

---

# Controllable and Lossless Non-Autoregressive End-to-End Text-to-Speech

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Some recent studies have demonstrated the feasibility of single-stage neural text-  
2       to-speech, which does not need to generate mel-spectrograms but generates the  
3       raw waveforms directly from the text. Single-stage text-to-speech often faces two  
4       problems: a) the one-to-many mapping problem due to multiple speech variations  
5       and b) insufficiency of high frequency reconstruction due to the lack of supervision  
6       of ground-truth acoustic features during training. To solve the a) problem and  
7       generate more expressive speech, we propose a novel phoneme-level prosody  
8       modeling method based on a variational autoencoder with normalizing flows  
9       to model underlying prosodic information in speech. We also use the prosody  
10      predictor to support end-to-end expressive speech synthesis. Furthermore, we  
11      propose the dual parallel autoencoder to introduce supervision of the ground-truth  
12      acoustic features during training to solve the b) problem enabling our model to  
13      generate high-quality speech. We compare the synthesis quality with state-of-  
14      the-art text-to-speech systems on an internal expressive English dataset. Both  
15      qualitative and quantitative evaluations demonstrate the superiority and robustness  
16      of our method for lossless speech generation while also showing a strong capability  
17      in prosody modeling.

## 18   1 Introduction

19   With the rapid development of deep learning, neural text-to-speech (TTS) systems can generate  
20   natural and high-quality speech and thus have drawn much attention in the machine learning and  
21   speech community. TTS is a task that aims at synthesizing raw speech waveforms from the given  
22   source text. Most previous neural TTS systems' pipelines are two-stage. The first stage is to generate  
23   intermediate speech representations (e.g., mel-spectrograms) autoregressively [36, 29, 24, 19] or  
24   non-autoregressively [27, 26, 14] from input text. The second stage is to synthesize speech waveforms  
25   from the generated intermediate speech representations using a vocoder [10, 22, 25, 38, 33]. These  
26   systems with two-stage pipelines can synthesize high-quality speech but still have drawbacks because  
27   they need sequential training or fine-tuning [15]. In addition, the use of predefined intermediate  
28   representations prevents further improvement in overall performance, as two system components can  
29   not be jointly trained and connected by learned intermediate representations.

30   Recently, several works (FastSpeech 2s [26], EATS [7], VITS [15]) have proposed parallel end-to-end  
31   TTS models that generate raw waveforms directly from input text in a single stage. FastSpeech 2s  
32   introduces explicit pitch and energy as mel-spectrogram decoder conditions to alleviate the one-to-  
33   many mapping problem in the TTS system. However, it needs to extract these handcrafted features in

34 advance, complicating the training pipeline. Moreover, FastSpeech 2s only models pitch and energy  
35 but does not disentangle other prosody features from the speech. EATS and VITS can synthesize  
36 high-quality speech, but they do not disentangle prosody information from speech, so they can not  
37 achieve prosody modeling and control.

38 To solve the problem that previous single-stage parallel end-to-end TTS models do not model the  
39 general prosody, we propose the **C**ontrollable and **L**Ossless **N**on-autoregressive **E**nd-to-end TTS  
40 (**CLONE**), which contains some carefully designed components to disentangle and model the general  
41 prosody from speech.

42 Firstly, to better solve the one-to-many mapping problem, we need to model the information variance  
43 other than text in speech. We propose an implicit phoneme-level prosody latent variable modeling  
44 instead of only explicitly modeling pitch and energy in FastSpeech 2s. The phoneme-level prosody  
45 latent variable models general prosody in the speech in a unified way without supervision. Specifically,  
46 we assume that prosody follows a normal distribution and use a variational autoencoder [16] (VAE)  
47 with normalizing flows [28, 6] to model it, which enhances the modeling ability of pure VAE and  
48 enables better modeling of prosody information that has extremely high variance. We propose a  
49 prosody predictor to predict the prior distribution of phoneme-level prosody latent variable from the  
50 input phoneme, which enables end-to-end synthesis as a TTS system.

51 In addition, we carefully study the problem of unsatisfactory high-frequency information generation  
52 in single-stage end-to-end speech synthesis, which is caused by the lack of the supervision of ground-  
53 truth acoustic features during training. To enhance the learned intermediate representation, we propose  
54 the dual parallel autoencoder (DPA) that consists of two parallel encoders (the acoustic encoder  
55 and the posterior wave encoder) and a wave decoder. DPA uses ground-truth linear spectrograms  
56 to regularize the learned intermediate representations for efficient learning. Besides, we introduce  
57 the multi-band discriminator (MBD) that significantly speeds up model convergence and improves  
58 generation quality. DPA and MBD enable CLONE to synthesize high-quality speech at the lossless  
59 high sampling rate (48 kHz) with better high-frequency information.

60 We conduct experiments on our private speech datasets. The results of extensive evaluations show that  
61 CLONE outperforms SOTA two-stage and single-stage TTS models [29, 26, 15] in terms of speech  
62 quality. In addition, CLONE can synthesize lossless speech at 48 kHz with better speech quality.  
63 Furthermore, we demonstrate that CLONE can model and control prosody by the phoneme-level  
64 prosody latent variable and generate speech with appropriate prosodic inflections. We attach audio  
65 samples generated by CLONE at <https://cloneneurips2022.github.io/CLONE/>.

## 66 2 Related Work

67 **Text-to-Speech** Text-to-Speech (TTS), which aims to synthesize intelligible and natural speech  
68 waveforms from the given text, has attracted much attention in recent years. Specifically, the neural  
69 network-based TTS models [36, 29, 27, 26, 22] have achieved tremendous progress. The quality of  
70 the synthesized speech is improved a lot and is close to that of the human counterpart. The previous  
71 prevalent methods are two-stage. The general pipeline of two-stage methods is: first, generate the  
72 acoustic features (e.g., mel-spectrograms) from text autoregressively [36, 29, 24, 19, 34] or non-  
73 autoregressively [27, 26, 14, 23], then synthesize the raw waveforms conditioned on the acoustic  
74 features [10, 22, 25, 18]. Recently, several single-stage end-to-end TTS models [26, 7, 15] have been  
75 proposed to generate raw waveform directly from the text. Among them, VITS [15] outperforms the  
76 two-stage models due to the advantages of learned intermediate speech representations obtained by  
77 fully end-to-end training. However, these single-stage methods have poor controllability over the  
78 prosody of synthesized speech. Specifically, EATS [7] and VITS cannot control prosody. FastSpeech  
79 2s [26] can only control some pre-defined prosodic features (i.e., pitch and energy), and these features  
80 are required to be extracted in advance. Unlike the aforementioned single-stage models, by utilizing  
81 a conditional VAE with normalizing flows, CLONE achieves high controllability over the general  
82 phoneme-level prosody of synthesized speech.

