# ONE-STEP IMAGE-FUNCTION GENERATION VIA CONSISTENCY TRAINING

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

Paper under double-blind review

#### ABSTRACT

Consistency models aim to deliver a U-Net generator to map noise to images directly and enable swift inference with minimal steps, even trained in isolation with consistency training mode. However, the U-Net generator requires heavy feature extraction layers for multi-level resolutions and learning convolution kernels with specific receptive fields, resulting in the challenge that consistency models suffer from heavy training resources and fail to generate images with any user-specific resolutions. In this paper, we first validate that training the original consistency model with a small batch size via consistency training mode is pretty unstable, which motivates us to investigate efficient and flexible consistency models. To this end, we propose to use a novel Transformer-based generator to generate continuous image functions, which can then be differentially rendered as images with arbitrary resolutions. We adopt implicit neural representations (INRs) to form such continuous functions, which help to decouple the resolution of generated images and the total amount of the parameters generated from the neural network. Extensive experiments on one-step image generation demonstrate that our method greatly improves the performance of consistency models with low training resources and also provides an efficient any-resolution image sampling process.

028 029 1

### 1 INTRODUCTION

031 Diffusion models Sohl-Dickstein et al. (2015); Song & Ermon (2019); 2020) have achieved remarkable efficacy in synthesizing various signals, including audio Kong et al. (2020); Chen et al. (2020). 033 image Dhariwal & Nichol (2021); Ramesh et al. (2022) and video Harvey et al. (2022); Ho et al. 034 (2022). However, diffusion models rely on an iterative sampling process, leading to slow generation. Consistency model Song et al. (2023) is an emerging family of diffusion models Ho et al. (2020); Song et al. (2020a) that directly map noise to data by maintaining point consistency on ODE trajectory. This unique characteristic enables consistency models to support rapid one-step gener-037 ation for high-quality samples by design. Consistency models Song et al. (2023) can be trained through two modes: consistency distillation from a pre-trained diffusion model or consistency training in isolation. Different from consistency distillation, consistency training does not rely on the 040 pre-trained diffusion model but utilizes an unbiased estimator to approximate the ground-truth score 041 function Song et al. (2023). Consequently, consistency training emerges as a more flexible and 042 convenient method for training consistency models and shows greater potential in generation tasks. 043

However, consistency models face challenges due to their substantial training resource requirement 044 and inflexible image generation with fixed resolution. Consistency models rely on a U-Net generator Ronneberger et al. (2015) to map noise to images, while the U-Net involves extensive convolution 046 to extract features and is proved to be less scalable than the Transformer-based generator in diffusion 047 models Peebles & Xie (2023). Our investigation reveals that training consistency models based on 048 a U-Net generator under the consistency training model requires a large batch size. Directly reducing the training batch size to accommodate limited training resources significantly diminishes the generation performance, leading to the generation of non-realistic images, as illustrated in Figure 1 051 (a). Besides, as depicted in Figure 1 (b), to generate a specific-resolution image, the U-Net needs to denoise a noisy image with the same resolution, causing consistency models can only generate 052 images with fixed resolution once trained on a dataset with specific-resolution images. Therefore, a more efficient and flexible generator is desired for consistency models.

056

058

060 061

062

063

064

065



Figure 1: (a) When we train the consistency models on the CelebA dataset with consistency training and a batch size of 4, the generation performance gradually diminishes. (b) The U-Net generates images with the same resolution as the noisy images. Therefore we need to train infinite separate generators at all resolutions if we want to sample images with arbitrary resolutions. (c) Our image function generator treats the images as continuous functions parameterized as the MLPs. With a single generator, it generates a fixed amount of MLP parameters as an image function, which can then be rendered as images with arbitrary resolutions.

066 067 068

090 091 092

093

094

096

097

098

099

102

103

In this paper, we propose a novel and efficient generator for consistency models that stabilizes the 069 consistency training process under low training resources and can generate images with any userspecific resolution in the inference phase. Specifically, instead of generating discrete grid represen-071 tations of images, we treat images as continuous functions and introduce a novel Transformer-based 072 generator Dosovitskiy et al. (2020) that predicts the continuous functions of images as intermedi-073 ates. These intermediates can then be differentially rendered into images with arbitrary resolution, 074 as depicted in Figure 1 (c). To represent the continuous functions of images, we leverage implicit 075 neural representations Xie et al. (2022) (INRs) that employ Multi-layer Perceptrons (MLPs) to map 076 the coordinates  $x \in \mathbb{R}^2$  to corresponding RGB values  $y \in \mathbb{R}^3$ . Compared to the U-Net generator, 077 our Transformer-based generator exhibits a more stable consistency training process of consistency 078 models with limited training resources. Moreover, the use of differential rendering enables the decoupling of the resolution of the generated images and the total amount of parameters generated 079 from the neural network, leading to efficient sampling of images at arbitrary resolutions.

081 Additionally, we empirically confirm that consistency training from scratch to generate functions 082 of high-resolution images is challenging and converges slowly. To enhance training efficiency, we 083 carefully design a novel architecture for the function generator, which consists of a feature extraction 084 module to denoise the noisy images, a function prediction module that predicts functions for clean images, and a non-learnable render module to obtain clean images. We then introduce an image 085 reconstruction pre-training task to pre-train the function prediction module, which compels the generator to convert a given clean image to its corresponding INR function, mitigating the optimization 087 challenges in consistency training when predicting the function of the clean image based solely on 088 a given noisy image. 089

- We summarize the contributions of this work as follows:
  - Toward efficient and flexible one-step generation for any-resolution images with consistency models, we propose a novel Transformer-based function generator, which first generates image functions with a single inference step, and then renders the image functions to produce images with arbitrary resolutions.
  - We carefully design a novel end-to-end architecture for the function generator that simplifies the generation of functions for clean images from given noisy images. This architecture involves a feature extraction module, a function prediction module, and a render module.
  - To enhance the consistency training efficiency, we introduce an image reconstruction pretraining task to fortify the function prediction module, enabling the function prediction module to predict functions of images giving clean images. This approach allows consistency training process to concentrate on the image-denoising task, resulting in more realistic image generation and faster convergence of the training process.
- Extensive experiments on one-step image generation under low training resources demonstrate that our methods can generate significantly more realistic images with much less training and inference resources than the original consistency models.



