

DEEP LEARNING

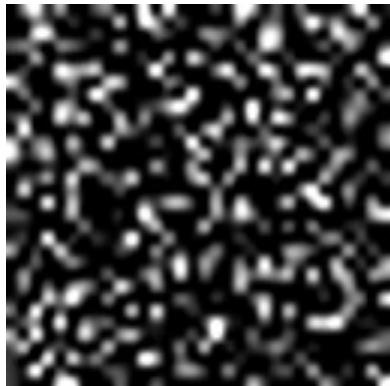
21. Diffusion Models

Naveen George
Supervised by
Dr. Konda Reddy Mopuri
Dept. of AI, IIT Hyderabad
Jan-May 2025

Diffusion models excel at representing and sampling from conditional distributions $p(x/y)$, enabling them to generate high-quality images that adhere to a given condition y.

Class conditioned Generation

$p(\text{image} | y = \text{Class Label})$



Class - 3



Class - 4

Text Conditioned Generation



Prompt: Roman city on top of a ridge, sci-fi illustration by Greg Rutkowski #sci-fi
detailed vivid colors gothic concept illustration by James Gurney and Zdzislaw
Beksiński vivid vivid colors concept illustration colorful interior



Prompt: cinematic photo of a red apple on a table in a classroom, on the blackboard are the words "go big or go home" written in chalk. Using Stable Diffusion

Video and Music Generation



[Source](#)

Image to Image



[Source](#)

Fast-Evolving Field

VAEs, 2013



GANs, 2014



PixelCNN, 2016



BigGAN, 2019



Imagen, 2022



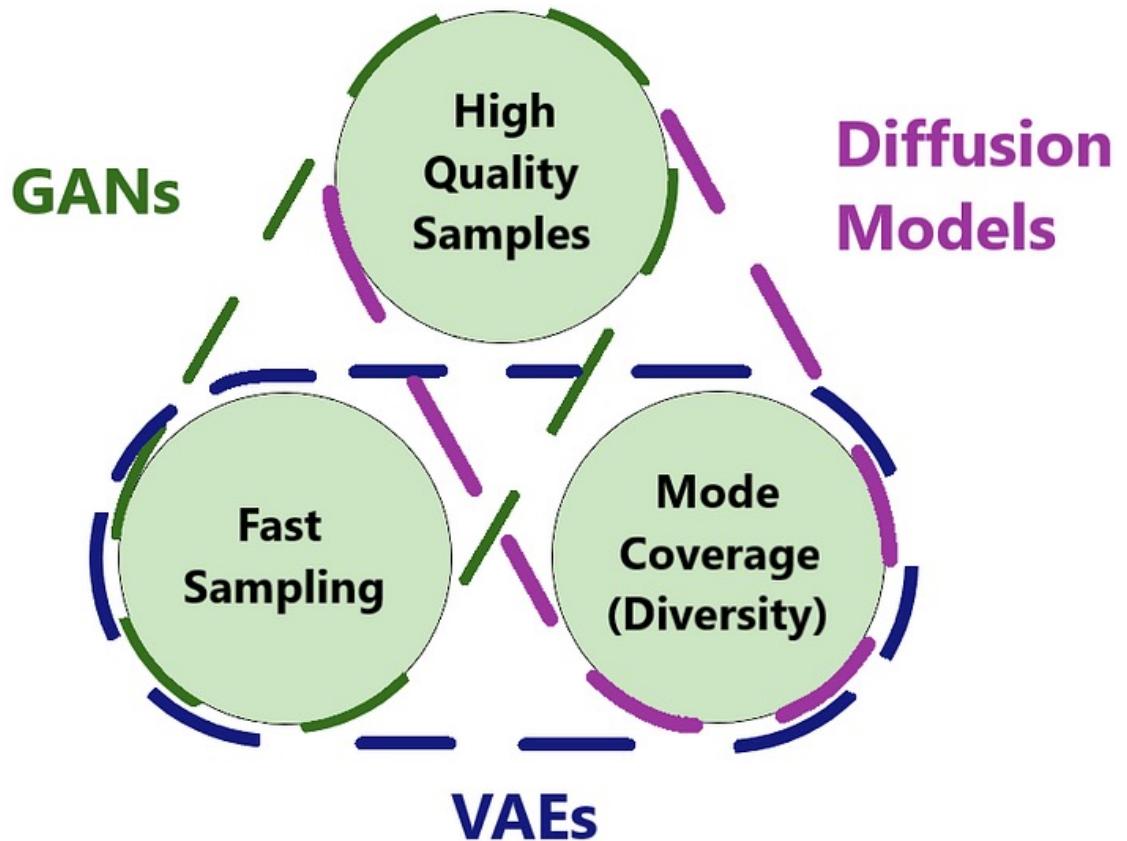
SORA 2024



Recap: Generative Modeling

- Learning a model of distribution from which we can generate new samples.
- Given $x \sim p(x)$ we want to learn $p_\theta(x) \cong p(x)$ (modeling)
- Then, we can generate new samples $x^* \sim p_\theta(x)$ (generation)
- Models you have seen before:
 - VAE
 - GAN

