

Lecture 17

Diffusion Models

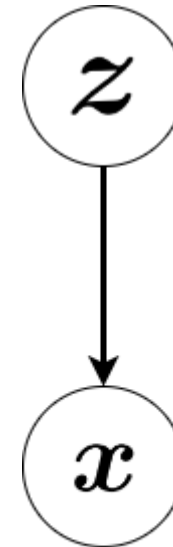
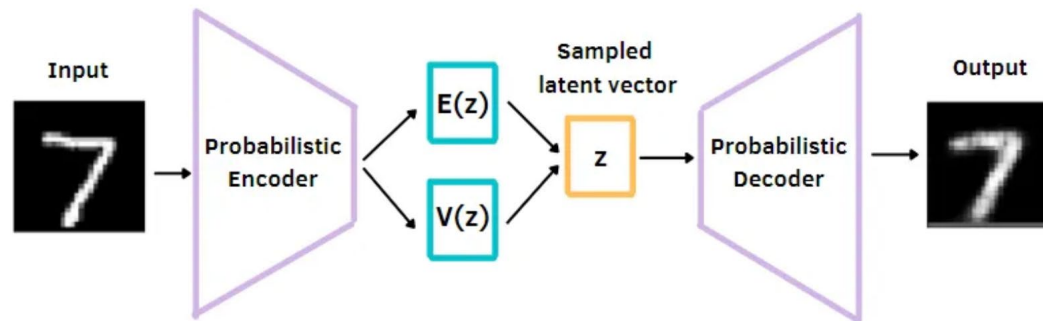
This presentation draws upon insights and information from various sources, including Kevin P. Murphy's 'Probabilistic Machine Learning: Advanced Topics,' Lilian Weng's 'What are Diffusion Models?', 'Denoising Diffusion Probabilistic Models' by Jonathan Ho, Ajay Jain, and Pieter Abbeel, and the tutorial on Denoising Diffusion Probabilistic Models by Jiaming Song, Chenlin Meng, and Arash Vahdat. The content presented is a synthesis of these resources, integrated and interpreted through my own understanding and analysis.

Diffusion Models

- Diffusion models take their inspiration from the principles of non-equilibrium thermodynamics
- Diffusion models use a Markov chain to gradually introduce noise into data, then learn to reverse this process to recreate the desired data from the noise.



Variational Autoencoder



VAE

Diffusion Models

$$x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$$

Data



Noise

Diffusion Models

$$x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$$

Data



Noise

$$x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$$

Diffusion Models vs VAE

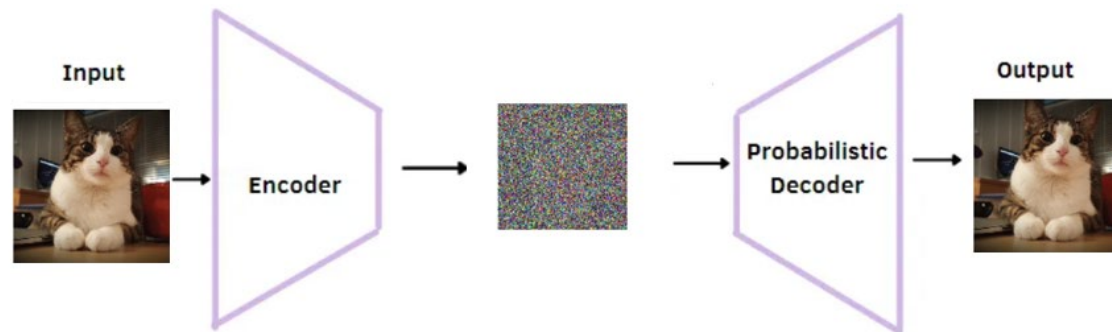
$$x_0 \longrightarrow x_1 \longrightarrow \dots \longrightarrow x_T$$

