

Pseudo Numerical Methods for Diffusion Models on Manifold

Pramook Khungurn

June 7, 2023

This note was written as I read the paper “Pseudo Numerical Methods for Diffusion Models on Manifolds” by Liu et al. [LRLZ22].

1 Introduction

- This paper is often abbreviated as “PNMD,” and the algorithm it proposes is called “PLMS.” The algorithm is one of the well-known solvers tailored to diffusion models. It is distributed with AUTOMATIC1111’s Stable Diffusion UI.
- The objective of the paper is to develop a fast algorithm for sampling a diffusion model.

2 Background

2.1 DDIM Sampler

- In this note, we follow the original DDPM formulation Ho et al. [HJA20].
- The forward and backward process is divided into multiple time steps. A time step is denoted by $t \in \{0, 1, \dots, T\}$.
- The noised sample at time step t is denoted by \mathbf{x}_t . So, \mathbf{x}_0 is the denoised sample from the data distribution, and \mathbf{x}_T should be distributed like $\mathcal{N}(\mathbf{0}, I)$.
- The forward process is controlled by the parameters $\{\beta_t\}_{t=1}^T$ and is given by:

$$\mathbf{x}_t \sim \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t I)$$

for $t = 1, 2, \dots, T$. We can then deduce that

$$\mathbf{x}_t \sim \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) I)$$

where $\alpha_t = 1 - \beta_t$, and $\bar{\alpha}_t = \prod_{u=1}^t \alpha_u$.

- According to Ho et al. [HJA20], the backward process can be formulated as:

$$p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t I)$$

where

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \frac{\beta_t \sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t, \quad (1)$$

$$\tilde{\beta}_t = \beta_t \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}. \quad (2)$$

- Song et al. later showed that the equation above is a special case of the following more general equation:

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\left(\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}}\right), \sigma_t^2 I\right)$$

for any $0 \leq \sigma_t < \sqrt{1 - \bar{\alpha}_t}$ [SME20]. When $\sigma_t = 0$, the backward process becomes deterministic once \mathbf{x}_T has been sampled according to $\mathcal{N}(\sqrt{\bar{\alpha}_T}\mathbf{x}_0, (1 - \bar{\alpha}_T)I)$.

- In a diffusion model, we either train a denoising model $\mathbf{x}_\theta(\mathbf{x}_t, t)$ that predicts \mathbf{x}_0 from \mathbf{x}_t , or we train the noise model $\xi_\theta(\mathbf{x}_t, t)$ that predicts the Gaussian noise $\xi \sim \mathcal{N}(\mathbf{0}, I)$ that was used to construct $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + (1 - \bar{\alpha}_t)\xi$. The models are related to each other via the following equation:

$$\xi_\theta(\mathbf{x}_t, t) = \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}}. \quad (3)$$

- The DDIM sampler uses (1) along with the assumption that $\sigma_t = 0$. The update equation is given by:

$$\begin{aligned} \mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_\theta(\mathbf{x}_t, t) + \sqrt{1 - \bar{\alpha}_{t-1}}\left(\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}}\right) \\ &= \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} - \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}}\right)\mathbf{x}_\theta(\mathbf{x}_t, t) \\ &= \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{(1 - \bar{\alpha}_t)\bar{\alpha}_{t-1} - (1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t})}\right)\mathbf{x}_\theta(\mathbf{x}_t, t) \\ &= \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{\bar{\alpha}_{t-1} - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t})}\right)\mathbf{x}_\theta(\mathbf{x}_t, t), \end{aligned} \quad (4)$$

or, equivalently,

$$\begin{aligned} \mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}}\left(\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\xi_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}\right) + \sqrt{1 - \bar{\alpha}_{t-1}}\xi_\theta(\mathbf{x}_t, t) \\ &= \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{\sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t} - \sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}}\right)\xi_\theta(\mathbf{x}_t, t) \\ &= \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t - \bar{\alpha}_{t-1}(1 - \bar{\alpha}_t)}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t})}\right)\xi_\theta(\mathbf{x}_t, t) \\ &= \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{\bar{\alpha}_t - \bar{\alpha}_{t-1}}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t})}\right)\xi_\theta(\mathbf{x}_t, t) \end{aligned} \quad (5)$$

