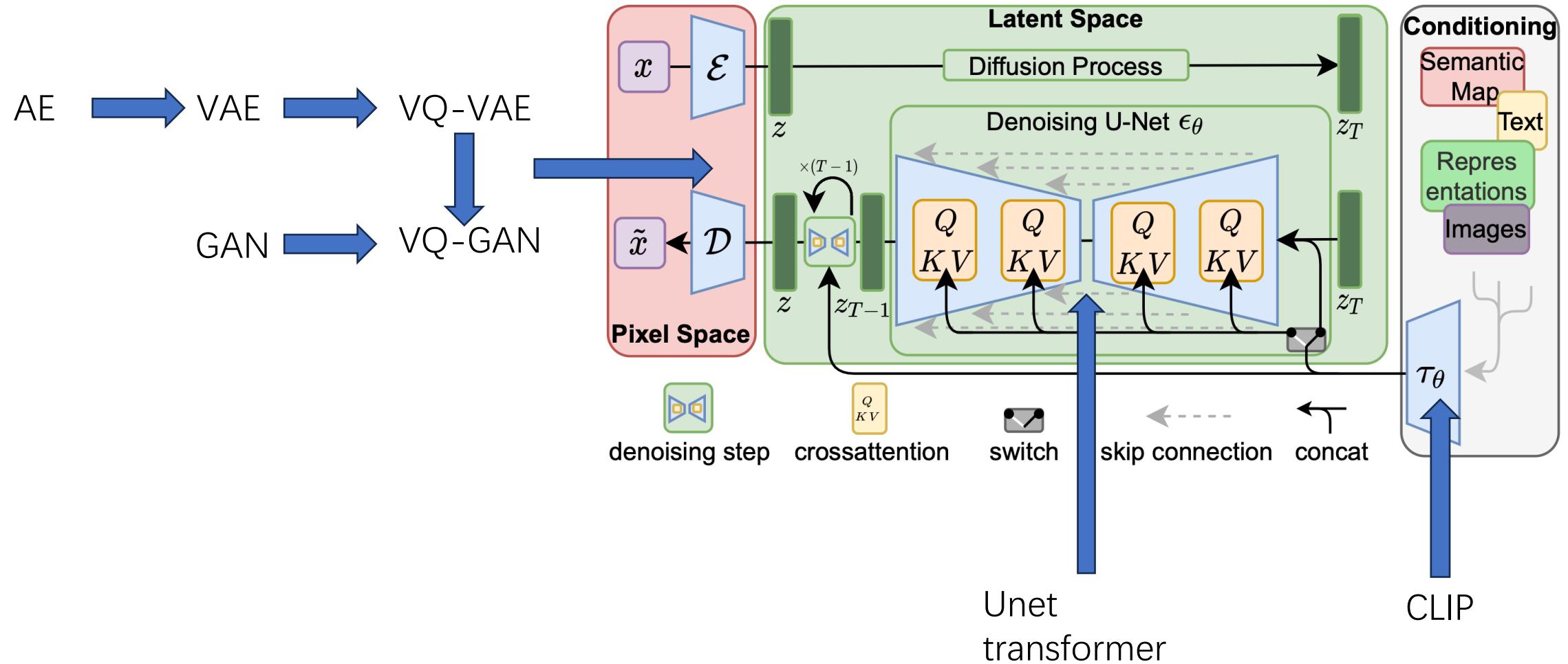
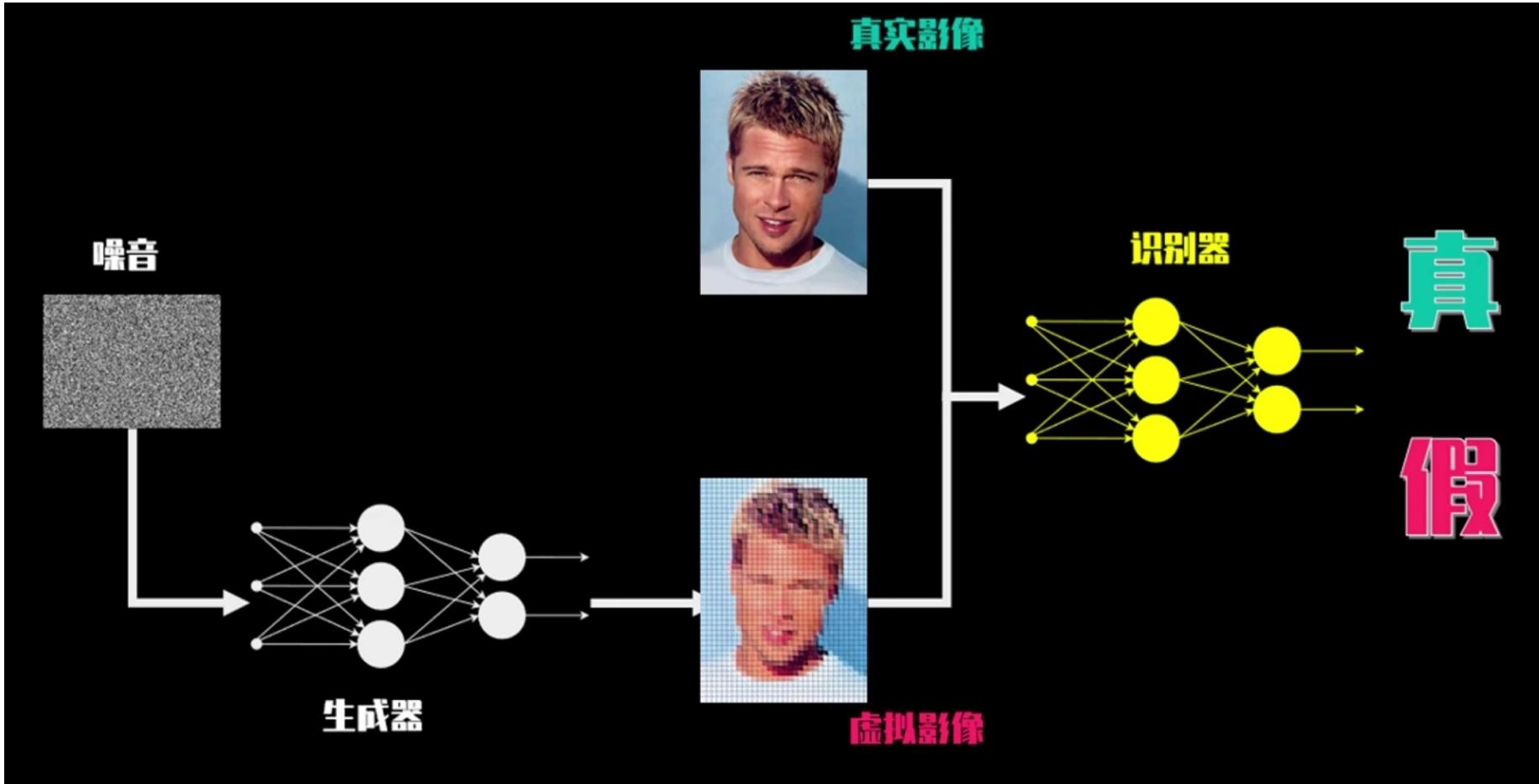


