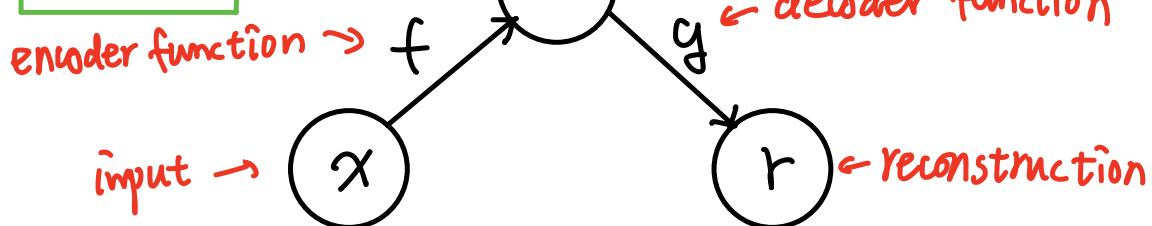


# Lecture - 02 VAE, DDPM, Discrete-VAE, D3PMs

- ① AE - Autoencoder
- ② VAE - Variational Autoencoder
- ③ DDPM - Denoising Diffusion Probabilistic Model
- ④ Discrete-VAE
- ⑤ D3PMs

## ① AE - Autoencoder

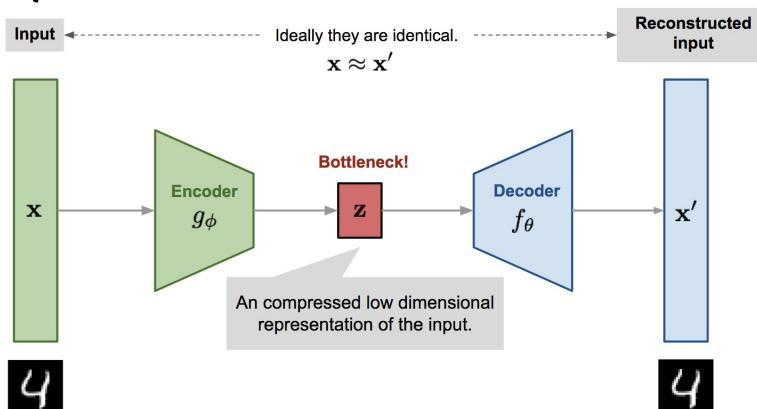
### o The Idea



training: attempt to copy its input to its output

$$r = g(f(x)) \Rightarrow \min_{f,g} L(x, g(f(x)))$$

application: dimension reduction, feature learning



• Encoder:

$$x \rightarrow z$$

• Decoder:

$$z \rightarrow x'$$

parameters =  $[\phi, \theta]$

- Ideas of autoencoder originated in 1980s, and later promoted by [Hinton and Salakhutdinov, 2006]
- Application to dimensionality reduction, like PCA
- If  $x$  are images:  $(x \approx f_\theta(g_\phi(x)))$

$$L_{AE}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - f_\theta(g_\phi(x^{(i)})))^2$$

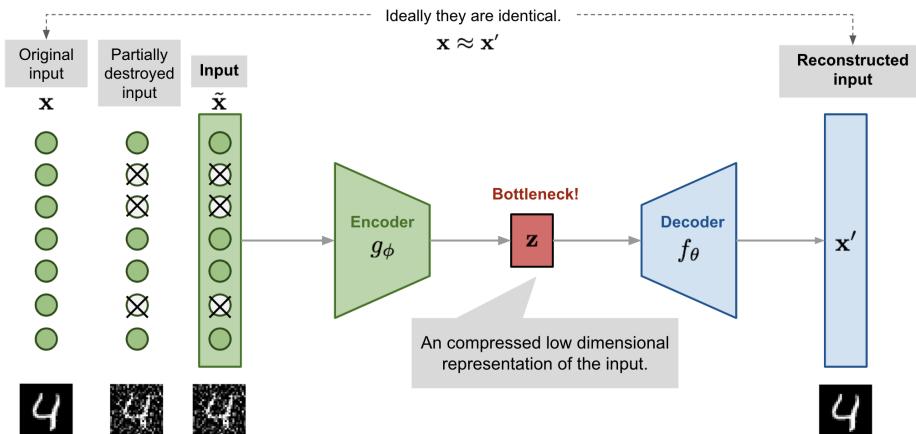
If  $x$  are categorical:

$$L_{AE}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k (-x_k^{(i)} \cdot \log \hat{x}_k^{(i)}),$$

$\hat{x}_k^{(i)}$  is the output of AE. } AE: Early overfitting

### Many variants:

- Denoising autoencoder: Add noise to  $x^{(i)}$



$$L_{DAE}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - f_\theta(g_\phi(\tilde{x}^{(i)})))^2,$$

$$\tilde{x}^{(i)} \sim M_0(\tilde{x}^{(i)} | x^{(i)})$$

- Sparse AE
- k-Sparse AE
- Contractive AE

Key limitations:

The latent space of AEs is unstructured: randomly sampling  $z$  produces meaningless outputs.

