

Diffusion Models in Depth: From a Theoretical and a Practical Perspective

Gregory Igor Sedykh

June 28th 2024



**UNIVERSITÉ
DE GENÈVE**

FACULTY OF SCIENCE
Department of Informatics

Bachelor thesis for the degree of Bachelor of Science in Computer Science
Supervised by Prof. Stéphane Marchand-Maillet

Abstract

This report aims to provide a detailed and comprehensive overview of the current state of diffusion models, a type of generative model that has gained popularity over the past few years.

The foundations of the theory behind diffusion models, particularly Denoising Diffusion Probabilistic Models (DDPMs), is seen in depth along with the most important improvements made to them.

To illustrate the theory and to be able to visualise the performance of diffusion models, a Jupyter Notebook is provided with a practical implementation of a simple DDPM using PyTorch. Overall, the reader will be able to understand the theory, visualise the model, be able to implement it, comprehend the improvements and the reasons behind them and observe the difference in performance between the different models, so as to have an understanding of the current progress made in the field of diffusion models.

Contents

1	Introduction	2
2	Denoising Diffusion Probabilistic Models	3
2.1	Forward Process	3
2.2	Reverse Process	5
2.2.1	Variational Lower Bound	5
2.2.2	Simplified training objective	11
2.2.3	Model architecture	13
3	Score-based formulation	14
3.1	Score matching	15
3.2	Sampling	16
3.3	Noise Conditional Score Networks	17
3.4	Equivalence between DDPM and NCSN	19
4	Improvements upon DDPMs	20
4.1	Improved likelihood	20
4.2	Denoising Diffusion Implicit Models (DDIMs)	22
4.3	Guidance	24
4.3.1	Classifier guidance	24
4.3.2	Classifier-free guidance	25
5	Further development	26
5.1	Cascaded Diffusion Models	26
5.2	Latent Diffusion Models	27
5.3	Diffusion Transformer	28
6	Implementation	29
7	Conclusion	31
A	Convolution of two Gaussian distributions	32

1 Introduction

Diffusion models have gained widespread popularity since 2020, as models such as *DALL-E*, *Stable Diffusion* [1, 16] and *Midjourney* have proven to be capable of generating high-quality images given a text prompt. Furthermore, OpenAI’s recent announcement of *Sora* has shown that diffusion models have now also become highly capable of generating minute long high-definition videos from a text prompt [2].

These models date back to 2015, where the idea of a diffusion model appeared, based on diffusion processes used in thermodynamics [20].

Denoising Diffusion Probabilistic Models (DDPMs) were a development of the original diffusion probabilistic model introduced in 2015 [4].

Subsequently, OpenAI improved upon the original DDPMs which did not have ideal log likelihoods [4] while also using less forward passes and therefore speeding up the sampling process [13].

The most recent progress done by OpenAI has allowed their diffusion models to obtain better metrics and better sample quality than Generative Adversarial Networks (GANs) which were previously considered the state-of-the-art in image generation [3].

The fairly recent apparition of diffusion models means not only that there is still a lot to be discovered about them, but also that progress is being made rapidly.

The theory behind diffusion models was mainly founded when Ho et al. [4] introduced their DDPMs in 2020, but many improvements have been made upon their work since then.

To understand what these models are and how they work, it is crucial to understand how DDPMs were developed, what choices were made when developing them and why these choices were made, as well as why changes were made to the original model and how they were made.

This report aims to provide an overview of the theory behind diffusion models, as well as a practical guide on how to implement a simple diffusion model using PyTorch, in order to compare what the theory shows us and what the practical implementation gives us.

2 Denoising Diffusion Probabilistic Models

A **diffusion model** is a generative model that consists of two Markov chains, one forward and one reverse.

Given an input (e.g. an image), the forward process will destroy the information in the image by gradually adding Gaussian noise at each step of the process [4].

The reverse process’ objective is to “invert” the forward process and learn how to do so, starting from a noisy uninformative image and step-by-step, estimating the structure of the information that was destroyed in the image at each step and regenerating it until the first step is reached, where we should obtain the original image [4].

2.1 Forward Process

The forward process is a Markov process that starts from the original image \mathbf{x}_0 and adds Gaussian noise during T steps which results in a more noisy image \mathbf{x}_t [4].

At each step $t \in [1, T]$, the noise added has variance $\beta_t \in \{\beta_1, \dots, \beta_T\}$.

This process is defined by q which is a probability distribution that takes as input an image \mathbf{x}_t and outputs its likelihood:

$$q : \mathbb{R}^D \rightarrow [0, 1] \quad (2.1)$$

where D is the data dimensionality (e.g. for a 64×64 RGB image, $D = 64 \times 64 \times 3$).

Pixels in the image \mathbf{x}_t are assumed to be mutually independent conditionally to the previous timestep.

This is reflected by the fact that the covariance matrix ($\beta_t \mathbf{I}$) is diagonal.

Formally, we obtain the following Markov process:

$$q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (2.2)$$

where:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (2.3)$$

Equation 2.3 gives us a single step forward: given the previous image \mathbf{x}_{t-1} , \mathbf{x}_t is obtained by sampling a D -dimensional Gaussian distribution with mean $\sqrt{1 - \beta_t} \mathbf{x}_{t-1}$ and variance $\beta_t \mathbf{I}$.

Equation 2.2 gives us the full forward process from the original image \mathbf{x}_0 to the final image \mathbf{x}_T .

The reparametrisation trick says that for a univariate Gaussian distribution where $z \sim p(z | x, \mu) = \mathcal{N}(x, \sigma^2)$, a valid reparametrisation would be $z = x + \sigma \epsilon$ where $\epsilon \sim \mathcal{N}(0, 1)$ is just noise [8].

In other words, sampling z conditionally to x from a Gaussian distribution is equivalent to getting z as x with Gaussian noise of variance σ^2 added to it. In this case, we can see that the image \mathbf{x}_t becomes more random because of the Gaussian noise ϵ_t added to \mathbf{x}_{t-1} at each step, since we are sampling:

$$\mathbf{x}_t = q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (2.4)$$

$$= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathbf{I} \epsilon_t \quad \epsilon_t \sim \mathcal{N}_D(0, \mathbf{I}) \quad (2.5)$$

It is important to note that as $T \rightarrow \infty$ and with a correct choice of β_t , \mathbf{x}_T will become a sample of an isotropic Gaussian distribution ($\mathcal{N}(0, \mathbf{I})$) [13, 20].

This is important for the reverse process, as it will allow us to take a sample $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ and reverse the forward process to obtain the original image \mathbf{x}_0 (however this cannot be done so simply, as seen in section 4) [13].

Ho et al. [4] use the cascade of the reparametrisation trick to be able to sample \mathbf{x}_t at any arbitrary step t of the forward process.

This property is due to the fact that the space of Gaussian distributions is closed under convolution. In other words, the addition of two samples from Gaussian distributions is the same as obtaining a sample from the convolution of the initial distributions (see Appendix A).

Let $\alpha_t = 1 - \beta_t$, then:

$$\begin{aligned}\mathbf{x}_t &\sim \mathcal{N}\left(\sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I}\right) \\ \mathbf{x}_t &= \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\epsilon_{t-1} \quad \epsilon_{t-1} \sim \mathcal{N}_D(0, \mathbf{I}) \\ &= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1}\end{aligned}$$

From this, we can apply it again to \mathbf{x}_{t-1} and obtain:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2}$$

Therefore, we get:

$$\mathbf{x}_t = \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{\alpha_t(1 - \alpha_{t-1})}\epsilon_{t-2} + \sqrt{1 - \alpha_t}\epsilon_{t-1}$$

We can write:

$$\mathbf{x}_t = \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \mathcal{N}(0, \alpha_t(1 - \alpha_{t-1})\mathbf{I}) + \sqrt{1 - \alpha_t}\epsilon_{t-1} \quad (2.6)$$

$$= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \mathcal{N}(0, \alpha_t(1 - \alpha_{t-1})\mathbf{I}) + \mathcal{N}(0, (1 - \alpha_t)\mathbf{I}) \quad (2.7)$$

The convolution of the two Gaussian distributions $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$ gives us a new Gaussian distribution with mean $\mu_1 + \mu_2$ and variance $\sigma_1^2 + \sigma_2^2$ (see Appendix A):

$$\mathbf{x}_t = \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \mathcal{N}(0, (\alpha_t(1 - \alpha_{t-1}) + 1 - \alpha_t)\mathbf{I}) \quad (2.8)$$

Which finally gives us:

$$\mathbf{x}_t = \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\bar{\epsilon}_{t-2}$$

Where $\epsilon_{t-1}, \epsilon_{t-2}, \dots \sim \mathcal{N}(0, \mathbf{I})$ and $\bar{\epsilon}_{t-2}$ is the convolution of the noise distributions ϵ_{t-1} and ϵ_{t-2} .

We can recursively apply this backwards until \mathbf{x}_0 .

Let $\bar{\alpha}_t = \prod_{i=0}^t \alpha_i$:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \quad (2.9)$$

This will give us a way to directly compute a noisy image \mathbf{x}_t at step t of the forward process, given the original image \mathbf{x}_0 (we ignore the noise $\epsilon \sim \mathcal{N}_D(0, \mathbf{I})$):

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (2.10)$$

We now also know that $1 - \bar{\alpha}_t$ is the equivalent variance of the noise added to \mathbf{x}_0 to obtain \mathbf{x}_t during the forward process [13].

Choosing the variances β_t is an important step when developing the forward process.

Ho et al. [4] chose to use a linear schedule starting from $\beta_1 = 10^{-4}$ and increasing linearly until $\beta_T = 0.02$ in order for them to be small compared to the pixel values that were in $[-1, 1]$.

However we will see in section 4 that this is not the best choice, as this destroys the images too quickly closer to the end of the process while adding little to sample quality, which made the model sub-optimal for 64×64 and 32×32 images [13].

Finally, Ho et al. [4] chose $T = 1000$ in order to be able to compare their model with Sohl et al.'s [20] model, which also used $T = 1000$ for most image experiments.

We will again see in section 4 that so many steps can make the sampling slow and that there exist better choices [13].

2.2 Reverse Process

As mentioned previously, with a correct choice of β_t , x_T will become the sample of an isotropic Gaussian distribution over \mathbb{R}^D as $T \rightarrow \infty$ [13, 20].

This should mean that we can take a sample $x_T \sim \mathcal{N}_D(0, I)$ and reverse the forward process to obtain the original image x_0 . However, this means that we need to use the entire dataset in order to figure out $q(x_{t-1}|x_t)$, which in practice cannot be done [13].

Therefore, an estimation $p_\theta(x_{t-1}|x_t)$ of $q(x_{t-1}|x_t)$ is found using a trained neural network (see 2.2.3), with θ the parameters of the neural network [13].

The reverse process is defined as follows [4]:

$$p_\theta(x_0, \dots, x_T) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \quad (2.11)$$

where

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \quad (2.12)$$

and

$$p(x_T) = \mathcal{N}_D(0, I) \quad (\text{isotropic Gaussian distribution})$$

Equation 2.11 gives us the full reverse process, starting from the final image x_T (sampled from noise) to the original image x_0 .

Equation 2.12 gives us the reverse process at step t , where we estimate the noise that was added to the image at step t of the forward process in order to find the image x_{t-1} .

Two parameters must be estimated to find the reverse process at step t (from x_t to x_{t-1}): the mean $\mu_\theta(x_t, t)$ and the variance $\Sigma_\theta(x_t, t)$.

At first, Ho et al. [4] used neural networks to estimate both the mean μ_θ and the variance Σ_θ , but explain that estimating the variance was not necessary as there was little difference between the estimated variance $\sigma_t^2 = \tilde{\beta}_t = \frac{1-\alpha_{t-1}}{1-\alpha_t}$ and the actual values $\sigma_t^2 = \beta_t$. Thus they set $\Sigma_\theta(x_t, t) = \sigma_t^2 I = \beta_t I$.

As for the mean, it is estimated using a neural network trained to optimise a variational lower bound (VLB) of the negative log likelihood $\mathbb{E}[-\log p_\theta(x_0)]$.

2.2.1 Variational Lower Bound

Computing the negative log likelihood $\mathbb{E}[-\log p_\theta(x_0)]$ requires knowing $p_\theta(x_0)$ which means going through T steps, which is computationally expensive.

