

# Denoising Diffusion Implicit Models

Jiaming Song, Chenlin Meng, Stefano Ermon  
ICLR 2021

- What is Diffusion Model?
- Denoising Diffusion Probabilistic Models
- Denoising Diffusion Implicit Model
- Experiments
- Conclusions

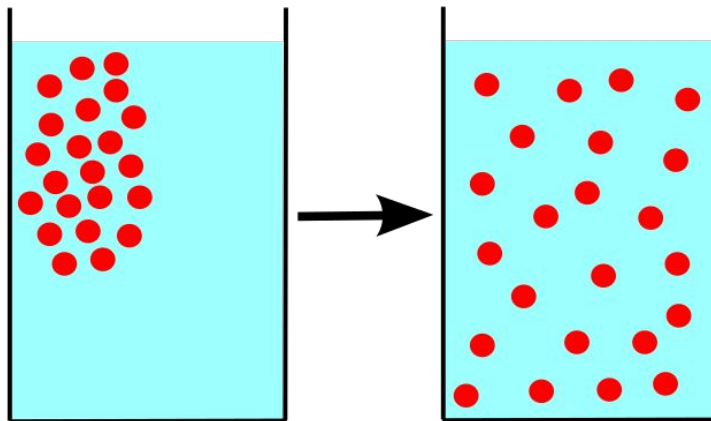
---

# What is Diffusion Model?

# ■ What is Diffusion Model?

## Diffusion

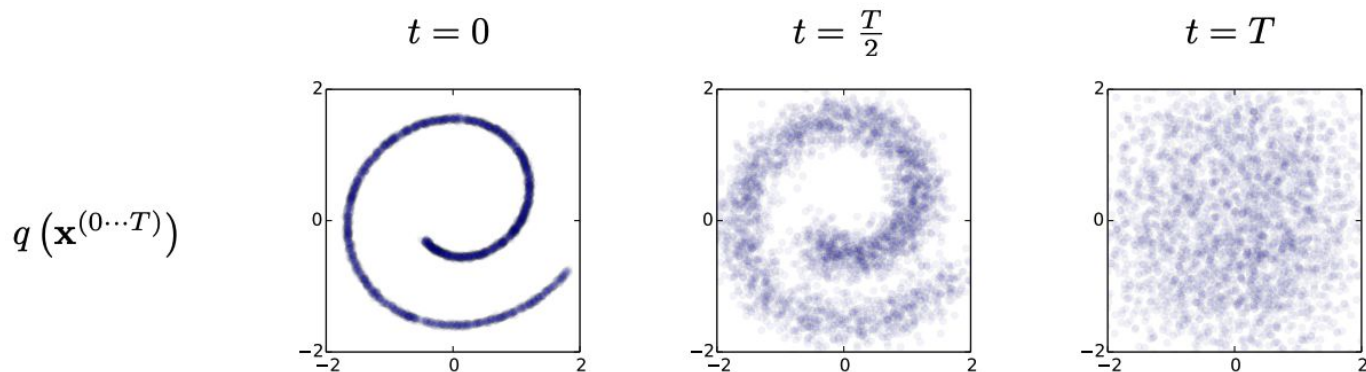
- Diffusion is the net movement of anything (for example, atoms, ions, molecules, energy) generally from a region of higher concentration to a region of lower concentration.



# What is Diffusion Model?

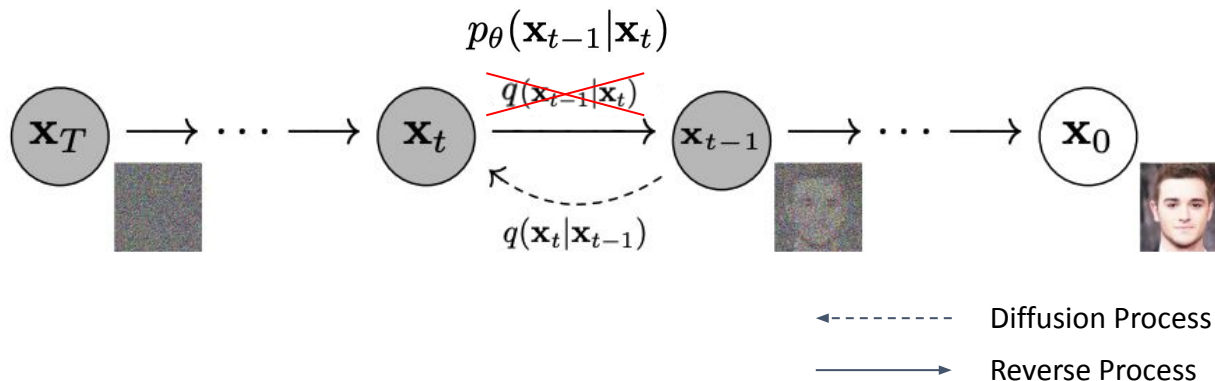
## Diffusion Process

- A (forward) diffusion process converts any complex data distribution into a simple, tractable, distribution.