Figure 2: The comparison of the consistency training efficiency with the U-Net generator and our function generator based on image functions, including the parameter amount (left, lower is better), the training GPU costs (middle, lower is better), and inference frame per second (FPS, right, higher is better). We can observe that the increment of the parameter amount of our generator is nearly negligible, while the training GPU cost and inference FPS of our function generator are more efficient and have better scalability than the U-Net generator.

- 2 RELATED WORK
- 125 126 127

119

120

121

122

123 124

Diffusion Models. Diffusion models Sohl-Dickstein et al. (2015); Song et al. (2020a) are emerging topics in the computer vision community. They train a generator to denoise the noise-corrupted data to estimate the score of data distribution and iteratively denoise the data point sampled from the noise distribution to generate new samples. Lots of work are done to accelerate the inference speed of diffusion models, such as faster ODE solver Song et al. (2020a); Lu et al. (2022a;b), predictor-corrector methods Song et al. (2020b) distillation methods Salimans & Ho (2022); Meng et al. (2023) Yin et al. (2024) and rectification Liu et al. (2022; 2023c).

135 **Consistency Models.** Consistency models Song et al. (2023) is a new type of diffusion model for 136 few-step sampling while maintaining good generation quality. They deliver consistency mapping to map any point in ODE trajectory to its origin, enabling one-step generation. Unlike GAN-based 137 generation models Goodfellow et al. (2014), consistency models do not rely on adversarial opti-138 mization and thus avoid the associated training difficulty. Consistency models can be trained either 139 by consistency distillation mode from a pre-trained diffusion model or by consistency training mode 140 with an unbiased estimator to approximate the ground-truth score function. Luo et al. Luo et al. 141 (2023) apply consistency distillation to train consistency models in latent space. In this work, we 142 focus on consistency training because it does not rely on an additional pre-trained diffusion model 143 and is more flexible and convenient for training. 144

Efficient Diffusion. Apart from classical U-Net Ronneberger et al. (2015), several works Bao et al. (2023a,b); Peebles & Xie (2023) successfully adopt Vision Transformer Dosovitskiy et al. (2020) as the generator in diffusion models. They Bao et al. (2023a); Peebles & Xie (2023) show that the Transformer generator enjoys better scalability and higher performance in generating highresolution images in latent diffusion models Rombach et al. (2022) than the U-Net generator. However, these works only explore the simple case of replacing the U-Net generator in the diffusion models with a ViT generator. On the contrary, our work is the first work to deliver a ViT generator to generate INR functions in the novel consistency models that support one-step generation.

152 Flexible Image Generation. Some works have targeted the problem of flexible image generation 153 with any resolution, either based on modifying the diffusion trajectory and model or based on a 154 patch-by-patch assembly manner. [Haji-Ali et al.] (2023); [Zhang et al.] (2023) focus on the flexible 155 sampling of U-Net-based diffusion models and enable iteratively generating images with specific 156 resolutions by decoupling the generation trajectory or dynamically adjusting the feature map size. 157 However, these methods rely on operating on the iterative sampling process, which makes them un-158 suitable for the single-step sampling process with consistency models. The patch-by-patch assem-159 bly manner works Chai et al. (2022); Lin et al. (2022) deliver a patch-by-patch strategy to generate patches and assemble the patches into a larger image with particular resolution. However, the output 160 of their generator is still of a fixed resolution, therefore they require multiple inference processes for 161 a larger image and cannot generate images with lower resolution.

162 **Implicit Neural Representations.** By mapping a coordinate to its corresponding quantity with a 163 neural network (e.g. MLP), INRs have shown great potential in representing complex continuous 164 functions for a lot of natural signals, such as time-serial signals Fons et al. (2022); Szatkowski et al., 165 images Sitzmann et al. (2020b); Skorokhodov et al. (2021); Liu et al. (2023a) and 3D scenes Park 166 et al. (2019); Sitzmann et al. (2020a); Liu et al. (2023b). A lot of works on how to generate INRs fast for unseen signals have been employed, including meta-learning Sitzmann et al. (2020a); Tancik 167 et al. (2021); Liu et al. (2023a); Finn et al. (2017) and hyper-network Chen & Wang (2022); Kim 168 et al. (2023); Zhang et al. (2024). 169

170 Diffusion Models Based on Implicit Neural Representations. Different from diffusion based on 171 explicit field Zhuang et al. (2023), when adopting diffusion models to high-resolution images or 172 complex 3D signals, existing novel research Dupont et al. (2022); Karnewar et al. (2023); Erkoc et al. (2023); Chen et al. (2024) firstly convert the signals to their corresponding INR functions and 173 then train a diffusion model on these INR functions, enabling generating complex signals efficiently. 174 However, the two-stage training process is inflexible and the error in the first representation stage 175 would greatly affect the performance of the second diffusion stage. On the contrary, our work 176 proposes an end-to-end training framework to generate INRs of signals and requires rendering for 177 only one time when generating one-step inference results. 178