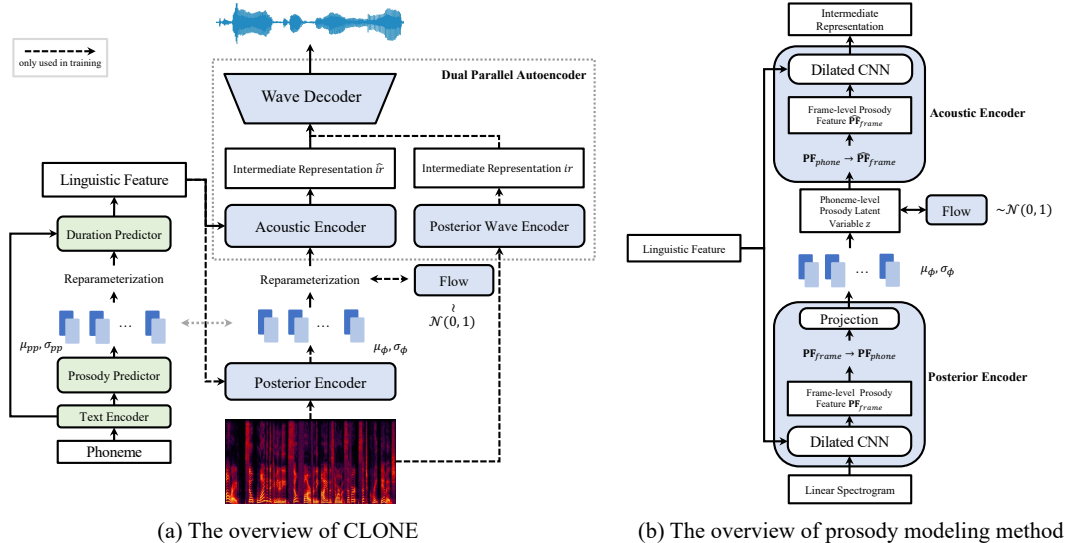


Figure 1: The overview of CLONE and prosody modeling method.

83 **Prosody Modeling** Many previous works have focused on learning underlying non-textual infor-  
 84 mation (e.g., style and prosody) in speech. In particular, some works [30, 37] introduce a reference  
 85 embedding to model style and prosody. [30] extracts a prosody embedding from a reference spec-  
 86 trogram, and [37] models a reference embedding as a weighted combination of a bank of learned  
 87 embeddings. Some works [1, 39] use the variational autoencoder (VAE) to model latent representa-  
 88 tions for styles and prosody of speech. Specifically, [32] uses multi-level VAE to model fine-grained  
 89 prosody at the phoneme and word level in an autoregressive way. [39] and [12] integrate VAE  
 90 with Tacotron 2 for better style modeling. Some works [8, 11] use GMM based mixture density  
 91 network to model prosodic information at phoneme and word levels. Unlike the previous models,  
 92 CLONE adopts a conditional VAE with normalizing flows for the phoneme-level prosody modeling  
 93 to a single-stage parallel TTS system. Inspired by [28, 5, 40] that improve the expressive capability  
 94 of prior and posterior distribution with normalizing flows, we add normalizing flows to enhance  
 95 the representation power of our prior distribution for better prosody modeling. Furthermore, the  
 96 prosody predictor enables CLONE to predict the prior distribution of prosody directly from the text  
 97 and control the prosody of generated speech during inference without the need for manually adjusting  
 98 the sampling points [39] or a reference speech [37] or other modality input (e.g., video [13]).

### 99 3 CLONE

#### 100 3.1 Overview

101 The overall model structure of CLONE shown in Figure 1a can be regarded as a conditional VAE.  
 102 Firstly, the posterior encoder converts the input spectrograms to a sequence of phoneme-level prosody  
 103 latent variable  $z$ . Then, the acoustic encoder transforms the prosody latent variable  $z$  into the learnable  
 104 intermediate representations conditioning on linguistic features. Finally, the wave decoder predicts the  
 105 waveforms from the learnable intermediate representations. The objective of CLONE is to maximize  
 106 the evidence lower bound (ELBO) of the intractable marginal log-likelihood of data  $\log_{\theta}(x | c)$ :

$$\text{ELBO} = \mathbb{E}_{q_{\phi}(z|x,c)} [\log p_{\theta}(x | z, c)] - D_{kl}(q_{\phi}(z | x, c) || p_{\theta}(z | c)), \quad (1)$$

107 where  $c$  and  $z$  denote the linguistic feature and the phoneme-level prosody latent variable respectively,  
 108  $q_{\phi}(z | x, c)$  is the approximate posterior distribution of  $z$  given a data point  $x$  and condition  $c$ ,  
 109  $p_{\theta}(x | z, c)$  is the likelihood function of  $x$  given  $z$  and  $c$ , and  $p_{\theta}(z | c)$  is the prior distribution of  $z$   
 110 given  $c$ .

111 The training loss of CLONE is the negative ELBO, which consists of the reconstruction loss  
 112 ( $-\mathbb{E}_{q_\phi(z|x,c)}[\log p_\theta(x|z,c)]$ ) and the KL divergence ( $D_{kl}(q_\phi(z|x,c)||p_\theta(z|c))$ ). The details  
 113 of the reconstruction loss and the KL divergence are described in Section 3.4.2 and Section 3.2,  
 114 respectively.