# What is Diffusion Model?

## Diffusion Model

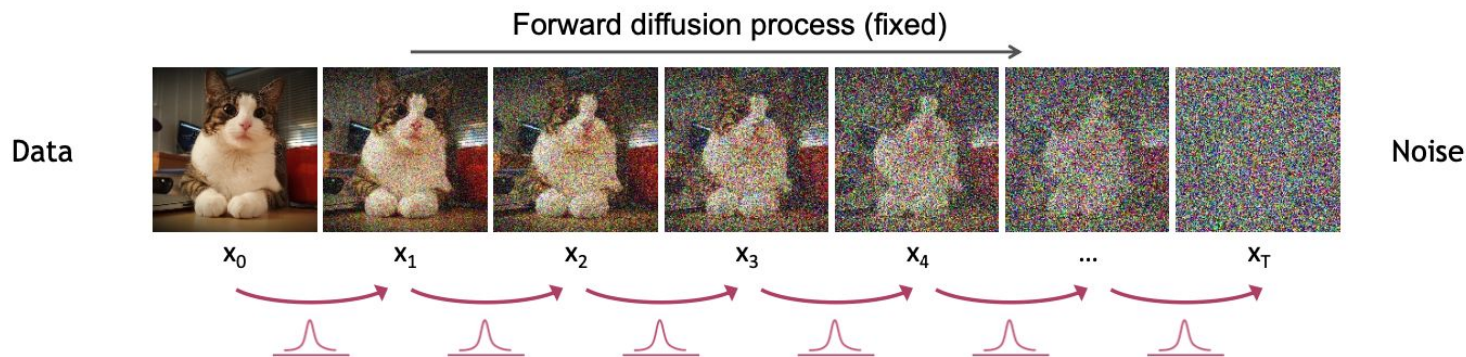


- Diffusion Process: gradually adds noise to input
- Reverse Process: learns to generate data by denoising

---

# Denoising Diffusion Probabilistic Models

## Diffusion Process



$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

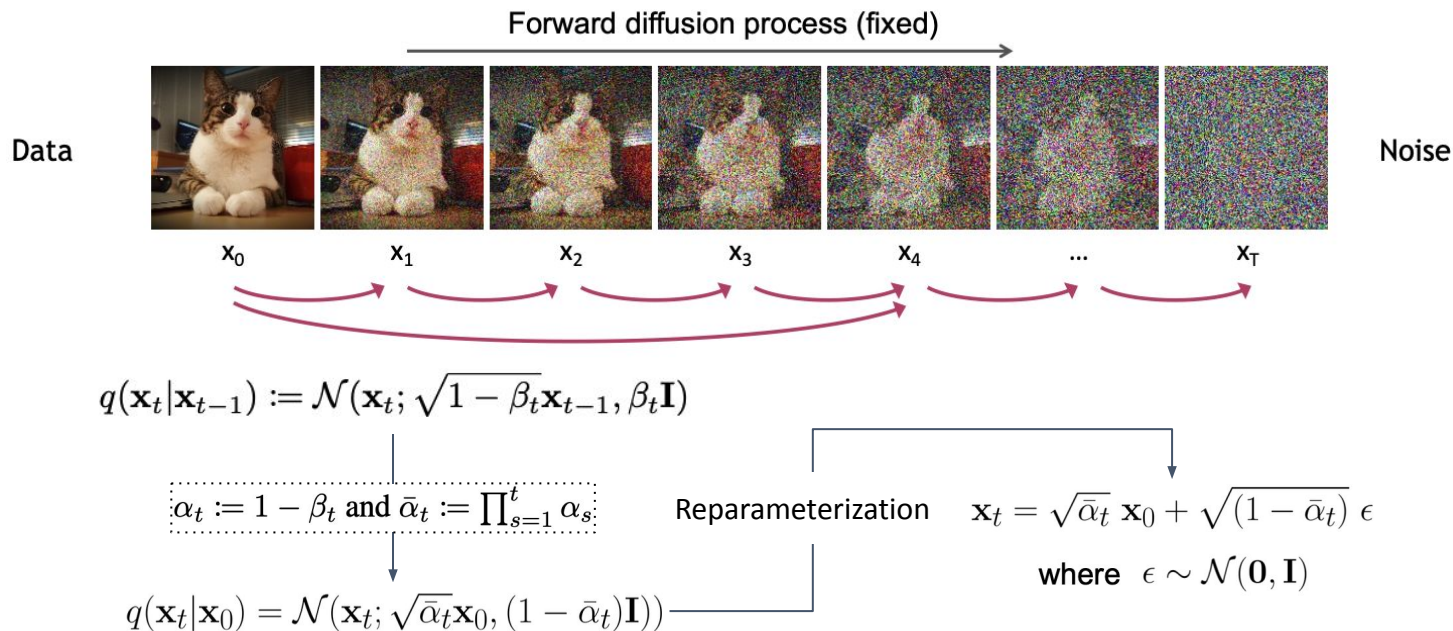
$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

$\beta_t$  increased linearly:  $\beta_1 = 10^{-4}$  to  $\beta_T = 0.02$



## Diffusion Process

- How to sample  $\mathbf{x}_t$ ?