① AE - Autoencoder

② VAE - Variational Autoencoder

③ DDPM - Denoising Diffusion Probabilistic Model

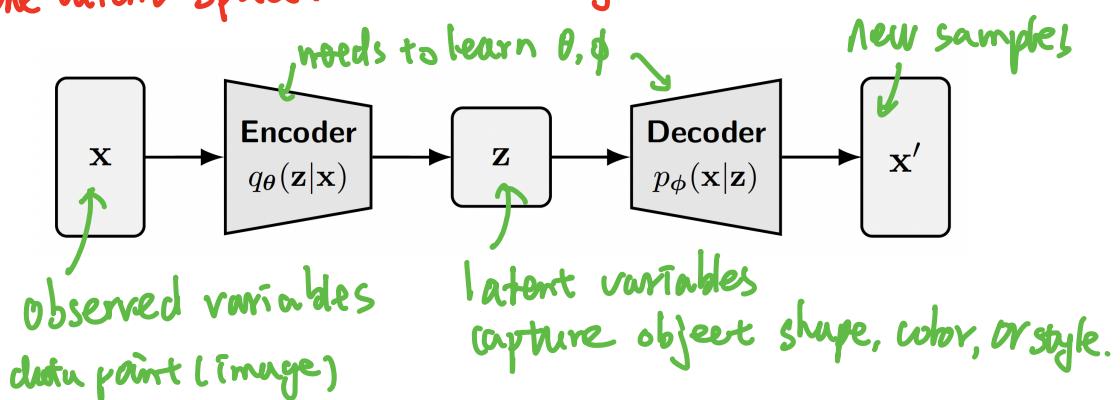
④ Discrete-VAE

⑤ D3PMs

○ New Idea : [Kingma and Welling, 2013]

Auto-Encoding Variational Bayes

From  $z \rightarrow p(z)$  : imposing a probabilistic structure on the latent space.  $\rightarrow$  a true generative model.



• generation of  $x'$ :

Step 1:  $z \sim p_{\text{prior}} := N(0, I)$

Step 2: decoding  $z$  via  $p_\phi$ :  $x' \sim p_\phi(x|z)$

latent-variable generative model:

$$p_\phi(x) = \underbrace{\int p_\phi(x|z) \cdot p(z) dz}_{\text{Intractable}} \Rightarrow \underset{\phi}{\text{goal:}} \max p_\phi(x)$$

Intractable, but can use variational method.

- construction of encoder (inference network)

given  $x$ , what latent code  $z$  could have produced it?

$$P_\phi(z|x) = \frac{P_\phi(x|z) \cdot p(z)}{P_\phi(x)} \leftarrow \text{intractable}$$

posterior

**Key:** = "variational" step in VAEs, replacing the intractable posterior with a tractable approximation!

Encoder  $q_\theta(z|x)$ : learnable proxy.

$$q_\theta(z|x) \approx P_\phi(z|x).$$

### Training via ELBO

#### Theorem.1 Evidence Lower Bound (ELBO)

For any data point  $x$ , the log-likelihood satisfies:

$$\log P_\phi(x) \geq \mathcal{L}_{\text{ELBO}}(\theta, \phi; x),$$

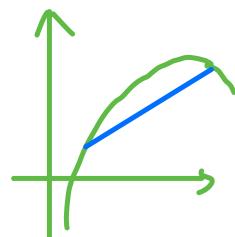
$$\mathcal{L}_{\text{ELBO}} = \underbrace{\mathbb{E}_{z \sim q_\theta(z|x)} [\log P_\phi(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\theta(z|x) || p(z))}_{\text{Latent Regularization}}$$

**Proof:**  $\log P_\phi(x) = \log \int P_\phi(x, z) dz$

$$= \log \int q_\theta(z|x) \cdot \frac{P_\phi(x, z)}{q_\theta(z|x)} dz$$

$$= \log \mathbb{E}_{z \sim q_\theta(z|x)} \left[ \frac{P_\phi(x, z)}{q_\theta(z|x)} \right]$$

$$\stackrel{\text{Jensen's Inequality}}{\geq} \mathbb{E}_{z \sim q_\theta(z|x)} \left[ \log \frac{P_\phi(x, z)}{q_\theta(z|x)} \right]$$



Note  $P_\phi(x, z) = P_\phi(x|z) \cdot P(z)$

$$\begin{aligned} \log P_\phi(x) &\geq \mathbb{E}_{z \sim q_\theta(z|x)} \left[ \log \frac{P_\phi(x|z) \cdot P(z)}{q_\theta(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_\theta} [\log P_\phi(x|z)] - D_{KL}(q_\theta(z|x) || p(z)) \end{aligned}$$

**Reconstruction:** encourages accurate recovery of  $x$  from  $z$

**Latent KL:** encourages the encoder  $q_\theta(z|x)$  to stay close to  $p_{\text{prior}}(z) \rightarrow$  enable meaningful generation.

