

Reproduction Of The Denoising Diffusion Probabilistic Model

AI6103 Project Report

Chen Zhitong (G2204728D),¹ Fu Ziming (G2205212F),²
Luo Yiyang (G2204807B),³ Wang Xinyue (G2204921G),⁴

s220166@e.ntu.edu.sg,¹ zfu009@e.ntu.edu.sg,² luoy0043@e.ntu.edu.sg,³ wang1913@e.ntu.edu.sg,⁴

Abstract

In this Project, we reproduce the model from the paper *Denoising Diffusion Probabilistic Models*, which presents an innovative image denoising approach employing diffusion processes. The method iteratively denoises images, leveraging a probabilistic model to estimate noise distribution and enable effective noise reduction. We implement the DDPM and train the model for the task of human face image synthesis. The project code is available at <https://github.com/MRTater/AI6103-Project>.

Introduction

THE *Denoising Diffusion Probabilistic Models* (DDPM) represents a significant advancement in generative modeling, introducing an approach based on diffusion processes for iterative denoising. This method has demonstrated state-of-the-art (SOTA) performance in various image processing tasks, including image synthesis, super-resolution, and inpainting, making it a versatile tool for diverse applications. Additionally, DDPM has inspired the development of other diffusion-based models, such as Stable Diffusion, further expanding its impact on the field.

In this Project, we focus on reproducing the DDPM based on our understanding of the original paper. We adapt the network architecture to suit our limited computational resources without compromising performance. Furthermore, we incorporate strategies such as learning rate scheduling, skip connections, better activations, etc. to enhance model performance. Our implementation targets the task of human face image synthesis, utilizing the *Synthetic Faces Dataset* (Beniaguev 2022) to train the model.

Background

Over the past few years, deep generative models have gained popularity due to their ability to learn complex data distributions, generate high-quality samples, and perform a wide range of tasks, such as image synthesis, inpainting, and translation. Examples of deep generative models include Variational Autoencoders (VAEs), Generative Adversarial Networks (GANs), and autoregressive models like Pixel-CNN.

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Despite their success, these models face challenges. GANs struggle with generating high-quality samples and avoiding mode collapse, while autoregressive models are computationally expensive and time-consuming. Consequently, researchers have sought alternative approaches to address these limitations.

In this context, denoising score matching has garnered attention, leading to the development of DDPM. This alternative generative model addresses some of the shortcomings of VAEs and GANs by using a diffusion process to learn a noise model within a continuous-time framework.

DDPMs have emerged as a promising alternative to traditional generative models due to their high-quality training samples and stability. The DDPM approach simulates a reverse diffusion process, transforming a simple noise distribution into a complex data distribution through a series of denoising steps. A neural network describes the denoising process, learning to predict the noise added at each step of the diffusion process.

Additionally, the diffusion process used in DDPM has applications in other domains, such as optimization and machine learning, leading to new algorithms and techniques. This highlights the significant impact DDPM has had on the machine learning field and its potential to drive further advancements in the future.

Dataset

We chose the Synthetic Faces High Quality (SFHQ) Part 1 dataset, as shown in Figure 1, for training our model. This dataset comprises 89,785 high-quality 1024×1024 curated face images created by "bringing to life" various artworks (paintings, drawings, 3D models) using a process similar to that described in a short Twitter thread. This process involves encoding the images into StyleGAN2 latent space and performing a minor manipulation that converts each image into a photo-realistic representation.

Data Preprocessing

Data preprocessing involves transforming raw data into a suitable range or scale that facilitates model learning. In our experiment, we performed several preprocessing steps on the data. Initially, we randomly selected 20,000 images from the dataset and resized them to 256×256 for our training



Figure 1: Preview of the Synthetic Faces dataset.

set. To further accelerate training, given limited computational resources, we resized the images to 64×64 and applied a random horizontal crop as data augmentation. The image data, which initially consisted of integers ranging from 0 to 255, was linearly scaled to the range of $[-1, 1]$ to ensure compatibility with the model.

