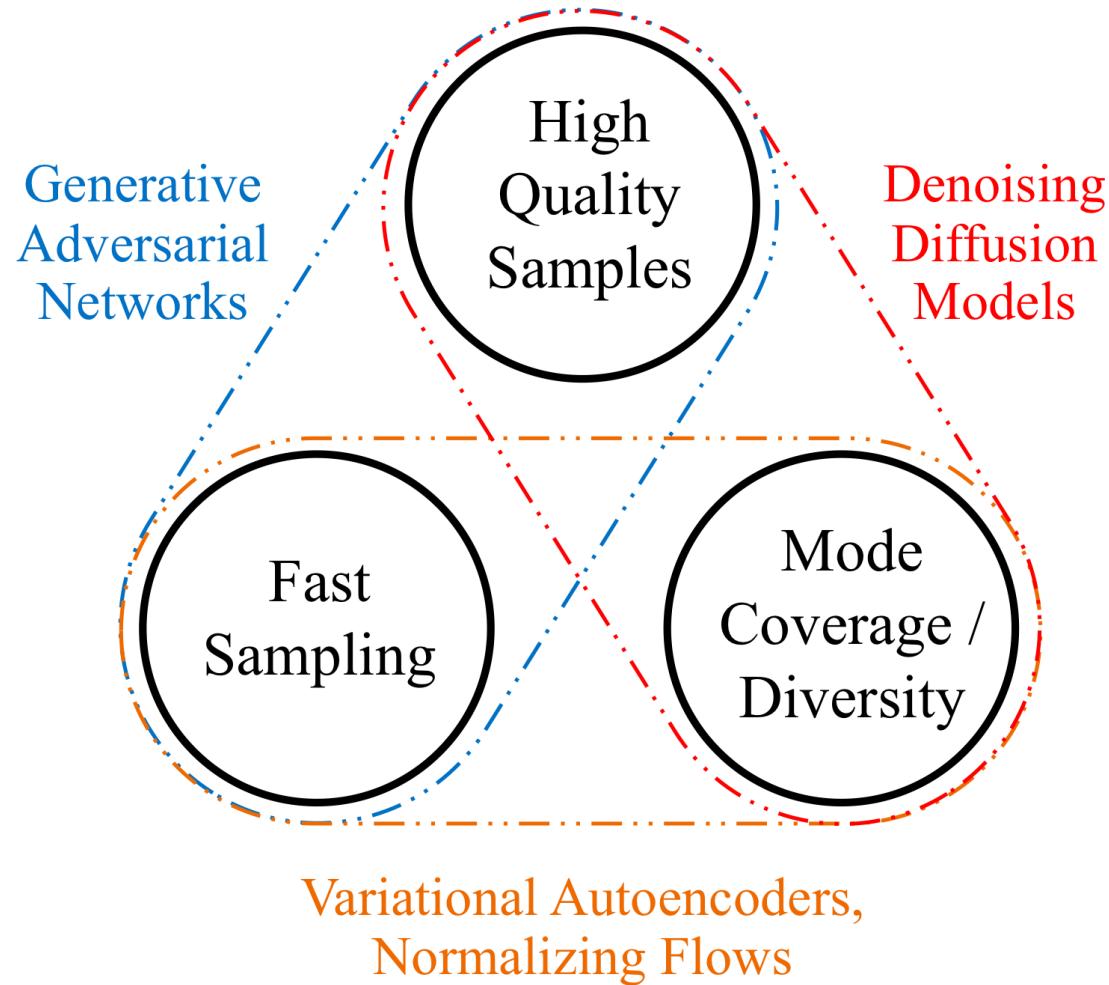


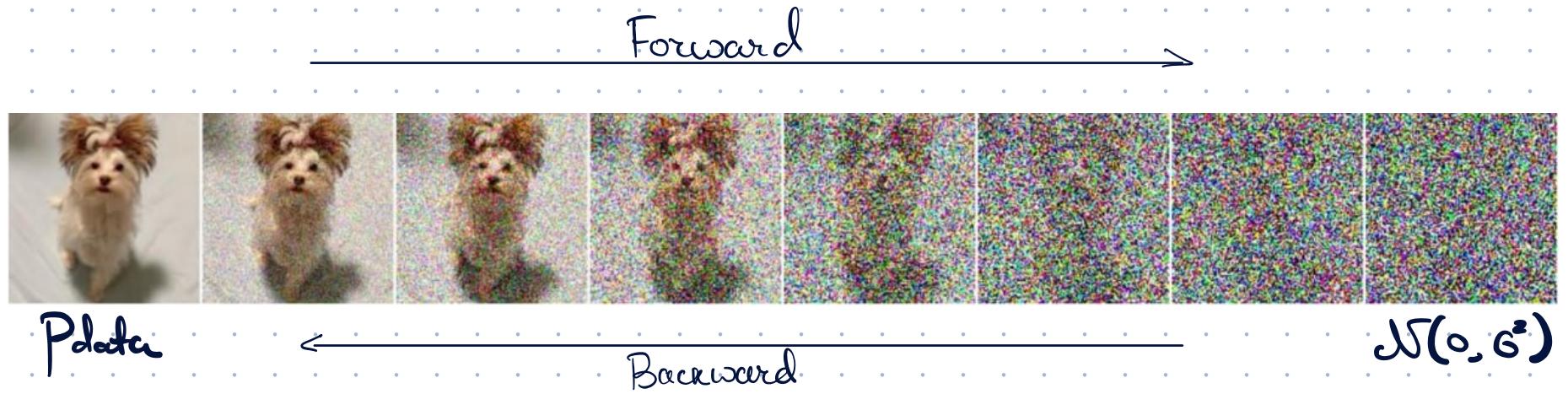
# Diffusion models

Denis Rakitin

# Generative trilemma

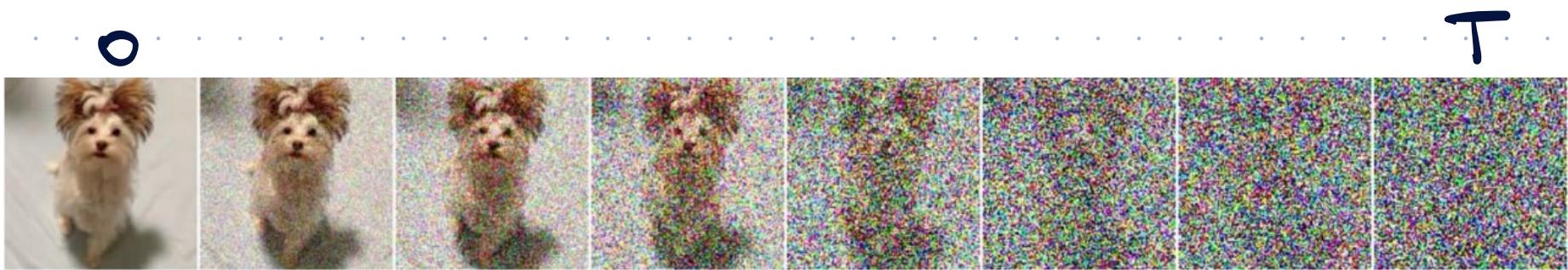


# Diffusion



Idea: define forward noising process and try to reverse it.

# Forward noising process



$x_0 \sim p_{\text{data}}$  ;  $x_t = x_{t-1} + g(t)\varepsilon_t$ ,  $\varepsilon_t \sim N(0, I)$   
(Variance exploding), Reg. c  $x_0 \dots x_{t-1}$ .

$x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_t$ ,  $\varepsilon_t \sim N(0, I)$   
Reg. c  $x_0 \dots x_{t-1}$ .

