

---

# Fast Sampling of Diffusion Models via Operator Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Diffusion models have found widespread adoption in various areas. However,  
2 sampling from them is still slow because it involves emulating a reverse stochastic  
3 process with hundreds-to-thousands of neural network evaluations. Inspired by the  
4 recent success of neural operators in accelerating differential equations solving, we  
5 approach this problem by solving the underlying neural differential equation from  
6 an operator learning perspective. We examine probability flow ODE trajectories  
7 in diffusion model and observe a compact energy spectrum that can be learned  
8 efficiently in Fourier space. With this insight, we propose diffusion Fourier neural  
9 operator (DFNO) with temporal convolution in Fourier space to parameterize the  
10 operator that maps initial condition to the solution trajectory. DFNO can apply to  
11 any diffusion models and generate high-quality samples in one step. Our method  
12 achieves the state-of-the-art clean FID of 5.9 (legacy FID 4.72) on CIFAR-10 using  
13 one network evaluation.

## 14 1 Introduction

15 Diffusion models, also known as score-based generative models, have emerged as a powerful genera-  
16 tive modeling framework in various areas. They have achieved state-of-the-art (SOTA) performance  
17 in many applications including image generation [1, 2, 3, 4], molecule generation [5], audio synthesis  
18 [6, 7] and model robustness [8, 9]. However, sampling from diffusion models is still slower than  
19 other generative models such as generative adversarial networks (GAN) [10] by several orders of  
20 magnitude. Accelerating the sampling process of diffusion models remains challenging but important  
21 in applying diffusion models in many downstream tasks. Many studies have worked on the fast  
22 sampling of diffusion models which can be summarized into two categories.

23 **Training-free sampling methods** focus on solving the reverse stochastic differential equation  
24 (SDE) or the corresponding probability flow ordinary differential equation (ODE) from a numerical  
25 perspective, which can be applied to any trained score model without extra training. SDE-based  
26 samplers often have better sample quality than ODE-based samplers but they are much slower and  
27 require hundreds if not thousands of function evaluations [1, 11]. ODE-based methods are fully  
28 deterministic and allow for larger steps in discretization. Existing studies work on reducing the  
29 approximation error with less steps [1, 11, 12, 13, 14, 15, 16] but still need more than 10 function  
30 evaluations to generate high-quality samples.

31 **Training-based sampling methods** require extra training including knowledge distillation [17, 18]  
32 and learning the noise schedule [19, 20]. Training-based methods work in the few-step regime with  
33 less than 10 steps. The current SOTA progressive distillation [18] reduces the number of steps down to  
34 4-8 without losing much sample quality. However, it requires progressive training from fine resolution  
35 to coarse resolution. It also breaks in the limit of one function evaluation. DDGAN [21] achieves

36 similar results as progressive distillation by leveraging conditional GAN to model the denoising  
 37 distributions or equivalently the reverse stochastic process, allowing for large denoising steps. LSGM  
 38 [2] accelerates sampling by encoding the data distribution into a smooth latent distribution that is  
 39 close to a Gaussian prior and obtains better image quality with 20 to 100 steps.

40 **Our contributions.** Inspired by the recent success of neural operators [22, 23, 24] in solving  
 41 differential equations, we propose to solve the probability flow ODE of diffusion models from an  
 42 operator learning perspective. We examine the characteristics of the ODE trajectories sampled from  
 43 trained diffusion models [11, 25] and observe a compact energy spectrum. With this observation, we  
 44 propose a diffusion Fourier neural operator (DFNO) with temporal convolution in the Fourier space  
 45 to obtain probability flow trajectories efficiently.