## Diffusion Process

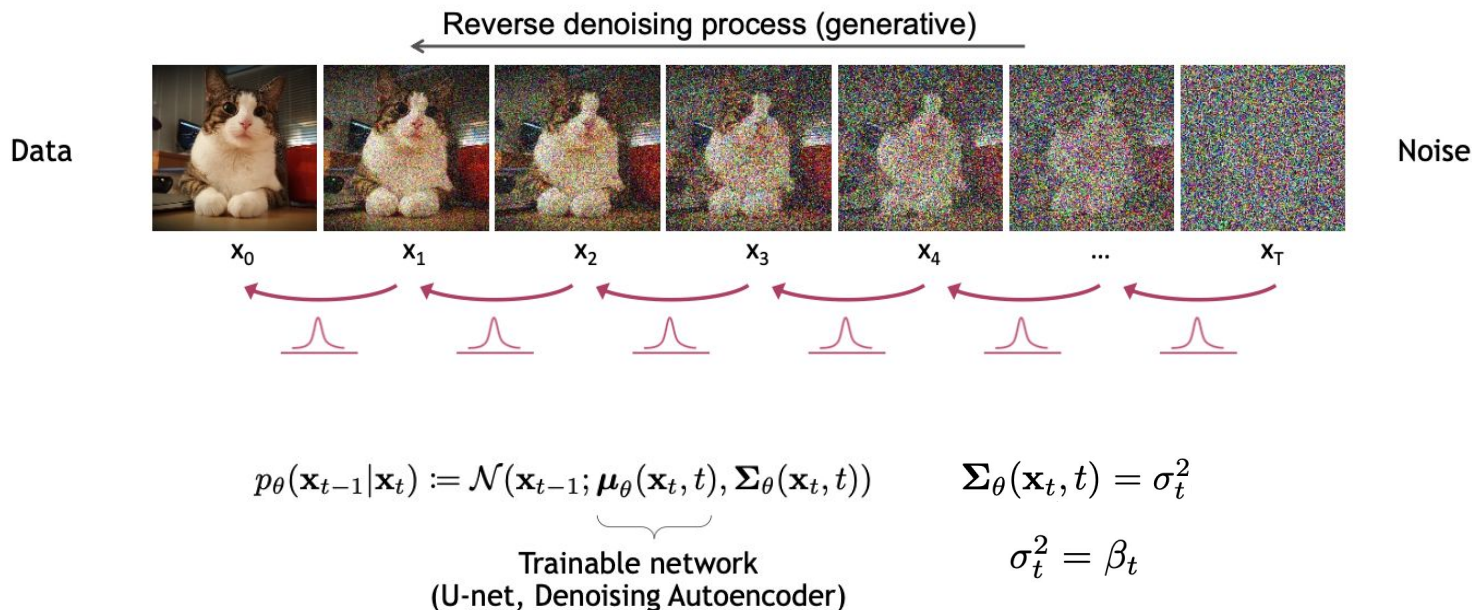
- Codes

```
alphas = 1. - betas
alphas_cumprod = torch.cumprod(alphas, dim=0)
register_buffer('sqrt_alphas_cumprod', torch.sqrt(alphas_cumprod))
register_buffer('sqrt_one_minus_alphas_cumprod', torch.sqrt(1. - alphas_cumprod))
```

```
@autocast(enabled = False)
def q_sample(self, x_start, t, noise = None):
    noise = default(noise, lambda: torch.randn_like(x_start))

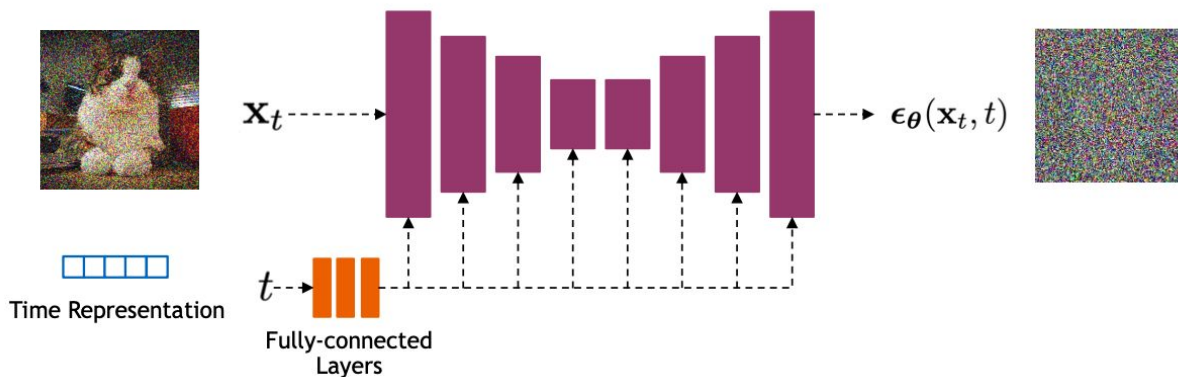
    return (
        extract(self.sqrt_alphas_cumprod, t, x_start.shape) * x_start +
        extract(self.sqrt_one_minus_alphas_cumprod, t, x_start.shape) * noise
    )
```

## Reverse Process

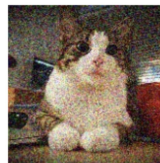


## Reverse Process with Diffusion Model

- Diffusion model consists of U-Net architectures with ResNet blocks and self-attention layers to represent  $\epsilon_{\theta}(\mathbf{x}_t, t)$ .



- We can generate  $\mathbf{x}_{t-1}$  image with  $\mathbf{x}_t$  and  $\epsilon_{\theta}(\mathbf{x}_t, t)$ .