Thus, Ho et al. [4] use a variational lower bound (VLB or more commonly known as ELBO) to estimate the negative log likelihood [20].

The variational lower bound given by Dederik et al. [8] is defined as:

$$\text{VLB} = \mathbb{E}_{q_\phi(\mathbf{x}|\mathbf{z})} [-\log q_\phi(\mathbf{z}|\mathbf{x}) + \log p_\theta(\mathbf{x}, \mathbf{z})] \quad (2.13)$$

$$= -D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z})) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] \quad (2.14)$$

where \mathbf{X}, \mathbf{Z} are jointly distributed random variables with distribution $p_\theta(\mathbf{x}, \mathbf{z})$, $\mathbf{x} \sim p_\theta$ and q_ϕ is any distribution.

This gives us:

$$\text{VLB} = \log p_\theta(\mathbf{x}_0) - D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)) \quad (2.15)$$

Since we need to minimise the loss with the negative log likelihood, we write:

$$-\text{VLB} = -\log p_\theta(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)) \quad (2.16)$$

We know that the Kullback-Leibler divergence D_{KL} is non-negative, which means that we can write:

$$-\log p_\theta(\mathbf{x}_0) \leq -\log p_\theta(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)) \quad (2.17)$$

$$= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (2.18)$$

$$= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{\frac{p_\theta(\mathbf{x}_0, \mathbf{x}_{1:T})}{p_\theta(\mathbf{x}_0)}} \right] \quad (2.19)$$

$$= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{\frac{p_\theta(\mathbf{x}_{0:T})}{p_\theta(\mathbf{x}_0)}} \right] \quad (2.20)$$

$$= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} + \log p_\theta(\mathbf{x}_0) \right] \quad (2.21)$$

$$= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \quad (2.22)$$

so we obtain that:

$$-\log p_\theta(\mathbf{x}_0) \leq \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \quad (2.23)$$

We can subsequently define L_{VLB} as the loss to be minimised:

$$L_{\text{VLB}} = \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \quad (2.24)$$

Ho et al. [4] reformulate the loss in order to reduce variance:

$$L_{\text{VLB}} = \mathbb{E}_q \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right] \quad (2.25)$$

$$= \mathbb{E}_q \left[\log \frac{\prod_{t \geq 1} q(x_t|x_{t-1})}{p(x_T) \prod_{t \geq 1} p_\theta(x_{t-1}|x_t)} \right] \quad (2.26)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \frac{\prod_{t \geq 1} q(x_t|x_{t-1})}{\prod_{t \geq 1} p_\theta(x_{t-1}|x_t)} \right] \quad (2.27)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \prod_{t \geq 1} q(x_t|x_{t-1}) - \log \prod_{t \geq 1} p_\theta(x_{t-1}|x_t) \right] \quad (2.28)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t \geq 1} \log q(x_t|x_{t-1}) - \sum_{t \geq 1} \log p_\theta(x_{t-1}|x_t) \right] \quad (2.29)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t \geq 1} \log q(x_t|x_{t-1}) - \log p_\theta(x_{t-1}|x_t) \right] \quad (2.30)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t \geq 1} \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)} \right] \quad (2.31)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} + \sum_{t \geq 2} \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)} \right] \quad (2.32)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} + \sum_{t \geq 2} \log \frac{q(x_{t-1}|x_t, x_0) q(x_t|x_0)}{q(x_{t-1}|x_0) p_\theta(x_{t-1}|x_t)} \right] \quad (2.33)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} + \sum_{t \geq 2} \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \sum_{t \geq 2} \log \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} \right] \quad (2.34)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} + \log \frac{q(x_T|x_0)}{q(x_1|x_0)} + \sum_{t \geq 2} \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} \right] \quad (2.35)$$

$$= \mathbb{E}_q \left[-\log(p(x_T)) + \log \frac{q(x_T|x_0)}{p_\theta(x_0|x_1)} + \sum_{t \geq 2} \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} \right] \quad (2.36)$$

$$= \mathbb{E}_q \left[-\log(p_\theta(x_0|x_1)) + \log \frac{q(x_T|x_0)}{p(x_T)} + \sum_{t \geq 2} \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} \right] \quad (2.37)$$

$$= \mathbb{E}_q \left[D_{\text{KL}}(q(x_T|x_0) \parallel p(x_T)) \right] \quad (2.38)$$

$$+ \sum_{t \geq 2} D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t)) - \log(p_\theta(x_0|x_1)) \right]$$

This formulation of the loss can be split up into 3 parts:

$$L_T = D_{KL} (q(x_T|x_0) \parallel p(x_T)) \quad (2.39)$$

$$L_{1:T-1} = \sum_{t \geq 2} D_{KL} (q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t)) \quad (2.40)$$

$$L_0 = -\log(p_\theta(x_0|x_1)) \quad (2.41)$$

$$\text{so that} \quad L_{VLB} = L_T + L_{1:T-1} + L_0 \quad (2.42)$$

It is worth noting that at equation (2.33), the terms of q are conditioned on x_0 ; this is in order to be able to easily compute the forward process posteriors, as they are much easier to find given the original image x_0 .

Ho et al. [4] define it as:

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I) \quad (2.43)$$

Using Bayes' theorem, we find:

$$q(x_{t-1}|x_t, x_0) = q(x_t|x_{t-1}, x_0) \frac{q(x_{t-1}|x_0)}{q(x_t|x_0)} \quad (2.44)$$

Assuming that $q(x_t|x_{t-1}, x_0) = \mathcal{N}(x_t; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I)$, then:

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I) \frac{\mathcal{N}(\sqrt{\tilde{\alpha}_{t-1}} x_0, (1 - \tilde{\alpha}_{t-1}) I)}{\mathcal{N}(\sqrt{\tilde{\alpha}_t} x_0, (1 - \tilde{\alpha}_t) I)} \quad (2.45)$$

$$= \frac{1}{\sqrt{2\pi(1 - \tilde{\alpha}_t)(1 - \tilde{\alpha}_{t-1})(\sqrt{1 - \beta_t})}} e^{-\frac{1}{2} \left(\frac{(x_t - \sqrt{1 - \beta_t} x_{t-1})^2}{\beta_t} + \frac{(x_{t-1} - \sqrt{\tilde{\alpha}_{t-1}} x_0)^2}{1 - \tilde{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\tilde{\alpha}_t} x_0)^2}{1 - \tilde{\alpha}_t} \right)} \quad (2.46)$$

Note: we rewrote $q(x_t|x_{t-1}, x_0)$ as $q(x_t|x_{t-1})$ (2.3) in equation (2.45) since x_0 adds no new information not already present in x_{t-1} .

We can ignore the constant factor in front of (2.46) and thus obtain:

$$\propto \exp \left(-\frac{1}{2} \left(\frac{x_t^2 - 2\sqrt{\alpha_t} x_{t-1} x_t + \alpha_t x_{t-1}^2}{\beta_t} + \frac{x_{t-1}^2 - 2\sqrt{\tilde{\alpha}_{t-1}} x_0 x_{t-1} + \tilde{\alpha}_{t-1} x_0^2}{1 - \tilde{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\tilde{\alpha}_t} x_0)^2}{1 - \tilde{\alpha}_t} \right) \right) \quad (2.47)$$

$$= \exp \left(-\frac{1}{2} \left(\frac{x_t^2}{\beta_t} - \frac{2\sqrt{\alpha_t} x_t x_{t-1}}{\beta_t} + \frac{\alpha_t x_{t-1}^2}{\beta_t} + \frac{x_{t-1}^2}{1 - \tilde{\alpha}_{t-1}} - \frac{2\sqrt{\tilde{\alpha}_{t-1}} x_0 x_{t-1}}{1 - \tilde{\alpha}_{t-1}} + \frac{\tilde{\alpha}_{t-1} x_0^2}{1 - \tilde{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\tilde{\alpha}_t} x_0)^2}{1 - \tilde{\alpha}_t} \right) \right) \quad (2.48)$$

$$= \exp \left(-\frac{1}{2} \left(x_{t-1}^2 \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \tilde{\alpha}_{t-1}} \right) + x_{t-1} \left(\frac{-2\sqrt{\alpha_t} x_t}{\beta_t} - \frac{2\sqrt{\tilde{\alpha}_{t-1}} x_0}{1 - \tilde{\alpha}_{t-1}} \right) + C(x_t, x_0) \right) \right) \quad (2.49)$$

where:

$$C(x_t, x_0) = \frac{x_t^2}{\beta_t} + \frac{\bar{\alpha}_{t-1}x_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{1 - \bar{\alpha}_t} \quad (2.50)$$

Considering $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=0}^t \alpha_i$, we can simplify the expression to:

$$= \exp \left(-\frac{1}{2} \left(x_{t-1}^2 \left(\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} + C(x_t, x_0) \right) \right) \quad (2.51)$$

$$= \exp \left(-\frac{1}{2} \left(x_{t-1}^2 \left(\frac{\alpha_t(1 - \bar{\alpha}_{t-1}) + (1 - \alpha_t)}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} + C(x_t, x_0) \right) \right) \quad (2.52)$$

$$= \exp \left(-\frac{1}{2} \left(x_{t-1}^2 \left(\frac{\alpha_t - \alpha_t \bar{\alpha}_{t-1} + 1 - \alpha_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} + C(x_t, x_0) \right) \right) \quad (2.53)$$

$$= \exp \left(-\frac{1}{2} \left(x_{t-1}^2 \left(\frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} + C(x_t, x_0) \right) \right) \quad (2.54)$$

$$= \exp \left(-\frac{1}{2} \left(\left(\frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left(x_{t-1}^2 - 2 \frac{\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}}}{\frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}} x_{t-1} + \frac{C(x_t, x_0)}{\frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}} \right) \right) \right) \quad (2.55)$$

$$= \exp \left(-\frac{1}{2} \left(\left(\frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left(x_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} + \frac{C(x_t, x_0)}{\frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}} \right) \right) \right) \quad (2.56)$$

As for the final term, since:

$$C(x_t, x_0) = \frac{x_t^2}{\beta_t} + \frac{\bar{\alpha}_{t-1}x_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{1 - \bar{\alpha}_t} \quad (2.57)$$

$$= \frac{x_t^2}{1 - \alpha_t} + \frac{\bar{\alpha}_{t-1}x_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{x_t^2 - 2\sqrt{\alpha_t}x_tx_0 + \alpha_tx_0^2}{1 - \bar{\alpha}_t} \quad (2.58)$$

We obtain:

$$C(x_t, x_0) / \frac{1 - \bar{\alpha}_{t-1}}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \quad (2.59)$$

$$= \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_{t-1}} \left(\frac{x_t^2}{1 - \alpha_t} + \frac{\bar{\alpha}_{t-1}x_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{x_t^2 - 2\sqrt{\alpha_t}x_tx_0 + \alpha_tx_0^2}{1 - \bar{\alpha}_t} \right) \quad (2.60)$$

$$= x_t^2 \left(\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} - \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)^2} \right) \quad (2.61)$$

$$+ x_0^2 \left(\frac{(1 - \alpha_t)(\bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} - \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{(1 - \bar{\alpha}_t)^2} \right)$$

$$+ \frac{2x_tx_0\sqrt{\alpha_t}(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)^2}$$

$$= x_t^2 \left(\frac{(1 - \bar{\alpha}_t)(1 - \bar{\alpha}_{t-1}) - (1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)^2} \right) \quad (2.62)$$

$$+ x_0^2 \left(\frac{(1 - \bar{\alpha}_t)(1 - \alpha_t)\bar{\alpha}_{t-1} - (1 - \alpha_t)(1 - \bar{\alpha}_{t-1})\bar{\alpha}_t}{(1 - \bar{\alpha}_t)^2} \right)$$

$$+ \frac{2x_tx_0\sqrt{\alpha_t}(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)^2}$$

$$= x_t^2 \frac{(1 - \bar{\alpha}_{t-1})^2 \alpha_t}{(1 - \bar{\alpha}_t)^2} + x_0^2 \frac{(1 - \bar{\alpha}_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} + \frac{2x_tx_0\sqrt{\alpha_t}(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)^2} \quad (2.63)$$

$$= \frac{(x_t\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) + x_0\sqrt{\bar{\alpha}_{t-1}}(1 - \bar{\alpha}_t))^2}{(1 - \bar{\alpha}_t)^2} \quad (2.64)$$

$$= \left(\frac{(x_t\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) + x_0\sqrt{\bar{\alpha}_{t-1}}(1 - \bar{\alpha}_t))}{(1 - \bar{\alpha}_t)} \right)^2 \quad (2.65)$$