# LDM预备知识

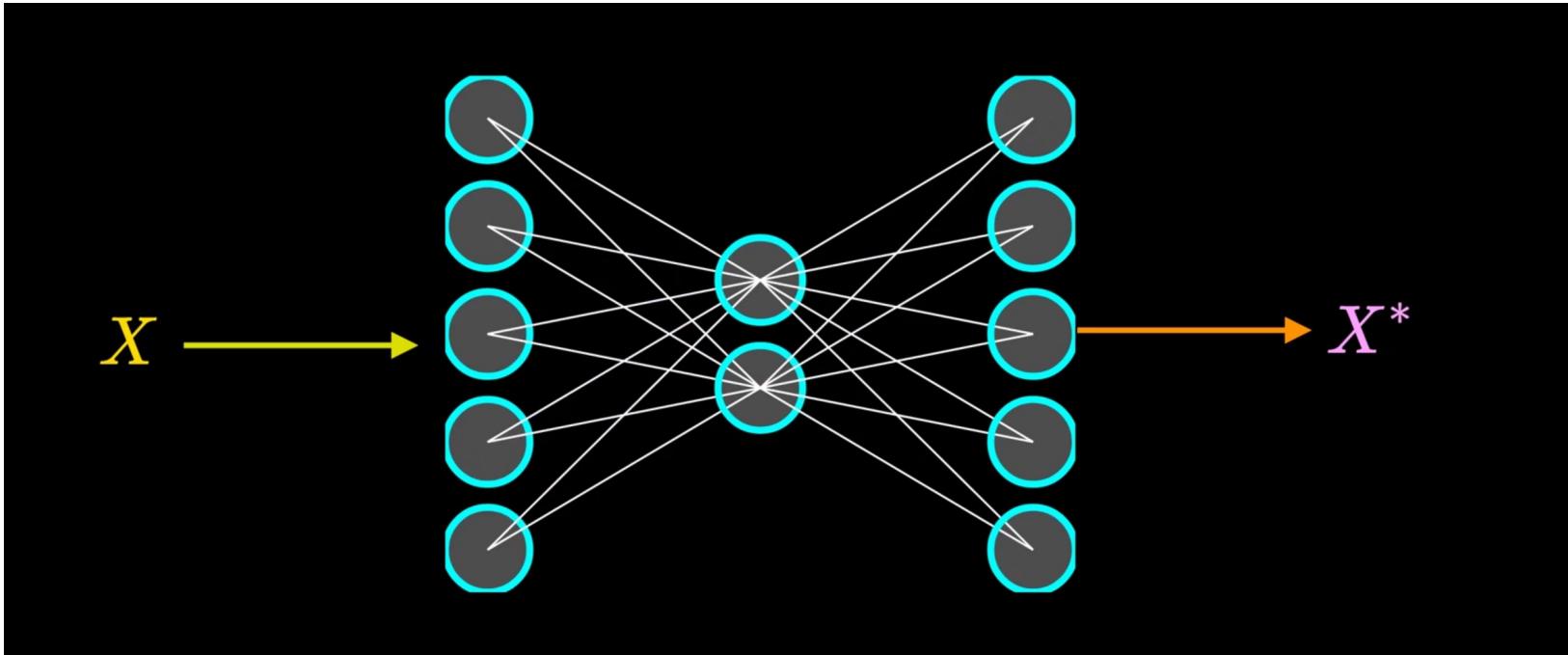
# LDM -Latent Diffusion Models -潜在扩散模型



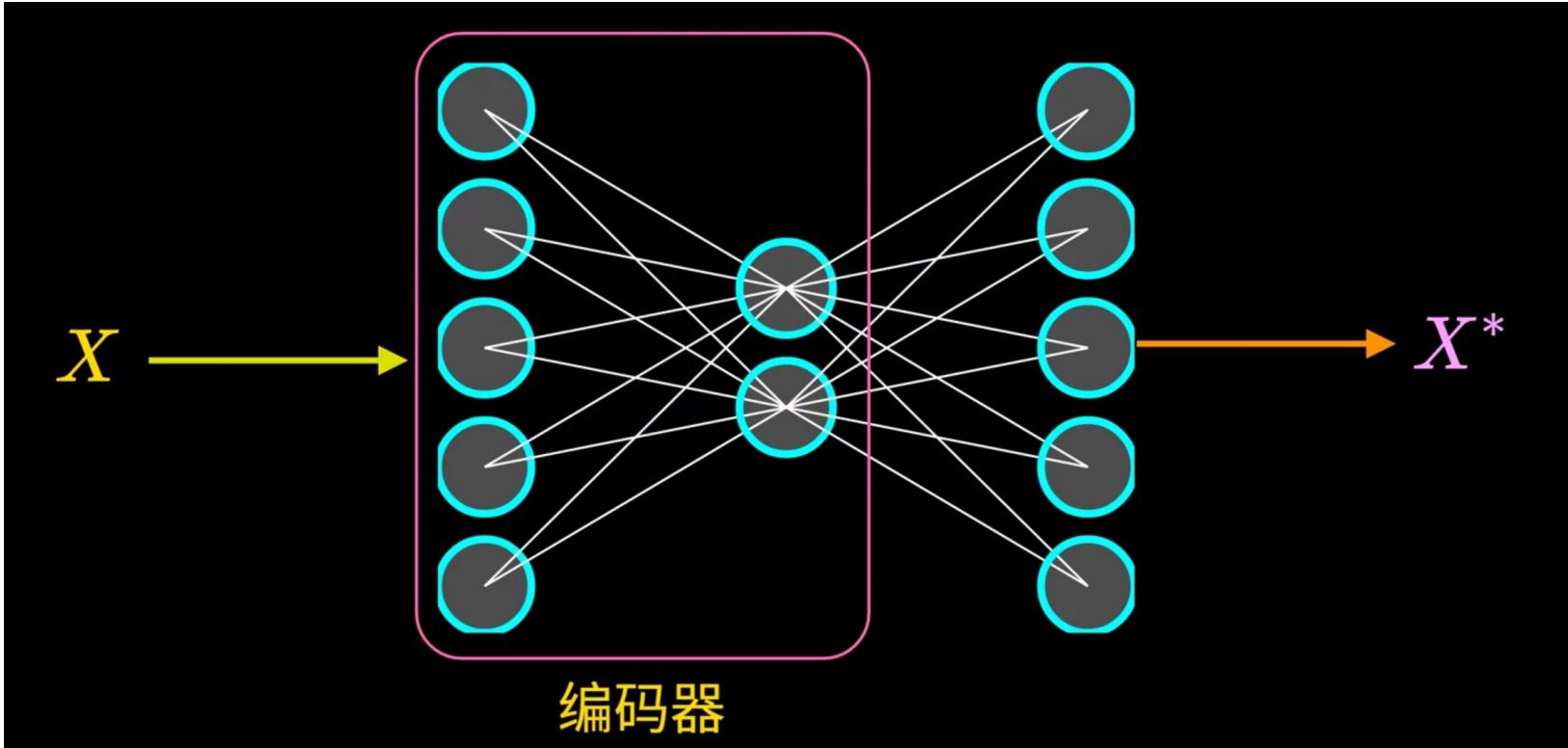
# GAN -Generative Adversarial Nets -生成对抗网络