- 46 • DFNO only takes one function evaluation to sample and has better generalization ability  
 47 than distillation methods. With the trajectory information guiding the sampling, DFNO  
 48 achieves the state-of-the-art FID of 5.9 for CIFAR-10 in the one-function-evaluation setting.
- 49 • Temporal convolution blocks in the Fourier space can learn a trajectory as a function of  
 50 time in the Fourier space efficiently. DFNO inherits the discretization invariant property  
 51 from the Fourier neural operator [23] over the temporal dimension. One can train DFNO  
 52 with high-resolution in time for stronger supervision and sample at low-resolution for fast  
 53 inference.
- 54 • Compared to the current SOTA progressive distillation [18], DFNO is easier to train and not  
 55 limited to specific time step scheme. This allows us to learn from a large class of ODE-based  
 56 samplers including training-based methods.

## 57 2 Background

58 We consider the general class of score-based generative models in a unified continuous-time frame-  
 59 work proposed by [11], which includes different variants of diffusion models [26, 25]. We will use  
 60 score-based models interchangeably with the diffusion models. Suppose the data distribution is  $p_{\text{data}}$ .  
 61 The forward pass is a diffusion process  $\{\mathbf{x}(t)\}$  starting from 0 to  $T$  can be expressed as

$$\mathrm{d}\mathbf{x} = f(\mathbf{x}, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w}, \quad (1)$$

62 where  $\mathbf{w}$  is the standard Wiener process, and  $f(\mathbf{x}, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$  and  $g(t) : \mathbb{R} \rightarrow \mathbb{R}$  are the drift  
 63 and diffusion coefficients respectively. Diffusion models choose  $f, g$  such that  $\mathbf{x}(0) \sim p_{\text{data}}$  and  
 64  $\mathbf{x}(T) \sim \mathcal{N}(0, \mathbf{I})$ . Song et al. [11] show that the following probability flow ODE produces the same  
 65 marginal distributions  $p_t(\mathbf{x})$ :

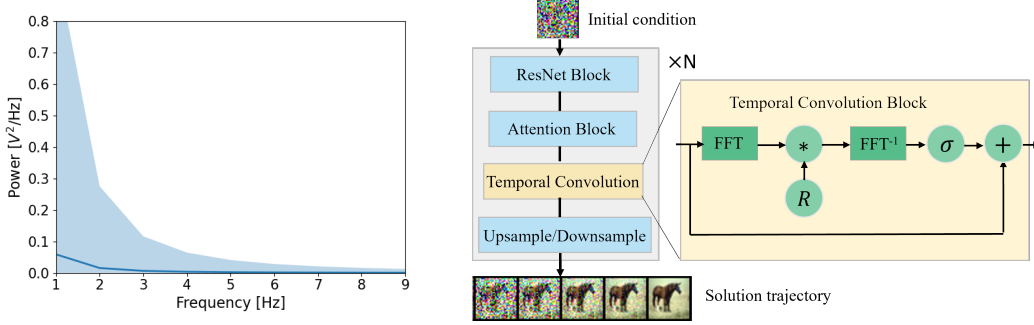
$$\mathrm{d}\mathbf{x} = f(\mathbf{x}, t)\mathrm{d}t - \frac{1}{2}g(t)^2\nabla_{\mathbf{x}}\log p_t(\mathbf{x})\mathrm{d}t. \quad (2)$$

66 The sampling process eventually becomes solving the probability flow ODE 2 from  $T$  to 0 given  
 67 the initial condition  $\mathbf{x}(T)$ . Furthermore,  $f(\mathbf{x}, t)$  often has the affine form  $f(\mathbf{x}, t) = f(t)\mathbf{x}$ , where  
 68  $f : \mathbb{R} \rightarrow \mathbb{R}$ . We can simplified the ODE 2 into a semi-linear ODE. Integrating both sides over time  
 69 gives the explicit form of solution for any  $t < s$ :

$$\mathbf{x}(t) = \phi(t, s)\mathbf{x}(s) - \int_s^t \phi(t, \tau) \frac{g(\tau)^2}{2} \nabla_{\mathbf{x}} \log p_{\tau}(\mathbf{x}) \mathrm{d}\tau, \quad (3)$$

