

Denoising Diffusion Probabilistic Models

August 10, 2022

Useful papers - [\[HJA20\]](#),

1 Introduction

This work is dedicated to generative modelling with *diffusion process*. As always our goal is to recover distribution of real data $p_{data}(x)$. For these purpose, authors proposed an approach based on two processes: *forward* and *reverse*. In the forward process we destroy our data step by step adding noise. We do a lot of small steps. In physics this calls *diffusion* (i'm not sure). In the reverse process we aim to obtain samples from noise, that is, we are going in the opposite direction. It is something like normalizing flows, but with significant differences: 1) we add noise to data (as we know adding noise in generative modelling is extremely necessary because of manifold hypothesis); 2) transformations can be arbitrary (as oppose to normalizing flows where we have to use invertible transformations). Lets discuss in detail both processes.

2 Forward process

As was said in forward (diffusion) process we add noise to data for several steps. This is pretty easy one, we do not need any neural networks to do it. Lets describe it more formally. Let $q(x_0)$ is data distribution. Diffusion process is a Markov Chain, thus we have:

- Sequence of random variables: x_0, x_1, \dots, x_T , where $x_0 \sim q(x_0), x_T \sim \text{noise}$.
- Proposal distribution: $q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{1-\beta_t}x_{t-1}, \beta_t I)$. This distribution is necessary for transformation from variable x_{t-1} to x_t . So, this is core of Markov Chain. Here β_t could be constant (that is, not to depend on time), or gradually increase from small one to 1. Proposal distribution is really important because it gives us the form of all x , i.e.:

$$x_{t+1} = \sqrt{1-\beta_{t+1}}x_t + \sqrt{\beta_{t+1}}\epsilon, \epsilon \sim \mathcal{N}(\epsilon|0, I), x_0 \sim q(x_0) \quad (1)$$

- Joint distribution: $q(x_0, \dots, x_T) = q(x_0) \prod_{t=1}^T q(x_t|x_{t-1})$. This is a markovity property, i.e., probability of next random variable depends only from previous random variable.

This three parts defines our diffusion process.

Lets see on some cool and necessary stuff. In future we will need $q(x_t|x_0)$. As we showed, $q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{1-\beta_t}x_{t-1}, \beta_t I)$. So, it is logical to assume that $q(x_t|x_0)$ is also normal. Thus, we need to obtain its parameters. Lets do some easy tricks.

$$\begin{aligned} x_{t+1} &= \sqrt{1-\beta_{t+1}}x_t + \sqrt{\beta_{t+1}}\epsilon \\ &= \sqrt{1-\beta_{t+1}} \left(\sqrt{1-\beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon \right) + \sqrt{\beta_{t+1}}\epsilon \\ &= \sqrt{1-\beta_{t+1}}\sqrt{1-\beta_t}x_{t-1} + \sqrt{1-\beta_{t+1}}\sqrt{\beta_t}\epsilon + \sqrt{\beta_{t+1}}\epsilon \end{aligned} \quad (2)$$

Then,

$$\begin{aligned} \mathbb{E} x_{t+1} &= \sqrt{1-\beta_{t+1}}\sqrt{1-\beta_t}x_{t-1} \\ \mathbb{D} x_{t+1} &= (1-\beta_{t+1})(1-\beta_t) + \beta_{t+1} \\ &= \beta_t + \beta_{t+1}(1-\beta_t) \\ &= -(-\beta_t - \beta_{t+1}(1-\beta_t) + 1) + 1 \\ &= 1 - (1-\beta_t)(1-\beta_{t+1}) \end{aligned} \quad (3)$$

So, we can say that

$$\begin{aligned} q(x_{t+1}|x_{t-1}) &= \mathcal{N}(x_{t+1}|\sqrt{1-\beta_{t+1}}\sqrt{1-\beta_t}x_{t-1}, 1 - (1-\beta_t)(1-\beta_{t+1})) \\ q(x_{t+1}|x_0) &= \mathcal{N}(x_{t+1}|\prod_{s=1}^{t+1} \sqrt{1-\beta_s}x_0, 1 - \prod_{s=1}^{t+1} (1-\beta_s)) \end{aligned} \quad (4)$$

