

Deep Generative Models: Diffusion Models

Fall Semester 2025

René Vidal

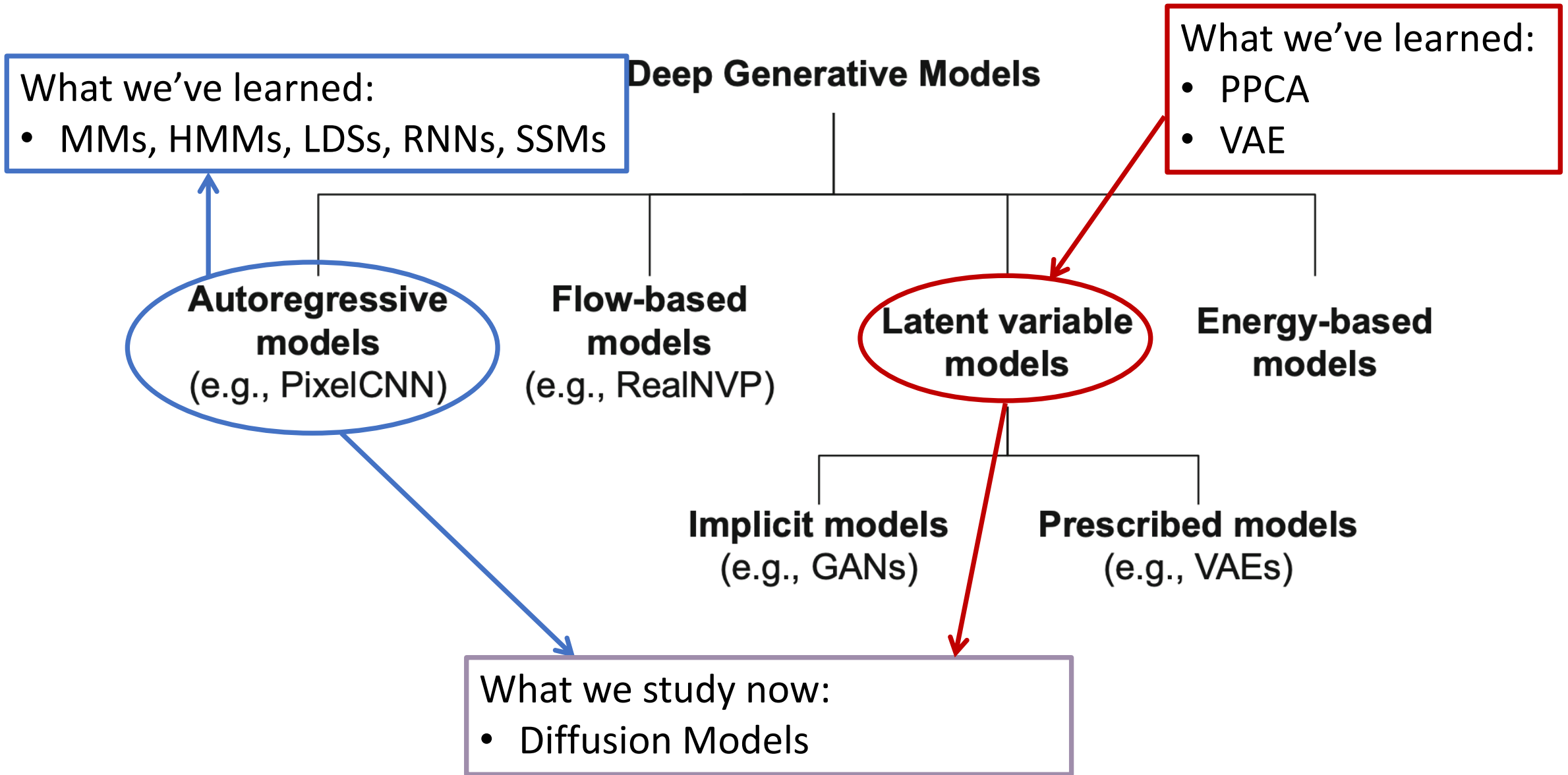
Director of the Center for Innovation in Data Engineering and Science (IDEAS)

Rachleff University Professor, University of Pennsylvania

Amazon Scholar & Chief Scientist at NORCE



Taxonomy of Generative Models

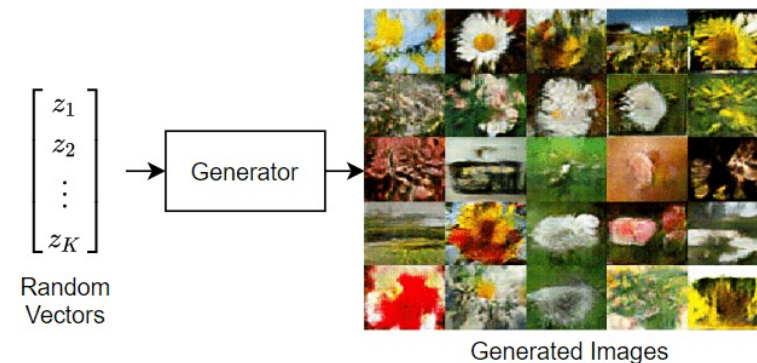
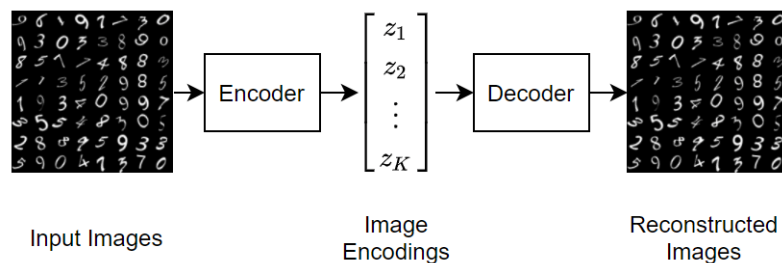


