

CL2: Text to Image Generation using AI

Demostrator: Yisong Shen

student: Haixiang Huang

CID: 02470313

Part A: Proof of DDPM

1. Concept of DDPM

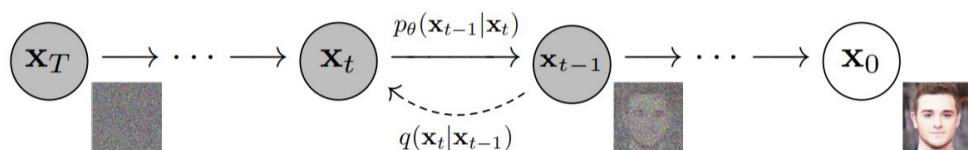


Figure 2: The directed graphical model considered in this work.

DDPM is a parameterized Markov chains

training using Variational inference to learn the transitions of this chain to reverse a diffusion process.

[learn $p_\theta(x_{t-1}|x_t)$ to denoise and generate image]

2. Diffusion Process.

2.1. One step Diffusion

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

$$x_t = \sqrt{1-\beta_t} x_{t-1} + \beta_t \epsilon, \text{ where } \epsilon \sim N(0, I)$$

N.B. use β_t to control the speed of noise adding.

2.2. multiple steps diffusion process

We have: $q(x_1:T | z_0) = \prod_{t=1}^T q(x_t | x_{t-1})$ [Markov chain]

We use: $\partial t = 1 - \beta_t$ and $\bar{\partial} t = \prod_{i=1}^t \partial i$

so, $x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon$

$$x_1 = \sqrt{1-\beta_1} z_0 + \sqrt{\beta_1} \epsilon \quad \text{substitute}$$

$$x_2 = \sqrt{1-\beta_2} x_1 + \sqrt{\beta_2} \epsilon$$

$$x_2 = \sqrt{1-\beta_1} \sqrt{1-\beta_2} z_0 + \sqrt{1-\beta_2} \sqrt{\beta_1} \epsilon + \sqrt{\beta_2} \epsilon$$

$$= \sqrt{1-\beta_1} \sqrt{1-\beta_2} z_0 + \sqrt{1-(1-\beta_2)(1-\beta_1)} \epsilon$$

⋮

$$x_t = \underbrace{\sqrt{1-\beta_1} \cdots \sqrt{1-\beta_t} z_0}_{\sqrt{\bar{\partial} t}} + \underbrace{\sqrt{1-(1-\beta_1) \cdots (1-\beta_t)}}_{\sqrt{1-\bar{\partial} t}} \epsilon.$$

$$\therefore x_t = \bar{x}_t + \bar{\epsilon}$$

$$q(x_t|x_0) = N(x_t; \bar{x}_t, \bar{I})$$

3. Reverse Process

Denoise

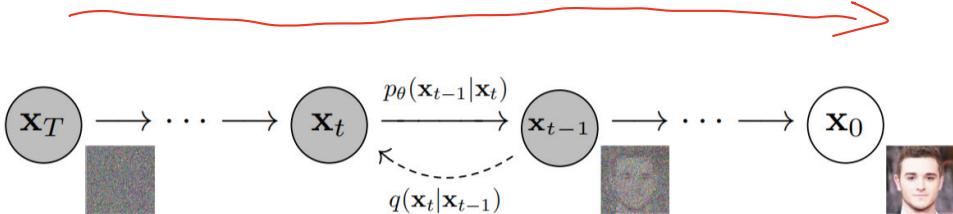


Figure 2: The directed graphical model considered in this work.

- ① We assume reverse process is a **Markov chain** and transitions of distribution is a **Gaussian**



- ② Parameters estimation problem

To learn $p_\theta(x_{t-1}|x_t)$ using neural Network.

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

Note: In DDPM paper, $\Sigma_\theta(x_t, t)$ is fixed for simplity.

