

CONTENTS

| | |
|---|----------|
| A Score matching framework | 2 |
| A.1 Forward SDE | 2 |
| A.1.1 VESDE and VPSDE | 3 |
| A.2 Backward SDE and ODE | 4 |
| A.3 Approximation of the score function | 5 |
| A.4 Summary | 5 |
| B Fast solvers | 6 |
| C Distillation techniques | 6 |
| D Analysis | 8 |
| D.1 Properties of the approximated score function | 8 |
| D.2 Properties of the distilled diffusion models | 9 |

Appendix

A SCORE MATCHING FRAMEWORK

This section is devoted to the derivation of the score based models framework, which can be divided into three parts: forward, backward processes and score approximation.

A.1 FORWARD SDE

Forward SDE can be written as

$$dx = f(x, t)dt + g(t)d\omega_t, \quad (1)$$

where ω_t is the standard Wiener process, which can be roughly understood as $d\omega_t \sim \sqrt{dt} \cdot \mathcal{N}(0, T)$. This equation is usually solved on the interval $t \in [0, 1]$. In the score matching it is assumed that the first term is linear $f(x, t) = f(t)x$. It can be shown that the perturbation kernels of this SDE have the Gaussian distribution:

$$p_t(x_t | x_0) = \mathcal{N}(x_t | \mu(x_0, t), \sigma^2(t)\mathbf{I}). \quad (2)$$

We use the notation $x_t = x(t)$ for simplicity. Here x_0 is the initial condition of the forward SDE (10). In terms of deep learning this is a sample from a dataset, $x_0 \sim p_{\text{data}}(x_0)$. The parameters, $\mu(x_0, t)$ and $\sigma^2(t)$, can be found using the Ito formula.

Proposition 1. *Parameters of the perturbation kernels (2) caused by the forward SDE (10) can be found from solution of the following ODEs:*

$$\begin{cases} d\mathbb{E}_x x = \mathbb{E}_x f(x, t)dt, \\ \mathbb{E}_x x(0) = x_0. \end{cases} \quad (3)$$

$$\begin{cases} d\mathbb{D}_x x = (2\mathbb{E}_x [xf(x, t) - f(x, t)\mathbb{E}_x x] + g(t)^2) dt, \\ \mathbb{D}_x x(0) = 0. \end{cases}$$

Proof. Let suppose that $\phi(x, t) : \mathbb{R}^n \rightarrow \mathbb{R}$ is an arbitrary twice differentiable function, where x is a solution of the (10). Thus, we can write the Ito formula for the ϕ :

$$\begin{aligned} d\phi &= \frac{\partial \phi}{\partial t} dt + \sum_{i=1}^n \frac{\partial \phi}{\partial x_i} dx_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \phi}{\partial x_i \partial x_j} dx_i dx_j + \dots \\ &= \frac{\partial \phi}{\partial t} dt + \sum_{i=1}^n \frac{\partial \phi}{\partial x_i} f^i dt + \frac{1}{2} g(t)^2 \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \phi}{\partial x_i \partial x_j} dt + g(t) \sum_{i=1}^n \frac{\partial \phi}{\partial x_i} d\omega_t^i, \end{aligned} \quad (4)$$

where f^i denotes the i -th component of the vector-valued function f . Note that we omitted the arguments for simplicity. Then we can take the expectation from both parts:

$$\begin{aligned} d\mathbb{E}_x \phi &= \mathbb{E}_x \left(\frac{\partial \phi}{\partial t} \right) dt + \frac{1}{2} g(t)^2 \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_x \left(\frac{\partial^2 \phi}{\partial x_i \partial x_j} \right) dt \\ &\quad + \sum_{i=1}^n \mathbb{E}_x \left(\frac{\partial \phi}{\partial x_i} f^i \right) dt + g(t) \sum_{i=1}^n \underbrace{\mathbb{E}_x \left(\frac{\partial \phi}{\partial x_i} d\omega_t^i \right)}_{\mathbb{E}_x \left(\frac{\partial \phi}{\partial x_i} \right) \mathbb{E}(d\omega_t^i) = 0} \\ &= \mathbb{E}_x \left(\frac{\partial \phi}{\partial t} \right) dt + \frac{1}{2} g(t)^2 \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_x \left(\frac{\partial^2 \phi}{\partial x_i \partial x_j} \right) dt + \sum_{i=1}^n \mathbb{E}_x \left(\frac{\partial \phi}{\partial x_i} f^i \right) dt. \end{aligned} \quad (5)$$

Let say that $\phi = x^k$, then

$$d\mathbb{E}_x x^k = \mathbb{E}_x \left(\underbrace{\frac{\partial x^k}{\partial t}}_0 \right) dt + \frac{1}{2} g(t)^2 \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_x \left(\underbrace{\frac{\partial^2 x^k}{\partial x_i \partial x_j}}_0 \right) dt + \sum_{i=1}^n \mathbb{E}_x \left(\frac{\partial x^k}{\partial x_i} f^i \right) dt = \mathbb{E}_x f^k dt. \quad (6)$$

