

# Text-to-image DDPMs

Giovani Tavares  
giovanitavares@usp.br

University of São Paulo — December 11, 2025

## 1 Goal

The goals of the present work are to first document the Deriving the training and sampling algorithms of Denoising Diffusion Probabilistic Models (DDPMs) and also use the insights obtained in such derivations to modify the DDPM implementation in [1] so that it becomes a *text-to-image* generative model, instead of a model that generates random images. In summary, [1] implements a DDPM that generates random images., and what I did was to modify it so that the images are not random, but instead it is possible to pass the class parameter to the sampling algorithm so it outputs an image from such class.

With the knowledge and insights obtained in the derivation of the training and sampling algorithms of DDPMs described in the first section, we will implement a conditional DDPM with FiLM layers [3] attached to this report. FiLM layers are what lets one marginalize the generated images distribution learned by the model on their classes, so that the result is a *text-to-image* DDPM.

The sections of the present work are:

- **What Are DDPMs:** Deriving the training and sampling algorithms of Denoising Diffusion Probabilistic Models (DDPMs) using the several references shown in the end of this work
- **Conditional DDPMs - A text-to-image Generative Model:** Deriving the training and sampling algorithms of a *conditional* DDPMs that lets one guide the generation of images of specific classes instead of random ones as in the original DDPM

## 2 What are DDPMs?

Denoising Diffusion Probabilistic Models (DDPMs) are models capable of predicting *noise* from a noisy input. By using such prediction, a sampling algorithm can be used to remove the noise from the input which results in a denoised output.

DDPMs are made of two processes: a **forward** and a **reversion** process. The former is responsible for gradually adding noise to a image by sampling from a normal distribution according to a Markov Chain. The latter removes added noise by sampling from another normal distribution. In simple terms, the training of DDPMs involve learning the reversion process' distribution's parameters.

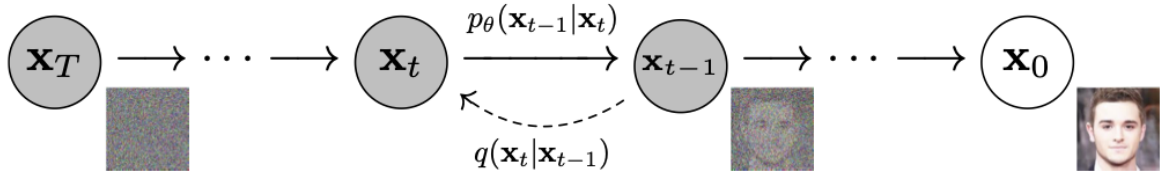


Figure 1: Image extracted from the original DDPM paper: "Denoising Diffusion Probabilistic Models" [2]

## 2.1 Forward Process $q$

The process of adding noise to an input image ( $x_0$ ) is a Markov Chain that generates a noisier image  $x_t$  from a less noisy image  $x_{t-1}$ . Hence,  $x_t$  represents the result of adding noise to  $x_{t-1}$  by transitioning it **once** in the Markov Chain.

From the original DDPM paper, we know that in the forward process, a sample  $x_t$  is produced by adding noise to a sample  $x_{t-1}$  according to a normal distribution defined below:

$$q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1 - \beta_t} \times x_{t-1}; \beta_t I) \quad (1)$$

**Ideally, one would need a single transition from  $x_0$  to get to  $x_t$ , with  $t > 0$ .**

The authors of [2] achieves such ideal scenario by defining a cumulative noise  $\alpha_t$  presented below:

**Definition 2.1** (Cumulative Noise). *The Cumulative Noise ( $\alpha_t$ ) added to an input  $x_0$  up to the  $t$ -th step is defined as:*

$$\alpha_t := 1 - \beta_t \quad (2)$$

$$\bar{\alpha}_t := \prod_{s=1}^t \alpha_s \quad (3)$$

which leads to (4)

$$q(x_t|x_0) := \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} \times x_0; (1 - \bar{\alpha}_t)I) \quad (5)$$

Definition 2.1 permits the sampling of  $x_t$  to be done in a single transition from  $x_0$ , instead of having to produce every intermediate image.

According to the original DDPM paper, the **Variance Schedule**  $\beta_1, \dots, \beta_T$  sequence that defines the noisy images distribution are held constant.