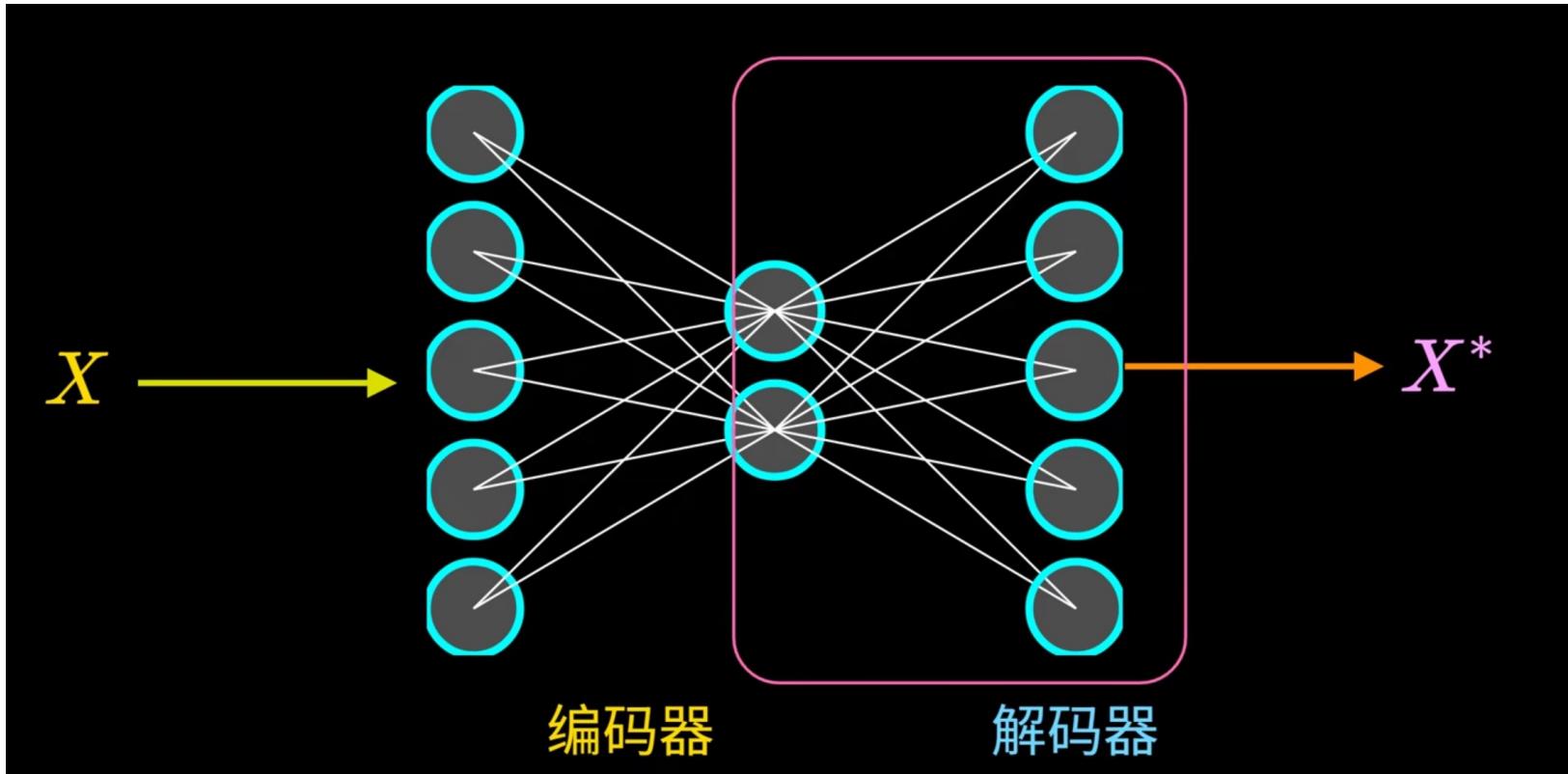
AE -autoencoder -自编码器



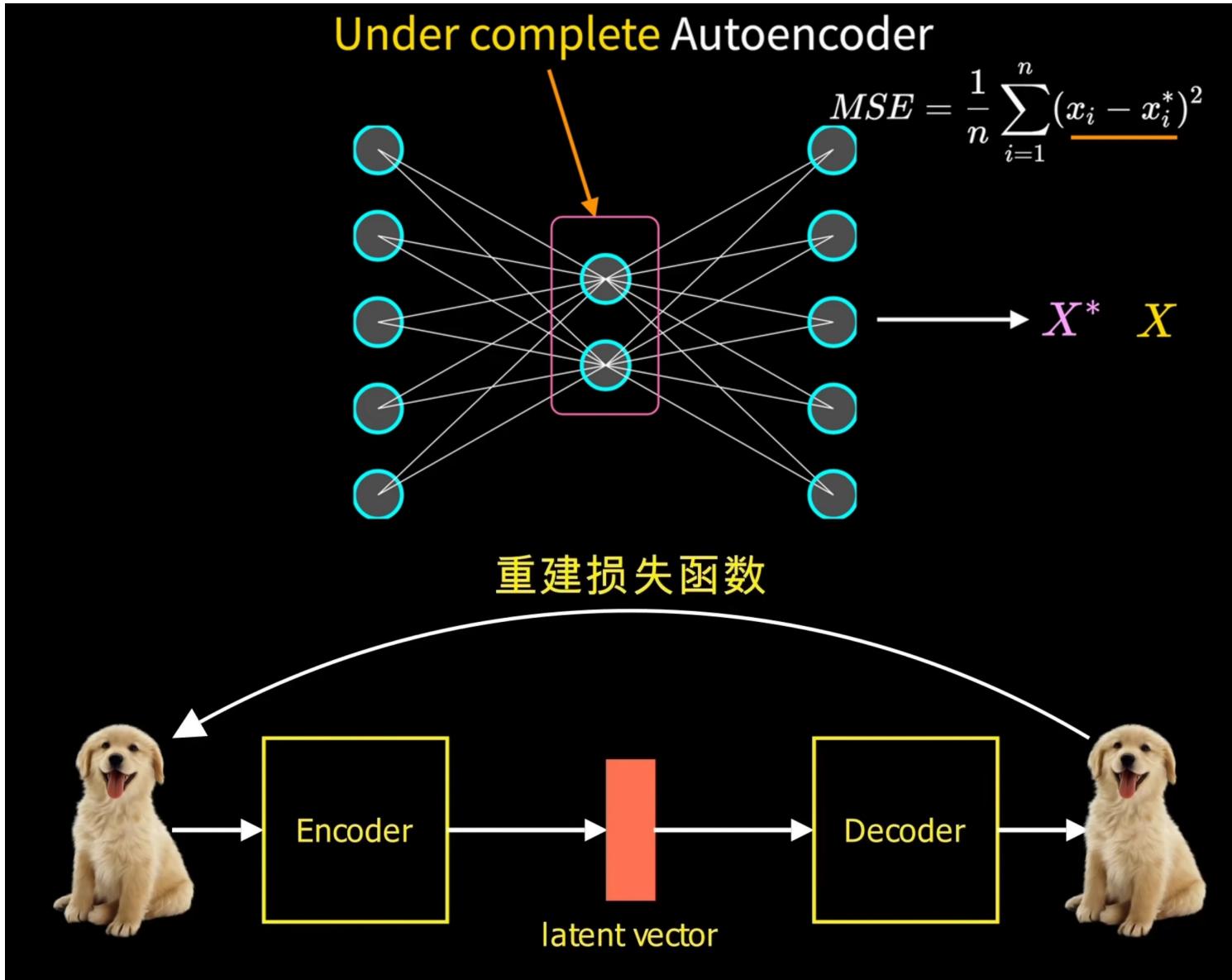
AE -autoencoder -自编码器