Let $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$, then $q(x_{t+1}|x_0) = \mathcal{N}(x_{t+1}|\sqrt{\bar{\alpha}_{t+1}}x_0, (1 - \bar{\alpha}_{t+1})I)$. It is possible to show that, $\lim_{t \rightarrow \infty} \bar{\alpha}_t = 0$. But then, if T is large:

$$q(x_T|x_0) \approx \mathcal{N}(x_T|0, I) \quad (5)$$

3 Reverse process

In reverse process we again have Markov Chain. But this process is more complicated because we need to recover data from noise. So, to solve such a difficult task our efforts are not enough we need neural networks. So, again we have

- Sequence of random variables: x_T, x_{T-1}, \dots, x_0 , where $x_0 \sim p_\theta(x_0)$, $x_T \sim \text{noise}$. Here we are going back through time.
- Proposal distribution: $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}|\mu_\theta(x_t, t), \sigma_t^2 I)$. Here $\mu_\theta(x_t, t)$ is neural network. In the papers authors fixed variance. Similar to the previous one, we can write:

$$x_{t-1} = \mu_\theta(x_t, t) + \sigma_t \epsilon, \epsilon \sim \mathcal{N}(\epsilon|0, I), x_T \sim \mathcal{N}(x_T|0, I) \quad (6)$$

- Joint distribution: $p_\theta(x_T, \dots, x_0) = p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$.

Knowing of proposal distribution gives us the ability to obtain real looked samples from noise.

4 Training

Now we have to discuss the training phase of the diffusion model. As always we want to maximize log likelihood

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}_{q(x_0)} \log p_{\theta}(x_0) \quad (7)$$

We can do the same trick like we did in VAE with ELBO:

$$\mathbb{E}_{p_{data}(x)} \log p_{\theta}(x) \geq \mathbb{E}_{p_{data}(x)} \int q(z|x) \log \frac{p_{\theta}(x, z)}{q(z|x)} dz \quad (8)$$

We can say that x_0 is observed variable and x_1, \dots, x_T is latent variables, i.e. $z = (x_1, \dots, x_T)$. Then ELBO for our case:

$$\begin{aligned} & \mathbb{E}_{q(x_0)} \log p_{\theta}(x_0) \\ & \geq \mathbb{E}_{q(x_0)} \int q(x_1, \dots, x_T | x_0) \log \frac{p_{\theta}(x_0, \dots, x_T)}{q(x_1, \dots, x_T | x_0)} dx_1 \dots dx_T \\ & = \int q(x_0) dx_0 \int q(x_1, \dots, x_T | x_0) \log \frac{p_{\theta}(x_0, \dots, x_T)}{q(x_1, \dots, x_T | x_0)} dx_1 \dots dx_T \\ & = \mathbb{E}_q(x_0, \dots, x_T) \log \frac{p_{\theta}(x_0, \dots, x_T)}{q(x_1, \dots, x_T | x_0)} \\ & = \mathbb{E}_q(x_0, \dots, x_T) \log \frac{p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1} | x_t)}{\prod_{t=1}^T q(x_t | x_{t-1})} \\ & = \mathbb{E}_q(x_0, \dots, x_T) \left[\log p_{\theta}(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_t | x_{t-1})} + \log \frac{p_{\theta}(x_0 | x_1)}{q(x_1 | x_0)} \right] \equiv \end{aligned} \quad (9)$$

Here we can see that we have something like KL-divergence, but times are reversed. We can use Bayes rule, but we will meet following problem:

$$q(x_t | x_{t-1}) = \frac{q(x_{t-1} | x_t) q(x_t)}{q(x_{t-1})} \quad (10)$$

We do not know $q(x_t)$ and $q(x_{t-1})$. To avoid this lets use Markov property:

$$\begin{aligned} & \equiv \mathbb{E}_q(x_0, \dots, x_T) \left[\log p_{\theta}(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_t | x_{t-1}, x_0)} + \log \frac{p_{\theta}(x_0 | x_1)}{q(x_1 | x_0)} \right] \\ & = \mathbb{E}_q(x_0, \dots, x_T) \left[\log p_{\theta}(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} + \log \frac{p_{\theta}(x_0 | x_1)}{q(x_1 | x_0)} \right] \\ & = \mathbb{E}_q(x_0, \dots, x_T) \left[\log p_{\theta}(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} + \log \frac{p_{\theta}(x_0 | x_1)}{q(x_1 | x_0)} \right] \\ & = \mathbb{E}_q(x_0, \dots, x_T) \left[\log p_{\theta}(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} + \log \frac{q(x_1 | x_0)}{q(x_T | x_0)} + \log \frac{p_{\theta}(x_0 | x_1)}{q(x_1 | x_0)} \right] \quad (11) \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_q(x_0, \dots, x_T) \left[\log \frac{p_\theta(x_T)}{q(x_T|x_0)} + \sum_{t=2}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} + \log p_\theta(x_0|x_1) \right] \\
&= -\mathbb{E}_{q(x_0, x_T)} \log \frac{q(x_T|x_0)}{p_\theta(x_T)} - \sum_{t=2}^T \mathbb{E}_{q(x_0, x_{t-1}, x_t)} \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \mathbb{E}_{q(x_0, x_1)} \log p_\theta(x_0|x_1) \equiv
\end{aligned} \tag{12}$$

Because of $q(x_0, x_{t-1}, x_t) = q(x_{t-1}|x_0, x_t)q(x_0, x_t)$

$$\begin{aligned}
&\equiv -\mathbb{E}_{q(x_0)} \mathbb{KL}(q(x_T|x_0)||p_\theta(x_T)) - \sum_{t=2}^T \mathbb{E}_{q(x_0, x_t)} \mathbb{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \mathbb{E}_{q(x_0, x_1)} \log p_\theta(x_0|x_1) \\
&= \mathbb{E}_q(x_0, \dots, x_T) \left[-\mathbb{KL}(q(x_T|x_0)||p_\theta(x_T)) - \sum_{t=2}^T \mathbb{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) + \log p_\theta(x_0|x_1) \right]
\end{aligned} \tag{13}$$

Looks cool. It is necessary to highlight one important thing. The first term is the same that we had in VAE (it was regularization term). But in VAE we had serious problem that posterior distribution could be different from prior. It is called **prior hole problem**. Here we do not have this problem because $\mathbb{KL}(q(x_T|x_0)||p_\theta(x_T)) = 0$. Thus, any sample from prior corresponds to real sample from data distribution.

Lets calculate the second KL-divergence.

$$\begin{aligned}
q(x_{t-1}|x_t, x_0) &= \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} \\
&= \frac{\mathcal{N}(x_t|\sqrt{1-\beta_t}x_{t-1}, \beta_t I) \mathcal{N}(x_{t-1}|\sqrt{\bar{\alpha}_{t-1}}x_0, (1-\bar{\alpha}_{t-1})I)}{\mathcal{N}(x_t|\sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I)} \\
&= \mathcal{N}(x_{t-1}|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) \\
\tilde{\mu}(x_t, x_0) &= \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}x_0 + \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}x_t, \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t
\end{aligned} \tag{14}$$

Then, KL-divergence is

$$\begin{aligned}
&\mathbb{KL}(\mathcal{N}(x_{t-1}|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)||\mathcal{N}(x_{t-1}|\mu_\theta(x_t, t), \sigma_t^2 I)) \\
&= \int \mathcal{N}(x_{t-1}|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) \log \frac{\mathcal{N}(x_{t-1}|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)}{\mathcal{N}(x_{t-1}|\mu_\theta(x_t, t), \sigma_t^2 I)} dx_{t-1} \\
&= \int \mathcal{N}(x_{t-1}|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) \log \frac{e^{\|x_{t-1}-\tilde{\mu}(x_t, x_0)\|^2}}{e^{\|x_{t-1}-\mu_\theta(x_t, t)\|^2}} + C \\
&= \int \mathcal{N}(x_{t-1}|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) \|\mu_\theta(x_t, t) - \tilde{\mu}(x_t, x_0)\|^2 + C \\
&= \frac{1}{2\sigma_t^2} \|\mu_\theta(x_t, t) - \tilde{\mu}(x_t, x_0)\|^2 + C
\end{aligned} \tag{15}$$

