatch of total target speech length within a certain range. 441

⁴⁴² The main limitations of our work and possible extensions include:

1) Multiple ASR guidance system are currently needed to generate satisfactory-quality speech, and our best system uses 12 ASR models for guidance. We use deliberately compact, small guidance ASR models (≈ 15 M parameters), thus in part alleviating the computational overhead in the inference stage. Our future works will focus on reducing the number of guidance models. Possible solutions include optimizing the guidance ASR model against adversarial attacks to improve the robustness of guidance.

2) Although we disentangle the training data of the unconditional SGM and the guidance ASR model, the guidance model must be trained with the data perturbed by the diffusion SDE. This prevents arbitrary pre-trained models from being directly used as the guidance model. In the future, it is necessary to explore using guidance models that do not require data-perturbed training but e.g. only fine-tuning with perturbed data.

454
3) Our proposed E2EASR guidance for TTS does not require commonly used phoneme duration
prediction models. We can generate speech with the input text and a total duration of the target
length. In our experiments, we empirically set the latter to the length of the ground truth. In future
work, the problem of generating natural speech robustly at a total input length needs to be explored.

4) This work focuses more on the utilization and optimization of guidance models but less on the
design of unconditional SGM. An important research direction is how to design an unconditional
SGM that is more effective when being guided.

462 5) Another potential extension for this work is cross-lingual zero-shot TTS. Using ASR models
463 trained on a target language may be still be able to guide an unconditional SGM trained on single
464 or multiple unlabeled speech data from other languages. This may solve the problem of lack of
465 high-quality minority languages or dialects TTS annotation.

466 467

468

References

- Syed Ammar Abbas, Thomas Merritt, Alexis Moinet, Sri Karlapati, Ewa Muszynska, Simon Slangen, Elia Gatti, and Thomas Drugman. Expressive, Variable, and Controllable Duration Modelling in TTS. In *Proc. Interspeech 2022*, pp. 4546–4550, 2022. doi: 10.21437/Interspeech. 2022-384.
- Brian D. O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, May 1982. ISSN 0304-4149. doi: 10.1016/0304-4149(82)90051-5.
- Dzmitry Bahdanau, Jan Chorowski, Dmitriy Serdyuk, Philémon Brakel, and Yoshua Bengio. End-to-end attention-based large vocabulary speech recognition. In 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4945–4949, March 2016. doi: 10.1109/ICASSP.2016.7472618.
- W. Chan, N. Jaitly, Q. Le, and O. Vinyals. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4960–4964, March 2016. doi: 10.1109/ICASSP.2016.7472621.
- Yung-Sung Chuang, Chi-Liang Liu, Hung-yi Lee, and Lin-shan Lee. SpeechBERT: An Audio-and-Text Jointly Learned Language Model for End-to-End Spoken Question Answering. In *Proc. Interspeech 2020*, pp. 4168–4172, 2020. doi: 10.21437/Interspeech.2020-1570.

524

525

526

527

486	Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. VoxCeleb2: Deep Speaker Recognition.
487	In <i>Proc. Interspeech 2018</i> , pp. 1086–1090, 2018. doi: 10.21437/Interspeech.2018-1929.
488	

- Johanes Effendi, Yogesh Virkar, Roberto Barra-Chicote, and Marcello Federico. Duration Modeling
 of Neural TTS for Automatic Dubbing. In *ICASSP 2022 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 8037–8041, May 2022. doi: 10.1109/
 ICASSP43922.2022.9747158.
- 493
 494
 494
 495
 Garofolo, John S., Graff, David, Paul, Doug, and Pallett, David. CSR-I (WSJ0) Complete, May 2007.
- A. Graves, A. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 6645–6649, May 2013. doi: 10.1109/ICASSP.2013.6638947.
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, pp. 369–376, New York, NY, USA, June 2006. Association for Computing Machinery. ISBN 978-1-59593-383-6. doi: 10.1145/1143844.1143891.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc., 2020.
- Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren. Prodiff: Progressive fast diffusion model for high-quality text-to-speech. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 2595–2605, 2022.
- Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score match *Journal of Machine Learning Research*, 6(4), 2005.
- Myeonghun Jeong, Hyeongju Kim, Sung Jun Cheon, Byoung Jin Choi, and Nam Soo Kim. Diff-tts:
 A denoising diffusion model for text-to-speech. *arXiv preprint arXiv:2104.01409*, 2021.
- Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-TTS: A Diffusion Model for Text-to-Speech via Classifier Guidance. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 11119–11133. PMLR, June 2022a.
- Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-TTS: A Generative Flow for
 Text-to-Speech via Monotonic Alignment Search. In *Advances in Neural Information Processing Systems*, volume 33, pp. 8067–8077. Curran Associates, Inc., 2020.
 - S. Kim, T. Hori, and S. Watanabe. Joint CTC-attention based end-to-end speech recognition using multi-task learning. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4835–4839, March 2017. doi: 10.1109/ICASSP.2017.7953075.
- Sungwon Kim, Heeseung Kim, and Sungroh Yoon. Guided-TTS 2: A Diffusion Model for High quality Adaptive Text-to-Speech with Untranscribed Data, May 2022b.
- 531 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint* 532 *arXiv:1412.6980*, 2014.
- Yuma Koizumi, Heiga Zen, Shigeki Karita, Yifan Ding, Kohei Yatabe, Nobuyuki Morioka, Michiel
 Bacchiani, Yu Zhang, Wei Han, and Ankur Bapna. LibriTTS-R: A Restored Multi-Speaker Textto-Speech Corpus. In *Proc. Interspeech 2023*, pp. 5496–5500, 2023. doi: 10.21437/Interspeech. 2023-1584.
- Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catanzaro, and Sungroh Yoon. BigVGAN: A Universal Neural Vocoder with Large-Scale Training. In *The Eleventh International Conference on Learning Representations*, September 2022.

