

DDPM

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 features
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 sampled from a decide distribution

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



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

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

$$= \sqrt{1-\beta_t} \cdot x_{t-1}$$

$$+ \sqrt{\beta_t} \cdot \epsilon,$$

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

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

Reverse:

① Start from the pure-noises image



② In each step :

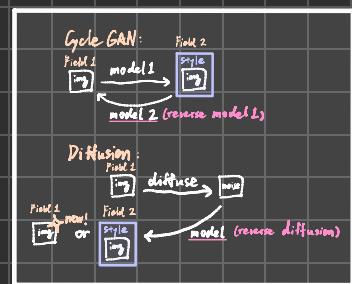
a. use x_t (to represent) & $\bar{x}^{(t)}$ (the ground-truth) known from 'Forward'
to learn $E_\theta(x_t, t)$ and get the distribution of x_{t-1}

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

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



③ Reach the
distribution of x_0 in
the end then sample
any new image



Forward: diffusion (inference)

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

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

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

Markov Chain Assumption

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

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

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

$$x_t = \int_{\beta_t}^{1-\beta_t} x_{t-1} + \int_{\beta_t}^{1-\beta_t} \epsilon^{(t)}$$

Sample once from $N(0, I)$

$$= \int_{1-\beta_t}^{\beta_t} (\int_{1-\beta_{t-1}}^{\beta_{t-1}} x_{t-2} + \int_{\beta_{t-1}}^{1-\beta_{t-1}} \epsilon^{(t-1)}) + \int_{\beta_t}^{1-\beta_t} \epsilon^{(t)} = \int_{1-\beta_t}^{\beta_t} (\int_{1-\beta_{t-1}}^{\beta_{t-1}} x_{t-2} + \int_{\beta_{t-1}}^{1-\beta_{t-1}} \beta_{t-1} \cdot \epsilon^{(t-1)} + \int_{\beta_t}^{1-\beta_t} \beta_t \cdot \epsilon^{(t)})$$

Sample twice from $N(0, I)$

$$= \int_{1-\beta_t}^{\beta_t} (\int_{1-\beta_{t-1}}^{\beta_{t-1}} x_{t-2} + \int_{\beta_{t-1}}^{1-\beta_{t-1}} \beta_{t-1} \cdot \epsilon^{(t-1)} + \int_{\beta_t}^{1-\beta_t} \beta_t \cdot \epsilon^{(t)})$$

N(0, I) times

$$= \dots = \sqrt{\pi_{t=1}^t} x_0 + \sqrt{1 - \pi_{t=1}^t} \cdot \epsilon^{(t)}$$