## Reverse Process with Predicted Noise

Bayes' Rule

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

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{x}_{t-1}) &:= \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \\ q(\mathbf{x}_{t-1} | \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0, (1 - \bar{\alpha}_{t-1}) \mathbf{I}) \\ q(\mathbf{x}_t | \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \end{aligned}$$

where

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{x}_{t-1}) &= \frac{1}{\sqrt{2\pi\beta_t}} \exp\left(-\frac{(\mathbf{x}_t - \sqrt{1 - \beta_t} \mathbf{x}_{t-1})^2}{2\beta_t}\right) \\ q(\mathbf{x}_t | \mathbf{x}_0) &= \frac{1}{\sqrt{2\pi(1 - \bar{\alpha}_t)}} \exp\left(-\frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0)^2}{2(1 - \bar{\alpha}_t)}\right) \\ q(\mathbf{x}_{t-1} | \mathbf{x}_0) &= \frac{1}{\sqrt{2\pi(1 - \bar{\alpha}_{t-1})}} \exp\left(-\frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0)^2}{2(1 - \bar{\alpha}_{t-1})}\right) \end{aligned}$$

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \approx \frac{1}{\sqrt{2\pi\tilde{\beta}_t}} \exp\left(-\frac{1}{2\tilde{\beta}_t} \left[ \mathbf{x}_{t-1} - \left( \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_t \right) \right]^2\right)$$

$\tilde{\beta}_t$   $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0)$

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

## Reverse Process with Predicted Noise

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\boldsymbol{\beta}}_t \mathbf{I}),$$

where  $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t$  and  $\tilde{\boldsymbol{\beta}}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$

Set as  $\beta_t$

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \boldsymbol{\epsilon} \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \longrightarrow \mathbf{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \times \mathbf{x}_t - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \times \boldsymbol{\epsilon}$$

► Above equation also can be used for predicted  $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$

## Reverse Process

```

@torch.inference_mode()
def p_sample(self, x, t: int, x_self_cond = None):
    b, *_ , device = *x.shape, self.device
    batched_times = torch.full((b,), t, device = device, dtype = torch.long)
    model_mean, _, model_log_variance, x_start = self.p_mean_variance(x = x, t = batched_times, x_self_cond = x_self_cond, clip_denoised = True)
    noise = torch.randn_like(x) if t > 0 else 0. # no noise if t == 0
    pred_img = model_mean + (0.5 * model_log_variance).exp() * noise
    return pred_img, x_start

@torch.inference_mode()
def p_sample_loop(self, shape, return_all_timesteps = False):
    batch, device = shape[0], self.device

    img = torch.randn(shape, device = device)
    imgs = [img]


    x_start = None

    for t in tqdm(reversed(range(0, self.num_timesteps)), desc = 'sampling loop time step', total = self.num_timesteps):
        self_cond = x_start if self.self_condition else None
        img, x_start = self.p_sample(img, t, self_cond)
        imgs.append(img)

    ret = img if not return_all_timesteps else torch.stack(imgs, dim = 1)

    ret = self.unnormalize(ret)
    return ret

```



## Diffusion Loss

$$\begin{aligned}
 & \mathbb{E}[-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\
 & \quad \text{Variational Bound} \\
 & \quad \text{Proof: Appendix A} \\
 & \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \| p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} + \underbrace{-\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right] \\
 & \quad \text{Denoising Matching} \\
 & \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[ \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2 \right] \\
 & \quad \downarrow \\
 & L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[ \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2 \right]
 \end{aligned}$$



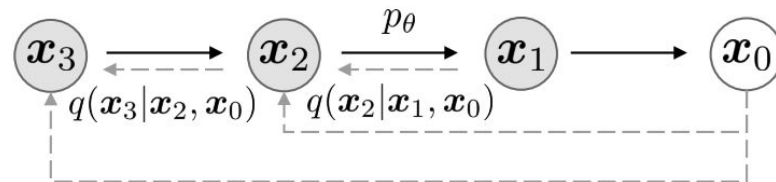
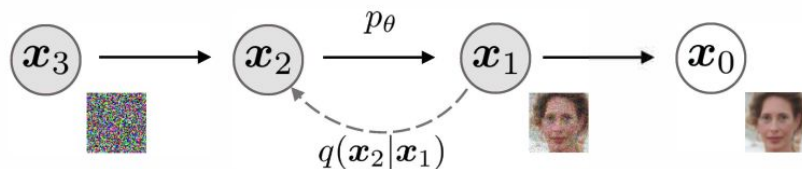
---

# Denoising Diffusion Implicit Model

## Main Idea

- Denoising diffusion probabilistic models (DDPMs) have achieved high quality image generation without adversarial training, yet they **require simulating a Markov chain for many steps in order to produce a sample.**
- Denoising diffusion implicit models (DDIMs) used **Non-Markovian processes** that can correspond to generative processes that are deterministic, giving rise to implicit models that **produce high quality samples much faster.**

## Non-Markovian Diffusion Process



$$q_{\sigma}(\mathbf{x}_{1:T}|\mathbf{x}_0) := q_{\sigma}(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^T q_{\sigma}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

$$\rightarrow \cancel{q(\mathbf{x}_1|\mathbf{x}_0)} q(\mathbf{x}_2|\mathbf{x}_1, \mathbf{x}_0) q(\mathbf{x}_3|\mathbf{x}_2, \mathbf{x}_0) \dots q(\mathbf{x}_T|\mathbf{x}_{T-1}, \mathbf{x}_0)$$