554

- 540 Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. Neural speech synthesis with 541 transformer network. In Proceedings of the AAAI conference on artificial intelligence, volume 33, 542 pp. 6706-6713, 2019. 543
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An ASR 544 corpus based on public domain audio books. In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5206–5210, April 2015. doi: 10.1109/ICASSP.2015. 546 7178964. 547
- Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-548 tts: A diffusion probabilistic model for text-to-speech. In International Conference on Machine 549 Learning, pp. 8599-8608. PMLR, 2021. 550
- Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech: 552 Fast, robust and controllable text to speech. Advances in neural information processing systems, 553 32, 2019a.
- Yi Ren, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Almost unsupervised text to 555 speech and automatic speech recognition. In *International conference on machine learning*, pp. 556 5410-5419. PMLR, 2019b.
- 558 Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech 559 2: Fast and high-quality end-to-end text to speech. In International Conference on Learning Representations, 2020. 560
- 561 Julius Richter, Simon Welker, Jean-Marie Lemercier, Bunlong Lay, and Timo Gerkmann. Speech 562 Enhancement and Dereverberation With Diffusion-Based Generative Models. IEEE/ACM Trans-563 actions on Audio, Speech, and Language Processing, 31:2351-2364, 2023. ISSN 2329-9290, 564 2329-9304. doi: 10.1109/TASLP.2023.3285241.
- Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi 566 Saruwatari. UTMOS: UTokyo-SaruLab System for VoiceMOS Challenge 2022. In Proc. Inter-567 speech 2022, pp. 4521–4525, 2022. doi: 10.21437/Interspeech.2022-439. 568
- 569 Takaaki Saeki, Soumi Maiti, Shinnosuke Takamichi, Shinji Watanabe, and Hiroshi Saruwatari. SpeechBERTScore: Reference-Aware Automatic Evaluation of Speech Generation Leveraging 570 NLP Evaluation Metrics, September 2024. 571
- 572 Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words 573 with Subword Units. In Katrin Erk and Noah A. Smith (eds.), Proceedings of the 54th Annual 574 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1715– 575 1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/ v1/P16-1162. 576
- 577 Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, 578 Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. Natural tts synthesis by con-579 ditioning wavenet on mel spectrogram predictions. In 2018 IEEE international conference on 580 acoustics, speech and signal processing (ICASSP), pp. 4779–4783. IEEE, 2018. 581
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep Unsupervised 582 Learning using Nonequilibrium Thermodynamics. In Proceedings of the 32nd International Con-583 ference on Machine Learning, pp. 2256–2265. PMLR, June 2015. 584
- 585 Yang Song and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data Distribu-586 tion. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019.
- 588 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and 589 Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. In 9th 590 International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. 592
- Jaesung Tae, Hyeongju Kim, and Taesu Kim. Editts: Score-based editing for controllable text-tospeech. In Interspeech 2022, pp. 421-425, 2022. doi: 10.21437/Interspeech.2022-6.