Data



Noise

$$x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$$



Diffusion Models vs VAE

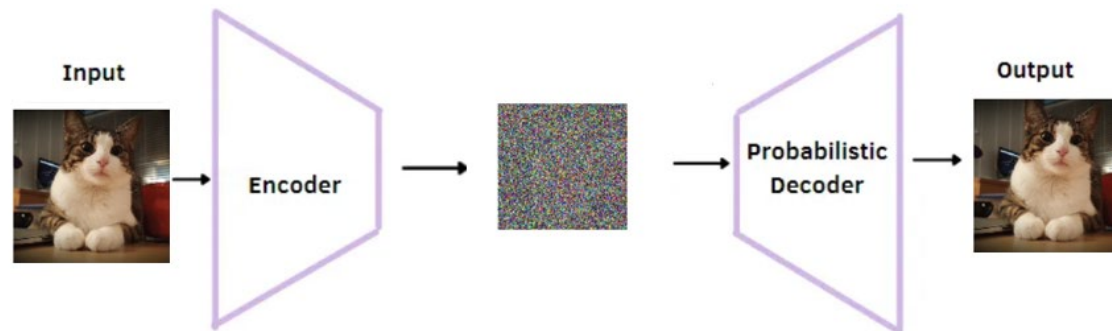
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

Data

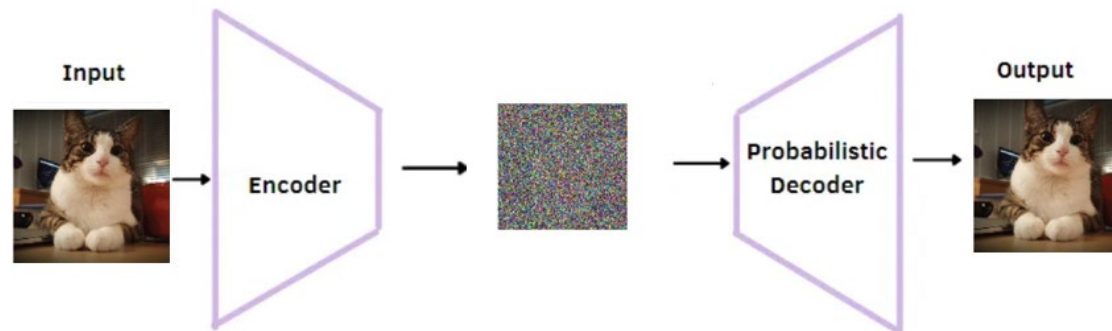
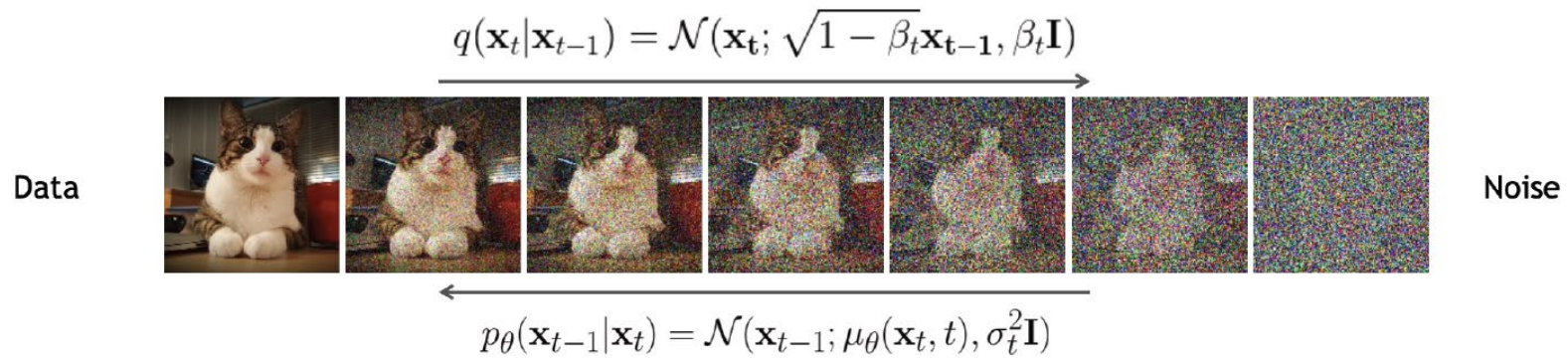


