

Diffusion Model

To generate new fake images being the same kind as the input images.

Forward: ① Given a ground-truth image to be learned to output

a mean to ensure feature
stability between generated &
ground-truth

each step's output provides a
snapshot for the reverse process
to learn

② Given a designed form to add steps of noises
while each noise is randomly
sampling from a decide distribution

a mean to introduce
stochasticity to later
model of the reverse process

$$X_t = \mu_t + \sigma_t \cdot \epsilon,$$

$$\mu_t = \sqrt{1 - \beta_t} \cdot X_{t-1}$$

$$\epsilon \sim N(0, I)$$

A designed form
A randomly sampled noise
(from an info decided distribution)

③ Get pure-noises image

meaning the entire map is of the decided distribution
(structured into, is completely lost, has no statistical difference
compared to any other pure-noise map of the same
distribution)

a standardised form
of the starting state of
the reverse process

Reverse : ① Start from the pure-noises image

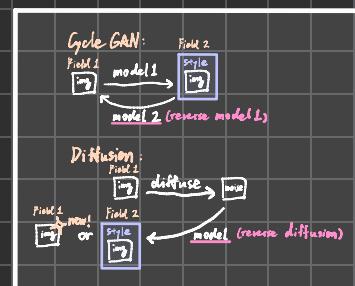
② In each step :

a. use X_t (to represent) & $\bar{E}^{(t)}$ (the ground-truth) known from 'Forward'
to learn $E_\theta(X_t, t)$ and get the distribution of X_{t-1}

b. sample a \tilde{X}_{t-1}

a mean to introduce
stochasticity to later
model of the reverse process

③ Reach the
distribution of X_0 in
the end then sample
any new image



Forward: embed noises (Diffusion)

Designed → ① Markov Chain of ② $q(x_t | x_{t-1}) = N(x_t; \sqrt{1 - \beta_t} \cdot x_{t-1}, \beta_t I)$

where ③ $\beta_t = S(t)$ (a Variance Schedule < diffusion rate > than $t \uparrow \Rightarrow \beta_t \uparrow$)

$$\Rightarrow q(x_1, \dots, x_T) = \frac{q(x_0)}{q(x_0)} = \frac{q(x_0) \cdot q(x_1 | x_0)}{q(x_0)} = \frac{q(x_0) \cdot q(x_1 | x_0) \cdot q(x_2 | x_0, x_1)}{q(x_0)} = \dots$$

Markov Chain Assumption

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

$$= \prod_{t=1}^T q(x_t | x_{t-1}) = \prod_{t=1}^T N(x_t; \sqrt{1 - \beta_t} \cdot x_{t-1}, \beta_t I) \Rightarrow \text{the Forward Proc. (Diffusion)}$$

\star Use the Multi-variance Method instead of the original N to be differentiable.

$$x_t = \underbrace{\int [1 - \beta_t] \cdot x_{t-1} + \int \beta_t \cdot E^{(t)}}_{\text{Initial noise}} \quad \text{Sample once from } N(0, I)$$

$$= \int [1 - \beta_t] \cdot (1 - \beta_{t-1}) \cdot x_{t-2} + \int \beta_{t-1} \cdot E^{(t-1)} + \underbrace{\int \beta_t \cdot E^{(t)}}_{\text{Embed noise from } N(0, I) \text{ times}} \quad \text{Sample twice from } N(0, I)$$

$$= \int [1 - \beta_t] \cdot (1 - \beta_{t-1}) \cdot x_{t-2} + \underbrace{\int [1 - \beta_t] \beta_{t-1} \cdot E^{(t-1)}}_{N(0, I) \text{ times}} + \int [1 - \alpha_t] \cdot x_{t-1} + \int [1 - \alpha_t] \beta_{t-1} \cdot E^{(t-1)} + \underbrace{\int \beta_t \cdot E^{(t)}}_{\text{Embed noise from } N(0, I) \text{ times}}$$