70 where  $\phi(t, s) = \exp\left(\int_s^t f(\tau)\mathrm{d}\tau\right)$ . The ODE is often solved using numerical techniques such as  
 71 Euler discretization [27] or multistep methods [14]. The score function  $\nabla_{\mathbf{x}}\log p_t(\mathbf{x})$  is usually  
 72 parameterized by  $\hat{\epsilon}_{\theta}(\mathbf{x}_t) \approx -\sigma_t\nabla_{\mathbf{x}}\log p_t(\mathbf{x})$ , where  $\sigma_t$  is the noise schedule [11, 25].

73 Neural operators [23, 22] are designed to solve the differential equations fast by learning a parametric  
 74 map between two Banach spaces. They are constructed as a stack of kernel integration layers where  
 75 the kernel function is parameterized by learnable weights. More details are in appendix A.2.



(a) Power spectrum of the ODE trajectories sampled from "DDPM++ cont. (VP)" model trained by [11] on CIFAR10. The mean is computed over all pixel locations and channels of randomly generated trajectories. Most power concentrates in the  $\leq 5$  Hz regime. The shaded region represents maximum and minimum power.

(b) Architecture of DFNO. Blue blocks are commonly used modules in diffusion models. The temporal convolution block first approximates the Fourier transform over temporal dimension via fast Fourier transform. It then multiply it by the filter  $R \in \mathbb{C}^{k \times d_{out} \times d_{in}}$  and do inverse Fourier transform. The shortcut connection restores the high-frequency information without extra cost. Temporal convolution blocks are introduced after every attention block.

Figure 1: Compact spectrum and DFNO architecture design.

### 76 3 Learning the trajectory as a function of time

#### 77 3.1 Problem statement

78 Given any initial condition  $\mathbf{x}(T) \sim \mathcal{N}(0, \mathbf{I})$ , our goal is to learn a neural operator that predicts the  
 79 probability flow trajectory  $\{\mathbf{x}(t)\}_s^0$  with time flowing from  $s$  to 0 defined in equation 3, where the end  
 80 point  $\mathbf{x}(0)$  is the data. Suppose the time domain  $D = [0, s], s > 0$ . Let  $\mathcal{A}$  be the finite-dimensional  
 81 space of initial condition,  $\mathcal{U} = \mathcal{U}(D; \mathbb{R}^d)$  denote the space of the target continuous time functions  
 82 with output value in  $\mathbb{R}^d$ . From the exact solution  $\mathbf{x}(t)$  in equation 3, we know the unique solution  
 83 operator  $G^* : \mathcal{A} \rightarrow \mathcal{U}$  exists and is a weighted integral operator of the score function. We build a  
 84 neural operator  $G_\theta$  parameterized by  $\theta$  to approximate the solution operator  $G^*$  by minimizing:

$$\min_{\theta} \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})} \mathcal{L}(G_\theta(\mathbf{x}_T) - G^*(\mathbf{x}_T)), \quad (4)$$

85 where  $\mathcal{L} : \mathcal{U} \rightarrow \mathbb{R}_+$  is some loss functional such as  $L^p$ -norm for some  $p \geq 1$ .

#### 86 3.2 Compact power spectrum

87 Learning the solution operator  $G^*$  is a challenging task in general. However, we observe that the  
 88 trajectory of probability flow ODEs has a compact energy spectrum over the time dimension and thus  
 89 can be learned efficiently in the Fourier space. Figure 1a visualizes the energy spectrum of ODE  
 90 trajectories sampled from the diffusion model "DDPM++ cont. (VP)" trained by [11] on CIFAR10.  
 91 Appendix A.1 explains the details of the power spectrum.