Diffusion Models

- Derivation of Diffusion Models (Today's Lecture)
 - Markov Hierarchical Variational Auto Encoders (MHVAE)
 - Diffusion Models are VAEs with Linear Gaussian Autoregressive latent space
 - Forward Process
 - Conditional Distributions for the Forward Process
 - Reverse Process
 - ELBO for Diffusion Models is a particular case of ELBO for VAEs with extra structure
 - Implementation Details
- Application of Diffusion Models (Next Lecture)
 - Stable Diffusion: Text-Conditioned Diffusion Model
 - ControlNet: Multimodal Control for Consistent Synthesis

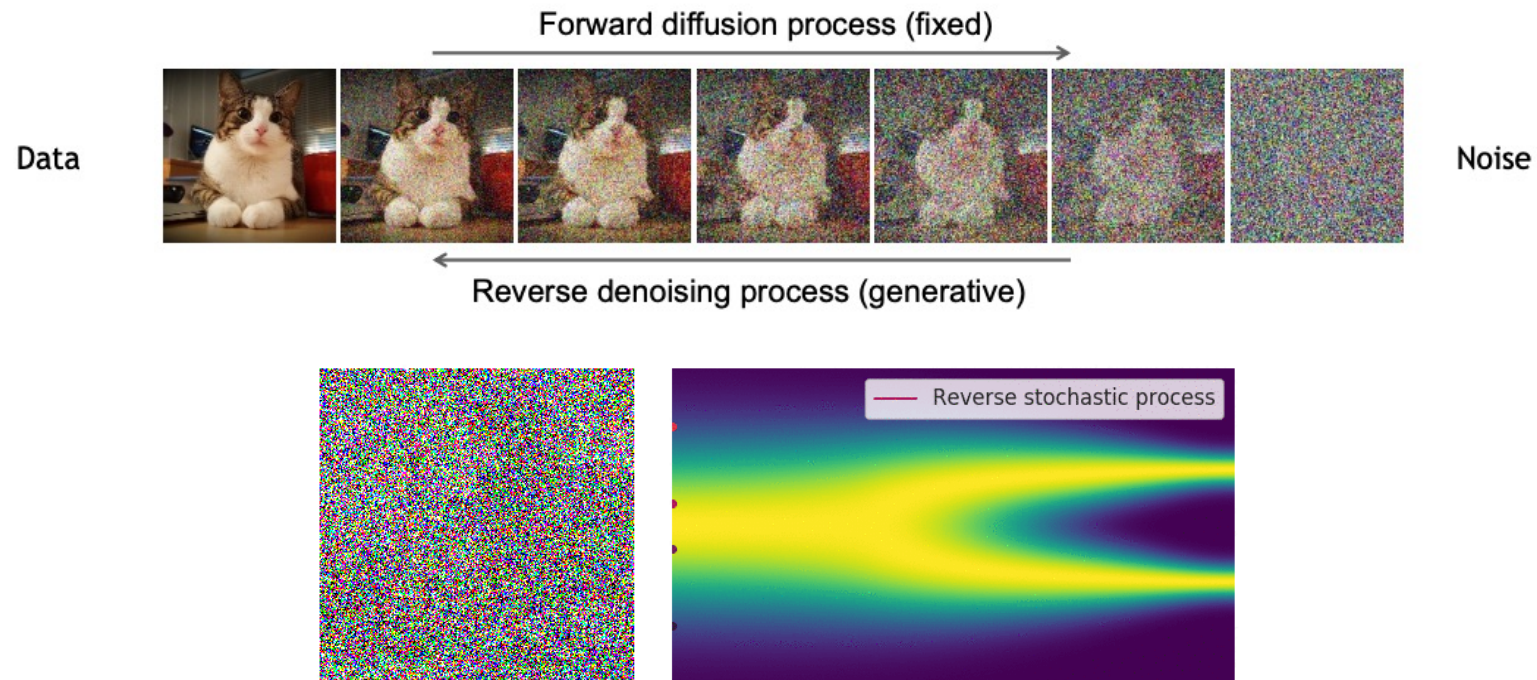
Diffusion Models

- The journey of generative models has evolved significantly in recent years.
- **Variational Autoencoders (VAEs)** introduce probabilistic modeling for latent representations but struggled with generating high-quality images.
- This led to the rise of **Generative Adversarial Networks (GANs)**, which leverage adversarial learning to produce high-quality, realistic outputs but suffered from issues like mode collapse and unstable training.
- The introduction of **Diffusion Models** achieve state-of-the-art results with superior stability and diversity in generated samples, particularly in multimodal image synthesis.

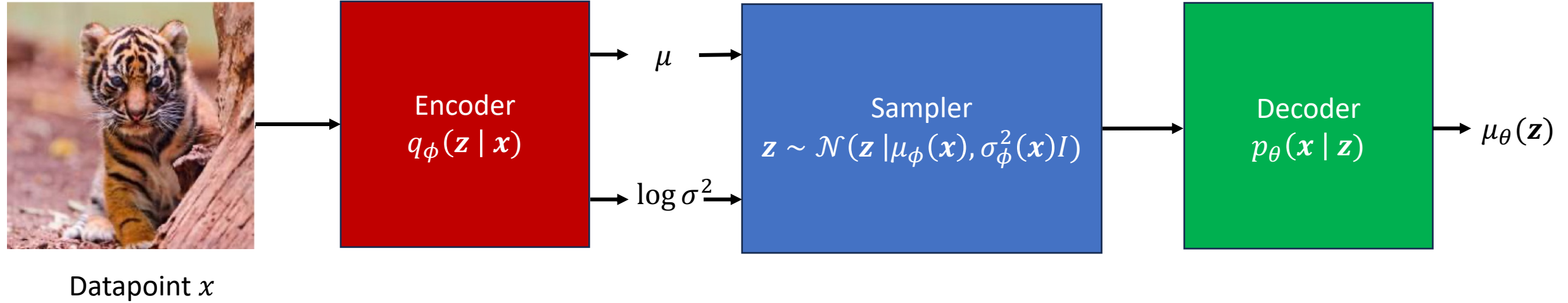


Diffusion Models

- A Latent Diffusion Model is a VAE with an autoregressive latent space
- The VAE encoder **maps data to noise** by gradually adding Gaussian noise to the input using a **diffusion process**.
- The VAE decoder **maps noise to data** by **learning how to reverse the forward diffusion process**. The reverse process predicts how to denoise.