(variance preserving, VP)

$$\text{Var } x_t = (1-\beta_t) \text{Var } x_{t-1} + \beta_t \cdot I$$

$$\text{Var } x_0 = I \rightarrow \text{Var } x_t = I \quad \forall t = 0 \dots T.$$

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

$$q_{t|0}(x_t|x_0) = \mathcal{N}(x_t |$$

$$x_t = c_t x_{t-1} + d_t \varepsilon_t = c_0 x_0 + \sum_{s=1}^t c_s \cdot \varepsilon_s$$

$$\mathbb{E}_{q_{t|0}(x_t|x_0)} x_t = \mathbb{E}(\sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_t) = \sqrt{1-\beta_t} \mathbb{E}_{q_{t-1|0}(x_{t-1}|x_0)} x_{t-1}$$

$$= \left( \prod_{s=1}^t \sqrt{1-\beta_s} \right) \cdot \mathbb{E}_{q_{0|0}(x_0|x_0)} x_0 = x_0 \cdot \prod_{s=1}^t \sqrt{1-\beta_s}$$

$$\underset{q_{t|0}(x_t|x_0)}{\text{Var}} \underset{q_{t|0}(x_t|x_0)}{X_t} = (1-\beta_t) \underset{q_{t-1|0}(x_{t-1}|x_0)}{\text{Var}} \underset{q_{t-1|0}(x_{t-1}|x_0)}{X_{t-1}} + \beta_t \cdot I$$

$$\begin{aligned} I - \underset{q_{t|0}(x_t|x_0)}{\text{Var}} \underset{q_{t|0}(x_t|x_0)}{X_t} &= I - \beta_t I - (1-\beta_t) \underset{q_{t-1|0}(x_{t-1}|x_0)}{\text{Var}} \underset{q_{t-1|0}(x_{t-1}|x_0)}{X_{t-1}} \\ &= (1-\beta_t) \left( I - \underset{q_{t-1|0}(x_{t-1}|x_0)}{\text{Var}} \underset{q_{t-1|0}(x_{t-1}|x_0)}{X_{t-1}} \right) \\ &= \dots = \left( \prod_{s=1}^t (1-\beta_s) \right) \cdot \left( I - \underset{q_{0|0}}{\text{Var}} \underset{q_{0|0}}{X_0} \right) = \left( \prod_{s=1}^t (1-\beta_s) \right) \cdot I \end{aligned}$$

$$\underset{q_{t|0}(x_t|x_0)}{\text{Var}} \underset{q_{t|0}(x_t|x_0)}{X_t} = I - \left( \prod_{s=1}^t (1-\beta_s) \right) I = \left( 1 - \prod_{s=1}^t (1-\beta_s) \right) \cdot I$$

$$q_{\text{HIO}}(x_t | x_0) = \left( \prod_{s=1}^t \sqrt{1-\beta_s} \right) \cdot x_0 \quad \text{Var } q_{\text{HIO}}(x_t | x_0) = \left( 1 - \prod_{s=1}^t (1-\beta_s) \right) I$$

" "