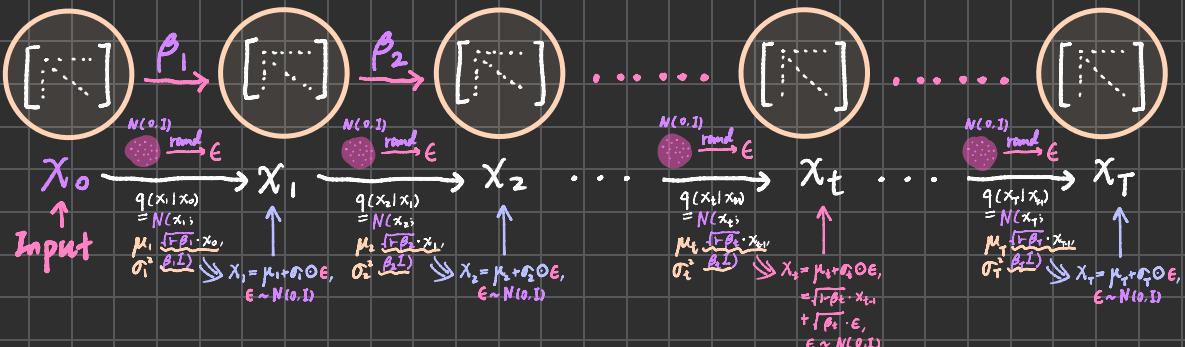
$$= \dots = \sqrt{\prod_{i=1}^t \alpha_i} \cdot x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \cdot E^{(t)}$$

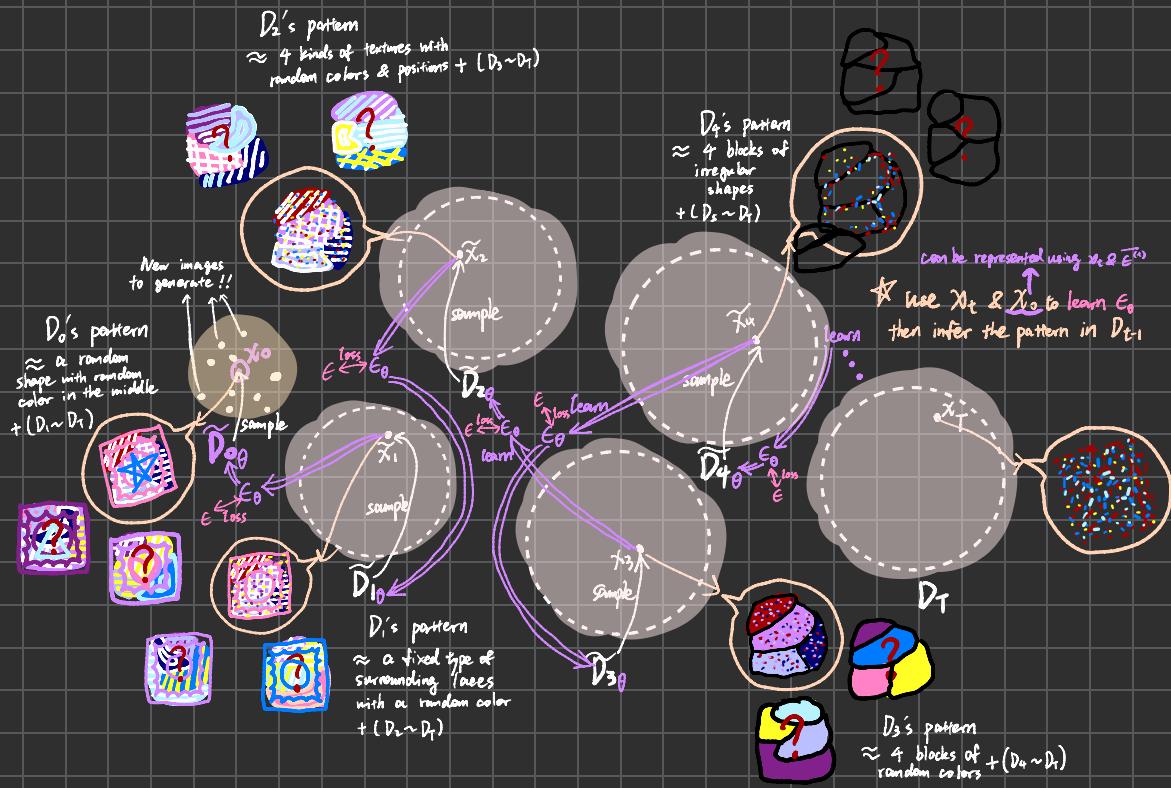
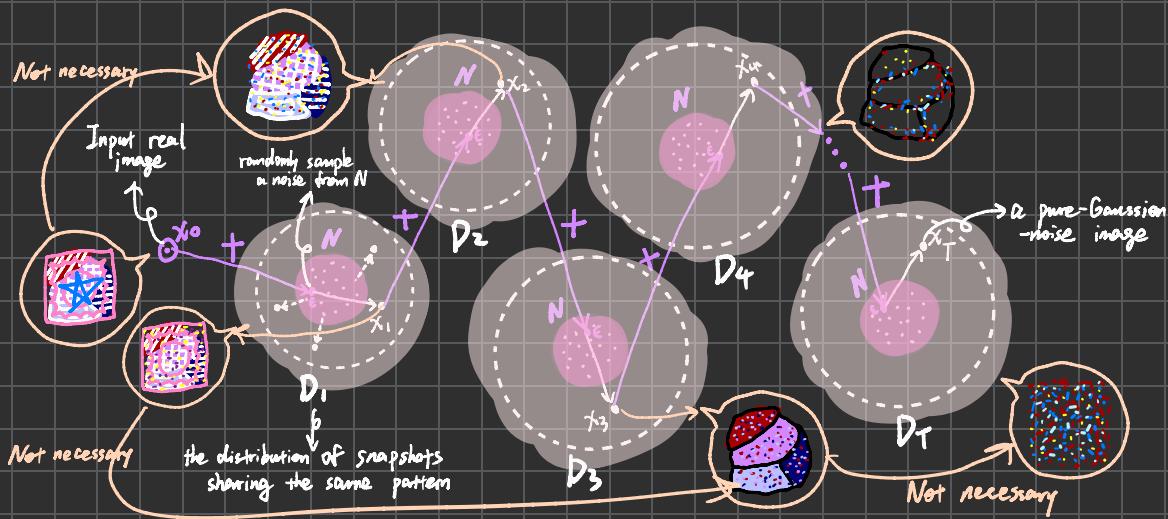
$$\overline{x}_t = \prod_{i=1}^t \alpha_i \cdot x_0 \quad \sqrt{\overline{x}_t} \xrightarrow{M} \quad \sqrt{1 - \overline{x}_t} \cdot E^{(t)} \xrightarrow{O}$$