### 115 3.2 Prosody Modeling

116 To model the prosody of speech more appropriately, we determine to model the phoneme-level  
 117 prosody rather than the frame-level prosody or the word-level prosody. Because the frame-level  
 118 prosody causes severe linguistic information leakages during training, and the granularity of the  
 119 word-level prosody is too large to reflect the details of the prosody well. Since the input linear  
 120 spectrogram is at the frame level, we need to convert the frame-level prosody feature to the phoneme  
 121 level. We use the obtained duration of phonemes to construct the hard alignment matrix representing  
 122 the correspondence between phonemes and spectrogram frames and convert the frame-level prosody  
 123 feature to the phoneme-level prosody feature by the matrix as follows:

$$\mathbf{PF}_{phone} = \text{diag}(\mathbf{s}) \cdot \mathbf{A} \cdot \mathbf{PF}_{frame}, \quad (2)$$

124 where  $\mathbf{PF}_{phone} \in \mathbb{R}^{T_p \times d}$  and  $\mathbf{PF}_{frame} \in \mathbb{R}^{T_f \times d}$  denote the phoneme-level prosody feature and  
 125 the frame-level prosody feature, respectively,  $\mathbf{A} \in \mathbb{R}^{T_p \times T_f}$  denotes the hard alignment matrix, and  
 126  $\mathbf{s} \in \mathbb{R}^{T_p}$  denotes the inverse of the duration of phonemes ( $s_i = 1 / \sum_{j=1}^{T_f} a_{ij}$ ).

127 The mean and standard deviation of the approximate posterior distribution  $q_\phi(z|x,c)$  are predicted  
 128 from the obtained phoneme-level prosody feature and linguistic feature. We obtain the linguistic  
 129 feature  $c$  by expanding the output of the text encoder according to the phoneme duration. The length  
 130 of the linguistic feature is the same as the number of frames in the spectrogram.

131 The formula for KL divergence is as follows:

$$\mathcal{L}_{kl} = D_{kl}(q_\phi(z|x,c)||p_\theta(z|c)) = \mathbb{E}_{q_\phi(z|x,c)}[\log q_\phi(z|x,c) - \log p_\theta(z|c)]. \quad (3)$$

132 Unlike traditional VAE, we assume that the approximate posterior distribution of the phoneme-level  
 133 prosody latent variable  $z$  is a normal distribution rather than a standard normal distribution, i.e.,  
 134  $q_\phi(z|x,c) = \mathcal{N}(z; \mu_\phi(x,c), \sigma_\phi(x,c))$ . As a TTS model, we want to control the phoneme-level  
 135 prosody explicitly. If  $z$  follows the standard normal distribution, the prosody variation is implicitly  
 136 determined in the random sampling process, which is not desired. Thus, a normal distribution with  
 137 variable mean and variance is a better choice. It enables the prosody variation to be contained in the  
 138 mean and variance of normal distribution so that the corresponding prosody can be determined by  
 139 predicting the mean and variance. In addition, compared with standard normal distribution, normal  
 140 distribution is more complex, which increases the prosody modeling ability to obtain diverse prosody  
 141 variation.

142 We also need to increase the expressiveness of the prior distribution to match the posterior distribution.  
 143 Therefore, we apply normalizing flows, which enable an invertible transformation from a simple  
 144 standard normal distribution into a more complex prior distribution following the rule of change-of-  
 145 variables:

$$p_\theta(z|c) = \mathcal{N}(f_\theta(z); \mathbf{0}, \mathbf{I}) \left| \det \frac{\partial f_\theta(z)}{\partial z} \right|, \quad (4)$$

146 where  $f_\theta$  denotes the normalizing flow. After the reparameterization of VAE, we get the phoneme-  
 147 level prosody latent variable  $z$  which represents the phoneme-level prosody of the speech. We need  
 148 to convert phoneme-level  $z$  to frame-level variable  $\widehat{\mathbf{PF}}_{frame}$  to match the length of the linguistic  
 149 feature as follows:

$$\widehat{\mathbf{PF}}_{frame} = \mathbf{A}^\top \cdot z, \quad (5)$$

150 where  $\mathbf{A}^\top$  is the transposed matrix of  $\mathbf{A}$ .

### 151 3.3 Prosody Predictor

152 We propose a prosody predictor to model the correspondence between phoneme and phoneme-level  
 153 prosody. The prosody predictor can predict the mean and variance of  $q_\phi(z|x,c)$  from the output

154 of the text encoder. Therefore, CLONE can predict phoneme-level prosody from text input in the  
 155 inference stage. During inference, there are three modes to generate highly natural speech with  
 156 suitable prosody (details in Section 4.4). The optimization goal of the prosody predictor is to minimize  
 157 the KL divergence between the predicted normal distribution and the posterior distribution. Therefore,  
 158 the training loss of the prosody predictor is as follows:

$$\mathcal{L}_{pp} = D_{kl}(\mathcal{N}(\mu_{pp}, \sigma_{pp}), \mathcal{N}(\mu_{\phi}, \sigma_{\phi})), \quad (6)$$

159 where  $\mu_{pp}$  and  $\sigma_{pp}$  are the mean and standard deviation predicted by the prosody predictor. Compared  
 160 with FastSpeech 2 to directly predict pitch and energy, we predict the distribution of phoneme-level  
 161 prosody, avoiding the one-to-many mapping problem.

162 It is worth noting that the duration predictor we used is an improved version of the duration predictor  
 163 in FastSpeech 2 [26]. Since prosodic information partly determines the phoneme duration, we use  $z$   
 164 as the condition of the duration predictor besides the output of the text encoder, i.e.,  $\hat{\mathbf{d}} = \mathcal{DP}(z, t) \in$   
 165  $\mathbb{R}^{T_p}$ , where  $\mathcal{DP}$  denotes the duration predictor, and  $t$  is the output of the text encoder. In this way,  
 166 CLONE can generate more stable and natural phoneme durations.