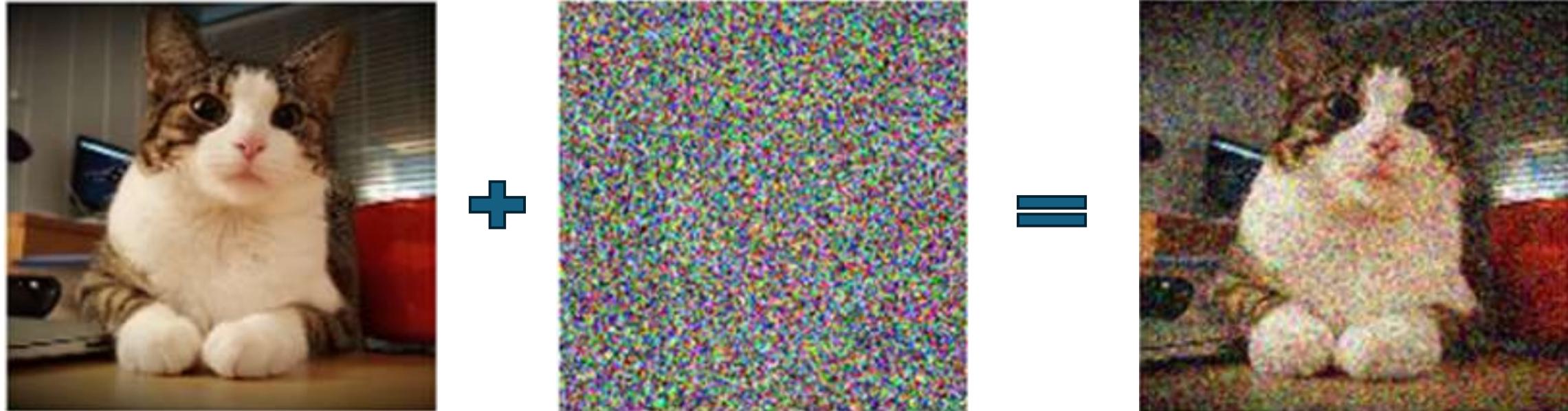
Comparison



Source

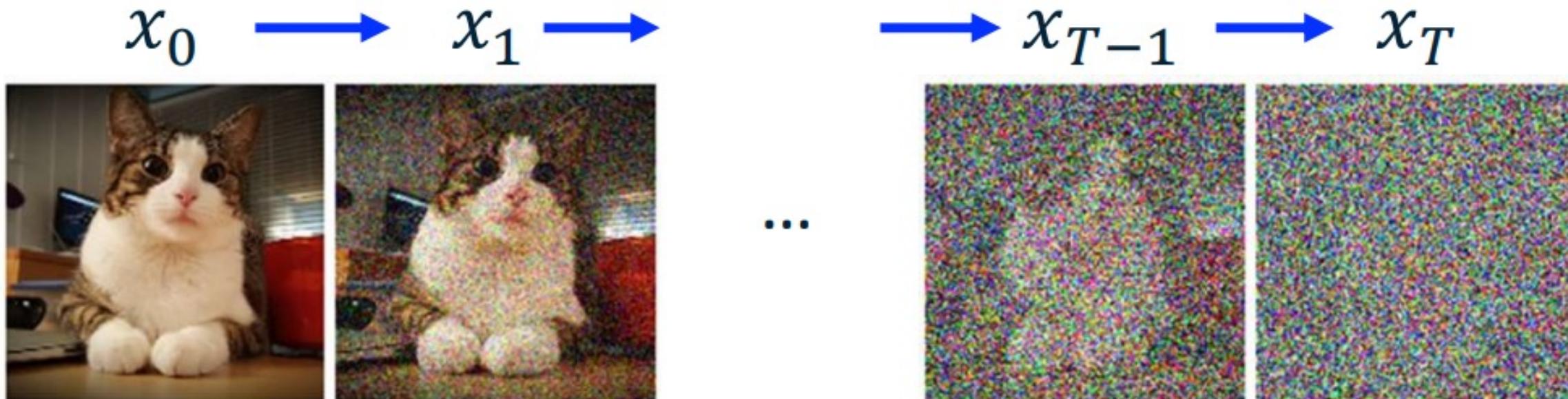
Diffusion Models

- Models data by gradually adding noise.
- Then learns to remove noise (denoise).



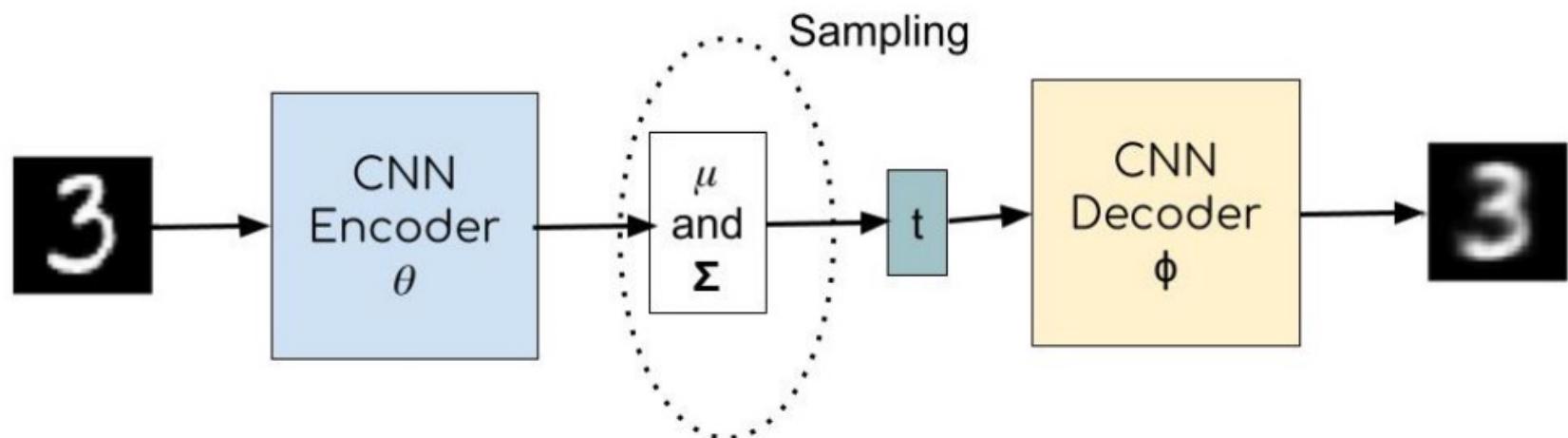
Diffusion Models

- Models data by gradually adding noise.
- Then learns to remove noise (denoise).



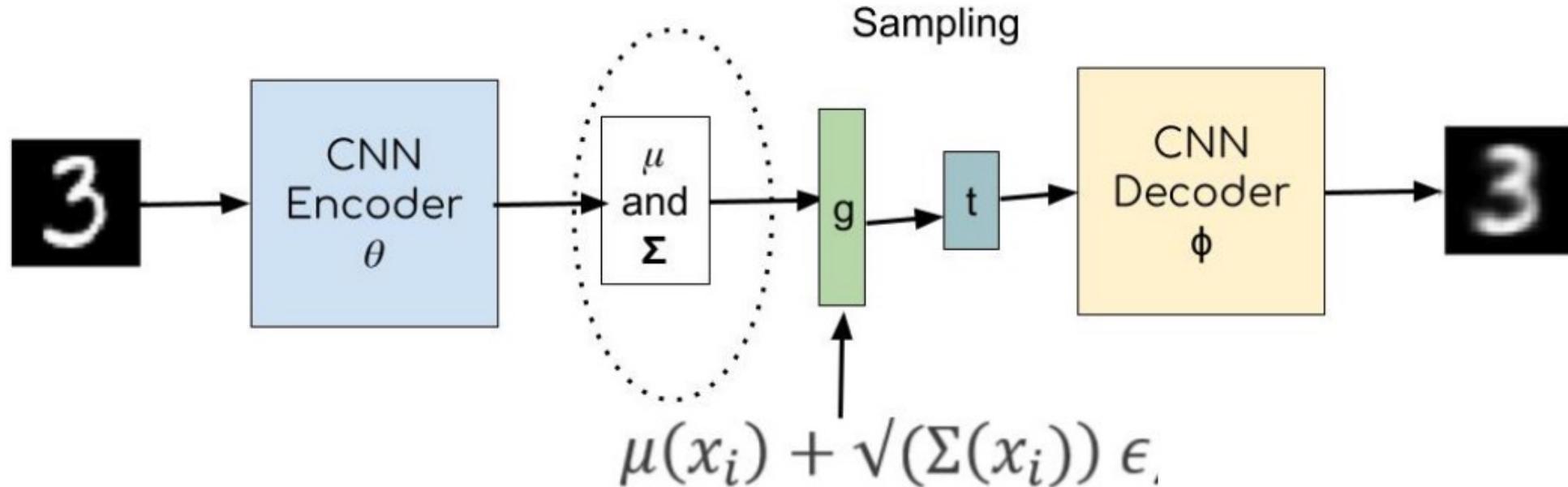
VAEs

Sampling is non-differentiable



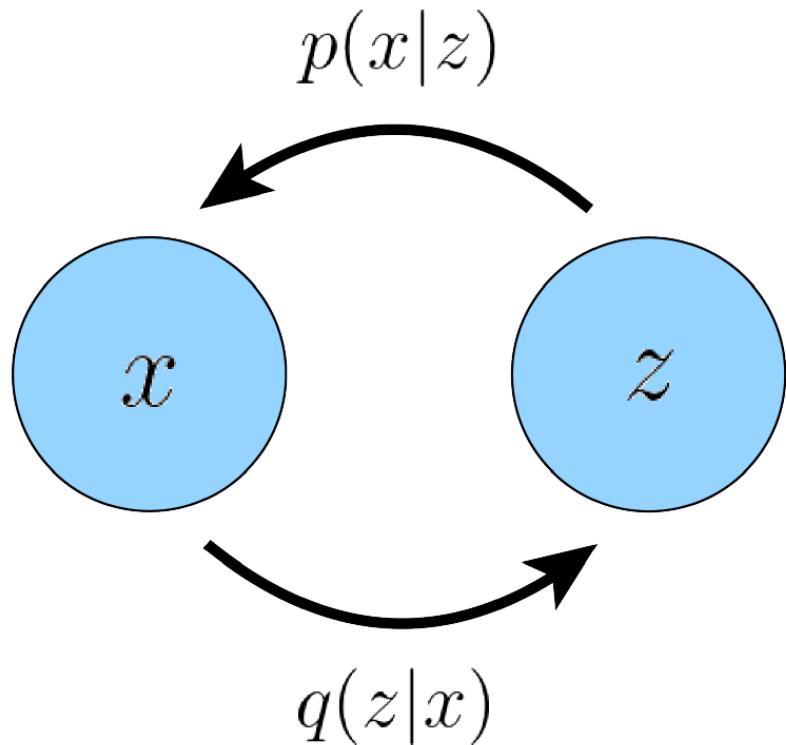
VAEs

So we use Reparameterization Trick



VAEs

- We can visualize it in this form -



Source

Evidence Lower BOund (ELBO)

$$\begin{aligned}
 \log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\
 &= \log \int \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z} | \mathbf{x})}{q_{\phi}(\mathbf{z} | \mathbf{x})} d\mathbf{z} \\
 &= \log \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] \\
 &\geq \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right]
 \end{aligned}$$