$\bar{x}_t = \pi_{t=1}^t x_0$

$\sqrt{\bar{x}_t} \xrightarrow{M}$

$\sqrt{1 - \bar{x}_t} \cdot \epsilon^{(t)} \xrightarrow{O}$

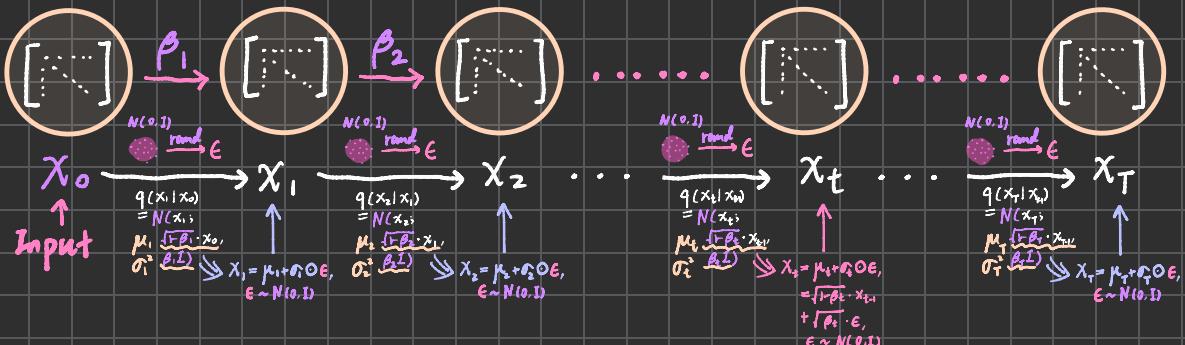
← Embedded noise from $N(0, I)$ t times

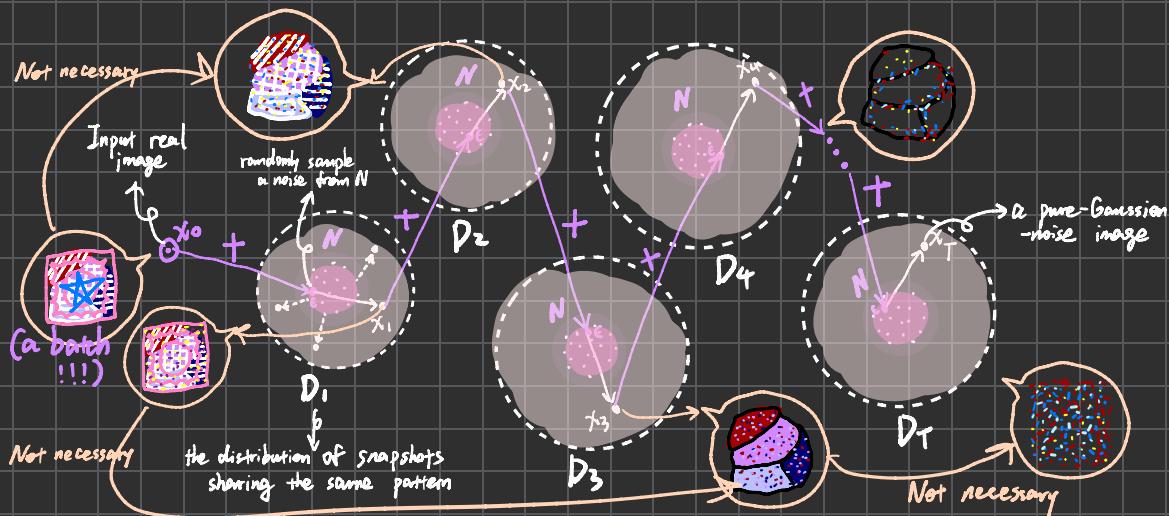
$$\Rightarrow q(x_t | x_0) = N(x_t; \sqrt{\bar{x}_t} \cdot x_0, 1 - \bar{x}_t) \rightarrow \text{For Reverse Process ("Snapshots")}$$

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!