Finally, we obtain the first equation in the (3)

$$d\mathbb{E}_{\mathbf{x}} \mathbf{x} = \mathbb{E}_{\mathbf{x}} \mathbf{f} dt \quad (7)$$

If we put $\phi = \mathbf{x}_k^2 - (\mathbb{E}_{\mathbf{x}} \mathbf{x}_k)^2$, then we will obtain the second equation.

$$\begin{aligned} d\mathbb{E}_{\mathbf{x}} [\mathbf{x}_k^2 - (\mathbb{E}_{\mathbf{x}} \mathbf{x}_k)^2] &= \mathbb{E}_{\mathbf{x}} \left(\underbrace{\frac{\partial [\mathbf{x}_k^2 - (\mathbb{E}_{\mathbf{x}} \mathbf{x}_k)^2]}{\partial t}}_{=-2\mathbb{E}_{\mathbf{x}} \mathbf{x}_k \cdot \mathbb{E}_{\mathbf{x}} \mathbf{f}_k} \right) dt + \underbrace{\frac{1}{2} g(t)^2 \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_{\mathbf{x}} \left(\frac{\partial^2 [\mathbf{x}_k^2 - (\mathbb{E}_{\mathbf{x}} \mathbf{x}_k)^2]}{\partial x_i \partial x_j} \right) dt}_{=2} \\ &\quad + \underbrace{\sum_{i=1}^n \mathbb{E}_{\mathbf{x}} \left(\frac{\partial [\mathbf{x}_k^2 - (\mathbb{E}_{\mathbf{x}} \mathbf{x}_k)^2]}{\partial x_i} \mathbf{f}^i \right) dt}_{2\mathbb{E}_{\mathbf{x}}(\mathbf{x}_k \mathbf{f}_k)}. \end{aligned} \quad (8)$$

Finally, we have the following

$$d\mathbb{D}_{\mathbf{x}} \mathbf{x} = (2\mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{f}(\mathbf{x}, t) - \mathbf{f}(\mathbf{x}, t) \mathbb{E}_{\mathbf{x}} \mathbf{x}] + g(t)^2) dt. \quad (9)$$

□

A.1.1 VESDE AND VPSDE

Two the most well known types of the forward SDE are VESDE (Variance Exploding SDE) and VPSDE (Variance Preserving SDE). They can be formulated as the following:

$$dx = f(t)xdt + g(t)d\omega_t, \quad (10)$$

VPSDE.

$$f(t) = \frac{1}{2} \frac{d \log \alpha_t}{dt}, \quad g(t) = \sqrt{\left(-\frac{d \log \alpha_t}{dt} \right)}, \quad (11)$$

where α_t is a bounded decreasing sequence from $\alpha_0 = 1$ to $\alpha_1 = 0$. The parameters of the perturbation kernels can be found from the (3)

$$\begin{aligned} d\mathbb{E}_{\mathbf{x}} \mathbf{x} &= \mathbb{E}_{\mathbf{x}} \frac{1}{2} \frac{d \log \alpha_t}{dt} \mathbf{x} dt = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \mathbf{x} \frac{d \alpha_t}{\alpha_t}, \\ \int_0^\tau \frac{d\mathbb{E}_{\mathbf{x}} \mathbf{x}}{\mathbb{E}_{\mathbf{x}} \mathbf{x}} &= \int_0^\tau \frac{1}{2} \frac{d \alpha_t}{\alpha_t} \rightarrow \mathbb{E}_{\mathbf{x}} \mathbf{x}(\tau) = \sqrt{\alpha_\tau} \mathbf{x}_0. \end{aligned} \quad (12)$$

$$\begin{aligned} d\mathbb{D}_{\mathbf{x}} \mathbf{x} &= \left(2\mathbb{E}_{\mathbf{x}} \left[\mathbf{x}^2 \frac{1}{2} \frac{d \log \alpha_t}{dt} - \frac{1}{2} \mathbf{x} \frac{d \log \alpha_t}{dt} \mathbb{E}_{\mathbf{x}} \mathbf{x} \right] - \frac{d \log \alpha_t}{dt} \right) dt, \\ d\mathbb{D}_{\mathbf{x}} \mathbf{x} &= \left(\frac{d \log \alpha_t}{dt} \mathbb{E}_{\mathbf{x}} \mathbf{x}^2 - \frac{d \log \alpha_t}{dt} (\mathbb{E}_{\mathbf{x}} \mathbf{x})^2 - \frac{d \log \alpha_t}{dt} \right) dt, \end{aligned} \quad (13)$$

$$d\mathbb{D}_{\mathbf{x}} \mathbf{x} = (\mathbb{D}_{\mathbf{x}} \mathbf{x} - 1) \frac{d \alpha_t}{\alpha_t} \rightarrow \int_0^\tau \frac{d\mathbb{D}_{\mathbf{x}} \mathbf{x}}{\mathbb{D}_{\mathbf{x}} \mathbf{x} - 1} = \int_0^\tau \frac{d \alpha_t}{\alpha_t},$$

$$\mathbb{D}_{\mathbf{x}} \mathbf{x}(\tau) - 1 = (\mathbb{D}_{\mathbf{x}} \mathbf{x}(0) - 1) \alpha_\tau \rightarrow \mathbb{D}_{\mathbf{x}} \mathbf{x}(\tau) = 1 - \alpha_\tau.$$