Bayes' Rule

$$q_{\sigma}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q_{\sigma}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q_{\sigma}(\mathbf{x}_t|\mathbf{x}_0)}{q_{\sigma}(\mathbf{x}_{t-1}|\mathbf{x}_0)}$$

$$\frac{q_{\sigma}(\mathbf{x}_1|\mathbf{x}_2, \mathbf{x}_0) \cancel{q_{\sigma}(\mathbf{x}_2|\mathbf{x}_0)}}{\cancel{q_{\sigma}(\mathbf{x}_1|\mathbf{x}_0)}} \times \frac{q_{\sigma}(\mathbf{x}_2|\mathbf{x}_3, \mathbf{x}_0) \cancel{q_{\sigma}(\mathbf{x}_3|\mathbf{x}_0)}}{\cancel{q_{\sigma}(\mathbf{x}_2|\mathbf{x}_0)}} \times \dots \times \frac{q_{\sigma}(\mathbf{x}_{T-1}|\mathbf{x}_T, \mathbf{x}_0) q_{\sigma}(\mathbf{x}_T|\mathbf{x}_0)}{\cancel{q_{\sigma}(\mathbf{x}_{T-1}|\mathbf{x}_0)}}$$

## Non-Markovian Diffusion Process

$$q_{\sigma}(\mathbf{x}_{1:T}|\mathbf{x}_0) := q_{\sigma}(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^T q_{\sigma}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

↓ Proof: Appendix B

$$q_{\sigma}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right)$$

$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left( \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } \mathbf{x}_0"} + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}_{\text{"direction pointing to } \mathbf{x}_t"} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

Based on DDPM

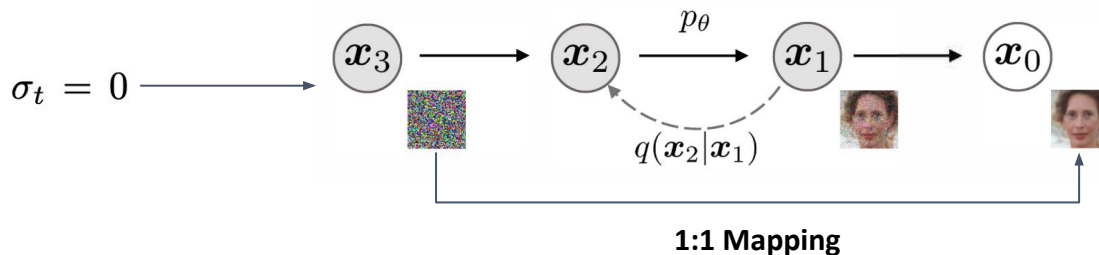
$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \times x_t - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \times \epsilon_{\theta}(x_t, t)$$

$$\epsilon_{\theta}(x_t, t) = \frac{x_t - \sqrt{\bar{\alpha}_t} \times x_0}{\sqrt{1 - \bar{\alpha}_t}}$$

$$x_t = \sqrt{\bar{\alpha}_t} \times x_0 - \sqrt{1 - \bar{\alpha}_t} \times \epsilon_{\theta}(x_t, t)$$

## Deterministic Generative Process

$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left( \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } \mathbf{x}_0"} + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}_{\text{"direction pointing to } \mathbf{x}_t"} + \cancel{\underbrace{\sigma_t \epsilon_t}_{\text{random noise}}}$$



## Accelerated Generation Processes

$$q_{\sigma}(\mathbf{x}_{1:T}|\mathbf{x}_0) := q_{\sigma}(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^T q_{\sigma}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

Proof: Appendix C.1

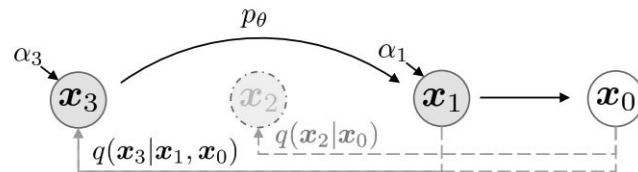


Figure 2: Graphical model for accelerated generation, where  $\tau = [1, 3]$ .

In the previous sections, the generative process is considered as the approximation to the reverse process; since of the forward process has  $T$  steps, the generative process is also forced to sample  $T$  steps. However, as the denoising objective  $L_1$  does not depend on the specific forward procedure as long as  $q_{\sigma}(\mathbf{x}_t|\mathbf{x}_0)$  is fixed, we may also consider forward processes with lengths smaller than  $T$ , which accelerates the corresponding generative processes without having to train a different model.

---

# Experiments

Table 1: CIFAR10 and CelebA image generation measured in FID.  $\eta = 1.0$  and  $\hat{\sigma}$  are cases of **DDPM** (although [Ho et al. \(2020\)](#) only considered  $T = 1000$  steps, and  $S < T$  can be seen as simulating DDPMs trained with  $S$  steps), and  $\eta = 0.0$  indicates **DDIM**.