Recall the Variational Autoencoder (VAE)



ELBO Objective

$$\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z}) - KL(q_\phi(\mathbf{z} | \mathbf{x}) \parallel p(\mathbf{z}))]$$

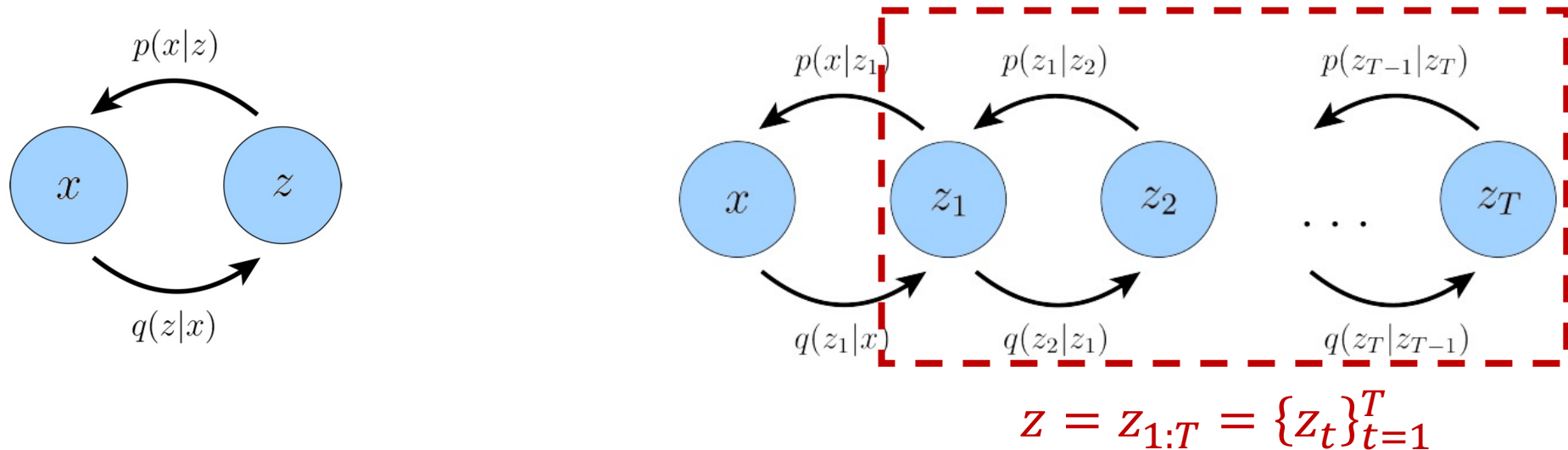
Recall the Evidence Lower Bound (ELBO)

- The ELBO is the sum of a reconstruction term and a prior matching term

$$\begin{aligned}\log p_{\theta}(x) &\geq \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(z)}{q_{\phi}(z|x)} \right] \\ &= \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}}\end{aligned}$$

Latent Diffusion Models as “Autoregressive VAEs”

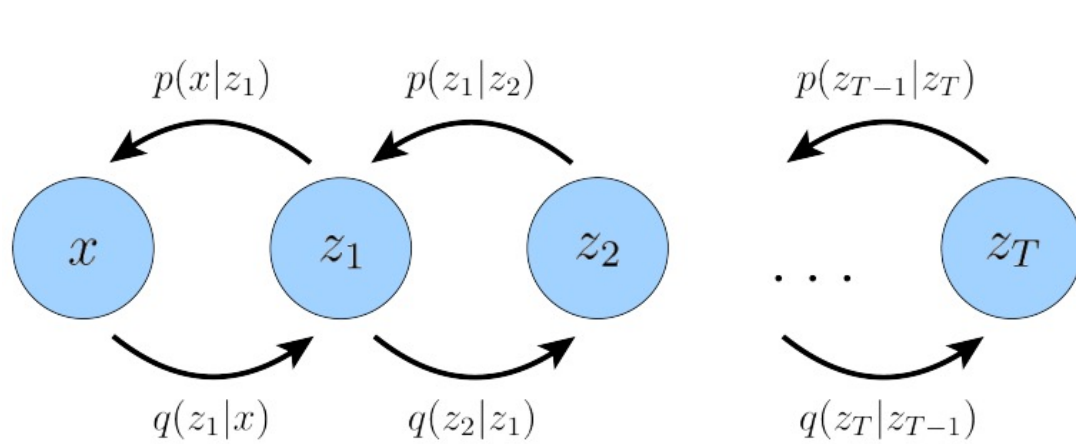
- A Latent Diffusion Model is as a **Markovian Hierarchical Variational Autoencoder (MHVAE)** with T hierarchical latents $\mathbf{z} = \mathbf{z}_{1:T} = \{z_t\}_{t=1}^T$ modeled by a Markov chain where each latent z_t is generated only from the previous latent z_{t+1} .



- What is the VAE encoder $q(\mathbf{z} | \mathbf{x})$ of a Diffusion Model ?
- What is the VAE decoder $p(\mathbf{x} | \mathbf{z})$ of a Diffusion Model ?
- What is the ELBO of a Diffusion Model ?