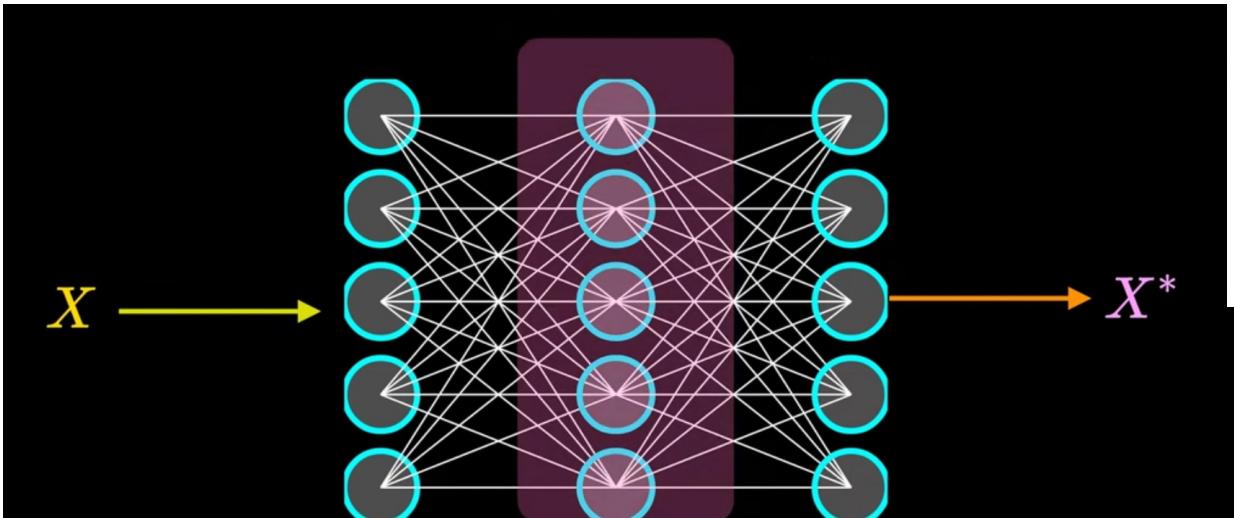
AE -autoencoder -自编码器



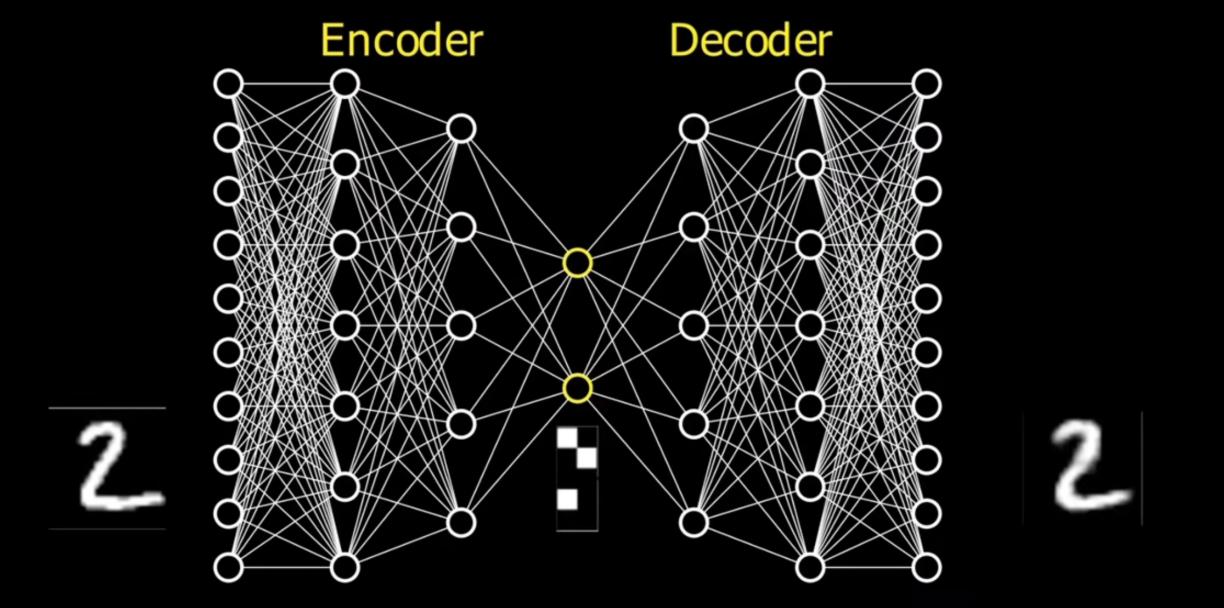
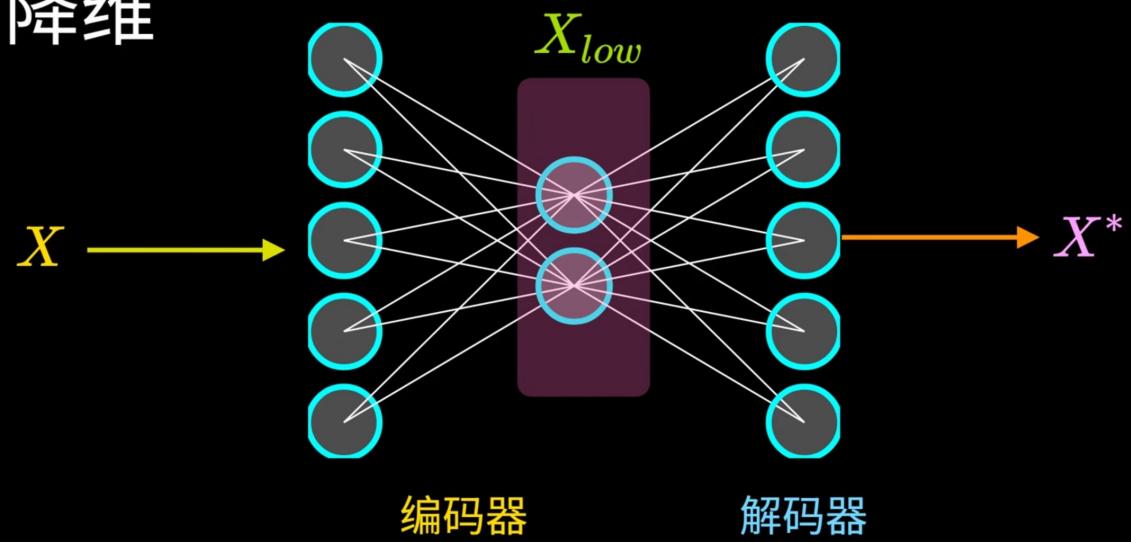
AE -autoencoder -自编码器