## 2.2 Reverse Process $p_\theta$

The Reverse Process is a Markov Chain with a Standard Gaussian initial state and Gaussian transition distribution  $p_\theta$  parametrized by mean  $\mu_\theta$  and variance  $\Sigma_\theta$ . **The core idea of the chain is that the transitions remove each a little bit of the noise from the initial state up until the noise-free state  $x_0$ .**

**Definition 2.2** (Reverse Process Transitions). *The Reverse Process is a Markov Chain with the following transitions:*

$$p(x_T) = \mathcal{N}(x_T; 0; 1) \quad (6)$$

$$p_\theta(x_{0:T}) := p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \quad (7)$$

$$p_\theta(x_{t-1}|x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t); \Sigma_\theta(x_t, t)) \quad (8)$$

What we ultimately want is to have a  $x_0$  that is as likely as possible. We can use the standard rule of probability to obtain a marginalization of  $p(x_0)$  using the latent variables  $\mathbf{x}_{1:T}$

**Definition 2.3** (Reverse Process Prior). *Ideally, the Reverse Process would let us sample from:*

$$p(\mathbf{x}_0) = \int p(\mathbf{x}_0, \mathbf{x}_{1:T}) d\mathbf{x}_{1:T} \quad (9)$$

$$(10)$$

From the defition above, we see that  $p(\mathbf{x}_0)$  is very complex due to its **multidimensionality**, which makes it intractable. That is why in DDPMs,  $p(\mathbf{x}_0)$  is never computed directly, but instead its lower bound.

### 2.2.1 Evidence Lower Bound / ELBO

The Evidence Lower Bound is a tight lower bound that limits  $\log(p(\mathbf{x}_0))$  from below. Hence, when  $p(\mathbf{x}_0)$  is intractable as in the case of DDPMs, one can always maximize such lower bound as a means to ensure that  $\log(p(\mathbf{x}_0))$  is as large as possible. Such lower bound is often called **ELBO** and will be demonstrated using two different approaches: the **Jensen's Inequality** and the **KL Divergence**.

**Definition 2.4** (Evidence Lower Bound - Jensen's Inequality). *Let's use the Rule Of Total Probability to find a lower bound for the log-likelihood function.*

$$\log[p_\theta(\mathbf{x}_0)] = \log \int_{\mathbf{x}_{1:T}} p(\mathbf{x}_0, \mathbf{x}_{1:T}) d\mathbf{x}_{1:T} \quad (11)$$

$$\log[p_\theta(\mathbf{x}_0)] = \log \int_{\mathbf{x}_{1:T}} p(\mathbf{x}_0, \mathbf{x}_{1:T}) \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \quad (12)$$

$$\log[p_\theta(\mathbf{x}_0)] = \log(\mathbb{E}_q \left[ \frac{p(\mathbf{x}_0, \mathbf{x}_{1:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]) \quad (13)$$

$$\text{the Jensen's inequality tells us} \quad (14)$$

$$f(\mathbb{E}(\mathbf{X})) \geq \mathbb{E}(f(\mathbf{X})) \quad (15)$$

$$\text{for any concave function } f. \quad (16)$$

$$\log \text{ is concave, hence:} \quad (17)$$

$$\log[p_\theta(\mathbf{x}_0)] \geq \mathbb{E}_q \left[ \log \frac{p(\mathbf{x}_0, \mathbf{x}_{1:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (18)$$

$$\text{If we define the Evidence Lower Bound } L \text{ as:} \quad (19)$$

$$L := \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_0, \mathbf{x}_{1:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (20)$$

$$\implies -\log[p_\theta(\mathbf{x}_0)] \leq L \quad (21)$$

**Definition 2.5** (Evidence Lower Bound - KL Divergence). *In order to reverse the forward process, we need that the forward process' distribution  $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$  is as close to  $p(\mathbf{x}_{1:T}|\mathbf{x}_0)$  as possible. We can use the Kullback-Leibler (KL) divergence between  $q$  and  $p$  ( $\mathbf{D}_{\text{KL}}$ ) to evaluate their difference as find a bound to  $\log[p_\theta(\mathbf{x}_0)]$ .*