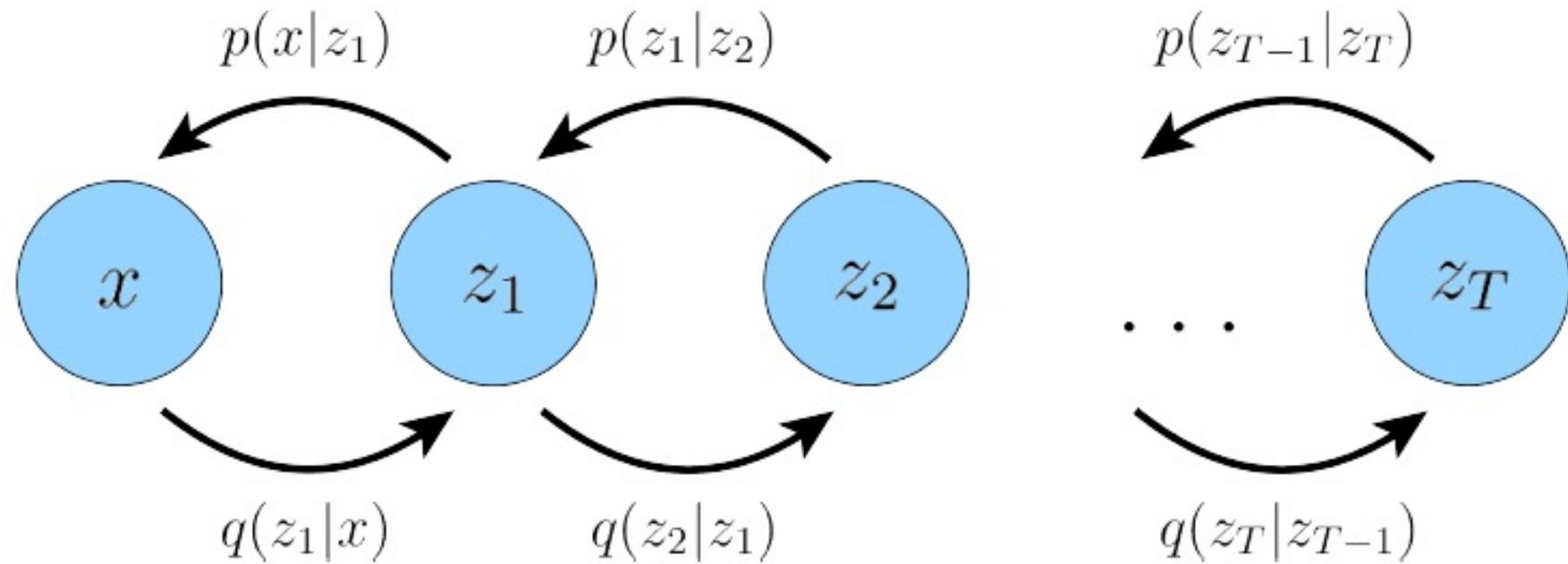
$$\log p(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z} | \mathbf{x}))$$

Evidence Lower BOund (ELBO)

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

Hierarchical VAEs

- A generalization of VAE to multiple hierarchies.
- A special case is called as Markovian HAVE (MHVAE); here, it is a Markov Chain.



[Source](#)

ELBO

$$p(x, z_{1:T}) = p(z_T)p(z_{T-1}|z_T)p(z_{T-2}|z_{T-1}, z_T) \dots p(x|z_1, z_2 \dots z_T)$$

Markovian Property: $p(z_t|z_{t+1}, z_{t+2} \dots z_T) = p(z_t|z_{t+1})$

$$q_\phi(z_{1:T}|x) = q_\phi(z_1|x)q_\phi(z_2|z_1, x) \dots q_\phi(z_T|z_1, z_2 \dots x)$$

Markovian Property: $q_\phi(z_t|z_{t-1}, z_{t-2} \dots z_1) = q_\phi(z_t|z_{t-1})$

$$p(\mathbf{x}, \mathbf{z}_{1:T}) = p(\mathbf{z}_T)p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}_1) \prod_{t=2}^T p_{\boldsymbol{\theta}}(\mathbf{z}_{t-1}|\mathbf{z}_t)$$

$$q_\phi(\mathbf{z}_{1:T}|\mathbf{x}) = q_\phi(\mathbf{z}_1|\mathbf{x}) \prod_{t=2}^T q_\phi(\mathbf{z}_t|\mathbf{z}_{t-1})$$

Extend ELBO

$$\begin{aligned}
 \log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T} \\
 &= \log \int \frac{p(\mathbf{x}, \mathbf{z}_{1:T}) q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x})}{q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x})} d\mathbf{z}_{1:T} \\
 &= \log \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x})} \right] \\
 &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x})} \right]
 \end{aligned}$$

Extend ELBO

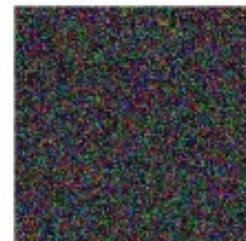
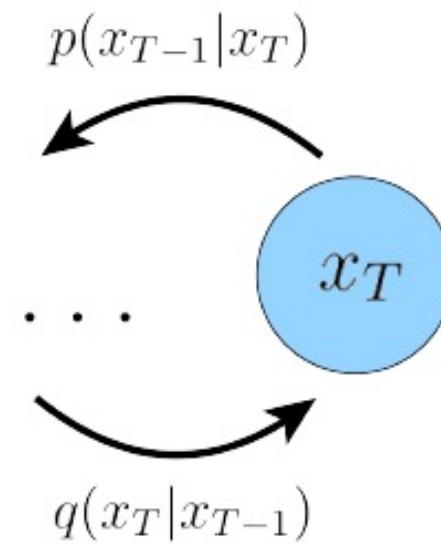
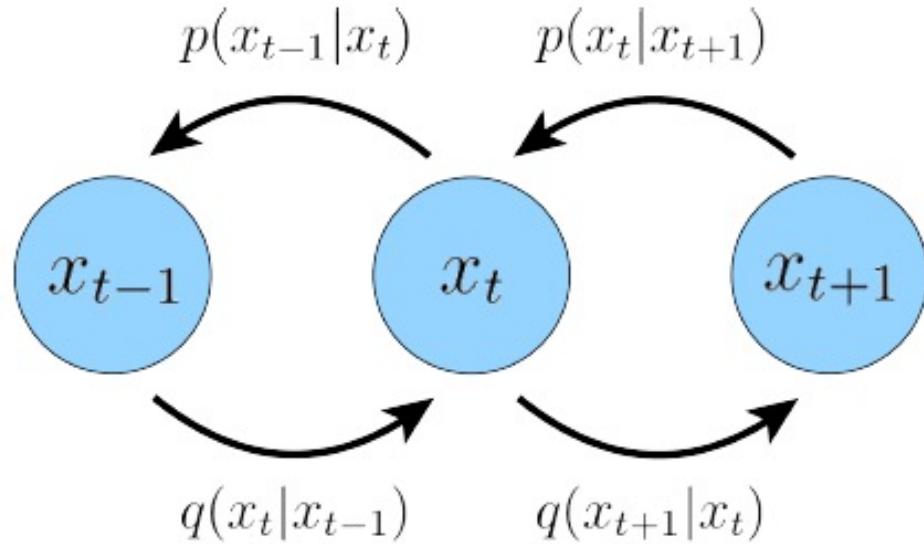
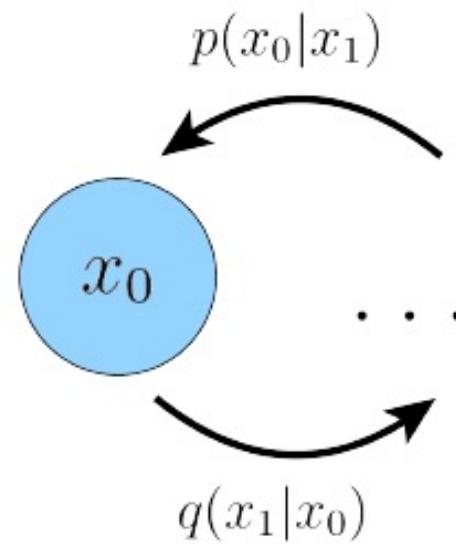
$$\mathbb{E}_{q_{\phi}(z_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z_{1:T})}{q_{\phi}(z_{1:T}|\mathbf{x})} \right] = \mathbb{E}_{q_{\phi}(z_{1:T}|\mathbf{x})} \left[\log \frac{p(z_T)p_{\theta}(\mathbf{x}|z_1) \prod_{t=2}^T p_{\theta}(z_{t-1}|z_t)}{q_{\phi}(z_1|\mathbf{x}) \prod_{t=2}^T q_{\phi}(z_t|z_{t-1})}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