MHVAE Encoder, Decoder, and ELBO

- A MHVAE is a VAE whose encoder and decoder are autoregressive models:



$$p_{\theta}(x, z_{1:T}) = p_{\theta}(z_T) p_{\theta}(x | z_1) \prod_{t=2}^T p_{\theta}(z_{t-1} | z_t)$$

$$q_{\phi}(z_{1:T} | x) = q_{\phi}(z_1 | x) \prod_{t=2}^T q_{\phi}(z_t | z_{t-1})$$

- Given this joint distribution and posterior, we can rewrite the ELBO for MHVAE as:

$$\mathbb{E}_{q_{\phi}(z_{1:T}|x)} \left[\log \frac{p_{\theta}(x, z_{1:T})}{q_{\phi}(z_{1:T} | x)} \right] = \mathbb{E}_{q_{\phi}(z_{1:T}|x)} \left[\log \frac{p_{\theta}(z_T) p_{\theta}(x | z_1) \prod_{t=2}^T p_{\theta}(z_{t-1} | z_t)}{q_{\phi}(z_1 | x) \prod_{t=2}^T q_{\phi}(z_t | z_{t-1})} \right]$$

Decomposition of the ELBO for a MHVAE

- Let us make the change of variables $x \rightarrow x_0$ and $z_{1:T} \rightarrow x_{1:T}$.

- **Theorem:** The ELBO for a MHVAE can be written as

$$\mathbb{E}_{q_\phi(x_{1:T}|x)} \left[\log \frac{p_\theta(x_T) p_\theta(x_0 | x_1) \prod_{t=2}^T p_\theta(x_{t-1} | x_t)}{q_\phi(x_1 | x_0) \prod_{t=2}^T q_\phi(x_t | x_{t-1})} \right] =$$
$$\underbrace{\mathbb{E}_{q_\phi(x_1|x_0)} [\log p_\theta(x_0 | x_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(x_T | x_0) \parallel p_\theta(x_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q_\phi(x_t|x_0)} \left[D_{\text{KL}}(q_\phi(x_{t-1} | x_t, x_0) \parallel p_\theta(x_{t-1} | x_t)) \right]}_{\text{score matching term}}$$

- **Proof (1/3):** Reversing q

$$q_\phi(x_t | x_{t-1}) = q_\phi(x_t | x_{t-1}, x_0) = \frac{q_\phi(x_{t-1}|x_t, x_0) q_\phi(x_t|x_0)}{q_\phi(x_{t-1}|x_0)}.$$

Decomposition of the ELBO for a MHVAE

- Proof (2/3): substituting q and using telescopic product to cancel factors

$$\begin{aligned}\log p(x) &\geq \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_T) p_\theta(x_0 | x_1) \prod_{t=2}^T p_\theta(x_{t-1} | x_t)}{q_\phi(x_1 | x_0) \prod_{t=2}^T q_\phi(x_t | x_{t-1})} \right] \\&= \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_T) p_\theta(x_0 | x_1)}{q_\phi(x_1 | x_0)} \prod_{t=2}^T \frac{p_\theta(x_{t-1} | x_t)}{\frac{q_\phi(x_{t-1} | x_t, x_0) q_\phi(x_t | x_0)}{q_\phi(x_{t-1} | x_0)}} \right] \\&= \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_T) p_\theta(x_0 | x_1)}{q_\phi(x_1 | x_0)} \prod_{t=2}^T \frac{p_\theta(x_{t-1} | x_t) q_\phi(x_1 | x_0)}{q_\phi(x_{t-1} | x_t, x_0) q_\phi(x_T | x_0)} \right] \\&= \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_T) p_\theta(x_0 | x_1) q_\phi(x_1 | x_0)}{q_\phi(x_1 | x_0) q_\phi(x_T | x_0)} \prod_{t=2}^T \frac{p_\theta(x_{t-1} | x_t)}{q_\phi(x_{t-1} | x_t, x_0)} \right]\end{aligned}$$

Decomposition of the ELBO for a MHVAE

- Proof (3/3): expanding into three terms and simplifying expectations

$$\begin{aligned}\log p(x) &\geq \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_T) p_\theta(x_0 | x_1)}{q_\phi(x_T | x_0)} \prod_{t=2}^T \frac{p_\theta(x_{t-1} | x_t)}{q_\phi(x_{t-1} | x_t, x_0)} \right] \\&= \mathbb{E}_{q_\phi(x_{1:T}|x_0)} [\log p_\theta(x_0 | x_1)] + \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_T)}{q_\phi(x_T | x_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_{t-1} | x_t)}{q_\phi(x_{t-1} | x_t, x_0)} \right] \\&= \mathbb{E}_{q_\phi(x_1|x_0)} [\log p_\theta(x_0 | x_1)] + \mathbb{E}_{q_\phi(x_T|x_0)} \left[\log \frac{p_\theta(x_T)}{q_\phi(x_T | x_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q_\phi(x_t|x_0)} \left[\log \frac{p_\theta(x_{t-1} | x_t)}{q_\phi(x_{t-1} | x_t, x_0)} \right] \\&= \underbrace{\mathbb{E}_{q_\phi(x_1|x_0)} [\log p_\theta(x_0 | x_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(x_T | x_0) || p_\theta(x_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q_\phi(x_t|x_0)} [D_{\text{KL}}(q_\phi(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t))]}_{\text{score matching term}}\end{aligned}$$

Why can we simplify expectations?

$$\int p(x_{2:T} | x_1) dx_{2:T} = 1$$

- For the first term:

$$\begin{aligned}\mathbb{E}_{q_\phi(x_{1:T}|x_0)}[\log(p_\theta(x_0 | x_1))] &= \int \log(p_\theta(x_0 | x_1)) q_\phi(x_{1:T} | x_0) dx_{1:T} \\ &= \int \log(p_\theta(x_0 | x_1)) q_\phi(x_1 | x_0) q_\phi(x_{2:T} | x_1) dx_1 dx_{2:T} \\ &= \int \log p_\theta(x_0 | x_1) q_\phi(x_1 | x_0) dx_1 = \mathbb{E}_{q_\phi(x_1|x_0)}[\log p_\theta(x_0 | x_1)]\end{aligned}$$

