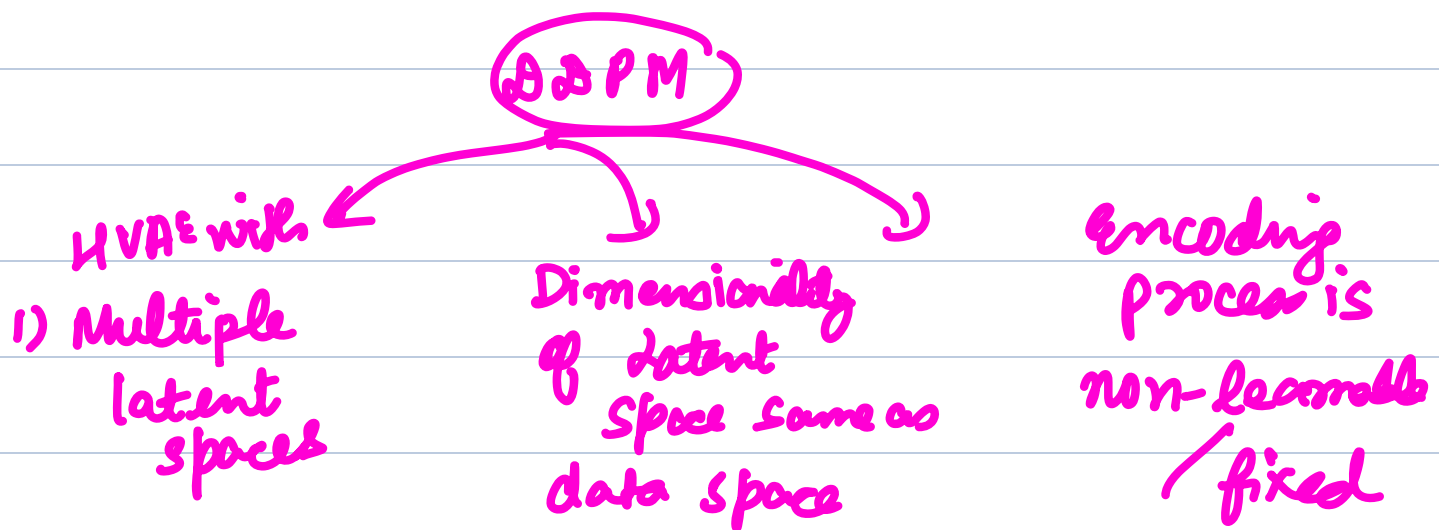
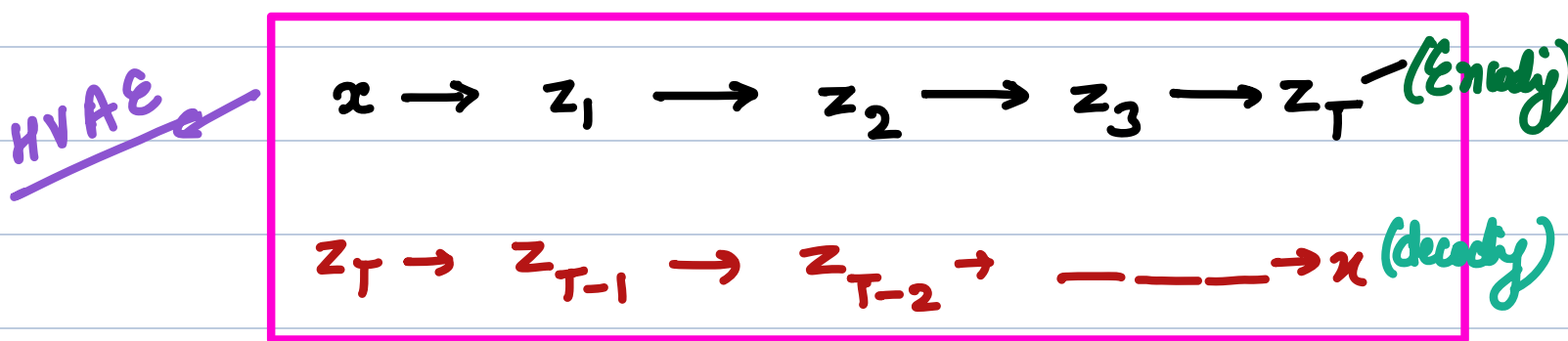
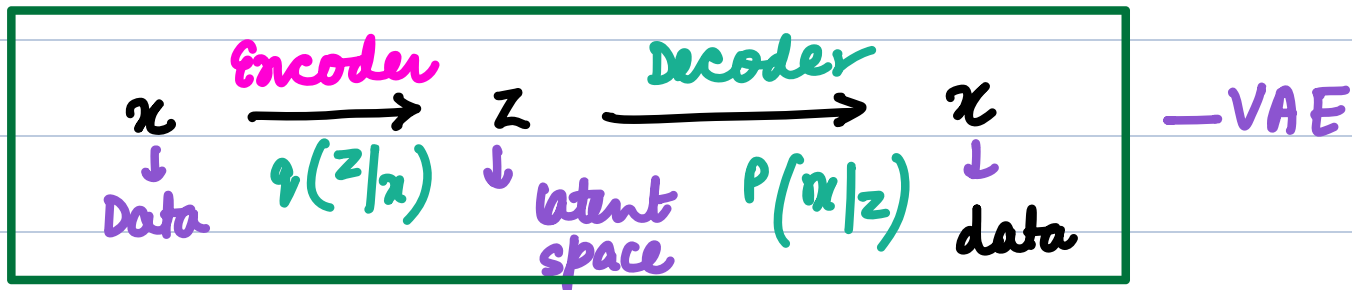