$S$	CIFAR10 ( $32 \times 32$ )					CelebA ( $64 \times 64$ )					
	10	20	50	100	1000	10	20	50	100	1000	
$\eta$	0.0	<b>13.36</b>	<b>6.84</b>	<b>4.67</b>	<b>4.16</b>	4.04	<b>17.33</b>	<b>13.73</b>	<b>9.17</b>	<b>6.53</b>	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	<b>3.17</b>	299.71	183.83	71.71	45.20	<b>3.26</b>	

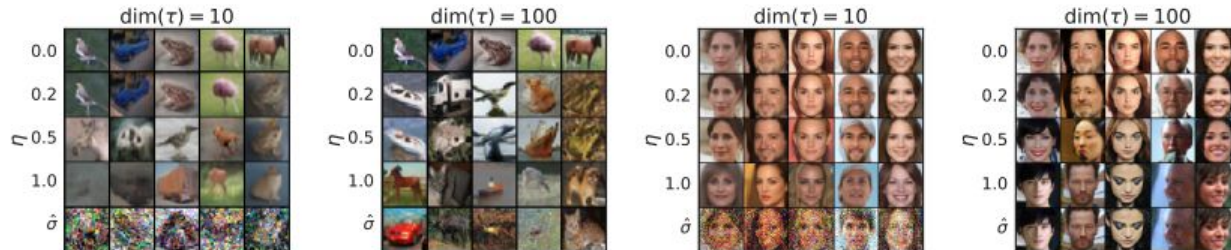
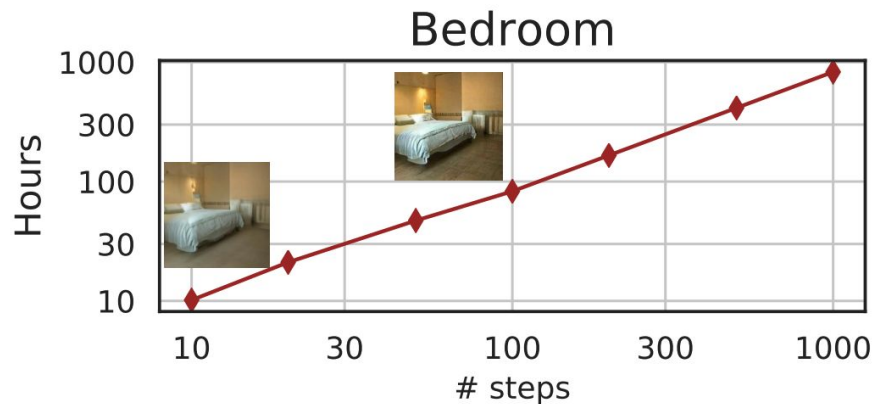
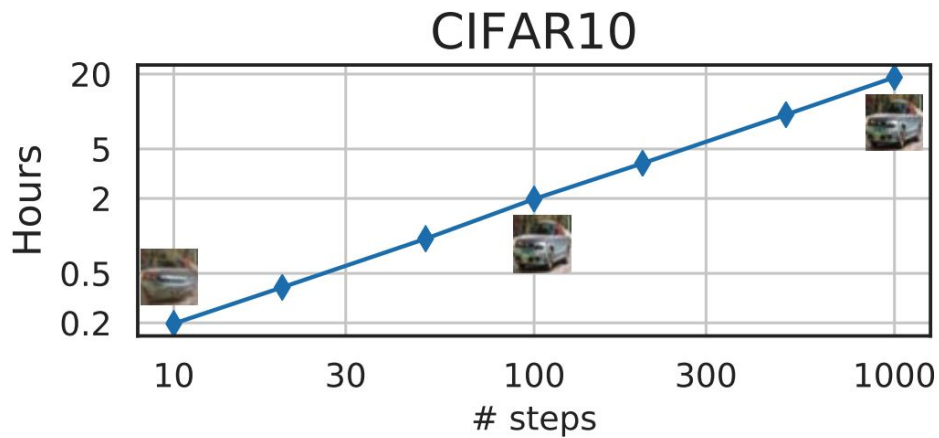


Figure 3: CIFAR10 and CelebA samples with  $\dim(\tau) = 10$  and  $\dim(\tau) = 100$ .