~~③~~

- We use **Variational inference** to approximate

$P\theta(x_{t-1} | x_t)$ to known posterior of diffusion process.

$$q(x_{t-1} | x_t, x_0) \left[q(x_t | x_{t-1}, x_0) = N(x_{t-1}; M(x_t, x_0), \bar{P}_t I) \right]$$

[proof below]

4. posterior of diffusion process.

$$\begin{aligned} q(x_{t-1} | x_t, x_0) &= \frac{q(x_{t-1}, x_t, x_0)}{q(x_t, x_0)} \\ &= q(x_t | x_{t-1}, x_0) \cdot \frac{q(x_{t-1} | x_0) q(x_0)}{q(x_t | x_0) q(x_0)} \\ &= q(x_t | x_{t-1}, x_0) \cdot \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} \\ &= q(x_t | x_{t-1}) \cdot \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} \quad [\text{Markov chain}] \end{aligned}$$

Due to $q(x_t | x_{t-1})$, $q(x_{t-1} | x_0)$, $q(x_t | x_0)$ all are Gaussian

$q(x_t | x_{t-1}) \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)}$ is also Gaussian.

$$\begin{aligned} &\propto \exp\left(-\frac{1}{2}\left(\frac{(x_t - \bar{x}_t x_{t-1})^2}{\beta_t} + \frac{(x_{t-1} - \bar{x}_{t-1} x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \bar{x}_t x_0)^2}{1 - \bar{\alpha}_t}\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(\underbrace{\left(\frac{\bar{\alpha}_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}\right)x_{t-1}^2}_{a} - \underbrace{\left(\frac{2\bar{\beta}_t}{\beta_t}x_t + \frac{2\bar{\beta}_{t-1}}{1 - \bar{\alpha}_{t-1}}x_0\right)x_{t-1}}_{b} + C\right)\right) \end{aligned}$$

We know that $ax^2 + bx + c = 0$.

$$\Rightarrow a(x - \frac{b}{2a})^2 + \frac{(4ac - b^2)}{4a} = 0$$

P.d.f of Gaussian is $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\Rightarrow \mu = \frac{b}{2a} \quad \sigma^2 = \frac{1}{a}$$

$$\Rightarrow \sigma^2 = 1 / \left(\frac{\partial}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t+1}} \right) = \underbrace{\frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_t}}_{C} \beta_t.$$

$$\begin{aligned} \hat{\mu}_t(x_t, x_0) &= \left(\frac{\sqrt{\bar{\alpha}_t}}{\beta_t} x_t + \frac{\sqrt{1 - \bar{\alpha}_t}}{1 - \bar{\alpha}_{t+1}} x_0 \right) / \left(\frac{\partial}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t+1}} \right) \\ &= \left(\frac{\sqrt{\bar{\alpha}_t}}{\beta_t} x_t + \frac{\sqrt{1 - \bar{\alpha}_t}}{1 - \bar{\alpha}_{t+1}} x_0 \right) \cdot \frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_t} \beta_t \\ &= \frac{\sqrt{\bar{\alpha}_t} (1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{1 - \bar{\alpha}_t} \beta_t}{1 - \bar{\alpha}_t} x_0. \end{aligned}$$

[multiple step forward] \uparrow
 $\hat{\mu}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_t)$

$$\Rightarrow \hat{\mu}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right)$$

$$\Rightarrow q(x_{t+1} | x_t, x_0) = N(x_{t+1}; \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right), \frac{1 - \bar{\alpha}_{t+1}}{1 - \bar{\alpha}_t} \beta_t)$$