- Information-theoretic Interpretation

Denote  $\begin{cases} \text{generative joint: } P_\phi(x, z) = p(z) \cdot P_\phi(x|z) \\ \text{inference joint: } q_\theta(x, z) = p_{\text{data}}(x) \cdot q_\theta(z|x) \end{cases}$

Recall our goal:  $D_{KL}(p_{\text{data}}(x) || P_\phi(x))$

We have:

$$D_{KL}(p_{\text{data}}(x) || P_\phi(x)) \leq D_{KL}(q_\theta(x, z) || P_\phi(x, z))$$

$\nearrow$   
chain Rule for KL Divergence

Proof:

$D_{KL}(q_\theta(x, z) || P_\phi(x, z))$  : Total error bound

$$= \mathbb{E}_{q_\theta(x, z)} \left[ \log \frac{P_{\text{data}}(x) \cdot q_\theta(z|x)}{P_\phi(x) \cdot P_\phi(z|x)} \right]$$

$$= \mathbb{E}_{p_{\text{data}}(x)} \left[ \log \frac{P_{\text{data}}(x)}{P_\phi(x)} + D_{KL}(q_\theta(z|x) || P_\phi(z|x)) \right]$$

$$\begin{aligned}
 // \text{Note : } \mathbb{E}_{q_\theta(z|x)} &= \mathbb{E}_{q_\theta(z|x) \cdot p_{\text{data}}(x)} \\
 &= \underbrace{D_{KL}(p_{\text{data}} \| p_\phi)}_{\text{True modeling error} \geq 0} + \underbrace{\mathbb{E}_{p_{\text{data}}(x)} [D_{KL}(q_\theta(z|x) \| p_\phi(z|x))]}_{\text{Inference error} \geq 0}
 \end{aligned}$$

$$\begin{aligned}
 \log p_\phi(x) &= \log \mathbb{E}_{z \sim q_\theta(z|x)} \left[ \frac{p_\phi(x, z)}{q_\theta(z|x)} \right] \\
 L_{ELBO}(\phi, \theta; x) &= \mathbb{E}_{z \sim q_\theta} [\log p_\phi(x|z)] - D_{KL}(q_\theta(z|x) \| p(z))
 \end{aligned}$$

↑   ↑  
 gap between log-likelihood      inference error.  
 and the ELBO proxy.

maximizing  $L_{ELBO} \Rightarrow$  reducing inference error.

$$D_{KL}(q_\theta(z|x) \| p_\phi(z|x)) = \int q_\theta(z|x) \log \frac{q_\theta(z|x)}{p_\phi(z|x)} dz$$

$$\begin{aligned}
 \text{Proof:} \quad &= \int q_\theta(z|x) \log \frac{q_\theta(z|x)p_\phi(x)}{p_\phi(z,x)} dz \\
 &= \int q_\theta(z|x) \left( \log p_\phi(x) + \log \frac{q_\theta(z|x)}{p_\phi(z,x)} \right) dz \\
 &\stackrel{*}{=} \log p_\phi(x) + \int q_\theta(z|x) \log \frac{q_\theta(z|x)}{p_\phi(z,x)} dz \\
 &= \log p_\phi(x) + \int q_\theta(z|x) \log \frac{q_\theta(z|x)}{p_\phi(x|z)p_\phi(z)} dz \\
 &= \log p_\phi(x) + \mathbb{E}_{z \sim q_\theta(z|x)} \left[ \log \frac{q_\theta(z|x)}{p_\phi(z)} - \log p_\phi(x|z) \right] \\
 &= \log p_\phi(x) + \underline{D_{KL}(q_\theta(z|x) \| p_\phi(z)) - \mathbb{E}_{z \sim q_\theta(z|x)} [\log p_\phi(x|z)]},
 \end{aligned}$$

where  $\stackrel{*}{=}$  is due to  $\int q_\theta(z|x) dz = 1$ .

$$\begin{aligned}
 &-L_{ELBO}(\phi, \theta; x) \\
 &= \log p_\phi(x) - L_{ELBO}(\phi, \theta; x)
 \end{aligned}$$

## ○ Gaussian VAE

- **Encoder:**  $q_{\theta}(z|x)$  is modeled as a Gaussian

$$q_{\theta}(z|x) := N(z; \mu_{\theta}(x), \text{diag}(\sigma_{\theta}^2(x))),$$

$$\text{with } \mu_{\theta}: \mathbb{R}^D \rightarrow \mathbb{R}^d, \sigma_{\theta}: \mathbb{R}^D \rightarrow \mathbb{R}_+^d$$

- **Decoder:** modeled as Gaussian

$$p_{\phi}(x|z) := N(x; \mu_{\phi}(z), \sigma^2 I),$$

$$\text{with } \mu_{\phi}(z): \mathbb{R}^d \rightarrow \mathbb{R}^D : NN$$

Under this model assumption

$$L_{\text{ELBO}} = \begin{cases} \text{Reconstruction term: } \mathbb{E}_{q_{\theta}(z|x)} [\log p_{\phi}(x|z)] \\ \text{Latent Regularization: } -D_{KL}[q_{\theta}(z|x) \parallel p_{\text{prior}}(z)] \end{cases}$$

- Reconstruction:

$$\mathbb{E}_{q_{\theta}(z|x)} [\log p_{\phi}(x|z)] = -\frac{1}{2\sigma^2} \mathbb{E}_{q_{\theta}(z|x)} [\|x - \mu_{\theta}(z)\|^2] + C$$

↓ ELBO objective

$$\min_{\theta, \sigma} \mathbb{E}_{q_{\theta}(z|x)} \left[ \frac{1}{2\sigma^2} \|x - \mu_{\theta}(z)\|^2 \right] + D_{KL}[q_{\theta}(z|x) \parallel p_{\text{prior}}(z)]$$