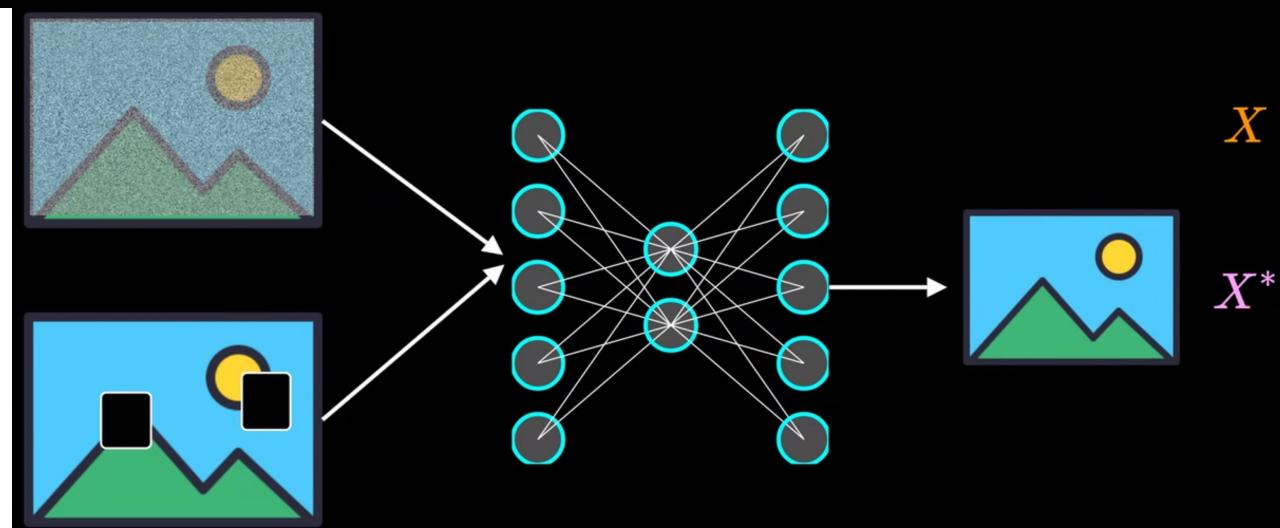
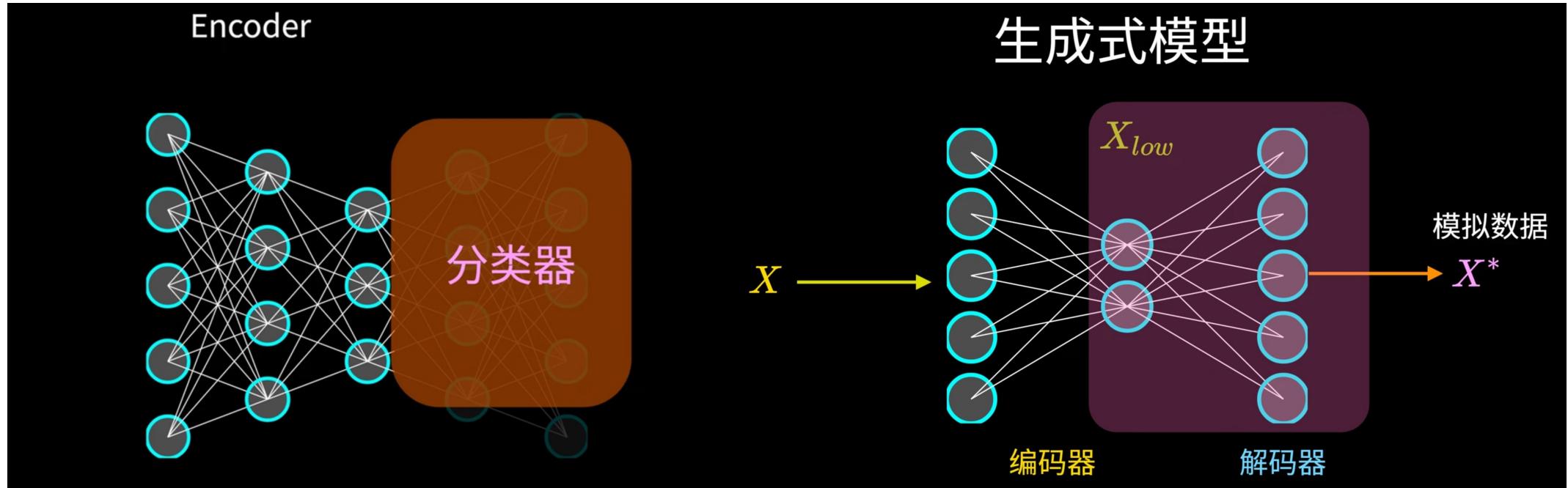
AE -autoencoder -自编码器



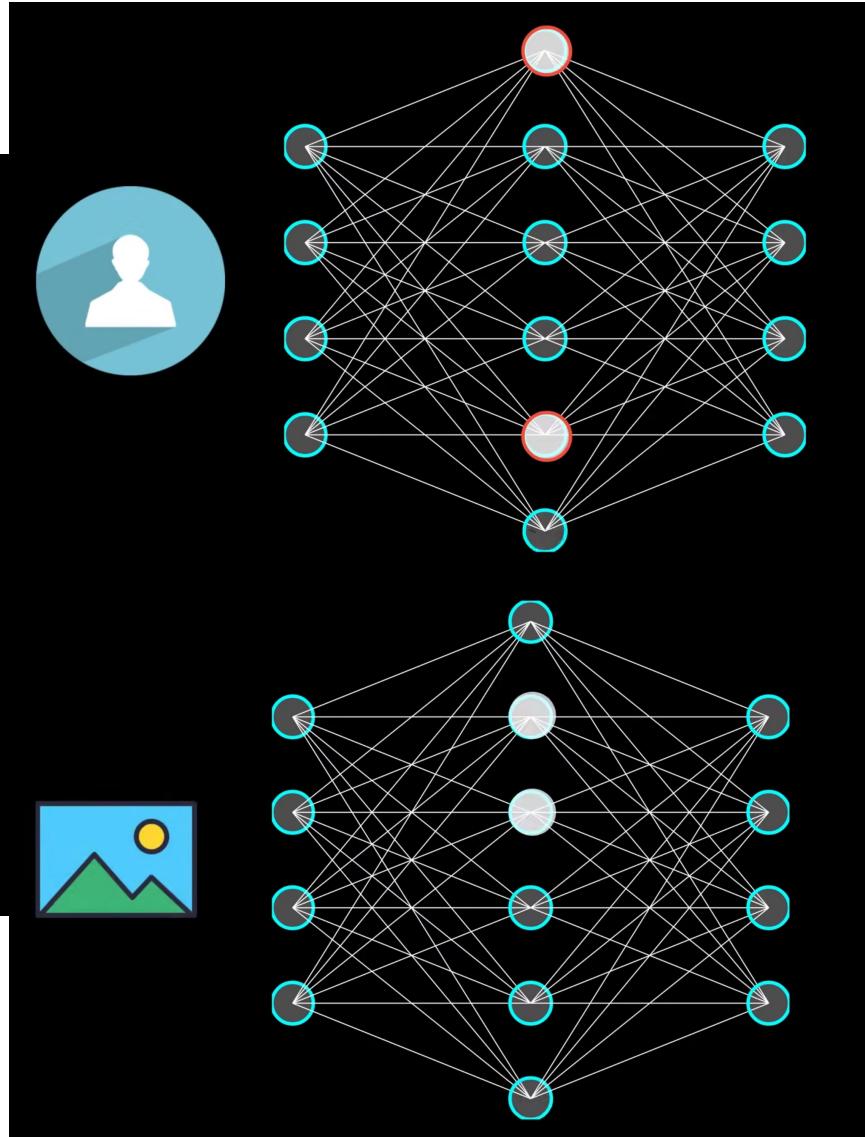
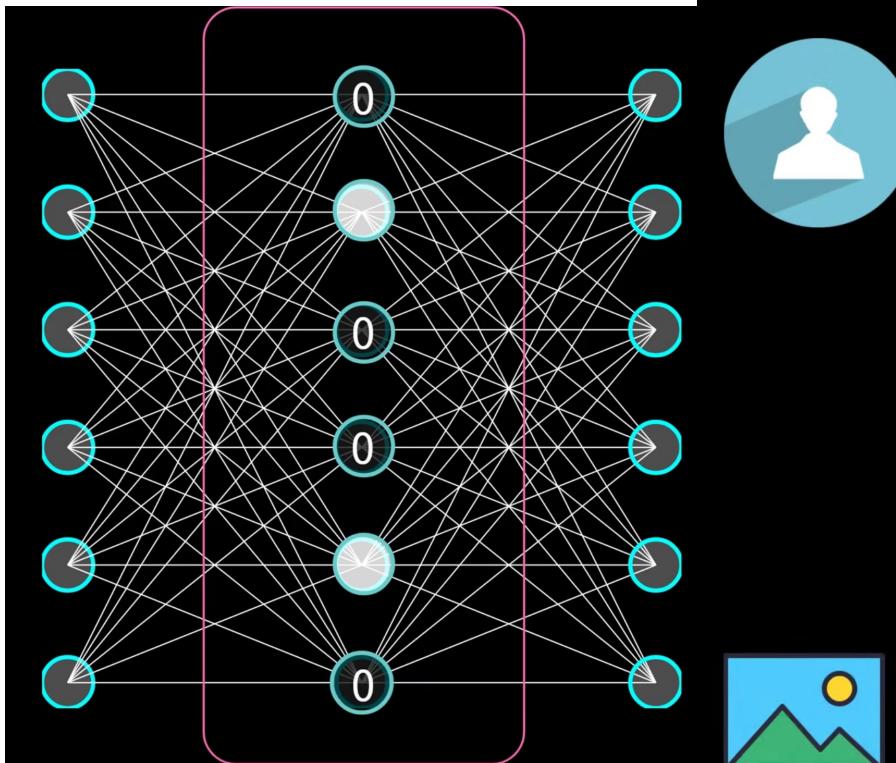
降维