2.2 ODEs for DDIM Sampling

- We start by turning the DDIM update equation (5) into an ODE. Subtracting \mathbf{x}_t from both sides of the equation, we have

$$\begin{aligned} \mathbf{x}_{t-1} - \mathbf{x}_t &= \frac{\sqrt{\bar{\alpha}_{t-1}} - \sqrt{\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}}\mathbf{x}_t + \left(\frac{\bar{\alpha}_t - \bar{\alpha}_{t-1}}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t})}\right)\xi_\theta(\mathbf{x}_t, t) \\ &= \frac{\bar{\alpha}_{t-1} - \bar{\alpha}_t}{\sqrt{\bar{\alpha}_t}(\sqrt{\bar{\alpha}_{t-1}} + \sqrt{\bar{\alpha}_t})}\mathbf{x}_t + \left(\frac{\bar{\alpha}_t - \bar{\alpha}_{t-1}}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t})}\right)\xi_\theta(\mathbf{x}_t, t) \\ &= (\bar{\alpha}_{t-1} - \bar{\alpha}_t)\left(\frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}(\sqrt{\bar{\alpha}_{t-1}} + \sqrt{\bar{\alpha}_t})} - \frac{\xi_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t})}\right). \end{aligned}$$

Now, we replace $t - 1$ with $t - \delta$:

$$\mathbf{x}_{t-\delta} - \mathbf{x}_t = (\bar{\alpha}_{t-\delta} - \bar{\alpha}_t)\left(\frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}(\sqrt{\bar{\alpha}_{t-\delta}} + \sqrt{\bar{\alpha}_t})} - \frac{\xi_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-\delta}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-\delta}}\sqrt{1 - \bar{\alpha}_t})}\right).$$

Dividing both sides by δ and taking the limit as $\delta \rightarrow 0$, we have

$$\begin{aligned} \lim_{\delta \rightarrow 0} \frac{\mathbf{x}_{t-\delta} - \mathbf{x}_t}{\delta} &= \lim_{\delta \rightarrow 0} \frac{\bar{\alpha}_{t-\delta} - \bar{\alpha}_t}{\delta} \left(\frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}(\sqrt{\bar{\alpha}_{t-\delta}} + \sqrt{\bar{\alpha}_t})} - \frac{\boldsymbol{\xi}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-\delta}}\sqrt{\bar{\alpha}_t} + \sqrt{\bar{\alpha}_{t-\delta}}\sqrt{1 - \bar{\alpha}_t})} \right) \\ &\quad - \frac{d\mathbf{x}_t}{dt} = -\bar{\alpha}'_t \left(\frac{\mathbf{x}_t}{2\bar{\alpha}_t} - \frac{\boldsymbol{\xi}_\theta(\mathbf{x}_t, t)}{2\bar{\alpha}_t\sqrt{1 - \bar{\alpha}_t}} \right) \\ \frac{d\mathbf{x}_t}{dt} &= \frac{\bar{\alpha}'_t}{2\bar{\alpha}_t} \left(\mathbf{x}_t - \frac{\boldsymbol{\xi}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \right). \end{aligned} \quad (6)$$

- Using (3), we have an equivalent ODE:

$$\frac{d\mathbf{x}_t}{dt} = \frac{\bar{\alpha}'_t}{2(1 - \bar{\alpha}_t)} \left(\frac{\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} - \mathbf{x}_t \right). \quad (7)$$

- Note that the above equation can also be derived by staring with (4).

$$\begin{aligned} &\mathbf{x}_{t-1} - \mathbf{x}_t \\ &= \frac{\sqrt{1 - \bar{\alpha}_{t-1}} - \sqrt{1 - \bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{x}_t + \left(\frac{\bar{\alpha}_{t-1} - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t})} \right) \mathbf{x}_\theta(\mathbf{x}_t, t) \\ &= \frac{\bar{\alpha}_t - \bar{\alpha}_{t-1}}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_t})} \mathbf{x}_t + \left(\frac{\bar{\alpha}_{t-1} - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t})} \right) \mathbf{x}_\theta(\mathbf{x}_t, t) \\ &= (\bar{\alpha}_{t-1} - \bar{\alpha}_t) \left(-\frac{\mathbf{x}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_t})} + \frac{\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t})} \right) \\ &= (\bar{\alpha}_{t-1} - \bar{\alpha}_t) \left(\frac{\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{\bar{\alpha}_t})} - \frac{\mathbf{x}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-1}} + \sqrt{1 - \bar{\alpha}_t})} \right). \end{aligned}$$

