

# DIFFUSION MODELS

Mollaev D.E.

May 9, 2023

# PROJECT DESCRIPTION

- Datasets - MNIST, CIFAR10
- Model - DDPM (Denoising Diffusion Probabilistics Models) with Unet
- Metrics - FID, INCEPTION SCORE

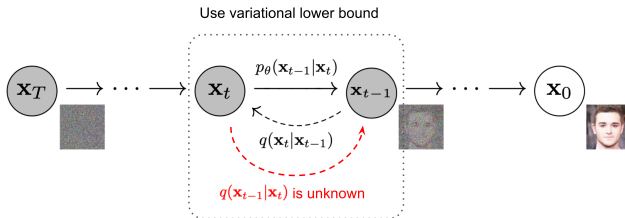
# DDPM

## FORWARD DIFFUSION PROCESS

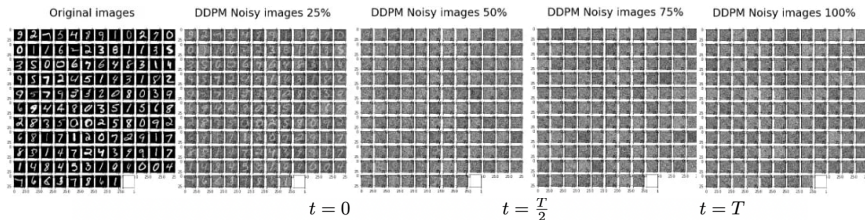
Let  $x_0 \sim q(x_0)$ , then we add small amount Gaussian noise to the sample in  $T$  steps. And we get sequence of noisy samples  $x_1, \dots, x_T$   
Then conditional probability

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} * x_{t-1}, b_t \mathcal{I})$$

Eventually when  $T \rightarrow \infty$ ,  $x_T$  is equivalent to an isotropic Gaussian distribution

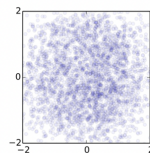
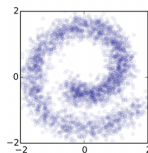
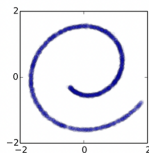


# VISUALIZATION OF FORWARD PROCESS



The forward trajectory

$$q(\mathbf{x}_{0:T})$$



# DDPM

## REVERSE DIFFUSION PROCESS

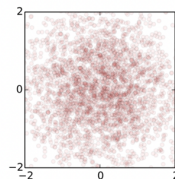
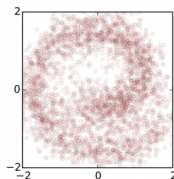
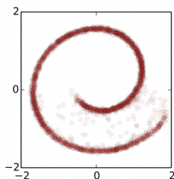
If we can reverse the forward process and sample from  $q(x_{t-1}|x_t)$ , we will be able to recreate the true sample from a Gaussian noise input,  $x_T \sim \mathcal{N}(0, \mathcal{I})$ .

But we cannot easily estimate  $q(x_{t-1}|x_t)$  because it needs to use the entire dataset and therefore we need to learn a model(Unet)  $p_\theta$  to approximate these conditional probabilities in order to run the reverse diffusion process.

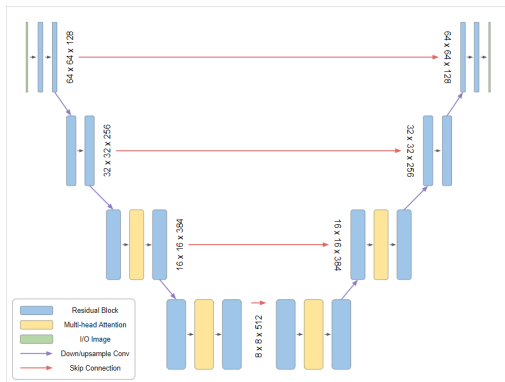
$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

The reverse trajectory

$p_\theta(\mathbf{x}_{0:T})$



[Link for visualization](#)