#### 92 3.3 Temporal convolution block in Fourier space

93 Based on the special characteristic of the ODE trajectory and the integral expression of the exact ODE  
 94 solution in equation 3, we build our temporal convolution block with Fourier integral operator  $\mathcal{K}$  to  
 95 efficiently model the trajectory. Given an integrable function  $u : D \rightarrow \mathbb{R}^{d_{in}}$ , the Fourier transform  
 96 operator  $\mathcal{F}$  is defined as

$$(\mathcal{F}u)_j(k) = \int_D u_j(t) \exp(-2\pi ikt) dt, \quad (5)$$

97 for  $j = 1, \dots, d_{in}$ , where  $i$  is the imaginary unit,  $k$  is the frequency. The Fourier integral operator  
 98  $\mathcal{K}_\phi$  parameterized by  $\phi$  is defined as

$$(\mathcal{K}_\phi u)(t) = \mathcal{F}^{-1}(R_\phi \cdot (\mathcal{F}u))(t) = \int_D (\mathcal{F}^{-1}R_\phi)(t)u(t)dt, \quad (6)$$

99 where  $R_\phi \in \ell^2(\mathbb{Z}; \mathbb{C}^{d_{\text{out}} \times d_{\text{in}}})$  is the Fourier transform of a kernel function parameterized by  $\phi$  that we  
 100 learn from data, and the second equality holds by the convolution theorem. Given an input function  
 101  $u$ , the temporal convolution layer  $\mathcal{P}$  is defined as

$$(\mathcal{P}u)(t) = u(t) + \sigma((\mathcal{K}_\phi u)(t)), \quad (7)$$

102 where  $\sigma$  is a point-wise activation function. Figure 1b demonstrates the implementation details of the  
 103 temporal convolution block.

### 104 3.4 Architecture of DFNO

105 As demonstrated in Figure 1b, the overall architecture of DFNO is similar to the UNet structure that  
 106 is commonly used in diffusion models. We introduce temporal convolution after every attention block.  
 107 Suppose the time resolution is  $M$ . The input noise is repeated  $M$  times in the first block and each  
 108 copy will be mixed with the corresponding time embedding in the ResNet block. More details are  
 109 provided in the appendix A.3.

## 110 4 Experiments

111 We use the Frechet inception distance (FID) [28] to evaluate the quality of generated images. FID is  
 112 computed between 50,000 generated images and CIFAR10 train set with the clean-fid library [29].  
 113 We use the checkpoint of "DDPM++ cont. (VP)" model by [11] trained on CIFAR10 [30]. 1 million  
 114 trajectories are sampled following the corresponding probability flow ODE of the variance preseving  
 115 (VP) SDE to train DFNO. The We use  $L^1$ -norm as the loss functional. Table 1 compares our results  
 116 against the recent training-free and training-based methods. Our method is clearly the best in the  
 117 one-function-evaluation setting, even better than most training-free methods with 10 steps.

Method	Model evaluations	FID score
DFNO (ours)	1	<b>5.92</b>
Knowledge Distillation [17]	1	9.36
Knowledge Distillation (our architecture)	1	8.06
Progressive Distillation [18]	1	9.12
	2	4.51
	4	3.00
GGDM + PRED + TIME [20]	5	13.77
	10	8.23
DDIM [27]	10	13.36
	20	6.84
DPM-solver-2 [15]	12	5.28
DPM-solver-3 [15]	12	6.03
3-DEIS [14]	5	16.09
	10	4.17

Table 1: Comparison of fast sampling methods on CIFAR-10 for diffusion models in the literature.

118 **Ablation study** The results of our ablation study are reported in the Table 2. We observe that 4-step  
 119 is much better than 1-step. Quadratic time step scheme and  $1/\sqrt{\sigma_t}$  loss weighting also improve the  
 120 sample quality. The setting of each ablation is explained in appendix A.4.

## 121 5 Conclusion

122 In this paper, we propose diffusion Fourier neural operator (DFNO), a training-based fast sampling  
 123 method for diffusion models. DFNO leverages the compact power spectrum characteristic of the  
 124 probability flow ODE trajectory and models it efficiently with Fourier integral operator. Experiments  
 125 show that our method achieves the best performance in one-function-evaluation setting.