Diffusion Forward Process

In the context of denoising diffusion probabilistic models (DDPM), the forward diffusion process, as described by Equation (1), is responsible for incrementally adding noise to the original image x_0 following a precise strategy at each timestep. This procedure aims to mimic the diffusion of the image across T timesteps, generating a sequence of increasingly noisy images that ultimately converge to Gaussian noise.

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1}) \quad (1)$$

The method of sampling noises is described by Equation (2). It is the conditional Gaussian noise with a mean that depends on the previous image and fixed variance, i.e., $(\beta_t I)$ in our implementation. The sequence of betas is called the noise schedule, which describes how much noise is added at each timestep.

$$q(x_t | x_{t-1}) = N\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right) \quad (2)$$

Since the sum of Gaussian distributions is still a Gaussian distribution, as proven by Equation (3), given the set of betas, we can precompute various terms that will be employed in the forward diffusion process, such as the cumulative product of alphas ($\alpha_t^{cumprod}$), and 1 minus the cumulative product of alphas ($1 - \alpha_t^{cumprod}$). These terms are utilized to control the variance and the noise introduced to the image at each timestep.

$$\begin{aligned} q(x_t | x_{t-1}) &= N\left(x_t, \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right) \\ &= \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon \\ &= \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon \\ &= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon \\ &= \sqrt{\alpha_t \alpha_{t-1} \alpha_{t-2}} x_{t-3} + \sqrt{1 - \alpha_t \alpha_{t-1} \alpha_{t-2}} \epsilon \\ &= \sqrt{\alpha_t \alpha_{t-1} \dots \alpha_1 \alpha_0} x_0 + \sqrt{1 - \alpha_t \alpha_{t-1} \dots \alpha_1 \alpha_0} \epsilon \\ q(x_t | x_0) &= N(x_t; \sqrt{\alpha_t} x_0, (1 - \bar{\alpha}_t) I) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \end{aligned} \quad (3)$$

Beta Scheduling

The choice of the beta schedule has a significant impact on the performance of the DDPM. We propose two beta schedules: linear and cosine.

The linear beta schedule is defined as Equation (4):

$$\beta_t = \text{start} + \frac{t(\text{end} - \text{start})}{T}, \quad (4)$$

where T is the total number of timesteps, and start and end are the initial and final values of the beta schedule.

The cosine beta schedule is defined as Equation(5):

$$\begin{aligned} \bar{\alpha}_t &= \frac{f(t)}{f(0)}, \quad f(t) = \cos\left(\frac{t/T + s}{1+s} \cdot \frac{\pi}{2}\right)^2, \\ \beta_t &= 1 - \frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}} \end{aligned} \quad (5)$$

The cosine beta schedule features a linear drop-off of $\bar{\alpha}_t$ in the middle while maintaining stable noise levels near $t = 0$ and $t = T$ as shown in Figure 2. This provides a more gradual noise addition compared to the linear schedule. A small offset s is introduced to avoid overly small β_t values near $t = 0$. The schedule is defined using the cosine function, ensuring smooth noise transitions. By clipping beta values to a maximum limit, it prevents excessive noise and improves denoising performance, especially for lower resolution images.

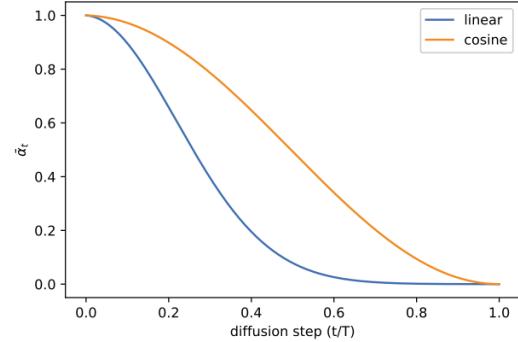


Figure 2: Linear and cosine beta schedules

Summary of the Forward Process

The forward diffusion process for a specific timestep t can be described as follows:

1. Sample Gaussian noise $\epsilon \sim \mathcal{N}(0, 1)$ with the same dimensions as the input image x_0 .
2. Compute the noisy image x_t as Equation(4):

$$x_t = \sqrt{\alpha_t^{cumprod}} x_0 + \sqrt{1 - \alpha_t^{cumprod}} \epsilon, \quad (6)$$

where $\alpha_t^{cumprod}$ and $1 - \alpha_t^{cumprod}$ govern the mean and variance of the noise introduced to the image, respectively. The noisy image is then clamped to the range $[-1, 1]$ to ensure valid pixel values.