- Similarly for the second term:

$$\mathbb{E}_{q_\phi(x_{1:T}|x_0)}\left[\log \frac{p_\theta(x_T)}{q_\phi(x_T | x_0)}\right] = \int \log\left(\frac{p_\theta(x_T)}{q_\phi(x_T | x_0)}\right) q_\phi(x_{1:T} | x_0) dx_{1:T} = \int \log\left(\frac{p_\theta(x_T)}{q_\phi(x_T | x_0)}\right) q_\phi(x_T | x_0) q_\phi(x_{1:T-1} | x_T, x_0) dx_{1:T}$$

ELBO for MHVAE

$$\int p(x_{2:T}) dx_{2:T} = 1$$

- Why can we simplify expectations? (1/2)

$$\mathbb{E}_{q_{\phi}(x_{1:T}|x_0)} \left[\log \left(\frac{p_{\theta}(x_{t-1} | x_t)}{q_{\phi}(x_{t-1} | x_t, x_0)} \right) \right] = \int \log \left(\frac{p_{\theta}(x_{t-1} | x_t)}{q_{\phi}(x_{t-1} | x_t, x_0)} \right) q_{\phi}(x_{1:T} | x_0) dx_{1:T}$$

$$= \int \log \left(\frac{p_{\theta}(x_{t-1}|x_t)}{q_{\phi}(x_{t-1}|x_t, x_0)} \right) \prod_{\tau=1}^T q_{\phi}(x_{\tau} | x_{\tau-1}) d\mathbf{x}_{1:T} \quad (\text{Markov property})$$

$$= \int \log \left(\frac{p_{\theta}(x_{t-1} | x_t)}{q_{\phi}(x_{t-1} | x_t, x_0)} \right) q_{\phi}(x_t | x_{t-1}) \prod_{\tau \neq t}^T q_{\phi}(x_{\tau} | x_{\tau-1}) d\mathbf{x}_{1:T}$$

The log term depends only on x_t, x_{t-1}, x_0 and thus marginalizing q_{ϕ} with respect to $d\mathbf{x}_{1:t-2}$

$$\int \prod_{\tau=1}^{t-2} q_{\phi}(x_{\tau} | x_{\tau-1}) d\mathbf{x}_{1:t-2} = q_{\phi}(x_{t-1} | x_0)$$

ELBO for MHVAE

$$\int p(x_{2:T}) dx_{2:T} = 1$$

- Why can we simplify expectations? (2/2)

$$\begin{aligned} \mathbb{E}_{q_\phi(x_{1:T}|x_0)} \left[\log \left(\frac{p_\theta(x_{t-1}|x_t)}{q_\phi(x_{t-1}|x_t, x_0)} \right) \right] &= \\ &= \int \log \left(\frac{p_\theta(x_{t-1} | x_t)}{q_\phi(x_{t-1} | x_t, x_0)} \right) q_\phi(x_t | x_{t-1}) q_\phi(x_{t-1} | x_0) dx_t dx_{t-1} \end{aligned}$$

Using the fact that $q_\phi(x_t | x_{t-1}) q_\phi(x_{t-1} | x_0) = q_\phi(x_{t-1}, x_t | x_0)$, it holds

$$\begin{aligned} &= \int \log \left(\frac{p_\theta(x_{t-1} | x_t)}{q_\phi(x_{t-1} | x_t, x_0)} \right) q_\phi(x_{t-1}, x_t | x_0) dx_t dx_{t-1} \\ &= \mathbb{E}_{q_\phi(x_t | x_0)} \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q_\phi(x_{t-1}|x_t, x_0)} \right] \\ &= \mathbb{E}_{q_\phi(x_t | x_0)} \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q_\phi(x_{t-1}|x_t, x_0)} \right] \end{aligned}$$

Interpretation of the ELBO for MHVAE

$$\begin{aligned} \log p(x) \\ \geq \underbrace{\mathbb{E}_{q_\phi(x_1|x_0)}[\log p_\theta(x_0 | x_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(x_T | x_0) || p_\theta(x_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q_\phi(x_t|x_0)} \left[D_{\text{KL}}(q_\phi(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t)) \right]}_{\text{score matching term}} \end{aligned}$$

$\mathbb{E}_{q_\phi(x_1|x_0)}[\log p_\theta(x_0 | x_1)]$ can be interpreted as a **reconstruction term**; like its analogue in the ELBO of a vanilla VAE. This term can be approximated and optimized using a Monte Carlo estimate.

$D_{\text{KL}}(q_\phi(x_T | x_0) || p_\theta(x_T))$ represents how **close the distribution of the final latent distribution is to the standard Gaussian prior**.

$\mathbb{E}_{q_\phi(x_t|x_0)} \left[D_{\text{KL}}(q_\phi(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t)) \right]$ is a **score matching term**.