# Experiments

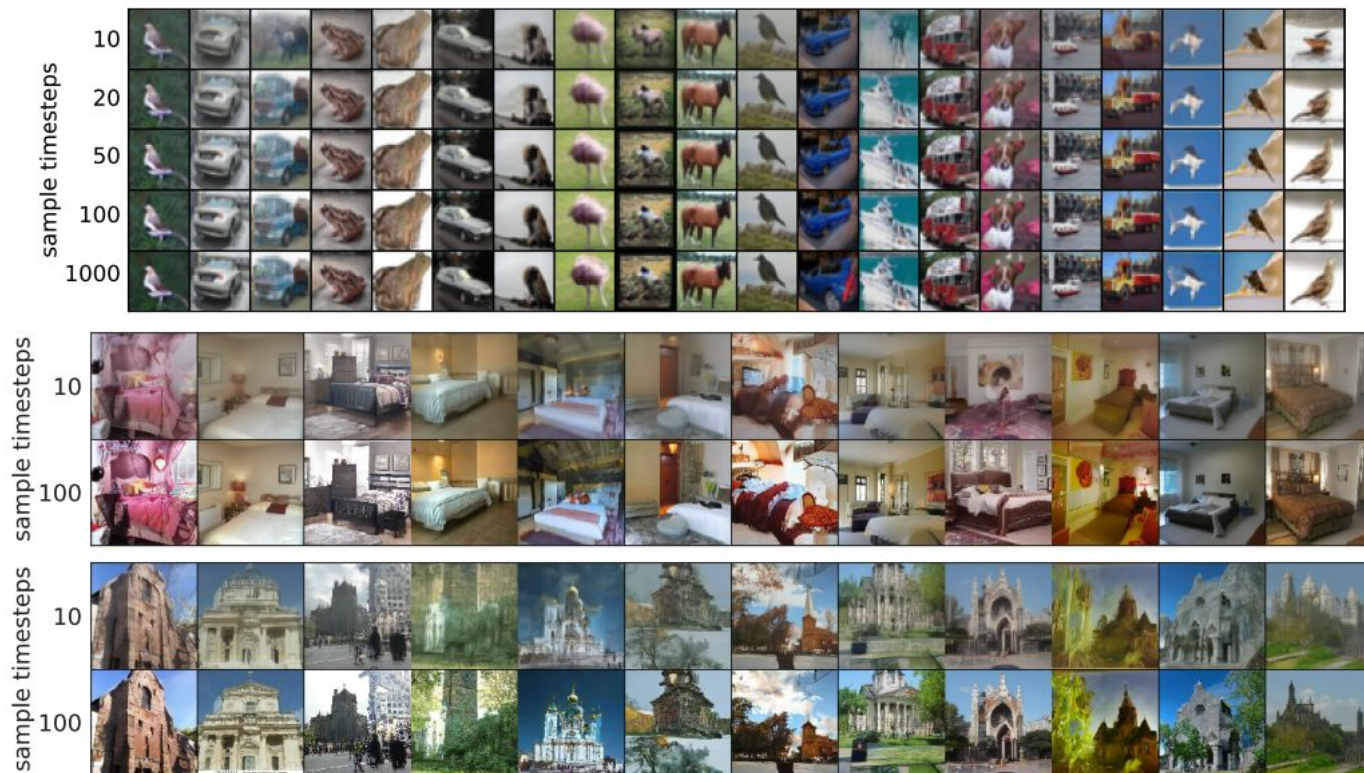


Figure 5: Samples from DDIM with the same random  $x_T$  and different number of steps.

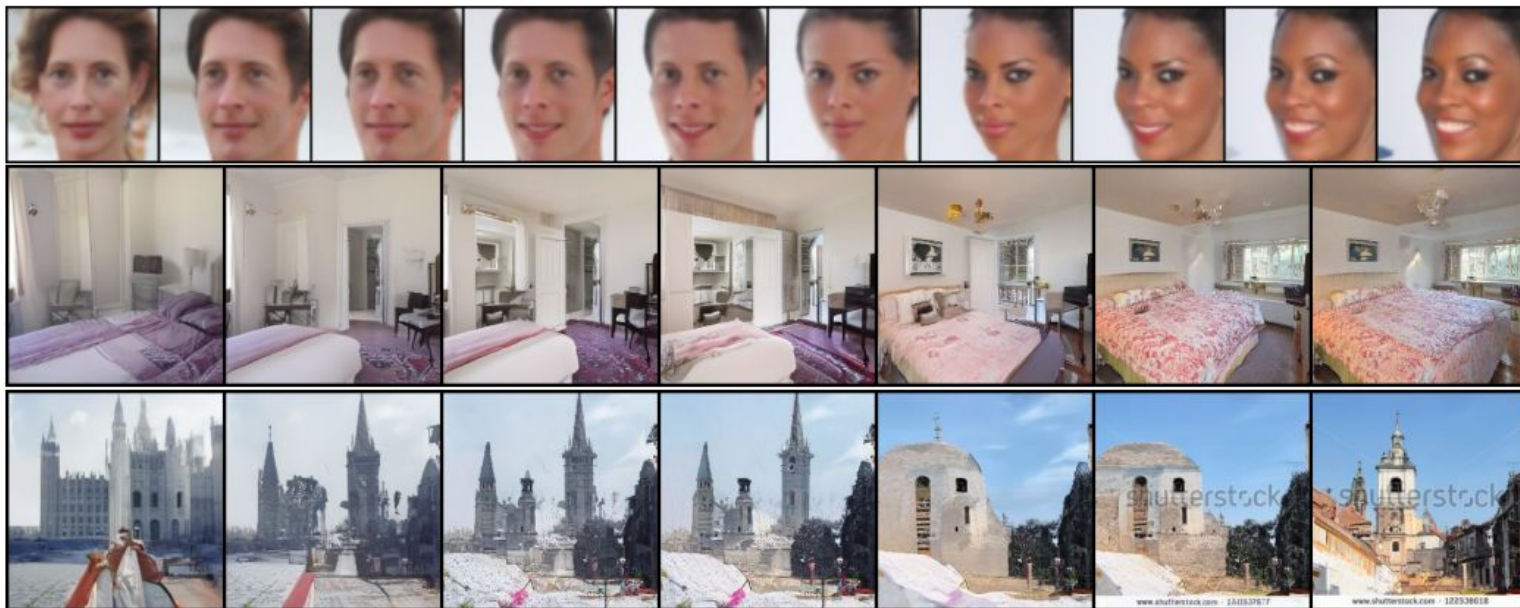


Figure 6: Interpolation of samples from DDIM with  $\dim(\tau) = 50$ .