This becomes

$$\mathbf{x}_{t-\delta} - \mathbf{x}_t = (\bar{\alpha}_{t-\delta} - \bar{\alpha}_t) \left(\frac{\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_t}\sqrt{\bar{\alpha}_{t-\delta}} + \sqrt{1 - \bar{\alpha}_{t-\delta}}\sqrt{\bar{\alpha}_t})} - \frac{\mathbf{x}_t}{\sqrt{1 - \bar{\alpha}_t}(\sqrt{1 - \bar{\alpha}_{t-\delta}} + \sqrt{1 - \bar{\alpha}_t})} \right),$$

which ultimately yields

$$\frac{d\mathbf{x}_t}{dt} = \frac{\bar{\alpha}'_t}{2(1 - \bar{\alpha}_t)} \left(\frac{\mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} - \mathbf{x}_t \right).$$

2.3 Simpler ODEs for DDIM Sampling

- Teng gave in his paper a simpler formulation of the ODE for DDIM sampling [WS23].
- Starting from (4), we have that

$$\begin{aligned} \mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\xi}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \boldsymbol{\xi}_\theta(\mathbf{x}_t, t) \\ \frac{\mathbf{x}_{t-1}}{\sqrt{\bar{\alpha}_{t-1}}} &= \frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}} + \left(\frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_{t-1}}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \right) \boldsymbol{\xi}_\theta(\mathbf{x}_t, t) \\ \frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}} - \frac{\mathbf{x}_{t-1}}{\sqrt{\bar{\alpha}_{t-1}}} &= \left(\frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_{t-1}}} \right) \boldsymbol{\xi}_\theta(\mathbf{x}_t, t). \end{aligned}$$

Replacing $t - 1$ with $t - \delta$, we have

$$\frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}} - \frac{\mathbf{x}_{t-\delta}}{\sqrt{\bar{\alpha}_{t-\delta}}} = \left(\frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_{t-\delta}}}{\sqrt{\bar{\alpha}_{t-\delta}}} \right) \boldsymbol{\xi}_\theta(\mathbf{x}_t, t).$$

Defining $\bar{\mathbf{x}}_t := \mathbf{x}_t / \sqrt{\bar{\alpha}_t}$ and $\gamma_t := \sqrt{1 - \bar{\alpha}_t} / \sqrt{\bar{\alpha}_t}$, we have that

$$\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t-\delta} = (\gamma_t - \gamma_{t-\delta}) \boldsymbol{\xi}_\theta(\mathbf{x}_t, t). \quad (8)$$

In other words,

$$\frac{\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t-\delta}}{\gamma_t - \gamma_{t-\delta}} = \boldsymbol{\xi}_\theta(\mathbf{x}_t, t).$$

Taking the limit as $\delta \rightarrow 0$, we have that

$$\frac{d\bar{\mathbf{x}}_t}{d\gamma_t} = \boldsymbol{\xi}_\theta(\mathbf{x}_t, t) \quad (9)$$

- Alternatively, we can start from (4).

$$\begin{aligned} \mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_\theta(\mathbf{x}_t, t) + \sqrt{1 - \bar{\alpha}_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \right) \\ \frac{\mathbf{x}_{t-1}}{\sqrt{1 - \bar{\alpha}_{t-1}}} &= \frac{\mathbf{x}_t}{\sqrt{1 - \bar{\alpha}_t}} + \left(\frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_{t-1}}} - \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} \right) \mathbf{x}_\theta(\mathbf{x}_t, t) \\ \frac{\mathbf{x}_t}{\sqrt{1 - \bar{\alpha}_t}} - \frac{\mathbf{x}_{t-1}}{\sqrt{1 - \bar{\alpha}_{t-1}}} &= \left(\frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} - \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_{t-1}}} \right) \mathbf{x}_\theta(\mathbf{x}_t, t). \end{aligned}$$