So, our first final objective for train diffusion model (we will minimize this objective)

$$\boxed{1. \mathcal{L}_{DDPM}(\theta) = \sum_{t=2}^T \mathbb{E}_{q(x_0)q(x_t|x_0)} \left[\frac{1}{2\sigma_t^2} \|\mu_\theta(x_t, t) - \tilde{\mu}(x_t, x_0)\|^2 \right]} \tag{16}$$

What can we say here? We train diffusion model in such a way that reverse proposal distribution $p_\theta(x_{t-1}|x_t)$ equals to $q(x_{t-1}|x_t, x_0)$. However, it is not clear how we can interpret the latter distribution. Because in forward process we define $q(x_{t-1}|x_t)$. I think that it is possible to say that $q(x_{t-1}|x_t, x_0)$ **is a true reverse proposal for samples from our dataset**. More precisely: we take some sample from data, x_0 , apply forward process for him, i.e. $q(x_t|x_{t-1})$, and most important - **we build reverse process for that sample using Bayes rule**. Finally, we want to approximate this true reverse process for all samples x_0 using our model proposal distribution $p_\theta(x_{t-1}|x_t)$. A fair question: why we need some approximation if we have true reverse proposal distribution? Answer: we have this proposal only for samples from dataset and on inference stage we will not have condition on x_0 . Informally, it is possible to say that in training phase we show how to transform from noise to data using samples from dataset. And on the inference, the model does it on its own without any x_0 .

Okay, lets go more deeper. We know that

$$\tilde{\mu}(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t \quad (17)$$

We know that our aim is samples x_0 . So what if we will directly learn our model for that aim. Then we can choose following parameterization:

$$\mu_\theta(x_t, t) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_\theta(x_t, t) + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t \quad (18)$$

Then, our second final objective

$$\boxed{2. \mathcal{L}_{DDPM}(\theta) = \sum_{t=2}^T \mathbb{E}_{q(x_0)q(x_t|x_0)} \frac{\bar{\alpha}_{t-1}\beta_t^2}{2\sigma_t^2(1 - \bar{\alpha}_t)^2} \|x_\theta(x_t, t) - x_0\|^2} \quad (19)$$

We can give following interpretation of this objective: we learn neural network to denoise. That is, we give noised sample x_t and want to reproduce real sample x_0 for which the noisy was received. We can derive more useful property, as we remember

$$\begin{aligned} x_t &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \epsilon \sim \mathcal{N}(\epsilon|0, I) \\ x_0 &= \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon) \end{aligned} \quad (20)$$

So, what if we say that we want to reproduce noise for which the real one was destroyed. Then,

$$x_\theta(x_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(x_t, t)) \quad (21)$$

Our, third final objective is

$$\boxed{3. \mathcal{L}_{DDPM}(\theta) = \sum_{t=2}^T \mathbb{E}_{p(\epsilon)q(x_t|x_0)} \frac{\beta_t^2}{2\sigma_t^2\alpha_t(1 - \bar{\alpha}_t)} \|\epsilon_\theta(x_t, t) - \epsilon\|^2} \quad (22)$$

Authors said that this objective is more preferable.

So, we got our training procedure. However, let's look on another cool stuff. We can remember objective function for denoising score matching:

$$\mathcal{L}_{dsm} = \sum_{i=1}^N \mathbb{E}_{q_{\sigma}(\hat{x}|x)p_{data}(x)} \frac{\sigma_i^2}{2L} [||\mathbf{nn}_{\theta}(\hat{x}, \sigma_i^2) - \nabla_{\hat{x}} \log q_{\sigma_i}(\hat{x}|x)||^2] \quad (23)$$

Looks pretty similar. Both has sum, similar expectations. We took q as $q_{\sigma}(\hat{x}|x) = \frac{1}{(2\pi)^{\frac{d}{2}}\sigma} \exp\left(-\frac{1}{2}\frac{||\hat{x}-x||^2}{\sigma^2}\right)$, then $\nabla_{\hat{x}} \log q_{\sigma}(\hat{x}|x) = \frac{1}{\sigma^2}(x - \hat{x})$

$$\mathcal{L}_{dsm} = \sum_{i=1}^N \mathbb{E}_{q_{\sigma}(\hat{x}|x)p(\epsilon)} \frac{1}{2L\sigma_i^2} [||\mathbf{nn}_{\theta}(\hat{x}, \sigma_i^2) - \epsilon||^2] \quad (24)$$