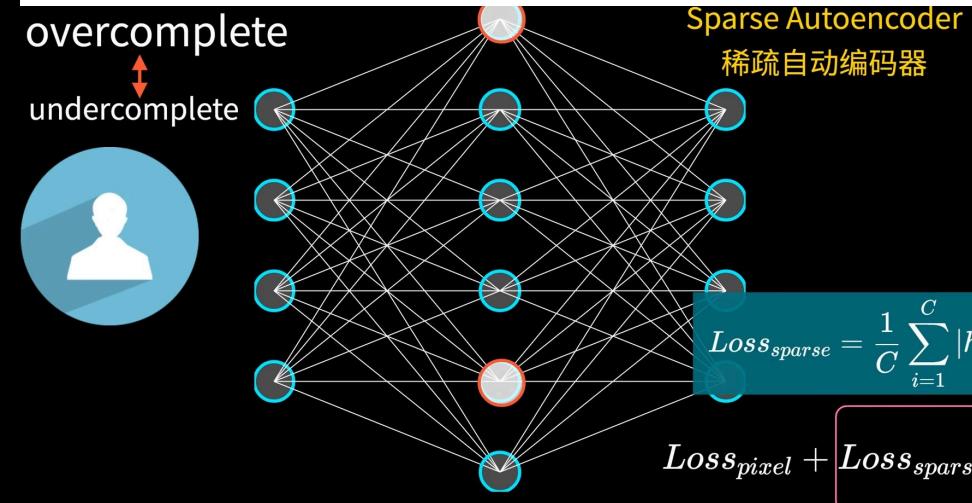
AE -autoencoder - 自编码器



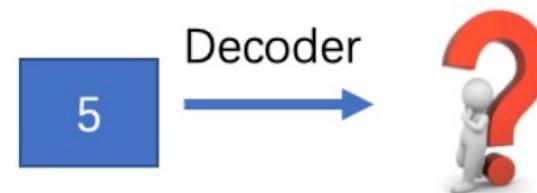
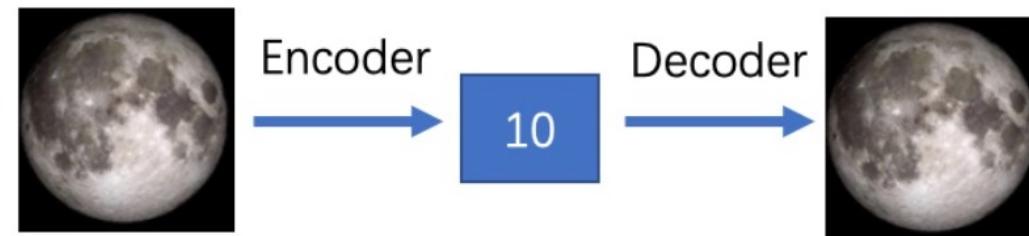
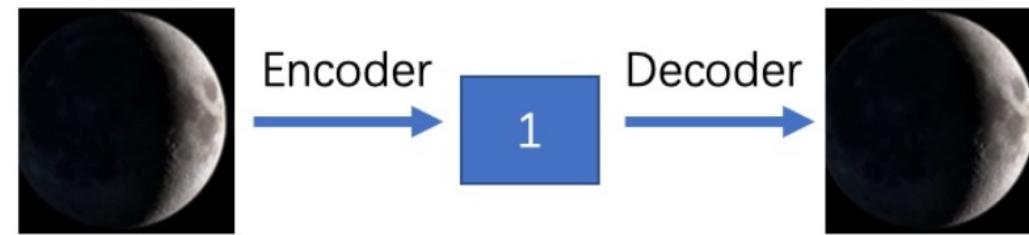
AE -autoencoder - 自编码器



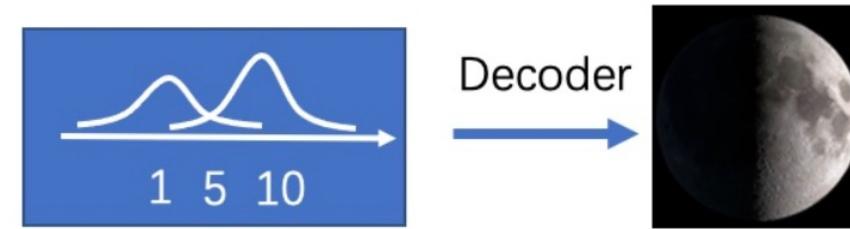
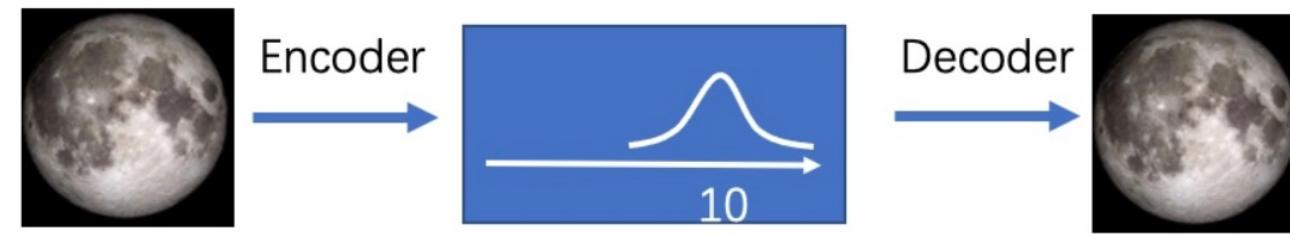
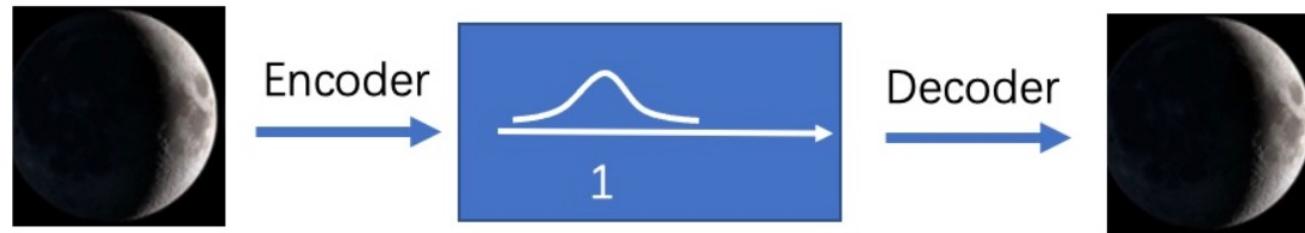
overcomplete  
↔  
undercomplete