Again, replacing $t - 1$ with $t - \delta$, we have

$$\frac{\mathbf{x}_t}{\sqrt{1 - \bar{\alpha}_t}} - \frac{\mathbf{x}_{t-\delta}}{\sqrt{1 - \bar{\alpha}_{t-\delta}}} = \left(\frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} - \frac{\sqrt{\bar{\alpha}_{t-\delta}}}{\sqrt{1 - \bar{\alpha}_{t-\delta}}} \right) \mathbf{x}_\theta(\mathbf{x}_t, t).$$

Define $\tilde{\mathbf{x}}_t := \mathbf{x}_t / \sqrt{1 - \bar{\alpha}_t}$ and $\omega_t = \sqrt{\bar{\alpha}_t} / \sqrt{1 - \bar{\alpha}_t} = \gamma_t^{-1}$. We have that the above equation can be rewritten as

$$\tilde{\mathbf{x}}_t - \tilde{\mathbf{x}}_{t-\delta} = (\omega_t - \omega_{t-\delta}) \mathbf{x}_\theta(\mathbf{x}_t, t). \quad (10)$$

Finally, taking the limit as $\delta \rightarrow 0$, we have

$$\frac{d\tilde{\mathbf{x}}_t}{d\omega_t} = \mathbf{x}_\theta(\mathbf{x}_t, t). \quad (11)$$

2.4 Numerical Methods

- A numerical method in this context is a way to numerically solve the initial value problem:

$$\begin{aligned} \frac{d\mathbf{y}}{ds} &= \mathbf{f}(\mathbf{y}, s) \\ \mathbf{y}(0) &= \mathbf{y}_0 \end{aligned}$$

where \mathbf{f} is an arbitrary function and \mathbf{y}_0 is a given value. By “numerically solving” the initial value problem, we mean to compute $\mathbf{y}(S)$ for any $S > 0$.

- In a typical setting, we would divide the interval $[0, S]$ into K equal subintervals. We let $\Delta s = S/K$. Numerical methods will compute $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \dots, \mathbf{y}_K$ that are supposed to approximate $\mathbf{y}(\Delta s), \mathbf{y}(2\Delta s), \mathbf{y}(3\Delta s), \dots, \mathbf{y}(K\Delta s)$, respectively. A numerical method generally gives how to compute \mathbf{y}_{k+1} from $\mathbf{y}_k, \mathbf{y}_k$, or any other terms that come before it.
- Here are some famous numerical methods. For brevity, let us s_k denote $k\Delta s$.

- Euler method.

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \mathbf{f}(\mathbf{y}_k, s_k).$$

- Heun’s method.

$$\begin{aligned}\mathbf{f}_1 &= \mathbf{f}(\mathbf{y}_k, s_k), \\ \mathbf{f}_2 &= \mathbf{f}(\mathbf{y}_k + \Delta s \mathbf{f}_1, s_{k+1}), \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \Delta s \left(\frac{\mathbf{f}_1 + \mathbf{f}_2}{2} \right).\end{aligned}$$

- Runge–Kutta method. The following is the 4th order version of the method.

$$\begin{aligned}\mathbf{f}_1 &= \mathbf{f}(\mathbf{y}_k, s_k), \\ \mathbf{f}_2 &= \mathbf{f}\left(\mathbf{y}_k + \frac{\Delta s}{2} \mathbf{f}_1, s_{k+1/2}\right), \\ \mathbf{f}_3 &= \mathbf{f}\left(\mathbf{y}_k + \frac{\Delta s}{2} \mathbf{f}_2, s_{k+1/2}\right), \\ \mathbf{f}_4 &= \mathbf{f}(\mathbf{y}_k + \Delta s \mathbf{f}_3, s_{k+1}), \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \Delta s \left(\frac{\mathbf{f}_1 + 2\mathbf{f}_2 + 2\mathbf{f}_3 + \mathbf{f}_4}{6} \right).\end{aligned}$$

- Linear multi-step method. Each of this is a collection of numerical methods. The i th method in the collection is referred to as the “ i th-order” method. For these methods, we define

$$\mathbf{f}_k = \mathbf{f}(\mathbf{s}_k, s_k).$$

Now, \mathbf{y}_{k+1} would be defined in terms of \mathbf{y}_k and \mathbf{f}_k , \mathbf{f}_{k-1} , and so on. One well-known linear multi-step method is the Adams–Bashforth method.