We can now replace this value into the original expression:

$$q(x_{t-1}|x_t, x_0) \propto \exp \left(-\frac{1}{2} \left(\left(\frac{1}{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_{t-1}}} \right) \right. \right. \quad (2.66)$$

$$\left(x_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} \right. \\ \left. \left. \left. + \left(\frac{(x_t\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) + x_0\sqrt{\bar{\alpha}_{t-1}}(1 - \bar{\alpha}_t))}{(1 - \bar{\alpha}_t)} \right)^2 \right) \right) \right) \quad (2.67)$$

$$= \exp \left(-\frac{1}{2} \left(\frac{1}{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_{t-1}}} \right) \right. \quad (2.68)$$

$$\left(x_{t-1} - \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} \right)^2 \left. \right)$$

which through normalisation gives:

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}\left(x_{t-1}; \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}, \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}I\right) \quad (2.69)$$

$$= \mathcal{N}\left(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I\right) \quad (2.70)$$

We have therefore obtained $\tilde{\mu}(x_t, x_0)$ and $\tilde{\beta}_t$ mentioned in Ho et al.'s [4] paper. The authors also rearrange $\tilde{\mu}(x_t, x_0)$ in order to remove the dependency on x_0 :

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon \quad (2.71)$$

$$\Leftrightarrow x_0 = \frac{x_t - \sqrt{1 - \alpha_t}\epsilon}{\sqrt{\alpha_t}} \quad (2.72)$$

$$\Rightarrow \tilde{\mu}(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\frac{x_t - \sqrt{1 - \alpha_t}\epsilon}{\sqrt{\alpha_t}}}{1 - \bar{\alpha}_t} \quad (2.73)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + (1 - \alpha_t)(x_t - \sqrt{1 - \alpha_t}\epsilon)\frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\alpha_t}}}{1 - \bar{\alpha}_t} \quad (2.74)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t} + \frac{(1 - \alpha_t)(x_t - \sqrt{1 - \alpha_t}\epsilon)\frac{1}{\sqrt{\alpha_t}}}{1 - \bar{\alpha}_t} \quad (2.75)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t + \frac{(1 - \alpha_t)}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}}x_t - \frac{(1 - \alpha_t)(\sqrt{1 - \alpha_t})}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}}\epsilon \quad (2.76)$$

$$= \left(\frac{\sqrt{\alpha_t}\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} + \frac{(1 - \alpha_t)}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}}\right)x_t - \frac{(1 - \alpha_t)(\sqrt{1 - \alpha_t})}{\sqrt{(1 - \bar{\alpha}_t)}\sqrt{(1 - \bar{\alpha}_t)}\sqrt{\alpha_t}}\epsilon \quad (2.77)$$

$$= \left(\frac{\alpha_t(1 - \bar{\alpha}_{t-1}) + (1 - \alpha_t)}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)}\right)x_t - \frac{(1 - \alpha_t)}{\sqrt{(1 - \bar{\alpha}_t)}\sqrt{\alpha_t}}\epsilon \quad (2.78)$$

$$= \left(\frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)}\right)x_t - \frac{1 - \alpha_t}{\sqrt{(1 - \bar{\alpha}_t)}\sqrt{\alpha_t}}\epsilon \quad (2.79)$$

$$= \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1 - \alpha_t}{\sqrt{(1 - \bar{\alpha}_t)}\sqrt{\alpha_t}}\epsilon \quad (2.80)$$

$$= \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{(1 - \bar{\alpha}_t)}}\epsilon\right) \quad (2.81)$$

Another important point is that the term L_T (2.39) can be omitted since q has no learnable parameters and together with the fact that it will almost become a Gaussian distribution, the Kullback-Leibler Divergence $D_{KL}(q(x_T|x_0) \parallel p(x_T))$ between a nearly Gaussian distribution and a Gaussian distribution $p(x_T) = \mathcal{N}(x_T; 0, I)$ will be close to 0 [4].

We thus find ourselves with:

$$L_{VLB} = \mathbb{E}_q \left[\sum_{t \geq 2} D_{KL}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t)) - \log(p_\theta(x_0|x_1)) \right] \quad (2.82)$$

2.2.2 Simplified training objective

Ho et al. [4] simplify the training objective by making several changes in order to obtain an easy quantity for minimising the loss.

First of all, $L_{1:T-1}$ can be written as such:

$$L_{1:T-1} = D_{KL} (q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t)) \quad (2.83)$$

$$= D_{KL} (\mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, t), \sigma_t^2 I) \parallel \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I)) \quad (2.84)$$

The Kullback-Leibler divergence between the two Gaussian distributions of dimension D is given by:

$$D_{KL} (\mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, t), \sigma_t^2 I) \parallel \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I)) \quad (2.85)$$

$$= \frac{1}{2} \left(\text{tr}((\sigma_t^2 I)^{-1}(\sigma_t^2 I)) - D + (\mu_\theta - \tilde{\mu}_t)^\top (\sigma_t^2 I)^{-1} (\mu_\theta - \tilde{\mu}_t) + \log \frac{\det(\sigma_t^2 I)}{\det(\sigma_t^2 I)} \right) \quad (2.86)$$

$$= \frac{1}{2} \left(D - D + (\mu_\theta - \tilde{\mu}_t)^\top (\sigma_t^2 I)^{-1} (\mu_\theta - \tilde{\mu}_t) + \log 1 \right) \quad (2.87)$$

$$= \frac{1}{2} \left((\mu_\theta - \tilde{\mu}_t)^\top (\sigma_t^2 I)^{-1} (\mu_\theta - \tilde{\mu}_t) \right) \quad (2.88)$$

$$= \frac{1}{2\sigma_t^2} \left((\mu_\theta - \tilde{\mu}_t)^\top (\mu_\theta - \tilde{\mu}_t) \right) \quad (2.89)$$

$$= \frac{1}{2\sigma_t^2} \|\mu_\theta - \tilde{\mu}_t\|_2^2 \quad (2.90)$$

However, Ho et al. [4] rewrite this formula so that it depends only on the noise ϵ :

$$\frac{1}{2\sigma_t^2} \|\mu_\theta - \tilde{\mu}_t\|_2^2 = \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon \right) - \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta \right) \right\|_2^2 \quad (2.91)$$

$$= \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} (\epsilon - \epsilon_\theta) \right\|_2^2 \quad (2.92)$$

$$= \frac{1}{2\sigma_t^2} \frac{1}{\alpha_t} \frac{(1-\alpha_t)^2}{1-\bar{\alpha}_t} \|\epsilon - \epsilon_\theta\|_2^2 \quad (2.93)$$

$$= \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1-\bar{\alpha}_t)} \|\epsilon - \epsilon_\theta\|_2^2 \quad (2.94)$$

Furthermore, they found that simply ignoring the factor in front gave better results, giving us:

$$L_{1:T-1} = \|\epsilon - \epsilon_\theta\|_2^2 \quad (2.95)$$

As for the second term L_0 , the authors define $p_\theta(x_0|x_1)$ as:

$$p_\theta(x_0|x_1) = \prod_{i=1}^D \int_{\delta_-(x_0^i)}^{\delta_+(x_0^i)} \mathcal{N}(x; \mu_\theta^i(x_1, 1), \sigma_1^2) dx \quad (2.96)$$

where $\delta_+(x)$ and $\delta_-(x)$ are defined as:

$$\delta_+(x) = \begin{cases} \infty & \text{if } x = 1 \\ x + \frac{1}{255} & \text{if } x < 1 \end{cases} \quad (2.97)$$

$$\delta_-(x) = \begin{cases} -\infty & \text{if } x = -1 \\ x - \frac{1}{255} & \text{if } x > -1 \end{cases} \quad (2.98)$$

with D the data dimensionality and i one coordinate from the image.

This is a decoder for the last step in order to get pixel values between $\{0, 1, \dots, 255\}$ from the values in $[-1, 1]$ used for the neural network [4, 13]

However the authors also ignore the entire term, leaving us with:

$$L_{\text{simple}}(\theta) = \mathbb{E}_{\mathbf{q}} [\|\epsilon - \epsilon_{\theta}\|_2^2] \quad (2.99)$$

2.2.3 Model architecture

Ho et al. [4] as well as the following papers [3, 13, 21] use a U-Net style architecture for the neural network to estimate the noise.

The U-Net architecture is a convolutional neural network (CNN) that was originally designed for image segmentation in biomedical imaging [17]. It consists of an encoder (or contracting path) and a decoder (or expansive path) that are connected by a bottleneck.

The encoder takes in the noisy image as input and repeatedly applies unpadded convolutions followed by max pooling layers in order to reduce the spatial dimensions of the image [17].

The bottleneck contains 3 unpadded convolutions each followed by a rectified linear unit (ReLU) activation function [17].

The decoder then applies up-convolutions followed by unpadded convolutions to increase the spatial dimensions of the image, with at each step a concatenation of the feature maps from the encoder [17].

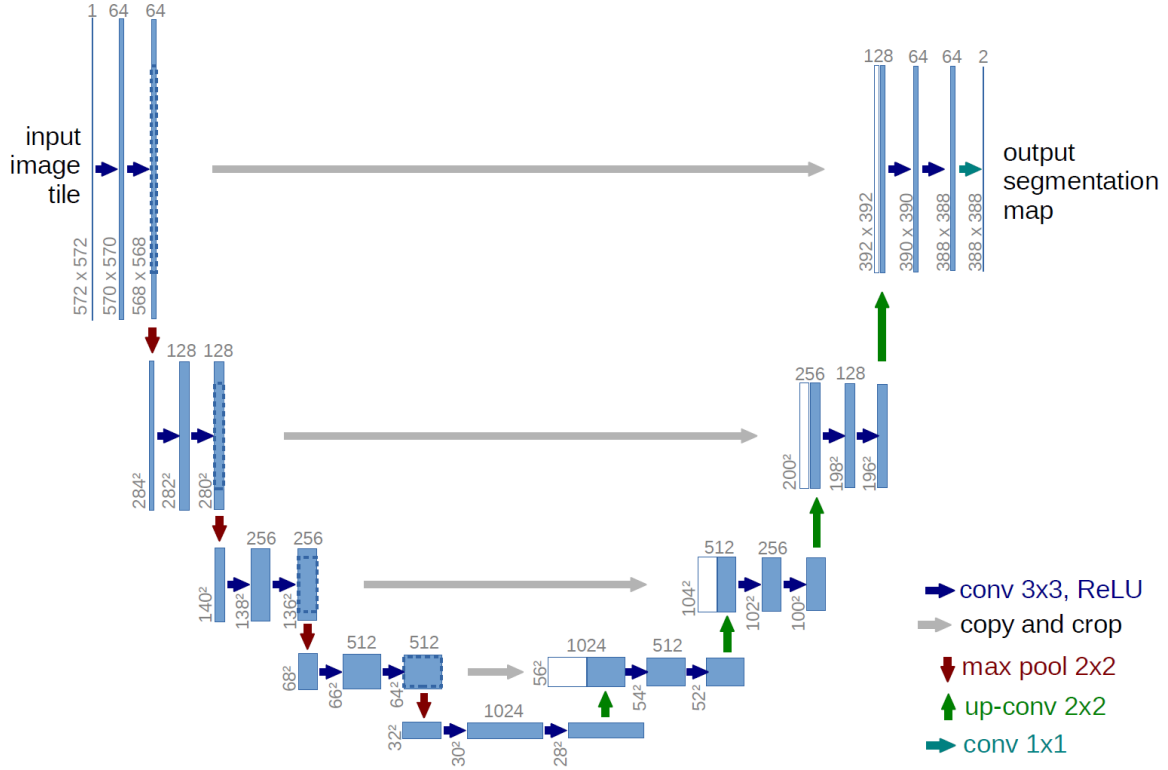


Figure 2.1: Original U-Net architecture [17]

However, Ho et al. [4] made some changes to the original U-Net architecture.