The forward diffusion process is applied iteratively for all T timesteps as shown in Figure 3, producing a sequence of increasingly noisy images. During training, the DDPM model learns to invert this process by predicting the noise added at each timestep, enabling it to synthesize new images by denoising Gaussian noise in the later part.



Figure 3: A preview of forward diffusion process by time steps T

Parametrized Backward Process

The backward process in denoising diffusion probabilistic models (DDPM) is responsible for reconstruct the image x_0 back from the noisy image at the final timestep x_T . This process, parametrized by θ , is described by Equation (9):

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

$x_{t-1} \approx x_t - \text{noise}$

The parametrized backward process is learned by the DDPM using a neural network, in our case, a simplified variant of the U-Net, along with sinusoidal position embeddings for the information of time steps. Notice that we only learn the mean, since the variance is fixes as mentioned.

U-Net

The U-Net serves as the core of the parametrized backward process in denoising diffusion probabilistic models. Its ability to capture both local and global information makes it ideal for denoising tasks.

The architecture maintains a large receptive field and preserves spatial information throughout the processing, enabling detailed image generation.

We use a simplified version of U-Net as shown in Figure 4, named SimpleUnet, consisting of an initial convolutional layer, a series of downsampling and upsampling blocks, and an output convolutional layer. The blocks incorporate time embedding and consist of convolutional layers, batch normalization layers, and activation functions.

The SimpleUnet is initialized with customizable settings like activation function (ReLU or SiLU), self-attention, and skip connections. Skip connections link corresponding downsampling and upsampling layers, propagating detailed spatial information to reconstruct denoised images more accurately. Additionally, there are skip connections within the blocks themselves, allowing for better information flow. The architecture has 5 downsampling and upsampling layers, with channels reaching up to 1024 in the middle layers. This structure captures hierarchical input image representations, including both high-level and low-level features.

Sinusoidal Position Embeddings

Incorporating temporal information is crucial in denoising diffusion probabilistic models (DDPM). To achieve this, we employ sinusoidal position embeddings, which provide a continuous and differentiable representation of time.

Sinusoidal position embeddings are computed using the following formulae:

$$\begin{aligned} P(k, 2i) &= \sin\left(\frac{k}{n^{2i/d}}\right) \\ P(k, 2i+1) &= \cos\left(\frac{k}{n^{2i/d}}\right) \end{aligned} \quad (8)$$

Where P represents the position embedding matrix, k is the time step, n is a constant (typically 10000), d is the dimensionality of the embeddings, and i is an index that ranges from 0 to $\frac{d}{2} - 1$. These equations generate sine and cosine functions with different frequencies, effectively encoding temporal information within the embeddings.

The resulting sinusoidal embeddings as shown in Figure 5, can be used within the DDPM, effectively encoding temporal information for the denoising process. The sinusoidal nature of these embeddings ensures that they can be easily combined with other representations, facilitating the learning process for the model.

Loss Function

In the DDPM, the loss function is used to measure the discrepancy between the actual noise and the predicted noise at each timestep during training. Here, we present two different loss functions: L1 and L2.

The L1 loss function is defined as follows:

$$\mathcal{L}1(\theta) = \sum_{t=0}^{T-1} \|\epsilon_t - \text{model}(x_t, t)\|_1, \quad (9)$$

The L2 loss function, also known as the Mean Squared Error (MSE) loss, can lead to smoother denoising results as it tends to penalize large deviations more heavily. It is defined similarly to the L1 loss:

$$\mathcal{L}2(\theta) = \sum_{t=0}^{T-1} \|\epsilon_t - \text{model}(x_t, t)\|_2^2, \quad (10)$$

where ϵ_t is the true noise, x_t is the noisy image, t is the timestep, and $\text{model}(x_t, t)$ is the predicted noise at timestep t .