179 180

181

182

#### **INR-BASED CONSISTENCY TRAINING** 3

#### 3.1 PRELIMINARIES

183 Consistency models Song et al. (2023) is a 184 new family of generative models that en-185 ables a few-step generation. The core idea 186 of CM is the PF-ODE Song et al. (2020b). 187 Denote the data distribution by  $p_{\text{data}}(\mathbf{x})$ , 188 and the perturbed distribution is presented 189 as  $p_{\sigma}(\mathbf{x}) = \int p_{\text{data}}(\mathbf{y}) \mathcal{N}(\mathbf{x} \mid \mathbf{y}, \sigma^2 \mathbf{I}) \, \mathrm{d}\mathbf{y}$ 190 if we add Gaussian noise  $\mathcal{N}(0, \sigma^2)$  with noise level  $\sigma$  to the data. Then, the PF-191 ODE presented in Karras et al. Karras 192 et al. (2022) is formulated as: 193

194 
$$\frac{\mathrm{d}\mathbf{x}}{\mathrm{d}\sigma} = -\sigma\nabla\log p_{\sigma}\left(\mathbf{x}\right) \approx -\sigma \boldsymbol{s}_{\phi}\left(\mathbf{x},\sigma\right),$$
196 (1)

where  $s_{\phi}(\mathbf{x}, \sigma) \approx \nabla \log p_{\sigma}(\mathbf{x})$  is the 197 score function Song et al. (2020b) of 198  $p_{\sigma}(\mathbf{x})$ . Here, as in Song et al. (2023); 199 Karras et al. (2022),  $\sigma$  is defined as  $\sigma \in$ 200  $[\sigma_{\min}, \sigma_{\max}]$ , where  $\sigma_{\min}$  is a small posi-201 tive number to ensure  $p_{\sigma_{\min}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$ 202 and  $\sigma_{\max}$  is a large positive number such 203 that  $p_{\sigma_{\max}}(\mathbf{x}) \approx \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 \mathbf{I}).$ 



Figure 3: Given a probability-flow ODE that smoothly converts data to noise, we learn a denoising neural network G to map any point (e.g.,  $\mathbf{x}_{\sigma}, \mathbf{x}_{\sigma'}$ , and  $\mathbf{x}_{\sigma_{\max}}$ ) on the ODE trajectory to the continuous function of the origin (e.g.,  $\beta_{\mathbf{x}_0}$ ) for generative modeling, which can then be differentially rendered as the original data with arbitrary resolutions.

204 Solving the PF-ODE from noise level  $\sigma$  to  $\sigma_{min}$  in Eq. 1 indeed establishes a bijective mapping 205 from a noisy data sample  $\mathbf{x}_{\sigma} \sim p_{\sigma}(\mathbf{x})$  to the real data sample  $\mathbf{x}_{\sigma_{\min}} \sim p_{\sigma_{\min}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$ . This mapping  $f^* : (\mathbf{x}_{\sigma}, \sigma) \mapsto \mathbf{x}_{\sigma_{\min}}$  is defined as a *consistency function* in Song *et al.* Song *et al.* (2023). 206 207 By the definition, the consistency function satisfies the boundary condition  $f^*(\mathbf{x}, \sigma_{\min}) = \mathbf{x}$ . To 208 approximate the consistency function with boundary condition, a consistency model  $f_{\theta}(\mathbf{x},\sigma)$  is 209 parameterized as: 210

$$\boldsymbol{f}_{\boldsymbol{\theta}}(\mathbf{x},\sigma) = c_{\text{skip}}\left(\sigma\right)\mathbf{x} + c_{\text{out}}\left(\sigma\right)\boldsymbol{F}_{\boldsymbol{\theta}}(\mathbf{x},\sigma),\tag{2}$$

211 where  $F_{\theta}(\mathbf{x}, \sigma)$  is a free-form denoising neural network with parameter  $\theta$ , while  $c_{skip}(\sigma)$  and  $c_{out}(\sigma)$ 212 are differential functions so that  $c_{\text{skip}} (\sigma_{\min}) = 1$  and  $c_{\text{out}} (\sigma_{\min}) = 0$ .

213 The consistency function has the property of *self-consistency*: the outputs are consistent for arbitrary 214 pairs of  $(\mathbf{x}_{\sigma}, \sigma)$  that belong to the same PF-ODE trajectory, which can be formulated as: 215

$$\boldsymbol{f}(\mathbf{x}_{\sigma},\sigma) = \boldsymbol{f}(\mathbf{x}_{\sigma'},\sigma'), \forall \sigma, \sigma' \in [\sigma_{\min},\sigma_{\max}].$$
(3)



Figure 4: The architecture for our image-function generator. It contains three modules: the feature extraction module that contains several encoder blocks to extract features on the input noisy images and gives data tokens as a condition to guide the generation of the image function; the function prediction module that contains several decoder blocks to predict the INR parameters for image functions with given data token from the previous module; the non-learnable render module that multiples the INR parameters with coordinates and finally render as images.

Therefore, consistency models can be trained by enforcing the self-consistency property with a consistency loss between the results denoised from the  $i^{th}$  noise level and the  $(i + 1)^{th}$  noise level:

$$\mathcal{L} = \mathbb{E} \left[ d \left( f_{\boldsymbol{\theta}} \left( \mathbf{x}_{\sigma_{i+1}}, \sigma_{i+1} \right), f_{\boldsymbol{\theta}^{-}} \left( \hat{\mathbf{x}}_{\sigma_{i}}^{\phi}, \sigma_{i} \right) \right) \right], \tag{4}$$

where  $d(\cdot, \cdot)$  is a metric function such as  $\ell_2$  metric or learned perceptual image patch similarity (LPIPS) metric Zhang et al. (2018) while  $\theta^- \leftarrow \mu \theta^- + (1 - \mu)\theta$  is the target model parameter updated with the exponential moving average (EMA) of the parameter  $\theta$  and EMA decay rate  $\mu$ .  $\hat{\mathbf{x}}_{\sigma_i}^{\phi}$  is derived from  $\mathbf{x}_{\sigma_{i+1}}$  by solving the PF-ODE in the reverse direction for a single step:

$$\hat{\mathbf{x}}_{\sigma_i} = \mathbf{x}_{\sigma_{i+1}} - (\sigma_i - \sigma_{i+1}) \sigma_{i+1} \nabla_{\mathbf{x}} \log p_{\sigma_{i+1}}(\mathbf{x}) \Big|_{\mathbf{x} = \mathbf{x}_{\sigma_{i+1}}}$$