**Proof:** For  $D$ -dimensional Gaussian with mean  $\mu$  and

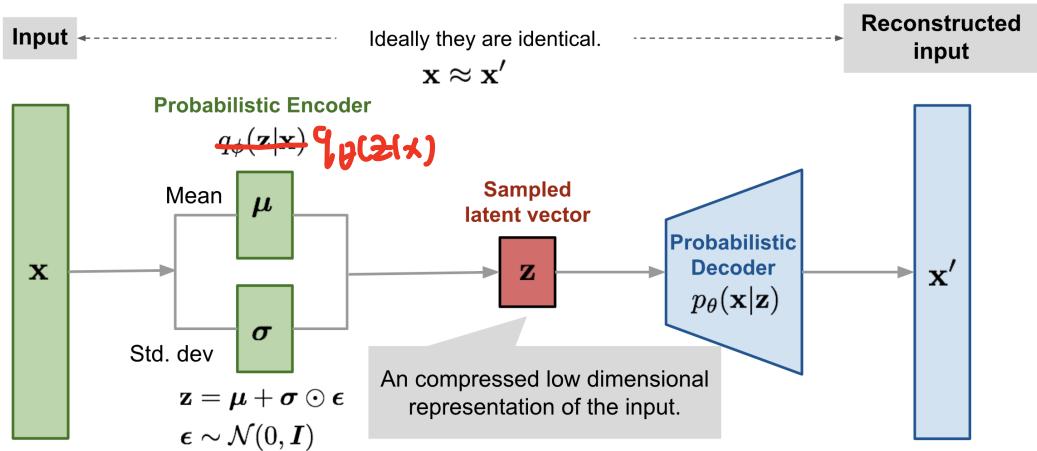
covariance  $\sigma^2 I$ :

$$p(x) = \frac{1}{(2\pi\sigma^2)^{D/2}} \cdot \exp \left( -\frac{1}{2\sigma^2} \|x - \mu\|^2 \right)$$

$$\log p(x) = -\frac{1}{2\sigma^2} \|x - \mu\|^2 - \frac{D}{2} \log(2\pi\sigma^2).$$

Replacing  $\mu$  with  $\mu_{\theta}(z)$  and  $C = -\frac{D}{2} \log(2\pi\sigma^2)$





- $q_\theta(z|x)$  For sample  $x^{(i)}$ :

$$z \sim q_\theta(z|x^{(i)}) = N(z; \mu_i, \sigma_i^2 \cdot I)$$

$$z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(0, I)$$

### ○ Drawbacks of Standard VAE

VAEs often produces blurry outputs

Consider:  $\left\{ \begin{array}{l} \text{fixed Gaussian encoder } q_{\text{enc}}(z|x) \\ \text{decoder } p_{\text{dec}}(x|z) = N(x; \mu_z, \sigma_z^2 I) \end{array} \right.$

$$\Rightarrow \text{ELBO: } \text{arg min}_\mu \mathbb{E}_{p_{\text{data}}(x)} [q_{\text{enc}}(z|x), \mathbb{E}_{p_{\text{data}}(x)} [\|x - \mu_z\|^2]]$$

$$\Rightarrow \mu^*(z) = \mathbb{E}_{q_{\text{enc}}(x|z)} [x] \quad \text{with}$$

$$q_{\text{enc}}(x|z) = \frac{q_{\text{enc}}(z|x) \cdot p_{\text{data}}(x)}{p_{\text{prior}}(z)}$$

$$\Rightarrow \mu^*(z) = \frac{\mathbb{E}_{p_{\text{data}}(x)} [q_{\text{enc}}(z|x) \cdot x]}{\mathbb{E}_{p_{\text{data}}(x)} [q_{\text{enc}}(z|x)]}$$

↑ missing point?

If two inputs  
 $x \neq x'$ , such that  
 $q_{\text{enc}}(\cdot|x)$   
 $q_{\text{enc}}(\cdot|x')$  overlap  
 $\Rightarrow \text{share } z \rightarrow \text{blurry}$

**Proof:**  $E_{x,z} [f(x, z)] = E_z [E_{x|z} [f(x, z)]]$ ,

$$\text{with } f(x, z) = \|x - \mu(z)\|^2$$

$$E_{\text{prior}} [E_{q(z|x)} [\cdot]] = E_{q(z)} [E_{q(x|z)} [\cdot]]$$

$$L(\mu) = E_{q(z)} [E_{q(x|z)} [\|x - \mu(z)\|^2]]$$

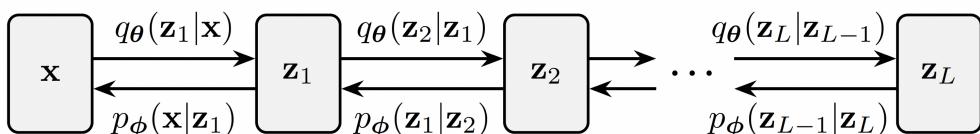
$$\arg \min_{\mu} L(\mu) = \arg \min_{\mu} E_{q(x|z)} [\|x - \mu(z)\|^2]$$

Each subproblem :

$$\min_{\mu(z)} E_{q(x|z)} [\|x - \mu(z)\|^2] \Leftrightarrow \mu^*(z) = E_{q(x|z)} [x]$$



### From VAE to Hierarchical VAEs



HVAEs introduce multiple layers of latent code  $z$

$$P_{\theta}(x, z_{1:L}) = p_{\phi}(x|z_1) \cdot \prod_{i=1}^L p_{\phi}(z_{i-1}|z_i) \cdot p(z_L)$$

$$P_{\text{HVAE}}(x) := \int P_{\theta}(x, z_{1:L}) d z_{1:L}$$

$$q_{\theta}(z_{1:L}|x) = q_{\theta}(z_1|x) \cdot \prod_{i=2}^L q_{\theta}(z_i|z_{i-1})$$