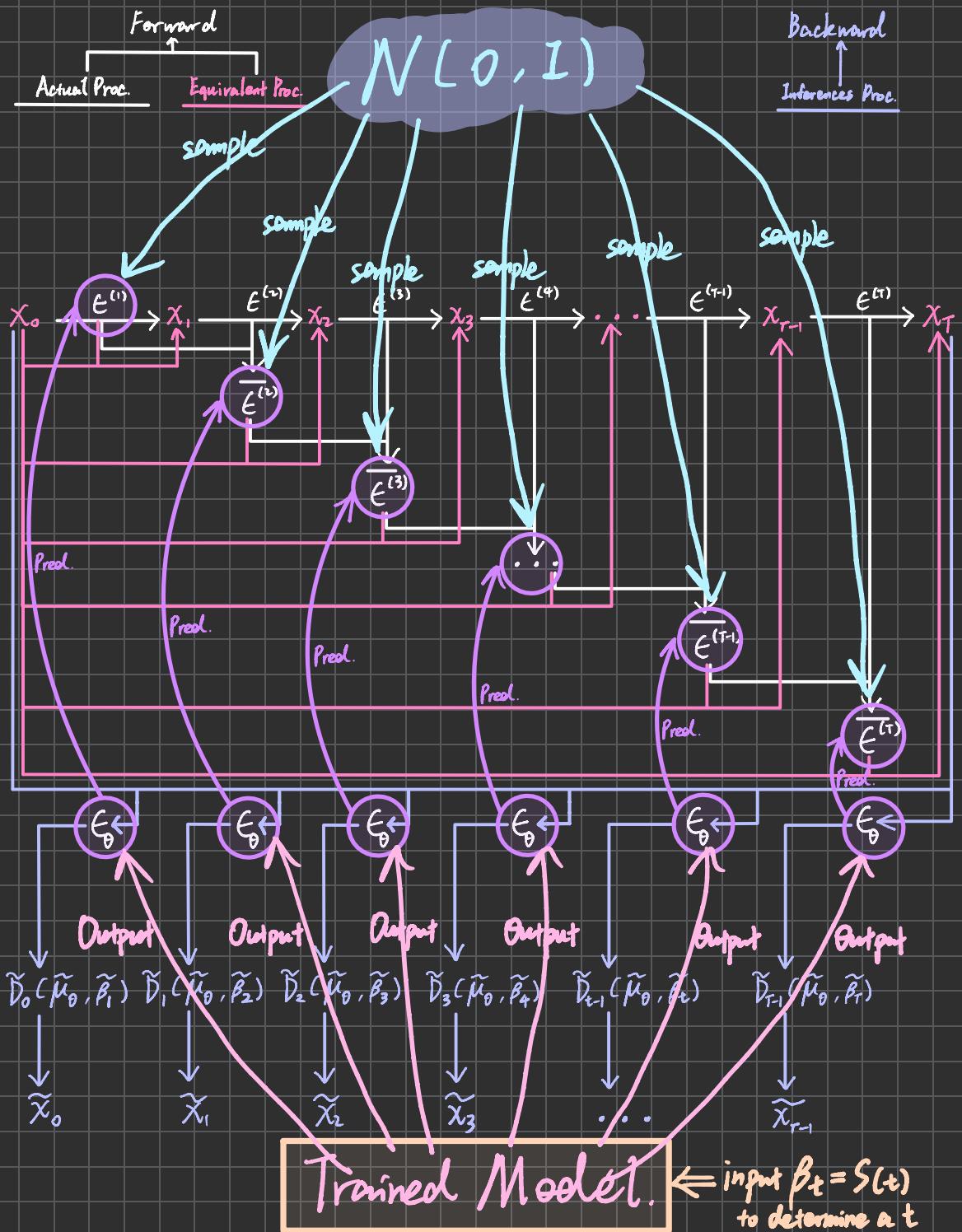
$$\Rightarrow q(x_t | x_0) = N(x_t; \sqrt{\overline{\alpha}_t} \cdot x_0, 1 - \overline{\alpha}_t) \xrightarrow{\text{For Reverse Process ("Snapshots")}}$$