To estimate the unknown score function  $\nabla_{\mathbf{x}} \log p_{\sigma_{i+1}}(\mathbf{x})$ , Song *et al.* Song *et al.* (2023) propose *consistency training* that employs an approximation  $\hat{\mathbf{x}}_{\sigma_i} = \mathbf{x} + \sigma_i \mathbf{z}$  with the same  $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$  and  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  to calculate  $\mathbf{x}_{\sigma_{i+1}} = \mathbf{x} + \sigma_{i+1}\mathbf{z}$ . Therefore, the consistency training objective  $\mathcal{L}_{\text{CT}}$  can be defined as:

$$\mathcal{L}_{\rm CT} = \mathbb{E}\left[d\left(\boldsymbol{f}_{\boldsymbol{\theta}}\left(\mathbf{x} + \sigma_{i+1}\mathbf{z}, \sigma_{i+1}\right), \boldsymbol{f}_{\boldsymbol{\theta}^{-}}\left(\mathbf{x} + \sigma_{i}\mathbf{z}, \sigma_{i}\right)\right)\right].$$
(5)

267

268

228

229

230

231

232

233 234

235

236 237 238

239

240 241

242 243 244

245

246

247

248

#### 3.2 GENERATE IMAGE FUNCTIONS AND SAMPLE ANY-RESOLUTION IMAGES

Our modification focuses on the free-form denoising neural network  $F_{\theta} : (\mathbf{x}, \sigma) \to \mathbf{x}_d$ , which is typically a U-Net model Song & Ermon (2019); Ronneberger et al. (2015) that gets noisy image  $\mathbf{x} \in \mathcal{R}^{C \times H \times W}$  and noise level  $\sigma \in \mathcal{R}$  as input. As shown in Figure 1 (b), the U-Net model applies heavy feature extraction at the different resolution of the image and finally outputs a denoised image  $\mathbf{x}_d \in \mathcal{R}^{C \times H \times W}$ . We find that the U-Net can only generate the output image with exactly the same resolution as the input image, which is not flexible enough to scale to high-resolution images with limited training resources.

Therefore, we seek a more efficient and more noise-robust method to generate high-resolution images Dupont et al. (2022); Rahaman et al. (2019); Skorokhodov et al. (2021). We show our pipeline in Figure 3. As in INR methods Sitzmann et al. (2020b); Liu et al. (2023a); Chen & Wang (2022), rather than discrete grid representation, we consider images as continuous functions, which can be parameterized as neural networks, e.g. MLPs Dupont et al. (2022). Specifically, we consider an image as a collection of paired coordinates and RGB values  $\{(\mathbf{c}, \mathbf{y})\}_{H \times W}$  and fit an MLP  $M_{\beta}$  with learnable parameter  $\beta$  to map the coordinates  $\mathbf{c} \in \mathcal{R}^2$  to its corresponding RGB values  $\mathbf{y} \in \mathcal{R}^3$ :

$$M_{\beta}(\mathbf{c}) = \mathbf{y}.\tag{6}$$

To obtain smooth super-resolution performance and eliminate artifacts, we follow Zhang et al. (2024) to apply variational coordinates to sample the coordinates c.



Figure 5: (a) & (b) The detailed implementation of encoder and decoder block (Norm is not shown).(c) Pre-training data flow. We set the noise level to 0, skip the feature extraction module, and optimize the Function Prediction Module to predict the image function for clean images with reconstruction loss. The black path is the data flow when consistency training the whole model while the red path is the pre-training of the function prediction module with the image reconstruction task.

Note that each image can be considered as an image function  $M_{\beta}$  and holds its own parameter  $\beta = \{\mathbf{W}^t, \mathbf{B}^t\}_{t=0}^{t=T-1}$  for an MLP with totally T layers. Generating a denoised image is equivalent to generating the INR parameter  $\beta$  of the image function corresponding to the denoised image:

$$\boldsymbol{G}_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\sigma}) = \boldsymbol{\beta},\tag{7}$$

where  $G_{\theta}$  is our proposed function generator and its architecture is introduced in Section 3.3

Once the parameter  $\beta$  is obtained, the denoised image  $\mathbf{x}_d$  can be reconstructed at arbitrary resolution by querying the RGB values with any specific coordinates  $\{\mathbf{c}\}_{H \times W}$ . As a result, we parameterize the free-form neural network  $F_{\theta}$  with a generator that generates image function as intermediate:

$$F_{\theta}(\mathbf{x},\sigma,\{\mathbf{c}\}) = M_{G_{\theta}(\mathbf{x},\sigma)}(\{\mathbf{c}\}) = \mathbf{x}_{d}.$$
(8)

300 Note that each layer of MLP can be formulated as:

$$\mathbf{c}^{t+1} = \operatorname{act}(\mathbf{W}^t \mathbf{c}^t + \mathbf{B}^t), \tag{9}$$

where act is the activation function. The whole forward process of  $M_{\beta}(\mathbf{c})$  is entirely differential if the activation function act is differential, which enables the end-to-end backward process after calculating consistency training loss, as shown in Eq. 5

Compared with the U-Net generator that generates images with the same resolution as the input images, our function generator only generates the INR parameters  $\beta$  of the image function that has a fixed size and does not scale with the input image resolution. When scaling to the larger resolution image, our pipeline only needs to query the image function with finer-grained coordinates {c}, which greatly decreases the parameters amount, training time, and memory cost of the whole pipeline and provides a flexible any-resolution image sampling process.