---

# Conclusions



- DDIMs can **produce high quality samples 10× to 50× faster** in terms of wall-clock time compared to DDPMs, allow us to trade off computation for sample quality applying non-Markovian processes.
- After **training with DDPM**, we can use **DDIM for inference**, because they have the same training procedure.



Thank you



---

# Appendix

## C.1 ACCELERATED SAMPLING PROCESSES

In the accelerated case, we can consider the inference process to be factored as:

$$q_{\sigma,\tau}(\mathbf{x}_{1:T}|\mathbf{x}_0) = q_{\sigma,\tau}(\mathbf{x}_{\tau_S}|\mathbf{x}_0) \prod_{i=1}^S q_{\sigma,\tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}, \mathbf{x}_0) \prod_{t \in \bar{\tau}} q_{\sigma,\tau}(\mathbf{x}_t|\mathbf{x}_0) \quad (52)$$

where  $\tau$  is a sub-sequence of  $[1, \dots, T]$  of length  $S$  with  $\tau_S = T$ , and let  $\bar{\tau} := \{1, \dots, T\} \setminus \tau$  be its complement. Intuitively, the graphical model of  $\{\mathbf{x}_{\tau_i}\}_{i=1}^S$  and  $\mathbf{x}_0$  form a chain, whereas the graphical model of  $\{\mathbf{x}_t\}_{t \in \bar{\tau}}$  and  $\mathbf{x}_0$  forms a star graph. We define:

$$q_{\sigma,\tau}(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_t}\mathbf{x}_0, (1 - \alpha_t)\mathbf{I}) \quad \forall t \in \bar{\tau} \cup \{T\} \quad (53)$$

$$q_{\sigma,\tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{\tau_{i-1}}}\mathbf{x}_0 + \sqrt{1 - \alpha_{\tau_{i-1}} - \sigma_{\tau_i}^2} \cdot \frac{\mathbf{x}_{\tau_i} - \sqrt{\alpha_{\tau_i}}\mathbf{x}_0}{\sqrt{1 - \alpha_{\tau_i}}}, \sigma_{\tau_i}^2 \mathbf{I}\right) \quad \forall i \in [S]$$

where the coefficients are chosen such that:

$$q_{\sigma,\tau}(\mathbf{x}_{\tau_i}|\mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_{\tau_i}}\mathbf{x}_0, (1 - \alpha_{\tau_i})\mathbf{I}) \quad \forall i \in [S] \quad (54)$$

i.e., the “marginals” match.



The corresponding “generative process” is defined as:

$$p_{\theta}(\mathbf{x}_{0:T}) := \underbrace{p_{\theta}(\mathbf{x}_T) \prod_{i=1}^S p_{\theta}^{(\tau_i)}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i})}_{\text{use to produce samples}} \times \underbrace{\prod_{t \in \bar{\tau}} p_{\theta}^{(t)}(\mathbf{x}_0|\mathbf{x}_t)}_{\text{in variational objective}} \quad (55)$$

where only part of the models are actually being used to produce samples. The conditionals are:

$$p_{\theta}^{(\tau_i)}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}) = q_{\sigma,\tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}, f_{\theta}^{(\tau_i)}(\mathbf{x}_{\tau_{i-1}})) \quad \text{if } i \in [S], i > 1 \quad (56)$$

$$p_{\theta}^{(t)}(\mathbf{x}_0|\mathbf{x}_t) = \mathcal{N}(f_{\theta}^{(t)}(\mathbf{x}_t), \sigma_t^2 \mathbf{I}) \quad \text{otherwise,} \quad (57)$$

where we leverage  $q_{\sigma,\tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}, \mathbf{x}_0)$  as part of the inference process (similar to what we have done in Section 3). The resulting variational objective becomes (define  $\mathbf{x}_{\tau_{L+1}} = \emptyset$  for conciseness):

$$J(\epsilon_{\theta}) = \mathbb{E}_{\mathbf{x}_{0:T} \sim q_{\sigma,\tau}(\mathbf{x}_{0:T})} [\log q_{\sigma,\tau}(\mathbf{x}_{1:T}|\mathbf{x}_0) - \log p_{\theta}(\mathbf{x}_{0:T})] \quad (58)$$

$$= \mathbb{E}_{\mathbf{x}_{0:T} \sim q_{\sigma,\tau}(\mathbf{x}_{0:T})} \left[ \sum_{t \in \bar{\tau}} D_{\text{KL}}(q_{\sigma,\tau}(\mathbf{x}_t|\mathbf{x}_0) \| p_{\theta}^{(t)}(\mathbf{x}_0|\mathbf{x}_t)) + \sum_{i=1}^L D_{\text{KL}}(q_{\sigma,\tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}, \mathbf{x}_0) \| p_{\theta}^{(\tau_i)}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i})) \right] \quad (59)$$

where each KL divergence is between two Gaussians with variance independent of  $\theta$ . A similar argument to the proof used in Theorem 1 can show that the variational objective  $J$  can also be converted to an objective of the form  $L_{\gamma}$ .