VESDE.

$$f(t) = 0, \quad g(t) = \sqrt{\frac{ds_t^2}{dt}}, \quad (14)$$

where s_t is an unbounded increasing sequence, $s_0 = 0, s_1 = c$

$$\begin{aligned} d\mathbb{E}_{\mathbf{x}} \mathbf{x} &= 0 \rightarrow \mathbb{E}_{\mathbf{x}} \mathbf{x}(\tau) = \mathbf{x}_0, \\ \int_0^\tau d\mathbb{D}_{\mathbf{x}} \mathbf{x} &= \int_0^\tau ds_t^2 \rightarrow \mathbb{D}_{\mathbf{x}} \mathbf{x}(\tau) = s_\tau^2. \end{aligned} \quad (15)$$

To summarize, the transition kernels take the following form:

$$p_t(\mathbf{x}_t | \mathbf{x}_0) = \begin{cases} \mathbf{VPSDE} : \mathcal{N}(\mathbf{x}_t | \sqrt{\alpha_t} \mathbf{x}_0, 1 - \alpha_t), & \alpha_0 = 1, \alpha_1 = 0, \\ \mathbf{VESDE} : \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, s_t^2), & s_0 = 0, s_1 = c. \end{cases} \quad (16)$$

Here we deal with the Variance-Exploding scheme. So, our main equations defining the noising process are the following:

$$\begin{aligned} dx &= \sqrt{\frac{ds_t^2}{dt}} d\mathbf{w}_t = d\mathbf{w}_{s_t^2}, \\ p_t(\mathbf{x}_t | \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, s_t^2). \end{aligned} \quad (17)$$

A.2 BACKWARD SDE AND ODE

A remarkable result from Anderson (1982) states that the reverse of a diffusion process is also a diffusion process, running backwards in time and given by the reverse-time SDE:

$$dx = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\mathbf{w}_t, \quad (18)$$

In the literature an explicit form of the score function is defined by noise, applied to a clean image, or clean image itself. We are working with the latter one, which can be obtained by Tweedie's formula.

Proposition 2. (*Tweedie's formula*) *The score function produced by Eq. (17) has the following form*

$$\begin{aligned} \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) &= -\frac{1}{s_t^2} (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}]), \\ \mathbb{E}[\mathbf{x}_0 | \mathbf{x}] &= \int_{\mathbb{R}^n} \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) d\mathbf{x}_0. \end{aligned} \quad (19)$$

Proof. Firstly, according to the chain rule, the following is truth

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = \frac{\nabla_{\mathbf{x}} p_t(\mathbf{x})}{p_t(\mathbf{x})}. \quad (20)$$

Then, we have to deal with $\nabla_{\mathbf{x}} p_t(\mathbf{x})$:

$$\begin{aligned} p_t(\mathbf{x}) &= \int_{\mathbb{R}^n} p_t(\mathbf{x} | \mathbf{x}_0) p_{\text{data}}(\mathbf{x}_0) d\mathbf{x}_0, \quad p_t(\mathbf{x} | \mathbf{x}_0) = \mathcal{N}(\mathbf{x} | \mathbf{x}_0, s_t^2), \\ \nabla_{\mathbf{x}} p_t(\mathbf{x}) &= \nabla_{\mathbf{x}} \int_{\mathbb{R}^n} p_t(\mathbf{x} | \mathbf{x}_0) p_{\text{data}}(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \int_{\mathbb{R}^n} p_{\text{data}}(\mathbf{x}_0) \nabla_{\mathbf{x}} p_t(\mathbf{x} | \mathbf{x}_0) d\mathbf{x}_0 \\ &= \int_{\mathbb{R}^n} p_{\text{data}}(\mathbf{x}_0) \left[2\pi s_t^2 \right]^{-\frac{n}{2}} \exp \frac{\|\mathbf{x} - \mathbf{x}_0\|_2^2}{-2s_t^2} d\mathbf{x}_0 \\ &= \int_{\mathbb{R}^n} p_{\text{data}}(\mathbf{x}_0) \left[2\pi s_t^2 \right]^{-\frac{n}{2}} \exp \frac{\|\mathbf{x} - \mathbf{x}_0\|_2^2}{-2s_t^2} \nabla_{\mathbf{x}} \left(\frac{\|\mathbf{x} - \mathbf{x}_0\|_2^2}{-2s_t^2} \right) d\mathbf{x}_0 \\ &= \int_{\mathbb{R}^n} p_{\text{data}}(\mathbf{x}_0) p_t(\mathbf{x} | \mathbf{x}_0) \nabla_{\mathbf{x}} \left(\frac{\|\mathbf{x} - \mathbf{x}_0\|_2^2}{-2s_t^2} \right) d\mathbf{x}_0 \\ &= \int_{\mathbb{R}^n} p_{\text{data}}(\mathbf{x}_0) p_t(\mathbf{x} | \mathbf{x}_0) \left[\frac{\mathbf{x}_0 - \mathbf{x}}{s_t^2} \right] d\mathbf{x}_0 \\ &= -\frac{1}{s_t^2} \left(\int_{\mathbb{R}^n} \mathbf{x} p_{\text{data}}(\mathbf{x}_0) p_t(\mathbf{x} | \mathbf{x}_0) d\mathbf{x}_0 - \int_{\mathbb{R}^n} \mathbf{x}_0 p_{\text{data}}(\mathbf{x}_0) p_t(\mathbf{x} | \mathbf{x}_0) d\mathbf{x}_0 \right) \\ &= -\frac{1}{s_t^2} \left(\mathbf{x} \int_{\mathbb{R}^n} p_t(\mathbf{x}, \mathbf{x}_0) d\mathbf{x}_0 - \int_{\mathbb{R}^n} \mathbf{x}_0 p_t(\mathbf{x}, \mathbf{x}_0) d\mathbf{x}_0 \right) \\ &= -\frac{1}{s_t^2} \left(\mathbf{x} p_t(\mathbf{x}) - \int_{\mathbb{R}^n} \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) p_t(\mathbf{x}) d\mathbf{x}_0 \right) \\ &= -\frac{p_t(\mathbf{x})}{s_t^2} \left(\mathbf{x} - \int_{\mathbb{R}^n} \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) d\mathbf{x}_0 \right). \end{aligned} \quad (21)$$