$$\sqrt{\alpha_t} \cdot x_0 \quad (1-\alpha_t) \cdot I$$

$$\alpha_t = \prod_{s=1}^t (1-\beta_s)$$

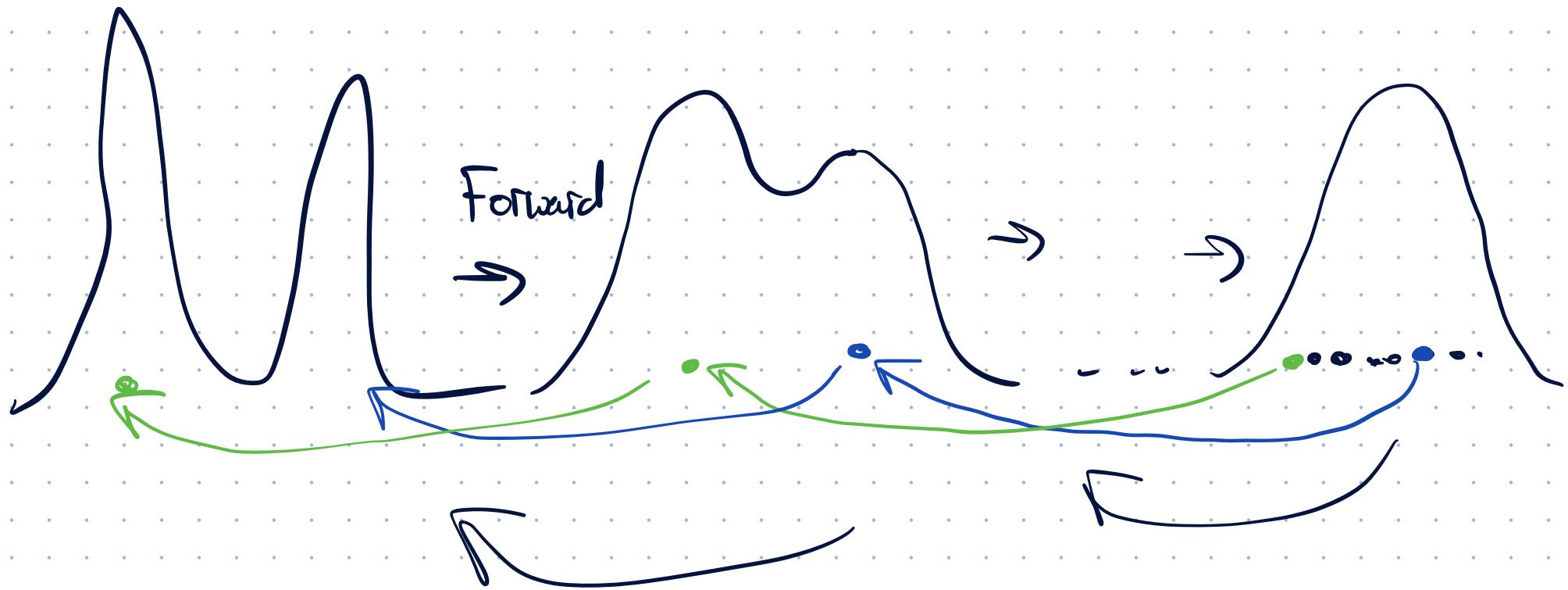
$$q_{\text{HIO}}(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\alpha_t} x_0, (1-\alpha_t) I)$$

$$q_{\text{TIO}}(x_T | x_0) \approx \mathcal{N}(x_T | 0, I)$$

$$(\text{depuis } T, \beta_t : \prod_{s=1}^T (1-\beta_s) \approx 0)$$

$$q_T(x_T) = \underbrace{\int q_{\text{TIO}}(x_T | x_0) q_0(x_0) dx_0}_{q_{0,T}(x_0, x_T)} \approx$$

$$\approx \int \mathcal{N}(x_T | x_0, I) q_0(x) dx_0 = \mathcal{N}(x_T | x_0, I).$$



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

$$q_{t|0}(x_t|x_0) = \mathcal{N}(x_t | \sqrt{\alpha_t} x_0, (\ell - \alpha_t) \mathbb{I}).$$

$$q_{t-1|t}(x_t \sim \mathcal{N}(0, \mathbb{I}) \quad x_{t-1} \sim q_{t-1|t}(x_{t-1}|x_t))$$

*ночью*

$$q_{t-1|t}(x_{t-1}|x_t) = \frac{q_{t|t-1}(x_t|x_{t-1}) q_{t-1}(x_{t-1})}{q_t(x_t)}$$

$$q_t(x_t) = \int q_{t|0}(x_t|x_0) p_{\text{data}}(x_0) dx_0$$

*ночью*      *сновие*

$q_{t+1|t}(x_{t+1}|x_t)$  — сконструирована распределение

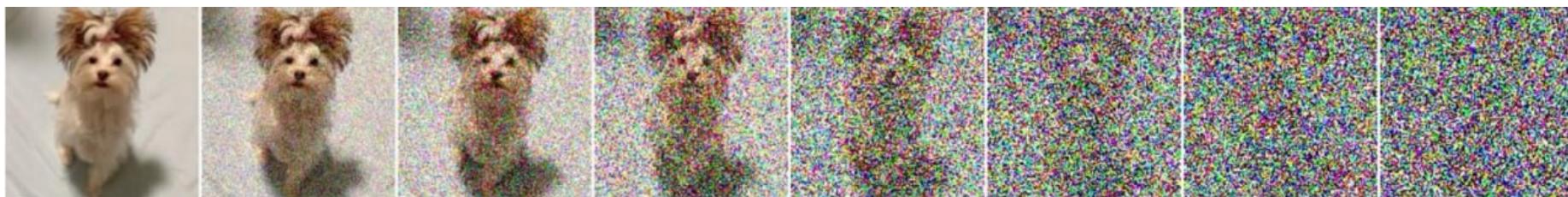
$$p_{t+1|t}^{\theta} (x_{t+1}|x_t) = \mathcal{N}(x_{t+1} | \mu_t^{\theta}(x_t), \sigma_t^2 I)$$