Noise

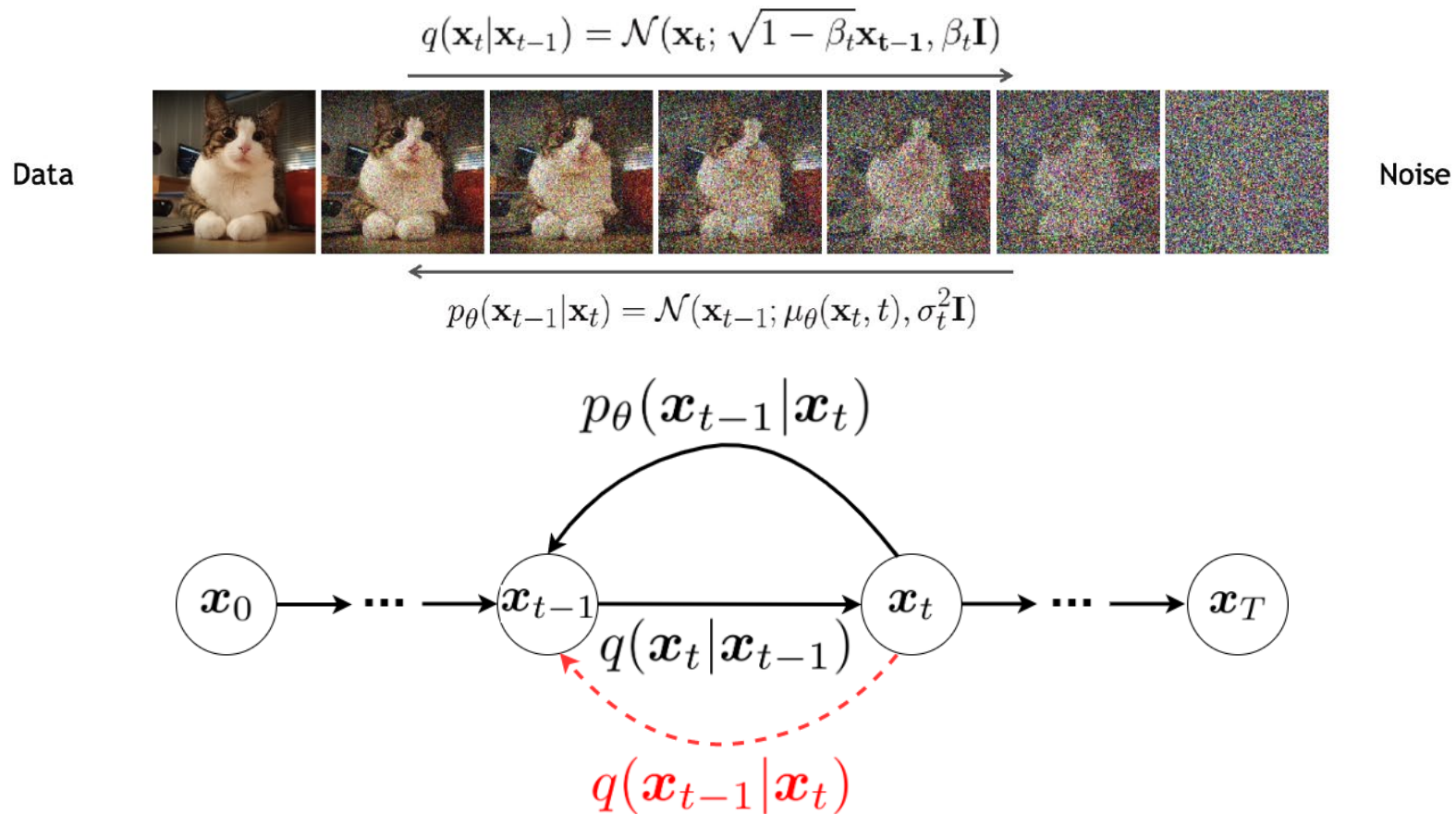
$$x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$$



Diffusion Models vs VAE



Denoising Diffusion Probabilistic Models (DDPMs)



The Forwards Encoder Process in DDPMs

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



The forward process is defined by the equation:

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

Where:

- ▶ β_t is a noise level parameter, part of a predetermined noise schedule.

$$\beta_1 < \beta_2 < \dots < \beta_T$$

$$\beta_t \in (0, 1)$$

Joint Distribution over Latent States in DDPMs

The joint distribution is expressed as:

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

Where:

- ▶ $x_{1:T}$ represents all latent states from time 1 to T.
- ▶ x_0 is the initial state.

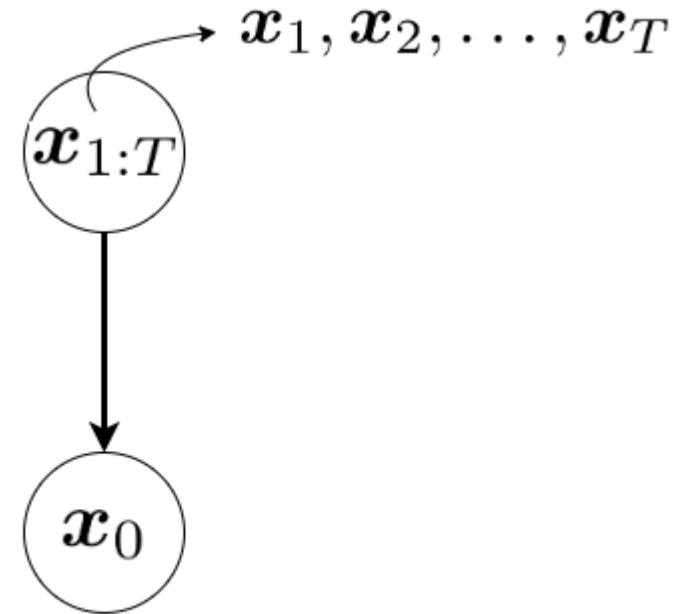
Joint Distribution over Latent States in DDPMs

The joint distribution is expressed as:

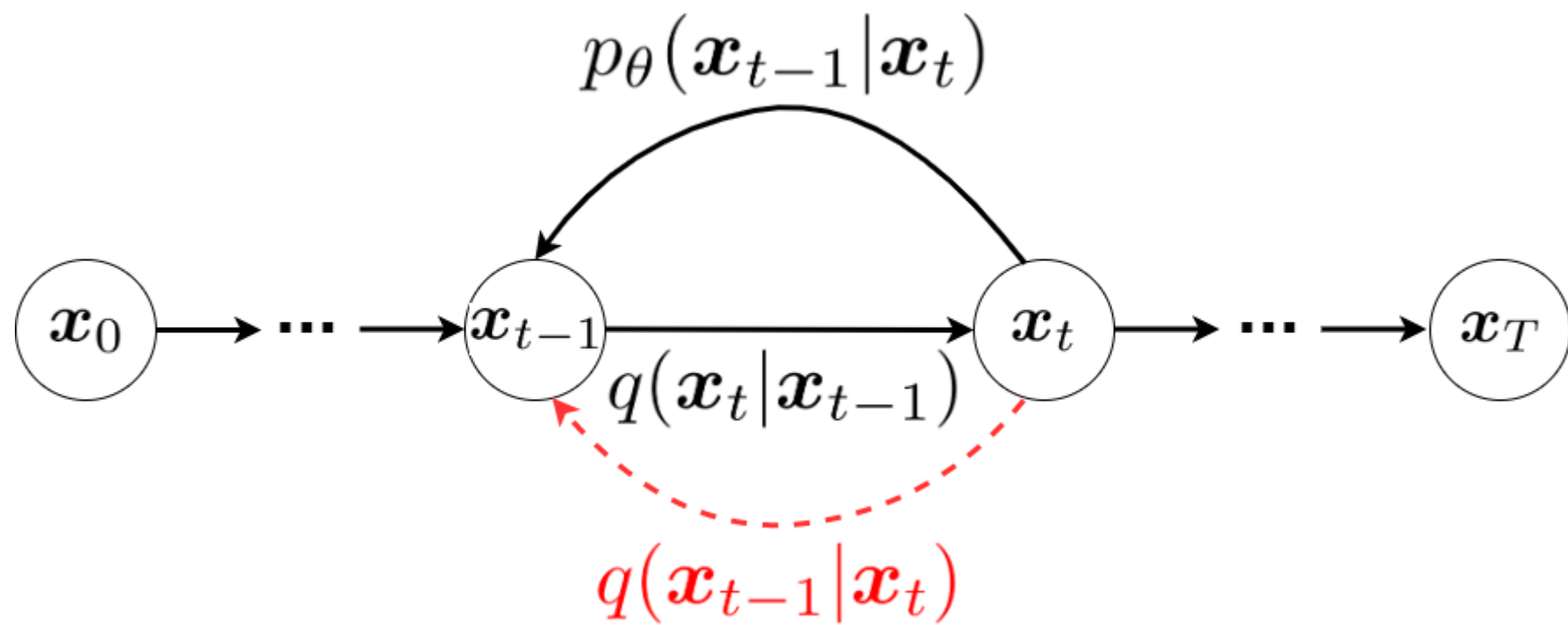
$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

Where:

- ▶ $x_{1:T}$ represents all latent states from time 1 to T.
- ▶ x_0 is the initial state.



Diffusion model





$$x_0 \longleftarrow x_1 \longleftarrow \cdots \longleftarrow x_T$$

The generator is modeled to reverse the diffusion process.

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

where the variance $\Sigma_{\theta}(x_t; t)$ is often set to $\sigma_t^2 I$.



$$x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$$

The joint distribution is given by

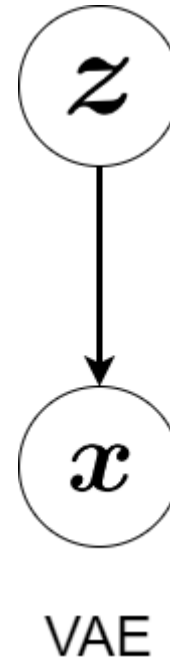
$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

with $p(x_T)$ being $\mathcal{N}(0; I)$.

DDPM vs VAE

- **In VAE**

- Observed variable x
- Hidden Variable z



DDPM vs VAE

- In VAE

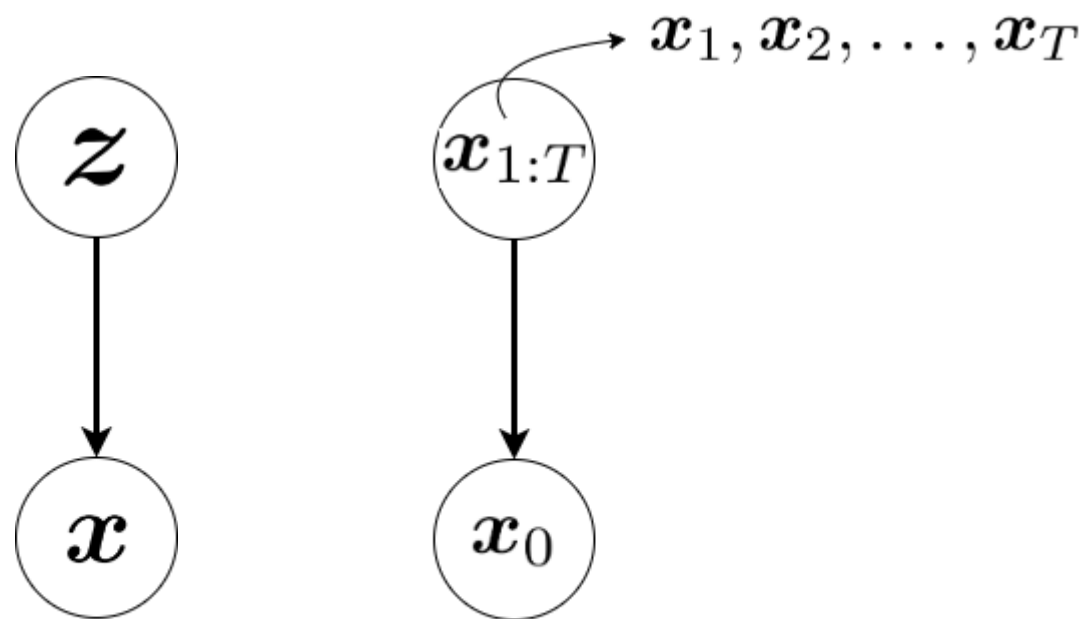
- Observed variable x
- Hidden Variable z



$$x_0 \longleftarrow x_1 \longleftarrow \dots \longleftarrow x_T$$

- In DDPM

- Observed variable x_0
- Hidden variable x_1, x_2, \dots, x_T



VAE

Diffusion model

Objective Functions in VAEs vs. DDPMs

VAEs:

Maximize the likelihood of the data:

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz$$

Objective Functions in VAEs vs. DDPMs

VAEs:

Maximize the likelihood of the data:

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz$$

Diffusion Models:

$$p_{\theta}(x_0) = \int p_{\theta}(x_0, x_{1:T}) dx_{1:T}$$

Variational Lower Bounds in VAEs vs. DDPMs

VAEs:

$$p_{\theta}(x) \geq \mathbb{E}_{q(z|x)}[\log p_{\theta}(x|z)] - D_{KL}(q(z|x) || p_{\theta}(z))$$

Variational Lower Bounds in VAEs vs. DDPMs

VAEs:

$$p_{\theta}(x) \geq \mathbb{E}_{q(z|x)}[\log p_{\theta}(x|z)] - D_{KL}(q(z|x) || p_{\theta}(z))$$

Diffusion Models:

$$p_{\theta}(x_0) \geq \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_{\theta}(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0) || p_{\theta}(x_{1:T}))$$

Variational Lower Bound

$$\mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0)||p_\theta(x_{1:T}))$$

Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0)||p_\theta(x_{1:T})) \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \int q(x_{1:T}|x_0) \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} dx_{1:T} \end{aligned}$$

Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0)||p_\theta(x_{1:T})) \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \int q(x_{1:T}|x_0) \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} dx_{1:T} \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \end{aligned}$$

Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0)||p_\theta(x_{1:T})) \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \int q(x_{1:T}|x_0) \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} dx_{1:T} \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \\ &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p_\theta(x_0|x_{1:T}) - \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \end{aligned}$$

Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0)||p_\theta(x_{1:T})) \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \int q(x_{1:T}|x_0) \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} dx_{1:T} \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \\ &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p_\theta(x_0|x_{1:T}) - \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \\ &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p_\theta(x_0|x_{1:T}) + \log \frac{p_\theta(x_{1:T})}{q(x_{1:T}|x_0)} \right] \end{aligned}$$

Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - D_{KL}(q(x_{1:T}|x_0)||p_\theta(x_{1:T})) \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \int q(x_{1:T}|x_0) \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} dx_{1:T} \\ &= \mathbb{E}_{q(x_{1:T}|x_0)}[\log p_\theta(x_0|x_{1:T})] - \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \\ &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p_\theta(x_0|x_{1:T}) - \log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{1:T})} \right] \\ &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p_\theta(x_0|x_{1:T}) + \log \frac{p_\theta(x_{1:T})}{q(x_{1:T}|x_0)} \right] \\ &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_0|x_{1:T})p_\theta(x_{1:T})}{q(x_{1:T}|x_0)} \right] = \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] \end{aligned}$$

$$\mathbb{E}_q \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right]$$

$$\begin{aligned} & \mathbb{E}_q \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] \\ = & \mathbb{E}_q \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{q(x_{1:T}|x_0)} \right] \end{aligned}$$

$$\begin{aligned} & \mathbb{E}_q \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] \\ &= \mathbb{E}_q \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{q(x_{1:T}|x_0)} \right] \\ &= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q[\log q(x_{1:T}|x_0)] \end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_q \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] \\
&= \mathbb{E}_q \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{q(x_{1:T}|x_0)} \right] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q[\log q(x_{1:T}|x_0)] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q \left[\log \prod_{t=1}^T q(x_t|x_{t-1}) \right]
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_q \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] \\
&= \mathbb{E}_q \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{q(x_{1:T}|x_0)} \right] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q[\log q(x_{1:T}|x_0)] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q \left[\log \prod_{t=1}^T q(x_t|x_{t-1}) \right] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \sum_{t=1}^T \mathbb{E}_q[\log q(x_t|x_{t-1})]
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_q \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] \\
&= \mathbb{E}_q \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{q(x_{1:T}|x_0)} \right] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q[\log q(x_{1:T}|x_0)] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \mathbb{E}_q \left[\log \prod_{t=1}^T q(x_t|x_{t-1}) \right] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q[\log p_\theta(x_{t-1}|x_t)] - \sum_{t=1}^T \mathbb{E}_q[\log q(x_t|x_{t-1})] \\
&= \mathbb{E}_q[\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right]
\end{aligned}$$

Decomposing the ELBO

- ▶ $\mathbb{E}_q[\log p(x_T)]$: Validates the model's calibration to the true noise at the last diffusion step, ensuring the reverse process starts accurately.
- ▶ $\sum_{t=1}^T \mathbb{E}_q \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right]$: Acts as a regularizer by comparing the model's backward predictions to the known forward diffusion, guiding the model to correct any discrepancies.

$$\mathcal{L}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right]$$

$$\mathcal{L}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right]$$

By the Markov property:

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

$$\mathcal{L}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right]$$

By the Markov property:

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

and by Bayes' rule:

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

$$\mathcal{L}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right]$$

By the Markov property:

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

and by Bayes' rule:

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

Plugging this equation into the ELBO, we get

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

$$\underbrace{\sum_{t=2}^T \log \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)}}_*$$

The term marked * is a telescoping sum, and can be simplified as follows:

$$\begin{aligned} * &= \log q(x_{T-1} | x_0) + \dots + \log q(x_2 | x_0) + \log q(x_1 | x_0) \\ &\quad - \log q(x_T | x_0) - \log q(x_{T-1} | x_0) - \dots - \log q(x_2 | x_0) \\ &= -\log q(x_T | x_0) + \log q(x_1 | x_0) \end{aligned}$$

$$\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \right. \\ \left. \underbrace{\sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} + \log \frac{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)}}_{*} \right]$$

$$* = -\log q(\mathbf{x}_T|\mathbf{x}_0) + \log q(\mathbf{x}_1|\mathbf{x}_0)$$

Hence the negative ELBO (variational upper bound) becomes

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

Revers Process

We need to learn a neural network to approximate the conditioned probability distributions in the reverse diffusion process,

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

$$q(\mathbf{x}_t \mid \mathbf{x}_0)$$

Recall that the forward process is defined by the equation:

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

We can sample \mathbf{x}_t at any arbitrary time step t in a closed form.

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

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

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

The reverse conditional probability is tractable when conditioned on \mathbf{x}_0 :

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

Using Bayes' rule, we have:

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

where $C(\mathbf{x}_t, \mathbf{x}_0)$ is some function not involving \mathbf{x}_{t-1} .

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

We can express the mean and variance in the following manner: (recall that $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^T \alpha_i$) :

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

Recall $\mathbf{x}_0 = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t)$

$$\begin{aligned} \tilde{\boldsymbol{\mu}}_t &= \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t) \\ &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right) \end{aligned}$$

KL for two Gaussian

$$D_{KL}(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t))$$

$$q(x) = \mathcal{N}(x; \mu_q, \Sigma_q) \quad \text{and} \quad p(x) = \mathcal{N}(x; \mu_p, \Sigma_p),$$

$$D_{KL}(q||p) = \frac{1}{2} \left(\text{tr}(\Sigma_p^{-1}\Sigma_q) + (\mu_p - \mu_q)^T \Sigma_p^{-1}(\mu_p - \mu_q) - k + \ln \left(\frac{\det \Sigma_p}{\det \Sigma_q} \right) \right)$$

for q

$$\tilde{\mu}_t(x_t, x_0) = \frac{1}{\sqrt{\alpha_t}} \left(x_+ - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

for p

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right)$$

Loss fnction

$$\begin{aligned}
 L_t &= E_{\mathbf{x}_0, \epsilon} \left[\frac{1}{2 \|\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)\|_2^2} \|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2 \right] \\
 &= E_{\mathbf{x}_0, \epsilon} \left[\frac{1}{2 \|\boldsymbol{\Sigma}_\theta\|_2^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right) - \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) \right\|^2 \right] \\
 &= E_{\mathbf{x}_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t (1 - \bar{\alpha}_t) \|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \right] \\
 &= E_{\mathbf{x}_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t (1 - \bar{\alpha}_t) \|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t, t)\|^2 \right]
 \end{aligned}$$

Loss function

- Ho et al. (2020) discovered empirically that the diffusion model produces higher-quality images when using a simplified objective, omitting the weighting term

$$\begin{aligned} L_{\text{simple}} &= E_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[\|\epsilon_t - \epsilon_{\theta}(\mathbf{x}_t, t)\|^2 \right] \\ &= E_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[\|\epsilon_t - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right] \end{aligned}$$

Training a DDPM model with L_{simple} .

```
1 while not converged do  
2    $\mathbf{x}_0 \sim q_0(\mathbf{x}_0)$   
3    $t \sim \text{Unif}(\{1, \dots, T\})$   
4    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
5   Take gradient descent step on  $\nabla_{\boldsymbol{\theta}} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$ 
```

Training a DDPM model with L_{simple} .

```
1 while not converged do
2    $\mathbf{x}_0 \sim q_0(\mathbf{x}_0)$ 
3    $t \sim \text{Unif}(\{1, \dots, T\})$ 
4    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5   Take gradient descent step on  $\nabla_{\boldsymbol{\theta}} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t)\|^2$ 
```

Sampling from a DDPM model.

```
1  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2 foreach  $t = T, \dots, 1$  do
3    $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) \right) + \sigma_t \boldsymbol{\epsilon}_t$ 
5 Return  $\mathbf{x}_0$ 
```
