

# Denoising Diffusion Probability Model

**Jonathan Ho**  
UC Berkeley  
jonathanho@berkeley.edu

**Ajay Jain**  
UC Berkeley  
ajayj@berkeley.edu

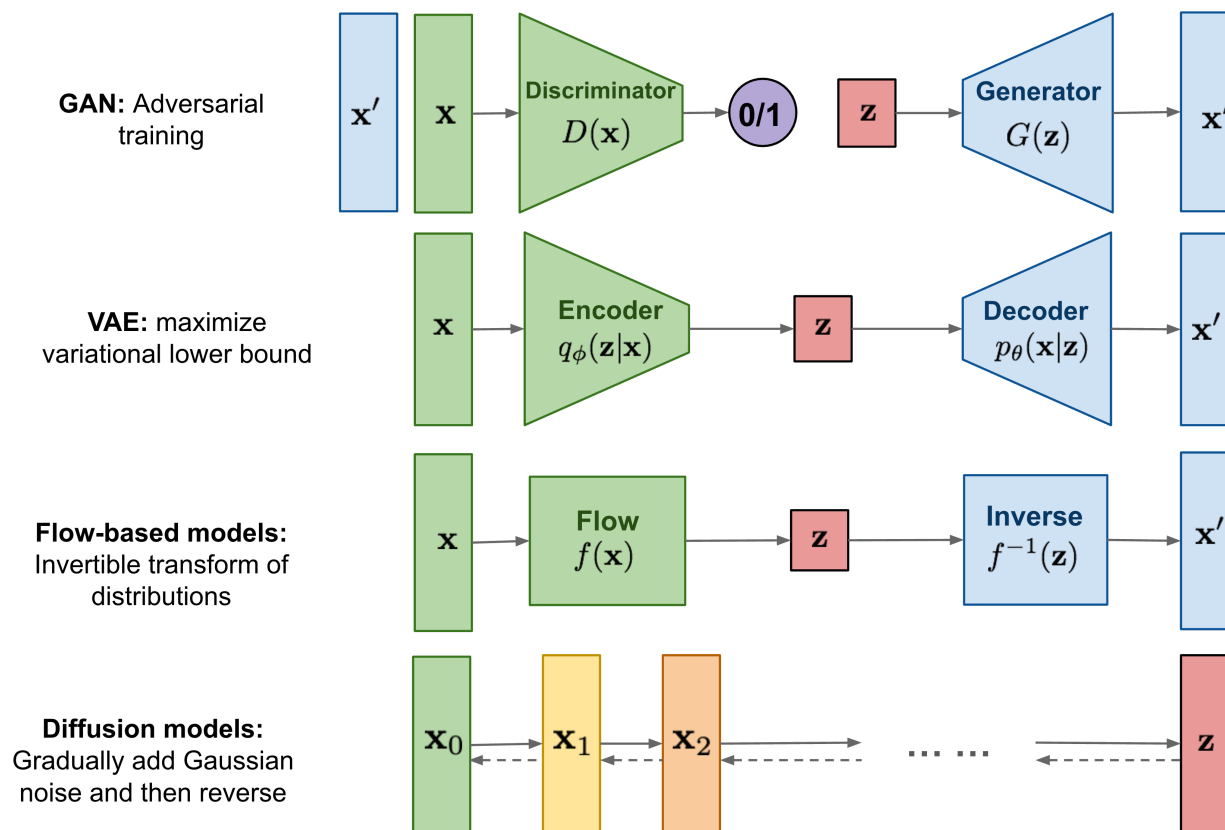
**Pieter Abbeel**  
UC Berkeley  
pabbeel@cs.berkeley.edu