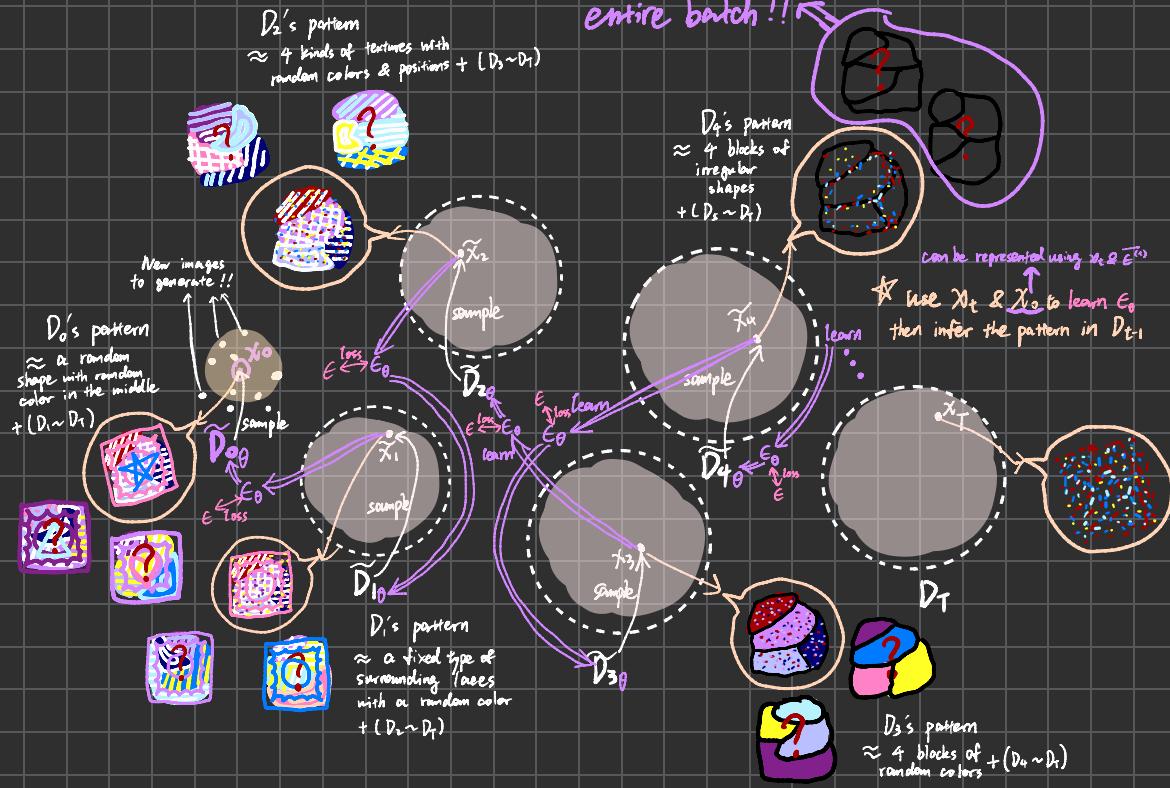
→ ① Record a \bar{x}_t ; ② Visual features of x_t .

$$\begin{aligned} & N(0, I) \\ & = \int_{(1-\beta_t)\beta_{t-1}}^{\beta_t} \epsilon^{(t)} + \int_{\beta_t}^{1-\beta_t} \beta_t \cdot \epsilon^{(t)} \\ & \vdots \\ & N(0, (1-\beta_{t-1})\beta_{t-2} + \dots + \beta_t \cdot \beta_1) \\ & = N(0, (1-\beta_{t-1})\beta_{t-2} + \dots + \beta_1) \\ & = \int_{(1-\beta_{t-1})\beta_{t-2} + \dots + \beta_1}^{1-\beta_t} \epsilon^{(t)} \\ & \vdots \\ & \int_{(1-\beta_1)\beta_0}^{1-\beta_t} \epsilon^{(t)} \\ & N(0, I) \end{aligned}$$





different patterns are learned from the entire batch !! ↗



Forward

$N(0, 1)$

?