Variational Diffusion Model

Assume 3 conditions with Markovian
HAVE -

- All Latent Dimensions are the same.
- Structure of latent encoder at each timestep is not learned, it is pre-defined as linear Gaussian model.
- Gaussian parameters of the latent encoders vary over time such that at final T it is standard Gaussian.





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

$$\begin{aligned}
\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
&= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2} \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2} \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \\
&= \dots \\
&= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \\
&= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 \\
&\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})
\end{aligned}$$



Latent samples from linear (top) and cosine (bottom) schedules respectively. The latents in the last quarter of the linear schedule are almost purely noise, whereas the cosine schedule adds noise more slowly.

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \quad (3^{\text{rd}} \text{ Condition Assumed})$$

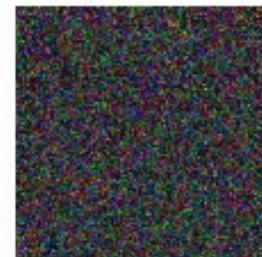
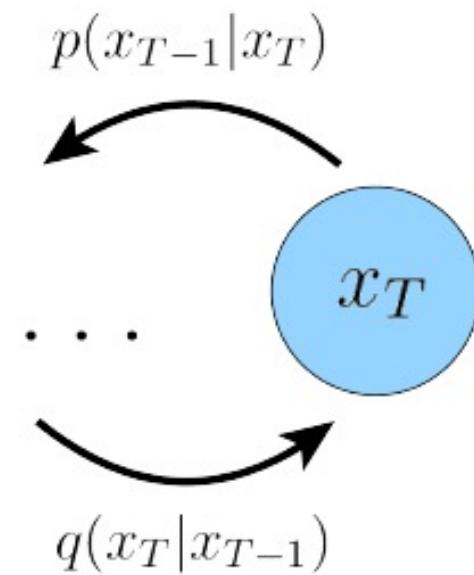
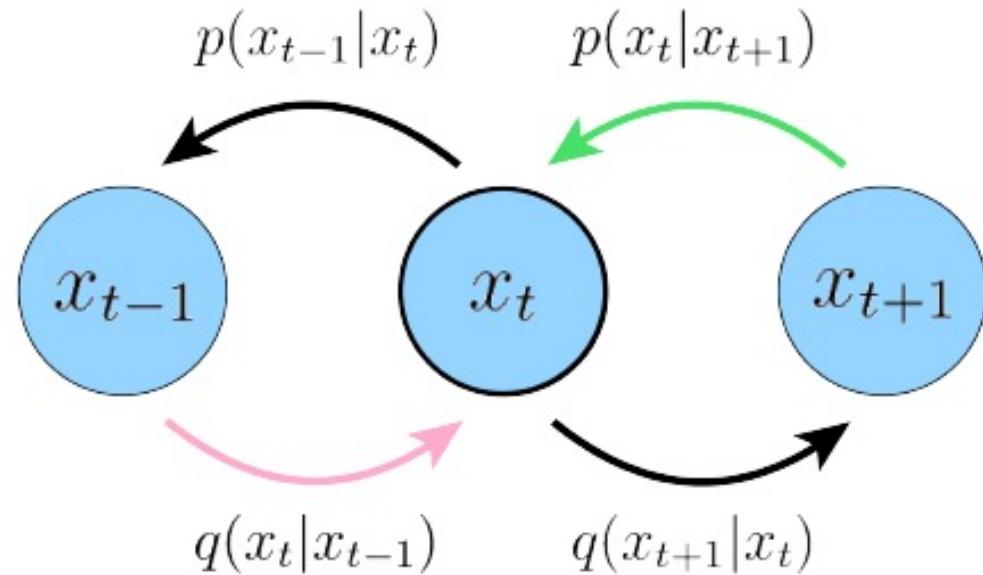
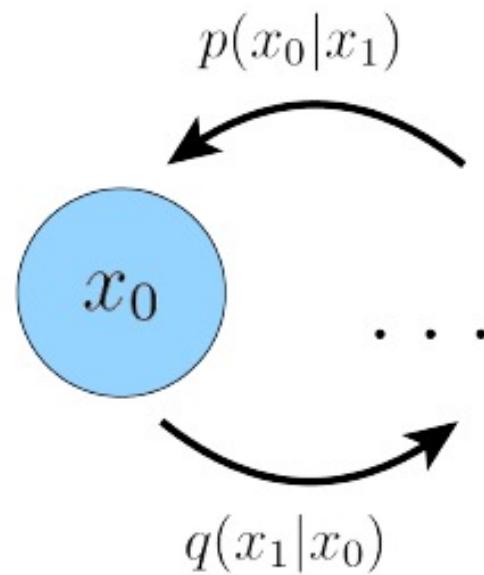
Now by using the 3 conditions which we have used we can simplify the ELBO to this form -

$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$



$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\&= \log \int \frac{p(\mathbf{x}_{0:T}) q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} d\mathbf{x}_{1:T} \\&= \log \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \\&\geq \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\boldsymbol{\theta}}(\mathbf{x}_0 | \mathbf{x}_1) \prod_{t=2}^T p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_T | \mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\boldsymbol{\theta}}(\mathbf{x}_0 | \mathbf{x}_1) \prod_{t=1}^{T-1} p_{\boldsymbol{\theta}}(\mathbf{x}_t | \mathbf{x}_{t+1})}{q(\mathbf{x}_T | \mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right]\end{aligned}$$