5. Variational inference

$$\begin{aligned}
 & \text{KL} (q(x_{1:T}|x_0) || p_\theta(x_{1:T}|x_0)) \\
 &= \mathbb{E}_{x_{1:T} \sim q} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})/p_\theta(x_0)} \right] \\
 &= \underbrace{\mathbb{E}_q \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right]}_{-\text{ELBO}} + \log p_\theta(x_0)
 \end{aligned}$$

$$\Rightarrow \text{KL} = -\text{ELBO} + \log p(x) \geq 0.$$

$$\therefore \log p_\theta(x_0) \geq \text{ELBO}$$

$$-\log p_\theta(x_0) \leq -\text{ELBO}$$

$$\underbrace{-\mathbb{E}_{q(x_0)} \log p_\theta(x_0)}_{\text{cross entropy loss.}} \leq \underbrace{\mathbb{E}_{q(x_{1:T})} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right]}_{\text{LVB.}}$$

MIN LOSS \Rightarrow MIN - ELBO. \Rightarrow MIN LVB

$$\begin{aligned}
 \text{LVB} &= \mathbb{E}_{q(x_{1:T})} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \\
 &= \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(x_t|x_{t-1})}{p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)} \right]
 \end{aligned}$$

$$\begin{aligned}
 q(x_t|x_{t-1}) &\stackrel{\text{softmax}}{=} \frac{q(x_t|x_{t-1}, x_0)}{\sum_{x_t} q(x_t|x_{t-1}, x_0)} \\
 &= \frac{q(x_t, x_{t-1}, x_0)}{q(x_{t-1}, x_0)} \\
 &= \frac{q(x_{t-1}|x_t, x_0) q(x_t, x_0)}{q(x_{t-1}) q(x_0)} \\
 &= \frac{q(x_{t-1}|x_t, x_0) q(x_t|x_0)}{q(x_{t-1})}
 \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=1}^T \log \frac{q(x_t | x_{t-1})}{p_\theta(x_{t-1} | x_t)} \right] \\
&= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_t | x_{t-1})}{p_\theta(x_{t-1} | x_t)} + \log \frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right] \\
&= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \cdot \frac{q(x_t | x_0)}{q(x_{t-1} | x_0)} \right) + \log \frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right] \\
&= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} + \underbrace{\sum_{t=2}^T \log \frac{q(x_t | x_0)}{q(x_{t-1} | x_0)}}_{= \log \frac{\prod q(x_t | x_0)}{\prod q(x_{t-1} | x_0)}} + \log \frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right] \\
&= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} + \log \frac{q(x_T | x_0)}{q(x_{T-1} | x_0)} + \log \frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right] \\
&= \mathbb{E}_q \left[\log \frac{q(x_T | x_0)}{p_\theta(x_T)} + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} - \log p_\theta(x_0 | x_1) \right] \\
&= D_{KL}(q(x_T | x_0) || p_\theta(x_T)) + \underbrace{\sum_{t=2}^T D_{KL}(q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t))}_{L_{T-1}} - \underbrace{\log p_\theta(x_0 | x_1)}_{L_0}.
\end{aligned}$$

L_T : all is Gaussian, is not related to optimization.

$$L_0: t=1, D_K(\underbrace{q(x_0 | x_t, x_0)}_{=1} || p_\theta(x_0 | x_1)) = \log p_\theta(x_0 | x_1)$$

It is specific L_{T-1}

So, we only care about L_{T-1}

* In this way, $\text{MIN } L_{VLB} \Rightarrow \text{MIN } D_{KL}(q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t))$

\Rightarrow approximate $P_\theta(x_{t+1}|x_t)$ to $q(x_{t+1}|x_t, x_0)$

$$\text{We have } KL(P, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

σ_p, σ_q is constant, only $\frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2}$ can contribute to KL

$$KL_t = \frac{1}{2\sigma_q^2} \left\| \tilde{m}_t(x_t, x_0) - m_\theta(x_t, t) \right\|^2$$

$\therefore \min KL \Rightarrow \min \tilde{m}_t \text{ and } m_\theta$.

$$\text{Also, } \tilde{m}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \frac{1-\bar{\alpha}_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_t)$$

x_t is known.