various



$\times \text{batchsize}$



$\times \text{batchsize}$

various



$\times \text{batchsize}$

various



$\times \text{batchsize}$

various



$\times \text{batchsize}$

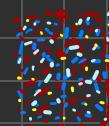
various



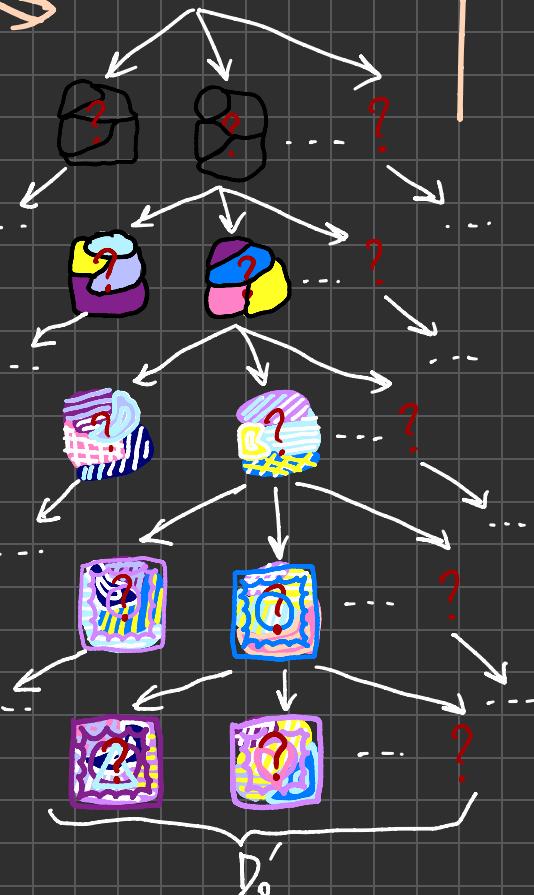
$\times \text{batchsize} \sim D_0$

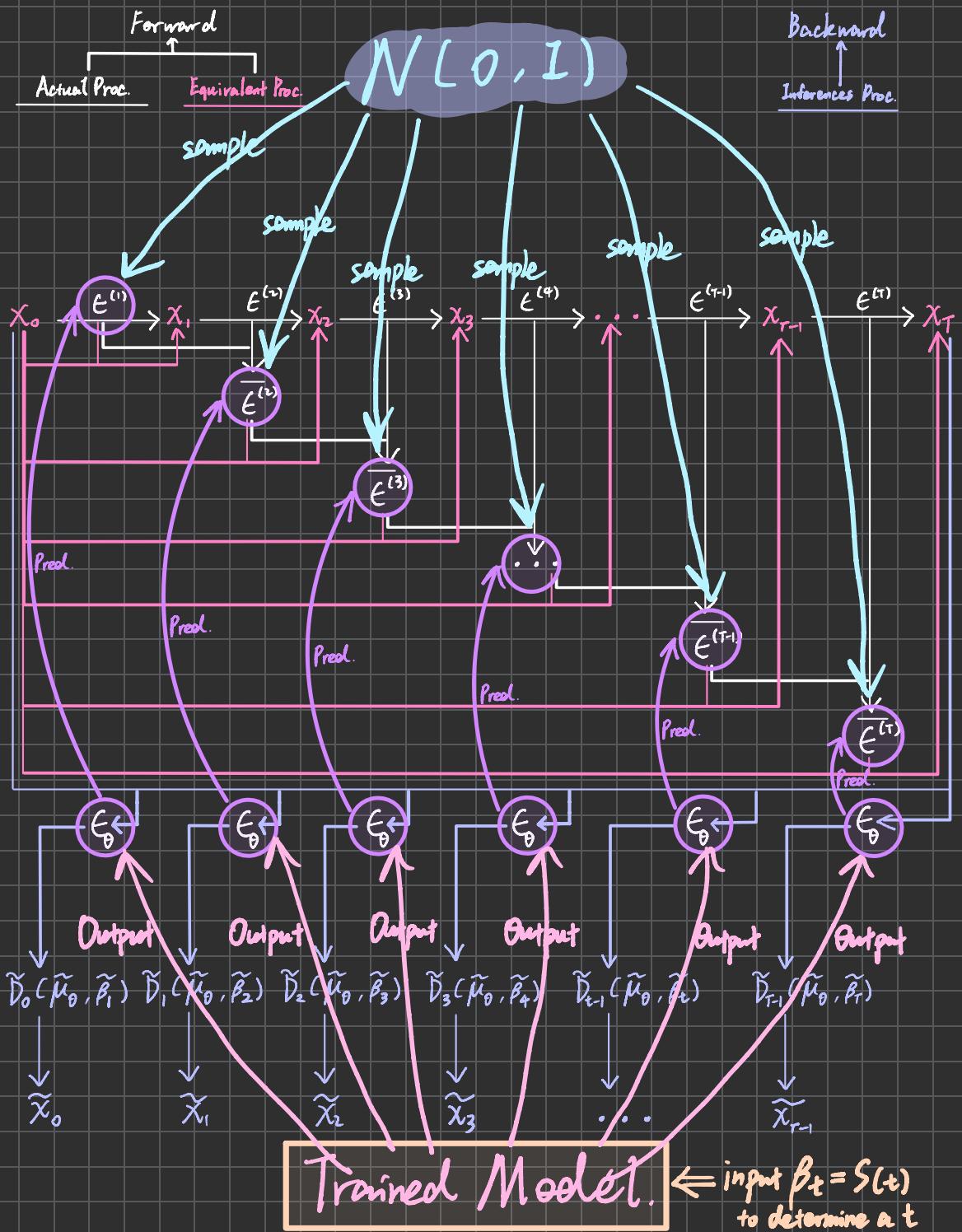
Backward

more t,
higher tree,
more
patterns!



larger batch,
more patterns!





Backward: denoise (generative)

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

$$\begin{aligned} \text{Assume that:} \\ p(x_{t-1}|x_t, x_{t-2}, \dots, x_T) \\ = p(x_{t-1}|x_t) \end{aligned}$$

diffusion time of small step size (τ) is the reversal of the forward process. It has the identical form as the forward process (Feller, 1949). Since $q(x^{(t)}|x^{(t-1)})$ is a Gaussian (binomial) distribution, and if σ_t^2 is small, then the variance of the noise is also small, i.e., the diffusion rate β_t can be made small.

⇒ Need to be learned!! Known from forward proc.: $x_0, \bar{x}_T, x_T, p_T \Rightarrow q(x_T|x_{t-1}), q(x_T|x_0), x_T$

target a destination to denoise eliminate the unknown $q(x)$ using Markov chain's assumption