## References

- 126
- 127 [1] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of  
128 diffusion-based generative models. *arXiv preprint arXiv:2206.00364*, 2022.
- 129 [2] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space.  
130 *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021.
- 131 [3] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis.  
132 *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- 133 [4] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.  
134 High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE*  
135 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- 136 [5] Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geo-  
137 metric diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*,  
138 2022.
- 139 [6] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile  
140 diffusion model for audio synthesis. In *ICLR*, 2021.
- 141 [7] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-  
142 tts: A diffusion probabilistic model for text-to-speech. In *International Conference on Machine*  
143 *Learning*, pages 8599–8608. PMLR, 2021.
- 144 [8] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar.  
145 Diffusion models for adversarial purification. In *International Conference on Machine Learning*  
146 *(ICML)*, 2022.
- 147 [9] Nicholas Carlini, Florian Tramèr, J Zico Kolter, et al. (certified!!) adversarial robustness for  
148 free! *arXiv preprint arXiv:2206.10550*, 2022.
- 149 [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
150 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications*  
151 *of the ACM*, 63(11):139–144, 2020.
- 152 [11] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and  
153 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*  
154 *preprint arXiv:2011.13456*, 2020.
- 155 [12] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In  
156 *ICLR*, 2021.
- 157 [13] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the  
158 optimal reverse variance in diffusion probabilistic models. In *International Conference on*  
159 *Learning Representations*, 2021.
- 160 [14] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential  
161 integrator. *arXiv preprint arXiv:2204.13902*, 2022.
- 162 [15] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A  
163 fast ode solver for diffusion probabilistic model sampling in around 10 steps. *arXiv preprint*  
164 *arXiv:2206.00927*, 2022.
- 165 [16] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models  
166 on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
- 167 [17] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for  
168 improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021.
- 169 [18] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models.  
170 In *International Conference on Learning Representations*, 2021.

- 171 [19] Max WY Lam, Jun Wang, Rongjie Huang, Dan Su, and Dong Yu. Bilateral denoising diffusion  
172 models. *arXiv preprint arXiv:2108.11514*, 2021.
- 173 [20] Daniel Watson, William Chan, Jonathan Ho, and Mohammad Norouzi. Learning fast samplers  
174 for diffusion models by differentiating through sample quality. In *International Conference on*  
175 *Learning Representations*, 2021.
- 176 [21] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma  
177 with denoising diffusion gans. In *International Conference on Learning Representations*, 2021.
- 178 [22] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya,  
179 Andrew Stuart, and Anima Anandkumar. Neural operator: Graph kernel network for partial  
180 differential equations. *arXiv preprint arXiv:2003.03485*, 2020.
- 181 [23] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya,  
182 Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differen-  
183 tial equations. *arXiv preprint arXiv:2010.08895*, 2020.
- 184 [24] Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Azizzadenesheli, Kaushik Bhattacharya,  
185 Andrew Stuart, and Anima Anandkumar. Neural operator: Learning maps between function  
186 spaces. *arXiv preprint arXiv:2108.08481*, 2021.
- 187 [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
188 *in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 189 [26] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-  
190 vised learning using nonequilibrium thermodynamics. In *International Conference on Machine*  
191 *Learning*, pages 2256–2265. PMLR, 2015.
- 192 [27] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
193 *preprint arXiv:2010.02502*, 2020.
- 194 [28] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
195 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*  
196 *neural information processing systems*, 30, 2017.
- 197 [29] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties  
198 in gan evaluation. In *CVPR*, 2022.
- 199 [30] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
200 2009.