312 313

282

283

284

286 287 288

289

290 291

292

294

295

296

297 298 299

301

302

3.3 INR GENERATOR DESIGN

In this part, we introduce the detailed architecture of our function generator. As presented in Figure 4, it contains three major modules: a feature extraction module, a function prediction module, and a render module.

**Feature Extraction Module.** The feature extraction module is utilized to extract features from noisy images and output the data tokens to guide the generation of the image function for the denoised images. We mainly follow DiT Peebles & Xie (2023) and deliver adaLN-Zero Transformer blocks to form a feature extraction encoder. As presented in Figure [5] (a), we adopt a linear layer to regress the dimension-wise scale and shift parameters  $\alpha$ ,  $\gamma$ , and  $\tau$  from noise level embedding.

**Function Prediction Module.** The function prediction module is designed for generating the INR parameters  $\beta$  for the image function based on the image feature extracted from the feature extraction

Dataset	Models	FID $(\downarrow)$	SFID $(\downarrow)$	IS (†)	P (†)	R (†)
Cifar10 32	CM-UNet	32.87	20.94	6.16	0.595	0.23
Cital 10-52	CM-Func	28.87	19.81	6.92	0.52	0.30
Calab A 64	CM-UNet	54.41	72.86	1.77	0.43	0.021
CEIEUA-04	CM-Func	29.49	45.10	2.08	0.78	0.094
Calab A 128	CM-UNet	89.46	158.64	1.40	0.46	0
CelebA-128	CM-Func	69.3	124.22	1.66	0.46	0.002
I SUN Church 128	CM-UNet	58.33	88.33	1.82	0.31	0.007
LSON Church-126	CM-Func	34.94	86.28	2.42	0.33	0.053
LSUN Classroom-128	CM-UNet	65.18	84.69	2.54	0.41	0.026
	CM-Func	57.96	76.31	2.83	0.50	0.033

Table 1: Table for the one-step image generation performance. The number in the name means the training resolution.

module. Before training, we randomly initialize the learnable INR tokens according to the shape of the parameters  $\beta$ . As shown in Figure [5] (b), we design a decoder block that adopts multi-head cross attention to fuse the data feature with INR tokens and predict the INR parameters for the image function of specific denoised images. Note that we follow Chen *et al.* [Chen & Wang] (2022) to use a grouping strategy to improve the efficiency and scalability of our function prediction module.

**Render Module.** After obtaining INR parameters  $\beta$  for the image functions, we need to differentially render images to ensure the whole pipeline is differential. As discussed before, we form a continuous image function as an MLP with parameter  $\beta$  corresponding to a specific image. With given resolution  $H \times W$ , we sample a coordinate list  $\{(\frac{i}{H}, \frac{j}{W})\}_{i,j}$  where  $i \in [0, H)$  and  $j \in [0, W)$ and query the RGB value at each coordinate, which finally can be reshaped to form a complete image. There is no learnable parameter in the render module since the INR parameter  $\beta$  is generated by the function prediction module.

#### 3.4 BENEFIT FROM IMAGE RECONSTRUCTION PRE-TRAINING

We notice that predicting the image function for a clean image from a noisy image is pretty difficult if not impossible, due to the large search space of the INR parameter  $\beta$ . Therefore, we propose an image reconstruction task to pre-train the function prediction module. We hope the feature extraction module focuses on transforming the noise image feature into the clean image feature while the function prediction module focuses on transforming the clean image feature into its image function.

Therefore, we design an image reconstruction pre-training task for the function prediction module. as shown in Figure [5](c). Specifically, during the image reconstruction pre-training task, we skip the feature extraction module and set the noise level  $\sigma$  as 0 so that the input image are neither perturbed with noise nor denoised by the feature extraction module. The pipeline is downgraded to a model (only the function prediction module contains learnable parameters) that transforms an input clean image to its corresponding image function with specific INR parameters  $\beta$ . We then optimize such a model with image reconstruction loss:

$$\mathcal{L}_{\text{rec}} = \text{MSE}(M_{\boldsymbol{G}_{\boldsymbol{\theta}}(\mathbf{x},\sigma=0)}\left(\{\mathbf{c}_{i}\}\right), \mathbf{x}), \tag{10}$$

<sup>366</sup> where MSE denotes the mean square error loss.

After pre-training the function prediction module, we train the whole model (including the feature extraction module and function prediction module) with the consistency training objective shown in Eq. 5, which enables our model to generate new images with given random noise. We will show that the image reconstruction pre-training task helps to make the consistency training process converge faster.

371372373

364 365

367

368

369

370

326 327 328

338

339

340

341

342

350

351

#### 4 EXPERIMENTS

374 375

In the experiment section, we first provide a detailed comparison between the original U-Net generator and our function generator in the one-step image generation task based on consistency models under the low-training resource, denoted as *CM-UNet* and *CM-Func* respectively. Then we provide

382

384

385 386

387 388

389

390

391 392

393

378



Figure 6: Results from (a) the CelebA dataset with resolution 64 (a) and from (b) the LSUN Church Outdoor with resolution 128. All corresponding images are generated from the same initial noise.

more discussion and ablation studies to illustrate the advantages of function generator when sampling images with different resolutions compared to other popular generators, such as UViT Bao et al. (2023a) and DiT Peebles & Xie (2023).