$\theta$  — параметры модели  
 $x_t$  — предыдущий вектор  
 $\mu_t^{\theta}(x_t)$  — предсказание модели  
 $\sigma_t^2$  — дисперсия предсказания

$$x_T \sim \mathcal{N}(0, I) \quad x_{t+1} = \mu_t^{\theta}(x_t) + \varepsilon_t \cdot \sigma_t$$

$$x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_t$$

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



$$x_{t+1} = \mu_t^{\theta}(x_t) + \sigma_t \hat{\varepsilon}_t \quad x_T \sim \mathcal{N}(0, I)$$

Forward:  $q_{0..T}(x_0..x_T) = q_0(x_0) \prod_{t=1}^T q_{t|t-1}(x_t|x_{t-1})$

\ Mapk.yemb

$p_{data}(x_0)$

Backward  $P_{0..T}^\theta(x_0..x_T) = \mathcal{N}(x_T | 0, I) \prod_{t=1}^T P_{t|t}^\theta(x_{t-1}|x_t)$

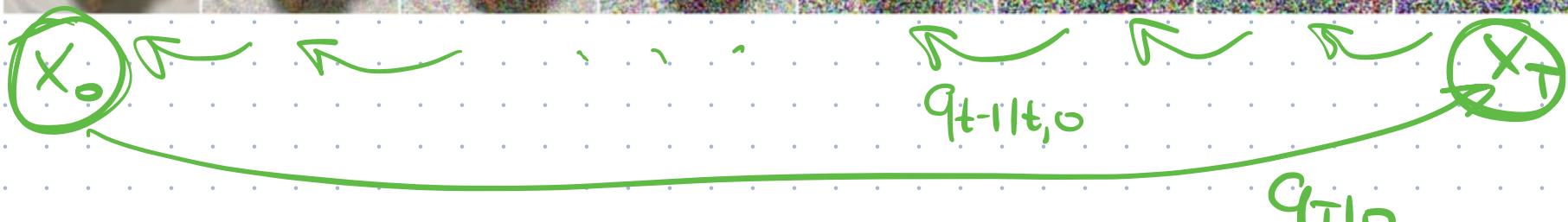
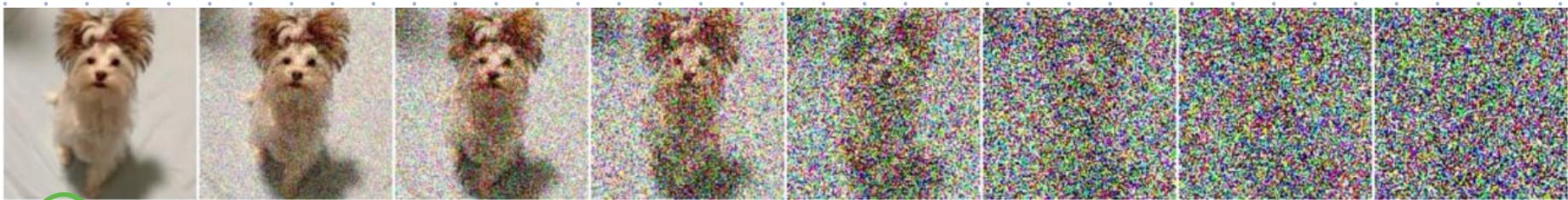
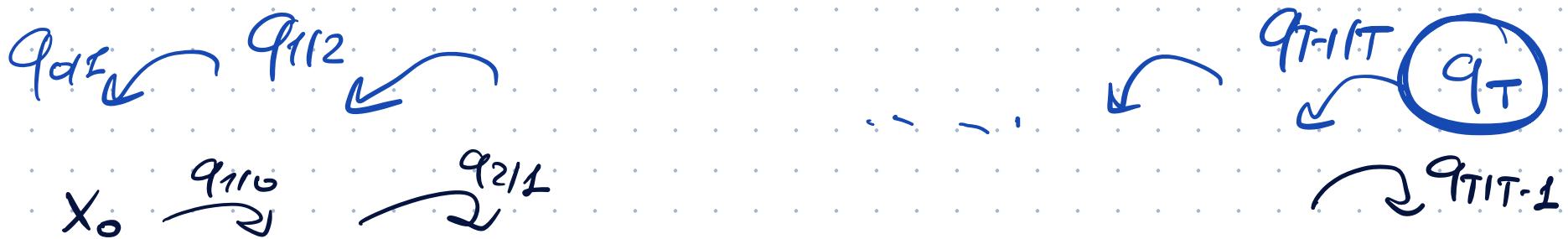
\ Mapk.yemb