First, they state that they used 4 feature map resolutions for their 32×32 resolution model, and 6 for the 256×256 resolution model. For each of these feature maps, they used two convolutional residual blocks similar of those of a Wide ResNet [30] consisting of a convolutional layer, a ReLU activation function, a convolutional layer, a group normalisation layer, and another ReLU activation function. The “residual” part of the block comes from the fact that the result at the end of the block on the encoder side is concatenated with the result at the end of the block on the decoder side in order to consider features that could have been lost during the downsampling [10].

They used group normalisation [29] rather than weight normalisation [19], which is a normalisation technique that divides the channels into groups and computes the mean and standard deviation for each group [4].

Finally, they state that they used self-attention layers from the Transformer architecture [25] at the 16×16 resolution between the two convolutional residual blocks [4].

Since the model needs to be able to predict the noise at any timestep t , the authors added a time embedding to the input of the network using the sinusoidal positional embedding from the Transformer architecture [25], which allows the parameters to be shared across all timesteps.

The sinusoidal positional embedding is given by [25]:

$$\text{PE}_{(\text{pos}, 2i)} = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \quad (2.100)$$

$$\text{PE}_{(\text{pos}, 2i+1)} = \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \quad (2.101)$$

where pos is the position and i is the dimension

Ho et al. [4] present two algorithms: one for training the model and one for sampling from the model:

Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0)$  ▷ Sample an image  $\mathbf{x}_0$  from the dataset
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ 
5:   Take gradient descent step on:
6:    $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\epsilon, t)\|^2$ 
7: until converged

```

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$  if  $t > 0$ , else  $\mathbf{z} = 0$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

3 Score-based formulation

The DDPM paper [4] allows us to predict either the next denoised image, the mean of the next denoised image or the noise itself.

However they mention another formulation of the DDPM model that is equivalent: a score-based generative modeling version

3.1 Score matching

This formulation uses score matching [7], either denoising score matching [26] or sliced score matching [24].

Essentially, score matching avoids optimizing the likelihood which can sometimes be difficult to find because of normalisation constants, and instead optimises a score function which is the gradient of the log-likelihood.

Suppose we have a probability density function $p_\theta(x)$ that models an i.i.d. dataset $\{x_i\}_{i=1}^N$ of the p.d.f. $p(x)$.

Let $f_\theta(x)$ be a real valued function which is an unnormalized density function and $Z_\theta = \int e^{-f_\theta(x)} dx$ be the normalisation constant such that $\int p_\theta(x) dx = 1$, where $p_\theta(x)$ is defined as:

$$p_\theta(x) = \frac{e^{-f_\theta(x)}}{Z_\theta} \quad (3.1)$$

However for a lot of models, Z_θ cannot be computed [12].

The score function¹ $s_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is defined as [23]:

$$s_\theta(x) = \nabla_x \log p_\theta(x) \quad (3.2)$$

If we try to train a model to optimise this score function, we get:

$$s_\theta(x) = \nabla_x \log p_\theta(x) \quad (3.3)$$

$$= \nabla_x \log \frac{e^{-f_\theta(x)}}{Z_\theta} \quad (3.4)$$

$$= \nabla_x (-f_\theta(x) - \log Z_\theta) \quad (3.5)$$

$$= -\nabla_x f_\theta(x) - \nabla_x \log Z_\theta \quad (3.6)$$

$$= -\nabla_x f_\theta(x) \quad (3.7)$$

Since Z_θ is a constant with respect to x , $\nabla_x \log Z_\theta = 0$ which means we don't need to compute the intractable normalisation constant Z_θ .

This means we could train the score network s_θ to match the true score function $s(x) = \nabla_x \log p(x)$. The model should then optimise the following Fisher divergence $J(\theta)$ [7, 12]:

$$J(\theta) = \frac{1}{2} \mathbb{E}_{x \sim p(x)} [\|s_\theta(x) - s(x)\|_2^2] \quad (3.8)$$

However, since we don't know the true score function $s(x)$ which is the score function of the true data distribution $p(x)$, another formulation was shown by Hyvärinen [7] to be equivalent:

$$J(\theta) = \mathbb{E}_{x \sim p(x)} \left[\text{tr}(\nabla_x s_\theta(x)) + \frac{1}{2} \|s_\theta(x)\|_2^2 \right] \quad (3.9)$$

This is the score matching objective we want to minimise [7].

It can also be computed using the datapoints with the following formula [7]:

$$J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \text{tr}(\nabla_x s_\theta(x_i)) + \frac{1}{2} \|s_\theta(x_i)\|_2^2 \quad (3.10)$$

¹In ML, the score function used is the Stein score, which is the data gradient of the log likelihood. However in statistics, the Fisher score is used and is the parameter gradient of the log likelihood $s_F(x) = \nabla_\theta \log p_\theta(x)$

A problem that arises with this formulation is that the Jacobian matrix $\nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x})$ is not suitable for deep neural networks as it requires backpropagation to compute all its diagonal elements to find the trace [23, 24]. To solve this problem, Song et al. [23] propose either denoising score matching or sliced score matching:

Denoising score matching [26] is a method that slightly noises a datapoint and then score matches the denoised datapoint. This is done by convolving datapoints sampled from the true data distribution \mathbf{p} with a noising kernel \mathbf{q}_{σ} such that $\mathbf{p}(\mathbf{x}) \approx \mathbf{q}_{\sigma}(\mathbf{x})$. Given a datapoint \mathbf{x} , we get a noised datapoint $\tilde{\mathbf{x}}$ by sampling from $\mathbf{q}_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$.

The score matching now estimates the score of the noised data distribution $\mathbf{q}_{\sigma}(\tilde{\mathbf{x}})$, which leads us to the following objective [23, 26]:

$$J(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{q}_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})\mathbf{p}(\mathbf{x})} [\|\mathbf{s}_{\theta}(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log \mathbf{q}_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})\|_2^2] \quad (3.11)$$

Sliced score matching [24] projects the vectors of the vector field $\mathbf{s}_{\theta}(\mathbf{x})$ onto random directions to obtain a scalar field. This will approximate the trace of the Jacobian matrix and it becomes much easier to compute [23].

We get the following objective [24]:

$$J(\theta) = \mathbb{E}_{\mathbf{p}_{\mathbf{v}}} \mathbb{E}_{\mathbf{x} \sim \mathbf{p}(\mathbf{x})} \left[\mathbf{v}^T \nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x}) \mathbf{v} + \frac{1}{2} \|\mathbf{s}_{\theta}(\mathbf{x})\|_2^2 \right] \quad (3.12)$$

The approximation of the Jacobian matrix $\mathbf{v}^T \nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x}) \mathbf{v}$ can be rewritten as [24]:

$$\mathbf{v}^T \nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x}) \mathbf{v} = \mathbf{v}^T \nabla_{\mathbf{x}} (\mathbf{v}^T \mathbf{s}_{\theta}(\mathbf{x})) \quad (3.13)$$

The gradient of the inner product $\nabla_{\mathbf{x}} (\mathbf{v}^T \mathbf{s}_{\theta}(\mathbf{x}))$ requires only one backpropagation pass and requires another inner product with \mathbf{v}^T to approximate the trace [24].

3.2 Sampling

Using the score function mentioned above, we can sample from the model by using Langevin dynamics [23, 27], which will give us the definition of score-based generative modeling [23].

Let $\mathbf{x}_0 \sim \pi(\mathbf{x})$ where $\pi(\mathbf{x})$ is an arbitrary distribution (such as a Gaussian distribution), $\epsilon > 0$ a fixed step size and $\mathbf{z}_t \sim \mathcal{N}_{\mathbf{D}}(0, \mathbf{I})$.

Using Stochastic Gradient Langevin Dynamics (SGLD) [23, 27], we can recursively sample using the following formula:

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\epsilon}{2} \nabla_{\mathbf{x}} \log \mathbf{p}(\mathbf{x}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t \quad (3.14)$$

When $T \rightarrow \infty$ and $\epsilon \rightarrow 0$, the distribution of \mathbf{x}_T will converge to the true data distribution $\mathbf{p}(\mathbf{x})$ therefore the sample will become one from the true data distribution [23, 27].

\mathbf{z}_t is used to add some Gaussian noise at each step because if there was no stochasticity at each step, the sampling would become deterministic and all the samples would essentially be the same.

Since we don't know the true data distribution $\mathbf{p}(\mathbf{x})$, we can use the score function $\mathbf{s}_{\theta}(\mathbf{x})$ trained to match the true score function $\mathbf{s}(\mathbf{x})$ to sample from the model.

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\epsilon}{2} \mathbf{s}_{\theta}(\mathbf{x}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t \quad (3.15)$$

There is however a problem with this method: the score function is well approximated at high density regions where there are a lot of data points but not at low density regions [23].

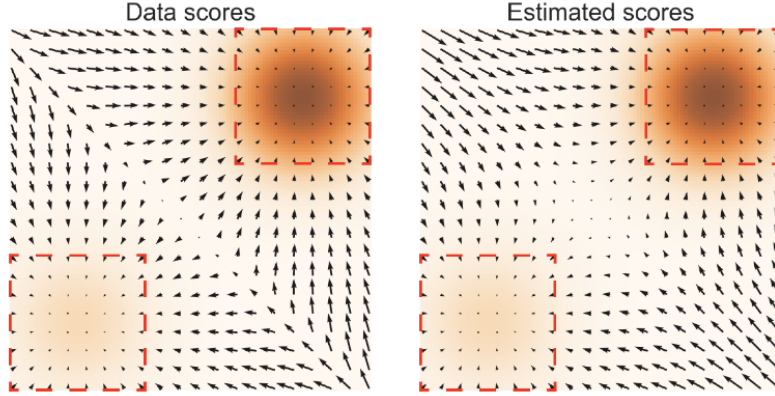


Figure 3.1: Data scores represent the real score function $\nabla_x \log p(x)$ and estimated scores represent the estimated score function $s_\theta(x)$. We can see that the regions in the red rectangles are the darker high density regions where the score function is well approximated, while the lighter lower density regions outside the red rectangles were not well estimated [23].

Song and Ermon [23] show this in figure 3.1 by comparing the real score function $\nabla_x \log p(x)$ and the estimated score function $s_\theta(x)$ for a mixture of two Gaussian distributions.

This means that when sampling using Langevin dynamics, the samples will ultimately be incorrect [23]. To correct this, Song and Ermon [23] use annealed Langevin dynamics.

The idea is to add noise to the data points in order to improve the score function estimation in low density regions [22]. However adding too much noise will make the score function estimation worse in high density regions [22].

Therefore, they control the amount of noise added using a scale of different standard deviations σ_i with $i = 1, 2, \dots, L$ such that $\sigma_1 < \sigma_2 < \dots < \sigma_L$ [22, 23].

To sample, the same Langevin dynamics method is used going from $i = L, L-1, \dots, 1$ with σ_i decreasing at each step. This is the annealed Langevin dynamics method they used, and the results can be seen in figure 3.2, where we clearly see that the annealed Langevin dynamics samples are much more similar to the original data samples than the simple Langevin dynamics samples [22, 23].

3.3 Noise Conditional Score Networks

Noise Conditional Score Networks (NCSN) are a score-based generative model that use a neural network to estimate the score function of the previously mentioned noised distribution and then use annealed Langevin dynamics to sample from the model [23].

Given the increasing sequence of standard deviations σ_i with $i = 1, 2, \dots, L$, the perturbed data distribution is defined as [23]:

$$q_{\sigma_i}(x) = \int p(t) N(x; t, \sigma_i^2) dt \quad (3.16)$$

Song and Ermon [23] then define the Noise Conditional Score Network s_θ that estimates the scores of

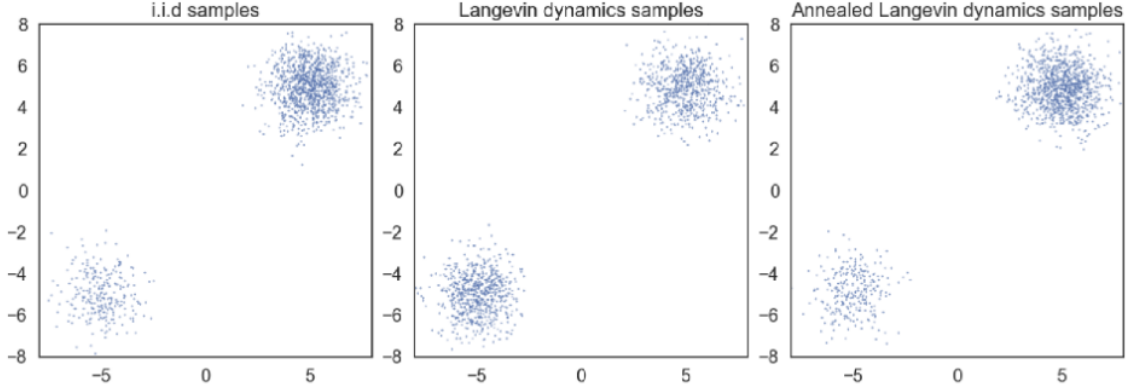


Figure 3.2: Left: Samples from data distribution. Middle: Samples from model using Langevin dynamics. Right: Samples from model using annealed Langevin dynamics [23].

all the perturbed data distributions $q_{\sigma_i}(\mathbf{x})$ given a data point \mathbf{x} :

$$s_{\theta} : \mathbb{R}^D \times \mathbb{R} \rightarrow \mathbb{R}^D \quad (3.17)$$

$$s_{\theta}(\mathbf{x}, \sigma) \approx \nabla_{\mathbf{x}} \log q_{\sigma}(\mathbf{x}) \quad \forall \sigma \in \{\sigma_1, \sigma_2, \dots, \sigma_L\} \quad (3.18)$$

To train the NCSN, the authors chose denoising score matching [26] as its goal is to estimate the score function of a noised data distribution while also being quicker than sliced score matching [23, 24].

Let $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2, \mathbf{I})$ be the noise distribution that noises the data point \mathbf{x} to get the noised data point $\tilde{\mathbf{x}}$.

Its score function is given by [23]:

$$\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) = -\frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \quad (3.19)$$

Let $\lambda(\sigma_i)$ be a positive coefficient function, often chosen to be $\lambda(\sigma_i) = \sigma_i^2$ [22, 23].

The denoising score matching objective is given by [23]:

$$\mathcal{L}(\theta, \{\sigma_i\}_{i=1}^L) = \frac{1}{2L} \sum_{i=1}^L \lambda(\sigma_i) \mathbb{E}_{\mathbf{p}}(\mathbf{x}) \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma_i^2 \mathbf{I})} \left[\left\| s_{\theta}(\tilde{\mathbf{x}}, \sigma_i) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma_i^2} \right\|_2^2 \right] \quad (3.20)$$

To sample from the model, Song and Ermon [23] define an algorithm that uses the previously mentioned annealed Langevin dynamics:

Algorithm 3 Sampling from NCSN

Require: $\{\sigma_i\}_{i=1}^L$, ϵ , T