During training, either the L1 or L2 loss can be used to optimize the model, depending on the desired denoising characteristics.

Summary of the Backward Process

The parametrized backward process in DDPM sythesize the original image x_0 from the noisy image x_T . The process involves a simplified U-Net architecture and sinusoidal position embeddings for encoding temporal information.

1. Perform batch training with noisy images x_T and corresponding time steps.

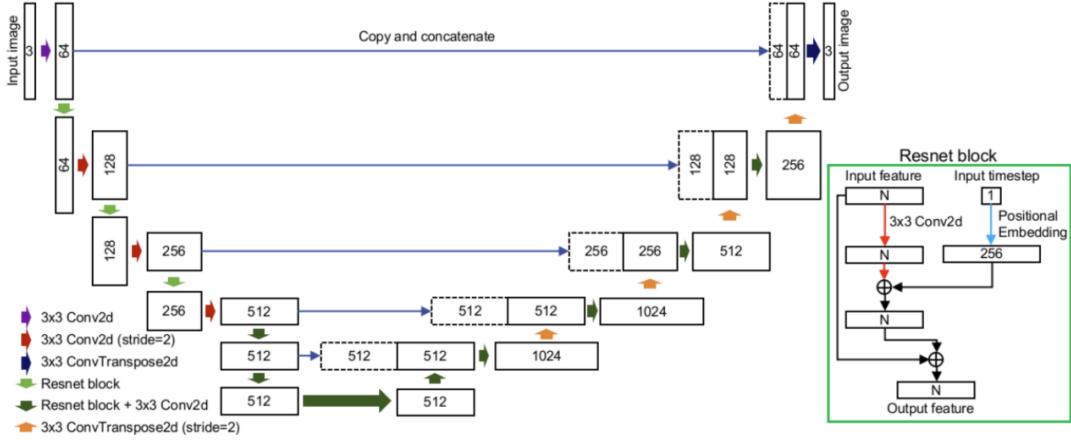


Figure 4: The SimpleUnet

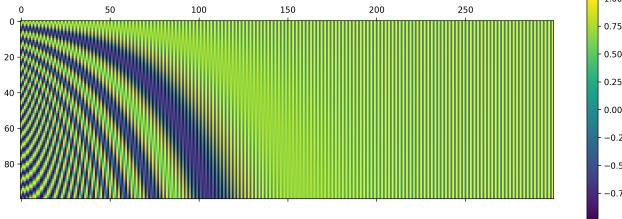


Figure 5: Position embeddings, $T=300$

2. Encode time steps using sinusoidal position embeddings to generate a position embedding matrix P .
3. Process input images using the SimpleUnet and feed in with the sinusoidal position embeddings to incorporate temporal information in the denoising process.
4. Iterate through the training process, adjusting model parameters to minimize the loss function and refining the denoising capability.
5. By following these steps, the parametrized backward process effectively synthesizes the original images from the noisy counterparts, improving the performance of the denoising diffusion probabilistic model.

Metric

Considering that our project focuses on the task of image synthesis, we choose the *Fréchet Inception Distance* (FID) as the metric for our task. FID is a widely-used metric for evaluating the quality of generated images. First introduced by Martin Heusel et al. in 2017, FID measures the similarity between distributions of generated and real images in the feature space of a pre-trained Inception-v3 neural network.

The principle behind FID for calculating image distance can be summarized as follows: given that a random variable follows a Gaussian distribution, the distribution can be determined by its mean and variance. Two distributions are considered identical if their means and variances are the same. In the case of multidimensional distributions, the covariance

matrix is used to measure the correlation between dimensions. Thus, FID employs the mean and covariance matrix to calculate the distance between two distributions.