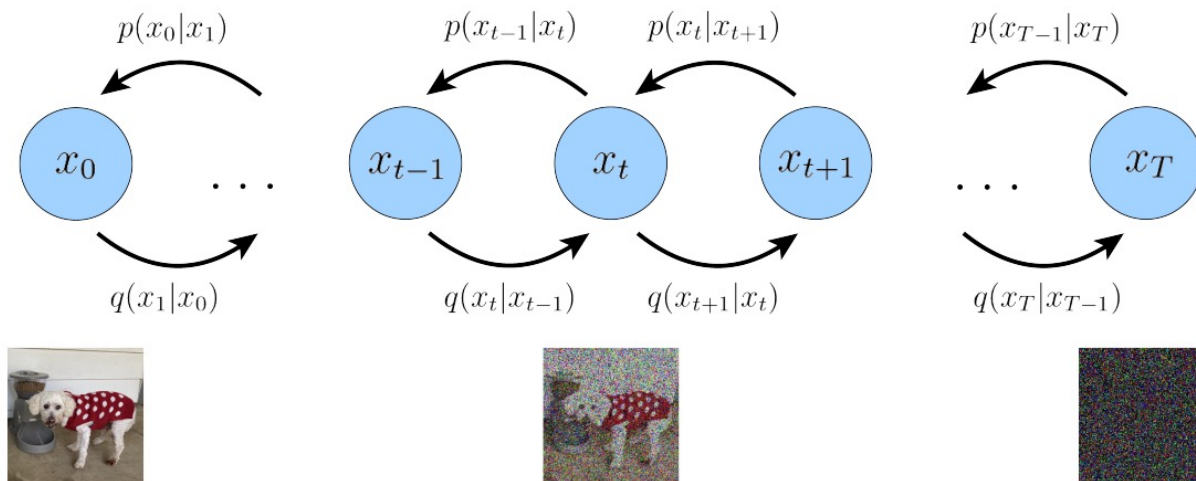
As we will show, the diffusion model learns the denoising transition step $p_\theta(x_{t-1} | x_t)$ as an approximation to the tractable, ground-truth denoising transition step $q_\phi(x_{t-1} | x_t, x_0)$.

Understand Diffusion Model from VAE Perspective

- **A Diffusion Model is an MHVAE:** $x_0 = x$ is the data and $x_{1:T} = z_{1:T}$ is the latent variable
- All latent variables have the same dimension as the dimension of the data
- The structure of the encoder $q_\phi(x_{1:T} | x_0) = \prod_{t=1}^T q_\phi(x_t | x_{t-1})$ is not learned, but it is pre-specified as a linear Gaussian model

$$q_\phi(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

- The parameter α_t is chosen such that $x_T \sim \mathcal{N}(x_T; 0, I)$ is a standard Gaussian



The Forward Process of Diffusion Model

Given the formulation of a single noising step:

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon; 0, I),$$

we can **recursively** derive the closed form for arbitrary noising steps:

$$\begin{aligned}\mathbb{E}[x_t \mid x_0] &= \mathbb{E}[\sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_t \mid x_0] \\ &= \sqrt{\alpha_t}\mathbb{E}[x_{t-1} \mid x_0] + \sqrt{1 - \alpha_t}\mathbb{E}[\epsilon_t] \\ &= \sqrt{\alpha_t}\mathbb{E}[x_{t-1} \mid x_0]\end{aligned}$$

That is:

$$\mathbb{E}[x_t \mid x_0] = \sqrt{\alpha_t}\sqrt{\alpha_{t-1}}\mathbb{E}[x_{t-2} \mid x_0] = \cdots = \sqrt{\bar{\alpha}_t}x_0$$

The Forward Process of Diffusion Model

The variance is given by

$$\begin{aligned}\text{Var}(x_t \mid x_0) &= \text{Var}(\sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_t \mid x_0) \\ &= \alpha_t \text{Var}(x_{t-1} \mid x_0) + (1 - \alpha_t)\text{Var}(\varepsilon_t) \\ &= \alpha_t \text{Var}(x_{t-1} \mid x_0) + (1 - \alpha_t)I\end{aligned}$$

That is:

$$\begin{aligned}\text{Var}(x_t \mid x_0) &= (1 - \alpha_t)I + \alpha_t(1 - \alpha_{t-1})I + \alpha_t\alpha_{t-1}(1 - \alpha_{t-2})I + \dots \\ &= (1 - \prod_{s=1}^t \alpha_s)I \\ &= (1 - \overline{\alpha_t})I\end{aligned}$$

The Forward Process of Diffusion Model

To summarize: $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t} \epsilon, \epsilon \sim \mathcal{N}(\epsilon; 0, I)$

That is: $x_0 = \frac{x_t - \sqrt{1 - \alpha_t} \epsilon}{\sqrt{\alpha_t}}$

The forward diffusion process can be seen as a paradigm that x_t is a linear Gaussian transformation of x_0 with scheduled randomness from a standard normal distribution.

We will use this for the reparameterization trick later.

ELBO for Diffusion Model: Score Matching Term

- To compute the third term, we need $q(x_{t-1} | x_t, x_0) = \frac{q(x_t | x_{t-1}, x_0)q(x_{t-1} | x_0)}{q(x_t | x_0)}$

- Letting $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, recall that $q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

- Therefore

$$\begin{aligned} q(x_{t-1} | x_t, x_0) &= \frac{q(x_t | x_{t-1}, x_0)q(x_{t-1} | x_0)}{q(x_t | x_0)} \\ &= \frac{\mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)\mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})I)}{\mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)} \\ &\propto \mathcal{N}\left(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \bar{\alpha}_{t-1}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}I}_{\Sigma_q(t)}\right) \end{aligned}$$

ELBO for Diffusion Model: Matching the Mean

- Recall KL divergence for Gaussians

$$D_{\text{KL}} \left(\mathcal{N}(x; \mu_x, \Sigma_x) \middle| \mathcal{N}(y; \mu_y, \Sigma_y) \right) = \frac{1}{2} \left[\log \frac{|\Sigma_y|}{|\Sigma_x|} - d + \text{tr}(\Sigma_y^{-1} \Sigma_x) + (\mu_y - \mu_x)^T \Sigma_y^{-1} (\mu_y - \mu_x) \right]$$