$$\mathbf{D}_{\text{KL}}[q(\mathbf{x}_{1:T}|\mathbf{x}_0)||p(\mathbf{x}_{1:T}|\mathbf{x}_0)] := \mathbb{E}_q[\log(q(\mathbf{x}_{1:T}|\mathbf{x}_0) - \log(p(\mathbf{x}_{1:T}|\mathbf{x}_0)))] \quad (22)$$

$$\text{using the Bayes' Rule we can write} \quad (23)$$

$$p(\mathbf{x}_{1:T}|\mathbf{x}_0) = \frac{p(\mathbf{x}_{1:T}, \mathbf{x}_0)}{p(\mathbf{x}_0)} \quad (24)$$

$$\implies \mathbf{D}_{\text{KL}}[q(\mathbf{x}_{1:T}|\mathbf{x}_0)||p(\mathbf{x}_{1:T}|\mathbf{x}_0)] = \mathbb{E}_q[\log(q(\mathbf{x}_{1:T}|\mathbf{x}_0) - \log(p(\mathbf{x}_{1:T}|\mathbf{x}_0)) + \log(p(\mathbf{x}_0)))] \quad (25)$$

$$\text{the prior of the latent variables does not depend on } q \quad (26)$$

$$\mathbf{D}_{\text{KL}}[q(\mathbf{x}_{1:T}|\mathbf{x}_0)||p(\mathbf{x}_{1:T}|\mathbf{x}_0)] = \mathbb{E}_q[\log(q(\mathbf{x}_{1:T}|\mathbf{x}_0) - \log(p(\mathbf{x}_{1:T}, \mathbf{x}_0))] + \log(p(\mathbf{x}_0)) \quad (27)$$

$$\implies \log(p(\mathbf{x}_0)) = \mathbf{D}_{\text{KL}}[q(\mathbf{x}_{1:T}|\mathbf{x}_0)||p(\mathbf{x}_{1:T}|\mathbf{x}_0)] - \mathbb{E}_q[\log(q(\mathbf{x}_{1:T}|\mathbf{x}_0) - \log(p(\mathbf{x}_{1:T}, \mathbf{x}_0))] \quad (28)$$

$$\log(p(\mathbf{x}_0)) = \mathbf{D}_{\text{KL}}[q(\mathbf{x}_{1:T}|\mathbf{x}_0)||p(\mathbf{x}_{1:T}|\mathbf{x}_0)] + \mathbb{E}_q[\log(p(\mathbf{x}_{1:T}, \mathbf{x}_0) - \log(q(\mathbf{x}_{1:T}|\mathbf{x}_0))] \quad (29)$$

$$\text{but } \mathbf{D}_{\text{KL}}[q(\mathbf{x}_{1:T}|\mathbf{x}_0)||p(\mathbf{x}_{1:T}|\mathbf{x}_0)] \geq 0 \quad (30)$$

$$\implies -\log(p(\mathbf{x}_0)) \leq \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_0, \mathbf{x}_{1:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (31)$$

$$\text{If we define the Evidence Lower Bound } L \text{ as:} \quad (32)$$

$$L := \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_0, \mathbf{x}_{1:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (33)$$

$$\implies -\log[p(\mathbf{x}_0)] \leq L \quad (34)$$

We have found upper bound  $L$  for the negative log-likelihood function that can be maximized in the forward process' training.

### 2.2.2 Noise Predictor Training

In order to use  $L$  as the loss function in our training, further algebraic manipulation must be performed

**Definition 2.6** (Noise Predictor’s Loss Derivation). *In order to build the Noise Predictor’s loss function, we need to remember that both forward and reverse processes are Markov Chains and use this fact to manipulate  $L$ .*

$$L = \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (35)$$

$$\text{The forward and reverse processes are Markov Chains, so} \quad (36)$$

$$\mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] = \mathbb{E}_q \left[ \log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \quad (37)$$

$$= \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (38)$$

$$= \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (39)$$

$$= \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) + \sum_{t=2}^T \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (40)$$

$$= \mathbb{E}_q \left[ -\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) - \mathbb{E}_q \left[ \sum_{t=2}^T \mathbf{D}_{\text{KL}} [q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)] \right] \quad (41)$$

We see that the first term is parameter free, because  $p(\mathbf{x}_T)$  is fixed and defined as a Gaussian, while  $q(\mathbf{x}_T|\mathbf{x}_0)$  is also Gaussian from the definition of the forward process. Hence, we are left with the second and third terms from  $L$ .