Denoising Diffusion Probabilistic models (DDPM)

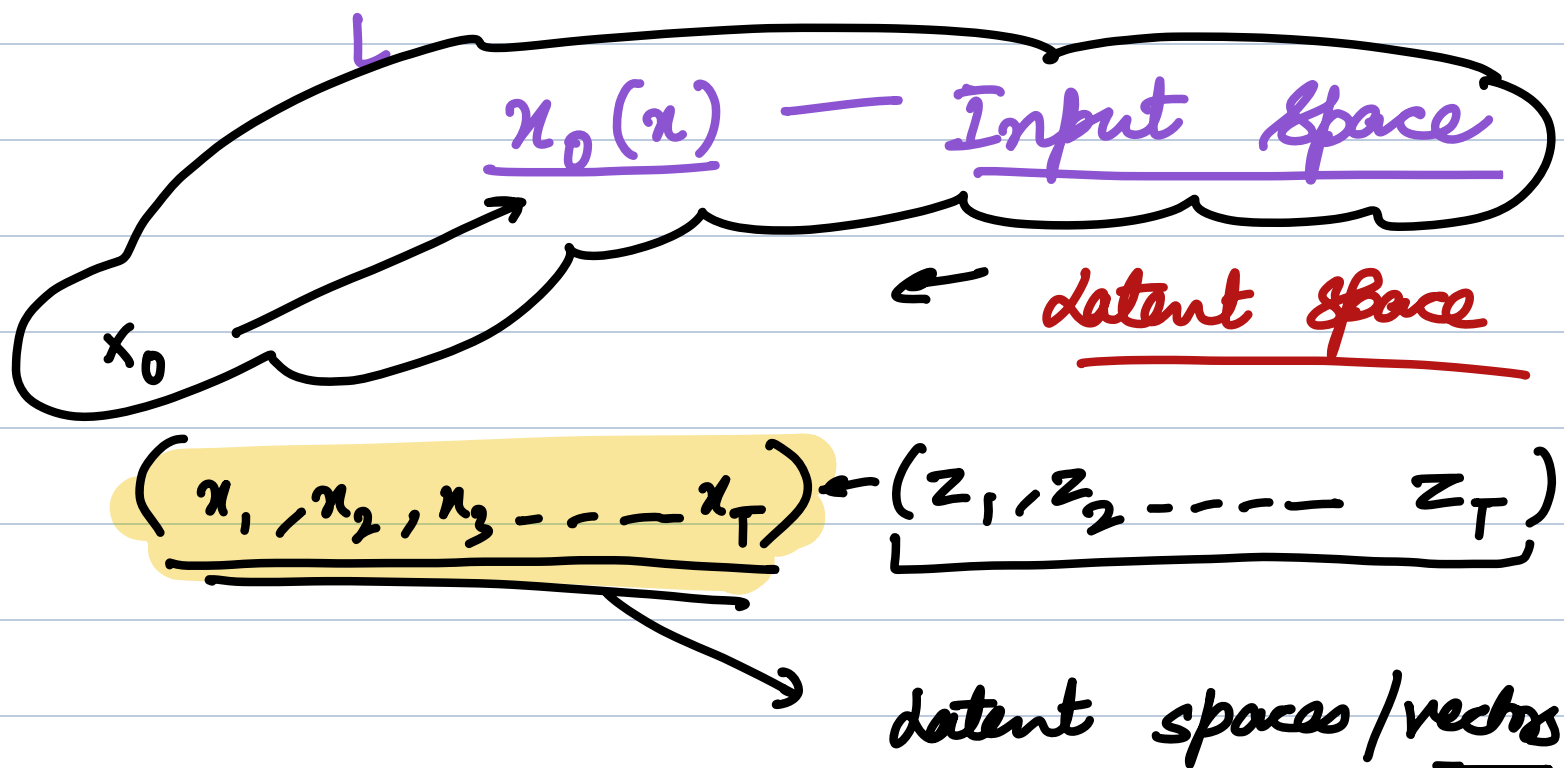
Given data $\mathcal{D} = \{x_0\} \sim P_{x_0}$ (data distribution)