HVAE's ELBO:

$$\begin{aligned}
\log p_{\text{HVAE}}(\mathbf{x}) &= \log \int p_{\phi}(\mathbf{x}, \mathbf{z}_{1:L}) d\mathbf{z}_{1:L} \\
&= \log \int \frac{p_{\phi}(\mathbf{x}, \mathbf{z}_{1:L})}{q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x})} q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x}) d\mathbf{z}_{1:L} \\
&= \log \mathbb{E}_{q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x})} \left[ \frac{p_{\phi}(\mathbf{x}, \mathbf{z}_{1:L})}{q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x})} \right] \\
&\geq \mathbb{E}_{q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x})} \left[ \log \frac{p_{\phi}(\mathbf{x}, \mathbf{z}_{1:L})}{q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x})} \right] \\
&=: \mathcal{L}_{\text{ELBO}}(\phi).
\end{aligned}$$

Substituting the factorized forms yields:

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{q_{\theta}(\mathbf{z}_{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{z}_L) \prod_{i=2}^L p_{\phi}(\mathbf{z}_{i-1}|\mathbf{z}_i) p_{\phi}(\mathbf{x}|\mathbf{z}_1)}{q_{\theta}(\mathbf{z}_1|\mathbf{x}) \prod_{i=2}^L q_{\theta}(\mathbf{z}_i|\mathbf{z}_{i-1})} \right].$$

### Observation 2 (Key observations from HVAE)

*Stacking layers allows the model to generate data progressively, starting with coarse details and adding finer ones at each step. This process makes it far easier to capture the complex structure of high-dimensional data.*

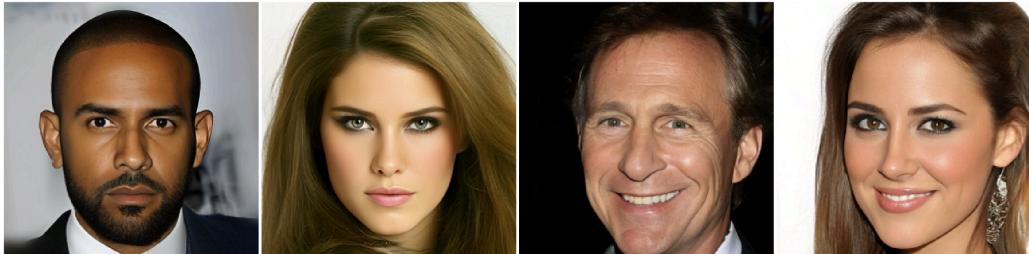
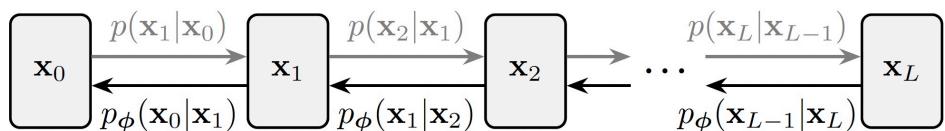


Figure 1: 256×256-pixel samples generated by NVAE, trained on CelebA HQ [28].

- Deeper Networks in a flat VAE are not enough.
- variational family is limited:
  - Greater network depth improves  $\mathbb{M}_{\theta}$  and  $\mathbb{D}_{\theta}$  but does not expand the family
  - if the decoder is too expensive, the model may suffer from posterior collapse

- ① AE - Autoencoder
- ② VAE - Variational Autoencoder
- ③ DDPM - Denoising Diffusion Probabilistic Model
- ④ Discrete-VAE
- ⑤ D3PMs



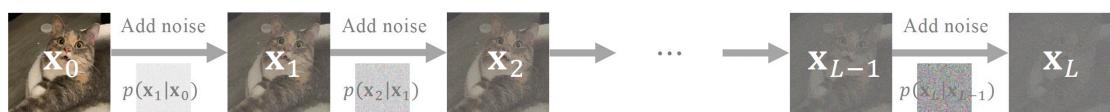
DDPMs { ①. Forward Pass (Fixed Encoder)  
 ②. Reverse Denoising Process (Learnable Decoder)

- DDPMs achieve remarkable stability and expressive power.

### ○ Forward Process (Fixed Encoder)

$$x_0 \sim p_{\text{data}} \longrightarrow x_L \sim p_{\text{prior}} := N(0, I)$$

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



Each step is governed by a fixed Gaussian transition kernel

$$p(x_i|x_{i-1}) := N(x_i; \sqrt{1-\beta_i^2}x_{i-1}, \beta_i^2 I),$$

$i=1, 2, \dots, L$  with  $x_0 \sim p_{\text{data}}$ .