```
1:  $\tilde{x}_0 \sim \pi(x)$ 
2: for  $i = 1, 2, \dots, L$  do
3:    $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ 
4:   for  $t = 1, 2, \dots, T$  do
5:      $z_t \sim \mathcal{N}(0, I)$ 
6:      $\tilde{x}_t \leftarrow \tilde{x}_{t-1} + \frac{\alpha_i}{2} s_\theta(\tilde{x}_{t-1}, \sigma_i) + \sqrt{\alpha_i} z_t$ 
7:   end for
8:    $\tilde{x}_0 \leftarrow \tilde{x}_T$ 
9: end for
10: return  $\tilde{x}_T$ 
```

3.4 Equivalence between DDPM and NCSN

Ho et al.'s DDPM paper [4] showed that the DDPM model is equivalent to the NCSN model with denoising score matching and annealed Langevin dynamics.

This can be seen by the fact that sampling algorithm of the DDPM paper (2) is in fact very similar to the sampling algorithm of the NCSN paper (3) but rather than estimating a score, the DDPM model estimates the noise [4].

Similarly, the score function $s_\theta(x, \sigma)$ can be rewritten in terms of the DDPM's forward process as [28]:

$$s_\theta(x, \sigma) \approx \nabla_x \log q_\sigma(x) \Leftrightarrow s_\theta(x_t, t) \approx \nabla_{x_t} \log q(x_t) \quad (\text{see 2.3}) \quad (3.21)$$

$$= \mathbb{E}_{q(x_0)} [\nabla_{x_t} \log q(x_t | x_0)] \quad (3.22)$$

$$= \mathbb{E}_{q(x_0)} [\nabla_{x_t} \log \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I)] \quad (3.23)$$

$$= \mathbb{E}_{q(x_0)} \left[\nabla_{x_t} \log e^{-\frac{1}{2} \left(\frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1 - \bar{\alpha}_t}} \right)^T \left(\frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1 - \bar{\alpha}_t}} \right)} \right] \quad (3.24)$$

$$= \mathbb{E}_{q(x_0)} [\Sigma^{-1} (x_t - \sqrt{\bar{\alpha}_t} x_0)] \quad (3.25)$$

$$= \mathbb{E}_{q(x_0)} \left[\Sigma^{-1} \left(x_t - \sqrt{\bar{\alpha}_t} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta}{\sqrt{\bar{\alpha}_t}} \right) \right) \right] \quad (3.26)$$

$$= \mathbb{E}_{q(x_0)} [\Sigma^{-1} (\sqrt{1 - \bar{\alpha}_t} \epsilon_\theta)] \quad (3.27)$$

$$= -\frac{1}{1 - \bar{\alpha}_t} \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta \quad (3.28)$$

$$= -\frac{\epsilon_\theta}{\sqrt{1 - \bar{\alpha}_t}} \quad (3.29)$$

We can see that this quantity is used in sampling algorithm (2) of the DDPM paper [4].

4 Improvements upon DDPMs

As mentioned previously, the DDPM paper [4] had certain choices that limited its performance. Several papers [13, 21] following it have improved upon the original DDPM in various ways.

4.1 Improved likelihood

Ho et al. [4] admit that their model, despite its high sample quality, did not have a competitive log-likelihood in comparison to other likelihood-based models such as VAEs and autoregressive models [13].

The *Improved Denoising Diffusion Probabilistic Models* paper published by Nichol et al. [13] in 2021 implemented certain changes in order to improve the log-likelihood of the model not only on simpler datasets such as CIFAR-10 [9] that the original DDPM was trained on, but also on more complex datasets such as ImageNet [8, 13, 14].

In order to improve the log-likelihood, Nichol et al. [13] started by investigating the fixed variance $\Sigma_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I} = \beta_t \mathbf{I}$ of the original DDPM and whether or not it could be worth learning it.

Recall that Ho et al. [4] saw no difference in sample quality when using a learned variance $\tilde{\beta}_t$ and therefore opted for using the fixed variance β_t .

Nichol et al. [13] came to the same conclusion, finding that as the number of steps T increased (the authors also mentioned that they used $T = 4000$ for their models rather than $T = 1000$ from the original paper), the difference between β_t and $\tilde{\beta}_t$ remained very small. However, they were still interested in the possibility of learning the variance in order to improve the log-likelihood [13].

For this, they used an interpolation approach, where they parametrised the variance as an interpolation between the fixed variance β_t and a learned variance $\tilde{\beta}_t$ in the log domain, with \mathbf{v} an output vector from the model containing one component per dimension:

$$\Sigma_\theta(\mathbf{x}_t, t) = \exp \left(\mathbf{v} \log \beta_t + (1 - \mathbf{v}) \log \tilde{\beta}_t \right) \quad (4.1)$$

As the simplified training objective L_{simple} used in the original DDPM paper [4] did not include the variance, Nichol et al. [13] introduced a new training objective L_{hybrid} , defined as:

$$L_{\text{hybrid}} = L_{\text{simple}} + \lambda L_{\text{VLB}} \quad (4.2)$$

with $\lambda = 0.001$ in order to keep L_{simple} as the main training objective [13].

At first, they found that the hybrid objective L_{hybrid} achieved a better log-likelihood and was easier to optimise than just L_{VLB} , contrary to what they had expected [13].

However, they obtained the best log-likelihoods by optimising L_{VLB} directly but with importance sampling, which is a technique where samples of a difficult distribution are taken from a distribution that is easier to compute and then the original distribution is approximated by weighted samples [13].

They defined it as:

$$L_{\text{VLB}} = \mathbb{E}_{t \sim \mathbf{p}_t} \left[\frac{L_t}{\mathbf{p}_t} \right] \quad \text{where } \mathbf{p}_t \propto \sqrt{\mathbb{E}[L_t^2]} \text{ and } \sum \mathbf{p}_t = 1 \quad (4.3)$$

Another change that Nichol et al. [13] made was to use a cosine noising schedule rather than the linear noising schedule used in the original DDPM paper [4].

They explain this change by the fact that towards the end of the forward process, the image is already very noisy and yet more noise is added while not really improving the model, particularly for 32×32 and 64×64 images [13].

In effect, the input images are destroyed by noise very quickly close to the start of the noising process.

Nichol and Dhariwal [13] found that up to 20% of the reverse process could be skipped while losing very little image quality and that therefore the last noising steps are redundant.

The cosine schedule they use adds noise at a slower rate throughout the process, in particular at the beginning and the end.

They define the cosine schedule as:

$$f(t) = \cos \left(\frac{\frac{t}{T} + s}{1 + s} \cdot \frac{\pi}{2} \right)^2 \quad (4.4)$$

$$\bar{\alpha}_t = \frac{f(t)}{f(0)} \quad (4.5)$$

$$\beta_t = 1 - \frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}} \quad (4.6)$$

T is the total number of timesteps, t is the timestep and s is a small offset used to prevent β_t from being too small when $t = 0$ in order to allow the noise ϵ to be predicted more accurately at the beginning of the process [13].

With these changes, Nichol et al. [13] were able to achieve a negative log-likelihood (NLL) of 2.94 bits/dim on CIFAR-10 [9] and 3.53 bits/dim on ImageNet [14], compared to the original DDPM paper [4] which had a NLL of 3.70 on CIFAR-10 and 3.77 on ImageNet [13].

Model	ImageNet	CIFAR
Glow (Kingma & Dhariwal, 2018)	3.81	3.35
Flow++ (Ho et al., 2019)	3.69	3.08
PixelCNN (van den Oord et al., 2016c)	3.57	3.14
SPN (Menick & Kalchbrenner, 2018)	3.52	-
NVAE (Vahdat & Kautz, 2020)	-	2.91
Very Deep VAE (Child, 2020)	3.52	2.87
PixelSNAIL (Chen et al., 2018)	3.52	2.85
Image Transformer (Parmar et al., 2018)	3.48	2.90
Sparse Transformer (Child et al., 2019)	3.44	2.80
Routing Transformer (Roy et al., 2020)	3.43	-
DDPM (Ho et al., 2020)	3.77	3.70
DDPM (cont flow) (Song et al., 2020b)	-	2.99
Improved DDPM (ours)	3.53	2.94

Figure 4.1: Comparison of NLLs of different models [13]

Moreover, Nichol et al. [13] were able to improve sampling speed due to the aforementioned changes, even though they trained their model with $T = 4000$ timesteps.

They were able to train their model with 4000 timesteps, but use only 100 timesteps for sampling which brought them near-optimal FID (Fréchet Inception Distance) scores when using the L_{hybrid} training objective [13].

They had tried to reduce the number of sampling steps with the original fixed variance versions from the DDPM paper [4] but found that the sample quality degraded a lot more [13].

4.2 Denoising Diffusion Implicit Models (DDIMs)

As mentioned previously, the DDPM paper [4] has a forward process and a reverse process which are both Markov processes.

While the forward process can be easily calculated for a specific timestep t with equation (2.10), the reverse process still requires 1000 timesteps in the case of Ho et al.'s [4] model, and Nichol et al.'s model requires 100 and 4000 timesteps for sampling and training respectively [13].

This is time-consuming and computationally expensive compared to GANs which only require one forward pass through the network rather than thousands to produce a sample [21].

To address this issue, Song et al. [21] introduced the Denoising Diffusion Implicit Models (DDIMs) in 2020.

The main objective of DDIMs was to accelerate sampling by making the forward process non-Markovian, which would allow the reverse process to require less iterations [21].

Song et al. [21] define a family \mathcal{Q} of forward processes. Each forward process in \mathcal{Q} is indexed by $\sigma = \{\sigma_1, \sigma_2, \dots, \sigma_T\} \in \mathbb{R}_{\geq 0}^T$, where σ_t is the standard deviation of the noise added at timestep t :

$$q_\sigma(x_{1:T}|x_0) = q_\sigma(x_T|x_0) \prod_{t=2}^T q_\sigma(x_{t-1}|x_t, x_0) \quad (4.7)$$

where

$$q_\sigma(x_T|x_0) = \mathcal{N}(\sqrt{\alpha_T}x_0, (1 - \alpha_T)I) \quad (4.8)$$

and for all $t > 1$:

$$q_\sigma(x_{t-1}|x_t, x_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}x_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{x_t - \sqrt{\alpha_t}x_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 I\right) \quad (4.9)$$

Song et al. [21] note that when $\sigma_t^2 = 0$, then x_{t-1} is fixed and known in advance.

Therefore, we can see that as Song et al. [21] state, the value of σ controls the stochasticity (or more simply, the randomness) of the forward process.

The authors use Bayes' rule similarly to the original DDPM paper [4] in equation (2.44), but instead to find the forward process $q_\sigma(x_t|x_{t-1}, x_0)$:

$$q_\sigma(x_t|x_{t-1}, x_0) = \frac{q_\sigma(x_{t-1}|x_t, x_0) q_\sigma(x_t|x_0)}{q_\sigma(x_{t-1}|x_0)} \quad (4.10)$$

Since $q_\sigma(x_t|x_{t-1}, x_0)$ means that x_t not only depends on x_{t-1} but also on x_0 , we see that the forward process is no longer Markovian, which was what the authors were aiming to achieve [21].

However, since the reverse process in the original DDPM paper [4] was an estimation of the forward process which was Markovian, a new reverse process was defined.

Their idea of the reverse process is to sample an image x_T like in the original DDPM paper [4], predict x_0 from x_T and then use this prediction of x_0 to sample x_{T-1} [21].

This is repeated until we reach x_1 and then x_0 is predicted from x_1 [21].

Similar to equation (2.72) from the DDPM paper, Song et al. [21] define the prediction of x_0 from x_t as:

$$f_\theta^{(t)}(x_t) = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon(x_t, t)}{\sqrt{\bar{\alpha}_t}} \quad (4.11)$$

Given this, the reverse process is defined as [21]:

$$\mathbf{p}_\theta(\mathbf{x}_T) = \mathcal{N}_D(0, \mathbf{I}) \quad (4.12)$$

$$\mathbf{p}_\theta^{(t)}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \begin{cases} \mathcal{N}(\mathbf{f}_\theta^{(1)}(\mathbf{x}_1), \sigma_1^2 \mathbf{I}) & \text{if } t = 1 \\ \mathbf{q}_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{f}_\theta^{(t)}(\mathbf{x}_t)) & \text{if } t > 1 \end{cases} \quad (4.13)$$

Therefore, \mathbf{x}_{t-1} is sampled from \mathbf{x}_t [21]:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sigma_t \boldsymbol{\epsilon}_t \quad (4.14)$$

In order to train the model, Song et al. [21] define the following *variational inference objective*:

$$J_\sigma(\boldsymbol{\epsilon}_\theta) = \mathbb{E}_{\mathbf{q}_\sigma(\mathbf{x}_{0:T})} [\log \mathbf{q}_\sigma(\mathbf{x}_{1:T}|\mathbf{x}_0) - \log \mathbf{p}_\theta(\mathbf{x}_{0:T})] \quad (4.15)$$

$$= \mathbb{E}_{\mathbf{q}_\sigma(\mathbf{x}_{0:T})} \left[\log \mathbf{q}_\sigma(\mathbf{x}_T|\mathbf{x}_0) + \sum_{t=2}^T \log \mathbf{q}_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) - \sum_{t=1}^T \log \mathbf{p}_\theta^{(t)}(\mathbf{x}_{t-1}|\mathbf{x}_t) - \log \mathbf{p}_\theta(\mathbf{x}_T) \right] \quad (4.16)$$

This would mean that for different choices of σ , we would have a different objective to optimise [21]. However, the authors [21] prove that J_σ is in fact equivalent to L_γ for certain choices of γ , where L_γ is defined as:

$$L_\gamma(\boldsymbol{\epsilon}_\theta) = \sum_{t=1}^T \gamma_t \mathbb{E}_{\mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0), \boldsymbol{\epsilon}_t \sim \mathcal{N}(0, \mathbf{I})} \left[\|\boldsymbol{\epsilon}_\theta^{(t)}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_t) - \boldsymbol{\epsilon}_t\|_2^2 \right] \quad (4.17)$$

where γ is a vector of length T with positive values that depend on α .

We easily see that when $\gamma = 1$, we get L_1 which is the training objective from the original DDPM paper [4, 21].

Song et al. prove that

$$\forall \sigma > 0, \exists \gamma \in \mathbb{R}_{\geq 0}^T \text{ and } \exists C \in \mathbb{R} \text{ s.t. } J_\sigma = L_\gamma + C$$

This means that the optimal solution for J_σ is the same as the optimal solution for L_1 , which means that the training objective for DDIMs can be kept the same as the one for DDPMs if the parameters being trained θ are not the same for all timesteps [21].

In order to speed up the sampling process, Song et al. [21] skip some timesteps in the reverse process. To do this, the forward process is redefined on a subset of timesteps of size S , $\{\mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_S}\}$ and such that $\mathbf{q}(\mathbf{x}_{\tau_i}|\mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_{\tau_i}} \mathbf{x}_0, (1 - \alpha_{\tau_i}) \mathbf{I})$

Since this is a subset of the timesteps, the training can be done for all timesteps and the sampling can be done on just the subset [21].

With the same amount of training timesteps $T = 1000$ and the same training objective as the DDPM paper [4], Song et al. [21] were able to achieve solid FID scores on CIFAR-10 [9] and CelebA [11] datasets with far fewer sampling timesteps, ranging from 10 to 100 timesteps [21].

A new hyperparameter $\eta \in \mathbb{R}_+$ is introduced to control the interpolation between the DDPM σ and the deterministic σ of the DDIM [21]:

$$\sigma_{\tau_i}(\eta) = \eta \sqrt{\frac{1 - \alpha_{\tau_{i-1}}}{1 - \alpha_{\tau_i}}} \sqrt{1 - \frac{\alpha_{\tau_i}}{\alpha_{\tau_{i-1}}}} \quad (4.18)$$

The results obtained by Song et al. [21] show that the DDIM outperformed the DDPM for 10, 20, 50, 100 and 1000 timesteps and was only slightly worse at 1000 timesteps and $\sigma_{\tau_i} = \hat{\sigma}_{\tau_i} = \sqrt{1 - \alpha_{\tau_i}/\alpha_{\tau_{i-1}}}$

S	CIFAR10 (32×32)					CelebA (64×64)					
	10	20	50	100	1000	10	20	50	100	1000	
η	0.0	13.36	6.84	4.67	4.16	4.04	17.33	13.73	9.17	6.53	3.51
	0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79	3.64
	0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09	4.28
	1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93	5.98
$\hat{\sigma}$	367.43	133.37	32.72	9.99	3.17	299.71	183.83	71.71	45.20	3.26	

Figure 4.2: Comparison of FID scores (lower is better) of DDPMs and DDIMs trained on S timesteps, $S < T$ and $\eta = 0$ is DDIM, $\eta = 1$ and $\hat{\sigma}$ are DDPMs [21]

4.3 Guidance

Recall that when sampling from a DDPM or a DDIM, we start with an image of Gaussian noise x_T and we denoise the image until we get some image x_0 that resembles an image from the original dataset. We may however be interested in sampling something specific; for instance, from some text. In this case it is helpful to guide the denoising process towards whatever we specified.

4.3.1 Classifier guidance

Classifier guidance was introduced by Dhariwal and Nichol in their *Diffusion Models Beat GANs on Image Synthesis* paper [3] in order to improve the samples generated.

To do this, they trained a classifier $p_\Phi(y|x_t, t)$ on noisy images [3] since pre-trained classifiers are generally trained on dataset images and not on noised images.

The reverse process is then guided towards a class y using the gradient $\nabla_{x_t} \log p_\Phi(y|x_t, t)$ [3]. The authors prove [3] that if we have a reverse process such as that from the DDPM $p_\theta(x_t, x_{t+1})$, to condition it on class y , we can sample from the following process instead:

$$p_{\theta, \Phi}(x_t | x_{t+1}, y) = Z p_\theta(x_t | x_{t+1}) p_\Phi(y | x_t) \quad (4.19)$$

where Z is a normalization constant.

Using the score-based formulation, we can easily find the new objective to train the model with classifier guidance [12].

Remember that the score-based generative model learned to estimate $\nabla_{x_t} \log q(x_t)$ (3.21).

If we want to condition the model on class y , we are interested in estimating $\nabla_{x_t} \log q(x_t | y)$ [12].

With Bayes' rule, we can rewrite this as [12]:

$$q(x_t | y) = \frac{q(y | x_t) q(x_t)}{q(y)} \quad (4.20)$$

$$\Leftrightarrow \log q(x_t | y) = \log \left[\frac{q(y | x_t) q(x_t)}{q(y)} \right] \quad (4.21)$$

$$\Leftrightarrow \nabla_{x_t} \log q(x_t | y) = \nabla_{x_t} \log q(x_t) + \nabla_{x_t} \log q(y | x_t) - \nabla_{x_t} \log q(y) \quad (4.22)$$

$$= \nabla_{x_t} \log q(x_t) + \nabla_{x_t} \log q(y | x_t) \quad (4.23)$$

We can see that $\nabla_{x_t} \log q(x_t)$ is the original unconditional score from the score-based generative model and $\nabla_{x_t} \log q(y | x_t)$ is the gradient of the log-likelihood of the class y given the image x_t [3, 12].

To control the strength of the guidance, the authors introduce a hyperparameter γ that weights the

classifier gradient [3, 12]:

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{y} | \mathbf{x}_t) \quad (4.24)$$

Trivially, when $\gamma = 0$, the model doesn't use the classifier and we recover the original score-based generative model.

Dhariwal and Nichol [3] gave two algorithms to train the model with classifier guidance: one for DDPMs and one that also works for DDIMs that uses the score-based generative model formulation:

Algorithm 4 Training a DDPM/DDIM with classifier guidance

Require: *Model* $\epsilon_\theta(\mathbf{x}_t, t)$, *Classifier* $\mathbf{p}_\Phi(\mathbf{y} | \mathbf{x}_t)$ with \mathbf{y} the class label