- Choose variance of p to match exactly variance of q

$$\begin{aligned} & D_{\text{KL}}(q(x_{t-1} | x_t, x_0) | p_{\theta}(x_{t-1} | x_t)) \\ &= D_{\text{KL}} \left(\mathcal{N}(x_{t-1}; \mu_q, \Sigma_q(t)) \middle| \mathcal{N}(x_{t-1}; \mu_{\theta}, \Sigma_q(t)) \right) \end{aligned}$$

$$= \frac{1}{2\sigma_q^2(t)} [|\mu_{\theta} - \mu_q|_2^2] = \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} [|\widehat{x}_{\theta}(x_t, t) - x_0|_2^2]$$

$$\sigma_q^2(t) = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}$$

- Choose mean of p to match form of mean of q

$$\mu_{\theta}(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\widehat{x}_{\theta}(x_t, t)}{1 - \bar{\alpha}_t} \quad \mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}$$

Reparameterization as an Alternative Form for ELBO

- Plugging our previous finding $\mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0}{\sqrt{\alpha_t}}$ into the denoising transition mean

$\mu_q(\mathbf{x}_t, \mathbf{x}_0)$, we have:

$$\begin{aligned} \mu_q(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \\ &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0}{\sqrt{\alpha_t}}}{1 - \bar{\alpha}_t} \\ &= \frac{1 - \bar{\alpha}_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \epsilon_0 \\ &= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \boxed{\frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \epsilon_0} \end{aligned}$$

- This inspires us to approximate the denoising transition mean as **choosing the mean of p to match q** : $\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \boxed{\frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \hat{\epsilon}_\theta(\mathbf{x}_t, t)}$

Progressive Denoising or Direct Reconstruction?

- The model predicts the noise to be removed in each step (i.e., denoising) by optimizing **score matching term**. This reduces to minimizing the difference between the predicted noise and the ground-truth schedule noise:

$$\begin{aligned} & \underset{\theta}{\operatorname{argmin}} D_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) \\ &= \underset{\theta}{\operatorname{argmin}} D_{\text{KL}}\left(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))\right) \\ &= \underset{\theta}{\operatorname{argmin}} D_{\text{KL}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) - \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0 \right\|_2^2 \right] \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} [\|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)\|_2^2] \end{aligned}$$



- Predicting \mathbf{x}_0 from a highly noisy \mathbf{x}_t in one step is complex because the signal is buried under significant noise, especially at large t . By predicting the noise at each step, the model progressively refines \mathbf{x}_t towards \mathbf{x}_0 , which makes the learning task more manageable (e.g., converges better / requires smaller network capacity).

Training and Sampling from Diffusion Model

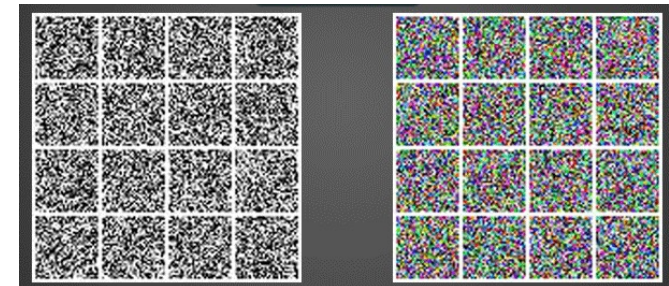
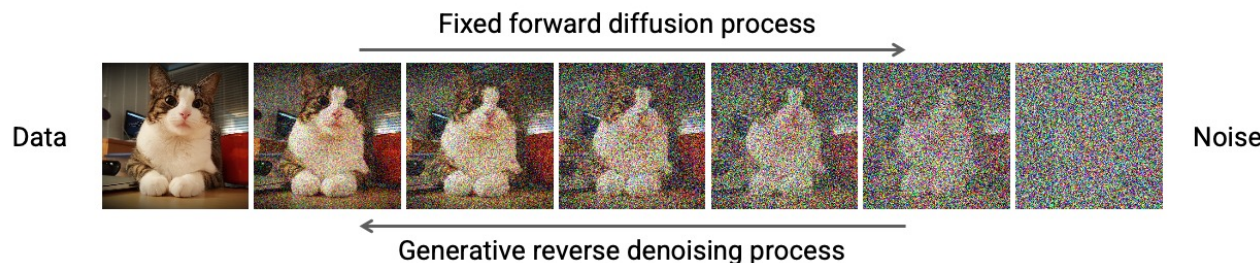
- [Ho et al., 2020] (DDPM) chooses to build the training procedure by performing SGD on the set of training images over timesteps.
- The sampling procedure iteratively executes the denoising process from a Gaussian initialization \mathbf{x}_T .

