

## Generative models: VAE and DM

$$p(x, y) \quad , \quad p(x, y) \geq 0 \quad \int p(x, y) dx dy = 1$$

$$\begin{array}{c} y - \text{disc.} \\ x - \text{cont.} \end{array} \quad \sum_y \int p(x, y) dx = 1$$

$$\text{Product rule: } p(x, y) = p(y|x)p(x) = p(x|y)p(y)$$

$$\text{Sum rule: } \int p(x, y) dy = p(x)$$

$$p(x, y, z) \quad \begin{array}{c} x - \text{observed} \\ y - \text{to predict} \\ z - \text{hidden} \end{array} \quad p(y|x) = \frac{p(x, y)}{p(x)} = \frac{\int p(x, y, z) dz}{\int p(x, \hat{y}, \hat{z}) d\hat{y} d\hat{z}}$$

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\mathcal{N}(x | \mu, \Sigma) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$\Sigma^T = \Sigma \geq 0$

$$\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2) \Rightarrow \mathcal{N}(x | \mu, \Sigma) = \prod_{i=1}^d \mathcal{N}(x_i | \mu_i, \sigma_i^2)$$

$$x \sim \mathcal{N}(x | \mu, \Sigma) \quad z \sim \mathcal{N}(z | 0, I), \quad x = \mu + L \cdot z \quad \Sigma = L \cdot L^T$$

$$x \sim \mathcal{N}(x | \mu, \Sigma) \quad z \sim \mathcal{N}(z | 0, I), \quad x = \mu + L \cdot z \quad L \succ (0)$$

$f$ -differentiable

$$\underset{\mu, \Sigma}{\min} E_{N(x|\mu, \Sigma)} f(x)$$

Init.  $\mu_1 \leftarrow$

Iter. until convergence:

$$z_1, \dots, z_B \sim N(z|0, I)$$

$$\text{Loss} = \frac{1}{B} \sum_j f(\mu + z_j)$$

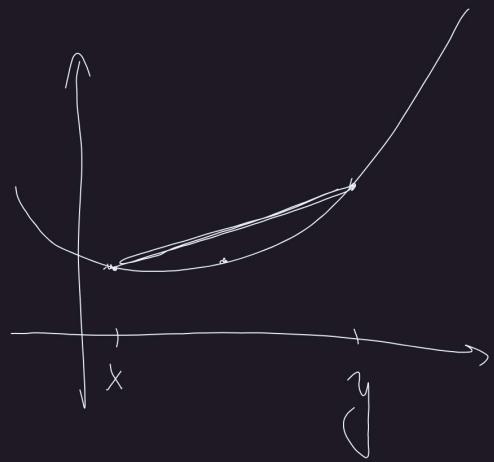
$$\mu, \Sigma = \text{opt\_step}(\nabla_{\mu, \Sigma} \text{Loss})$$

$$\nabla_{\mu} E_{N(x|\mu, \Sigma)} f(x) = \nabla_{\mu} E_{N(z|0, I)} f(\mu + z) =$$

$$= E_{N(z|0, I)} \nabla_{\mu} f(\mu + z) \approx \frac{1}{B} \sum_{j=1}^B \nabla_{\mu} f(\mu + z_j)$$

$$z_j \sim N(z|0, I)$$

$$f - \text{convex} \iff f(\lambda x + (1-\lambda)y) \leq f(x) + (1-\lambda)f(y) \quad \forall x, y \quad \forall \lambda \in (0, 1)$$