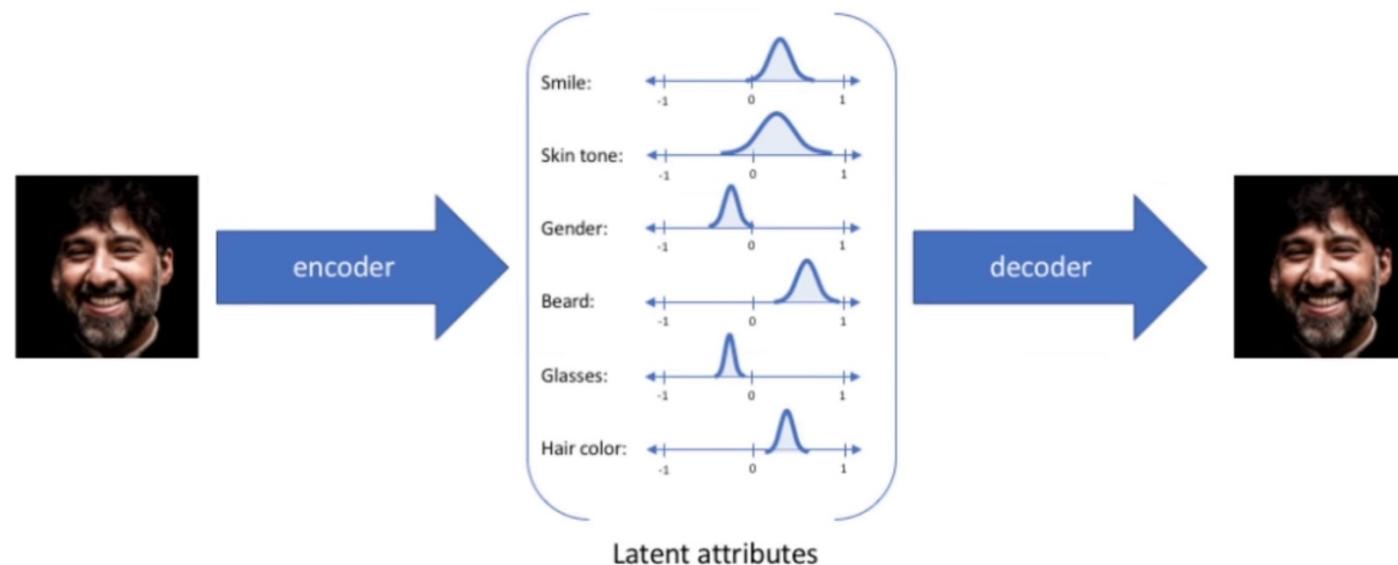
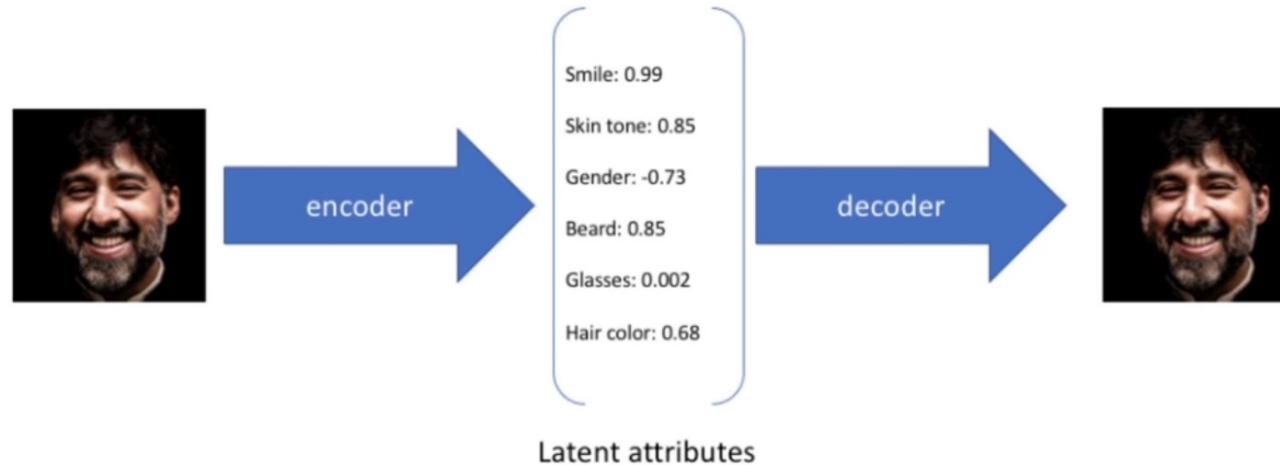
# VAE -Variational Autoencoder -变分自编码器