### 167 3.4 Waveform Generation

#### 168 3.4.1 Motivation

169 Some end-to-end TTS studies [26, 7, 15] focus on generating raw waveforms directly from phonemes  
 170 recently. These single-stage TTS systems usually generate learned intermediate representations from  
 171 phonemes and then synthesize raw waveforms from the learned intermediate representations. Unlike  
 172 the mel-spectrogram used in two-stage methods, the intermediate representation in single-stage  
 173 methods is predicted by the model without training supervision, subject to prediction errors and  
 174 over-smoothness. The lack of supervision of intermediate representations during training expands the  
 175 search space of the single-stage model, resulting in the model being more challenging to optimize,  
 176 which is reflected in the poor modeling ability for high-frequency information in our experiments. To  
 177 narrow the search space of waveform generation in single-stage models, we introduce ground-truth  
 178 speech signals. We design an autoencoder architecture called Dual Parallel Autoencoder (DPA) to  
 179 regularize the learned intermediate representation. In addition, to further enhance the quality of the  
 180 generated waveforms, we propose a new discriminator called Multi-Band Discriminator (MBD).  
 181 MBD divides the waveform into multiple bands so that our model can separately supervise the  
 182 low-frequency and high-frequency parts of the audio, improving the overall quality of the synthesized  
 183 speech.

#### 184 3.4.2 Dual Parallel Autoencoder

185 In the dual parallel autoencoder, two parallel encoders, namely the acoustic encoder and the posterior  
 186 wave encoder, generate intermediate representation, and a dual training method is used to optimize  
 187 them. The acoustic encoder transforms  $z$  into the predicted intermediate representations  $\hat{ir}$  condition-  
 188 ing on linguistic features, and the posterior wave encoder transforms linear spectrograms into  
 189 the intermediate representations  $ir$ . The wave decoder acts as the decoder of DPA and generates the  
 190 waveform from both intermediate representations. In practice, we concatenate  $ir$  and  $\hat{ir}$  in the batch  
 191 dimension to get  $ir_{concat}$  and send  $ir_{concat}$  into the wave decoder to produce the waveform  $\hat{w}$ . For  
 192 dual training, we compute the mean absolute error (MAE)  $\mathcal{L}_{ir}$  between  $ir$  and  $\hat{ir}$ , so that  $\hat{ir}$  gets  
 193 the supervision of ground-truth acoustic features from  $ir$ , and  $\hat{ir}$  is regularized by  $ir$ , which assists  
 194 CLONE in learning intermediate representation efficiently. Please note that we only use the acoustic  
 195 encoder without the posterior wave encoder during inference.

196 To calculate the reconstruction loss, we convert  $\hat{w}$  to mel-spectrogram  $m_1$  and calculate MAE with  
 197 ground-truth mel-spectrogram  $m_{gt}$ . To make  $ir$  and  $\hat{ir}$  focus on the information at the human  
 198 voice frequency band for better prosody modeling, we use a one-layer linear network to predict the  
 199 mel-spectrogram  $m_2$  from  $ir_{concat}$ . Therefore, the whole reconstruction loss  $\mathcal{L}_{recon}$  is:

$$\mathcal{L}_{recon} = \text{MAE}(m_{gt}, m_1) + \text{MAE}(m_{gt}, m_2). \quad (7)$$

### 200 3.4.3 Discriminator

201 We use the popular adversarial training approach as HiFi-GAN [17] to improve the resolution of  
202 synthesized speech. The following equations describe the loss of adversarial training:

$$\mathcal{L}_{advD} = \mathbb{E}_{(w, \hat{w})} [((D(w) - 1)^2 + (D(\hat{w}))^2)], \quad (8)$$

$$\mathcal{L}_{advG} = \mathbb{E}_w [(D(\hat{w}) - 1)^2] + \beta \mathcal{L}_{fm}, \quad \mathcal{L}_{fm} = \mathbb{E}_{(w, \hat{w})} \left[ \sum_{l=1}^L \frac{1}{N_l} \text{MAE}(D_l(w), D_l(\hat{w})) \right], \quad (9)$$

203 where  $\mathcal{L}_{advD}$ ,  $\mathcal{L}_{advG}$  and  $\mathcal{L}_{fm}$  denote the loss of the discriminator, the loss of the wave decoder, and  
204 the loss of the feature map, respectively.  $L$  denotes the total number of layers in the discriminators.  
205  $D_l$  and  $N_l$  denote the features and the number of features in the  $l$ -th layer of the discriminator,  
206 respectively.  $\beta$  is the coefficient of the feature map loss  $\mathcal{L}_{fm}$ , and we set it to 0.1 in our experiments.

207 We use two discriminators for joint adversarial training, namely MPD in HiFi-GAN and MBD. MBD  
208 uses the pseudo quadrature mirror filter bank (Pseudo-QMF) to divide the waveform into  $N$  sub-  
209 bands. These  $N$  sub-bands with one full-band are respectively sent to  $N + 1$  scale discriminators in  
210 MelGAN [18]. By applying different discriminators on different frequency bands of the synthesized  
211 audio, MBD can significantly enhance the generation quality of high-frequency parts, allowing  
212 CLONE to generate high-fidelity and lossless audio at a high sample rate. A similar idea is used  
213 in [21]. However, our method differs in that we send different frequency bands into different  
214 discriminators.

### 215 3.5 Loss Function

216 The total loss of CLONE can be expressed as follows:

$$\mathcal{L} = \lambda_1 * \mathcal{L}_{recon} + \lambda_2 * \mathcal{L}_{kl} + \lambda_3 * \mathcal{L}_{ir} + \lambda_4 * \mathcal{L}_{pp} + \lambda_5 * \mathcal{L}_{dur} + \lambda_6 * \mathcal{L}_{advG}, \quad (10)$$

217 where  $\lambda_{[1 \rightarrow 6]}$  represent coefficients of different components of the total loss.