$\{\beta_i\}_{i=1}^L$ : pre-determined or learned,  $\beta_i \uparrow$  and  $\beta_i \in (0, 1)$

Define  $\lambda_i := \sqrt{1 - \beta_i^2}$ . Based on the definition of  $p(x_i | x_{i-1})$ ,

$$x_i = \lambda_i \cdot x_{i-1} + \beta_i \cdot \varepsilon_i, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, I)$$

$\curvearrowleft$  controlled noise

[This is from the fact of sampling from a Gaussian]

If  $X \sim N(\mu, \sigma^2 I)$ , the the sample can be generated as:  $X = \mu + \sigma \cdot \varepsilon, \quad \varepsilon \sim N(0, I)$ . ]

Applying the transition kernels recursively:

$$x_1 = \lambda_1 x_0 + \beta_1 \varepsilon_1$$

$$x_2 = \lambda_2 x_1 + \beta_2 \varepsilon_2$$

$$= \lambda_2 (\lambda_1 x_0 + \beta_1 \varepsilon_1) + \beta_2 \varepsilon_2$$

$$= \lambda_2 \lambda_1 x_0 + \lambda_2 \beta_1 \varepsilon_1 + \beta_2 \varepsilon_2$$

$$x_3 = \lambda_3 x_2 + \beta_3 \varepsilon_3$$

$$\vdots \\ = \lambda_3 \lambda_2 \lambda_1 x_0 + \lambda_3 \lambda_2 \beta_1 \varepsilon_1 + \lambda_3 \beta_2 \varepsilon_2 + \beta_3 \varepsilon_3$$

$$\Rightarrow x_i = \left( \prod_{j=1}^i \lambda_j \right) \cdot x_0 + \underbrace{\sum_{k=1}^i \left( \beta_k \cdot \prod_{j=k+1}^i \lambda_j \right) \cdot \varepsilon_k}_{\sum_{k=1}^i \beta_k \cdot \frac{\bar{\lambda}_i}{\bar{\lambda}_k}} \quad \text{Note, they are i.i.d.}$$

$$\text{Define } \bar{\lambda}_i = \prod_{j=1}^i \lambda_j$$

$$\text{To show } \text{Var} \left[ \sum_{k=1}^i \beta_k \cdot \frac{\bar{\lambda}_i}{\bar{\lambda}_k} \varepsilon_k \right] = 1 - \bar{\lambda}_i^2 \quad \text{telescoping}$$

$$\text{Note } 1 - \bar{\lambda}_i^2 = 1 - \prod_{j=1}^i \lambda_j^2 \quad \text{and} \quad \begin{cases} 1 - \lambda_i^2 = \beta_i^2 \\ 1 - \lambda_i^2 \lambda_{i-1}^2 = (1 - \lambda_i^2) + (1 - \lambda_{i-1}^2) \cdot \lambda_i^2 \\ 1 - \lambda_i^2 \lambda_{i-1}^2 \lambda_{i-2}^2 = (1 - \lambda_i^2) + (1 - \lambda_{i-1}^2) \cdot \lambda_i^2 + (1 - \lambda_{i-2}^2) \cdot \lambda_{i-1}^2 \cdot \lambda_i^2 \end{cases}$$

$$\sum_{k=1}^i \left( \beta_k \cdot \frac{\bar{\lambda}_i}{\bar{\lambda}_k} \right)^2 = \sum_{k=1}^i (1 - \lambda_k^2) \cdot \frac{\bar{\lambda}_i^2}{\bar{\lambda}_k^2}$$

So, the mean is  $\bar{\lambda}_i \cdot x_0$  and variance is  $(1 - \bar{\lambda}_i^2) \cdot I$

Perturbation Kernel and Prior Distribution:

$$p_i(x_i|x_0) = N(x_i; \bar{\alpha}_i \cdot x_0, (1-\bar{\alpha}_i^2) \cdot I),$$

$$\bar{\alpha}_i = \prod_{k=1}^i \alpha_k = \prod_{k=1}^i \sqrt{1-\beta_k^2} = \prod_{k=1}^i \alpha_k$$

Sampling:  $x_i = \bar{\alpha}_i \cdot x_0 + \sqrt{1-\bar{\alpha}_i^2} \cdot \xi, \xi \sim N(0, I)$

Since  $\{\beta_i\}_{i=1}^L \downarrow 0$  and  $\beta_i \in (0, 1)$ ,  $\beta_L \rightarrow 1$ .

$$\Rightarrow \sqrt{1-\beta_L^2} \rightarrow 0, \bar{\alpha}_i \rightarrow 0 \Rightarrow \sqrt{1-\bar{\alpha}_i^2} \rightarrow 1.$$

$p_L(x_L|x_0) \rightarrow N(0, I)$  as  $L \rightarrow \infty$ ,

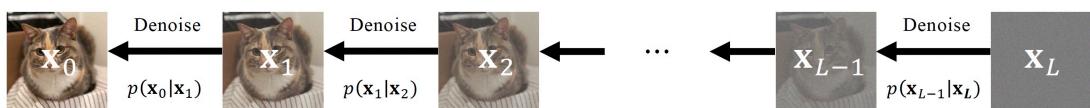
which motivates the choice of the Prior as:

$p_{\text{prior}} := N(0, I) \leftrightarrow \text{No reliance on data } x_0.$

### Reverse Denoising Process (Learnable Decoder)

$x_L \sim p_{\text{prior}} \xrightarrow{\text{denoise}} x_0 : \text{meaningful data point}$

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



key Q: Can we effectively approximate these reverse transition kernels  $p(x_{i-1}|x_i)$  especially when  $x_i \sim p_i(x_i)$  is complex?