So, finally we have the following:

$$\begin{aligned}\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) &= \frac{\nabla_{\mathbf{x}} p_t(\mathbf{x})}{p_t(\mathbf{x})} = -\frac{1}{s_t^2} \left(\mathbf{x} - \int_{\mathbb{R}^n} \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) d\mathbf{x}_0 \right), \\ \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) &= -\frac{1}{s_t^2} (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}]).\end{aligned}\quad (22)$$

□

The backward SDE has the following form, taking into account the obtained formula and Variance-Exploding scheme

$$d\mathbf{x} = 2(\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}]) \frac{ds_t}{s_t} + d\mathbf{w}_{s_t^2}, \quad (23)$$

where $d\mathbf{w}_{s_t^2}$ means the standard Wiener process with changed variable, roughly speaking, $d\mathbf{w}_{s_t^2} \sim \sqrt{ds_t^2} \cdot \mathcal{N}(0, 1)$.

However, there is a backward ODE, providing the same marginal probability densities, p_t , as the backward SDE, which in practice usually gives better results in terms of quality and efficiency.

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2} g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt, \quad (24)$$

or in the case of VE scheme and Tweedie's formula

$$d\mathbf{x} = (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}]) \frac{ds_t}{s_t}. \quad (25)$$

A.3 APPROXIMATION OF THE SCORE FUNCTION

The conditional expectation, $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}]$, cannot be calculated analytically since the conditional distribution, $p_t(\mathbf{x}_0 | \mathbf{x})$, is unavailable. Thus, we have to utilize a neural network to approximate it. To this end, a least square regression problem is solved.

$$\begin{aligned}\mathcal{L}(\theta) &= \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_t(\mathbf{x}|\mathbf{x}_0)} \|\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}, t)\|_2^2, \\ \mathcal{L}(\theta) &\rightarrow \min_{\theta}\end{aligned}\quad (26)$$

Then, we approximate the conditional expectation using a trained neural network $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}] \approx \mathbf{x}_{\theta^*}(\mathbf{x}, t)$.

A.4 SUMMARY

Below we summarize three main parts of the score-based models.

| | |
|--|--|
| | 1. Forward process |
| | $dx = d\mathbf{w}_{s_t^2},$ |
| | $p_t(\mathbf{x}_t \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t \mathbf{x}_0, s_t^2).$ |
| | 2. Backward process |
| | $d\mathbf{x} = -\frac{1}{2} \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) ds_t^2,$ |
| | $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\frac{1}{s_t^2} (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 \mathbf{x}]).$ |
| | 3. Score approximation |
| | $\mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_t(\mathbf{x} \mathbf{x}_0)} \ \mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}, t)\ _2^2 \rightarrow \min_{\theta},$ |
| | $\mathbb{E}[\mathbf{x}_0 \mathbf{x}] \approx \mathbf{x}_{\theta^*}(\mathbf{x}, t).$ |

B FAST SOLVERS

In this paper we consider two the most popular types of fast solvers, EDM and DPM-Solver, which are currently the state-of-the-art in the image generation.

DPM-Solver. DPM-Solver is heavily based on the well-known in the ODE literature exponential integrators methods. Firstly, we can rewrite Eq. (24) using Proposition 2, reparameterization through the noise $\epsilon(x, t) = \frac{1}{s_t} (\mathbf{x} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}])$ and the semi-linear structure of the diffusion ODE.

$$\begin{aligned} d\mathbf{x} &= \mathbf{f}(\mathbf{x}, t)dt - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})dt \\ &= f(t)\mathbf{x}dt + \frac{g(t)^2}{2s_t}\epsilon(x, t)dt. \end{aligned} \quad (28)$$

The exponential integrator method enter an auxiliary variable, in general case called as *transition matrix*, $\phi(s, t)$. Denoting s, τ as the starting and ending points, respectively.