## 218 4 Experiment

### 219 4.1 Dataset

220 We used proprietary English speech datasets, including a single-speaker dataset and a multi-speaker  
221 dataset. The single-speaker dataset contains 11,176 utterances with a total audio length of 10 hours  
222 at both 24 kHz and 48 kHz. We use 9,000 utterances as the training set, 100 utterances as the  
223 validation set, and the remaining data as the test set. The multi-speaker speech data contains five  
224 English speakers (two males and three females) with a total audio length of 22 hours. We evaluate the  
225 high-quality generation capability of CLONE on the single-speaker dataset and the prosody transfer  
226 capability of CLONE on the multi-speaker dataset.

### 227 4.2 Data Preprocessing

228 We convert the text sequences into phoneme sequences following [2, 19, 26]. In experiments, we use  
229 80-dimensional mel-spectrograms and linear spectrograms with 2048 filters. For the audio at 24 kHz,  
230 the hop size is 300, and the window size is 1200. For the audio at 48 kHz, the hop size is 300, and the  
231 window size is 2048. All use the Hann window.

### 232 4.3 Model Configuration

233 Our text encoder uses a stack of 6 feed-forward transformer blocks [35], and the prosody predictor  
234 consists of 4 WaveNet residual blocks (dilated CNN) [22], which consists of layers of dilated  
235 convolutions with a gated activation unit and skip connection. The duration predictor consists of a  
236 2-layer 1D convolutional network with ReLU activation, each followed by layer normalization [3]  
237 and the dropout layer [31], and an extra linear layer with ReLU activation to output a scalar, which is

238 the predicted phoneme duration. The posterior encoder consists of 8 blocks of dilated CNN, and the  
 239 normalizing flow is a stack of affine coupling layers [6] consisting of 4 blocks of dilated CNN. The  
 240 acoustic encoder is composed of 8 blocks of dilated CNN. The posterior wave encoder consists of 8  
 241 blocks of dilated CNN. The structure of the wave decoder is the same as HiFi-GAN V1. To match  
 242 our hop size, we change the upsampling rate to 5, 5, 4, 3 and change the upsampling kernel size to  
 243 15, 15, 12, 9. The detailed hyper-parameters of CLONE are listed in the appendix.

#### 244 4.4 Training and Inference

245 We train our model on 4 NVIDIA Tesla V100 32G GPUs with the batch size of 16 on each GPU for  
 246 500k steps. We use the AdamW [20] optimizer with  $\beta_1 = 0.8$  and  $\beta_2 = 0.99$ . The learning rate of  
 247 CLONE is fixed at  $2e-4$ . The discriminator uses the same optimization settings, with a fixed learning  
 248 rate of  $1e-4$ . As VITS, to reduce training time and GPU memory usage, we employ a windowed  
 249 generator for training, randomly sampling segments of intermediate representation with a window  
 250 size of 32 frames as input to the wave decoder.

251 CLONE has three modes of inference, and the difference lies in how the prosody latent variable is  
 252 calculated. (a) Use the prosody predictor to predict prosody information based on text information,  
 253 and the input of the model is only text information, just like a typical TTS model. (b) The prosody  
 254 latent variable is obtained through the inverse flow transformation, where the inputs of CLONE are  
 255 textual information and the sampling value of the standard normal distribution. (c) The prosody  
 256 information is extracted by the posterior encoder from the input linear spectrogram. In this mode, the  
 257 model inputs are text information and the reference linear spectrogram. Besides, we test the inference  
 258 speed of CLONE in mode (a) on an NVIDIA Tesla V100 32G GPU and compare it to that of VITS.  
 259 The RTF of CLONE and VITS are 0.012 and 0.017, respectively, indicating that our model has a  
 260 comparable inference speed to the SOTA single-stage parallel TTS model.

#### 261 4.5 Evaluation

##### 262 4.5.1 MOS Evaluation

263 We conduct the MOS (mean opinion score) evaluation to measure the synthesis quality of different  
 264 models. We use each model to synthesize the same 30 utterances<sup>1</sup>, and let 15 English native  
 265 speakers evaluate them. We compare with the state-of-the-art (SOTA) autoregressive TTS model  
 266 Tacotron 2 [29] with GMM-based attention mechanisms [4], the SOTA non-autoregressive TTS model  
 267 FastSpeech 2 [26], and the SOTA single-stage TTS model VITS [15]. The vocoder for Tacotron  
 268 2 and FastSpeech 2 is TFGAN [33], as TFGAN has better generation robustness than HiFi-GAN  
 269 in our experiments. The above models all generate audio at 24 kHz. We evaluate two kinds of  
 270 CLONE, namely 24 kHz CLONE and 24 kHz CLONE without MBD, of which the discriminators  
 271 are the same as those used in HiFi-GAN. The results of the MOS evaluation are shown in Table 1. It  
 272 can be seen that CLONE surpasses other SOTA models, indicating that our phoneme-level prosody  
 273 modeling method and the introduction of DPA enable the model to synthesize highly natural speech  
 274 with appropriate prosodic inflections. In addition, CLONE is better than CLONE without MBD,  
 275 demonstrating the effectiveness of MBD for generating higher quality speech.

Table 1: The MOS result with 95% confidence intervals of different models.

| Method               | MOS           | CI           |
|----------------------|---------------|--------------|
| Tacotron 2 + TFGAN   | 4.0717        | $\pm 0.0578$ |
| FastSpeech 2 + TFGAN | 4.0983        | $\pm 0.0579$ |
| VITS                 | 4.0108        | $\pm 0.0601$ |
| 24 kHz CLONE w/o MBD | 4.1000        | $\pm 0.0606$ |
| 24 kHz CLONE         | <b>4.1567</b> | $\pm 0.0528$ |

<sup>1</sup>To test the synthesis quality and robustness of the model and avoid data leakage, we use 30 general utterances outside the dataset for testing.

276 Furthermore, we conduct CMOS evaluations on 30 utterances with ground-truth recording audio.  
 277 We compare 24 kHz VITS, 24 kHz CLONE, and 48 kHz CLONE with 48 kHz ground-truth audio,  
 278 respectively, as shown in Table 2. We find that 48 kHz CLONE can generate the audio close to the 48  
 279 kHz ground-truth audio and outperform both the 24 kHz CLONE and 24 kHz VITS, demonstrating  
 280 that CLONE can synthesize high-fidelity 48 kHz audio.