$KL(q_{0..T}(x_0..x_T) || P_{0..T}^\theta(x_0..x_T)) \rightarrow \min_\theta$

"

$$\mathbb{E}_{q_{0..T}(x_0..x_T)} \log \frac{q_{0..T}(x_0..x_T)}{P_{0..T}^\theta(x_0..x_T)}$$

# Model and Training



$$q_{t-1|t} = \frac{q_{t|t-1} q_{t-1}}{q_t}$$

$$q_{t+1|t,0} = \frac{\overbrace{q_{t|t-1}}^N \overbrace{q_{t+1|0}}^N}{q_{t|0} - N}$$

$$KL\left(q_{0:T}(x_0 \dots x_T) \parallel p_{0:T}^{\theta}(x_0 \dots x_T)\right) \rightarrow \min_{\theta} \sum_T$$

$$\mathbb{E}_{q_{0..T}(x_0..x_T)} \log \frac{q_{0..T}(x_0..x_T)}{p_{0..T}^\theta(x_0..x_T)} = \mathbb{E}_{q_{0..T}} \log \frac{q_T(x_T) \prod_{t=1}^T q_{t-1|t}(x_{t-1}|x_t)}{\mathcal{N}(x_T|0, I) \prod_{t=1}^T p_{t-1|t}^\theta(x_{t-1}|x_t)}$$

$$= \overline{E}_{q_{0:T}} \log \frac{q_0(x_0) q_{T|0}(x_T|x_0)}{\mathcal{N}(x_T|0, I) P_{0|1}^0(x_0|x_1)} \prod_{t=2}^T q_{t-1|t,0}(x_{t-1}|x_t, x_0)$$

$$P_{0|1}^{\Theta}(x_0|x_1) \approx S_{x_1}(x_0)$$

$$\approx \underset{q_{0..T}(x_0..x_T)}{\mathbb{E}} \log \frac{\prod_{t=2}^T q_{t-1|t,0}(x_{t-1}|x_t, x_0)}{\prod_{t=2}^T p_{t-1|t}^\theta(x_{t-1}|x_t)} =$$

$$= \sum_{t=2}^T \underset{q_{0..T}(x_0..x_T)}{\mathbb{E}} \log \frac{q_{t-1|t,0}(x_{t-1}|x_t, x_0)}{p_{t-1|t}^\theta(x_{t-1}|x_t)} =$$

$$= \sum_{t=2}^T \underset{q_{0,t-1,t}(x_0, x_{t-1}, x_t)}{\mathbb{E}} \log \frac{q_{t-1|t,0}(x_{t-1}|x_t, x_0)}{p_{t-1|t}^\theta(x_{t-1}|x_t)} =$$

$$= \sum_{t=2}^T \underset{q_{0,t}(x_0, x_t)}{\mathbb{E}} \underset{q_{t-1|0,t}(x_{t-1}|x_0, x_t)}{\mathbb{E}} \log \frac{q_{t-1|t,0}(x_{t-1}|x_t, x_0)}{p_{t-1|t}^\theta(x_{t-1}|x_t)} =$$

$$= \sum_{t=2}^T \underset{q_{0,t}(x_0, x_t)}{\mathbb{E}} \text{KL}\left(q_{t-1|t,0}(x_{t-1}|x_t, x_0) \parallel p_{t-1|t}^\theta(x_{t-1}|x_t)\right)$$