Δ The snapshots are not for exactly anchoring the targets in Reversing !!
 They're only the clues for the reverse model to land in their own distributions unknown to the model! → ① Record a $\overline{E}^{(t)}$; ② Visual features of x_t .

$$\begin{aligned} N(0, I) &= \frac{1}{(1 - \beta_t) \beta_{t-1} + \beta_t^2} \cdot E^{(t)} + \beta_t \cdot E^{(t)} \\ N(0, I) &= N(0, (1 - \beta_t) \beta_{t-1} + \beta_t^2) \\ &= \frac{1}{(1 - \beta_t) \beta_{t-1} + \beta_t^2} \cdot N(0, I) \\ &\Downarrow \\ &= \frac{1}{(1 - \beta_t) \beta_{t-1} + \beta_t^2} \cdot E^{(t)} \\ &N(0, I) \end{aligned}$$







Backward: denoise (inferences)

Designed → Markov Chain

Assumed → $q(x_{t-1}|x_t) = N(x_{t-1}; \tilde{\mu}_t^?, \tilde{\beta}_t^? I)$ because β_t is small enough.

⇒ Need to be learned!! Known from forward proc.: $x_0, \tilde{e}_0, x_T, p_T \Rightarrow q(x_t|x_{t+1}), q(x_t|x_0), x_t$

$$\Rightarrow q(x_{t-1}|x_t, x_0) = \frac{q(x_{t-1}, x_t, x_0)}{q(x_t, x_0)} = \frac{q(x_{t-1}|x_t, x_0)}{q(x_t, x_0)} = \frac{q(x_{t-1}|x_{t-1}, x_0)}{q(x_t, x_0)} = \frac{q(x_{t-1}|x_{t-1})}{q(x_t, x_0)} \cdot \frac{q(x_{t-1}, x_0)}{q(x_t, x_0)} \rightarrow N(x_{t-1}; (I - \tilde{\beta}_{t-1}^2) \cdot x_0, \tilde{\beta}_{t-1}^2 I) \rightarrow N(x_{t-1}; (\tilde{\mu}_{t-1}^? + \tilde{\beta}_{t-1}^? \cdot x_0), \tilde{\beta}_{t-1}^? I)$$