As it can be seen, these objective are same (without some constants). That's why diffusion models are called score based models. Here we have to make important remark. When score based models is being derived, this noise adding with different magnitude **was motivated**. I mean that we solved problems with Langevin dynamics using different levels of noise. And as a result, diffusion models were obtained. BUT in diffusion models we are working with noise without any motivation. We just use noise because the diffusion process requires it. I think the combination of both articles (about score-based models and diffusion models) should provide an explanation of what's going on (but I still don't have a full understanding :)).

5 Summary

Lets make most important conclusions:

- Diffusion model it is models for generating images based on **forward** and **reverse** processes.
- At training stage we **build forward and reverse processes for samples from dataset using Bayes rule**. Then we train a neural network to study this reverse process. After network is trained it will be able to do reverse process without condition on x_0 . **In inference we build reverse process for arbitrary starting point**.
- We showed that **score based models equivalent to diffusion models**.

A fair question: **why this model works and works good?**

- We have no prior hole problem. This is the main problem of VAE which prevents from being the perfect model
- We have no mode collapse. This is the main problem of GAN.
- We are working with noise. Noise is very important for real data because of the manifold hypothesis.

- Our model is not restricted. We can create any architecture in diffusion models unlike normalizing flows models. So we can build very deep and very flexible models, which is really important for generating high quality samples.

TRAINING

$$\begin{aligned}\theta^{k+1} &= \theta^k - \nabla_{\theta} [\mathbb{E}_{\mathcal{U}(t)q(x_0)p(\epsilon)} \|\epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) - \epsilon\|^2] \\ \mathcal{U}(t) &= \text{Uniform}(\{1, \dots, T\}), p(\epsilon) = \mathcal{N}(\epsilon|0, I)\end{aligned}\tag{25}$$

INFERENCE

At each step we need to sample $x_t \sim p_{\theta}(x_t|x_{t+1}) = \mathcal{N}(x_t|\mu_{\theta}(x_{t+1}, t+1), \sigma_{t+1}^2 I)$. We remember that

$$\begin{aligned}\mu_{\theta}(x_t, t) &= \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_{\theta}(x_t, t) + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t \\ x_{\theta}(x_t, t) &= \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_{\theta}(x_t, t)) \\ \mu_{\theta}(x_t, t) &= \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_{\theta}(x_t, t)\end{aligned}\tag{26}$$

Then $p_{\theta}(x_t|x_{t+1}) = \mathcal{N}\left(x_t|\frac{1}{\sqrt{\alpha_{t+1}}}(x_{t+1} - \frac{1 - \alpha_{t+1}}{\sqrt{1 - \bar{\alpha}_{t+1}}}\epsilon_{\theta}(x_{t+1}, t)), \sigma_{t+1}^2 I\right)$. And we have following sampling procedure (if $\sigma_t^2 = \beta_t$):

1. $x_T \sim \mathcal{N}(\epsilon|0, I)$
2. $x_{T-1} = \frac{1}{\sqrt{\alpha_T}}\left(x_T - \frac{1 - \alpha_T}{\sqrt{1 - \bar{\alpha}_T}}\epsilon_{\theta}(x_T, T)\right) + \sqrt{1 - \alpha_T}z, z \sim \mathcal{N}(z|0, I)$
- ...
3. $x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_{\theta}(x_t, t)\right) + \sqrt{1 - \alpha_t}z, z \sim \mathcal{N}(z|0, I)$
- ...
4. $x_0 = \frac{1}{\sqrt{\alpha_1}}\left(x_1 - \frac{1 - \alpha_1}{\sqrt{1 - \bar{\alpha}_1}}\epsilon_{\theta}(x_1, 1)\right)$

(Note for me: $\epsilon_{\theta}(x_t, t)$ - noise that affected the x_t). $\alpha_t = 1 - \beta_t, \bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$, $T = 1000, \beta_1 = 10^{-4}, \beta_T = 0.02$.

References

- [HJA20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.