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  
        $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta} \left( \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2$   
6: until converged
```

Algorithm 2 Sampling

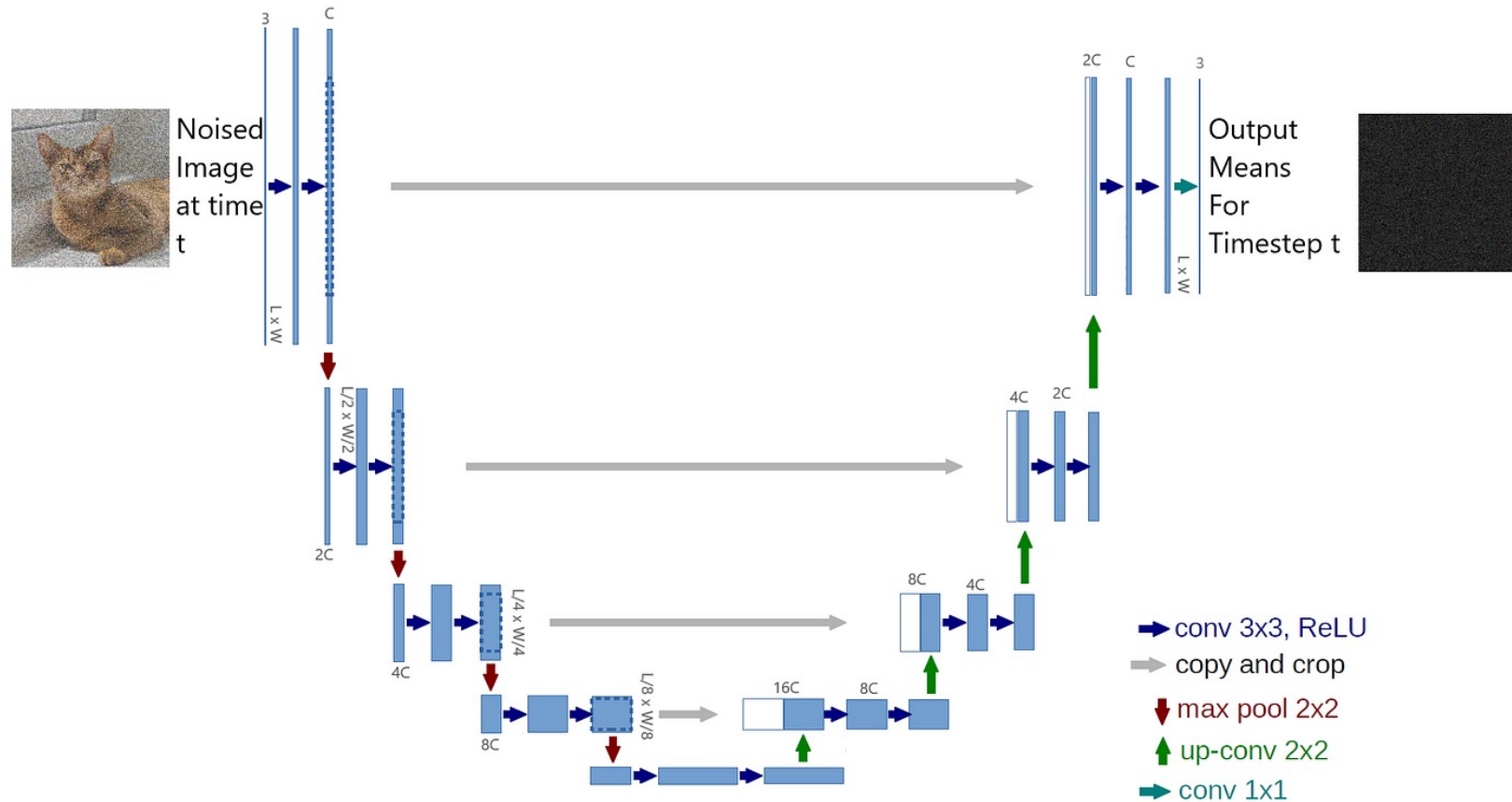
```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
2: for  $t = T, \dots, 1$  do  
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$   
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1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



Implementation (DDPM)

DDPM uses U-Net with residual connection and self-attention layers to represent $\epsilon_{\theta}(\mathbf{x}_t, t)$.

The time representation is conditioned in the U-Net as **sinusoidal positional embeddings** or **Fourier features**.



Implementation (DDPM)

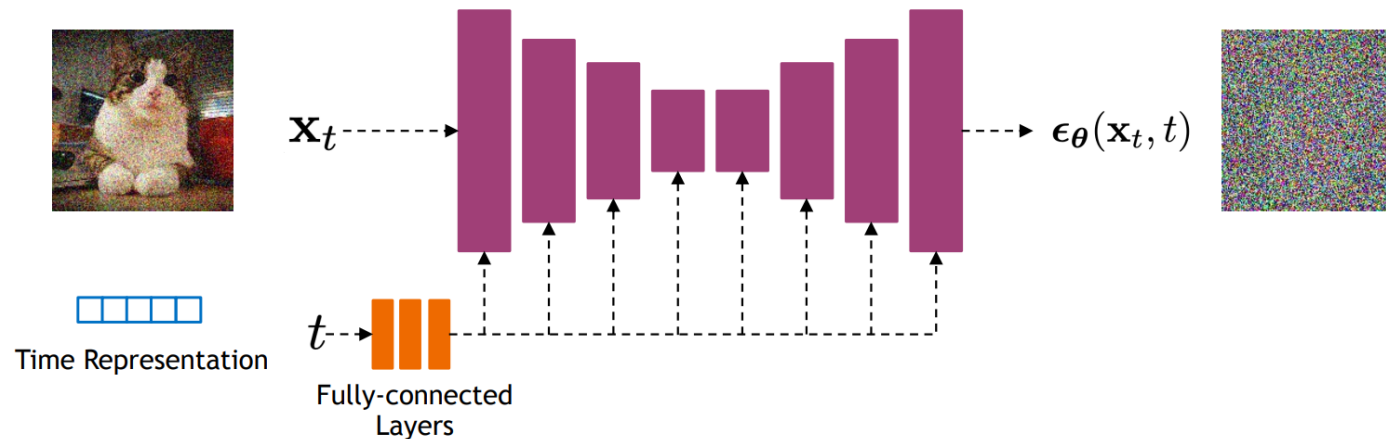
Scheduler for beta (β_t) indicates a predefined sequence of noise variances for each timestep t .

- Linear Schedule: β_t increases linearly from a small initial value to a maximum value.
- Cosine Schedule: Uses a cosine function to define β_t for smoother transitions.

Alpha Terms (α_t and $\bar{\alpha}_t$) is then derived from the beta:

- $\alpha_t = 1 - \beta_t$
- $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$

Creating Training Data as the forward diffusion (noising) process is simulated by adding Gaussian noise to images according to the noise schedule. For each training image x_0 and timestep t , we generate a noisy image x_t using the closed-form equation: $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \epsilon \sim \mathcal{N}(\epsilon; 0, I)$



Implementation

- Samples of DDPM