$$\begin{cases} d(\phi(t, \tau)\mathbf{x}) = \phi(t, \tau) \frac{g(t)^2}{2s_t} \epsilon(x, t)dt, \\ \frac{\partial \phi(r, t)}{\partial r} = f(r)\phi(r, t), \phi(t, t) = 1. \end{cases} \quad (29)$$

$$\begin{cases} \int_s^\tau d(\phi(t, \tau)\mathbf{x}) = \int_s^\tau \phi(t, \tau) \frac{g(t)^2}{2s_t} \epsilon(x, t)dt, \\ \phi(s, t) = e^{\int_s^t f(r)dr}. \end{cases} \quad (30)$$

Solving the above system of equations, one can obtain the following:

$$\mathbf{x}_\tau = e^{\int_s^\tau f(r)dr} \mathbf{x}_s + \int_s^\tau \left(e^{\int_t^\tau f(r)dr} \frac{g(t)^2}{2s_t} \epsilon(x, t) \right) dt. \quad (31)$$

Transition matrix can be calculated analytically for a chosen scheme (VP or VE), however an exact form of $\epsilon(x, t)$ is unknown, which makes the second integral intractable. The authors of DPM-solver use the Taylor expansion of $\epsilon(x, t)$ w.r.t t .

EDM. This framework proposes to reformulate the previous ODE/SDE (24, 18) using $\mathbf{f}(\mathbf{x}, t) = 0$ and $g(t) = \sqrt{2t}$.

$$\begin{aligned} d\mathbf{x} &= -t \nabla_{\mathbf{x}} \log p_t(\mathbf{x})dt, \\ d\mathbf{x} &= -t \nabla_{\mathbf{x}} \log p_t(\mathbf{x})dt + \beta(t)t^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})dt + \sqrt{2\beta(t)}tdt. \end{aligned} \quad (32)$$

To approximate the above equations authors utilize Heun's 2nd order method and the following schedule design:

$$t_{i < N} = \left(t_{\max}^{1/\rho} + \frac{i}{N-1} \left(t_{\min}^{1/\rho} - t_{\max}^{1/\rho} \right) \right)^\rho, \quad t_N = 0. \quad (33)$$

This choice was motivated by truncation error analysis. Here N is the number of sampling steps, $\rho = 7$ and $t_{\max} = 80, t_{\min} = 0.002$.

At this moment, the described solvers are the state-of-the-art in the field of image generation. Thus, in our experiments we adhere to the settings described in the corresponding papers.

C DISTILLATION TECHNIQUES

In this work we consider three the most popular distillation techniques: knownledge, progressive and consistency distillation.

Consistency distillation. This type of distillation introduces a new type of generative models called *consistency models*. Assume that $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x}_0)$ is a sample from a real data distribution and $p_t(\mathbf{x}_t | \mathbf{x}_0)$ is transition kernels caused by the forward process (??). Then consistency model is $f(\mathbf{x}_t, t) = \mathbf{x}_0, \forall t$.

Transition from diffusion to consistency model occurs due to the self-consistency property.

$$f(\mathbf{x}_{t_n}, t_n) = f(\mathbf{x}_{t_{n-1}}, t_{n-1}) = \dots = f(\mathbf{x}_{t_0}, t_0). \quad (34)$$

In practice we enforce this property solving the following optimization problem:

$$\mathcal{L}^{\text{CD}}(\theta) = \mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0) p_{t_n}(\mathbf{x} | \mathbf{x}_0) q_{t_{n-1}}^\phi(\hat{\mathbf{x}} | \mathbf{x})} \|\mathbf{x}_\theta(\hat{\mathbf{x}}, t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 \rightarrow \min_{\theta}, \quad (35)$$

where $q_{t_{n-1}}^\phi(\hat{\mathbf{x}} | \mathbf{x})$ is a single step of ODE solver, from t_n to t_{n-1} , using a fixed pretrained diffusion model with parameters ϕ , $\theta^- = \text{EMA}(\theta)$ and $\mathbf{x}_\theta(\mathbf{x}, t_n)$ is initialized from a pretrained diffusion model that approximates the conditional expectation 2. Moreover, the consistency model is parameterized as $\mathbf{x}_\theta(\mathbf{x}, t_1) = \mathbf{x}$.

It was shown that, under some regularity conditions, the difference between the estimated and real consistency models becomes infinitely small.

$$\sup_{n, \mathbf{x}} \|\mathbf{x}_\theta(\mathbf{x}, t_n) - f(\mathbf{x}, t_n)\| = O(\Delta t^p), \quad (36)$$

where $p + 1$ is the local error of ODE solver and $\Delta t = \max_n |t_{n+1} - t_n|$.