```

1: Sample  $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ 
2: for  $t = T, T-1, \dots, 1$  do
3:    $\hat{\epsilon} \leftarrow \epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{\mathbf{x}_t} \mathbf{p}_\Phi(\mathbf{y} | \mathbf{x}_t)$ 
4:    $\mathbf{x}_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{\mathbf{x}_t - \hat{\epsilon} \sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \right) + \hat{\epsilon} \sqrt{1 - \bar{\alpha}_{t-1}}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

4.3.2 Classifier-free guidance

Despite the impressive performance of the *ablated diffusion model* (ADM) that used classifier guidance, Ho and Salimans made an improvement by using classifier-free guidance [5].

Classifier guidance requires a classifier to be trained on noisy images. This means that not only does it require training a classifier in addition to the generative model, but it also requires the classifier to be trained on noisy images for which pre-trained classifiers do not generally exist [5].

Therefore, Ho and Salimans [5] introduced a classifier-free guidance method which does not require a trained classifier at all.

The authors use an unconditional diffusion model $q(\mathbf{x}_t)$ and a conditional diffusion model $q(\mathbf{x}_t | \mathbf{y})$ where \mathbf{y} is the class label [5, 12].

Starting from equation (4.24):

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{y} | \mathbf{x}_t) \quad (4.25)$$

$$= \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \gamma (\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y}) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)) \quad (4.26)$$

$$= \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y}) - \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) \quad (4.27)$$

$$= \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y}) + \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) - \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) \quad (4.28)$$

$$= \gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y}) + (1 - \gamma) \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) \quad (4.29)$$

$\gamma \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{y})$ is the score of the conditional model and $(1 - \gamma) \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$ is the score of the unconditional model [12].

However it isn't necessary to train a conditional model and an unconditional model as we can just train the conditional model $q(\mathbf{x}_t | \mathbf{y})$ and simply ignore the conditioning information \mathbf{y} (by setting it to null, for example) to get the unconditional model $q(\mathbf{x}_t)$ [5, 12].