- * The first-order method is the Euler method.

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \mathbf{f}_k.$$

- * The 2nd-order method is given by:

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \left(\frac{3\mathbf{f}_k - \mathbf{f}_{k-1}}{2} \right).$$

- * The 3rd-order method is given by

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \left(\frac{23\mathbf{f}_k - 16\mathbf{f}_{k-1} + 5\mathbf{f}_{k-2}}{12} \right).$$

- * The 4th-order method is given by

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \left(\frac{55\mathbf{f}_k - 59\mathbf{f}_{k-1} + 37\mathbf{f}_{k-2} - 9\mathbf{f}_{k-3}}{24} \right).$$

3 Method

3.1 Problems with Conventional Numerical Methods

- Once we have an ODE like (6) or (7), we can apply numerical methods discussed previously.

- However, there are difficulties in applying numerical methods to (6) and (7).
 - Both ODEs require the term $\bar{\alpha}'_t$. Nevertheless, t is discrete in case of the original DDPM of Ho et al. [HJA20], and so how are we supposed to find the derivative of a discrete function to begin with? We can do something about it by trying to make $\bar{\alpha}_t$ a continuous function, so this is just a nuisance.
 - The RHSs of the ODEs blow up near $t = 0$ and $t = T$.
- On the other hand, the ODEs of Section 2.3, which are (9) and (11), do not have any of the above problems. There are no derivate terms. Moreover, the RHSs do not blow up.
- The problem with these equations is that the derivatives on the LHSs are with respect to γ_t and ω_t , not t . So, the numerical methods cannot be applied directly. However, we will see immediately in the next section how to circumvent this problem.
- Theoretical problems aside, Salimans and Ho observed that using higher-order numerical methods on (6) or (7) can introduce noise when the number of steps are small [SH22].
- The paper suggests that the above noise problem is due to the fact that diffusion models are typically trained to denoise \mathbf{x}_t of the form $\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\xi}$ where $\boldsymbol{\xi}$ is a univariate Gaussian noise. So, there is a “manifold” where the noised samples belong. If the trajectory of a numerical method strays too far away from this manifold, then the results would be off. The thesis is that, if we use the update rule based on (6) or (7), it becomes easy for the trajectory to veer off course.

3.2 Pseudo Numerical Methods

- The paper observes that the numerical methods in Section 2.4 can be divided into two steps.
 - **Gradient step.** This involves evaluating the function \mathbf{f} at several points and then combining the values somehow to compute and estimate of $d\mathbf{y}/ds$.
 - **Transfer step.** This involves taking the gradient from the previous step and update the variable \mathbf{y} .
- Let us take the 2nd-order Adams–Bashforth method

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \left(\frac{3\mathbf{f}_k - \mathbf{f}_{k-1}}{2} \right)$$

as an example. Its gradient stepn is give by

$$\bar{\mathbf{f}} = \frac{3\mathbf{f}_k - \mathbf{f}_{k-1}}{2},$$

and its transfer step is given by

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \Delta s \bar{\mathbf{f}}. \tag{12}$$

Note that the equation of the transfer step above is the same for all numerical methods in Section 2.4.

- There’s nothing wrong with the gradient step, but the transfer step (12) only works when the timesteps are of the form $s_k = k\Delta s$, which is a linear function of k . However, for ODEs like (9) and (9), the timesteps are $\gamma_t = \sqrt{1 - \bar{\alpha}_t}/\sqrt{\bar{\alpha}_t}$ and $\omega_t = \sqrt{\bar{\alpha}_t}/\sqrt{1 - \bar{\alpha}_t}$, which are not linear functions of t .