To approximate  $p(x_{i-1}|x_i)$ , we use  $p_\phi(x_{i-1}|x_i)$  and minimize

$$(*) \quad \min_{\phi} [E_{p_i(x_i)} [ D_{KL}(p(x_{i-1}|x_i) || p_\phi(x_{i-1}|x_i)) ]]$$

Note  $p(x_{i-1}|x_i) = p(x_i|x_{i-1}) \cdot \frac{p_{i-1}(x_{i-1})}{p_i(x_i)}$  ↗ intractable

as  $p_i(x_i) = \int p_i(x_i|x_0) \cdot \underbrace{p_{\text{data}}(x_0)}_{\text{no closed-form}} dx_0$

**Key Idea:** we condition the reverse transition on a clean data sample  $x$ .  $\rightarrow$  make the kernel tractable.

$$\Rightarrow p(x_{i-1}|x_i, x) = \frac{p(x_i|x_{i-1}, x)}{\frac{p(x_i|x_{i-1})}{p(x_{i-1}|x)}}$$

(Markov property of  $p_{\text{data}}(x_{i-1}, x) = p(x_i|x_{i-1})$  from the forward process assumption)

$\Leftrightarrow$  [ Marginal KL Minimization  $\Leftrightarrow$  Conditional KL Min.]

### Theorem 2.2.1: Equivalence Between Marginal and Conditional KL Minimization

The following equality holds:

$$\begin{aligned} & \mathbb{E}_{p_i(\mathbf{x}_i)} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i))] \\ &= \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{p(\mathbf{x}_i|\mathbf{x})} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i))] + C, \end{aligned} \quad (2.2.3)$$

where  $C$  is a constant independent of  $\phi$ . Moreover, the minimizer of Equation (2.2.3) satisfies

$$p^*(\mathbf{x}_{i-1}|\mathbf{x}_i) = \mathbb{E}_{p(\mathbf{x}|\mathbf{x}_i)} [p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x})] = p(\mathbf{x}_{i-1}|\mathbf{x}_i), \quad \mathbf{x}_i \sim p_i.$$

#### Proof. Derivation of Equation (2.2.3).

We start by expanding the right-hand side expectation:

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_i)} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i))] \\ &= \int \int p(\mathbf{x}_0, \mathbf{x}_i) \mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)) d\mathbf{x}_0 d\mathbf{x}_i. \end{aligned}$$

By the definition of KL divergence,

$$\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)) = \int p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} d\mathbf{x}_{i-1}.$$

Substituting this into the expectation, we have

$$\int \int \int p(\mathbf{x}_0, \mathbf{x}_i) p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} d\mathbf{x}_{i-1} d\mathbf{x}_0 d\mathbf{x}_i.$$

Using the chain rule of probability,

$$p(\mathbf{x}_0, \mathbf{x}_i) = p(\mathbf{x}_i)p(\mathbf{x}_0|\mathbf{x}_i),$$

we rewrite the integral as

$$\int p(\mathbf{x}_i) \int p(\mathbf{x}_0|\mathbf{x}_i) \int p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} d\mathbf{x}_{i-1} d\mathbf{x}_0 d\mathbf{x}_i.$$

This allows us to express the expectation in nested form:

$$\mathbb{E}_{p(\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} \right] \right] \right].$$

Next, we apply the decomposition of the logarithm:

$$\log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} = \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p(\mathbf{x}_{i-1}|\mathbf{x}_i)} + \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)}.$$

Substituting this back into the expectation gives two terms:

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)}{p(\mathbf{x}_{i-1}|\mathbf{x}_i)} \right] \right] \right] \\ & + \mathbb{E}_{p(\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} \right] \right] \right]. \end{aligned}$$

Since the second logarithmic term does not depend on  $\mathbf{x}_0$ , by the law of total probability

$$\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} \right] \right] = \mathbb{E}_{p(\mathbf{x}_{i-1}|\mathbf{x}_i)} \left[ \log \frac{p(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i)} \right].$$

Similarly, the first term is the KL divergence

$$\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_i)} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \| p(\mathbf{x}_{i-1}|\mathbf{x}_i))].$$

Putting it all together, we obtain the decomposition:

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_i)} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i))] \\ & = \mathbb{E}_{p(\mathbf{x}_i)} \left[ \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_i)} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x}_0) \| p(\mathbf{x}_{i-1}|\mathbf{x}_i))] \right] \\ & + \mathbb{E}_{p(\mathbf{x}_i)} [\mathcal{D}_{\text{KL}}(p(\mathbf{x}_{i-1}|\mathbf{x}_i) \| p_\phi(\mathbf{x}_{i-1}|\mathbf{x}_i))]. \end{aligned}$$

**Proof of Optimality.** To prove:

$$p^*(\mathbf{x}_{i-1}|\mathbf{x}_i) = p(\mathbf{x}_{i-1}|\mathbf{x}_i) = \mathbb{E}_{p(\mathbf{x}|\mathbf{x}_i)} [p(\mathbf{x}_{i-1}|\mathbf{x}_i, \mathbf{x})], \quad \mathbf{x}_i \sim p_i.$$

The first identity follows from the fact that the KL divergence  $\mathcal{D}_{\text{KL}}(p \| p_\phi)$  is minimized when  $p^* = p$ , assuming the parameterization is sufficiently expressive. The second identity follows directly from the law of total probability.

■

Conditioning to obtain a tractable objective, forms  
the foundation of DDPMs.