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

$$q_{t+1|t,0}(x_{t+1}|x_t, x_0) = \frac{q_{t|t-1}(x_t|x_{t-1}) q_{t+1|0}(x_{t+1}|x_0)}{q_{t|0}(x_t|x_0)}$$

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

$$q_{t|0} = \mathcal{N}(x_t | \sqrt{\alpha_t} x_0, (1-\alpha_t) I)$$

$$q_{t+1|0} = \mathcal{N}(x_{t+1} | \sqrt{\alpha_{t+1}} x_0, (1-\alpha_{t+1}) I)$$

$$\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \mathbf{G}^2 \mathbf{I}) = \frac{1}{(2\pi G^2)^{d/2}} \exp\left(-\frac{1}{2G^2} \|\mathbf{z}-\boldsymbol{\mu}\|^2\right)$$

$$= \text{const.} \cdot \exp(-\text{kb.} \varphi - \epsilon(\mathbf{z}, \boldsymbol{\mu}))$$

$$\begin{aligned}
 q_{t-1|t,0}(x_{t-1}|x_t, x_0) &= \\
 \frac{\exp(-\text{rb.}(x_{t-1}, x_t)) \exp(-\text{rb.}(x_{t-1}, x_0))}{\exp(-\text{rb.}(x_t, x_0))} \cdot \text{const} \\
 &= \text{const.} \cdot \exp(-\text{rb.}(x_{t-1}, x_t, x_0)) \\
 &= \mathcal{N}(x_{t-1} | A_t x_t + B_t x_0, C_t^2 I).
 \end{aligned}$$

$$= \sum_{t=2}^T \mathbb{E}_{q_{0,t}(x_0, x_t)} \text{KL} \left( q_{t-1|t,0}(x_{t-1}|x_t, x_0) \parallel p_{t-1|t}^\theta(x_{t-1}|x_t) \right)$$

$$p_{t-1|t}^\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1} | \mu_t^\theta(x_t), C_t^2 I)$$

$$q_{t-1|t,0}(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1} | A_t x_t + B_t x_0, C_t^2 I).$$

$$= \sum_{t=2}^T \mathbb{E}_{q_{0,t}(x_0, x_t)} w_t \| \mu_t^\theta(x_t) - (A_t x_t + B_t x_0) \|^2$$

$$\mu_t^\theta(x_t) = A_t x_t + B_t D_t^\theta(x_t)$$

$$= \sum_{t=2}^T \mathbb{E}_{q_{0,t}(x_0, x_t)} \frac{\hat{w}_t \| D_t^\theta(x_t) - x_0 \|^2}{p_{\text{data}}(x_0) q_{t|0}(x_t|x_0)} \rightarrow \min_{\theta}$$

Условия : batch-size  $B$

- $X_0^{(1)} \dots X_0^{(B)} \sim p_{\text{data}}$
- $t^{(1)} \dots t^{(B)} \sim \text{Uniform dist.} \text{ (Haus. } U\{2, \dots, T\})$
- $\varepsilon^{(1)} \dots \varepsilon^{(B)} \sim \text{reg. } \mathcal{N}(0, I)$
- $\alpha^{(1)} \dots \alpha^{(B)} : \alpha^{(i)} = \alpha_{t^{(i)}} = \prod_{s=1}^{t^{(i)}} (1 - \beta_s)$
- $X_{\text{noisy}}^{(i)} = \sqrt{\alpha^{(i)}} X_0^{(i)} + \sqrt{1 - \alpha^{(i)}} \varepsilon^{(i)}$ .
- $\frac{1}{B} \sum_{i=1}^B \| D_{t^{(i)}}^{\theta} (X_{\text{noisy}}^{(i)}) - X_0^{(i)} \|^2, \text{backward}()$