NeurIPS, 2020

Miao's Group - Paper Reading



# Image Generation Models



# Image Generation: GAN

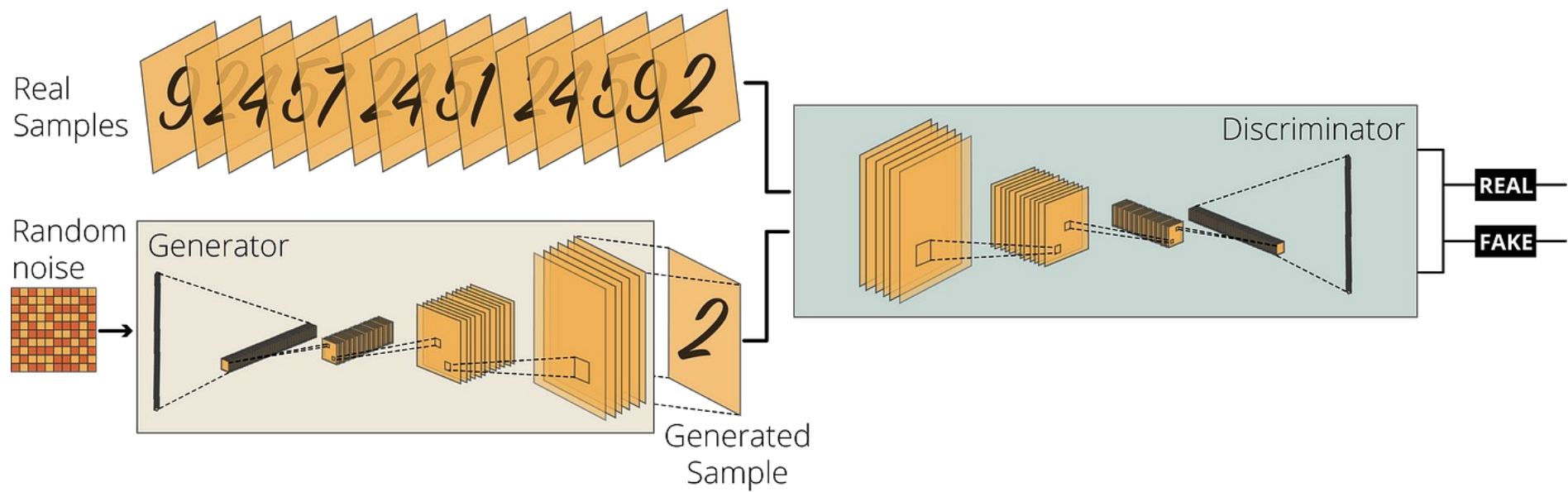
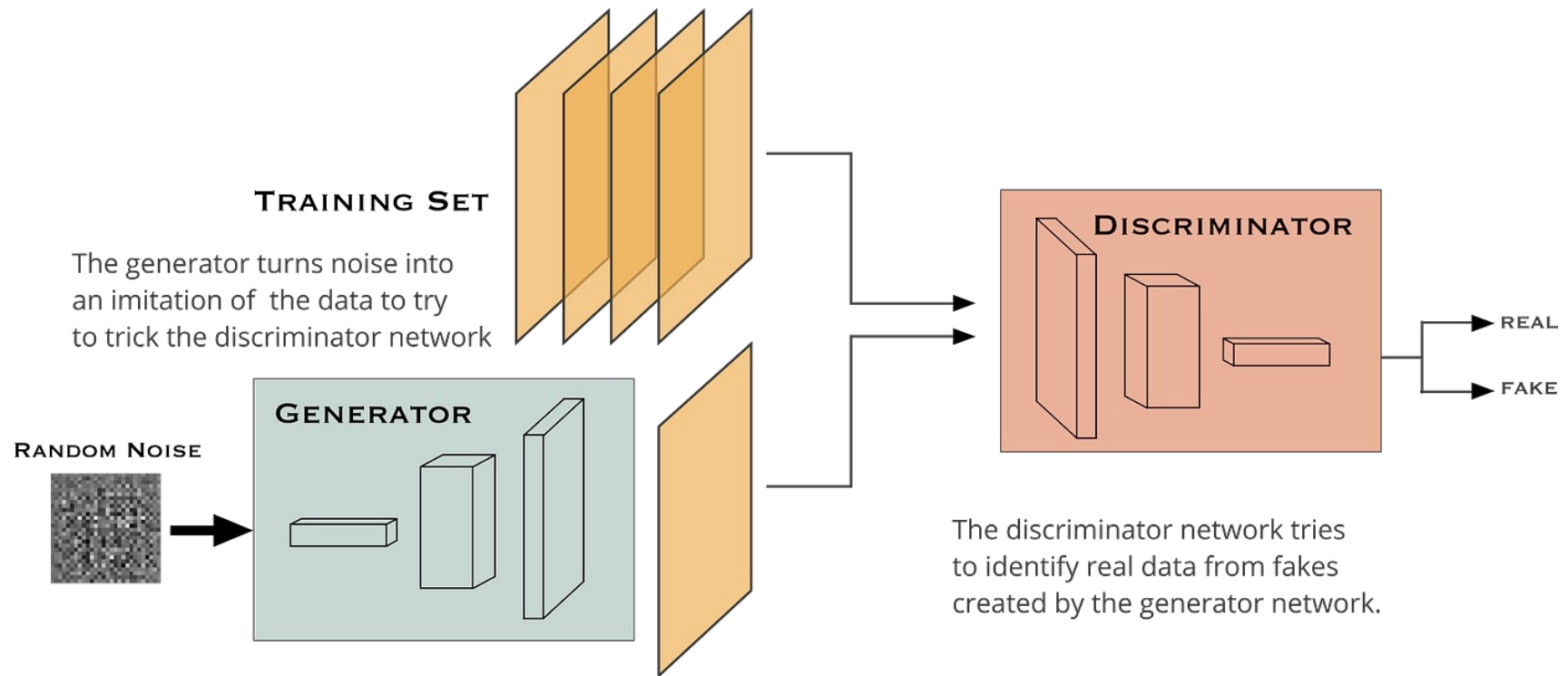