The FID score can be calculated using the following steps:

1. Extract feature vectors from the real and generated images using a pre-trained Inception-v3 neural network.
2. Calculate the mean and covariance of the feature vectors separately for the real and generated images.
3. Compute the Fréchet distance between the two distributions.

The mathematical formula for FID is expressed as Equation 11. Here, μ_x and μ_y represent the feature-wise means of the real and generated images, and Σ_x and Σ_y denote the covariance matrices of the real and generated feature vectors.

$$\text{FID}(x, g) = \|\mu_x - \mu_g\|_2^2 + \text{Tr} \left(\Sigma_x + \Sigma_g - 2 (\Sigma_x \Sigma_g)^{\frac{1}{2}} \right) \quad (11)$$

By improving the FID score, our model aims to generate images that are more similar to the real images in terms of their feature distributions, resulting in higher quality synthesized images.

Experiments

For all experiments, we set $T = 300$. We chose forward process variances as constants, increasing linearly from $\beta_1 = 10^{-4}$ to $\beta_T = 0.02$, consistent with the original paper's settings. The size of the sampled images is 64×64 .

We employed a U-Net backbone as the primary network, illustrated in Figure 4. Parameters are shared across time, specified to the network using sinusoidal position embeddings from the Transformer architecture.

In addition, we implemented a cosine β schedule for further investigation of the diffusion model.

Due to computational constraints, we omitted the self-attention layer as it was deemed too computationally expensive.

Epochs	FID
200	233.36
300	220.25
450	198.03

Table 1: SFHQ Results, measured with FID

Sampling Process

The sampling process in our experiments involves reversing the diffusion process. By starting with pure noise as input, the model iteratively removes noise until a realistic image is generated. In the context of the diffusion model, the sampling process is essentially equivalent to denoising.

We sampled 64×64 images using our diffusion model in our experiments. Our implementation achieved a generation speed of 25.7 seconds per batch, with each batch containing 128 images. Table 1 presents the FID scores on the Synthetic Faces High Quality (SFHQ) part 1 dataset, where the FID score decreases with the number of training epochs. Figure 6 illustrates the loss decreasing with the number of training epochs, indicating that our model's ability to generate high-quality images improves as training progresses. Some generated results are showcased in Figure 7.

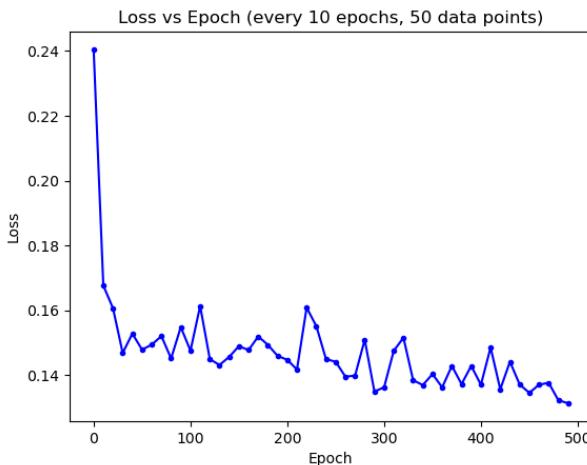


Figure 6: Loss vs Epoch, T=300

Beta Scheduler

For our experiments, we employed both linear and cosine β schedules to investigate their impact on the model's performance.

Our results revealed that the cosine schedule outperformed the linear schedule in terms of sample quality and training stability, as illustrated in Figure 8. Specifically, the cosine schedule yielded lower FID scores compared to the linear schedule.

In conclusion, we determined that the cosine β schedule is a more effective approach for training our diffusion model than the linear schedule. Its smoother and more gradual transition between values facilitates enhanced performance and stability during the training process.



Figure 7: SFHQ samples, FID=198.03

Activations	FID
SILU	233.36
RELU	235.13

Table 2: Comparison between activations, measured by FID, epoch=200

Activation Function

We conducted experiments to investigate the impact of two different activation functions on the performance of our model: SiLU and ReLU as shown in Table 2.

Our results indicate that the choice of activation function does not significantly influence the model's performance. Specifically, we observed that employing SiLU as the activation function resulted in a slight improvement in the model's accuracy, while ReLU demonstrated a marginally better convergence speed.