- 594 Shuai Wang, Zhengyang Chen, Bing Han, Hongji Wang, Chengdong Liang, Binbin Zhang, Xu Xi-595 ang, Wen Ding, Johan Rohdin, Anna Silnova, et al. Advancing speaker embedding learning: 596 Wespeaker toolkit for research and production. Speech Communication, 162:103104, 2024. 597 Yixiang Wang, Jiqiang Liu, Xiaolin Chang, Ricardo J. Rodríguez, and Jianhua Wang. DI-AA: An 598 interpretable white-box attack for fooling deep neural networks. *Information Sciences*, 610:14– 32, September 2022. ISSN 0020-0255. doi: 10.1016/j.ins.2022.07.157. 600 601 Yuxuan Wang, R.J. Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J. Weiss, Navdeep Jaitly, 602 Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, Quoc Le, Yannis Agiomyrgiannakis, Rob Clark, and Rif A. Saurous. Tacotron: Towards end-to-end speech synthesis. In Interspeech 603 2017, pp. 4006–4010, 2017. doi: 10.21437/Interspeech.2017-1452. 604 605 Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson 606 Enrique Yalta Soplin, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala, 607 and Tsubasa Ochiai. ESPnet: End-to-End Speech Processing Toolkit. In Interspeech 2018, pp. 608 2207-2211. ISCA, September 2018. doi: 10.21437/Interspeech.2018-1456. 609 Hossein Zeinali, Shuai Wang, Anna Silnova, Pavel Matějka, and Oldřich Plchot. But system descrip-610 tion to voxceleb speaker recognition challenge 2019. arXiv preprint arXiv:1910.12592, 2019. 611 612 Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J. Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 613 LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech. In Proc. Interspeech 2019, 614 pp. 1526-1530, 2019. doi: 10.21437/Interspeech.2019-2441. 615 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. BERTScore: Eval-616 uating Text Generation with BERT. In International Conference on Learning Representations, 617 September 2019. 618 619 620 APPENDIX А 621 622 DETAILS OF ASR GUIDANCE MODELS A.1 623 624 A total of twelve ASR guidance models were trained on the WSJ dataset (referred to as WSJ-ASR 1 to 12). The WSJ-ASR models utilize a Transformer encoder-decoder architecture, with each model 625 featuring 2048 hidden units and 4 attention heads. 626 627 All ASR guidance models were trained using the Adam optimizer (Kingma & Ba, 2014), with an 628 initial learning rate of 0.001 and 5,000 warm-up steps. Each model was trained for 100 epochs. We 629 employed a joint CTC-attention training framework (Kim et al., 2017), where the loss weights for the CTC and attention objectives were empirically set to 0.3 and 0.7, respectively. We use byte-pair 630 encoding (BPE) tokens for all models. Further architectural details of the ASR guidance models are 631 provided in Table 4. 632 633 A.2 INSTRUCTIONS FOR SUBJECTIVE EVALUATIONS 634 635 N-MOS: Your task is to judge the Naturalness of the speech you hear in relation to the reference 636 speech. Please concentrate on the fluidity and naturality of the interaction as well as the expres-637 siveness of the speakers regardless of meaning. 638 M-MOS: Your task is to judge the **Meaningfulness** of the speech you hear in relation to the reference 639 speech. Please focus on whether the sequence of words is identical to the reference speech. 640 641 S-MOS: Your task is to judge the Speaker Similarity of the speech you hear in relation to the 642 reference speech. Please concentrate on the speaker similarity regardless of speech quality, nat-643 uralness, and meaning. 644 645 A.3 ABLATION STUDIES ON TARGET LENGTH 646
- Figure 2 shows two speech syntheses with a longer or shorter length than the ground truth length. Although our methods do not include an explicit duration model, they are still robust in generating

648						
649	Table 4: Details of ASR models					
650		# token	# encoder lavers	# decoder lavers	# param (M)	
651		" token	" encoder layers	" decoder layers		
652	1	100	6	3	14.9	
653	2	200	6	3	14.9	
655	3	300	6	3	15.0	
054	4	100	6	3	14.9	
655	5	200	6	3	14.9	
656	6	300	6	3	15.0	
657	7	100	8	4	19.1	
658	8	200	8	4	19.2	
659	9	300	8	4	19.2	
660	10	100	12	6	27.5	
661	11	200	12	6	27.6	
662	12	300	12	6	27.7	

speech based on the target length. For the longer length, it will try to speak slower and generate silence at the beginning and end of the speech; for the shorter length, it will try to talk faster. It is worth noting that our method has no control over whether to speak slower or add silence at the beginning or end of the audio when the target speech length is too long. The specific approach to control it needs to be studied in the future. Our current model is not good at producing perfect silence because the unconditional diffusion model's training data is well segmented and contains fewer data that start or end with a longer duration. If this problem is solved in future work, the model can be more robust against longer input speech lengths.



Figure 2: Speech synthesized with 1.2 times or 0.85 times length of the ground truth length.