Table 2: The CMOS comparison to evaluate speech generation quality at high sample rates. The audio synthesized by 48 kHz CLONE, 24 kHz CLONE, and 24 kHz VITS is compared with 48 kHz ground-truth audio, respectively.

| Method                    | CMOS           |
|---------------------------|----------------|
| 48 kHz Ground-truth Audio | 0              |
| 48 kHz CLONE              | <b>-0.1712</b> |
| 24 kHz CLONE              | -0.2393        |
| 24 kHz VITS               | -0.3621        |

281 **4.5.2 Prosody Modeling**

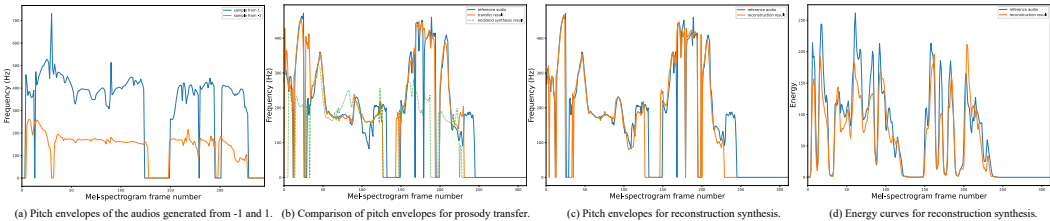


Figure 2: Visualization of prosody modeling. (a) shows the variation in the pitch of the audio generated from all  $-1$  and all  $1$  sampling, respectively. (b) shows the pitch comparison of the reference audio, transfer result, and end-to-end synthesis result. (c) and (d) show the pitch and energy of the reconstruction result and reference audio, respectively.

282 **Prosody Variation** We infer CLONE in mode (b), i.e., sampling values in the standard normal  
 283 distribution and obtaining the phoneme-level prosody latent variable through the inverse flow trans-  
 284 formation. We find that the prosody of the generated speech has a very high variance by sampling  
 285 different values in the standard normal distribution, which further proves the prosody modeling  
 286 capability of CLONE. Figure 2a visualizes the variation in the pitch of the audio generated by setting  
 287 sample values to all  $-1$  and all  $1$ , respectively. It can be seen that a significant prosody variation  
 288 can be achieved by adjusting sample values in the standard normal distribution of the inverse flow  
 289 transformation.

290 **Prosody Transfer** We test the prosody transfer performance of CLONE. Firstly, we train a multi-  
 291 speaker CLONE<sup>2</sup>. Then we infer CLONE in mode (c). We input the reference speech of speaker 1 to  
 292 the posterior encoder and use the speaker embedding of speaker 2 and the duration of the reference  
 293 speech to synthesize the transfer result. The goal is to use the timbre of speaker 2 to synthesize audio  
 294 with the same prosody as the reference audio of speaker 1. We also use the end-to-end synthesis  
 295 result with the timbre of speaker 2 for comparison. In Figure 2b, the pitch envelope of the transfer  
 296 result is very similar to that of the reference audio and quite different from that of the end-to-end  
 297 synthesis result. Besides, we calculate the MAE on pitch and energy of the prosody transfer result  
 298 and the end-to-end synthesis result (both are synthesized using the duration of the reference audio)  
 299 with the reference audio as ground truth, as shown in Table 3. It can be seen that the MAE of the  
 300 prosody transfer result is significantly smaller than that of the end-to-end synthesis result. Above two  
 301 experiments prove the effectiveness of prosody transfer.

<sup>2</sup>To implement multi-speaker TTS, we add speaker embedding to prosody predictor, duration predictor, posterior encoder, and acoustic encoder.



Table 3: The mean absolute error of pitch and energy of different methods.

| Method                                      | Pitch MAE    | Energy MAE   |
|---|--------------|--------------|
| end-to-end synthesis                        | 79.38        | 35.91        |
| prosody transfer synthesis                  | <b>41.44</b> | <b>31.92</b> |
| reconstruction synthesis by CLONE with HAPE | <b>32.56</b> | <b>16.34</b> |
| reconstruction synthesis by CLONE with SAPE | 64.83        | 23.70        |

302 **Prosody Reconstruction** CLONE can extract the prosody of the reference audio through the  
 303 posterior encoder and then reconstruct the audio as a conditional VAE. We draw the pitch and energy  
 304 curves of the reference audio and the reconstruction result, as shown in Figure 2c and Figure 2d. We  
 305 find that the two curves in each figure are very similar, indicating that CLONE can accurately extract  
 306 and reconstruct the prosody (e.g., pitch and energy) of the reference audio.

### 307 4.5.3 Ablation Study

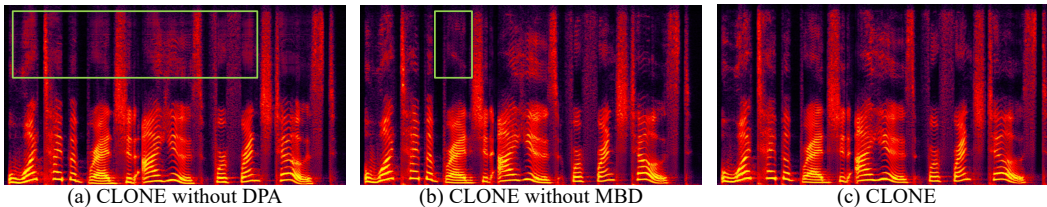


Figure 3: The spectrograms generated by three different CLONE, (a) CLONE without DPA, (b) CLONE without MBD, and (c) complete CLONE.

308 **Waveform Generation** To further investigate the improvement of our method on waveform gen-  
 309 eration, we conduct ablation experiments. Firstly, we conduct a MOS evaluation on the audio  
 310 synthesized by the CLONE with and without MBD, as shown in Table 1. We find that MBD enhances  
 311 the quality of synthesized audio. In addition, we plot the spectrograms of the synthesized audio, as  
 312 shown in Figure 3. Figure 3a shows that without the DPA, the synthesized audio suffers obvious  
 313 over-smoothness at high frequency, and Figure 3b shows that without MBD, the high-frequency  
 314 details of synthesized audio are insufficient. These demonstrate that DPA significantly weakens the  
 315 over-smoothness of the high frequency, and MBD further enhances the high-frequency details.