After the model is distilled, stochastic sampling is used to produce images, one step of which has the following form:

$$\begin{aligned} \mathbf{x}_{t_n} &\sim p_t(\mathbf{x}_{t_n} | \mathbf{x}_0), \\ \mathbf{x}_0 &= \mathbf{x}_\theta(\mathbf{x}, t_n). \end{aligned} \quad (37)$$

Progressive distillation. This strategy consists of several stages at each of which distilled model aims to halve the number of steps as a teacher model. A single stage can be written as the following:

$$\mathcal{L}^{\text{PD}}(\theta) = \mathbb{E}_{n \sim \mathcal{U}(3, N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0) p_{t_n}(\mathbf{x} | \mathbf{x}_0) q_{t_{n-2}}^\phi(\hat{\mathbf{x}} | \mathbf{x})} \|\mathbf{x}_\theta(\mathbf{x}, t_n) - \mathbf{x}'\|_2^2 \rightarrow \min_{\theta}. \quad (38)$$

Here $q_{t_{n-2}}^\phi$ is two steps ($t_n \rightarrow t_{n-1} \rightarrow t_{n-2}$) of the DDIM sampler using the teacher model with parameters ϕ . Concretely,

$$\begin{aligned} q_{t_{n-2}}^\phi(\hat{\mathbf{x}} | \mathbf{x}) &= \delta(\hat{\mathbf{x}} - \mathbf{x}_{t_{n-2}}(\mathbf{x})), \\ \mathbf{x}_{t_{n-2}}(\mathbf{x}) &= \mathbf{x}_\phi(\mathbf{x}_{t_{n-1}}) + s_{t_{n-2}} \frac{\mathbf{x}_{t_{n-1}} - \mathbf{x}_\phi(\mathbf{x}_{t_{n-1}})}{s_{t_{n-1}}}, \\ \mathbf{x}_{t_{n-1}} &= \mathbf{x}_\phi(\mathbf{x}) + s_{t_{n-1}} \frac{\mathbf{x} - \mathbf{x}_\phi(\mathbf{x})}{s_{t_n}}, \end{aligned} \quad (39)$$

where \mathbf{x}_ϕ is the teacher model. Note that to obtain this discretization we use VE scheme (27) and the Euler integration rule. The target for distilled model, \mathbf{x}' , can be derived from the consideration that during distillation we aim to approximate two steps of the DDIM sampler in only one step.

$$\begin{aligned} \mathbf{x}'_{t_{n-2}} &= \mathbf{x}' + s_{t_{n-2}} \frac{\mathbf{x} - \mathbf{x}'}{s_{t_n}} = \mathbf{x}_{t_{n-2}}, \\ \rightarrow \mathbf{x}' &= \frac{\mathbf{x}_{t_{n-2}} - s_{t_{n-2}}/s_{t_n} \mathbf{x}}{1 - s_{t_{n-2}}/s_{t_n}}. \end{aligned} \quad (40)$$

After the model is distilled, the DDIM solver is used to generate samples.

Knowledge distillation. Knowledge distillation is the most simplest technique among considered, but at the same time the most computationally expensive. In this case the optimization problem formulates as follows:

$$\mathcal{L}^{\text{KD}}(\theta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0) p_{t_N}(\mathbf{x} | \mathbf{x}_0) q_{t_0}^\phi(\hat{\mathbf{x}} | \mathbf{x})} \|\mathbf{x}_\theta(\mathbf{x}, t_N) - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\theta}. \quad (41)$$

Here $q_{t_0}^\phi(\hat{\mathbf{x}} | \mathbf{x})$ is N steps of ODE solver, from t_N to 0. Beyond computational complexity, the distilled model loses its ability to multi-step generation which leads to a quality degradation. To fix

this problem, we consider the following modification (Multi-Step Knowledge Distillation) of the knowledge distillation:

$$\mathcal{L}^{\text{MSKD}}(\theta) = \mathbb{E}_{n \sim \mathcal{U}(2, N)} p_{\text{data}}(\mathbf{x}_0) p_{t_N}(\mathbf{x} | \mathbf{x}_0) q_{t_0}^\phi(\hat{\mathbf{x}} | \mathbf{x}) q_{t_n}^\theta(\mathbf{x}' | \mathbf{x}) \|\mathbf{x}_\theta(\mathbf{x}', t_n) - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\theta}. \quad (42)$$

In this case we additionally train the model on the intermediate timesteps, t_n , and not only on the last one, t_N . This modification allows us to utilize ODE solver at the inference phase. In our experiments we use 1-6 steps of DDIM sampler.

D ANALYSIS

D.1 PROPERTIES OF THE APPROXIMATED SCORE FUNCTION

To understand the properties of the optimized network and, as a consequence, the approximated score function, we can obtain the optimal solution of Eq. (26) analytically.

Proposition 3. (*Karras*) *The optimal solution to the least squares regression problem Eq. (26) has the following form:*

$$\mathbf{x}_\theta(\mathbf{x}, t) = \frac{\sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) \mathbf{x}_0^j}{\sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2)}. \quad (43)$$

Proof. We assume that $p_{\text{data}}(\mathbf{x}_0)$ is an experimental distribution, consisting of N samples, $\mathbf{X} = \{\mathbf{x}_0^1, \dots, \mathbf{x}_0^N\}$. In other words

$$p_{\text{data}}(\mathbf{x}_0) = \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{x}_0 - \mathbf{x}_0^j). \quad (44)$$