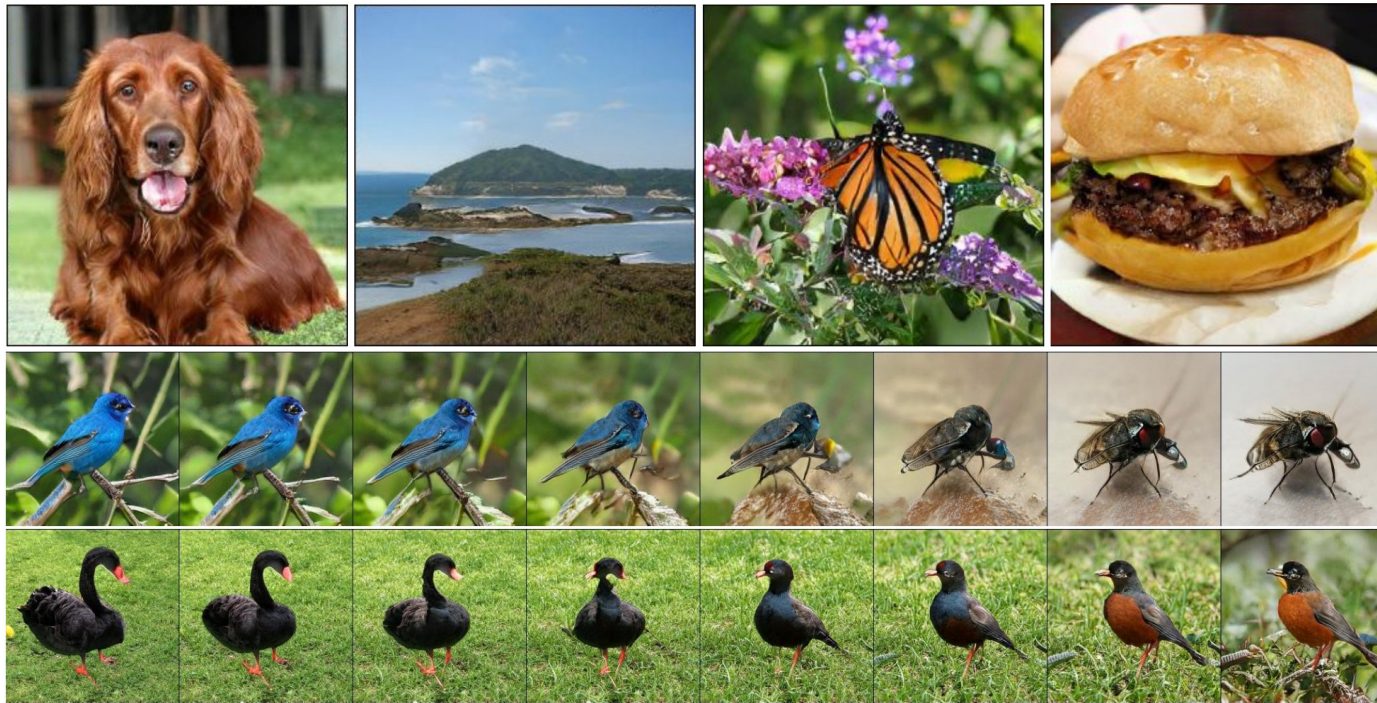
Figure adopted from: <https://towardsdatascience.com/image-generation-in-10-minutes-with-generative-adversarial-networks-c2afc56bfa3b>

# Image Generation: GAN



# Image Generation: GAN

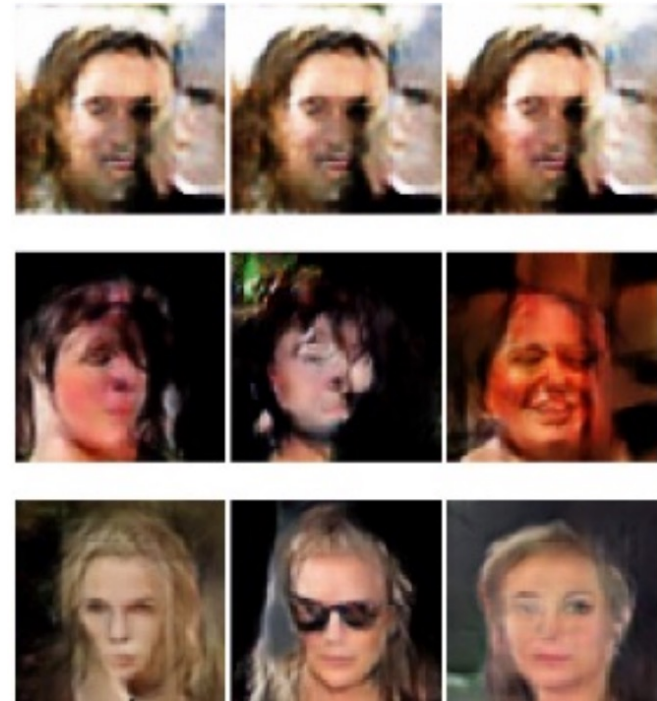
BigGAN: High-resolution image (train on 128x128 ImageNet images, generate 128/256/512 images)



A. Brock, et al. "Large Scale GAN Training for High Fidelity Natural Image Synthesis," ICLR, 2019.

## Problems of GANs

- Difficult to train due to the adversarial nature of the problem formulation
- Outputs lack diversity
- Mode collapse
- Vanishing gradients
- Problem learning multimodal distribution

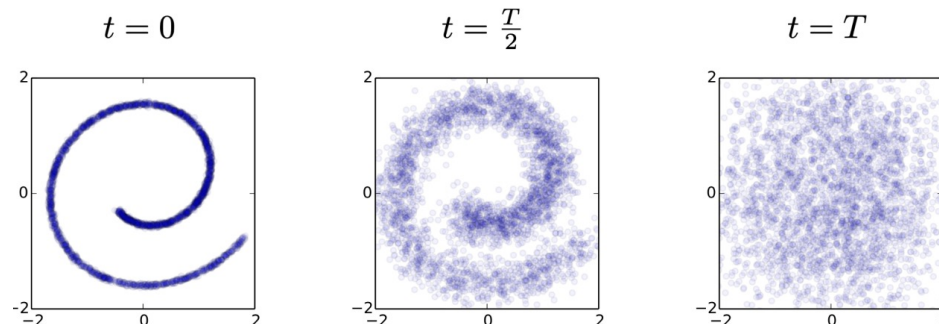




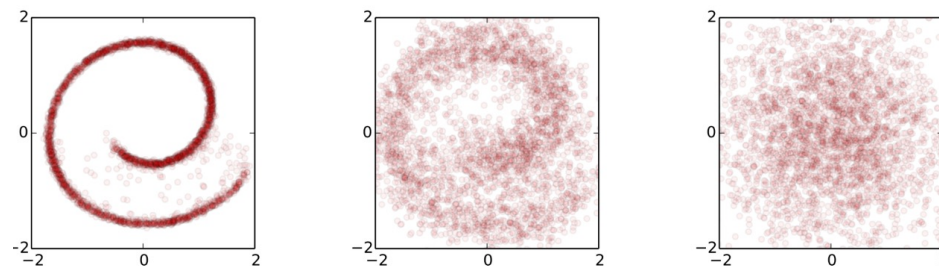
# Image Generation: Diffusion Models

- The process originates from **probabilistic likelihood estimation**, and take inspiration from physical phenomenon
- Central idea comes from the **thermodynamics of gas molecules**, whereby the molecules diffuse from high density to low density areas
- This movement is often referred in physics literature as the increase of entropy