More specifically, we can conclude that maximizing ELBO ( $L$ ) is equivalent to minimizing the KL-Divergence between  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$  and  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$

We know that both distributions are Gaussians, which makes computing the KL Divergence between them easier if we know their mean and variance. We will begin by calculating such moments for  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ .

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \quad (42)$$

$$\text{we know the } q \text{ distribution from Definition 2.1, hence} \quad (43)$$

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \text{ is a product of known Gaussians over another known Gaussian that lets us define} \quad (44)$$

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\bar{\alpha}_t}\mathbf{x}_t + (1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{(1 - \bar{\alpha}_t)} \quad (45)$$

$$\Sigma_q(t) = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I} \quad (46)$$

$$\implies q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_q(\mathbf{x}_t, \mathbf{x}_0); \Sigma_q(t)) \quad (47)$$

We have just defined the  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$  distribution. Now let’s move on to the  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  distribution.

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t); \Sigma_\theta(t)) \quad (48)$$

$$\text{the reverse process variance is defined as the ground truth variance of the forward process:} \quad (49)$$

$$\Sigma_\theta(t) = \Sigma_q(t) \quad (50)$$

$$\text{we are only left with the distribution’s mean } \mu_\theta \quad (51)$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t); \Sigma_q(t)) \quad (52)$$

Equation 50 makes it much easier to calculate  $\mathbf{D}_{\text{KL}} [q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)]$ :

**Now we know we are trying to compute the KL Divergence between two Gaussians with the exact same variance.**

For that, there is the following result that arises from the definition of such divergence

$$d_1(x) = \mathcal{N}(\mu_1, \sigma^2) \quad (53)$$

$$d_2(x) = \mathcal{N}(\mu_2, \sigma^2) \quad (54)$$

$$\text{The KL divergence } D_{KL}(d_1|d_2) \text{ is given by:} \quad (55)$$

$$D_{KL}(d_1|d_2) = \frac{(\mu_1 - \mu_2)^2}{2\sigma^2} \quad (56)$$

$$\text{Hence,} \quad (57)$$

$$\mathbf{D}_{KL}[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] = \mathbf{D}_{KL}(\mathcal{N}(\mathbf{x}_{t-1}; \mu_q(\mathbf{x}_t, \mathbf{x}_0); \Sigma_q(t)), \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t); \Sigma_q(t))) \quad (58)$$

$$= \frac{1 - \bar{\alpha}_t}{2(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \|(\mu_q - \mu_\theta)_2\|^2 \quad (59)$$

As our goal is to minimize the KL Divergence, from equation 59 we see that such goal comes down to basically minimizing the difference  $\mu_q - \mu_\theta$ , i.e.,

**We just need to minimize the difference between the means of the reverse and forward processes' distributions.** We need to define the reverse process' distribution mean ( $\mu_\theta$ ) prediction by taking a look at the forward process' one ( $\mu_q$ ).

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}\mathbf{x}_t + (1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{(1 - \bar{\alpha}_t)} \quad (60)$$

$$\text{we can define the prediction} \quad (61)$$

$$\hat{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \mu_\theta(\mathbf{x}_t) \quad (62)$$

$$= \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}\mathbf{x}_t + (1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_\theta}{(1 - \bar{\alpha}_t)} \quad (63)$$

In equation 63 we see that we are using our reverse process model's prediction  $\mathbf{x}_\theta$  in the prediction of its distribution's mean, which let's us rewrite  $\|(\mu_\theta - \mu_q)_2\|^2$  in terms of  $\mathbf{x}_0$  and  $\mathbf{x}_\theta$  which leaves us with the following for the KL Divergence:

$$\mathbf{D}_{KL}(\mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t); \Sigma_q(t)), \mathcal{N}(\mathbf{x}_{t-1}; \mu_q(\mathbf{x}_t, \mathbf{x}_0); \Sigma_q(t))) = \frac{(1 - \bar{\alpha}_t)(\bar{\alpha}_{t-1})}{2(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \|(\mathbf{x}_\theta - \mathbf{x}_0)_2\|^2 \quad (64)$$