Then,

$$\begin{aligned} \mathcal{L}(\mathbf{x}_\theta) &= \mathbb{E}_{t \sim \mathcal{U}(0, 1)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)} \int_{\mathbb{R}^n} \mathcal{N}(\mathbf{x} | \mathbf{x}_0, s_t^2) \|\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}, t)\|_2^2 d\mathbf{x} \\ &= \frac{1}{N} \sum_{j=1}^N \int_0^1 \mathcal{U}(t | 0, 1) dt \int_{\mathbb{R}^n} \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) \|\mathbf{x}_0^j - \mathbf{x}_\theta(\mathbf{x}, t)\|_2^2 d\mathbf{x} \\ &= \int_0^1 dt \int_{\mathbb{R}^n} \frac{1}{N} \sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) \|\mathbf{x}_0^j - \mathbf{x}_\theta(\mathbf{x}, t)\|_2^2 d\mathbf{x} \\ &= \int_0^1 dt \int_{\mathbb{R}^n} \mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t) d\mathbf{x}, \end{aligned} \quad (45)$$

$\mathcal{L}(\theta)$ can be optimized independently for given \mathbf{x} and t . That is, we have to find the following minimum:

$$\mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t) \rightarrow \min_{\mathbf{x}_\theta}. \quad (46)$$

Taking into account the convexity of the problem

$$\begin{aligned} \nabla_{\mathbf{x}_\theta} \mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t) &= 0, \\ \nabla_{\mathbf{x}_\theta} \mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t) &= \nabla_{\mathbf{x}_\theta} \left(\frac{1}{N} \sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) \|\mathbf{x}_0^j - \mathbf{x}_\theta(\mathbf{x}, t)\|_2^2 \right) \\ &= \frac{1}{N} \sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) \nabla_{\mathbf{x}_\theta} \|\mathbf{x}_0^j - \mathbf{x}_\theta(\mathbf{x}, t)\|_2^2 \\ &= \frac{-2}{N} \sum_{j=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_0^j, s_t^2) [\mathbf{x}_0^j - \mathbf{x}_\theta(\mathbf{x}, t)]. \end{aligned} \quad (47)$$

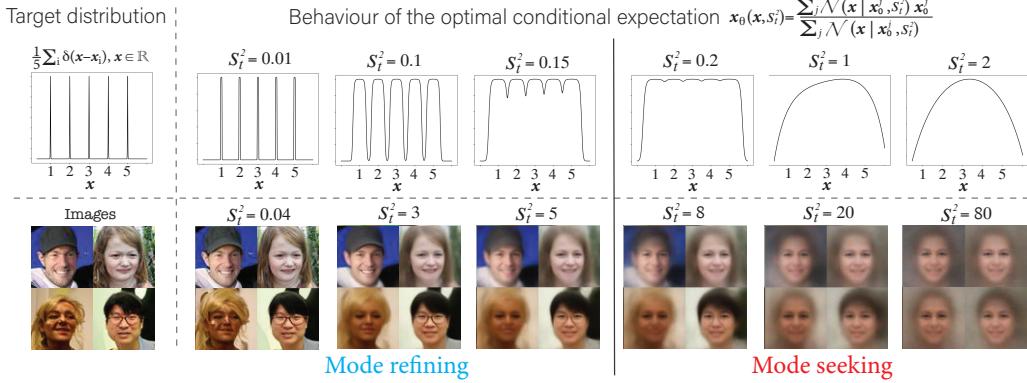


Figure 1: The behaviour of the optimal conditional expectation, $\mathbb{E}[\mathbf{x}_0|\mathbf{x}]$, for the toy distribution (top row) and images (bottom row). The behaviour leads to a rough border which separates two parts – **mode refining** and **mode seeking**. We argue that the mode seeking part **inevitably** requires significant sampler steps number, which blocks them from efficient sampling (1-10 steps).

$$\begin{aligned} \frac{-2}{N} \sum_{j=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_t^2) [\mathbf{x}_0^j - \mathbf{x}_\theta(\mathbf{x}, t)] &= 0, \\ \mathbf{x}_\theta(\mathbf{x}, t) &= \frac{\sum_{j=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_t^2) \mathbf{x}_0^j}{\sum_{j=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_t^2)}. \end{aligned} \quad (48)$$

□

This result gives us the form of the optimal score function for the empirical distribution, $p_{\text{data}}(\cdot)$.

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\frac{1}{s_t^2} \left(\mathbf{x} - \frac{\sum_{j=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_t^2) \mathbf{x}_0^j}{\sum_{j=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_t^2)} \right). \quad (49)$$

To better understand its properties we can visualize the corresponding conditional expectation. Figure 1 demonstrates the behaviour of $\mathbb{E}[\mathbf{x}_0|\mathbf{x}]$ with respect to the different s_t (we denote it as σ for simplicity). In this example we consider two distributions:

- Toy, in which the target variable is one dimensional and the target distribution is the sum of five delta functions.
- Images, where the target variable is a sample from the FFHQ-64 dataset.

As it can be seen, there is a rough border which separates two parts. We call them as the mode refining and seeking. The reason is that the latter demonstrates common features of a entire distribution while the former deals with particular modes of this distribution.