5 Further development

5.1 Cascaded Diffusion Models

Training a diffusion model on high-resolution images is computationally expensive and therefore time-consuming.

To address this, Ho et al. [6] introduced the cascaded diffusion model (CDM) which is essentially a pipeline of diffusion models where the first one generates a low-resolution image and the subsequent ones are trained to upsample the image to a higher resolution.

This method is much faster than training a single diffusion model on high-resolution images as most of the important features for sampling are already present in the low-resolution images which are fast to train, whereas training at bigger resolutions doesn't add much more detail considering the computation costs [6].

As mentioned, the first model is trained to generate an image as we have seen before. The subsequent models are trained using what Ho et al. call conditioning augmentation, a type of data augmentation, where the model is given a lower-resolution image and is trained on modified versions of the image to generate a higher-resolution and overall higher-quality image [6].

For low-resolution images, the authors use added Gaussian noise and use Gaussian blurring for higher resolution images [6].

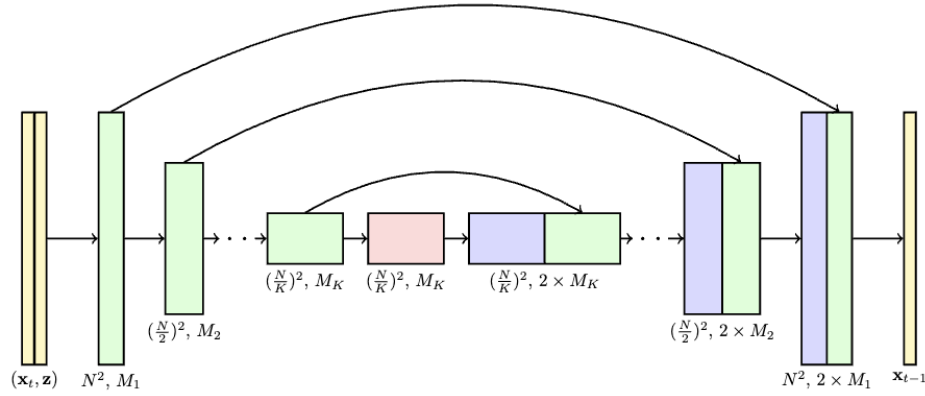


Figure 5.1: U-Net used for each of the model of a CDM pipeline: the upsampling models are given an image to upsample x_t and an upsampled low-resolution image z . For each of the blocks $1, \dots, K$, the image is upsampled/downsampled by 2, the channels are given by the channel multipliers M_1, \dots, M_K . [6]

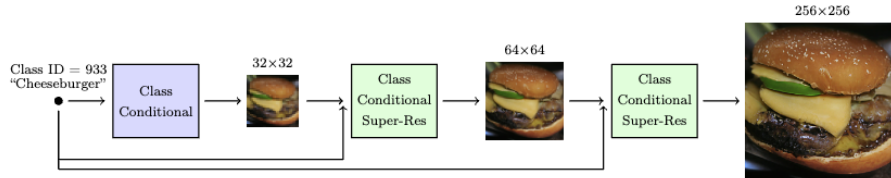


Figure 5.2: Example of a cascading diffusion pipeline for generating a 256×256 resolution image conditioned on a class [6]

As a result, the authors achieve an FID of 4.88 on ImageNet 256×256 with a CDM, where the first model generates a 32×32 image, the second model upsamples it to 64×64 and the third model upsamples it to 256×256 [6].

As an example, *Imagen* [18] from Google is a popular cascaded diffusion model that can efficiently generate high-resolution images.

5.2 Latent Diffusion Models

Latent Diffusion Models (LDMs) were introduced by Rombach et al. in their paper *High-Resolution Image Synthesis with Latent Diffusion Models* [16].

Rather than do the forward and reverse processes in pixel space, the authors propose to do them in latent space in order to speed up the computation by reducing the dimensionality of the data [16].

In order to do this, a Variational Autoencoder (VAE) [8] is used to first encode the image x into a latent representation z . The forward process is then applied to the latent representation z until we obtain z_T [16]. As for the reverse process, the model (a U-Net in their paper) will denoise the latent representation z_T until we obtain \tilde{z} which is then passed through the decoder of the VAE to obtain the denoised image \tilde{x} [16]. This means that given an image x , the encoder \mathcal{E} will encode x into a latent representation z :

$$z = \mathcal{E}(x) \quad (5.1)$$

Once the forward process and reverse process are done in the latent space, the decoder \mathcal{D} will decode the denoised latent representation \tilde{z} into the denoised image \tilde{x} :

$$\tilde{x} = \mathcal{D}(\tilde{z}) \quad (5.2)$$

From the training objective (2.99), we get a new training objective for LDMs [16]:

$$\mathcal{L}_{\text{LDM}} = \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0, I), t} [\|\epsilon - \epsilon_\theta(z_t, t)\|_2^2] \quad (5.3)$$

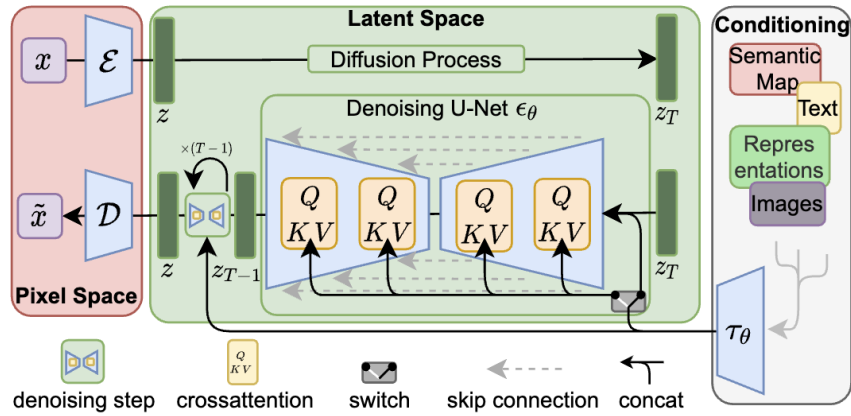


Figure 5.3: Latent Diffusion processes [16]

In terms of results, Rombach et al. [16] managed to achieve an FID of 3.60 on ImageNet 256×256 with 271 days of training on a Nvidia V100 GPU, which is better than the 3.85 FID achieved by Dhariwal and Nichol [3] with their ADM-G and ADM-U (guidance and upsampling) models which required 349 days of training on a V100 [16].

An example of a current LDM is *Stable Diffusion* [1, 16], which has become very popular nowadays and is easily available for use to the general public.

5.3 Diffusion Transformer

Recently, a new alternative to the U-Net architecture has been proposed by Peebles and Xie in their paper *Scalable Diffusion Models with Transformers* [15].

As the name suggests, the Diffusion Transformer (DiT) architecture is based on the Transformer architecture [25] used for Natural Language Processing (NLP).

The architecture builds upon Latent Diffusion Models (LDMs) [16], since it takes as input the latent representation z .

The first step is called *patchify*: given that we have a spatial input, this input is to be converted into tokens. Positional encodings are subsequently added to the tokens to give the model information about the spatial structure of the image [15].

This information is then passed through DiT blocks (see figure 5.4) before going through layer normalisation.

Finally, the output is passed through a linear layer and then reshaped back in order to return the noise and the covariance matrix Σ [15].

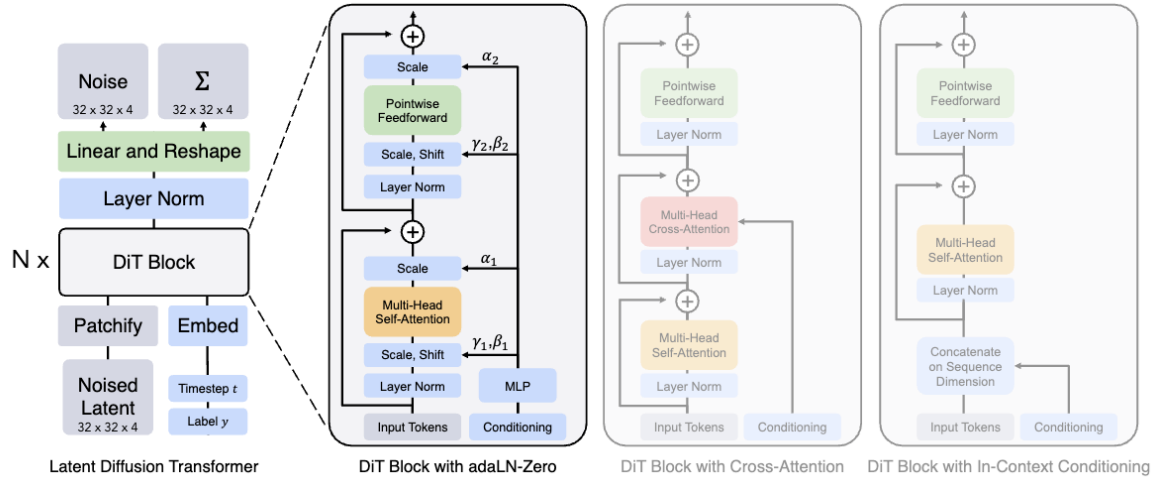


Figure 5.4: Diffusion Transformer architecture [15]

Compared to the previous experiments on class-conditioned ImageNet 256×256 , the DiT model has the best FID of them all with 2.27 [15], showing that the DiT brings significant improvement to sample quality.

Newer models such as *Sora* [2] and *Stable Diffusion 3* use this new DiT architecture and have achieved impressive visual results.

6 Implementation

A simple implementation of a DDPM for the MNIST dataset in PyTorch is given in the Jupyter notebook `DMID.ipynb`. The code aims to be as close as possible to the original DDPM paper [4]. It can be run on a CPU or a CUDA-enabled GPU.

The code provides an explanation and an implementation of a U-Net architecture similar to the original one.

Ho et al. used a self-attention block not only between the two ResNet-like blocks in the bottleneck, but also at the encoder and decoder stages at 16×16 resolution image steps. They trained the model as such for CIFAR-10, CelebA-HQ and LSUN datasets [4].

`DMID.ipynb` contains a model without self-attention blocks as this visually gave worse results on the MNIST dataset, possibly because of overfitting.

The forward process uses the linear noising schedule β and the reverse process uses the same 1000 timesteps from the original DDPM paper [4], however the code is easily modifiable to use a different noising schedule (e.g. cosine noising schedule) and a different number of timesteps. Similarly, a learning rate of 10^{-4} was used as it was found to work well for the MNIST dataset but this can also be changed to the original 2×10^{-4} learning rate.

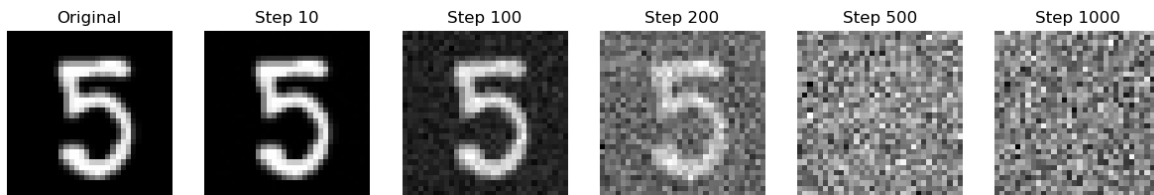


Figure 6.1: Noised images at different timesteps of the forward process

The model contains 7'339'297 parameters, much less than Ho et al.'s CIFAR-10 model which has 35.7 mln parameters and the CelebA-HQ model with its 114 mln parameters [4].

On an Apple M2 Pro chip using MPS, training takes around 70s per epoch for 1000 timesteps and 32×32 MNIST images and sampling 1 image takes around 13.4s.

This means that the model is not so complex but still takes a long time to train and sample from without significant resources such as a TPU v3-8 used by Ho et al. [4].

The code is available at <https://github.com/gregorysedykh/diffusion-models-in-depth>.



Figure 6.2: Samples generated by the model

7 Conclusion

Throughout this report, we analysed Denoising Diffusion Probabilistic Models (DDPMs) in detail by covering their mathematical foundation, which includes the two Markovian processes and the simplified training objective. Additionally, we explored the U-Net architecture employed in training these models.

The score-based generative model was also analysed, particularly its use of score matching to estimate the score of the data distribution, and demonstrated its equivalence to DDPMs.