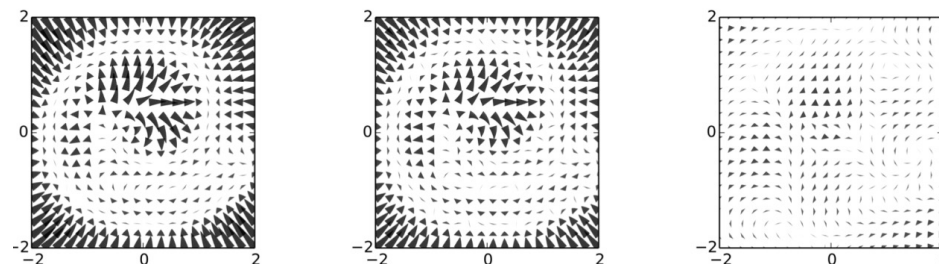
The forward trajectory  
 $q(\mathbf{x}_{0:T})$



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

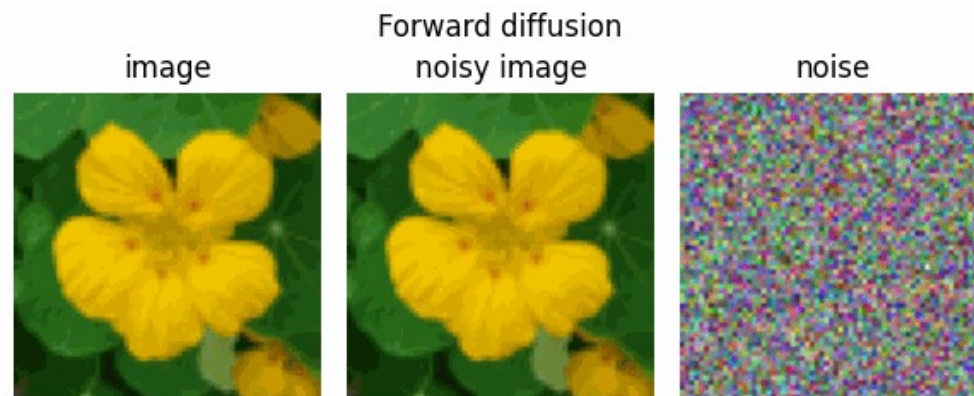


The drifting term  
 $\mu_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_t$



# Denoising Diffusion Probability Model

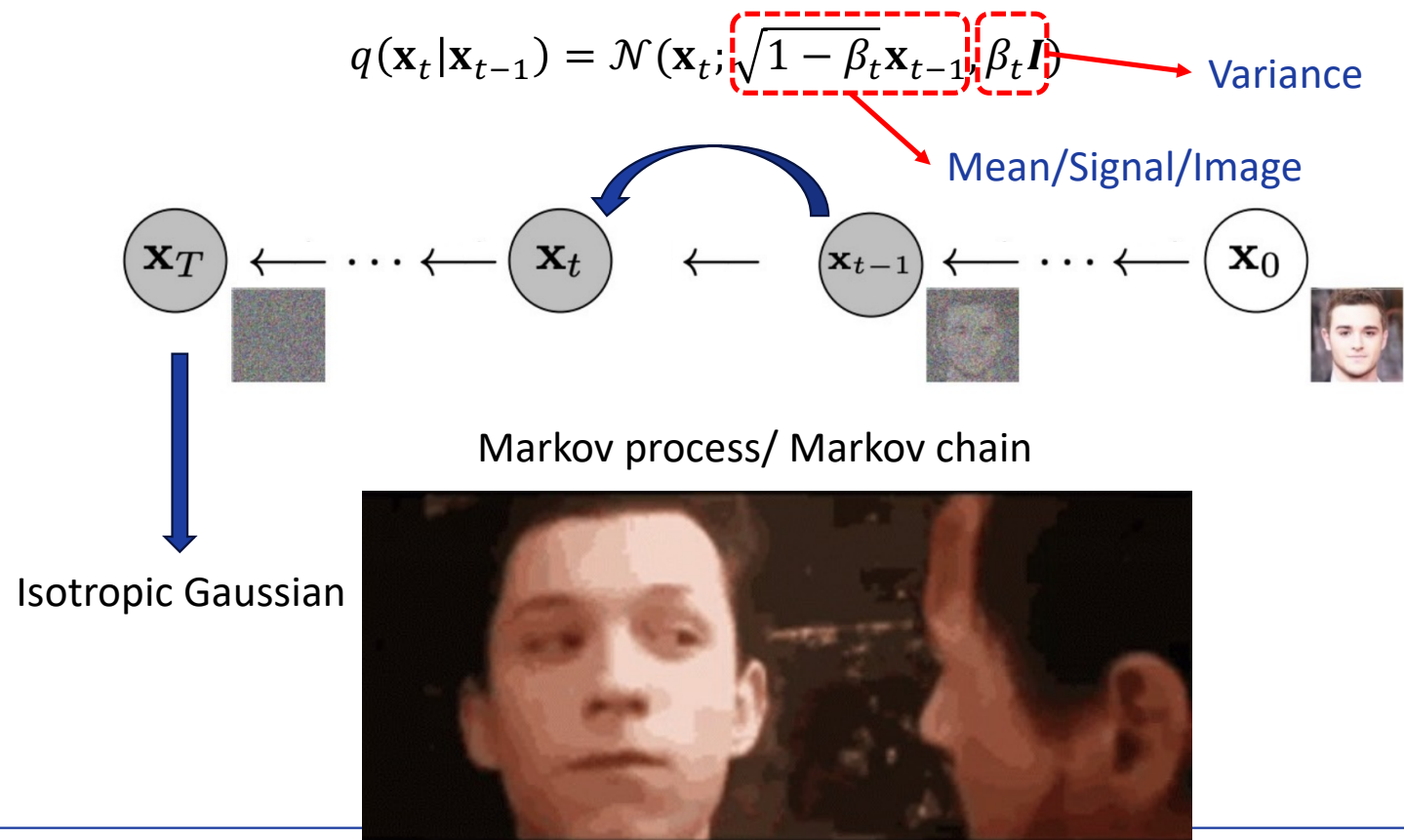
**Statistical definition:** *“Diffusion is the process of transforming a complex distribution into a predefined simpler one”*