4.1 Setting

Datasets. We show the performance for the unconditional image generation on two popular image datasets: CelebA Liu et al. (2015) with resolution 64 and 128 and LSUN dataset Yu et al. (2015) with resolution 128. We also show that our pipeline can be applied to the class-conditional image generation with the Cifar10 dataset Krizhevsky et al. (2009). We mainly evaluate the models by sampling images with corresponding resolutions.

**Hyper-parameters.** Following the setting of training the consistency models as in Song *et al.* [Song et al.] (2023), we use the following default hyper-parameters for all datasets unless otherwise stated:  $\sigma_{\min} = 0.002$ ,  $\sigma_{\max} = 80$ . We use a low batch-size training strategy to evaluate the performance of consistency models under the low-training resource setting. We set the number of encoder blocks N = 8 and the number of decoder blocks M = 6. We adopt MLP with a depth of 5 and width of 256, with ReLU activation, and positional embedding as our image function. For detailed hyper-parameters for consistency models and optimization processes, we present more details in Appendix [A].

For the pre-training task, we use the same dataset as the consistency training task and follow Chen *et al.* Chen & Wang (2022) to optimize the models with Adam optimizer Kingma & Ba (2014) and learning rate 1e - 4 for 30 epochs. All results are reported for models with the pre-training task unless otherwise discussed.

Metric. We first compare the efficiency of consistency models with the U-Net generator and function generator. Then we report the quantitative generation results according to Frechet Inception Distance (FID) Heusel et al. (2017), Sliding Fréchet Inception Distance (sFID) Szegedy et al. (2016), Inception Score (IS) Salimans et al. (2016), Precision (P) Kynkäänniemi et al. (2019), and Recall (R) Kynkäänniemi et al. (2019). We follow Song et al. (2023) and Dhariwal & Nichol (2021) to generate 50000 images for credible scores.

We also define a new metric, the total denoising distance, to reflect the denoising quality of our generators during training. The total denoising distance is simply calculated as the L2 difference between the denoised result from the noisy image and the original clean image:

$$d_{\text{denoising}}^{i} = |\boldsymbol{f}_{\theta} (\mathbf{x}_{\sigma_{i}}, \sigma_{i}) - \mathbf{x}_{0}|$$

Though this metric does not reflect the generation performance in the test phase, it means the denoising ability of the generator during training. More details for this metric are in Appendix B

4.2 Results

420

421

422

423 424

425

Model efficiency. We quantitatively compare the efficiency of our function generator and the base-line U-Net generator on the unconditional image generation task to show that our model is more efficient. We report the total parameters, the GPU cost, and the inference FPS in Figure 2. We verify that our model has fewer total parameters, less training GPU cost, and higher inference FPS.

431 Apart from the absolute quantitative value, we also observe that when increasing the resolution of images (from Cifar10-32 to CelebA-64 and then to CelebA-128), the increment of the total param-



Figure 7: (a) Visual results are generated from the two initial noises by models with different optimization iterations. (b) The convergence curve for the consistency training. The results indicate that the U-Net is unstable while our model is much more stable when keeping training.



Figure 8: Top Left: Visual results of optimizing models with or without pre-training task for 50,000 iterations. Bottle Left: Alabtion of Pre-training task on CelebA-64 with 400,000 training iterations. Right: Denoising ability during training. Pre-training task greatly improves the denoising ability.

eter and GPU cost of our model is much less than the U-Net or even negligible, which verifies that our model has better scalability than U-Net when scaling to image with higher resolution.

**Generation performance.** We demonstrate the quantitative generation performance in Table 1. We can observe that the generation performance on the Cifar10, CelebA, and LSUN datasets beats the baseline U-Net in terms of all evaluation metrics, while our model has much fewer model parameters and enjoys faster training, faster inference, and flexible employment.

We also present some visual generation results based on the CelebA dataset with resolution 64 and the LSUN dataset with resolution 128 in Figure 6. We can clearly observe that given exactly the same initial noise, our model can generate much more realistic than the baseline consistency model with the U-Net generator. More visual and quantitative results are presented in Appendix D

4.3 DISCUSSIONS

Training oscillates for Consistency models with U-Net generator. We go deeper into the failing case when training consistency models with U-Net and low batch size 32 and show the results in Figure 7. We find that during the training of the consistency model with U-Net with a low batch size (32 for one single A6000 GPU), the generation results tend to be corrupted and fail to be denoised. On the contrary, the results from the consistency model with our function generator are much more consistent during training and we can find that these results are clean. We also present the consistency training objective convergence curve for the training consistency model in Figure 7(b). The result shows that the consistency training process with the U-Net generator is pretty unstable and oscillating. In contrast, the process with our function generator is much more stable and consistent. 

Pre-training improves the convergence of the consistency training with the function generator. We compare the total denoising distance metric of training consistency models with our function generator with or without pre-training task and show the result in Figure 8. We also show the generation performance at 400,000 iterations on the CelebA-64 dataset by ablating the pre-training task in the Bottle Left Table in Figure 8. The results show that our pre-training reconstruction task greatly improves the denoising ability of our function generator during the early stage of consistency training and leads to better generation performance. We also verify that our pre-training task is very efficient and costs much less time compared with the consistency training process. More quantitative results for the efficiency of the pre-training task are shown in Appendix C

	Efficience		Accuracy			
Models	Any-Resolution	Multi-Resolution		sFID ( $\downarrow$ )	IS (†)	P (†)
widdels	Sampling	Sampling FPS ( <sup>†</sup> )	[ I'ID (↓)			
CM-UNet	X	83.64	54.41	72.86	1.77	0.43
CM-UViT	X	83.58	40.14	41.47	1.90	0.68
CM-DiT	X	96.61	25.52	37.27	2.00	0.81
<b>CM-Func</b>	$\checkmark$	447.02	29.49	45.10	2.08	0.78

Table 2: Table for efficiency and accuracy of different generators on CelebA-64 dataset.