We state that ODE/SDE samplers of the backward equation have to utilize significant number of steps to get through the mode seeking part. This fact blocks them from achieving high quality on small steps (1-10). In other words, samplers firstly need to "grop" objects in the mixture and then tune its fine-grained details. We believe that for efficient sampling the mode seeking part should be completed in a few steps (1-4). However, the optimal solution of the least square regression problem Eq. (26) does not allow us to make it possible.

D.2 PROPERTIES OF THE DISTILLED DIFFUSION MODELS

Similar to diffusion models we can obtain the optimal form of distilled models.

Proposition 4. *The optimal solution to the consistency distillation problem Eq. (35) has the following form:*

$$\mathbf{x}_\theta(\mathbf{x}, t) = \mathbf{x}_0^k. \quad (50)$$

Proof. Assuming

$$p_{\text{data}}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j). \quad (51)$$

Then,

$$\begin{aligned} \mathcal{L}^{\text{CD}}(\mathbf{x}_\theta) &= \mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_{t_n}(\mathbf{x}|\mathbf{x}_0)q_{t_{n-1}}^\phi(\hat{\mathbf{x}}|\mathbf{x})} \|\mathbf{x}_\theta(\hat{\mathbf{x}}, t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 \\ &= \mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_{t_n}(\mathbf{x}|\mathbf{x}_0)} \int_{\mathbb{R}^n} \delta(\hat{\mathbf{x}} - \mathbf{x}_{t_{n-1}}(\mathbf{x})) \|\mathbf{x}_\theta(\hat{\mathbf{x}}, t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 d\hat{\mathbf{x}}, \end{aligned} \quad (52)$$

where we use $q_{t_{n-1}}^\phi(\hat{\mathbf{x}}|\mathbf{x}) = \delta(\hat{\mathbf{x}} - \mathbf{x}_{t_{n-1}}(\mathbf{x}))$ and note $\mathbf{x}_{t_{n-1}}(\mathbf{x})$ as one step of ODE solver with parameters ϕ , from t_n to t_{n-1} , starting from \mathbf{x} .

$$\begin{aligned} \mathcal{L}^{\text{CD}}(\mathbf{x}_\theta) &= \mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)p_{t_n}(\mathbf{x}|\mathbf{x}_0)} \|\mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 \\ &= \mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{\text{data}}(\mathbf{x}_0)} \int_{\mathbb{R}^n} \mathcal{N}(\mathbf{x}|\mathbf{x}_0, s_{t_n}^2) \|\mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 d\mathbf{x} \\ &= \frac{1}{N} \sum_{n=2}^N \frac{1}{M} \sum_{j=1}^M \int_{\mathbb{R}^n} \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_{t_n}^2) \|\mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 d\mathbf{x} \\ &= \frac{1}{N} \sum_{n=2}^N \int_{\mathbb{R}^n} d\mathbf{x} \frac{1}{M} \sum_{j=1}^M \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_{t_n}^2) \|\mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 \\ &= \frac{1}{N} \sum_{n=2}^N \int_{\mathbb{R}^n} d\mathbf{x} \mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t_n). \end{aligned} \quad (53)$$

$\mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t_n)$ can be optimized independently for given \mathbf{x} and n .

$$\begin{aligned} \nabla_{\mathbf{x}_\theta} \mathcal{L}(\mathbf{x}_\theta, \mathbf{x}, t_n) &= \nabla_{\mathbf{x}_\theta} \left(\frac{1}{M} \sum_{j=1}^M \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_{t_n}^2) \|\mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1}) - \mathbf{x}_\theta(\mathbf{x}, t_n)\|_2^2 \right) = 0, \\ \mathbf{x}_\theta(\mathbf{x}, t_n) &= \frac{\sum_{j=1}^M \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_{t_n}^2) \mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1})}{\sum_{j=1}^M \mathcal{N}(\mathbf{x}|\mathbf{x}_0^j, s_{t_n}^2)}. \end{aligned} \quad (54)$$

Taking into account that $\mathbf{x} \sim p_{t_n}(\mathbf{x}|\mathbf{x}_0)$, we can assume $\mathbf{x} = \mathbf{x}_0^k + s_{t_n} \mathbf{z}$, $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|0, 1)$. Then,

$$\mathbf{x}_\theta(\mathbf{x}, t_n) = \mathbf{x}_\theta(\mathbf{x}_{t_{n-1}}(\mathbf{x}), t_{n-1}). \quad (55)$$

Let $n = 2$.

$$\mathbf{x}_\theta(\mathbf{x}, t_2) = \mathbf{x}_\theta(\mathbf{x}_{t_1}(\mathbf{x}), t_1) \stackrel{(1)}{=} \mathbf{x}_{t_1}(\mathbf{x}) = \mathbf{x}_{t_1}(\mathbf{x}_0^k + s_{t_2} \mathbf{z}) \stackrel{(2)}{=} \mathbf{x}_0^k, \quad (56)$$

where (1) is due to parameterization of the consistency model and (2) because t_1 is the starting point ($s_{t_1} = 0$). \square

As it can be seen, distilled models predict distribution modes directly without iterative denoising. This leads to the acceleration in the low-step regime. However, it is much easier for neural network to learn iterative denoising since this task is simpler compared to direct prediction. Thus, diffusion models demonstrate better quality of generated distribution.