# Denoising Diffusion Probability Model

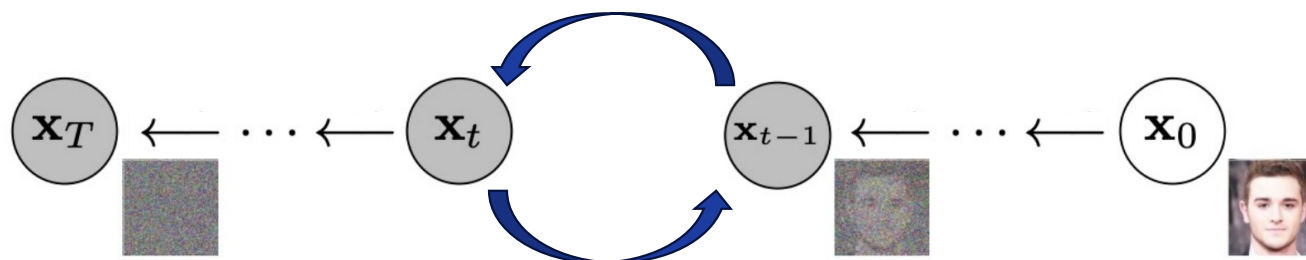
## Forward Diffusion Process



# Denoising Diffusion Probability Model

## Reverse Diffusion Process

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



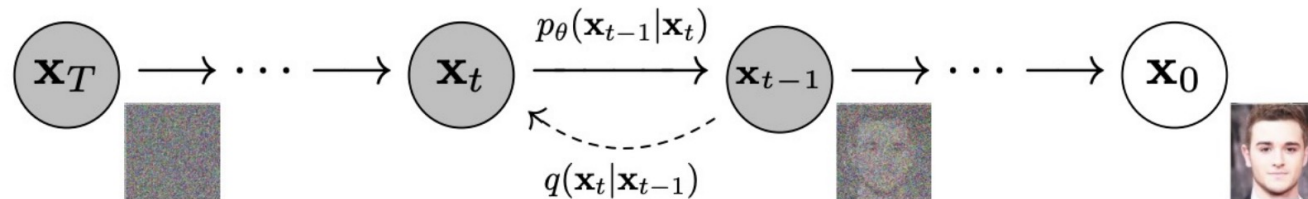
$$q(\mathbf{x}_{t-1}|\mathbf{x}_t) \propto q(\mathbf{x}_{t-1})q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad \text{Intractable} \rightarrow \text{Machine Learning}$$



Neural Network:

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

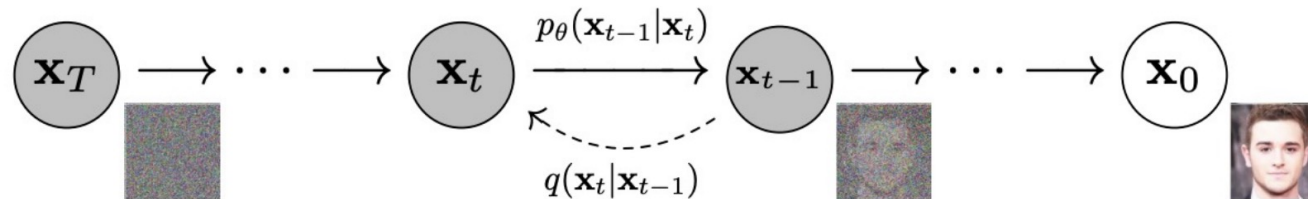
# Denoising Diffusion Probability Model



**Key concept behind:**

- Denoising diffusion models; **two steps**: 1) forward, 2) reverse / reconstruction
- Build a learning model which can learn the **systematic decay** of information due to noise
- Reverse the process and therefore, **recover** the information back from the noise