$$\begin{aligned} \Rightarrow q(x_{t-1}|x_t, x_0) &= \frac{q(x_{t-1}|x_t, x_0)}{q(x_t|x_0)} = \frac{q(x_{t-1}|x_{t-1}, x_0)}{q(x_t|x_0)} = q(x_{t-1}|x_{t-1}) \cdot q(x_t|x_0) \\ &= \frac{1}{\sqrt{2\pi \cdot \beta_{t-1}}} \cdot e^{-\frac{(x_{t-1}-\tilde{\mu}_{t-1})^2}{2\beta_{t-1}}} \cdot \frac{1}{\sqrt{2\pi \cdot (\beta_t + \alpha_t)}} \cdot e^{-\frac{(x_t-\tilde{\mu}_t)^2}{2(\beta_t + \alpha_t)}} = \frac{1}{\sqrt{2\pi \cdot \beta_{t-1}}} \cdot e^{-\frac{1}{2} \left[\frac{(x_{t-1}-\tilde{\mu}_{t-1})^2}{\beta_{t-1}} - \frac{(x_t-\tilde{\mu}_t)^2}{\beta_t + \alpha_t} + C(x_{t-1}, x_t) \right]} \end{aligned}$$

$$\Leftrightarrow \begin{cases} \frac{1}{\beta_{t-1}} + \frac{\alpha_t}{\beta_t} = \frac{1}{\beta_t} \\ \frac{x_{t-1} - \tilde{\mu}_{t-1}}{\beta_{t-1}} + \frac{x_t - \tilde{\mu}_t}{\beta_t} = \frac{\tilde{\mu}_t}{\beta_t} \end{cases} \Leftrightarrow \begin{cases} \tilde{\sigma}^2 = \tilde{\beta}_t = \frac{\beta_t(1-\beta_{t-1})}{\beta_t + \alpha_t(1-\beta_{t-1})} = \beta_t(1-\tilde{\alpha}_{t-1}) \\ \tilde{\mu}_t = \frac{\beta_t x_{t-1} + (1-\beta_{t-1})x_t}{\beta_t + \alpha_t(1-\beta_{t-1})} \end{cases}$$