## 201 **A Appendix**

### 202 **A.1 Energy spectrum**

203 The discrete-time Fourier transform of the signal  $x(t)$  with period  $T$  is given by

$$X_j = \sum_{i=1}^N x(t_i) \exp\left(-\frac{2\pi}{T} j t_i\right), \quad (8)$$

204 where  $t_i = \frac{iT}{N}$ .  $\frac{j}{T}$  is the frequency.  $j$  is called the frequency mode. Let  $\Delta = \frac{1}{N}$  be the time step. The  
205 spectrum is defined as the product of Fourier transform of  $x$  with its conjugate:

$$S_j = \frac{2\Delta^2}{T} X_j X_j^*, \quad (9)$$

206 where  $X_j^*$  is the complex conjugate. In practice, the statistics are computed over all pixel locations  
207 and channels of randomly generated trajectories.  $T = 1$  and the sampling frequency is 1000 Hz to  
208 avoid aliasing. We observe that most power concentrates in the regime where the frequency mode is  
209 less than 5.

210 **A.2 Background: neural operators**

211 Let  $\mathcal{A}$  and  $\mathcal{U}$  be two Banach spaces and  $G : \mathcal{A} \rightarrow \mathcal{U}$  be a non-linear map. Suppose we have a finite  
 212 collection of data  $\{a_i, u_i\}_{i=1}^N$  where  $a_i \sim \mu$  are i.i.d. samples from the distribution  $\mu$  supported on  $\mathcal{A}$   
 213 and  $u_i = G(a_i)$ . Neural operators aim to learn  $G_\phi$  parameterized by  $\phi$  to approximate  $G$  from the  
 214 observed data by minimizing the empirical risk given by

$$\min_{\phi} \mathbb{E}_{a \sim \mu} \|G(a) - G_\phi(a)\|_{\mathcal{U}} \approx \min_{\phi} \frac{1}{N} \sum_{i=1}^N \|u_i - G_\phi(a_i)\|_{\mathcal{U}}. \quad (10)$$

215 The architecture of neural operators is constructed as a stack of kernel integration layers where the  
 216 kernel function is parameterized by learnable weights.

217 **A.3 Architecture detail**

218 The overall architecture is similar to the UNet structure used in diffusion models. On top of that,  
 219 we introduce temporal convolution after each attention layer and replace all the Conv2d layers  
 220 with Conv3d. So the intermediate feature map will have an additional time dimension compared to  
 221 standard UNet. Suppose the output time resolution is  $M$ . The input noise has a shape of  $(B, C, H, W)$   
 222 where  $B$  is the batch size,  $C$  is the number of channels,  $H, W$  are the height and width. In the first  
 223 convolution block, the input noise will be repeated  $M$  times as the initial trajectory. We add the time  
 224 embeddings of  $M$  steps in each ResNet block.

225 **A.4 Ablation detail**

226 For all results in the same column, we keep all the other settings the same except for the control  
 227 variable. The left column reports the ablation on the temporal resolution with default setting of  
 228 uniform time steps and  $1/\sqrt{\sigma_t}$  loss weighting. The middle column compares uniform and quadratic  
 229 time step scheme with the default setting of resolution 4 and  $1/\sqrt{\sigma_t}$  loss weighting. The right column  
 compares different loss weightings with the default setting of resolution 4 and uniform time steps.

Resolution	FID score	Time step scheme	FID score	Loss weighting	FID score
1	7.74	uniform	6.21	uniform weights	7.8
4	6.21	quadratic	5.92	$1/\sqrt{\sigma_t}$	7.2
6	6.17				

Table 2: Left: ablation on the temporal resolution; Middle: ablation on the time step scheme; Right: ablation on the loss weighting. The numbers are clean FID score [29].

230

231 **A.5 Model complexity and inference Time**

232 "DDPM++cont. (VP)" architecture has 106,632,579 parameters. The DFNO with resolution 4 has  
 233 114,562,435 parameters. The computation cost of one model evaluation of DFNO is 2 times as that  
 234 of the standard score model when the batchsize is 2 on 16G-V100.