# Denoising Diffusion Probability Model



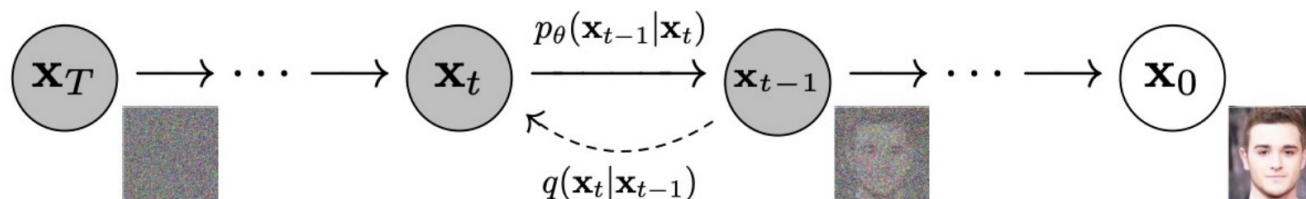
**Forward diffusion:**

- Gaussian noise is **introduced successively** until the data becomes all noise

**Reverse diffusion:**

- Undoes the noise by **learning the conditional probability densities** using a neural network

# Denoising Diffusion Probability Model

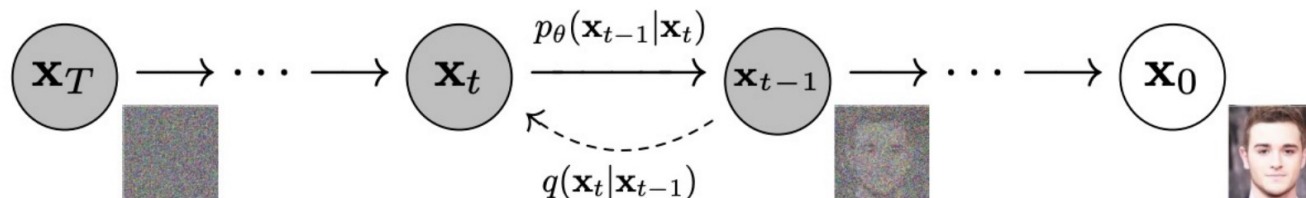


## Forward step:

- We can formally define the forward diffusion process as a Markov Chain
- Starting with the initial data point, we add **Gaussian noise** for  **$T$  successive steps**, and obtain a set of noisy samples
- The prediction of probability density at **time  $t$**  is only dependent on the immediate predecessor at **time  $t - 1$**
- The **conditional probability density** can be computed as:

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

# Denoising Diffusion Probability Model



## Forward step:

- The **conditional probability density** can be computed as:

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

- The **complete distribution** of the whole process can then be computed as follows:

$$q(\mathbf{x}_{0:T}|\mathbf{x}_0) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

- The **mean** and **variance** of the density function depends on a parameter  $\beta_t$ , which is a hyper parameter whose value can either be taken as a **constant** throughout the process or can be **gradually changed** in the successive steps



# Denoising Diffusion Probability Model

## Forward step:

- The above derivation is enough to predict the successive states
- We would like to sample at any given time interval  $t$  without going through all the intermediary steps
- This allows an efficient implementation
- Re-formulation:

$$\alpha_t = 1 - \beta_t, \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}\end{aligned}$$

## Conditional probability density:

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

## Complete distribution:

$$q(\mathbf{x}_{0:T} | \mathbf{x}_0) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

$$\longrightarrow q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

# Denoising Diffusion Probability Model

## Reverse step:

- The reverse process requires the estimation of probability density at an earlier time step given the current state of the system
- This means estimating the  $q(\mathbf{x}_t|\mathbf{x}_{t-1})$  when  $t = T$
- Thereby generating data sample from isotropic Gaussian noise
- The estimation of previous state from the current state requires the **knowledge of all the previous gradients** which we can't obtain without having a learning model that can predict such estimates

## Conditional probability density:

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

## Complete distribution:

$$q(\mathbf{x}_{0:T}|\mathbf{x}_0) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

## Re-formulation:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

# Denoising Diffusion Probability Model

## Reverse step:

- The estimation of previous state from the current state requires the **knowledge of all the previous gradients** which we can't obtain without having a learning model that can predict such estimates
- Solution
  - Train a neural network that estimates the  $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$  based on **learned weights**  $\theta$  and the current state at time  $t$
  - Formulated as:

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

$$p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

## Conditional probability density:

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

## Complete distribution:

$$q(\mathbf{x}_{0:T}|\mathbf{x}_0) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

## Re-formulation:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

# Denoising Diffusion Probability Model

Reverse step:

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

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

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

# Denoising Diffusion Probability Model

$$\alpha_t = 1 - \beta_t, \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

Reverse step:

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

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}\right)\mathbf{x}_{t-1}^2 - \left(\frac{2\sqrt{\alpha_t}}{\beta_t}\mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}}\mathbf{x}_0\right)\mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0)\right)\right)$$

$$\tilde{\beta}_t = 1/\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}\right) = 1/\left(\frac{\alpha_t - \bar{\alpha}_t + \beta_t}{\beta_t(1 - \bar{\alpha}_{t-1})}\right) = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t$$

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \left(\frac{2\sqrt{\alpha_t}}{\beta_t}\mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}}\mathbf{x}_0\right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}\right)$$

$$= \left(\frac{2\sqrt{\alpha_t}}{\beta_t}\mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}}\mathbf{x}_0\right) \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\mathbf{x}_0 = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}_t\right)$$

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_t \\ \mathbf{x}_0 &= \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_t)\end{aligned}$$

# Denoising Diffusion Probability Model

Reverse step:

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

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right)$$

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right)$$



# Denoising Diffusion Probability Model

## Reverse step:

- The parameterization for:

- Mean function:

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right)$$

- Variance function (fixed as  $\boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$ ,  $\sigma_t^2 = \beta_t$  or  $\sigma_t^2 = \tilde{\beta}_t$ ):