We only need to train neural network to predict ϵ_t

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

$$\Rightarrow KL = \frac{1}{2\sigma_q^2} \left\| \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \frac{1-\bar{\alpha}_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_t) - \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \frac{1-\bar{\alpha}_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t)) \right\|^2$$

$$= \frac{(1-\bar{\alpha}_t)^2}{2\bar{\alpha}_t(1-\bar{\alpha}_t)\sigma_q^2} \left\| \epsilon_t - \epsilon_\theta(x_t, t) \right\|^2$$

C.

$$\because x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon_t$$

$$\therefore KL = C \cdot \left\| \epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon_t, t) \right\|^2$$

In DDPM original paper, author find $C=1$ doesn't

affect the result for simplicity.

$$\therefore \text{MIN DKL}(q(x_{t+1} | x_t, x_0) || p_\theta(x_{t+1} | x_t)) \Rightarrow \text{MIN} \| \epsilon_t - \epsilon_0 \|$$

We train NN to predict ϵ_t .

Based Above, we know:

$$\text{MIN Loss} \Rightarrow \text{MIN LVB} \Rightarrow \text{MIN} \| \epsilon_t - \epsilon_0 \|$$

$$\therefore \text{Loss} = \epsilon - \epsilon_0 (\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon_t)^2$$

Finally:

We get two algorithm: Training and Sampling;

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  $\mathbf{x}_0$ 
```

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

$$p_\theta(x_{t+1} | x_t) = \mathcal{N}\left(x_{t+1}; \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right), \frac{1-\bar{\alpha}_t}{1-\bar{\alpha}_t} \beta_t\right)$$

Part B: Implementation of DDPM

1. Set up hyperparameters

1.1 Hyperparameters

```
1 # new hyperparams
2 TIMESTEPS = 1000
3 beta_0 = 0.0001
4 beta_T = 0.02
5 LEARNING_RATE = 0.001
6 # generate alphas,betas,alpha_bars using my unique name hhx to store
7 hhx_betas = torch.linspace(beta_0, beta_T, TIMESTEPS, device = DEVICE)
8 hhx_alphas = 1.0 - hhx_betas
9 hhx_alpha_bars = torch.cumprod(hhx_alphas, axis=0)
```

$\beta_0 \sim \beta_T \in [0.02, 10^{-4}]$ based on DDPM paper.

$$\begin{cases} \bar{\alpha} = 1 - \beta_t \\ \bar{\alpha} = \prod_{i=1}^t \alpha_i \end{cases}$$

2. Diffusion process.

2.1 One step diffusion process.

```
def fixed_one_step_forward(xt_minus_1, t):
    """
    fixed one step of a forward process according to one step formula in DDPM
    Args:
        xt_minus_1(torch.Tensor): The input image tensor in t-1 step (previous image)
        t (int): The current time step.

    Returns:
        xt (torch.Tensor): The output image tensor in t step after one step of t
    """
    noise = torch.randn_like(xt_minus_1, device = DEVICE) →  $\epsilon \sim N(0, I)$ 
    beta_t = hhx_betas[t]
    xt = torch.sqrt(1.0 - beta_t) * xt_minus_1 + torch.sqrt(beta_t) * noise →  $X_t = \sqrt{1-\beta_t} X_{t-1} + \beta_t \epsilon_t$ 
    return xt
```

$$= \sqrt{\bar{\alpha}_t} X_{t-1} + \sqrt{1-\bar{\alpha}_t} \epsilon_t$$

2.2. multiple steps. diffusion process.

```

def hhx_forward_diffusion(x_0, t):
    """
    Perform the forward diffusion process based on the given input clean image x_0.
    This function computes x[t] based on x[0], using the diffusion process defined
    in the paper.

    Args:
    - x_0 (torch.Tensor): The input clean image tensor. Expected to be a 4D tensor
    | | | | (batch_size, Channels, Img_size, Img_size).
    - t (torch.Tensor): The timestep tensor.