Jensen inequality:

$$f\left(\sum_{i=1}^m d_i x_i\right) \leq \sum_{i=1}^m d_i f(x_i)$$

$$d_i \geq 0 \quad \sum_{i=1}^m d_i = 1$$

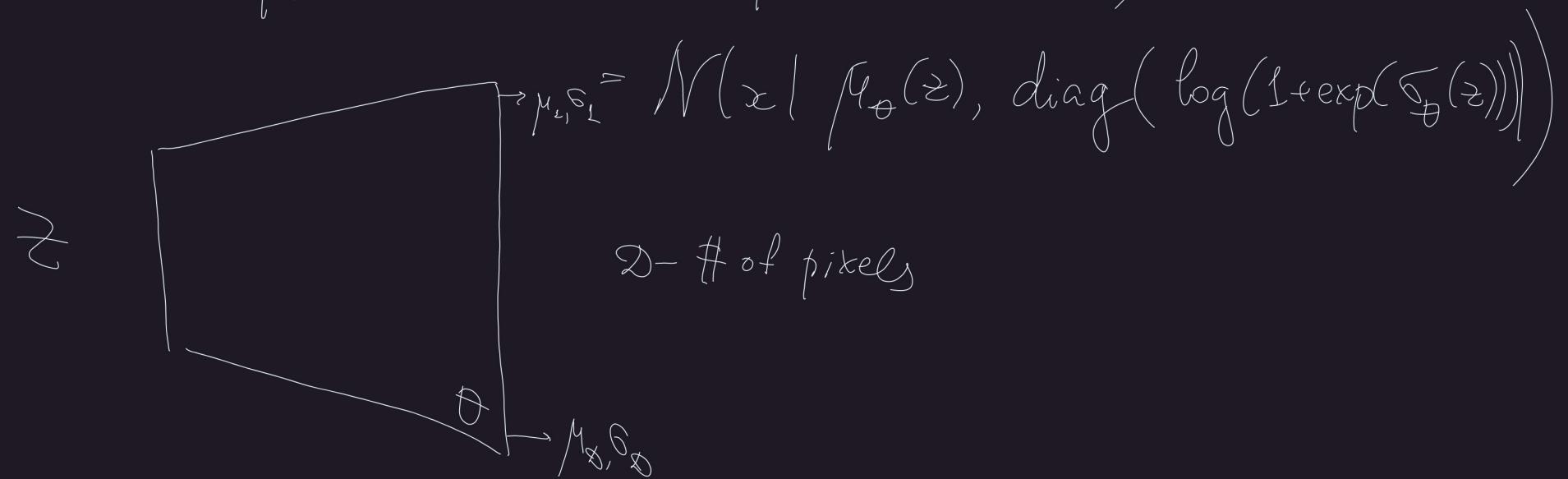
$$f\left(\int \alpha(t) x(t) dt\right) \leq \int \alpha(t) f(x(t)) dt$$

$$\alpha(t) \geq 0 \quad \int \alpha(t) dt = 1$$

VAE

$$z \sim p(z) = \mathcal{N}(z | 0, I)$$

$$x | z \sim p_\theta(x|z) = \mathcal{N}(x | \mu_\theta(z), \Sigma_\theta(z)) =$$