(use the $\tilde{\mu}_t$ as actual input) ; $\tilde{\mu}_t = \frac{\beta_t x_{t-1} + (1-\beta_{t-1})x_t}{\beta_t + \alpha_t(1-\beta_{t-1})}$; $\tilde{\mu}_t = \frac{\beta_t x_{t-1} + (1-\beta_{t-1})x_t}{\beta_t + \alpha_t(1-\beta_{t-1})} + \frac{x_t - \sqrt{\beta_t \cdot \alpha_t}}{\sqrt{\beta_t \cdot \alpha_t}}$

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

& $\because \tilde{E}(x) = \frac{x_t - \sqrt{\beta_t \cdot \alpha_t}}{\sqrt{\beta_t \cdot \alpha_t}}$ is sampled & x_0 is given

∴ To be learnable, use $E_\theta(x_t, t)$ as an equivalence to the single random $E(x)$

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

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

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

Only used when Denoising! (Not in training)

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



$$\begin{aligned} \text{OUT} \downarrow \quad \tilde{x}_0 &\leftarrow \tilde{\beta}_1 \quad \tilde{x}_1 &\leftarrow \tilde{\beta}_2 \quad \tilde{x}_2 &\leftarrow \cdots \quad \tilde{x}_{t-1} &\leftarrow \tilde{\beta}_t \quad x_T &\leftarrow \cdots \\ \tilde{x}_0 &= \tilde{\beta}_0 + \tilde{\beta}_0 \cdot E^*, \quad \tilde{x}_1 = \tilde{\beta}_1 + \tilde{\beta}_1 \cdot E^*, \quad \tilde{x}_2 = \tilde{\beta}_2 + \tilde{\beta}_2 \cdot E^*, \quad \vdots \quad \tilde{x}_{t-1} = \tilde{\beta}_{t-1} + \tilde{\beta}_{t-1} \cdot E^*, \quad x_T = \tilde{\beta}_t + \tilde{\beta}_t \cdot E^* \\ &E^* \sim N(0, I) \quad E^* \sim N(0, I) \quad E^* \sim N(0, I) \quad \dots \quad E^* \sim N(0, I) \quad E^* \sim N(0, I) \end{aligned}$$

$$\Rightarrow P_\theta(x_0|T) = P_\theta(x_1|T) \cdot P_\theta(x_2|x_1, T) \cdot P_\theta(x_3|x_2, T) \cdot \dots \cdot P_\theta(x_{t-1}|x_{t-2}, T) \cdot P_\theta(x_t|x_{t-1}, T) \cdot P_\theta(x_{t+1}|x_t, T) \cdot \dots \cdot P_\theta(x_T|x_{T-1}, T)$$

$$= P_\theta(x_1) \cdot P(x_{T-1}|x_T) \cdot P(x_{T-2}|x_{T-1}, x_T) \cdot \dots \cdot P(x_0|x_1, T)$$

$q(x_0)$

$$= P(x_1) \cdot \prod_{t=1}^T P_\theta(x_{t-1}|x_t) = P(x_1) \prod_{t=1}^T N(x_{t-1}; \tilde{\mu}_\theta(x_t, t), \tilde{\beta}_\theta I) \Rightarrow \text{the Reverse Proc (Denoise)}$$

$N(x_1; 0, I)$
(Pure noise)

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: $G_\theta(x_t, t)$ equivalent to $E^{(t)} \sim N(t; 0, I)$

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

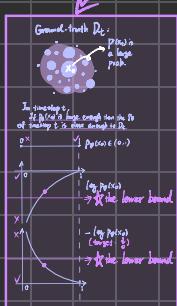
Practical target: $P_\theta(x_0) \rightarrow p(x_0) \rightarrow$ In actual Distribution of x_0

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

More Practical target:

$\min L$ where --

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



Actual vs Cross-entropy loss between P_θ & p prior values of x_0 :
 $\Rightarrow H(p, q) = -\sum_x p(x) \log q(x)$
 $= -(1 \cdot \log p(x_0)) + \dots + \log p(x_{n-1})$
 $= -\log p(x_0)$