Several improvements to DDPMs were discussed, notably Denoising Diffusion Implicit Models (DDIMs) and the guidance methods, both playing a crucial role in simultaneously improving sample quality and speeding up the sampling process.

Subsequently, we explored the present and the future of diffusion models, especially the Diffusion Transformer (DiT) architecture which is now becoming increasingly popular for models available to the public.

Finally, a straightforward implementation of a DDPM for the MNIST dataset implemented using Python and PyTorch was presented, together with some samples generated by the model.

A Convolution of two Gaussian distributions

Equation (2.7) had given us:

$$x_t = \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \mathcal{N}(0, \alpha_t (1 - \alpha_{t-1}) I) + \mathcal{N}(0, (1 - \alpha_t) I)$$

We can find the distribution of the sum of two independant 1D random variables $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ as the convolution of two Gaussian distributions.

Let $Z = X + Y$ with X and Y independent.

$$\begin{aligned} f_Z(z) &= \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma_X^2}} \exp\left[-\frac{(x-\mu_X)^2}{2\sigma_X^2}\right] \frac{1}{\sqrt{2\pi\sigma_Y^2}} \exp\left[-\frac{(z-x-\mu_Y)^2}{2\sigma_Y^2}\right] dx \\ &= \frac{1}{2\pi\sigma_X\sigma_Y} \int_{-\infty}^{\infty} \exp\left[-\frac{(x-\mu_X)^2}{2\sigma_X^2}\right] \exp\left[-\frac{(z-x-\mu_Y)^2}{2\sigma_Y^2}\right] dx \\ &= \frac{1}{2\pi\sigma_X\sigma_Y} \int_{-\infty}^{\infty} \exp\left[-\frac{(x-\mu_X)^2}{2\sigma_X^2} - \frac{(z-x-\mu_Y)^2}{2\sigma_Y^2}\right] dx \\ &= \frac{1}{2\pi\sigma_X\sigma_Y} \int_{-\infty}^{\infty} \exp\left[-\frac{x^2 - 2\mu_X x + \mu_X^2}{2\sigma_X^2} - \frac{z^2 + x^2 - 2xz - 2z\mu_Y + 2\mu_Y x + \mu_Y^2}{2\sigma_Y^2}\right] dx \\ &= \frac{1}{2\pi\sigma_X\sigma_Y} \int_{-\infty}^{\infty} \exp\left[-\frac{\sigma_Y^2(x^2 - 2\mu_X x + \mu_X^2)}{2\sigma_X^2\sigma_Y^2} - \frac{\sigma_X^2(z^2 + x^2 - 2xz - 2z\mu_Y + 2\mu_Y x + \mu_Y^2)}{2\sigma_X^2\sigma_Y^2}\right] dx \\ &= \frac{1}{2\pi\sigma_X\sigma_Y} \int_{-\infty}^{\infty} \exp\left[-\frac{\sigma_Y^2 x^2 - 2\mu_X \sigma_Y^2 x + \mu_X^2 \sigma_Y^2 + \sigma_X^2 z^2 + \sigma_X^2 x^2 - 2xz\sigma_X^2 + 2z\mu_Y \sigma_X^2 - 2\mu_Y x \sigma_X^2 + \sigma_X^2 \mu_Y^2}{2\sigma_X^2\sigma_Y^2}\right] dx \\ &= \frac{1}{2\pi\sigma_X\sigma_Y} \int_{-\infty}^{\infty} \exp\left[-\frac{(\sigma_X^2 + \sigma_Y^2) x^2 - 2x(\mu_X \sigma_Y^2 + z\sigma_X^2 - \mu_Y \sigma_X^2) + \sigma_X^2 (z - \mu_Y)^2 + \mu_X^2 \sigma_Y^2}{2\sigma_X^2\sigma_Y^2}\right] dx \end{aligned}$$

Let $\sigma_Z = \sqrt{\sigma_X^2 + \sigma_Y^2}$.

$$\begin{aligned} &= \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sigma_Z}} \frac{1}{\sqrt{2\pi} \sigma_Z} \int_{-\infty}^{\infty} \exp\left[-\frac{x^2 - 2x \frac{\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X}{\sigma_Z^2} + \frac{\sigma_X^2 (z - \mu_Y)^2 + \mu_X^2 \sigma_Y^2}{\sigma_Z^2}}{2 \left(\frac{\sigma_X \sigma_Y}{\sigma_Z}\right)^2}\right] dx \\ &= \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sigma_Z}} \frac{1}{\sqrt{2\pi} \sigma_Z} \int_{-\infty}^{\infty} \exp\left[-\frac{\left(x - \frac{\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X}{\sigma_Z^2}\right)^2 - \left(\frac{\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X}{\sigma_Z^2}\right)^2 + \frac{\sigma_X^2 (z - \mu_Y)^2 + \mu_X^2 \sigma_Y^2}{\sigma_Z^2}}{2 \left(\frac{\sigma_X \sigma_Y}{\sigma_Z}\right)^2}\right] dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sigma_Z}} \exp\left[-\frac{\left(x - \frac{\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X}{\sigma_Z^2}\right)^2}{2 \left(\frac{\sigma_X \sigma_Y}{\sigma_Z}\right)^2}\right] dx \end{aligned}$$

$$\begin{aligned}
& \frac{1}{\sqrt{2\pi}\sigma_Z} \exp \left[-\frac{\sigma_Z^2 (\sigma_X^2 (z - \mu_Y)^2 + \sigma_Y^2 \mu_X^2) - (\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X)^2}{2\sigma_Z^2 (\sigma_X \sigma_Y)^2} \right] dx \\
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sigma_Z}} \exp \left[-\frac{\left(x - \frac{\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X}{\sigma_Z^2} \right)^2}{2 \left(\frac{\sigma_X \sigma_Y}{\sigma_Z} \right)^2} \right] \frac{1}{\sqrt{2\pi}\sigma_Z} \exp \left[-\frac{(z - (\mu_X + \mu_Y))^2}{2\sigma_Z^2} \right] dx \\
&= \frac{1}{\sqrt{2\pi}\sigma_Z} \exp \left[-\frac{(z - (\mu_X + \mu_Y))^2}{2\sigma_Z^2} \right] \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sigma_Z}} \exp \left[-\frac{\left(x - \frac{\sigma_X^2 (z - \mu_Y) + \sigma_Y^2 \mu_X}{\sigma_Z^2} \right)^2}{2 \left(\frac{\sigma_X \sigma_Y}{\sigma_Z} \right)^2} \right] dx \\
&= \frac{1}{\sqrt{2\pi}\sigma_Z} \exp \left[-\frac{(z - (\mu_X + \mu_Y))^2}{2\sigma_Z^2} \right] \\
&= \mathcal{N}(z; \mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)
\end{aligned}$$

Given this, we can refer back to our case and therefore we find:

$$\begin{aligned}
x_t &= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \mathcal{N}(0, \alpha_t (1 - \alpha_{t-1}) I) + \mathcal{N}(0, (1 - \alpha_t) I) \\
&= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \mathcal{N}(0, \alpha_t ((1 - \alpha_{t-1}) + 1 - \alpha_t) I)
\end{aligned}$$

References

- [1] Stability AI. URL: <https://github.com/Stability-AI/stablediffusion>.
- [2] Tim Brooks et al. “Video generation models as world simulators”. In: (2024). URL: <https://openai.com/research/video-generation-models-as-world-simulators>.
- [3] Prafulla Dhariwal and Alex Nichol. *Diffusion Models Beat GANs on Image Synthesis*. 2021. arXiv: [2105.05233 \[cs.LG\]](#).
- [4] Jonathan Ho, Ajay Jain, and Pieter Abbeel. *Denoising Diffusion Probabilistic Models*. 2020. arXiv: [2006.11239 \[cs.LG\]](#).
- [5] Jonathan Ho and Tim Salimans. *Classifier-Free Diffusion Guidance*. 2022. arXiv: [2207.12598 \[cs.LG\]](#).
- [6] Jonathan Ho et al. *Cascaded Diffusion Models for High Fidelity Image Generation*. 2021. arXiv: [2106.15282 \[cs.CV\]](#).
- [7] Aapo Hyvarinen. *Estimation of Non-Normalized Statistical Models by Score Matching*. 2005. URL: <https://jmlr.org/papers/volume6/hyvarinen05a/hyvarinen05a.pdf>.
- [8] Diederik P Kingma and Max Welling. *Auto-Encoding Variational Bayes*. 2022. arXiv: [1312.6114 \[stat.ML\]](#).
- [9] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. “CIFAR-10 (Canadian Institute for Advanced Research)”. In: (2024). URL: <http://www.cs.toronto.edu/~kriz/cifar.html>.
- [10] Zhitong Lai et al. *Rethinking Skip Connections in Encoder-decoder Networks for Monocular Depth Estimation*. 2022. arXiv: [2208.13441 \[cs.CV\]](#).
- [11] Ziwei Liu et al. “Deep Learning Face Attributes in the Wild”. In: *Proceedings of International Conference on Computer Vision (ICCV)*. Dec. 2015.
- [12] Calvin Luo. *Understanding Diffusion Models: A Unified Perspective*. 2022. arXiv: [2208.11970 \[cs.LG\]](#).
- [13] Alex Nichol and Prafulla Dhariwal. *Improved Denoising Diffusion Probabilistic Models*. 2021. arXiv: [2102.09672 \[cs.LG\]](#).
- [14] Aaron van den Oord et al. *Conditional Image Generation with PixelCNN Decoders*. 2016. arXiv: [1606.05328 \[cs.CV\]](#).
- [15] William Peebles and Saining Xie. *Scalable Diffusion Models with Transformers*. 2023. arXiv: [2212.09748 \[cs.CV\]](#).
- [16] Robin Rombach et al. *High-Resolution Image Synthesis with Latent Diffusion Models*. 2022. arXiv: [2112.10752 \[cs.CV\]](#).
- [17] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. *U-Net: Convolutional Networks for Biomedical Image Segmentation*. 2015. arXiv: [1505.04597 \[cs.CV\]](#).
- [18] Chitwan Saharia et al. *Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding*. 2022. arXiv: [2205.11487 \[cs.CV\]](#).
- [19] Tim Salimans and Diederik P. Kingma. *Weight Normalization: A Simple Reparameterization to Accelerate Training of Deep Neural Networks*. 2016. arXiv: [1602.07868 \[cs.LG\]](#).
- [20] Jascha Sohl-Dickstein et al. *Deep Unsupervised Learning using Nonequilibrium Thermodynamics*. 2015. arXiv: [1503.03585 \[cs.LG\]](#).
- [21] Jiaming Song, Chenlin Meng, and Stefano Ermon. *Denoising Diffusion Implicit Models*. 2022. arXiv: [2010.02502 \[cs.LG\]](#).

-
- [22] Yang Song. *Blog post: Generative Modeling by Estimating Gradients of the Data Distribution*. 2021. URL: <https://yang-song.net/blog/2021/score/>.
 - [23] Yang Song and Stefano Ermon. *Generative Modeling by Estimating Gradients of the Data Distribution*. 2020. arXiv: [1907.05600 \[cs.LG\]](#).
 - [24] Yang Song et al. *Sliced Score Matching: A Scalable Approach to Density and Score Estimation*. 2019. arXiv: [1905.07088 \[cs.LG\]](#).
 - [25] Ashish Vaswani et al. *Attention Is All You Need*. 2023. arXiv: [1706.03762 \[cs.CL\]](#).
 - [26] Pascal Vincent. *A Connection Between Score Matching and Denoising Autoencoders*. 2010. URL: https://www.iro.umontreal.ca/~vincentp/Publications/smdae_techreport.pdf.
 - [27] Max Welling and Yee Whye Teh. *Bayesian Learning via Stochastic Gradient Langevin Dynamics*. 2011. URL: <https://www.stats.ox.ac.uk/~teh/research/compstats/WelTeh2011a.pdf>.
 - [28] Lilian Weng. “What are diffusion models?” In: *lilianweng.github.io* (July 2021). URL: <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>.
 - [29] Yuxin Wu and Kaiming He. *Group Normalization*. 2018. arXiv: [1803.08494 \[cs.CV\]](#).
 - [30] Sergey Zagoruyko and Nikos Komodakis. *Wide Residual Networks*. 2017. arXiv: [1605.07146 \[cs.CV\]](#).