496 Flexible sampling image resolution 497 leads to a more efficient sampling 498 process. A comprehensive experiment on the CelebA dataset shown in 499 Table 2 is conducted to evaluate the 500 efficiency and generation quality. We 501 evaluate the efficiency by measuring the 502 FPS to sample 10000 image signals, each of which requires 3 resolutions 504  $(32 \times 32, 64 \times 64, 128 \times 128, \text{ totally})$ 505 30000 images). Only our function gen-506 erator supports sampling any-resolution 507 images with one model, while other gen-508 erators require training multiple separate 509 models at each resolution. Therefore, the multi-resolution sample FPS for those 510 generators that generates fixed-resolution 511



Figure 9: After training on the CelebA-64 dataset, our pipeline supports sampling at any resolution while the results at all resolutions remain clean and realistic.

images is calculated as  $\frac{30000}{\sum_{i=1}^{10000} T_{32}^i + \sum_{i=1}^{10000} T_{64}^i + \sum_{i=1}^{10000} T_{128}^i}$ , where  $T_k^i$  is the average time that 512 513 generates  $i^{th}$  image with resolution k. Since the function generator only needs to generate one image function for each image, the multi-resolution sample FPS for our function generator is 514 calculated as  $\frac{30000}{\sum_{i=1}^{10000} (T^i + T_{R32} + T_{R64} + T_{R128})}$ , where  $T^i$  is the average time that generates  $i^{th}$  image 515 functions, and  $T_{R32}, T_{R64}, T_{R128}$  are the time to render the image function to image with 32/64/128 516 517 resolutions (which is nearly negligible compared to  $T^{i}$ ). The results indicate that our function generator achieves a much higher multi-resolution sampling FPS. In contrast, other generator 518 needs to deliver different models to separately denoise input noisy images, which leads to a slower 519 sampling process. 520

For generation quality, we find that our function generator greatly beats U-Net and is comparable
with Transformer-based generators. We show the generation results with different sampling resolutions in Figure 9. The results indicate that the sampling results at all resolutions remain clean and
realistic. We also find that even trained with lower-resolution images, our function generator is still
better than the U-Net generator. In addition, we find that our image functions have a better interpolation performance than the linear interpolation. See more quantitative and visualization results in
Appendix D.2 and Appendix D.3

528 529

#### 5 CONCLUSION

530 In this paper, we explore efficient consistency training generation for consistency models. We ob-531 serve that the U-Net generator is resource-intensive and inflexible. Reducing batch size for single 532 GPU use harms performance. To address this, we introduce a Transformer generator that generates 533 image functions parameterized as INR and then renders them into images of any resolution. We de-534 sign a new end-to-end function predictor that simplifies generating clean image functions from noisy images. Additionally, we enhance training efficiency with an image reconstruction pre-training task. 536 Extensive experiments show our method produces more realistic images with fewer resources than 537 the original consistency models and also provides an efficient any-resolution image sampling process. A limitation of this paper lies in that INRs for image functions are global representations of 538 signals, and have poor representation ability of the local semantic information, which makes our pipeline hard for generation with finer details.

## 540 REFERENCES

547

559

- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *CVPR*, pp. 22669–22679, June 2023a.
- Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang
  Su, and Jun Zhu. One transformer fits all distributions in multi-modal diffusion at scale. *arXiv preprint arXiv:2303.06555*, 2023b.
- Lucy Chai, Michael Gharbi, Eli Shechtman, Phillip Isola, and Richard Zhang. Any-resolution training for high-resolution image synthesis. In *European Conference on Computer Vision*, pp. 170– 188. Springer, 2022.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wave grad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*, 2020.
- 553
   554
   555
   Yinbo Chen and Xiaolong Wang. Transformers as meta-learners for implicit neural representations. In *ECCV*, pp. 170–187. Springer, 2022.
- Yinbo Chen, Oliver Wang, Richard Zhang, Eli Shechtman, Xiaolong Wang, and Michael Gharbi.
   Image neural field diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8007–8017, 2024.
  - Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NIPS*, 34: 8780–8794, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
  Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
  image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Emilien Dupont, Hyunjik Kim, SM Eslami, Danilo Rezende, and Dan Rosenbaum. From data to functa: Your data point is a function and you can treat it like one. *arXiv preprint arXiv:2201.12204*, 2022.
- Ziya Erkoç, Fangchang Ma, Qi Shan, Matthias Nießner, and Angela Dai. Hyperdiffusion: Generating implicit neural fields with weight-space diffusion. *arXiv preprint arXiv:2303.17015*, 2023.
- 572 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
   573 of deep networks. In *ICML*, pp. 1126–1135. PMLR, 2017.
- 574
  575
  576
  576
  576
  577
  576
  577
  577
  578
  578
  579
  579
  570
  570
  570
  571
  572
  573
  574
  574
  574
  574
  574
  575
  576
  577
  577
  578
  578
  578
  578
  579
  579
  579
  570
  570
  570
  571
  572
  574
  574
  575
  575
  576
  577
  577
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
   Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Moayed Haji-Ali, Guha Balakrishnan, and Vicente Ordonez. Elasticdiffusion: Training-free arbitrary size image generation, 2023.
- William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible
  diffusion modeling of long videos. *NIPS*, 35:27953–27965, 2022.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
   Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NIPS*, 30, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NIPS*, 33: 6840–6851, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
   Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
   video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