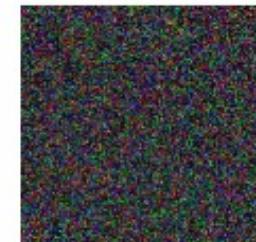
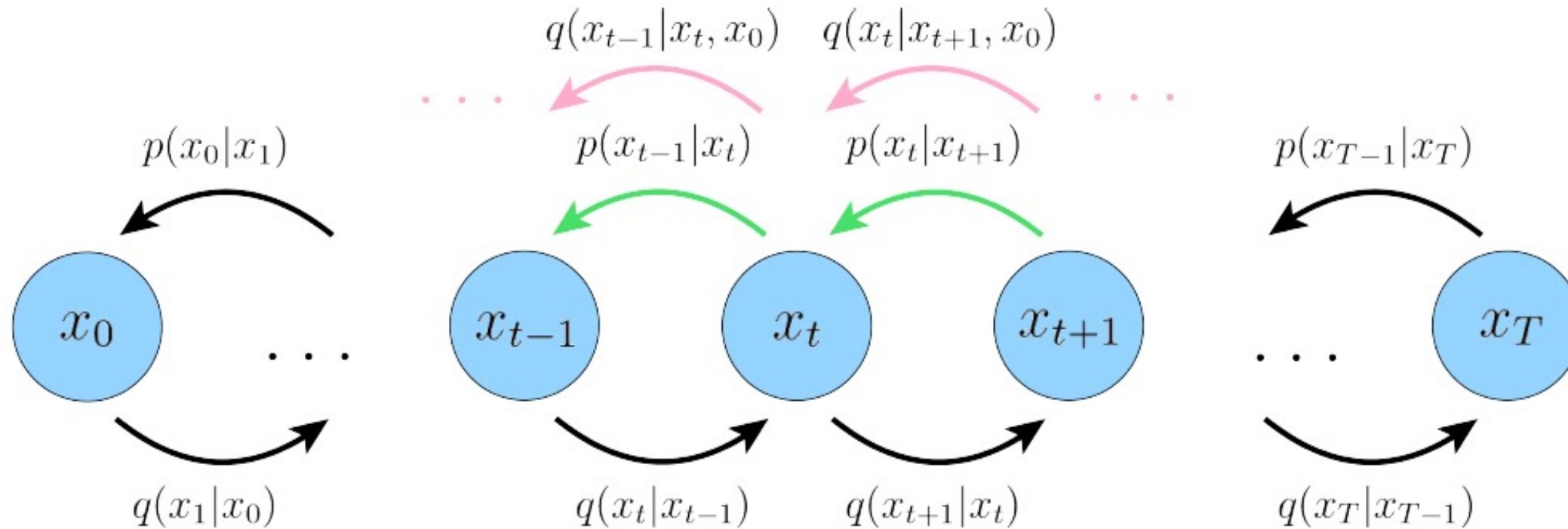
$$\begin{aligned}
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)p_{\theta}(x_0|x_1)}{q(x_T|x_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \prod_{t=1}^{T-1} \frac{p_{\theta}(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)}{q(x_T|x_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^{T-1} \log \frac{p_{\theta}(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)}{q(x_T|x_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)] + \mathbb{E}_{q(\mathbf{x}_{T-1}, \mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(x_T)}{q(x_T|x_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] \\
&= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{\text{KL}}(q(x_T|x_{T-1}) \parallel p(x_T))]}_{\text{prior matching term}} \\
&\quad - \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(x_t|x_{t-1}) \parallel p_{\theta}(x_t|x_{t+1}))]}_{\text{consistency term}}
\end{aligned}$$



Source

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

$$\begin{aligned}
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)p_{\theta}(x_0|x_1)}{q(x_1|x_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(x_{t-1}|x_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(x_t|x_0)}{q(x_{t-1}|x_0)}} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)p_{\theta}(x_0|x_1)}{q(x_1|x_0)} + \log \frac{q(x_1|x_0)}{q(x_T|x_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)p_{\theta}(x_0|x_1)}{q(x_T|x_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(x_T)}{q(x_T|x_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] \\
&= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(x_T)}{q(x_T|x_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] \\
&= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(x_0|x_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(x_T|x_0) \parallel p(x_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]}_{\text{denoising matching term}}
\end{aligned}$$



[Source](#)

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

On simplifying -

$$\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)})$$

In order to match approximate denoising transition step $p_{\theta}(x_{t-1} / x_t)$ to ground-truth denoising transition step $q(x_{t-1} / x_t, x_0)$ as closely as possible, we can also model it as a Gaussian.

On simplifying -

$$\begin{aligned}
 & \arg \min_{\theta} D_{\text{KL}}(q(x_{t-1} | x_t, x_0) \| p_{\theta}(x_{t-1} | x_t)) \\
 &= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(x_{t-1}; \mu_q, \Sigma_q(t)) \| \mathcal{N}(x_{t-1}; \mu_{\theta}, \Sigma_q(t))) \\
 &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_{\theta} - \mu_q\|_2^2 \right]
 \end{aligned}$$

Different Interpretation for VDM

$$\begin{aligned}
 & \arg \min_{\theta} D_{\text{KL}}(q(x_{t-1} | x_t, x_0) \| p_{\theta}(x_{t-1} | x_t)) \\
 &= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(x_{t-1}; \mu_q, \Sigma_q(t)) \| \mathcal{N}(x_{t-1}; \mu_{\theta}, \Sigma_q(t))) \\
 &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \left[\|\epsilon_0 - \hat{\epsilon}_{\theta}(x_t, t)\|_2^2 \right]
 \end{aligned}$$

Some works have shown that predicting noise resulted in better performance.

Pseudocode

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
6:      $\nabla_{\theta} \|\mathbf{x}_0 - \hat{\mathbf{x}}_{\theta}(\mathbf{x}_0 + \alpha_t \epsilon, t)\|^2$ 
7: until converged
  
```

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$ :
3:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\epsilon = 0$ 
4:    $\mathbf{x}_{t-1} = \hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) + \alpha_{t-1} \epsilon$ 
5: end for
6: return  $\mathbf{x}_0$ 
  
```

Practical Details

- For images, $\epsilon_\theta(x_t, t)$ is typically implemented via the U-Net architecture.
- Time input t is discrete integer – usually handled via (learnable) embeddings
- Can tune forward process: number of steps T , variance schedule β_t