The author's of the DDPM paper mention that equation 64 can be used as the loss function to train the reverse process model. On the other hand, we now that the forward process actually predict the noise that was added to an input  $\mathbf{x}_t$  instead of predicting  $\mathbf{x}_\theta$  directly. This means that the loss function must account for the error prediction somehow. This is achieved by further analysing  $\mathbf{x}_0$  and  $\mathbf{x}_\theta$  and remembering how  $\mathbf{x}_\theta$  was defined in the forward process.

$$q(x_t|x_0) := \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} \times x_0; (1 - \bar{\alpha}_t)\mathbb{I}) \quad (65)$$

$$\text{which let's us write} \quad (66)$$

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \quad (67)$$

$$\Rightarrow \mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon}{\sqrt{\bar{\alpha}_t}} \quad (68)$$

$$\text{for a Standard Gaussian Noise } \epsilon. \quad (69)$$

$$\text{We can now define our prediction } \hat{\epsilon} = \epsilon_\theta \quad (70)$$

$$\mathbf{x}_\theta = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta}{\sqrt{\bar{\alpha}_t}} \quad (71)$$

$$\Rightarrow \frac{(1 - \bar{\alpha}_t)(\bar{\alpha}_{t-1})}{2(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \|(\mathbf{x}_\theta - \mathbf{x}_0)_2\|^2 = \frac{(1 - \bar{\alpha}_t)(\bar{\alpha}_{t-1})}{2(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \|(\epsilon_\theta - \epsilon)_2\|^2 \quad (72)$$

Even though equation 72 could be used directly as the loss function for the noise predictor, the DDPM paper authors mention that optimizing  $\|(\epsilon_\theta - \epsilon)_2^2\|$  without the scaling factor with the cumulative noise  $\alpha_t$  is enough. Hence, we have finally defined a function to be minimized for the noise predictor training and hence write its algorithm.

---

**Algorithm 1:** Noise Predictor Training

---

- 1: **repeat**
  - 2:  $\mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0)$  ▷ Sample image from training set
  - 3:  $\mathbf{x}_0 \sim \text{Uniform}(\{1, \dots, T\})$  ▷ Sample the step of the Forward Process Markov Chain
  - 4:  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  ▷ Sample standard gaussian noise to be added to the input
  - 5:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$  ▷ Forward Process/ Generating Noisy Image
  - 6: Take Gradient Descent Step on  $\nabla_\theta(\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|)$
  - 7: **until** converged
- 

### 2.2.3 Sampling Algorithm Derivation

Now that we have defined a way to predict an input image  $\mathbf{x}_t$ 's noise, we need a way to use such prediction to reconstruct the original de-noised image  $\mathbf{x}_0$ , i.e., we need a way to sample from  $p(\mathbf{x}_0)$ . To do so, let's recall how we have defined the  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  distribution.

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t); \Sigma_q(t)) \quad (73)$$

$$\mu_\theta(\mathbf{x}_t) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}\mathbf{x}_t + (1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_\theta}{(1 - \bar{\alpha}_t)} \quad (74)$$

$$\mathbf{x}_\theta = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta}{\sqrt{\bar{\alpha}_t}} \quad (75)$$

$$\implies \mu_\theta(\mathbf{x}_t) = \frac{\mathbf{x}_t}{\sqrt{\alpha_t}} - \frac{(1 - \alpha_t)(\sqrt{1 - \bar{\alpha}_t})}{(1 - \bar{\alpha}_t)(\sqrt{\alpha_t})}\epsilon_\theta = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta\right) \quad (76)$$

$$\Sigma_\theta(t) = \Sigma_q(t) = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I} \quad (77)$$

We now have defined  $\mathbf{x}_{t-1}$ 's mean and variance given  $\mathbf{x}_t$  which let's us write it as (78)

$$\mathbf{x}_{t-1} = \mu_\theta(\mathbf{x}_t) + \sqrt{\Sigma_\theta(t)}\mathbf{z} \quad (79)$$

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (80)$$

We have now a way to generate  $\mathbf{x}_{t-1}$  from  $\mathbf{x}_t$ . This means that if we have  $\mathbf{x}_1$  we can generate  $\mathbf{x}_0$ . This let's us finally define our sampling algorithm:

---

**Algorithm 2:** Sampling

---

- 1: **repeat**
  - 2:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  ▷ Sample random noisy image
  - 3:  $\mathbf{T} \sim \text{Uniform}(\{1, \dots, 1000\})$  ▷ Sample random length of the Denoising Chain. Max chain size of 1000 was set arbitrarily
  - 4: **for**  $t = \mathbf{T}, \dots, 1$  **do**
  - 5:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$  else  $\mathbf{z} = \mathbf{0}$
  - 6:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)\right) + \sqrt{\Sigma_\theta(t)}\mathbf{z}$  ▷ Sampling  $\mathbf{x}_{t-1}$  from  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$
  - 7: **end for**
  - 8: **return**  $\mathbf{x}_0$
-

### 3 Conditional DDPMs - A text-to-image Generative Model

In a standard Denoising Diffusion Probabilistic Model (DDPM), the neural network  $\epsilon_\theta(x_t, t)$  receives a noisy sample  $x_t$  and a timestep  $t$  and learns to predict the noise  $\epsilon$  that produced  $x_t$  from the clean image  $x_0$ . A *conditional* DDPM extends this parametrization to

$$\epsilon_\theta(x_t, t, c),$$

where  $c$  is some conditioning signal, in our case it is the class label turned into a text embedding. It is such that the denoising process is guided by it.

#### 3.1 FiLM Conditioning

Feature-wise Linear Modulation (FiLM) [3] introduces conditioning by modulating the intermediate feature maps of the noise prediction (in our case it is implemented as a U-Net) using an affine transformation derived from the conditioning embedding. Given an intermediate activation  $h$  and a conditioning vector  $c$ , a FiLM layer produces channel-wise scale and shift parameters  $\gamma(c)$  and  $\beta(c)$ , and applies the transformation:

$$\text{FiLM}(h \mid c) = \gamma(c) \odot h + \beta(c),$$

where  $\odot$  denotes channel-wise multiplication. This modulation is applied after a the normalization layer, which in our case it is as GroupNorm (GN):

$$h' = \gamma(c) \odot \text{GN}(h) + \beta(c).$$

FiLM allows the conditioning signal to influence the forward and reverse processes *at every layer* of the U-Net, making the model significantly more expressive than simple concatenation of inputs or adding the embedding only at the input level. While the diffusion process itself (the forward noising and reverse denoising equations) does not change, the neural network that parameterizes the reverse distribution becomes condition-dependent through FiLM.

#### 3.2 Effect on Training

The DDPM training objective remains the “simple” noise prediction loss:

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{x_0, t, \epsilon, c} \left[ \|\epsilon - \epsilon_\theta(x_t, t, c)\|^2 \right],$$

but the U-Net now incorporates FiLM in each residual block so that the predicted noise depends on the conditioning. The rest of the DDPM training algorithm is unchanged.

#### 3.3 Effect on Sampling

During sampling, the reverse process uses the conditional model  $\epsilon_\theta(x_t, t, c)$  inside each transition step. The DDPM equations for sampling are unaffected; only the network that predicts the mean of the reverse distribution depends on  $c$ .

### 3.4 Algorithms

---

#### Algorithm 3: FiLM-Conditional DDPM Training

---

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0)$  ▷ Sample clean image
- 3:  $\mathbf{c} \sim \mathbf{q}(\mathbf{c}|\mathbf{x}_0)$  ▷ Sample conditioning signal
- 4:  $\mathbf{t} \sim \text{Uniform}(\{1, \dots, T\})$  ▷ Sample timestep
- 5:  $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$  ▷ Sample Gaussian noise
- 6:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$  ▷ Forward (noising) process
- 7: Take gradient step on

$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(x_t, t, c)\|^2,$$

where  $\epsilon_{\theta}$  is a U-Net whose residual blocks apply FiLM modulation:

$$\text{GN}(h) \xrightarrow{\text{FiLM}(c)} \gamma(c) \odot \text{GN}(h) + \beta(c).$$