## Sampling

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

$$x_{t-1} \sim P_{t-1|t}^{\Theta}(x_{t-1}|x_t) = \mathcal{N}\left(x_{t-1} | \mu_t^{\Theta}(x_t), C_t^2 I\right)$$
$$= A_t x_t + B_t D_t^{\Theta}(x_t)$$

$$x_{t-1} = A_t x_t + B_t D_t^{\Theta}(x_t) + C_t \varepsilon_t$$

DDPM

$$\sum_{t=2}^T \omega_t \mathbb{E} \| D_t^\epsilon(x_t) - x_0 \|^2 =$$

$q_{0,t}(x_0, x_t)$

$$q_{0,t}(x_0, x_t) =$$

$$= q_0(x_0) \mathcal{N}(\epsilon|_0, I)$$

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon$$

$x_t =$

$$= \sum_{t=2}^T \mathbb{E}_{\substack{\omega_t \\ q_0(x_0) \mathcal{N}(\varepsilon|_0, I)}} \| D_t^\circ(x_t) - x_0 \|^2 =$$

$$= \sum_t \mathbb{E}_{\substack{\omega_t \\ q_0(x_0) \mathcal{N}(\varepsilon|_0, I)}} \| D_t^\circ(x_t) - \frac{x_t - \sqrt{1-\alpha_t} \varepsilon}{\sqrt{\alpha_t}} \|^2 =$$

$$D_t^\circ(x_t) = \frac{x_t - \sqrt{1-\alpha_t} \varepsilon_t^\circ(x_t)}{\sqrt{\alpha_t}}$$

$$\textcircled{=} \sum_t \mathbb{E}_{\substack{\omega_t \\ q_0(x_0) N(\varepsilon | 0, I)}} \left\| \varepsilon_t^o (\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon) - \varepsilon \right\|^2$$

$$L_{\text{simple}} = \sum_{t=1}^T \mathbb{E}_{\substack{\omega_t \\ q_0(x_0) N(\varepsilon | 0, I)}} \left\| \varepsilon_t^o (\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon) - \varepsilon \right\|^2.$$

## Forward transitions

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

Show  $q_{t|0}(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\alpha_t} x_0, (1-\alpha_t)I)$ ,

$$\alpha_t = \prod_{s=1}^t (1-\beta_s).$$

## $\epsilon$ -prediction

Let  $\sum_{t=1}^T \omega_t \mathbb{E}_{q_{\theta, t}^{(x_0, x_t)}} \| D_t^\theta(x_t) - x_0 \|^2 \rightarrow \min_{\theta}$

be the training functional of DDPM.

Show that it is equivalent to

$$\sum_{t=1}^T \hat{\omega}_t \mathbb{E}_{q_\theta(x_0) \mathcal{N}(\varepsilon | 0, I)} \| \varepsilon_t^\theta (\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon) - \varepsilon \|^2 \rightarrow \min_{\theta}$$

## KL between Gaussians

Let  $x \sim \underbrace{\mathcal{N}(\mu_1, \sigma_1^2)}_{P_1}; y \sim \underbrace{\mathcal{N}(\mu_2, \sigma_2^2)}_{P_2}; x, y \in \mathbb{R}$ .

Show that  $\text{KL}(P_1 \parallel P_2) = \frac{1}{2} \left( \log \frac{\sigma_2^2}{\sigma_1^2} + \frac{\sigma_1^2}{\sigma_2^2} + \frac{(\mu_1 - \mu_2)^2}{\sigma_2^2} - 1 \right)$ .

In particular, if  $\hat{x} \sim \mathcal{N}(\mu_1, \sigma_1^2 I)$ ,  $\hat{y} \sim \mathcal{N}(\mu_2, \sigma_2^2 I)$ ,

then  $\text{KL}(P_1 \parallel P_2) = C \cdot \|\mu_1 - \mu_2\|^2 + f(\sigma_1, \sigma_2)$ .