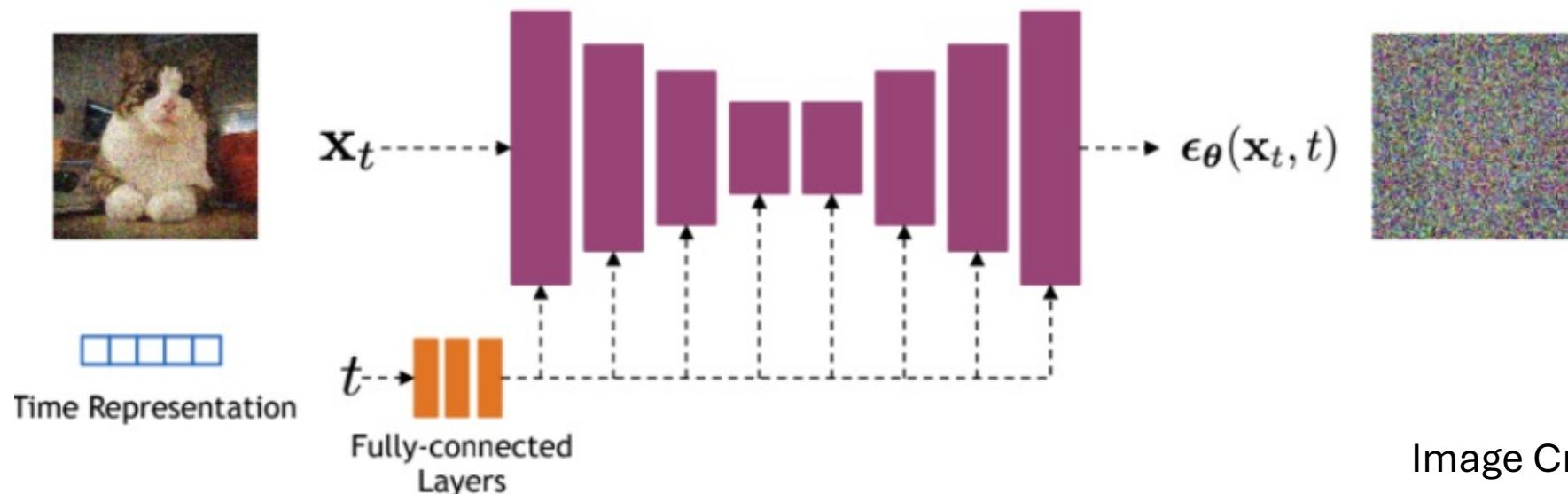
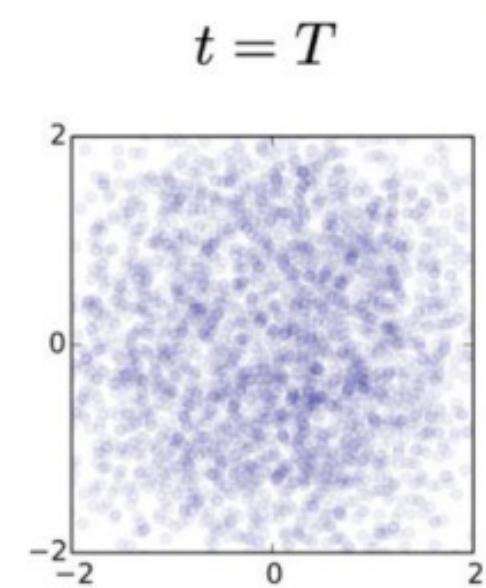
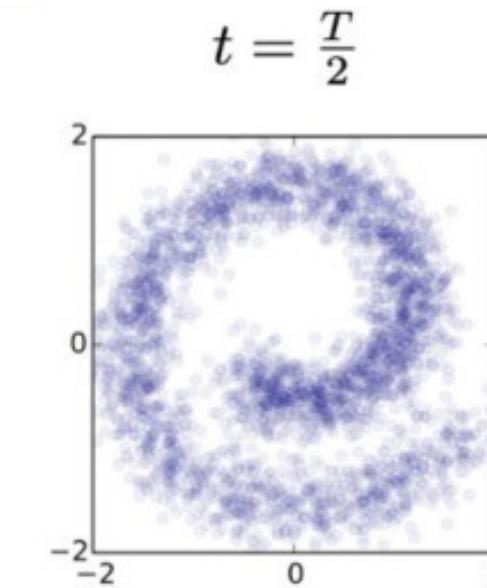
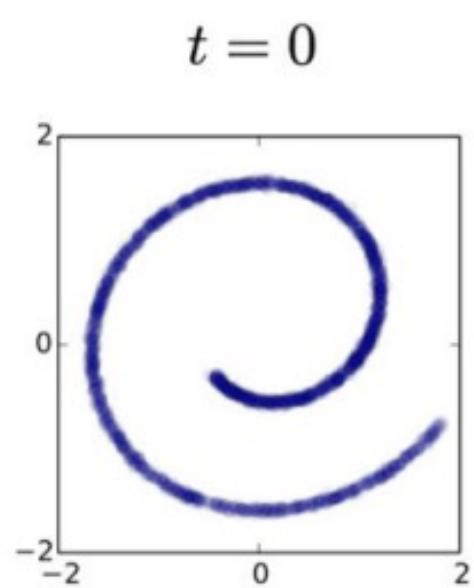
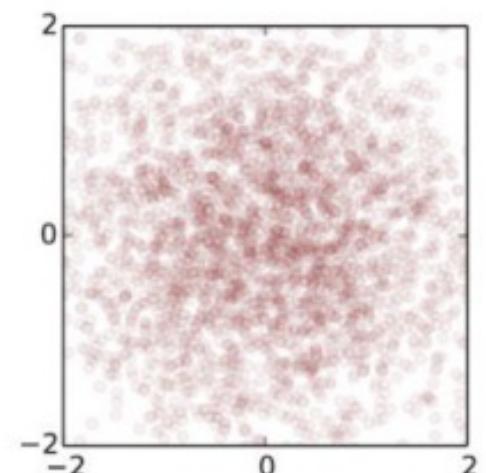
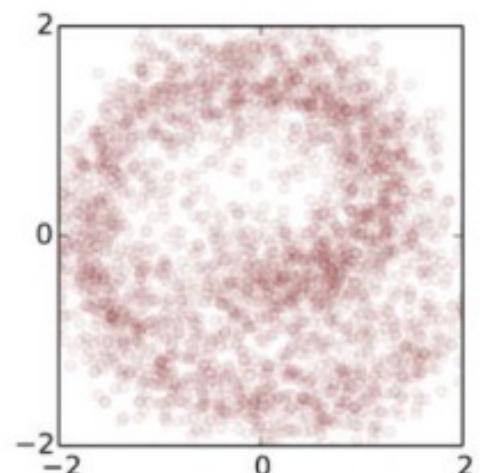
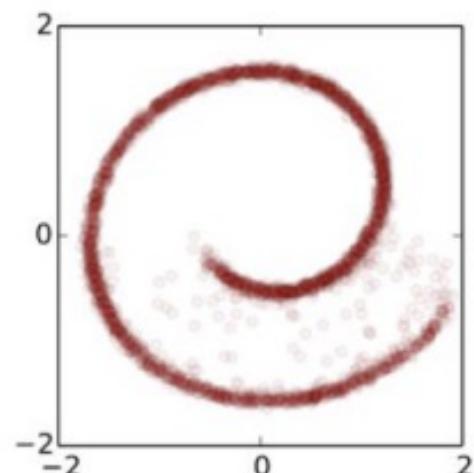


Image Credit: Arash Vahdat

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

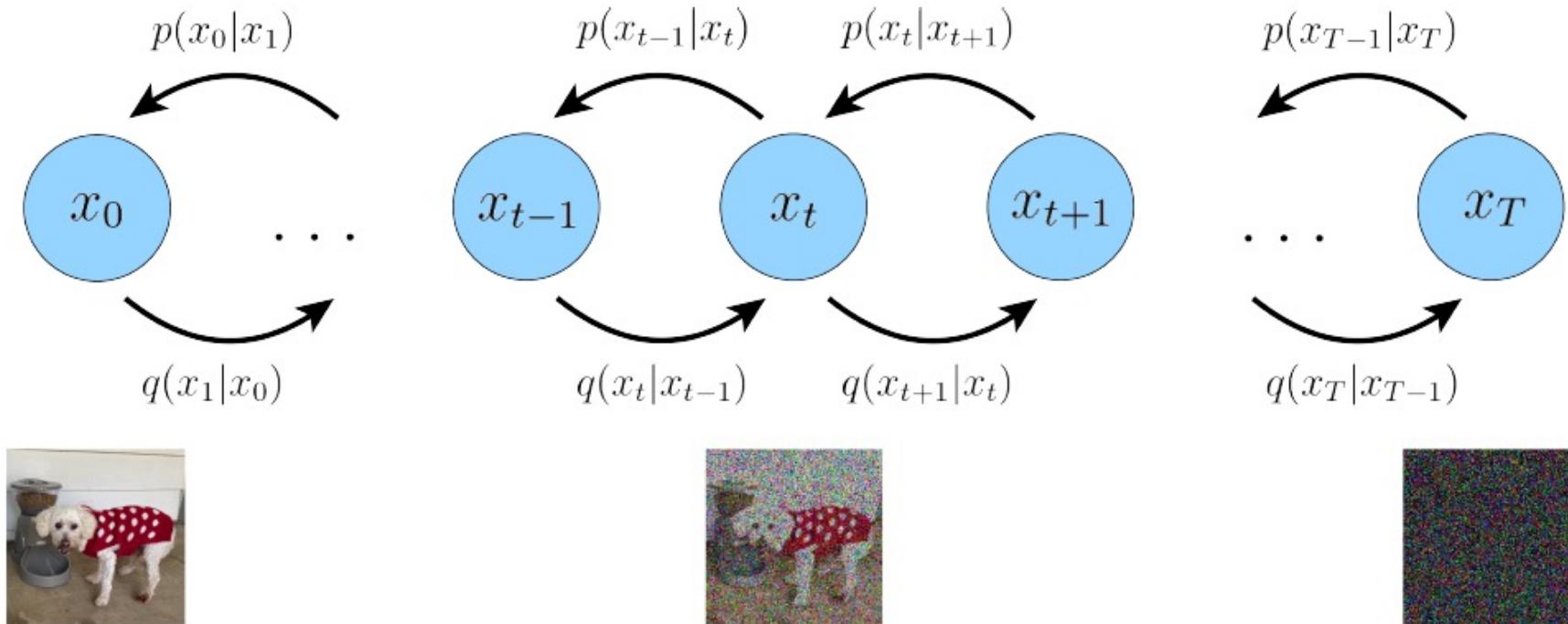


The reverse trajectory
 $p_\theta(\mathbf{x}_{0:T})$



DDPM

- Slow sampling due to many steps.
- High computational cost.
- Limited control over generation process.