- 8: **until** converged
- 

---

#### Algorithm 4: FiLM-Conditional DDPM Sampling

---

- 1:  $\mathbf{c}$  given ▷ Condition for generation
  - 2:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, I)$  ▷ Initial noise sample
  - 3: **for**  $\mathbf{t} = T, \dots, 1$  **do**
  - 4:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I)$  **if**  $t > 1$  **else**  $\mathbf{z} = \mathbf{0}$
  - 5:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t, \mathbf{c}) \right) + \sqrt{\Sigma_{\theta}(\mathbf{t})} \mathbf{z}$
  - 6: **end for**
  - 7: **return**  $\mathbf{x}_0$
- 

### 3.5 Results of the PyTorch Implementation of Conditional DDPM

The original DDPM implementation in [1] defines the noise prediction model as a U-Net with generic residual blocks. In order to turn this implementation into a *text-to-image* model, we change its residual layer to a FiLM residual layer.

In the notebook attached to this work, there are more details about the implementation. Essentially, what was done was substituting the **ResBlock** class from [1] with the **FiLMResBlock** class defined in the section 4 of the notebook. Also, the training and sampling algorithms were modified to reproduced the algorithms defined in the previous section (3.4).

Below there are three samples of the images generated by the conditional DDPM.

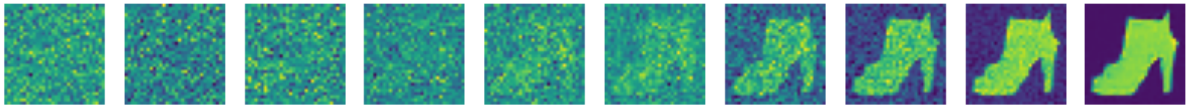


Figure 2: Heel generated by the text-to-image DDPM

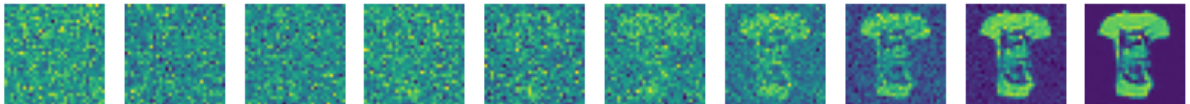


Figure 3: T-shirt generated by the text-to-image DDPM



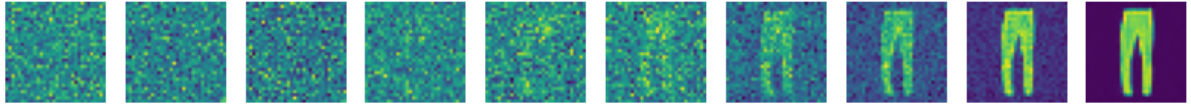


Figure 4: Pants generated by the text-to-image DDPM

## 4 Self evaluation

I did a deep study I did on the DDPM paper [2] and the resources that explain its Evidence Lower Bound, which resulted in my presentation of this work's first section to the Deep Learning class' students this semester along with my own implementation of [1].

Motivated by that presentation, I was not satisfied with simply reproducing [1]. I wanted to go further and implement a text-to-image generative model. This way, I studied about FiLM with [3] and was able to modify [1] to become a text-to-image generative model, which I consider one of my greatest achievements in deep learning up until today, given that it was an original idea.

This way, even though I think there is still room to improve the present work so that the trained text-to-image generated more diverse images from each class, I believe this work is a success.

## References

- [1] Nicholas DiSalvo. Diffusion model from scratch in pytorch (ddpm). <https://towardsdatascience.com/diffusion-model-from-scratch-in-pytorch-ddpm-9d9760528946/>, Jul 2024. Towards Data Science article, accessed 2025-12-11.
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *CoRR*, abs/2006.11239, 2020.
- [3] Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer, 2017.
- [4] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR*, abs/1503.03585, 2015.
- [5] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics, 2015.
- [6] Jake Tae. A step up with variational autoencoders. <https://jaketae.github.io/study/vae/>, 2019. Accessed: 2025-12-11.
- [7] Jake Tae. From elbo to ddpm. <https://jaketae.github.io/study/elbo/>, 2021. Accessed: 2025-12-11.
- [8] YouTube. Neural networks [7.8] : Deep learning – variational bound. <https://www.youtube.com/watch?v=pStDscJh2Wo>, n.d. Accessed: 2025-12-11.