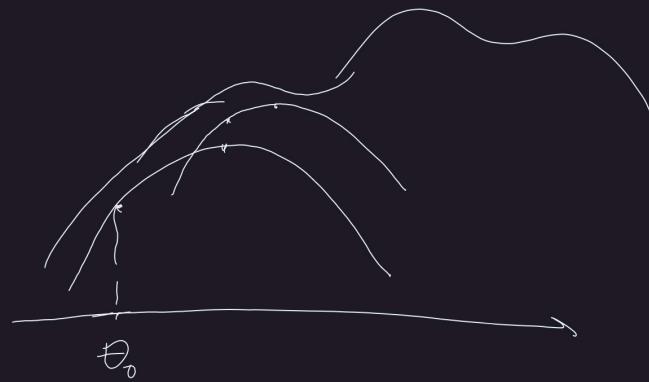
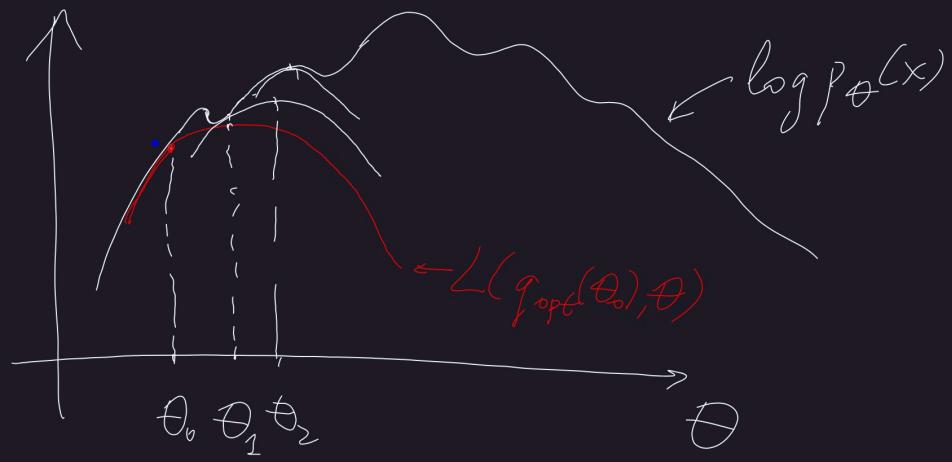


$$P_{\theta}(x, z) = P_{\theta}(x|z) \cdot p(z)$$

$$\mathbb{E}_{x \sim p_{\text{data}}(x)} \log p_{\theta}(x) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \underbrace{\log \int p_{\theta}(x, z) dz}_{\rightarrow \max_{\theta}} \rightarrow \max_{\theta}$$

$$\log p_{\theta}(x) = \log \int p_{\theta}(x, z) dz \stackrel{\substack{\forall q: q(z) \geq 0 \\ \int q(z) dz = 1}}{=} \log \int q(z) \cdot \frac{p_{\theta}(x, z)}{q(z)} dz \stackrel{\substack{\text{can't be computed} \\ \text{Jensen inequality}}}{\geq}$$

$$\geq \int q(z) \log \frac{p_{\theta}(x, z)}{q(z)} dz = \mathbb{E}_{q(z)} \log \frac{p_{\theta}(x, z)}{q(z)} \stackrel{\substack{\rightarrow \max_{\theta, q} \\ \text{ELBO}}}{=} \mathcal{L}(q, \theta)$$



$$q(z|\lambda) \approx p(z|x)$$

$$\frac{1}{N} \sum_{i=1}^N \log p_\theta(x_i) \geq \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q(z_i|\lambda)} \log \frac{p_\theta(x_i|z_i)}{q(z_i|\lambda)} =$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q(z_i|\lambda)} \log \frac{p_\theta(x_i|z_i)p(z_i)}{q(z_i|\lambda)} =$$

①  $q(z_i|\lambda) = N(z_i | \tilde{\mu}_i, \text{diag}(\tilde{\Sigma}_i))$   
 $\forall i = 1 \dots N$