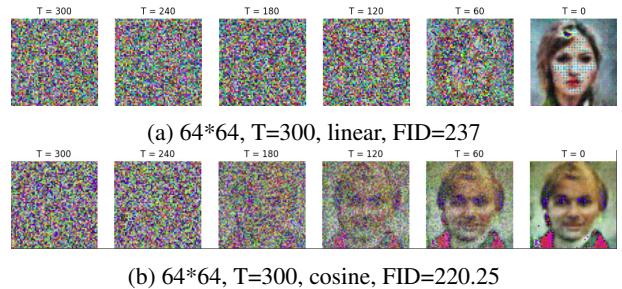


Figure 8: Comparison of synthesized images via backward process with different beta scheduler: (a) linear; (b) cosine.

Result Analysis

Our model was only able to generate images with an FID score of around 200, which is less realistic than the results reported in the original paper. Moreover, some noise and color jittering can be observed in the generated images, indicating that the model's denoising ability could be further improved. There are several factors that might have negatively impacted our model's performance.

One such factor is that, in contrast to the original paper, which employed CelebA-HQ models with 114 million parameters (Ho, Jain, and Abbeel 2020), we utilized a much shallower network consisting of only 62 million parameters. Shallower networks, with fewer layers and parameters, can constrain the model's learning ability. In the context of our diffusion model, this means that the model may exhibit inferior performance in learning noise distributions, thereby generating lower-quality images. Self-attention layers facilitate the model's focus on crucial information and relationships, which can enhance the accuracy of the diffusion process. Eliminating self-attention layers may impair the model's capacity to capture features of noise distribution and the relationships within different feature maps, ultimately diminishing the model's denoising ability.

Furthermore, we opted for a smaller value of $T = 300$, as opposed to the original paper's $T = 1000$ (Ho, Jain, and Abbeel 2020). Although this choice facilitated training, it adversely affected the generation process. Employing a smaller T value can impact the diffusion results in several ways:

- **Reduced Quality.** Using a smaller T value truncates the diffusion process, leading to lower-quality image or data generation. The resulting images or data may contain artifacts or appear blurry, which might be undesirable in certain applications.
- **Limited Detail.** A smaller T value also restricts the amount of detail in the generated images or data. Consequently, the generated images or data may lack fine-grained detail compared to those produced with a larger T value.
- **Limited Diversity.** The diversity of the generated images or data might also be compromised when using a smaller T value. Images or data generated with a smaller T value may be less diverse and exhibit a narrower range of features or characteristics.

Additionally, we selected a relatively modest epoch setting, which could have potentially affected the model's performance. Initially, we assumed that the epoch settings would not significantly influence the experiment's outcome beyond a certain timestep (200 in our experiment), considering that the loss remained relatively stable. However, upon examining the results, we observed that the FID of the generated images varied considerably compared to the change in loss, suggesting that longer epochs might be necessary to optimize the model's performance. This could be attributed to the limited learning capacity of our model.

In summary, our model's ability to generate realistic images was limited, with generated images only attaining an

FID score of approximately 200 and exhibiting noise and color jittering. Factors that adversely affected the model's performance included using a shallower network with fewer parameters than the original paper, a smaller T value, and a relatively modest epoch setting.

Conclusion

In summary, our reimplementations of the DDPM has generated image samples comparable to those in the original paper, demonstrating its effectiveness as a powerful generative model. Although our implementation was adapted to accommodate limited computational resources, it has underscored the significant potential of the DDPM for generative tasks. Furthermore, recent advances such as Stable Diffusion have bolstered the performance of diffusion-based models in this domain, highlighting the potential of incorporating DDPM into various generative models and machine learning systems. As research into DDPM's capabilities progresses, we expect to uncover novel applications and deepen our understanding of this promising approach, paving the way for further advancements in the field.

References

- Beniaguev, D. 2022. Synthetic Faces High Quality (SFHQ) part 1.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33: 6840–6851.