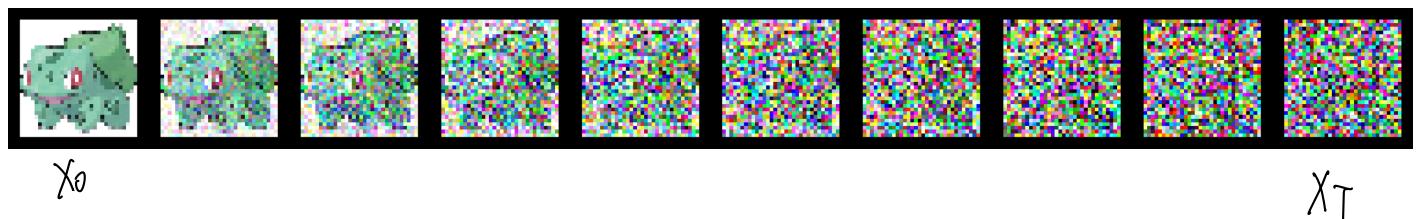
    Returns:
    - torch.Tensor: The tensor x[t] after the diffusion process.
    - torch.Tensor: The noise tensor generated during the process.
    """
    # Generate a sample from the normal distribution
    noise = torch.randn_like(x_0)
    # must be 4 dim batch_size X Channels X Img_size X Img_size
    if not x_0.dim() == 4:
        assert(False)
    # Get the mean and variance. Choose sqrt_alphas_bar_t and sqrt_one_minus_alpha
    # values based on the timestep.
    sqrt_alphas_bar_t = torch.sqrt(torch.gather(hhx_alpha_bars, 0, t))
    sqrt_one_minus_alphas_bar = torch.sqrt(1-torch.gather(hhx_alpha_bars, 0, t))
    batch_dim = x_0.size(0)

    # Compute x[t] based on x[0]
    x_t = sqrt_alphas_bar_t.view((batch_dim, 1, 1, 1)) * x_0 + \
          sqrt_one_minus_alphas_bar.view(batch_dim, 1, 1, 1) * noise
    return x_t, noise

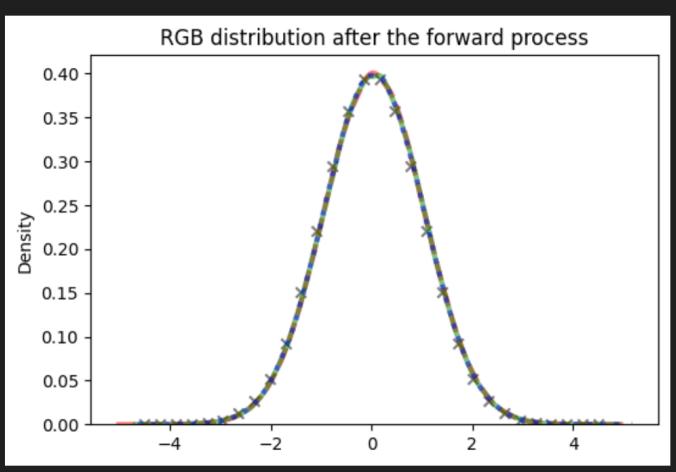
```

$$\bar{x}_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \mathcal{E}$$

forward Visualize :



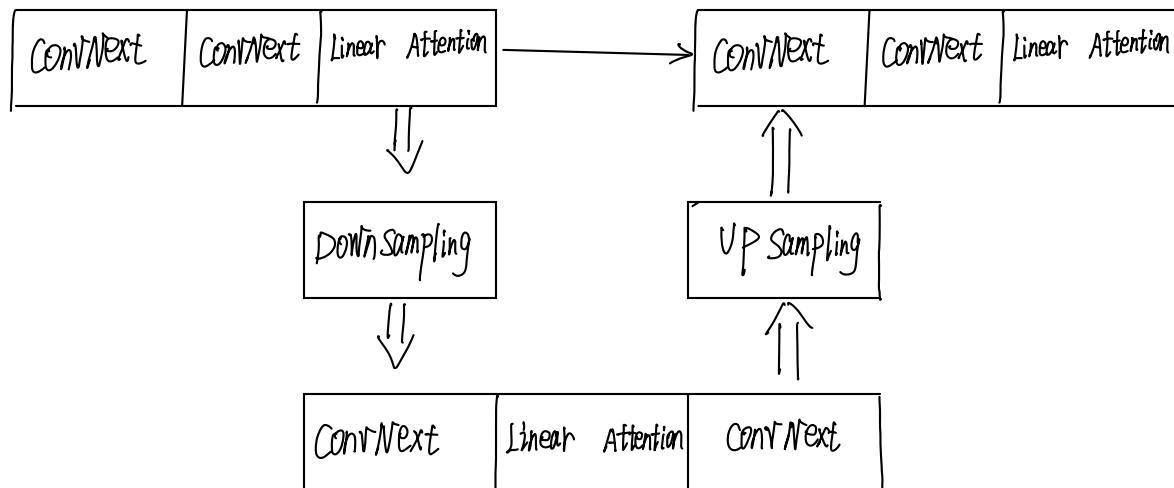
Processing a total of 919 images in the dataset.



RGB statistics follows standard Gaussian distribution

3 Diffusion Backward Noise predict model

[UNet based on ConvNextBlock]



2f. One Step Sampling