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$$

## Current state estimation:

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

# Denoising Diffusion Probability Model

**Loss:**

$$\begin{aligned} -\log p_{\theta} &\leq -\log p_{\theta}(\mathbf{x}_0) + D_{KL} \left( (q(\mathbf{x}_{1:T}|\mathbf{x}_0) || p_{\theta}(\mathbf{x}_{1:T}|\mathbf{x}_0)) \right) \\ &= -\log p_{\theta}(\mathbf{x}_0) + \mathbb{E}_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})/p_{\theta}(\mathbf{x}_0)} \right] \\ &= -\log p_{\theta}(\mathbf{x}_0) + \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} + \log p_{\theta}(\mathbf{x}_0) \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\ &= L_{VLB} \end{aligned}$$

# Denoising Diffusion Probability Model

**Loss:**

$$\begin{aligned} L_{VLB} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_q \left[ \log \frac{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} + \log \frac{q(\mathbf{x}_1 | \mathbf{x}_0)}{p_\theta(\mathbf{x}_0 | \mathbf{x}_1)} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \left( \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} \cdot \frac{q(\mathbf{x}_t | \mathbf{x}_0)}{q(\mathbf{x}_{t-1} | \mathbf{x}_0)} \right) + \log \frac{q(\mathbf{x}_1 | \mathbf{x}_0)}{p_\theta(\mathbf{x}_0 | \mathbf{x}_1)} \right] \end{aligned}$$

# Denoising Diffusion Probability Model

**Loss:**

$$\begin{aligned} L_{VLB} &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \left( \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \cdot \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right) + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \\ &= \mathbb{E}_q \left[ D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)||p_\theta(\mathbf{x}_T)) + \sum_{t=2}^T D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \end{aligned}$$

# Denoising Diffusion Probability Model

**Loss:**

$$L_{VLB} = \mathbb{E}_q \left[ \underbrace{D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)||p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_t} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right]_{L_0}$$

$$L_{VLB} = L_T + L_{T-1} + \cdots + L_0$$

# Denoising Diffusion Probability Model

Loss:

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

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right)$$

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right)$$

$$L_t = D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \frac{1}{2 \|\boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)\|^2} \|\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)\|^2$$

$$L_t^{simple} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \boldsymbol{\epsilon}} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\|^2 = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \boldsymbol{\epsilon}} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\|^2$$



# Denoising Diffusion Probability Model

Algorithm:

---

## 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  $\mathbf{x}_0$ 
```

---

# Denoising Diffusion Probability Model

Generation Ability:

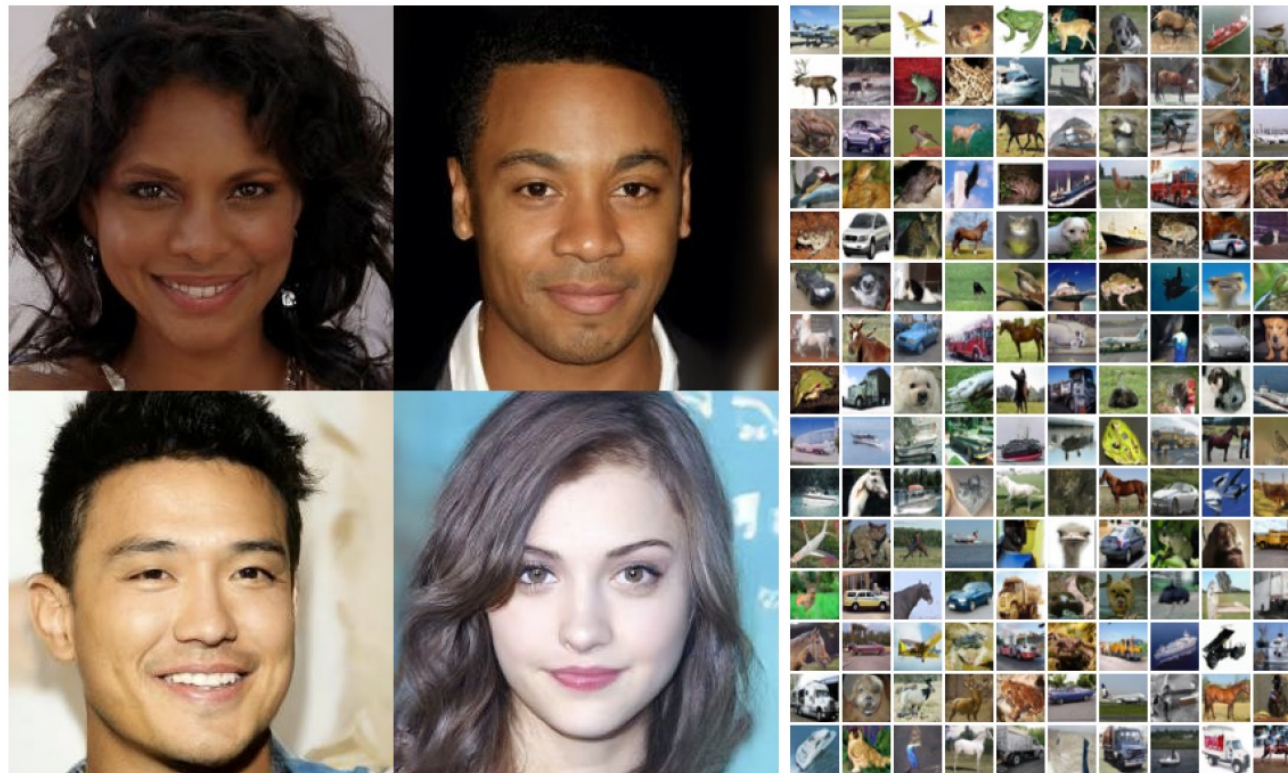


Figure 1: Generated samples on CelebA-HQ  $256 \times 256$  (left) and unconditional CIFAR10 (right)

# Denoising Diffusion Probability Model

## Generation Ability:

Table 1: CIFAR10 results. NLL measured in bits/dim.

Model	IS	FID	NLL Test (Train)
<b>Conditional</b>			
EBM [11]	8.30	37.9	
JEM [17]	8.76	38.4	
BigGAN [3]	9.22	14.73	
StyleGAN2 + ADA (v1) [29]	<b>10.06</b>	<b>2.67</b>	
<b>Unconditional</b>			
Diffusion (original) [53]			$\leq 5.40$
Gated PixelCNN [59]	4.60	65.93	3.03 (2.90)
Sparse Transformer [7]			<b>2.80</b>
PixelIQN [43]	5.29	49.46	
EBM [11]	6.78	38.2	
NCSNv2 [56]		31.75	
NCSN [55]	$8.87 \pm 0.12$	25.32	
SNGAN [39]	$8.22 \pm 0.05$	21.7	
SNGAN-DDLS [4]	$9.09 \pm 0.10$	15.42	
StyleGAN2 + ADA (v1) [29]	$9.74 \pm 0.05$	3.26	
Ours ( $L$ , fixed isotropic $\Sigma$ )	$7.67 \pm 0.13$	13.51	$\leq 3.70$ (3.69)
<b>Ours (<math>L_{\text{simple}}</math>)</b>	$9.46 \pm 0.11$	<b>3.17</b>	$\leq 3.75$ (3.72)

# Denoising Diffusion Probability Model

Beat GAN:

Model	FID	sFID	Prec	Rec
<b>LSUN Bedrooms 256×256</b>				
DCTransformer <sup>†</sup> [42]	6.40	6.66	0.44	<b>0.56</b>
DDPM [25]	4.89	9.07	0.60	0.45
IDDPM [43]	4.24	8.21	0.62	0.46
StyleGAN [27]	2.35	6.62	0.59	0.48
ADM (dropout)	<b>1.90</b>	<b>5.59</b>	<b>0.66</b>	0.51

<b>LSUN Horses 256×256</b>				
StyleGAN2 [28]	3.84	6.46	0.63	0.48
ADM	2.95	<b>5.94</b>	0.69	<b>0.55</b>
ADM (dropout)	<b>2.57</b>	6.81	<b>0.71</b>	<b>0.55</b>

<b>LSUN Cats 256×256</b>				
DDPM [25]	17.1	12.4	0.53	0.48
StyleGAN2 [28]	7.25	<b>6.33</b>	0.58	0.43
ADM (dropout)	<b>5.57</b>	6.69	<b>0.63</b>	<b>0.52</b>

<b>ImageNet 64×64</b>				
BigGAN-deep* [5]	4.06	3.96	<b>0.79</b>	0.48
IDDPM [43]	2.92	<b>3.79</b>	0.74	0.62
ADM	2.61	<b>3.77</b>	0.73	0.63
ADM (dropout)	<b>2.07</b>	4.29	0.74	<b>0.63</b>

Model	FID	sFID	Prec	Rec
<b>ImageNet 128×128</b>				
BigGAN-deep [5]	6.02	7.18	<b>0.86</b>	0.35
LOGAN <sup>†</sup> [68]	3.36			
ADM	5.91	<b>5.09</b>	0.70	<b>0.65</b>
ADM-G (25 steps)	5.98	7.04	0.78	0.51
ADM-G	<b>2.97</b>	<b>5.09</b>	0.78	0.59

<b>ImageNet 256×256</b>				
DCTransformer <sup>†</sup> [42]	36.51	8.24	0.36	<b>0.67</b>
VQ-VAE-2 <sup>†‡</sup> [51]	31.11	17.38	0.36	0.57
IDDPM <sup>‡</sup> [43]	12.26	5.42	0.70	0.62
SR3 <sup>†‡</sup> [53]	11.30			
BigGAN-deep [5]	6.95	7.36	<b>0.87</b>	0.28
ADM	10.94	6.02	0.69	0.63
ADM-G (25 steps)	5.44	5.32	0.81	0.49
ADM-G	<b>4.59</b>	<b>5.25</b>	0.82	0.52

<b>ImageNet 512×512</b>				
BigGAN-deep [5]	8.43	8.13	<b>0.88</b>	0.29
ADM	23.24	10.19	0.73	<b>0.60</b>
ADM-G (25 steps)	8.41	9.67	0.83	0.47
ADM-G	<b>7.72</b>	<b>6.57</b>	0.87	0.42

# Reference

- [What are Diffusion Models? | Lil'Log \(lilianweng.github.io\)](https://lilianweng.github.io/lil-log/)
- <https://medium.com/@marija.jegorova/a-quick-dive-into-diffusion-models-part-1-fundamentals-of-the-ddpms-4de1304d7b6f>
- <https://towardsdatascience.com/image-generation-in-10-minutes-with-generative-adversarial-networks-c2afc56bfa3b>
- Sohl-Dickstein, et al. "Deep unsupervised learning using nonequilibrium thermodynamics." International conference on machine learning. PMLR, 2015.
- Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." NIPS 2020



**Q&A**







# Thanks

## Denoising Diffusion Probability Model

**Jonathan Ho**  
UC Berkeley  
jonathanho@berkeley.edu

**Ajay Jain**  
UC Berkeley  
ajayj@berkeley.edu

**Pieter Abbeel**  
UC Berkeley  
pabbeel@cs.berkeley.edu

NeurIPS, 2020

Miao's Group - Paper Reading



 向乾龙

 时间 : 2023.11.28