$$= \frac{1}{\sqrt{2\pi \cdot \beta_t}} \cdot e^{-\frac{(x_t - (\tilde{\mu}_t^? + \tilde{\beta}_t^? x_0))^2}{2\beta_t^2}} \cdot \frac{1}{\sqrt{2\pi \cdot (1-\beta_t)}} \cdot e^{-\frac{(x_{t-1} - (\tilde{\mu}_{t-1}^? + \tilde{\beta}_{t-1}^? x_0))^2}{2(1-\beta_t)^2}} = \frac{1}{\sqrt{1+\beta_t}} \cdot e^{-\frac{1}{2} \left[\frac{(x_t - (\tilde{\mu}_t^? + \tilde{\beta}_t^? x_0))^2}{\beta_t^2} + \frac{(x_{t-1} - (\tilde{\mu}_{t-1}^? + \tilde{\beta}_{t-1}^? x_0))^2}{(1-\beta_t)^2} \right]}$$

$$\frac{\tilde{e}_{t-1}}{\sqrt{1+\beta_t}} = \tilde{e}_0 \quad \tilde{e}_0 \cdot e^{-\frac{1}{2} \left[\left(\frac{1}{1-\beta_t} + \frac{\alpha_t}{\beta_t} \right) \cdot x_0^2 + x_{t-1}^2 - 2 \left(\frac{x_t / \beta_t}{1-\beta_t} + \frac{x_{t-1} / \beta_t}{1-\beta_t} \right) \cdot x_0 \cdot x_{t-1} + C(x_{t-1}, x_0) \right]} \rightarrow \text{Can be represented with } \tilde{e}_0, x_0, x_{t-1} \text{ but not enough!} \rightarrow \text{Model } q \text{ with } p_0 \text{ for learning}$$

$$\Leftrightarrow \begin{cases} \frac{1}{1-\beta_t} + \frac{\alpha_t}{\beta_t} = \frac{1}{\delta^2} \\ \frac{x_t / \beta_t}{1-\beta_t} + \frac{x_{t-1} / \beta_t}{1-\beta_t} = \frac{\tilde{\mu}_t^?}{\delta^2} \end{cases} \Leftrightarrow \begin{cases} \tilde{\sigma}^2 = \tilde{\beta}_t = \frac{\beta_t (1-\beta_t)}{\beta_t + \alpha_t (1-\beta_t)} = \frac{\beta_t (1-\beta_t)}{1-\beta_t} \\ \tilde{\mu}_t^? = \frac{\beta_t x_0 \sqrt{\delta^2} + (1-\beta_t) x_{t-1} \sqrt{\delta^2}}{\beta_t + \alpha_t (1-\beta_t)} = \frac{\beta_t x_0 + (1-\beta_t) x_{t-1} + \sqrt{1-\beta_t} \cdot \epsilon_t}{1-\beta_t} \end{cases}; \tilde{\mu}_t^? \rightarrow \text{Model } q \text{ with } p_0 \text{ for learning}$$

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

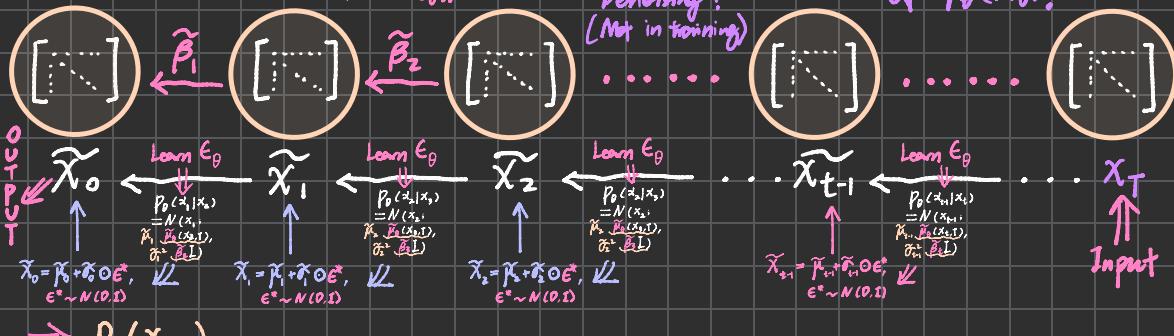
& ∵ $(\tilde{e}_0) = \frac{x_t - \sqrt{\tilde{\beta}_t} \cdot x_0}{\sqrt{1-\tilde{\beta}_t}}$ is sampled & x_0 is given
 \therefore To be learnable, use $E_\theta(x_t|t)$ as an equivalence to the single random $\tilde{e}(t)$

$$\Rightarrow P_\theta(x_{t-1}|x_t) := N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta I) \xrightarrow{\text{learn}} q(x_{t-1}|x_t) \xrightarrow{\text{Markov Chain Assumption}} q(x_{t-1}|x_t, x_0)$$