316 **Prosody Extraction** CLONE uses a hard alignment prosody extractor (HAPE) to extract prosody,  
 317 which improves the accuracy of prosody extraction. To verify this, we compare HAPE with the soft  
 318 attention prosody extractor (SAPE) [9, 32] where the text encoder output is to query the ground-truth  
 319 linear spectrogram by soft attention. We compute the MAE on pitch and energy of the above two  
 320 methods in Table 3. It can be seen that the MAE of HAPE is smaller than that of SAPE, indicating  
 321 that HAPE can extract prosody more accurately.

## 322 5 Conclusion

323 In this work, we propose a single-stage TTS system called CLONE that can directly generate lossless  
 324 waveforms from the text in parallel. Specifically, we design a phoneme-level prosody modeling  
 325 method based on a variational autoencoder with normalizing flows and a prosody predictor to solve the  
 326 one-to-many mapping problem better and support end-to-end expressive speech synthesis. Besides,  
 327 the dual parallel autoencoder is introduced to solve the problem of lacking supervision of ground-truth  
 328 acoustic features during training, which allows the single-stage model to generate lossless speech.  
 329 Our experiments demonstrate that CLONE outperforms existing SOTA single-stage and two-stage  
 330 TTS models in speech quality while performing strong controllability over prosody.

331 **References**

- 332 [1] Kei Akuzawa, Yusuke Iwasawa, and Yutaka Matsuo. Expressive speech synthesis via modeling  
333 expressions with variational autoencoder. *Proc. Interspeech 2018*, pages 3067–3071, 2018.
- 334 [2] Sercan Ö Arik, Mike Chrzanowski, Adam Coates, Gregory Diamos, Andrew Gibiansky, Yong-  
335 guo Kang, Xian Li, John Miller, Andrew Ng, Jonathan Raiman, et al. Deep voice: Real-time  
336 neural text-to-speech. In *International Conference on Machine Learning*, pages 195–204.  
337 PMLR, 2017.
- 338 [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*  
339 *arXiv:1607.06450*, 2016.
- 340 [4] Eric Battenberg, R.J. Skerry-Ryan, Soroosh Mariooryad, Daisy Stanton, David Kao, Matt  
341 Shannon, and Tom Bagby. Location-relative attention mechanisms for robust long-form speech  
342 synthesis. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and*  
343 *Signal Processing (ICASSP)*, pages 6194–6198, 2020.
- 344 [5] Xi Chen, Diederik P Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya  
345 Sutskever, and Pieter Abbeel. Variational lossy autoencoder. *arXiv preprint arXiv:1611.02731*,  
346 2016.
- 347 [6] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In  
348 *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April*  
349 *24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 350 [7] Jeff Donahue, Sander Dieleman, Mikolaj Binkowski, Erich Elsen, and Karen Simonyan. End-  
351 to-end adversarial text-to-speech. In *International Conference on Learning Representations*,  
352 2020.
- 353 [8] Chenpeng Du and Kai Yu. Rich prosody diversity modelling with phone-level mixture density  
354 network. *arXiv preprint arXiv:2102.00851*, 2021.
- 355 [9] Isaac Elias, Heiga Zen, Jonathan Shen, Yu Zhang, Ye Jia, Ron J Weiss, and Yonghui Wu. Parallel  
356 tacotron: Non-autoregressive and controllable tts. In *ICASSP 2021-2021 IEEE International*  
357 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5709–5713. IEEE,  
358 2021.
- 359 [10] Daniel Griffin and Jae Lim. Signal estimation from modified short-time fourier transform. *IEEE*  
360 *Transactions on acoustics, speech, and signal processing*, 32(2):236–243, 1984.
- 361 [11] Yiwei Guo, Chenpeng Du, and Kai Yu. Unsupervised word-level prosody tagging for control-  
362 lable speech synthesis. In *ICASSP 2022-2022 IEEE International Conference on Acoustics,*  
363 *Speech and Signal Processing (ICASSP)*, pages 7597–7601. IEEE, 2022.
- 364 [12] Wei-Ning Hsu, Yu Zhang, Ron J Weiss, Heiga Zen, Yonghui Wu, Yuxuan Wang, Yuan Cao,  
365 Ye Jia, Zhifeng Chen, Jonathan Shen, et al. Hierarchical generative modeling for controllable  
366 speech synthesis. In *International Conference on Learning Representations*, 2018.
- 367 [13] Chenxu Hu, Qiao Tian, Tingle Li, Wang Yuping, Yuxuan Wang, and Hang Zhao. Neural dubber:  
368 Dubbing for videos according to scripts. In *Advances in Neural Information Processing Systems*,  
369 2021.
- 370 [14] Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-tts: A generative flow  
371 for text-to-speech via monotonic alignment search. *arXiv preprint arXiv:2005.11129*, 2020.
- 372 [15] Jaehyeon Kim, Jungil Kong, and Juhee Son. Conditional variational autoencoder with adversar-  
373 ial learning for end-to-end text-to-speech. In *International Conference on Machine Learning*,  
374 pages 5530–5540. PMLR, 2021.