## ○ DDPM Sampling:

- $\{\xi_{\phi^x}(x_i, i) \mid i=1, 2, \dots, L\} \leftarrow \text{after training}$

- $x_L \sim P_{\text{prior}} = N(0, I)$  : sampling procedure

$$x_L \xrightarrow{P_{\phi^x}(x_{L-1} | x_L)} x_{L-1} \xrightarrow{P_{\phi^x}(x_{L-2} | x_{L-1})} x_{L-2} \dots \xrightarrow{P_{\phi^x}(x_0 | x_1)} x_0$$

$$x_{i-1} \leftarrow \frac{1}{\bar{\alpha}_i} \left( x_i - \frac{1-\bar{\alpha}_i^2}{\sqrt{1-\bar{\alpha}_i^2}} \xi_{\phi^x}(x_i, i) \right) + \sigma(i) \cdot \epsilon_i.$$

$\underbrace{\xi_{\phi^x}(x_i, i)}$

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

- Interpretation of DDPM's sampling:

$x_i \leftarrow \text{interpolation between } \underline{x_i} \text{ and } \underline{\xi_{\phi^x}(x_i, i) + \sigma(i) \cdot \epsilon_i}$

$$\Leftrightarrow x_i = A \cdot \underline{x_i} + B \cdot \underline{\xi_{\phi^x}(x_i, i) + \sigma(i) \cdot \epsilon_i}$$

Proof: By  $x_i = \bar{\alpha}_i \cdot x_{\phi^x}(x_i, i) + \sqrt{1-\bar{\alpha}_i^2} \xi_{\phi^x}(x_i, i)$

$$\Rightarrow x_{\phi^x}(x_i, i) = \frac{x_i - \sqrt{1-\bar{\alpha}_i^2} \cdot \xi_{\phi^x}(x_i, i)}{\bar{\alpha}_i}$$

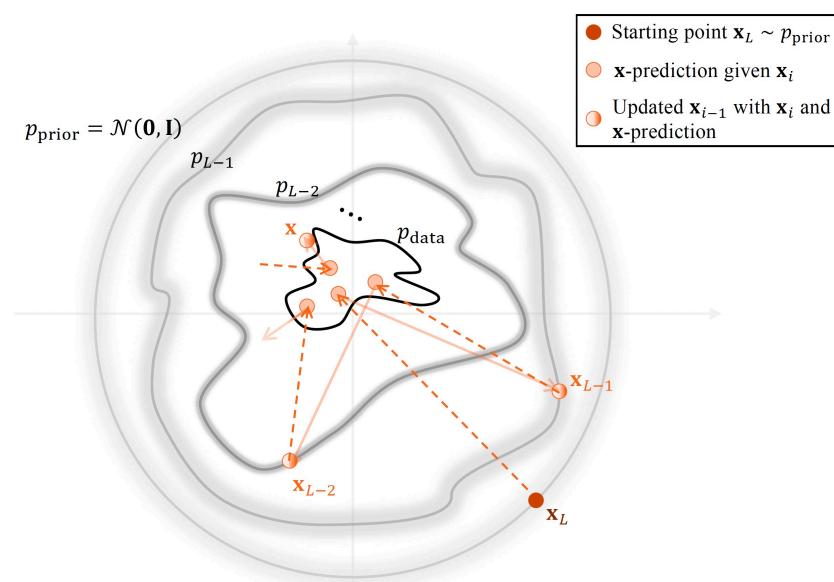
Rearranging to isolate  $\xi_{\phi^x}$ :

$$\xi_{\phi^x}(x_i, i) = \frac{x_i - \bar{\alpha}_i \cdot x_{\phi^x}(x_i, i)}{\sqrt{1-\bar{\alpha}_i^2}}$$

$$\Rightarrow M_{\phi^x}(x_i, i) = \underbrace{\left( \frac{1}{\bar{\alpha}_i} - \frac{1-\bar{\alpha}_i^2}{\bar{\alpha}_i(1-\bar{\alpha}_i^2)} \right)}_A \cdot x_i + \underbrace{\left( \frac{(1-\bar{\alpha}_i^2) \cdot \bar{\alpha}_i}{\bar{\alpha}_i(1-\bar{\alpha}_i^2)} \right)}_B \cdot x_{\phi^x}(x_i, i)$$

DDPM can be viewed as an iterative denoising process:

- (1) : Estimating clean data  $x_{\phi^x}(x_i, i)$  from current noisy input  $x_i$ .
- (2) : Sampling less noisy latent  $x_{i-1}$  via the update rule using this clean estimate.



### ○ Remarks of Sampling via DDPM:

- (1) : Blurry predictions at high noise levels

$$x^*(x_i, i) = \mathbb{E}[x_0 | x_i], x_i \sim p_i$$

the model predict the expected clean data given  $x_i$

- (2) : Early steps set the global structure  
Later steps add fine detail

- (3) : Slow sampling :

$L = 1000$  steps in practice

- (4) In practice,  $p_\phi(x_{i-1}|x_i)$  modeled as Gaussian

$\left\{ \begin{array}{l} \text{small } \beta_i \rightarrow \text{true reverse distribution closer to Gaussian} \\ \text{large } \beta_i \rightarrow \text{induce strong non-Gaussian.} \end{array} \right.$

References:

- [1] <https://lilianweng.github.io/posts/2018-08-12-vae/>