★ Use the Multi-varian Method instead of the original N to be differentiable:

And $\tilde{x}_{t-1} = \mu_\theta(x_t, t) + \sqrt{\Sigma_\theta} \cdot \epsilon^*$
 Learned! Contain E_θ !! Still sampled from $N(0, I)$!

Ready for the Variational lower bound of $P_\theta(x_0)$!



$$= p(x_0) \cdot p(x_1|x_0) \cdot p(x_2|x_1, x_0) \cdots p(x_{t-1}|x_{t-2}, x_0) \cdots p(x_T|x_{t-1}, x_0) = p(x_0) \cdot \prod_{t=1}^T p_\theta(x_t|x_{t-1}, x_0) = p(x_0) \prod_{t=1}^T N(x_t; \tilde{\mu}_\theta(x_t, t), \tilde{\beta}_\theta(t) I) \Rightarrow \text{the Reverse Proc. (Denoise)}$$

$N(x_t; 0, I)$
 (Pure noise)

★ Assume that:
 $p(x_{t-1}|x_t, x_{t-1}, \dots, x_T) = p(x_{t-1}|x_t)$

★ diffusion limit of small step size δ the reversal of the forward process (Feller, 1949). Since $x^{(k)}(t+\delta t)$ is a Gaussian (binomial) distribution, and δt is small, then $x^{(k)}(t+\delta t)$ is close to a normal distribution. The longer the trajectory the smaller the diffusion rate δ can be made.

↳ <https://arxiv.org/pdf/2202.09770.pdf>

★ The Multi-varian method
 from forward proc.

$$x_t = \frac{x_0 + \sqrt{1-\beta_t} \cdot \epsilon_t}{\sqrt{\beta_t}}$$

$$\sum \theta := \frac{x_t - \frac{x_0 + \sqrt{1-\beta_t} \cdot \epsilon_t}{\sqrt{\beta_t}}}{\sqrt{\beta_t}} \rightarrow \text{Loss}$$

(as actually cap value
 $\leq \tilde{\beta}_t$) \rightarrow learned

$$\mu_\theta := x_t - \frac{\tilde{\beta}_t \cdot \epsilon_t}{\sqrt{1-\tilde{\beta}_t}}$$

\downarrow
 Loss

Training.

Assumption: $q(x_t)$ is Gaussian

Learning: $P_\theta(x_t)$ learns the noise in q_t as a function

Ground-truth: D_t has $p(x_t) \approx 1$, which is the target of P_θ

Initial target: $E_\theta(x_t, t)$ equivalent to $\epsilon^{(t)} \sim N(t; 0, I)$

(△ No ways to calculate loss directly because equivalent is not exact "equal")

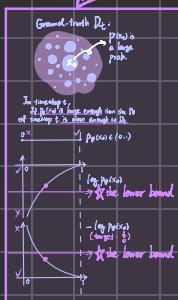
Practical target: $P_\theta(x_0) \rightarrow p(x_0) \rightarrow$ In actual Distribution of x_0

(△ Still difficult to optimize P_θ directly / Impossible to calculate prior prob.)

More Practical target:

$\min L$ where --

$$\begin{aligned} \mathbb{E}[-\log P_\theta(x_0)] &\leq L := L_{\text{nb}} = \mathbb{E}_q[-\log \frac{P_\theta(x_0, t)}{q(x_{0:T}|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_T) - \sum_{t=1}^T \log \frac{P_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})}] \end{aligned}$$



Acting like a Cross-entropy loss because
 $P_\theta(x_t)$ & $q(x_t|x_{t-1})$ are both probability distributions
 $\Rightarrow H(p(x_t)) = -\sum_x p(x_t) \log p(x_t)$
 $= -(\log p(x_t))$
 $= -\log p(x_t)$

means that using the simple Gaussian q to calculate this lower bound by means of which to optimize the black-boxed ground-truth prior p .