598

- Animesh Karnewar, Andrea Vedaldi, David Novotny, and Niloy J Mitra. Holodiffusion: Training a 3d diffusion model using 2d images. In *CVPR*, pp. 18423–18433, 2023.
  - Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusionbased generative models. *NIPS*, 35:26565–26577, 2022.
- Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5148–5157, 2021.
- Chiheon Kim, Doyup Lee, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Generalizable implicit neural representations via instance pattern composers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11808–11817, 2023.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
   2009.
- Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *NIPS*, 32, 2019.
- Chieh Hubert Lin, Hsin-Ying Lee, Yen-Chi Cheng, Sergey Tulyakov, and Ming-Hsuan Yang. Infin itygan: Towards infinite-pixel image synthesis. In *International Conference on Learning Representations*, 2022.
- Ke Liu, Feng Liu, Haishuai Wang, Ning Ma, Jiajun Bu, and Bo Han. Partition speeds up learning
   implicit neural representations based on exponential-increase hypothesis. In *ICCV*, pp. 5474–5483, 2023a.
- Ke Liu, Ning Ma, Zhihua Wang, Jingjun Gu, Jiajun Bu, and Haishuai Wang. Implicit neural distance optimization for mesh neural subdivision. In *ICME*, pp. 2039–2044, 2023b. doi: 10.1109/ICME55011.2023.00349.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei
  Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019.
- Kingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- Kingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, and Qiang Liu. Instaflow: One step is enough
  for high-quality diffusion-based text-to-image generation. *arXiv preprint arXiv:2309.06380*, 2023c.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.
   In *ICCV*, pp. 3730–3738, 2015.
- 637
  638
  639
  640
  640
  641
  641
  642
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
  644
- 641 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast
  642 solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*,
  643 2022b.
- Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- 647 Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *CVPR*, pp. 14297–14306, 2023.

648 649 650	Anish Mittal, Rajiv Soundararajan, and Alan C. Bovik. Making a "completely blind" image qual- ity analyzer. <i>IEEE Signal Processing Letters</i> , 20(3):209–212, 2013. doi: 10.1109/LSP.2012. 2227726.
651 652 653 654	Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In <i>CVPR</i> , pp. 165–174, 2019.
655 656 657	William Peebles and Saining Xie. Scalable diffusion models with transformers. In <i>ICCV</i> , pp. 4195–4205, 2023.
658 659 660	Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In <i>ICML</i> , pp. 5301–5310. PMLR, 2019.
661 662	Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text- conditional image generation with clip latents. <i>arXiv preprint arXiv:2204.06125</i> , 1(2):3, 2022.
663 664 665	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>CVPR</i> , pp. 10684–10695, 2022.
666 667	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedi- cal image segmentation. In <i>MICCAI</i> , pp. 234–241. Springer, 2015.
668 669 670	Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. <i>arXiv</i> preprint arXiv:2202.00512, 2022.
671 672	Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. <i>NIPS</i> , 29, 2016.
673 674 675	Vincent Sitzmann, Eric Chan, Richard Tucker, Noah Snavely, and Gordon Wetzstein. Metasdf: Meta-learning signed distance functions. <i>NIPS</i> , 33:10136–10147, 2020a.
676 677	Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. <i>NIPS</i> , 33:7462–7473, 2020b.
678 679 680	Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial generation of continuous images. In <i>CVPR</i> , pp. 10753–10764, 2021.
681 682	Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In <i>ICML</i> , pp. 2256–2265. PMLR, 2015.
683 684 685	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. <i>arXiv</i> preprint arXiv:2010.02502, 2020a.
686 687	Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. <i>NIPS</i> , 32, 2019.
689 690	Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. <i>NIPS</i> , 33:12438–12448, 2020.
691 692 693 694	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. <i>arXiv preprint arXiv:2011.13456</i> , 2020b.
695 696	Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. <i>arXiv preprint arXiv:2303.01469</i> , 2023.
697 698 699 700	Filip Szatkowski, Karol J Piczak, Przemysław Spurek, Jacek Tabor, and Tomasz Trzciński. Hy- pernetworks build implicit neural representations of sounds. In <i>Joint European Conference on</i> <i>Machine Learning and Knowledge Discovery in Databases</i> .
704	Christian Szegedy Vincent Vanhoucke Sergey Joffe Jon Shlens and Zhigniew Woing Pathinking

701 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pp. 2818–2826, 2016.

702	Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T
703	Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representa-
704	tions. In CVPR, pp. 2846–2855, 2021.
705	

- Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 2555–2563, 2023.
- Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. In *Computer Graphics Forum*, volume 41, pp. 641–676. Wiley Online Library, 2022.
- Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1191–1200, 2022.
- Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6613–6623, 2024.
- Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv* preprint arXiv:1506.03365, 2015.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pp. 586–595, 2018.
- Shen Zhang, Zhaowei Chen, Zhenyu Zhao, Zhenyuan Chen, Yao Tang, Yuhao Chen, Wengang Cao, and Jiajun Liang. Hidiffusion: Unlocking high-resolution creativity and efficiency in low-resolution trained diffusion models. *arXiv preprint arXiv:2311.17528*, 2023.
- Shuyi Zhang, Ke Liu, Jingjun Gu, Xiaoxu Cai, Zhihua Wang, Jiajun Bu, and Haishuai Wang.
  Attention beats linear for fast implicit neural representation generation. *arXiv preprint arXiv:2407.15355*, 2024.
- Peiye Zhuang, Samira Abnar, Jiatao Gu, Alex Schwing, Joshua M Susskind, and Miguel Angel
   Bautista. Diffusion probabilistic fields. In *The Eleventh International Conference on Learning Representations*, 2023.
- 737
- 738 739

741 742

- 743
- 744 745
- 746
- 747

- 750
- 751
- 752
- 753
- 754 755