- To make the transfer step works for (9) and (11), we rewrite

$$\mathbf{y}_{k+1} = \mathbf{y}_k + (s_{k+1} - s_k)\bar{\mathbf{f}}.$$

Note that this equation is very similar to (8) and (10), which we reproduce (in rewritten forms) below:

$$\begin{aligned}\bar{\mathbf{x}}_{t-\delta} &= \bar{\mathbf{x}}_t + (\gamma_{t-\delta} - \gamma_t)\boldsymbol{\xi}_\theta(\mathbf{x}_t, t), \\ \tilde{\mathbf{x}}_{t-\delta} &= \tilde{\mathbf{x}}_t + (\omega_{t-\delta} - \omega_t)\mathbf{x}_\theta(\mathbf{x}_t, t).\end{aligned}$$

So, the transfer steps for (9) and (11) of the form

$$\begin{aligned}\bar{\mathbf{x}}_{t-\delta} &= \bar{\mathbf{x}}_t + (\gamma_{t-\delta} - \gamma_t)\hat{\boldsymbol{\xi}}, \\ \tilde{\mathbf{x}}_{t-\delta} &= \tilde{\mathbf{x}}_t + (\omega_{t-\delta} - \omega_t)\hat{\mathbf{x}}\end{aligned}$$

where $\hat{\boldsymbol{\xi}}$ is some linear combination of values obtained by evaluating $\boldsymbol{\xi}_\theta$ at several points, and $\hat{\mathbf{x}}$ is also some linear combination of values obtained by evaluating \mathbf{x}_θ at several points.

- Before we go ahead, let us take stock. What does the transfer step like

$$\bar{\mathbf{x}}_{t-\delta} = \bar{\mathbf{x}}_t + (\gamma_{t-\delta} - \gamma_t)\hat{\boldsymbol{\xi}}$$

actually mean? Does it mean that we have to compute in terms of $\bar{\mathbf{x}}_t = \mathbf{x}_t/\sqrt{\bar{\alpha}_t}$ of \mathbf{x}_t ? If so, this would not be very nice because this value would blow up near $t = 0$. The good news is that we can still do everything in terms of \mathbf{x}_t because the above equation is just another form of the DDIM update rule (5). In other words, after rearranging, it becomes:

$$\mathbf{x}_{t-\delta} = \sqrt{\bar{\alpha}_{t-\delta}} \left(\frac{\mathbf{x}_t - \hat{\boldsymbol{\xi}}\sqrt{1-\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \right) + \hat{\boldsymbol{\xi}}\sqrt{1-\bar{\alpha}_{t-\delta}}.$$

For convenience, let us define

$$\phi_\xi(\mathbf{x}_t, \hat{\boldsymbol{\xi}}, t, s) := \sqrt{\bar{\alpha}_s} \left(\frac{\mathbf{x}_t - \hat{\boldsymbol{\xi}}\sqrt{1-\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \right) + \hat{\boldsymbol{\xi}}\sqrt{1-\bar{\alpha}_s}.$$

So, we may write the transfer step as:

$$\mathbf{x}_{t-\delta} = \phi_\xi(\mathbf{x}_t, \hat{\boldsymbol{\xi}}, t, t-\delta).$$

- Similarly, the transfer step

$$\tilde{\mathbf{x}}_{t-\delta} = \tilde{\mathbf{x}}_t + (\omega_{t-\delta} - \omega_t)\hat{\mathbf{x}}$$

is equivalent to

$$\mathbf{x}_{t-\delta} = \phi_{\mathbf{x}}(\mathbf{x}_t, \hat{\mathbf{x}}, t, s)$$

where

$$\phi_{\mathbf{x}}(\mathbf{x}_t, \hat{\mathbf{x}}, t, s) := \sqrt{\bar{\alpha}_s}\hat{\mathbf{x}} + \sqrt{1-\bar{\alpha}_s} \left(\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}}{\sqrt{1-\bar{\alpha}_t}} \right).$$

- Because the transfer steps above are not the same as the transfer step of conventional numerical methods, the paper calls methods that use them “pseudo numerical.”
 - However, I think this nomenclature is not really necessary. The formulation in Section 2.3 tells us that these methods are just conventional numerical methods applied to different ODEs.