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 .

$$\begin{aligned} &= \mathbb{E}_q[-\log P_\theta(x_T) - \sum_{t>0} \log \frac{P_\theta(x_t|x_0)}{q(x_t|x_0)} - \log \frac{P_\theta(x_0|t)}{q(x_0|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_T) - (\sum_{t>0} \log \frac{P_\theta(x_t|x_0)}{q(x_t|x_0)} + \log \frac{q(x_t|x_0)}{P_\theta(x_t|x_0)}) - \log \frac{P_\theta(x_0|t)}{q(x_0|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_T) - \sum_{t>0} \log \frac{P_\theta(x_t|x_0)}{p(x_t|x_0)} - \log \frac{q(x_t|x_0)}{P_\theta(x_t|x_0)}] \\ &= \mathbb{E}_q[-\log P_\theta(x_T) + \sum_{t>0} \log \frac{q(x_t|x_0)}{p(x_t|x_0)} + \log \frac{q(x_0|x_0)}{P_\theta(x_0|x_0)}] \\ &= \mathbb{E}_q[\sum_{t>0} \log \frac{q(x_t|x_0)}{p(x_t|x_0)} + \log \frac{q(x_0|x_0)}{P_\theta(x_0|x_0)} - \log P_\theta(x_0|x_0)] \end{aligned}$$

variance reduction
using KL div.

$$L_T = \mathbb{E}_q[D_{KL}(q(x_t|x_0) || p(x_t))] + \sum_{t>0} D_{KL}(q(x_{t+1}|x_t, x_0) || p(x_{t+1}|x_t)) - \log P_\theta(x_0|x_0)$$

L_T (constant)
 L_{t+1} (to be learnt)
 L_0 (modeled separately with a discrete decoder derived from $N(x_0; \mu_0(x_0), \Sigma_0(x_0))$)

$$\begin{aligned} L_T &= P_\theta(x_0) - \frac{1}{2} \left(\log \frac{\Sigma_0}{\alpha_t} - \dim + \ln \left(\frac{\Sigma_0}{\alpha_t} \right) + (\mu_0 - \mu_t)^T \Sigma_0^{-1} (\mu_0 - \mu_t) \right) \\ P_\theta &= \frac{1}{\sqrt{\alpha_t}} \left(x_0 - \frac{\mu_0 - \mu_t}{\sqrt{\alpha_t}} \right) \\ \mu_\theta &= \frac{1}{\sqrt{\alpha_t}} \left(x_0 - \frac{\mu_0 - \mu_t(x_0, t)}{\sqrt{1-\alpha_t}} \right) \\ \text{Multi-variate method for } x_t \\ (x_t = x_0 + (x_0, \epsilon)) \end{aligned}$$

$$E_{x_0, \epsilon} \left[\frac{1}{2 \|\Sigma_0(x_0, t)\|_F^2} \| \tilde{\mu}_t(x_0, \epsilon) - \mu_t(x_0, t) \|^2 \right]$$

$$E_{x_0, \epsilon} \left[\frac{1}{2 \|\Sigma_0(x_0, t)\|_F^2} \left(x_0 - \frac{\mu_0 - \mu_t}{\sqrt{1-\alpha_t}} \right) - \frac{1}{\sqrt{\alpha_t}} \left(x_0 - \frac{\mu_0 - \mu_t(x_0, t)}{\sqrt{1-\alpha_t}} \right) \right]^2$$

$$E_{x_0, \epsilon} \left[\frac{(1-\alpha_t)^2}{2 \alpha_t (1-\alpha_t)} \|\Sigma_0(x_0, t)\|_F^2 \right] \| E - E_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \epsilon, t) \|^2$$

(to speed up training)



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

$$\int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p_\theta(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log \left[\prod_{t=1}^T \frac{p_\theta(\mathbf{x}^{(t)} | \mathbf{x}^{(0:t-1)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(0:t-1)})} \right] \quad (11)$$