$$\begin{aligned} &= \mathbb{E}_q[-\log P_\theta(x_t) - \sum_{t=1}^T \log \frac{P_\theta(x_t|x_{t-1})}{q(x_t|x_{t-1})} - \log \frac{P_\theta(x_0|x_0)}{q(x_0|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_t) - (\sum_{t=1}^T \log \frac{P_\theta(x_t|x_{t-1})}{q(x_t|x_{t-1})} + \log \frac{q(x_t|x_{t-1})}{q(x_t|x_t)}) - \log \frac{P_\theta(x_0|x_0)}{q(x_0|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_t) - \sum_{t=1}^T \log \frac{P_\theta(x_t|x_{t-1})}{q(x_t|x_{t-1})} - \log \frac{\prod_{t=1}^T q(x_t|x_{t-1})}{\prod_{t=1}^T q(x_t|x_t)} \cdot \frac{P_\theta(x_0|x_0)}{q(x_0|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_t) + \sum_{t=1}^T \log \frac{q(x_t|x_{t-1})}{P_\theta(x_t|x_{t-1})} + \log \frac{q(x_t|x_t)}{P_\theta(x_t|x_t)}] \\ &= \mathbb{E}_q[\sum_{t=1}^T \log \frac{q(x_t|x_{t-1})}{P_\theta(x_t|x_{t-1})} + \log \frac{q(x_t|x_t)}{P_\theta(x_t|x_t)} - \log P_\theta(x_0|x_0)] \end{aligned}$$

variance reduction
using KL div.

$$L_T = \underbrace{\mathbb{E}_q[D_{KL}(q(x_t|x_0) || p_\theta(x_t))] + \sum_{t=1}^T \mathbb{E}_q[D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t)) - \log p_\theta(x_0|x_0)]}_{\text{constant}} \quad L_{t+1} = \underbrace{\mathbb{E}_q[D_{KL}(q(x_t|x_0) || p_\theta(x_t))]}_{\text{to be learned}} \quad L_0 = \underbrace{\text{monitored separately with a discrete derivative derived from } N(x_0; \mu_0(x_0), \Sigma_0(x_0))}_{\text{involving } \sum_{t=1}^T \mathbb{E}_q[D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))]}$$

$$\begin{aligned} L_t &= \frac{P_\theta(x_0) - \frac{1}{2} \left(\log \frac{\mu_0}{\mu_t} + \text{div} + \text{tr}(\Sigma_0^{-1} \Sigma_t) + (\mu_0 - \mu_t)^T \Sigma_0^{-1} (\mu_0 - \mu_t) \right)}{E_{x_0, \epsilon} \left[\frac{1}{2 \| \Sigma_0(x_t, t) \|_2^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \|^2 \right]} \\ \mu_t &= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\mu_0 - \mu_t}{\sqrt{1 - \alpha_t}} \right) \\ \mu_0 &= \frac{1}{\sqrt{\alpha_0}} \left(x_t - \frac{\mu_0 - \mu_t}{\sqrt{1 - \alpha_t}} \right) \\ \text{Multi-variate normal for } x_t &= \frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \alpha_t) \| \Sigma_0(x_t, t) \|_2^2} \| \epsilon - \epsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon, t) \|^2 \end{aligned}$$

2.4. Training
Training amounts to maximizing the model log likelihood,

$$\begin{aligned} &\int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p_\theta(\mathbf{x}^{(0)}) \quad (10) \\ &= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log \frac{\int d\mathbf{x}^{(1:T)} p_\theta(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})}{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} \quad (11) \end{aligned}$$

which has a lower bound provided by Jensen's inequality:

$$\log \frac{\int d\mathbf{x}^{(1:T)} p_\theta(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})}{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} \geq \log \frac{\int d\mathbf{x}^{(1:T)} p_\theta(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})}{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} + \log \frac{\int d\mathbf{x}^{(1:T)} p_\theta(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})}{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} \cdot \log \frac{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})}{\int d\mathbf{x}^{(1:T)} q(\mathbf{x}^{(1:T)} | \mathbf{x}^{(0)})} \quad (12)$$

As described in Appendix B, for our diffusion trajectory this reduces to,
 $L \geq K$

$N(0, I)$