- Armed with the new transfer steps, we can turn any numerical methods in Section 2.4 pseudo numerical.
 - The pseudo numerical Euler method is just the DDIM sampling method.

$$\begin{aligned}\widehat{\xi} &= \xi_{\theta}(\mathbf{x}_t, t), \\ \mathbf{x}_{t-1} &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}, t, t-1)\end{aligned}$$

- Pseudo Huen's method.

$$\begin{aligned}\widehat{\xi}_1 &= \xi_{\theta}(\mathbf{x}_t, t) \\ \widehat{\mathbf{x}}_1 &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}_1, t, t-1) \\ \widehat{\xi}_2 &= \xi_{\theta}(\widehat{\mathbf{x}}_1, t-1) \\ \widehat{\xi} &= \frac{\widehat{\xi}_1 + \widehat{\xi}_2}{2} \\ \mathbf{x}_{t-1} &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}, t, t-1)\end{aligned}$$

- Pseudo Runge-Kutta method (PRK).

$$\begin{aligned}\widehat{\xi}_1 &= \xi_{\theta}(\mathbf{x}_t, t) \\ \widehat{\mathbf{x}}_1 &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}_1, t, t-0.5) \\ \widehat{\xi}_2 &= \xi_{\theta}(\widehat{\mathbf{x}}_1, t-0.5) \\ \widehat{\mathbf{x}}_2 &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}_2, t, t-0.5) \\ \widehat{\xi}_3 &= \xi_{\theta}(\widehat{\mathbf{x}}_2, t-0.5) \\ \widehat{\mathbf{x}}_3 &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}_3, t, t-1) \\ \widehat{\xi}_4 &= \xi_{\theta}(\widehat{\mathbf{x}}_3, t-1) \\ \widehat{\xi} &= \frac{\widehat{\xi}_1 + 2\widehat{\xi}_2 + 2\widehat{\xi}_3 + \widehat{\xi}_4}{6} \\ \mathbf{x}_{t-1} &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}, t, t-1)\end{aligned}$$

- Pseudo linear multi-step method (PLMS). Let $\widehat{\xi}_t = \xi_{\theta}(\mathbf{x}_t, t)$.
 - * Second order (PLMS2).

$$\begin{aligned}\widehat{\xi} &= \frac{3\widehat{\xi}_t - \widehat{\xi}_{t+1}}{2} \\ \mathbf{x}_{t-1} &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}, t, t-1)\end{aligned}$$

- * Third order (PLMS3).

$$\begin{aligned}\widehat{\xi} &= \frac{23\widehat{\xi}_t - 16\widehat{\xi}_{t+1} + 5\widehat{\xi}_{t+2}}{12} \\ \mathbf{x}_{t-1} &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}, t, t-1)\end{aligned}$$

- * Fourth order (PLMS4).

$$\begin{aligned}\widehat{\xi} &= \frac{55\widehat{\xi}_t - 59\widehat{\xi}_{t+1} + 37\widehat{\xi}_{t+2} - 9\widehat{\xi}_{t+3}}{24} \\ \mathbf{x}_{t-1} &= \phi_{\xi}(\mathbf{x}_t, \widehat{\xi}, t, t-1)\end{aligned}$$

3.3 Putting It All Together

- The methods we covered in the last section are ways to generate \mathbf{x}_{t-1} from \mathbf{x}_t , \mathbf{x}_{t+1} , and so on. However, a complete sampling algorithm must generate \mathbf{x}_0 from scratch.
- The paper’s main contribution is a sampling algorithm where most of the sampling steps is PLMS4. However, PLMS4 can be used to compute \mathbf{x}_{t-1} only when $\hat{\xi}_t$, $\hat{\xi}_{t+1}$, $\hat{\xi}_{t+2}$, and $\hat{\xi}_{t+3}$ have been computed. As a result, we cannot use it to compute \mathbf{x}_{T-1} , \mathbf{x}_{T-2} , and \mathbf{x}_{T-3} . A sampling algorithm thus must specify how to compute these values.
- There are two implementations of the sampling algorithm using PLMS.
 - The one proposed by Liu et al. in the paper. It uses pseudo Runge–Kutta to compute \mathbf{x}_{T-1} , \mathbf{x}_{T-2} , and \mathbf{x}_{T-3} .
 - The one documented in Teng’s paper [WS23]. It uses PLMS methods of lower orders to compute \mathbf{x}_{T-1} , \mathbf{x}_{T-2} , and \mathbf{x}_{T-3} .
- Let’s start with Teng’s implementation because it’s the easier both to understand and to implement.

Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, I)$.

for $t \leftarrow T, T-1, \dots, 1$ **do**

$\hat{\xi}_t \leftarrow \xi_\theta(\mathbf{x}_t, t)$

if $t = T$ **then**

$\hat{\xi} \leftarrow \hat{\xi}_t$

else if $t = T-1$ **then**

$\hat{\xi} \leftarrow \frac{1}{2}(3\hat{\xi}_t - \hat{\xi}_{t+1})$

else if $t = T-2$ **then**

$\hat{\xi} \leftarrow \frac{1}{12}(23\hat{\xi}_t - 16\hat{\xi}_{t+1} + 5\hat{\xi}_{t+2})$

else

$\hat{\xi} \leftarrow \frac{1}{24}(55\hat{\xi}_t - 59\hat{\xi}_{t+1} + 37\hat{\xi}_{t+2} - 9\hat{\xi}_{t+3})$

end if

$\mathbf{x}_{t-1} \leftarrow \phi_\xi(\mathbf{x}_t, \hat{\xi}, t, t-1)$

end for

return \mathbf{x}_0 .

- Liu et al.’s implementation requires PRK as a subroutine.

PRK(\mathbf{x}, t, s)

begin

$\hat{\xi}_1 \leftarrow \xi_\theta(\mathbf{x}, t)$

$\hat{\mathbf{x}}_1 \leftarrow \phi_\xi(\mathbf{x}, \hat{\xi}_1, t, (t+s)/2)$

$\hat{\xi}_2 \leftarrow \xi_\theta(\hat{\mathbf{x}}_1, (t+s)/2)$

$\hat{\mathbf{x}}_2 \leftarrow \phi_\xi(\mathbf{x}, \hat{\xi}_2, t, (t+s)/2)$

$\hat{\xi}_3 \leftarrow \xi_\theta(\hat{\mathbf{x}}_2, (t+s)/2)$

$\hat{\mathbf{x}}_3 \leftarrow \phi_\xi(\mathbf{x}, \hat{\xi}_3, t, s)$

$\hat{\xi}_4 \leftarrow \xi_\theta(\hat{\mathbf{x}}_3, s)$

$\hat{\xi} \leftarrow \frac{1}{6}(\hat{\xi}_1 + 2\hat{\xi}_2 + 2\hat{\xi}_3 + \hat{\xi}_4)$

$\hat{\mathbf{x}} \leftarrow \phi_\xi(\mathbf{x}, \hat{\xi}, t, s)$

```

    return  $\hat{\mathbf{x}}, \hat{\xi}$ 
end

```

- Liu et al.'s implementation is as follows.

```

Sample  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, I)$ .
for  $t \leftarrow T, T-1, \dots, 1$  do
    if  $t \in \{T, T-1, T-2\}$  then
         $\mathbf{x}_{t-1}, \hat{\xi}_t \leftarrow \text{PRK}(\mathbf{x}_t, t, t-1)$ 
    else
         $\hat{\xi}_t \leftarrow \xi_\theta(\mathbf{x}_t, t)$ 
         $\hat{\xi} \leftarrow \frac{1}{24}(55\hat{\xi}_t - 59\hat{\xi}_{t+1} + 37\hat{\xi}_{t+2} - 9\hat{\xi}_{t+3})$ 
         $\mathbf{x}_{t-1} \leftarrow \phi_\xi(\mathbf{x}_t, \hat{\xi}, t, t-1)$ 
    end if
end for
return  $\mathbf{x}_0$ .

```

- Liu et al.'s implementation might not be as versatile as one might hope. The thing is that it requires evaluating ξ_θ at $t = T-1/2, T-3/2$, and $T-5/2$. The results might not be accurate on models that are trained on discrete times.

References

- [HJA20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *CoRR*, abs/2006.11239, 2020.
- [LRLZ22] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds, 2022.
- [SH22] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *CoRR*, abs/2202.00512, 2022.
- [SME20] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2020.
- [WS23] Suttisak Wizadwongsa and Supasorn Suwajanakorn. Accelerating guided diffusion sampling with splitting numerical methods, 2023.