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

$$L \geq \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log [p_\theta(\mathbf{x}^{(0:t-1)}) + \sum_{t=1}^T \log \frac{p_\theta(\mathbf{x}^{(t)} | \mathbf{x}^{(0:t-1)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(0:t-1)})}]$$

As described in Appendix B, for our diffusion trajectory this reduces to

$$L > K \quad \text{min} \quad \text{(lower - higher)}$$

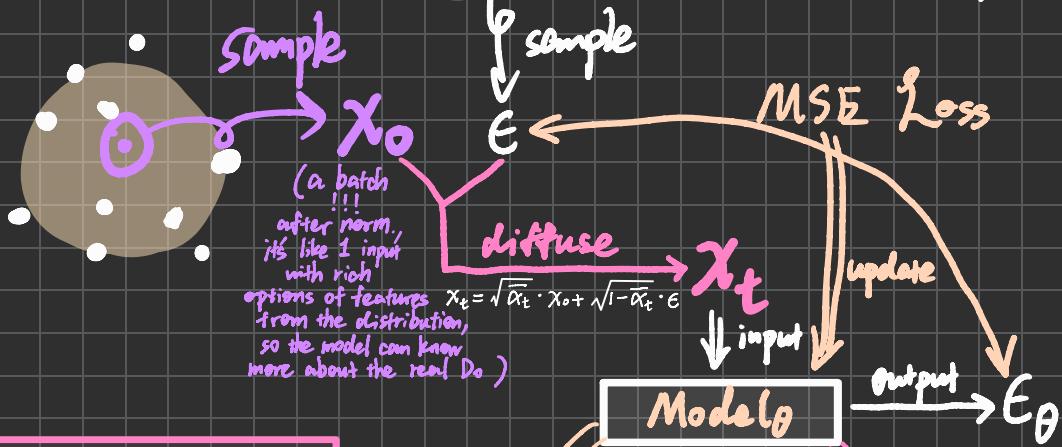
(13)

模型上扩散的目的是：

对原图进行从高层特征到低层特征的分解，从而可以在多层特征的分布下各自扩充特征类型（学习更多 ϵ ），最终叠加效果产生的 X_0 就更多样

训练分时间步的目的是：
通过分时间步独立学习减少误差率减小总误差

Ground-truth (D_0)



TRAINING

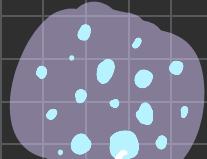
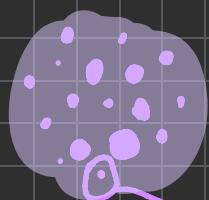
★ Ready for current
Denoising !!
(may be updated before,
or just the initial one)



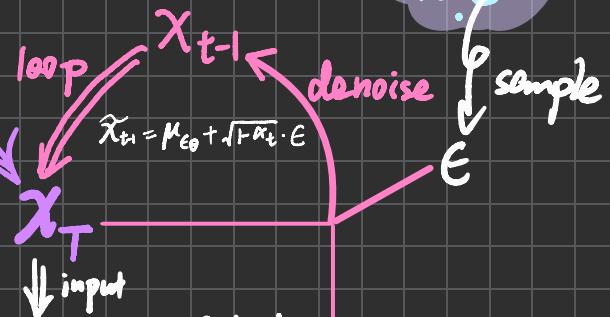
U-net!!
(latent = X_t 's latent implied by t)

$N(0, I)$

$N(0, I)$



SAMPLING



Model ϕ

Output $\rightarrow E_\theta$

↑ input

β_t

$\rightarrow U\text{-net}!!$

(latent = x_t 's latent implied by t)

★ Ready for any t
Denoising !!

from $t=T-1$ to $t=2$

$\uparrow \beta_t = S(t)$

Denoise Decoder

x_1

output
 \downarrow
 x_0

time steps