Denoising Diffusion Implicit Models (DDIM)

- Accelerates sampling speed significantly.
- Deterministic sampling (no randomness in reverse).
- Retains quality with fewer steps.

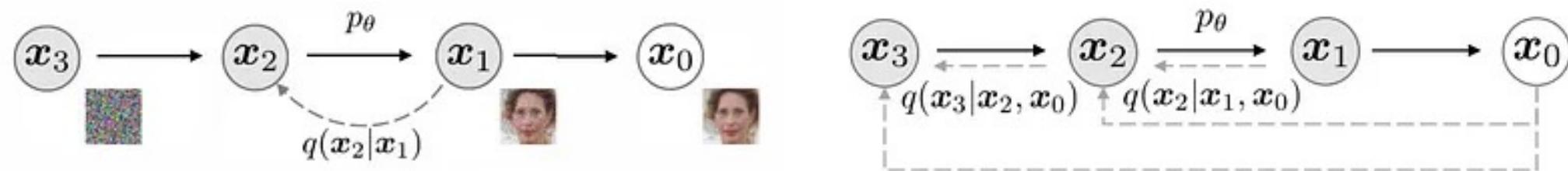


Figure 1: Graphical models for diffusion (left) and non-Markovian (right) inference models.

Denoising Diffusion Implicit Models (DDIM)

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta^{(t)}(x_t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \epsilon_\theta^{(t)} + \sigma_t \epsilon_t$$

Make it deterministic -

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta^{(t)}(x_t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \epsilon_\theta^{(t)}$$

Guidance

- Till now, we focused on $p(x)$ – *Unconditional distribution*
- With Conditional distribution $p(x|y)$ – Control image generation

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

Classifier Guidance

- Uses external classifier to guide generation.
- Conditions output on specific labels or features.
- Improved accuracy for conditional generation tasks.



[Source](#)

Classifier Guidance

$$\begin{aligned}
 \nabla \log p(\mathbf{x}_t|y) &= \nabla \log \left(\frac{p(\mathbf{x}_t)p(y|\mathbf{x}_t)}{p(y)} \right) \\
 &= \nabla \log p(\mathbf{x}_t) + \nabla \log p(y|\mathbf{x}_t) - \nabla \log p(y) \\
 &= \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y|\mathbf{x}_t)}_{\text{adversarial gradient}}
 \end{aligned}$$

Algorithm 1 Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $p_\phi(y|x_t)$, and gradient scale s .

Input: class label y , gradient scale s
 $x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
for all t from T to 1 **do**
 $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$
 $x_{t-1} \leftarrow$ sample from $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$
end for
return x_0

[Source](#)

Classifier-Free Guidance

- Removes dependency on external classifiers.
- Here we train $p(x/y)$ with conditional dropout.
- Simpler training and more flexible conditioning.

[Source](#)



Classifier-Free Guidance

$$\nabla \log p(y|\mathbf{x}_t) = \nabla \log p(\mathbf{x}_t|y) - \nabla \log p(\mathbf{x}_t)$$

$$\begin{aligned}\nabla \log p(\mathbf{x}_t|y) &= \nabla \log p(\mathbf{x}_t) + \gamma (\nabla \log p(\mathbf{x}_t|y) - \nabla \log p(\mathbf{x}_t)) \\ &= \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{x}_t|y) - \gamma \nabla \log p(\mathbf{x}_t) \\ &= \underbrace{\gamma \nabla \log p(\mathbf{x}_t|y)}_{\text{conditional score}} + \underbrace{(1 - \gamma) \nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}}\end{aligned}$$

Algorithm 1 Joint training a diffusion model with classifier-free guidance

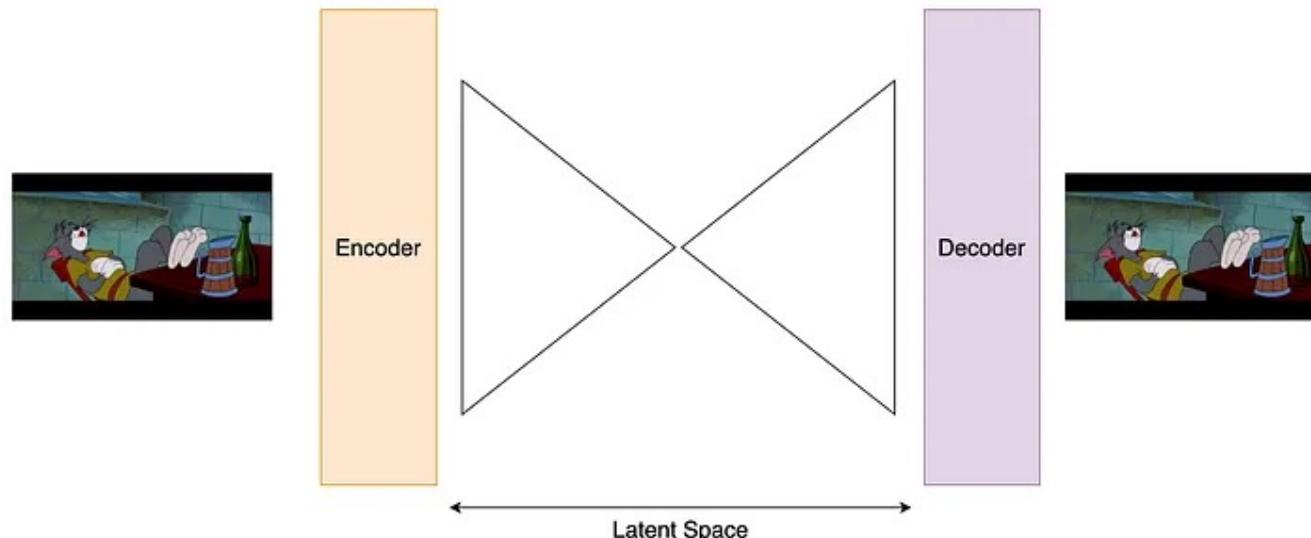
Require: p_{uncond} : probability of unconditional training

- 1: **repeat**
 - 2: $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$ ▷ Sample data with conditioning from the dataset
 - 3: $\mathbf{c} \leftarrow \emptyset$ with probability p_{uncond} ▷ Randomly discard conditioning to train unconditionally
 - 4: $\lambda \sim p(\lambda)$ ▷ Sample log SNR value
 - 5: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 6: $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \epsilon$ ▷ Corrupt data to the sampled log SNR value
 - 7: Take gradient step on $\nabla_\theta \|\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \epsilon\|^2$ ▷ Optimization of denoising model
 - 8: **until** converged
-

[Source](#)

Latent Diffusion Models (LDM)

- Operates in latent space instead of pixel space.
- Reduces computational resources dramatically.
- High-quality outputs with efficient training.



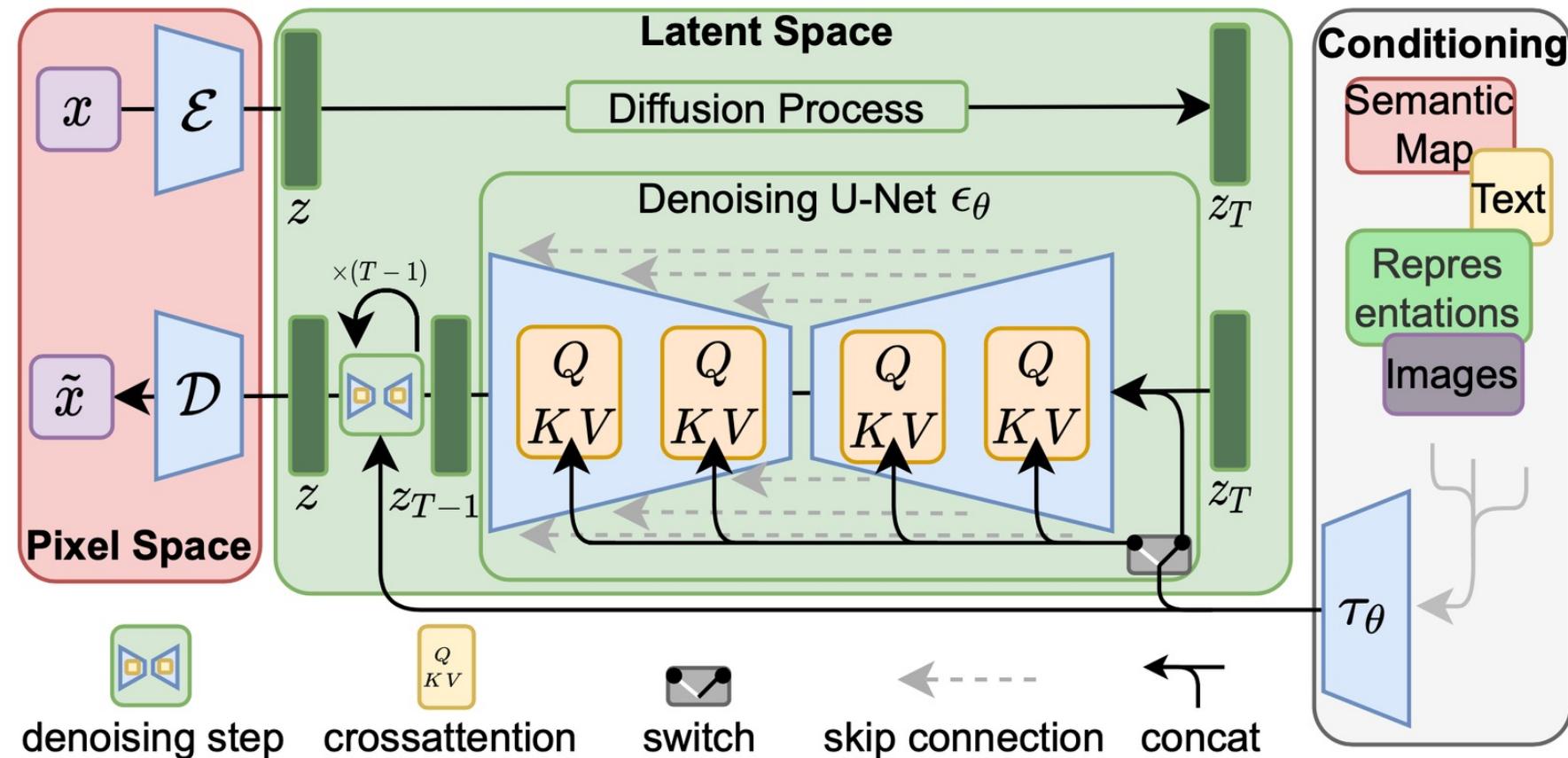
Latent Diffusion Models (LDM)

- *Denoising Model: Time-conditioned U-Net with cross-attention for flexible conditioning (e.g., class labels, semantic maps).*
- Cross Attention Mechanism – Integrates different Modalities

Text/Conditional Diffusion Models

- Conditions generation on textual input.
- Enables text-to-image and text-to-video tasks.
- Prominent example: Stable Diffusion.

Latent Diffusion Models (LDM)

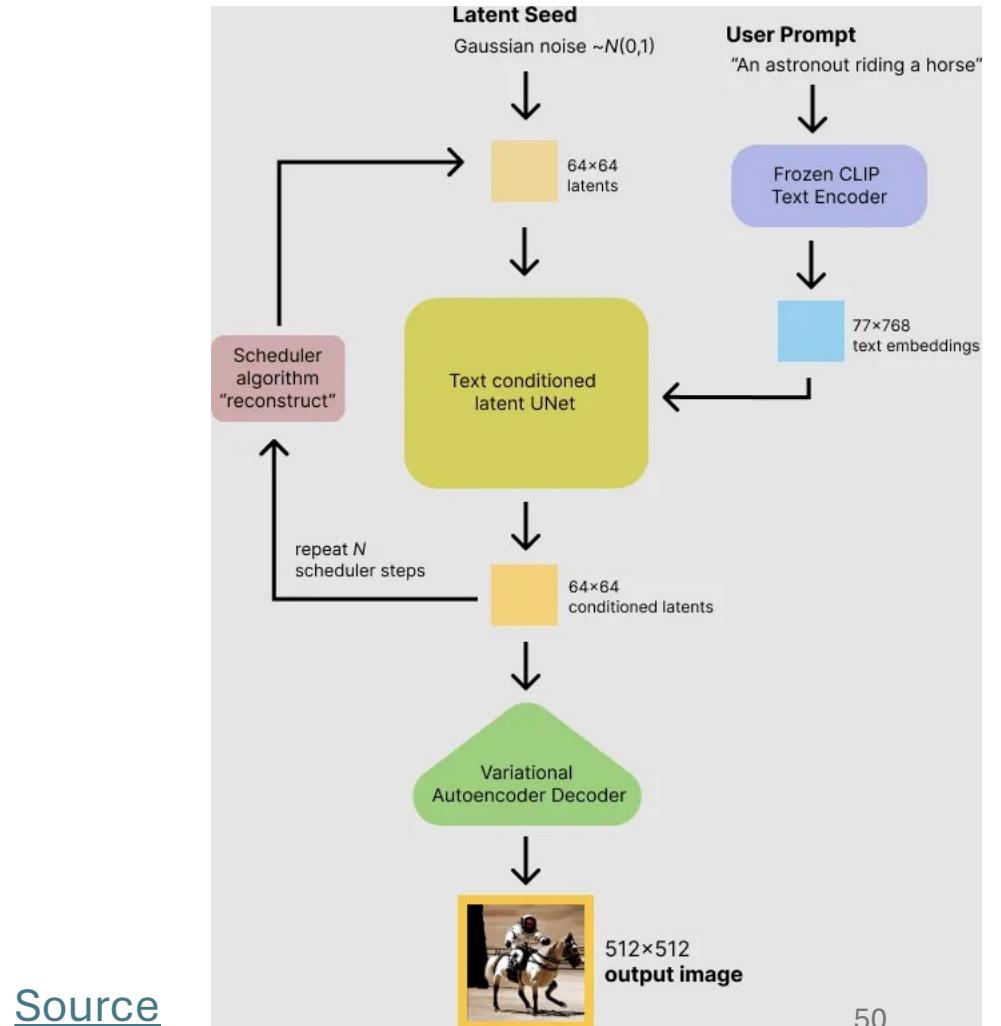


Source

Latent Diffusion Models (LDM)

Mainly 3 main Components in Latent Diffusion Models :

- An autoencoder (VAE)
- An U-Net
- A text encoder. Eg: CLIP



[Source](#)

THANK YOU