- 375 [16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
376 *arXiv:1312.6114*, 2013.
- 377 [17] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks  
378 for efficient and high fidelity speech synthesis. *Advances in Neural Information Processing*  
379 *Systems*, 33:17022–17033, 2020.
- 380 [18] Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose  
381 Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron Courville. Melgan: Generative  
382 adversarial networks for conditional waveform synthesis. *arXiv preprint arXiv:1910.06711*,  
383 2019.
- 384 [19] Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. Neural speech synthesis  
385 with transformer network. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
386 volume 33, pages 6706–6713, 2019.
- 387 [20] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International*  
388 *Conference on Learning Representations*, 2018.
- 389 [21] Ahmed Mustafa, Nicola Pia, and Guillaume Fuchs. Stylemelgan: An efficient high-fidelity  
390 adversarial vocoder with temporal adaptive normalization. In *ICASSP 2021-2021 IEEE Interna-*  
391 *tional Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6034–6038.  
392 IEEE, 2021.
- 393 [22] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex  
394 Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative  
395 model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- 396 [23] Kainan Peng, Wei Ping, Zhao Song, and Kexin Zhao. Non-autoregressive neural text-to-speech.  
397 In *International conference on machine learning*, pages 7586–7598. PMLR, 2020.
- 398 [24] Wei Ping, Kainan Peng, Andrew Gibiansky, Sercan O Arik, Ajay Kannan, Sharan Narang,  
399 Jonathan Raiman, and John Miller. Deep voice 3: 2000-speaker neural text-to-speech. *Proc.*  
400 *ICLR*, pages 214–217, 2018.
- 401 [25] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. Waveglow: A flow-based generative network  
402 for speech synthesis. In *ICASSP 2019-2019 IEEE International Conference on Acoustics,*  
403 *Speech and Signal Processing (ICASSP)*, pages 3617–3621. IEEE, 2019.
- 404 [26] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech  
405 2: Fast and high-quality end-to-end text to speech. In *International Conference on Learning*  
406 *Representations*, 2021.
- 407 [27] Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech:  
408 Fast, robust and controllable text to speech. In *Advances in Neural Information Processing*  
409 *Systems*, 2019.
- 410 [28] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In  
411 *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- 412 [29] Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang,  
413 Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. Natural tts synthesis by  
414 conditioning wavenet on mel spectrogram predictions. In *2018 IEEE International Conference*  
415 *on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4779–4783. IEEE, 2018.
- 416 [30] RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron  
417 Weiss, Rob Clark, and Rif A Saurous. Towards end-to-end prosody transfer for expressive speech  
418 synthesis with tacotron. In *international conference on machine learning*, pages 4693–4702.  
419 PMLR, 2018.

- 420 [31] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
421 Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine*  
422 *learning research*, 15(1):1929–1958, 2014.
- 423 [32] Guangzhi Sun, Yu Zhang, Ron J Weiss, Yuan Cao, Heiga Zen, and Yonghui Wu. Fully-  
424 hierarchical fine-grained prosody modeling for interpretable speech synthesis. In *ICASSP*  
425 *2020-2020 IEEE international conference on acoustics, speech and signal processing (ICASSP)*,  
426 pages 6264–6268. IEEE, 2020.
- 427 [33] Qiao Tian, Yi Chen, Zewang Zhang, Heng Lu, Linghui Chen, Lei Xie, and Shan Liu. Tfgan:  
428 Time and frequency domain based generative adversarial network for high-fidelity speech  
429 synthesis. *arXiv preprint arXiv:2011.12206*, 2020.
- 430 [34] Rafael Valle, Kevin J Shih, Ryan Prenger, and Bryan Catanzaro. Flowtron: an autoregressive  
431 flow-based generative network for text-to-speech synthesis. In *International Conference on*  
432 *Learning Representations*, 2020.
- 433 [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
434 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*,  
435 2017.
- 436 [36] Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly,  
437 Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: Towards end-to-end  
438 speech synthesis. *arXiv preprint arXiv:1703.10135*, 2017.
- 439 [37] Yuxuan Wang, Daisy Stanton, Yu Zhang, RJ-Skerry Ryan, Eric Battenberg, Joel Shor, Ying  
440 Xiao, Ye Jia, Fei Ren, and Rif A Saurous. Style tokens: Unsupervised style modeling, control  
441 and transfer in end-to-end speech synthesis. In *International Conference on Machine Learning*,  
442 pages 5180–5189. PMLR, 2018.
- 443 [38] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel wavegan: A fast waveform gen-  
444 eration model based on generative adversarial networks with multi-resolution spectrogram. In  
445 *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing*  
446 *(ICASSP)*, pages 6199–6203. IEEE, 2020.
- 447 [39] Ya-Jie Zhang, Shifeng Pan, Lei He, and Zhen-Hua Ling. Learning latent representations for style  
448 control and transfer in end-to-end speech synthesis. In *ICASSP 2019-2019 IEEE International*  
449 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6945–6949. IEEE,  
450 2019.
- 451 [40] Zachary Ziegler and Alexander Rush. Latent normalizing flows for discrete sequences. In  
452 *International Conference on Machine Learning*, pages 7673–7682. PMLR, 2019.

## 453 Checklist

- 454 1. For all authors...
- 455 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
456 contributions and scope? [Yes]
- 457 (b) Did you describe the limitations of your work? [Yes] See appendix.
- 458 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
459 appendix.
- 460 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
461 them? [Yes]
- 462 2. If you are including theoretical results...
- 463 (a) Did you state the full set of assumptions of all theoretical results? [Yes]

- 464 (b) Did you include complete proofs of all theoretical results? [Yes]
- 465 3. If you ran experiments...
- 466 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
467 mental results (either in the supplemental material or as a URL)? [Yes] We give the  
468 demo website URL of this work, where you can find results generated by our proposed  
469 method. We also attach important code snippets in the supplementary material, which  
470 are helpful for reproducing the main experimental results.
- 471 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
472 were chosen)? [Yes] Training details see Section 4, hyperparameters see appendix.
- 473 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
474 ments multiple times)? [Yes]
- 475 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
476 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.4.
- 477 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 478 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 479 (b) Did you mention the license of the assets? [Yes]
- 480 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 481 (d) Did you discuss whether and how consent was obtained from people whose data you're  
482 using/curating? [Yes] We use internal datasets.
- 483 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
484 information or offensive content? [N/A]
- 485 5. If you used crowdsourcing or conducted research with human subjects...
- 486 (a) Did you include the full text of instructions given to participants and screenshots, if  
487 applicable? [No]
- 488 (b) Did you describe any potential participant risks, with links to Institutional Review  
489 Board (IRB) approvals, if applicable? [N/A]
- 490 (c) Did you include the estimated hourly wage paid to participants and the total amount  
491 spent on participant compensation? [No]