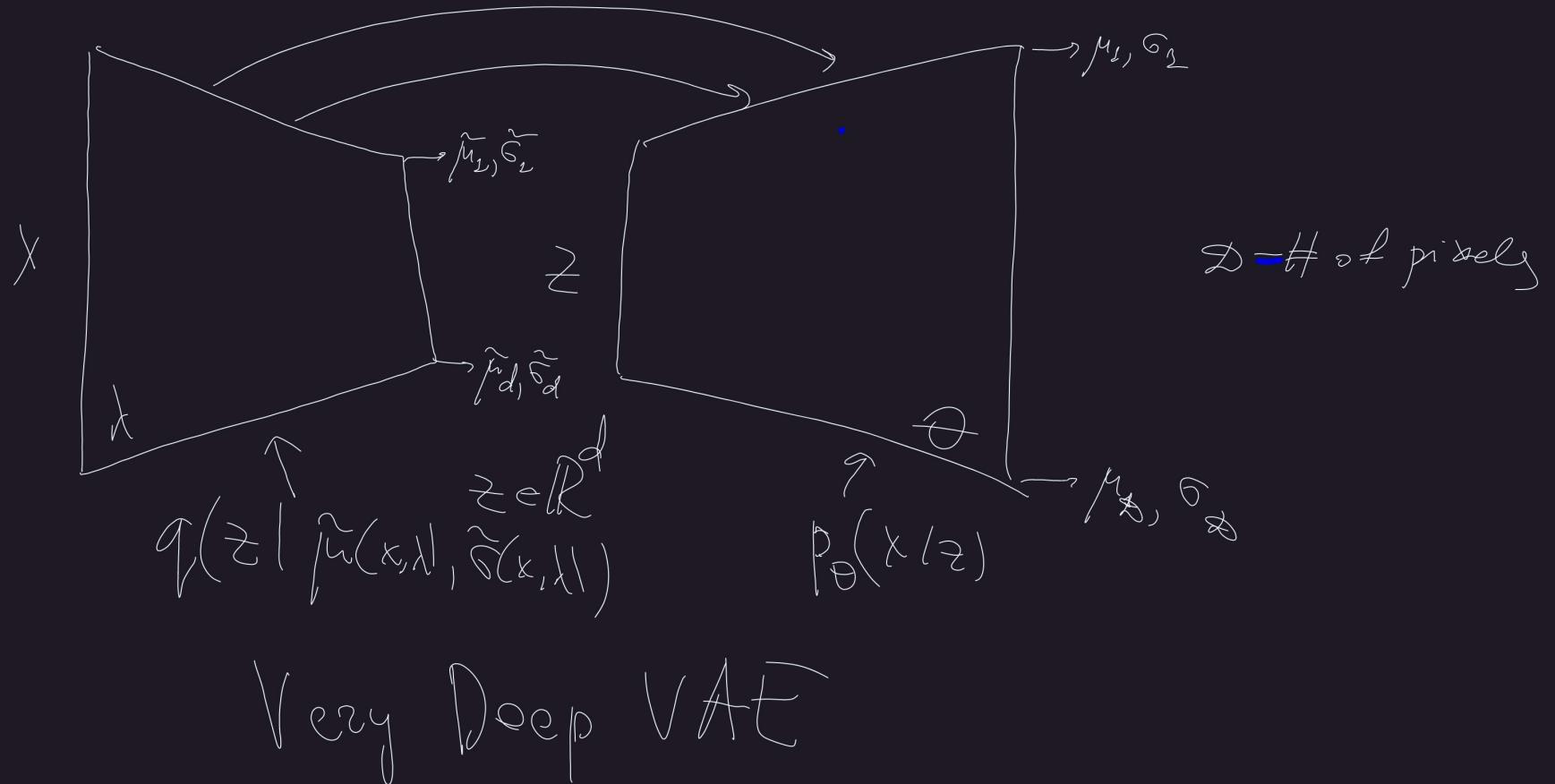
$$= \frac{1}{N} \sum_{i=1}^N \left[ \mathbb{E}_{q(z_i|\lambda)} \log p_\theta(x_i|z_i) + \mathbb{E}_{q(z_i|\lambda)} \log \frac{p(z_i)}{q(z_i)} \right]$$

②  $q(z_i|\lambda) = N(z_i | \tilde{\mu}(x_i|\lambda), \text{diag}(\tilde{\Sigma}(x_i|\lambda)))$

$$= \frac{1}{N} \sum_{i=1}^N \left[ \left( \mathbb{E}_{q(z_i|\lambda)} \log p_\theta(x_i|z_i) \right) - \text{KL}\left(q(z_i|\lambda) \parallel p(z_i)\right) \right]$$

reconstruction term regularization

$\rightarrow \max_{\theta, \lambda}$



$$\frac{1}{N} \sum_{i=1}^N E_{\mathcal{N}(\varepsilon_i | 0, I)} \log p_\theta(x_i \mid \underbrace{\tilde{\mu}(x_i, \lambda) + \tilde{\sigma}(x_i, \lambda) \cdot \varepsilon_i}_{z_i}) - KL\left(\tilde{\mu}(x_i, \lambda), \tilde{\sigma}(x_i, \lambda)\right)$$