(Goal) \rightarrow To learn to sample from $\underline{P_{x_0}}$

(special case of Hierarchical VAE)



$\text{VAE} \rightarrow q_\theta(z/x) \rightarrow \text{latent} \checkmark$
 both encoding & decoding learnt.
 $\text{S2PM} \rightarrow q(z/x) - \text{fixed / not learnt}$
 only decoding is learned



Forward Process

$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow x_3 \dots x_T$

input

latent / vectors / spaces corresponding to x_0

$$x_1 = \underbrace{(\sqrt{\alpha_1})}_{\text{input}} x_0 + \underbrace{(\sqrt{1-\alpha_1})}_{\text{noise}} \epsilon_1$$

$\alpha = \text{scalar}$

where $\epsilon_1 \sim N(0, I)$

$$x_2 = \sqrt{\alpha_2} x_1 + \sqrt{1 - \alpha_2} \epsilon_2$$

where $\epsilon_2 \sim N(0, I)$

$$x_3 = \sqrt{\alpha_3} x_2 + \sqrt{1 - \alpha_3} \epsilon_3$$

$$x_t = \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon_t$$

$\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_T$ are fixed scalars $\in [0, 1]$

$$\sqrt{\alpha_1} x_0 + \sqrt{1 - \alpha_1} \epsilon_1 \rightarrow \sqrt{\alpha_2} x_1 + \sqrt{1 - \alpha_2} \epsilon_2$$

$$x_T \dots x_3$$

Forward Process

$$q(x_t | x_{t-1}) = \mathcal{N} \left(x_t ; \underbrace{\sqrt{\alpha_t} x_{t-1}}_{\text{Mean}}, \underbrace{(1-\alpha_t)I}_{\text{Variance}} \right)$$

Conditional forward encoding.

output

β_t : small variance schedule

$$\alpha_t + \beta_t = 1$$

$$\boxed{\beta_t = 1 - \alpha_t}$$

Define the Model —

$$p_{\theta}(\underbrace{x_0}_{\text{data}}, \underbrace{x_1, x_2, \dots, x_T}_{\text{latent variables}}) = p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_t | x_{t-1})$$

Reverse Process

$$p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1} | x_t)$$

Mean

where

$$p_\theta(x_{t-1}|x_t) \triangleq N(x_{t-1}, \mu_\theta(x_t), \Sigma_\theta(x_t))$$

output variance

ELBO optimization
↓ VAE

VAE

$$\log p_\theta(x) = \log \int p_\theta(x, z) dz$$

$$\log \approx J_\theta(q_\phi) = E_{q_\phi(z|x)} \log \frac{p_\theta(x, z)}{q_\phi(z|x)}$$

evidence
lower bound
(ELBO)

In DDPM

$$J_\theta(q)^{\text{DDPM}} = E \log p_\theta(x_0, x_1, \dots, x_T)$$

$q(x_1, x_2, \dots, x_T | x_0)$ $q(x_1, x_2, \dots, x_T | x_0)$

$$x_{1:T} = (x_1, x_2, \dots, x_T)$$

$$x_{0:T} = (x_0, x_1, x_2, \dots, x_T)$$



$$\mathcal{J}_\theta(q)^{\text{DDPM}} = E \log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)}$$

Optimizing the ELBO for DDPM

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

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

Similarly,

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

$$q(x_{1:T} | q_0) = q\left(\frac{x_1}{x_0}\right) q(x_2 | x_1, x_0)$$

$$q(x_3 | x_2, x_1, x_0)$$

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

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

encoding
Forward
encoding