```
def denoise_one_step_sample(model, x, timestep):
    beta_t = hhx_betas[timestep]
    alpha_t = hhx_alphas[timestep]
    alpha_bar_t = hhx_alpha_bars[timestep]
    sqrt_one_minus_alpha_bar_t = torch.sqrt(1-alpha_bar_t)
    sqrt_recip_alphas_t = torch.sqrt(1.0 / alpha_t)

    timestep_tensor = torch.tensor([timestep], device=DEVICE)
    model_mean = sqrt_recip_alphas_t * (
        x - beta_t * model(x, timestep_tensor) / sqrt_one_minus_alpha_bar_t)

    alpha_bar_prev = F.pad(hhx_alpha_bars[:-1], (1, 0), value=1.0)

    # calculations for posterior q(x_{t-1} | x_t, x_0)
    posterior_variance = hhx_betas * (1. - alpha_bar_prev) / (1. - hhx_alpha_bars)

    if timestep == 0:
        return model_mean
    else:
        posterior_variance_t = posterior_variance[timestep]
        noise = torch.randn_like(x)
        return model_mean + torch.sqrt(posterior_variance_t) * noise
```

$$x_{t+1} = \frac{1}{\beta t} \left(x_t - \frac{1-\bar{\alpha}t}{1-\bar{\alpha}t} E_0(x_t, t) \right) + \frac{1-\bar{\alpha}t+1}{1-\bar{\alpha}t} \beta t E$$

$E \sim N(0, I)$

5. training

```
def train_epoch(model, dataloader, optimizer):
    loss_list = list()
    for x in dataloader:
        x= x.to(DEVICE)
        batch_size = x.shape[0]
        optimizer.zero_grad()
        t = torch.randint(0, TIMESTEPS - 1, (batch_size,), device = DEVICE)
        x_t, noise = hhx_forward_diffusion(x, t) ← forward
        noise_pred = model(x_t, t) ← reverse predict
        noise = noise.squeeze()
        # print("noise")
        # print(noise.shape)
        # print("noise_pred")
        # print(noise_pred.shape)
        loss = F.mse_loss(noise_pred, noise) ← loss = ||Et - Eθ(xt, t)||
        loss.backward()
        loss_list.append(loss.item())
        optimizer.step()

    return np.mean(loss_list)
```

training visualize:

Loss in epoch 3300 is: 0.017955930132804245
<Figure size 600x400 with 0 Axes>



APPendix: Project git record:

