

UNIT 5

Diffusion Models



Johannes Brandstetter

Institute for Machine Learning

Copyright statement:

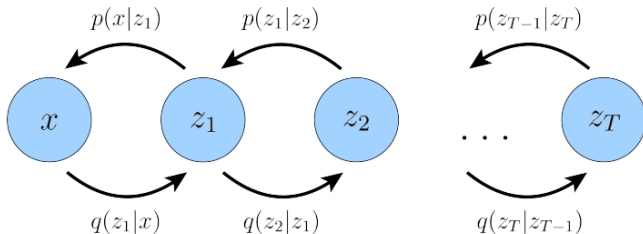
This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(\mathbf{x}) &= \log p(\mathbf{x}) \times \int q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\&= \int \log p(\mathbf{x}) \times q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\&= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})] \\&= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right] \\&= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \times \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\&= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]}_{\text{ELBO}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \right]}_{\mathbb{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))}\end{aligned}$$

Hierarchical VAE

- HVAE extends to multiple hierarchies over latent variables
- Latent variables are interpreted as generated from other higher-level, i.e., more abstract latents
- Special case: **Markovian HVAE** \Rightarrow generative process is a Markov chain



Markovian HVAE

- The generative process is modeled as a Markov chain:
each latent z_t is generated only from the previous latent z_{t+1}
- Joint distribution:

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

- The encoder needs to model:

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

ELBO of Hierarchical VAE

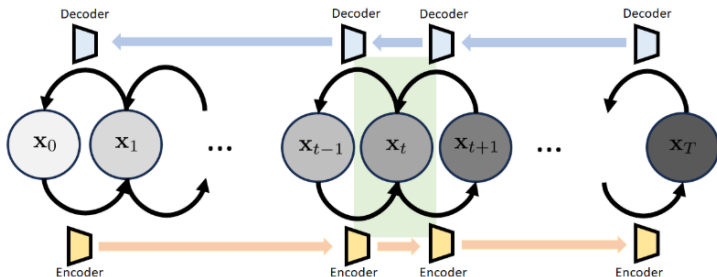
$$\begin{aligned}\log p(\mathbf{x}) &= \log p(\mathbf{x}) \times \int q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})d\mathbf{z}_{1:T} \\&= \int \log p(\mathbf{x}) \times q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})d\mathbf{z}_{1:T} \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} [\log p(\mathbf{x})] \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})}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

From HVAE to Diffusion

- The latent dimension is exactly equal to the data dimension
- The structure of the latent encoder at each timestep is not learned; it is pre-defined as a linear Gaussian model. In other words, it is a Gaussian distribution **centered around the output of the previous timestep**
- The Gaussian parameters of the latent encoders vary over time in such a way that the distribution of the latent at **final timestep T is a standard Gaussian**

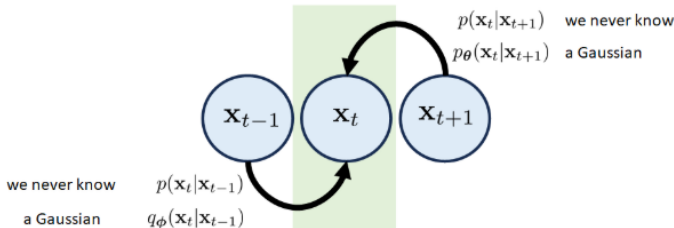
Diffusion through the lense of VAEs

- The input image is x_0 and the white noise is x_T
- The intermediate states x_1, \dots, x_{T-1} are latent variables
- The **forward transition** from x_{t-1} to x_t is analogous to the encoder, the **backward transition** from x_t to x_{t-1} is analogous to the decoder
- Input and output dimensions of the encoder/decoder are identical



Transition blocks

- Holds for $t = 1, \dots, T - 1$
- Just like in a VAE, the encoder (forward) transition distribution $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ is never accessible
 - We approximate it by a Gaussian $q_\phi(\mathbf{x}_t | \mathbf{x}_{t-1})$
- backward (decoder) transition goes from \mathbf{x}_{t+1} to \mathbf{x}_t
 - We approximate it by another Gaussian $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$, but this time we need to estimate the (means of the) distribution by a neural net



Initial and final block

- **Initial block:** Since the diffusion process starts at \mathbf{x}_0 , we only need to worry about the backward transition $p(\mathbf{x}_0|\mathbf{x}_1)$
 - Since $p(\mathbf{x}_0|\mathbf{x}_1)$ is never accessible, we approximate it by a Gaussian $p_\theta(\mathbf{x}_0|\mathbf{x}_1)$, where we again estimate the mean by a neural network
- **Final block:** There is only a forward transition from \mathbf{x}_{T-1} to \mathbf{x}_T , which is approximated by a Gaussian $q_\phi(\mathbf{x}_T|\mathbf{x}_{T-1})$

Transition distribution

- In a denoising diffusion probabilistic model, the transition distribution $q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})$ is defined as:

$$q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}) \stackrel{\text{def}}{=} \mathcal{N}(\sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

- To proof this transition distribution, we start with two scalars $a, b \in \mathbb{R}$, and define the transition distribution as:

$$q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(a\mathbf{x}_{t-1}, b^2\mathbf{I})$$

The equivalent sampling step is $\mathbf{x}_t = a\mathbf{x}_{t-1} + b\epsilon_{t-1}$

- For a random variable $X \sim \mathcal{N}(\mu, \sigma^2)$, drawing X from this distribution is equivalent to $X = \mu + \sigma\eta$ where $\eta \sim \mathcal{N}(0, 1)$

Transition distribution

$$\begin{aligned}x_t &= ax_{t-1} + b\epsilon_{t-1} \\&= a(ax_{t-2} + b\epsilon_{t-2}) + b\epsilon_{t-1} \\&= a^2x_{t-2} + ab\epsilon_{t-2} + b\epsilon_{t-1} \\&\vdots \\&= a^t x_0 + b \underbrace{[\epsilon_{t-1} + a\epsilon_{t-2} + a^2\epsilon_{t-3} + \dots + a^{t-1}\epsilon_0]}_{=w_t}\end{aligned}$$

$\mathbb{E}[x_t] = a^t \mathbb{E}[x_0] \Rightarrow$ for $a < 1$, $\mathbb{E}[x_t]$ goes towards zero for $t \rightarrow \infty$

$$\begin{aligned}\text{Cov}[x_t] &= 0 + \text{Cov}(w_t) \\&= b^2 (\text{Cov}(\epsilon_{t-1}) + a^2 \text{Cov}(\epsilon_{t-2}) + \dots + (a^{t-1})^2 \text{Cov}(\epsilon_0)) \mathbf{I} \\&= b^2 \frac{1 - a^{2t}}{1 - a^2} \mathbf{I}\end{aligned}$$

Transition distribution

- We want $\text{Cov}[\mathbf{x}_t] \underset{t \rightarrow \infty}{=} \frac{b^2}{1 - a^2} \mathbf{I} = \mathbf{I}$
 - I.e., the distribution $\text{Cov}[\mathbf{x}_t]$ should approach $\mathcal{N}(0, \mathbf{I})$
- Setting $a = \sqrt{\alpha}$ this yields $b = \sqrt{1 - \alpha}$
- This gives:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha} \mathbf{x}_{t-1} + \sqrt{1 - \alpha} \boldsymbol{\epsilon}_{t-1} \\ q_\phi(\mathbf{x}_t | \mathbf{x}_{t-1}) &= \mathcal{N}(\sqrt{\alpha} \mathbf{x}_{t-1}, (1 - \alpha) \mathbf{I})\end{aligned}$$

- If a scheduler is used, α is replaced by α_t

Conditional distribution

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \underbrace{\sqrt{\alpha_t} \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}}_{\mathbf{w}_1}\end{aligned}$$

- Sum of two Gaussians remains a Gaussian \Rightarrow we just need to re-calculate its new covariance: $\text{Cov}(\mathbf{w}_1) = \left[(\sqrt{\alpha_t} \sqrt{1 - \alpha_{t-1}})^2 + (\sqrt{1 - \alpha_t})^2 \right] \mathbf{I} = [1 - \alpha_t \alpha_{t-1}] \mathbf{I}$

Conditional distribution

■ Therefore:

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \\ &= \sqrt{\alpha_t \alpha_{t-1} \alpha_{t-2}} \mathbf{x}_{t-3} + \sqrt{1 - \alpha_t \alpha_{t-1} \alpha_{t-2}} \boldsymbol{\epsilon}_{t-3} \\ &\vdots \\ &= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \end{aligned}$$

■ Consequently, for $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 ,$$

■ The distribution $q_\phi(\mathbf{x}_t | \mathbf{x}_0)$ can be written as:

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

Evidence lower bound

$$\begin{aligned}\log p(\mathbf{x}_0) &= \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\ &= \log \int p(\mathbf{x}_{0:T}) \frac{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\ &= \log \int q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0) \frac{p(\mathbf{x}_{0:T})}{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\ &= \log \mathbb{E}_{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\frac{p(\mathbf{x}_{0:T})}{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \geq \mathbb{E}_{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]\end{aligned}$$

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

Evidence lower bound

$$\mathbb{E}_{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p(\mathbf{x}_0|\mathbf{x}_1)}{q_\phi(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] = \underbrace{\mathbb{E}_{q_\phi(\mathbf{x}_1|\mathbf{x}_0)} [\log p(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction}} - \underbrace{\mathbb{E}_{q_\phi(\mathbf{x}_T, \mathbf{x}_{T-1}|\mathbf{x}_0)} [\mathbb{D}_{\text{KL}}(q_\phi(\mathbf{x}_T|\mathbf{x}_{T-1})||p(\mathbf{x}_T))]}_{\text{prior matching}}$$

$$\begin{aligned} \mathbb{E}_{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{\prod_{t=1}^{T-1} p(\mathbf{x}_t|\mathbf{x}_{t+1})}{\prod_{t=1}^{T-1} q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] &= \sum_{t=1}^{T-1} \mathbb{E}_{q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_t|\mathbf{x}_{t+1})}{q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\ &= \sum_{t=1}^{T-1} \mathbb{E}_{q_\phi(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_t|\mathbf{x}_{t+1})}{q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\ &= \underbrace{\sum_{t=1}^{T-1} \mathbb{E}_{q_\phi(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} \left[\mathbb{D}_{\text{KL}}(q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})||p(\mathbf{x}_t|\mathbf{x}_{t+1})) \right]}_{\text{Consistency}} \end{aligned}$$

- We replace $p(\mathbf{x}_0|\mathbf{x}_1)$ with a neural network $p_\theta(\mathbf{x}_0|\mathbf{x}_1)$ and $p(\mathbf{x}_t|\mathbf{x}_{t+1})$ with a neural network $p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})$

Bayes' theorem on consistency term

- Problem: $q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})$ and $p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})$ run in opposite directions
- Solution via Bayes' theorem:

$$q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}) = \frac{q_\phi(\mathbf{x}_{t-1}|\mathbf{x}_t)q_\phi(\mathbf{x}_t)}{q(\mathbf{x}_{t-1})}$$

- We add a conditioning on \mathbf{x}_0 :

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

Bayes' theorem on consistency term

$$q_{\phi}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{\mathcal{N}(\sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I}) \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})}$$

- Minimizing this distribution (without proof) yields the mean of the distribution, setting the second derivative to zero yields the standard deviation:

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t}\mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t}\mathbf{x}_0$$

$$\Sigma_q(t) = \frac{(1 - \alpha_t)(1 - \sqrt{\hat{\alpha}_{t-1}})}{1 - \bar{\alpha}_t}\mathbf{I}$$

- $q_{\phi}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is completely characterized by \mathbf{x}_t and \mathbf{x}_0 .
No NN required to estimate the mean and variance!

New ELBO

$$\begin{aligned} \text{ELBO}_{\phi, \theta}(\mathbf{x}) = & \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)]}_{\text{reconstruction}} - \underbrace{\mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{x}_T | \mathbf{x}_{T-1}) || p(\mathbf{x}_T))}_{\text{new prior matching}} \\ & - \underbrace{\sum_{t=2}^T \mathbb{E}_{q_{\phi}(\mathbf{x}_t | \mathbf{x}_0)} [\mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))]}_{\text{new consistency}} \end{aligned}$$

- **Reconstruction**: we are maximizing the likelihood
- **Prior matching**: nothing to train
- **Consistency**: Instead of asking a forward transition to match with a reverse transition, we use q_{ϕ} to construct a reverse transition and use it to match with p_{θ}

Consistency term

$$q_\phi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\underbrace{\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)}_{\text{known}}, \underbrace{\sigma_q^2(t)\mathbf{I}}_{\text{known}}\right)$$
$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}\left(\underbrace{\boldsymbol{\mu}_\theta(\mathbf{x}_t)}_{\text{not known}}, \underbrace{\sigma_q^2(t)\mathbf{I}}_{\text{known}}\right)$$

- The KL divergence is simplified to:

$$\begin{aligned} \mathbb{D}_{\text{KL}}(q_\phi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &= \mathbb{D}_{\text{KL}}\left(\mathcal{N}(\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0), \sigma_q^2(t)\mathbf{I}) || \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_t), \sigma_q^2(t)\mathbf{I})\right) \\ &= \frac{1}{2\sigma_q^2(t)} \|\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t)\|^2 \end{aligned}$$

ELBO for training

- We need to find a network that minimizes:

$$\frac{1}{2\sigma_q^2(t)} \|\mu_q(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t)\|^2, \text{ where}$$

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_0$$

- To make it more convenient, we choose:

$$\mu_\theta(\mathbf{x}_t) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \underbrace{\hat{\mathbf{x}}_\theta(\mathbf{x}_t)}_{\text{network}}$$

- This yields:

$$\begin{aligned} & \frac{1}{2\sigma_q^2(t)} \|\mu_q(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t)\|^2 \\ &= \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\hat{\mathbf{x}}_\theta(\mathbf{x}_t) - \mathbf{x}_0\|^2 \end{aligned}$$

ELBO for training

- For DDPM training, we use:

$$\begin{aligned}\text{ELBO}_\theta(\mathbf{x}) &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction}} \\ &\quad - \underbrace{\sum_{t=2}^T \mathbb{E}_{q_\phi(\mathbf{x}_t|\mathbf{x}_0)} [\mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{consistency}} \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction}} \\ &\quad - \sum_{t=2}^T \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\hat{\mathbf{x}}_\theta(\mathbf{x}_t) - \mathbf{x}_0\|^2\end{aligned}$$

- We use:

$$\begin{aligned}\log p_\theta(\mathbf{x}_0|\mathbf{x}_1) &= \log \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_1), \sigma_q^2(1)\mathbf{I}) \propto -\frac{1}{2\sigma_q^2(1)} \|\boldsymbol{\mu}_\theta(\mathbf{x}_1) - \mathbf{x}_0\|^2 \\ &= -\frac{1}{2\sigma_q^2(1)} \|\hat{\mathbf{x}}_\theta(\mathbf{x}_1) - \mathbf{x}_0\|^2\end{aligned}$$

ELBO for training

- Using $\alpha_0 = 1$ and $\hat{\alpha}_1 = \alpha_1$, we obtain:

$$\text{ELBO}_\theta = - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\hat{\mathbf{x}}_\theta(\mathbf{x}_t) - \mathbf{x}_0\|^2 \right]$$

- Therefore, the training objective minimizes:

$$\theta^* = \operatorname{argmin} \sum_{t=1}^T \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \|\hat{\mathbf{x}}_\theta(\mathbf{x}_t) - \mathbf{x}_0\|^2$$

- Ignoring the constants and expectations, the objective is:

$$\operatorname{argmin}_{\theta} \|\hat{\mathbf{x}}_\theta(\mathbf{x}_t) - \mathbf{x}_0\|$$

DDPM training

- $\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}$: We are not denoising any random noisy image. Instead, we the noisy image is fully defined as:
$$\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$
- $\frac{1}{\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2}$: controls the relative emphasis on each denoising loss. However, only minor impact in practice, usually dropped
- $\sum_{t=1}^T$: the summation can be replaced by a uniform distribution $t \sim \text{Uniform}[1, T]$

DDPM training

- For every image x_0 in your training set:
 - Repeat the following steps until convergence
 - Pick a random time stamp $t \sim \text{Uniform}[1, T]$
 - Draw a sample $x_t \sim q(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$
 - Take gradient descent step on:

$$\nabla_{\theta} \|\hat{x}_{\theta}(x_t) - x_0\|^2$$

DDPM inference

- We start with white noise vector $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$
- For $t = T, T - 1, \dots, 1$:
 - Calculate $\hat{\mathbf{x}}_\theta(\mathbf{x}_t)$ using our trained denoiser
 - Update \mathbf{x}_t following:

$$\begin{aligned}\mathbf{x}_{t-1} &\sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \\ &= \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \hat{\mathbf{x}}_\theta(\mathbf{x}_t) + \sigma_q(t)\mathbf{z}\end{aligned}$$

- $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$
- Remember, we chose:

$$\mu_\theta(\mathbf{x}_t) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \underbrace{\hat{\mathbf{x}}_\theta(\mathbf{x}_t)}_{\text{network}}$$

DDPM training - based on noise prediction

- Reformulate diffusion to predict the noise instead of the signal:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0 \Rightarrow \mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0}{\sqrt{\bar{\alpha}_t}}$$

- Substituting into

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_0$$

gives (after several steps):

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0$$

- Thus, we design our mean estimator to match this form:

$$\mu_\theta(\mathbf{x}_t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_\theta(\mathbf{x}_t)$$

DDPM training - based on noise prediction

- In Ho et al, the new training objective becomes:

$$\text{ELBO}_{\theta} = - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\hat{\epsilon}_{\theta}(\epsilon_t) - \epsilon_0\|^2 \right]$$

- For every image \mathbf{x}_0 in your training set:
 - Repeat the following steps until convergence
 - Pick a random time stamp $t \sim \text{Uniform}[1, T]$
 - Draw a sample $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$
 - Take gradient descent step on:

$$\nabla_{\theta} \|\hat{\epsilon}_{\theta}(\mathbf{x}_t) - \epsilon_0\|^2$$

DDPM inference - based on noise prediction

- The inference step can be derived through:

$$\begin{aligned}\mathbf{x}_{t-1} &\sim p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_t), \sigma_q^2(t)\mathbf{I}) \\ &= \boldsymbol{\mu}_{\theta}(\mathbf{x}_t) + \sigma_q^2(t)\mathbf{z} \\ &= \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t) + \sigma_q(t)\mathbf{z}\end{aligned}$$

- We start with white noise vector $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$

- For $t = T, T - 1, \dots, 1$:

- Calculate $\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t)$ using our trained denoiser
- Update \mathbf{x}_t following:

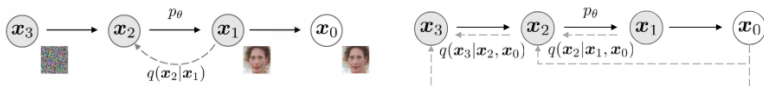
$$\begin{aligned}\mathbf{x}_{t-1} &\sim p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) \\ &= \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t) + \sigma_q(t)\mathbf{z}\end{aligned}$$

DDPMs are slow

- One of the most prevalent drawbacks of DDPM is that they need a large number of iterations to generate a reasonably good looking samples.
- If the reverse diffusion process intrinsically requires many steps to converge, then it will take us many denoising steps. Therefore, to speed up the computing, it is necessary to reduce the number of iterations.

Generalization of DDPMs

- Generalization of DDPMs via a non-Markovian forward process



- Song et al propose a family of backward processes:

$$q_\sigma(x_{t-1}|x_t, x_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}x_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{x_t - \alpha_t x_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right)$$

- σ_t modulates the stochasticity of the process:

□ If $\sigma_t = \sqrt{\frac{(1 - \alpha_{t-1})}{(1 - \alpha_t)}} \sqrt{1 - \frac{\alpha_t}{\alpha_{t-1}}}$, we obtain the original DDPM

□ If $\sigma_t = 0$, the forward process becomes Markovian and the backward process becomes deterministic; the same original noise leads to the same image

Derivation of DDIMs (sketch)

- The original DDPM transition probabilities are:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}\left(\mathbf{x}_t|\sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I}\right)$$
$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}\right)$$

- One important observation is that the transition probability $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ follows a Markov chain
- The Markov property ensures that the system is memoryless, i.e., once we know \mathbf{x}_{t-1} , we know \mathbf{x}_t
- Downside: a Markov chain can take many steps to converge

Derivation of DDIMs (sketch)

- Choice: $\alpha_t = \frac{\alpha_t}{\alpha_{t-1}}$
- I.e., $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}\left(\mathbf{x}_t|\sqrt{\frac{\alpha_t}{\alpha_{t-1}}}\mathbf{x}_{t-1}, (1 - \frac{\alpha_t}{\alpha_{t-1}})\mathbf{I}\right)$
- The product simplifies to $\bar{\alpha}_t = \prod_{i=1}^t \frac{\alpha_i}{\alpha_{i-1}} = \alpha_t$
- Thus, $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t|\sqrt{\alpha_t}\mathbf{x}_0, (1 - \alpha_t)\mathbf{I}\right)$
- $\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\epsilon$ and $\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1}}\epsilon$
- $\Rightarrow \mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1}}\left(\frac{\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_0}{\sqrt{1 - \alpha_t}}\right)$
- We have replaced the Gaussian by an estimate

Derivation of DDIMs (sketch)

- Given:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_t}\mathbf{x}_0, (1 - \alpha_t)\mathbf{I}\right)$$

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

- Can we ensure that

$$q(\mathbf{x}_{t-1}|\mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_0, (1 - \alpha_{t-1})\mathbf{I}\right)?$$

- The solution is a family of backward processes:

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

DDIM

- Special case when $\sigma_t = 0$ for all t
- The forward process becomes deterministic given x_{t-1} and x_0 , except for $t = 1$
- Additionally, no random noise σ_q in the reverse process
- The resulting model becomes an **implicit probabilistic model**, where samples are generated from latent variables with a fixed procedure (from x_T to x_0)
- This is named **denoising diffusion implicit model** (DDIM), because it is an implicit probabilistic model trained with the DDPM objective (despite the forward process no longer being a diffusion)

Generalized training objective

- Training objectives are equivalent for any value of σ_t :

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_0, q(\mathbf{x}_t|\mathbf{x}_0)} [\mathbb{D}_{\text{KL}}(q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] \\ &= \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_0 \left[\frac{\|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t)\|^2}{2\sigma_t^2 \alpha_t} \right] \end{aligned}$$

- A model trained for the original DDPM process can be used for any process of the family

Generalized sampling method

- Sampling becomes:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \hat{\epsilon}_\theta(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\mathbf{x}_t) + \sigma_t \epsilon_t$$

- Different choices of σ_t result in different generative processes, all for the same model ϵ_θ
- **Accelerated generative processes**: we define the forward process not on all the latent variables $\mathbf{x}_{1:T}$, but only on a subset $\{\mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_S}\}$, where τ is an increasing sub-sequence of $[1, \dots, T]$ of length S
 - $q(\mathbf{x}_{\tau_i} | \mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_{\tau_i}} \mathbf{x}_0, (1 - \alpha_{\tau_i}) \mathbf{I})$ matches the “marginals”
 - This achieves significant speed-ups when the length of the sampling trajectory is much smaller than T

Relevance to neural ODEs

- We can rewrite the sampling to:

$$\frac{\mathbf{x}_{t-\Delta t}}{\sqrt{\alpha_{t-\Delta t}}} = \frac{\mathbf{x}_t}{\sqrt{\alpha_t}} + \left(\sqrt{\frac{1 - \alpha_{t-\Delta t}}{\alpha_{t-\Delta t}}} - \sqrt{\frac{1 - \alpha_t}{\alpha_t}} \right) \hat{\epsilon}_{\theta}(\mathbf{x}_t)$$

- With enough discretization steps, we can reverse the generation process (going from $t = 0$ to T)
 - This encodes \mathbf{x}_0 to \mathbf{x}_T and simulates the reverse of the ODE
 - We can use DDIM to obtain encodings of the observations (as the form of \mathbf{x}_T)
 - Useful for other downstream applications that require latent representations of a model

Literature

■ Important papers:

- ☐ First 2015 diffusion paper
- ☐ Jonathan Ho - DDPM
- ☐ Jiaming Song - DDIM
- ☐ Latent diffusion
- ☐ Hooeboom - look at the appendix Step-by-Step Diffusion: An Elementary Tutorial