VAE training

Init.  $\theta, \lambda$

Iterations until convergence:

$\rightarrow \max_{\theta, \lambda}$

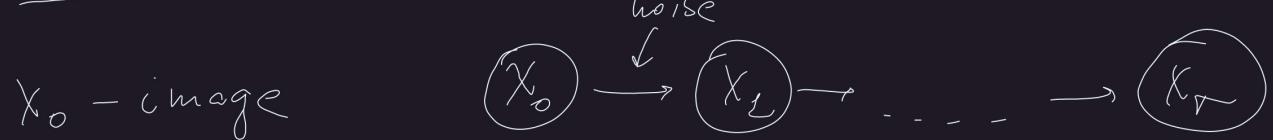
$$x_1 \dots x_B \sim P_{\text{data}}(\omega)$$

$$\varepsilon_1, \dots, \varepsilon_B \sim \mathcal{N}(\varepsilon | 0, I)$$

$$\text{Loss} = \frac{1}{B} \sum_{j=1}^B \left[ \log p_\theta(x_j \mid \tilde{\mu}(x_j, \lambda) + \tilde{\sigma}(x_j, \lambda) \varepsilon_j) - KL\left(\tilde{\mu}(x_j, \lambda), \tilde{\sigma}(x_j, \lambda)\right) \right]$$

$$\theta, \lambda \leftarrow \text{opt-step}\left(\nabla_{\theta, \lambda} \text{Loss}\right)$$

# Diffusion model



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

$$q(x_1, x_2, \dots, x_T | x_0) = q(x_1 | x_0) q(x_2 | x_1) \dots q(x_T | x_{T-1})$$

$$x_t = \sqrt{1-\beta_t} x_{t-1} + \varepsilon_0 \sqrt{\beta_t}, \quad \varepsilon \sim \mathcal{N}(\varepsilon | 0, I)$$

$$\text{Var}(x_{t,i}) = (1-\beta_t) \underbrace{\text{Var}(x_{t-1,i})}_{\frac{1}{1}} + \beta_t \underbrace{\text{Var}(\varepsilon_i)}_{\frac{\beta}{1}} = 1$$

## Linear Gaussian Model

$$x \sim \mathcal{N}(x | \mu, \Sigma)$$

$$y|x \sim \mathcal{N}(y | Ax, \Gamma)$$

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

$$\left( q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\mathcal{D}_t} x_0, (1 - \mathcal{D}_t) I) \right)$$

$$q(x_T | x_0) q(x_{T-1} | x_T, x_0) q(x_{T-2} | x_{T-1}, x_0) \dots q(x_1 | x_2, x_0)$$

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1} | \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I)$$

$$x, y \sim \mathcal{N}$$

$$x|y \sim \mathcal{N}$$

$$y \sim \mathcal{N}(y | A\mu, \Gamma + A\Sigma A^\top)$$

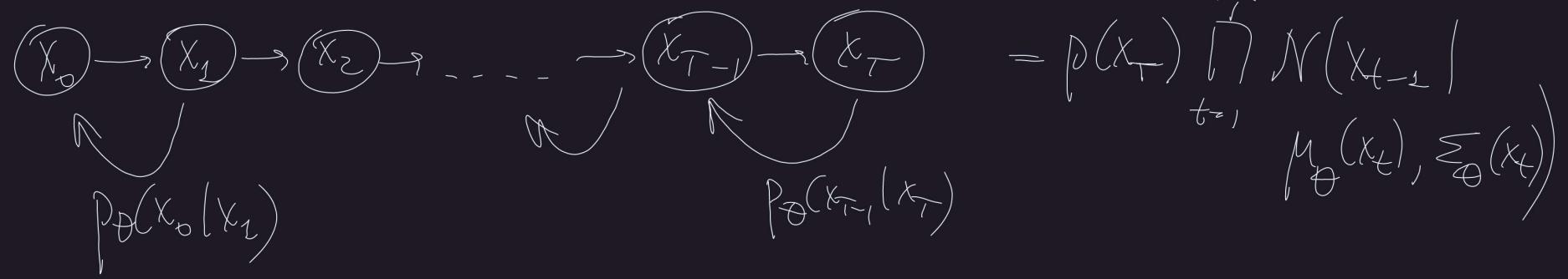
$$\mathcal{D}_t = \prod_{s=1}^t (1 - \beta_s)$$