---

**Algorithm 1** Training

---

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
        $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$ 
6: until converged
```

---

---

**Algorithm 2** Sampling

---

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

---

**FID(Frechlet Inception Distance)** is a performance metric that calculates the distance between the feature vectors of real images and the feature vectors of generate images

$$d^2((\mathbf{m}, \mathbf{C}), (\mathbf{m}_w, \mathbf{C}_w)) = \|\mathbf{m} - \mathbf{m}_w\|_2^2 + \text{Tr}(\mathbf{C} + \mathbf{C}_w - 2(\mathbf{C}\mathbf{C}_w)^{1/2})$$

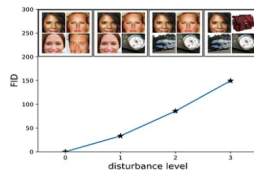
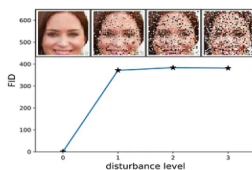
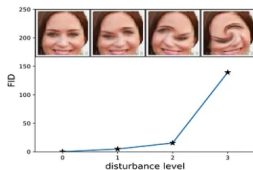
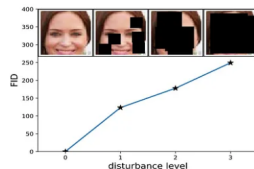
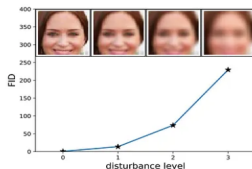
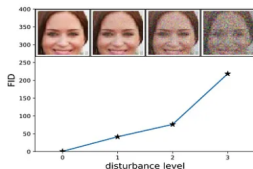
How to calculate FID?

- 1 Use the Inception V2 pre-trained model to extract the feature vectors of real images and generated images by the generator
- 2 Calculate the feature-wise mean of the feature vectors generated in step 1
- 3 Generate the covariance matrices of the feature vectors —  $\mathbf{C}, \mathbf{C}_w$
- 4 Calculate trace
- 5 Calculate the squared difference of the mean vectors calculated in step 2
- 6 Finally, add the output of step 4 and step 5



# METRICS

## FID



**The Inception Score (IS)** is an objective performance metric, used to evaluate the quality of generated images or synthetic images. It measures how realistic and diverse the output images are.

It measures two things:

- **Diversity** (Variety) — How diverse the generated images are — The entropy of the overall distribution should be high.
- **Quality** (Goodness) — How good the generated images are — Low entropy with high predictability is required.

$$\text{IS}(G) = \exp \left( \mathbb{E}_{\mathbf{x} \sim p_a} D_{KL}(p(y|\mathbf{x}) \parallel p(y)) \right)$$

- **Conditional Probability Distribution** —  $p(y|x)$ . It should be highly predictable and with low entropy. Here  $y$  is the set of labels and  $x$  is the image.
- **Marginal Probability Distribution** —  $p(y)$

$$\int_z p(y|x = G(z))dz$$

Here,  $G(z)$  is the generated image by the generator model when provided with a latent vector. If the data distribution for  $y$  is uniform with high entropy, then the synthetic images will be diverse.

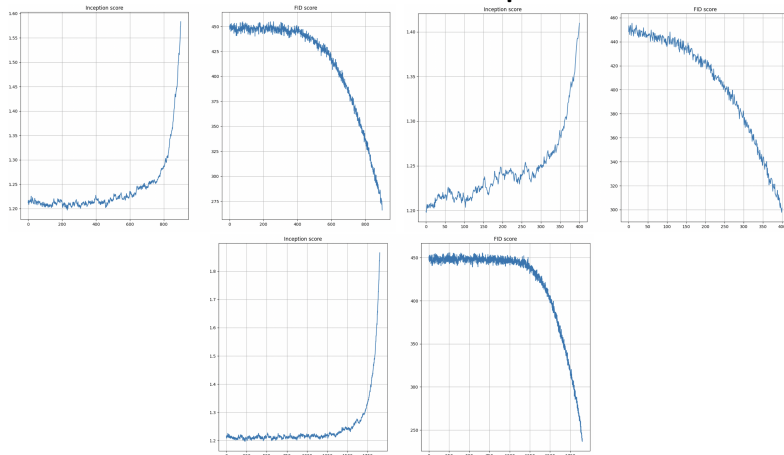
How to calculate IS?

- 1 Pass the generated images through the Inception model to get the conditional label distribution  $p(y|x)$
- 2 Calculate the marginal probability distribution  $p(y)$
- 3 Calculate the KL Divergence between  $p(y)$  and  $p(y|x)$
- 4 Calculate the sum over classes and take the average of outputs over images
- 5 Finally, take the exponential of the averaged value.

# RESULTS

I learned three generator with 500, 1000, 2000 steps on datasets CIFAR10 and show plot of metrics.

**1000, 500, 2000 steps:**



- Link for visualization(1000 steps on CIFAR10)
- Link for visualization(500 steps on CIFAR10)
- Link for visualization(2000 steps on CIFAR10)

## Links:

- Denoising Diffusion Probabilistic Models
- Lil'Log What are Diffusion Models?
- Medium: Generating images with DDPMs: A PyTorch Implementation
- Medium: A Very Short Introduction to Inception Score(IS)
- Medium: A Very Short Introduction to Frechlet Inception Distance(FID)
- Special course on Mechmath MSU: Introduction to Machine and Deep Learning Theory