$$q(x_{1:T} | x_0)$$

$$q(x_{t-1} | x_t, x_0)$$

$$\mathcal{N}(x_T | \phi, I)$$

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



$$p_{\theta}(x_0, x_1, \dots, x_T)$$

$$\log p_{\theta}(x_0) = \log \int p(x_0, \dots) dx_{1:T} \geq E_{q(x_{1:T}|x_0)} \log \frac{p_{\theta}(x_0, \dots)}{q(x_{1:T}|x_0)}$$

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

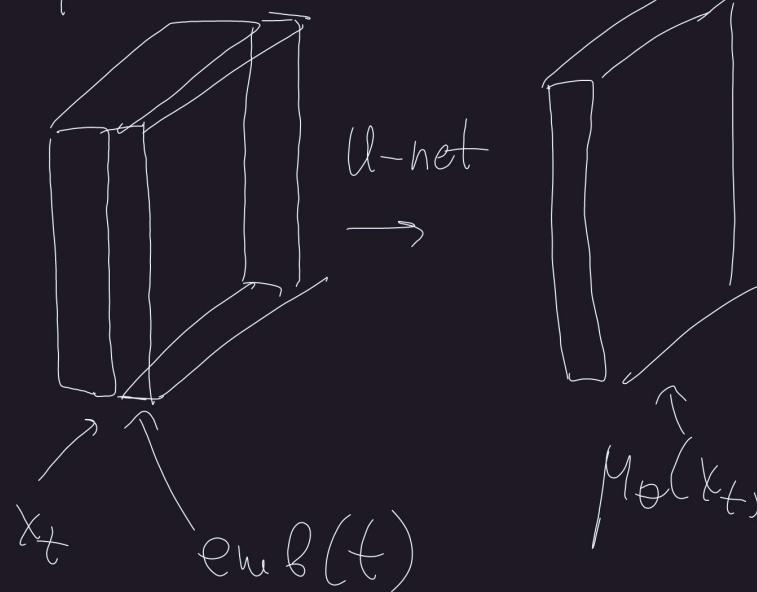
$$= E_q \left[ \log p(x_T) + \sum_{t=1}^T \log \frac{p_{\theta}(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right] \quad \text{=} \quad$$

$$= \left\{ \begin{aligned} q(x_{t-1}|x_t, x_0) &= \frac{q(x_{t-1}, x_t|x_0)}{q(x_t|x_0)} = \cancel{q(x_t|x_{t-1}, x_0)} \cancel{q(x_{t-1}|x_0)} \\ &\Rightarrow q(x_t|x_0) \\ &= \cancel{q(x_t|x_{t-1})} \times \cancel{\frac{q(x_t|x_0)}{q(x_{t-1}|x_0)}} \end{aligned} \right\}$$

$$\begin{aligned}
& \Leftarrow E_q \left[ \log p(x_T) + \sum_{t=2}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \cdot \underbrace{\left( \frac{q(x_{t-1}|x_0)}{q(x_t|x_0)} \right)}_{\log q(x_{t-1}|x_0) - \log q(x_t|x_0)} + \log \frac{p_\theta(x_0|x_1)}{q(x_1|x_0)} \right] = \\
& \quad \log q(x_{T-1}|x_0) - \log q(x_T|x_0) + \log q(x_{T-2}|x_0) - \\
& \quad - \log q(x_{T-1}|x_0) + \dots + \log q(x_1|x_0) - \log q(x_2|x_0) \\
& = E_q \left[ \log \frac{p(x_T)}{q(x_T|x_0)} + \sum_{t=2}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_{t-1}, x_0)} + \log p_\theta(x_0|x_1) \right] = \text{reconst. term} \\
& = - \text{KL}\left( q(x_T|x_0) \parallel p(x_T) \right) - \sum_{t=2}^T E_{q(x_t|x_0)} [\text{KL}\left( q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t) \right)] + E_{q(x_1|x_0)} [\log p_\theta(x_0|x_1)] \\
& \quad \text{regul. term} \quad \rightarrow \max_{\theta}
\end{aligned}$$

$$p_{\theta}(x_{t+1} | x_t) = \mathcal{N}(x_{t+1} | \mu_{\theta}(x_t, t), \sigma_t^2 I)$$

$$\sigma_t^2 = \beta_t \quad (\text{fixed!})$$



$$KL(q(x_{t+1} | x_t, x_0) || p_{\theta}(x_{t+1} | x_t)) = \\ \mathcal{N}(x_{t+1} | \tilde{\mu}(x_t, x_0), \tilde{\sigma}_t^2 I)$$

$$\mu_{\theta}(x_t, t) = \text{const} + \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_{\theta}(x_t, t) \|^2 \\ = x_t(x_0, \varepsilon)$$

$$E_{q(x_t | x_0)} KL \dots = E_{\mathcal{N}(\varepsilon | 0, I)} \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(\underbrace{\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon}_{= x_t(x_0, \varepsilon)}, t) - \mu_{\theta}(\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon, t) \|^2 \oplus$$

$$\text{② } |E_{\mathcal{N}(\varepsilon|0, I)} \frac{1}{2\zeta_t^2} \left( \left| \frac{\beta_t}{\sqrt{\lambda_t}} (x_t(x_0, \varepsilon) - \frac{\beta_t}{\sqrt{1-\lambda_t}} \varepsilon) - \mu_\theta(x_t(x_0, \varepsilon), t) \right|^2 \right| =$$

$$\mu_\theta(x_t, t) := \frac{1}{\sqrt{\lambda_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\lambda_t}} \varepsilon_\theta(x_t, t) \right)$$

$$= |E_{\mathcal{N}(\varepsilon|0, I)} \frac{\beta_t^2}{2\zeta_t^2 \lambda_t(1-\lambda_t)} \| \varepsilon - \varepsilon_\theta(\underbrace{\sqrt{\lambda_t} x_0 + \sqrt{1-\lambda_t} \varepsilon}_x, t) \|^2|$$

$$\mathcal{L}_{\text{Simple}}(\theta) = |E_{t, x_0, \varepsilon} \| \varepsilon - \varepsilon_\theta(\sqrt{\lambda_t} x_0 + \sqrt{1-\lambda_t} \varepsilon, t) \|^2|$$

DDPM training:

Init.  $\theta$

Iter. until convergence:

$$x_0 \sim p_{\text{data}}(x_0)$$

$$t \sim \text{Uniform}(1, 2, \dots, T)$$

$$\varepsilon \sim \mathcal{N}(\varepsilon | 0, I)$$

$$\text{Loss} = \left\| \varepsilon - \varepsilon_\theta \left( \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon, t \right) \right\|^2$$

$$\theta = \text{opt-step}(\nabla_\theta \text{Loss})$$

Sampling:

$$x_T \sim \mathcal{N}(x_T | 0, I)$$

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

$$\zeta \sim \mathcal{N}(\zeta | 0, I)$$

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon_\theta(x_t, t) \right) + \zeta \cdot \sigma_t$$

$$+ \zeta \cdot \sigma_t$$

Output:  $x_0$