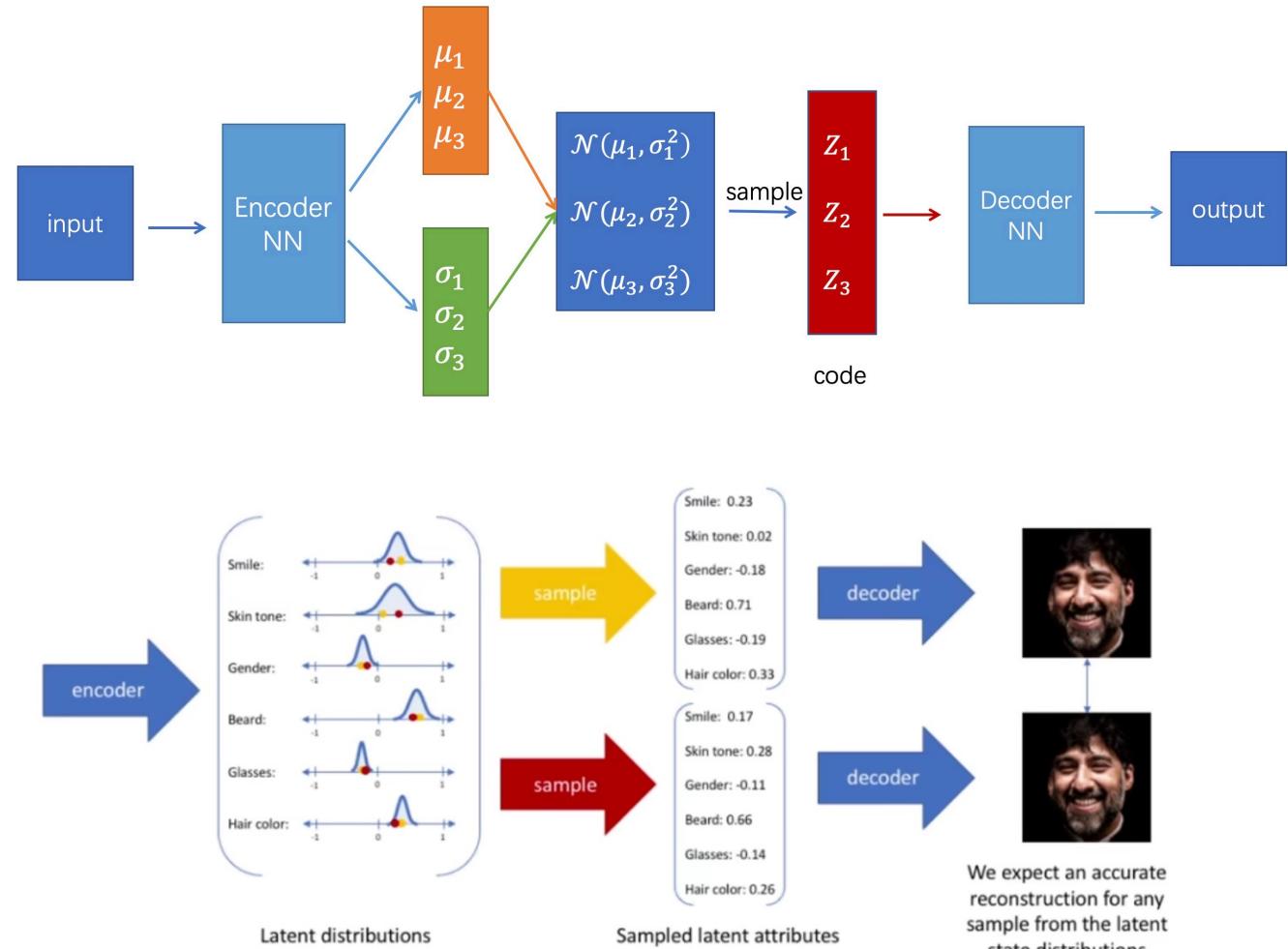
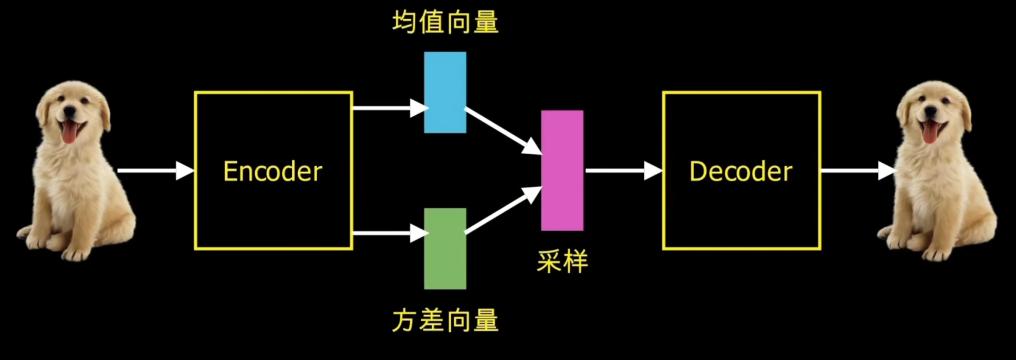
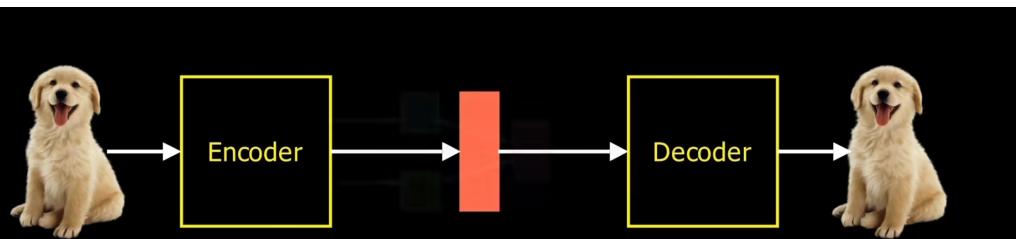
# VAE -Variational Autoencoder -变分自编码器



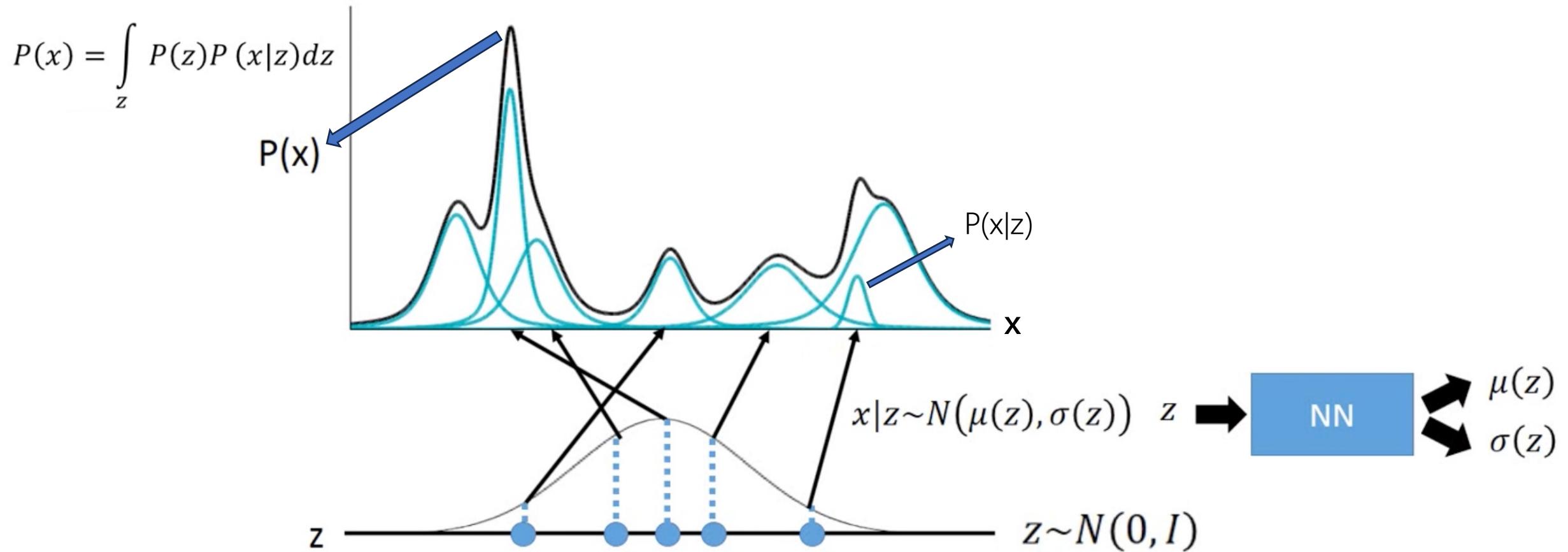
# VAE -Variational Autoencoder -变分自编码器



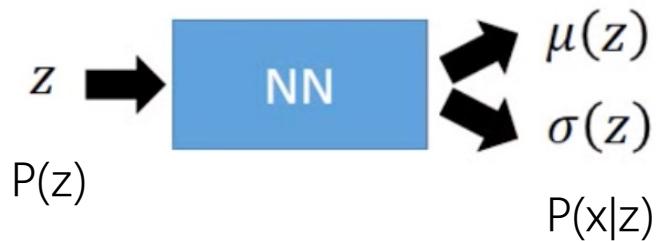
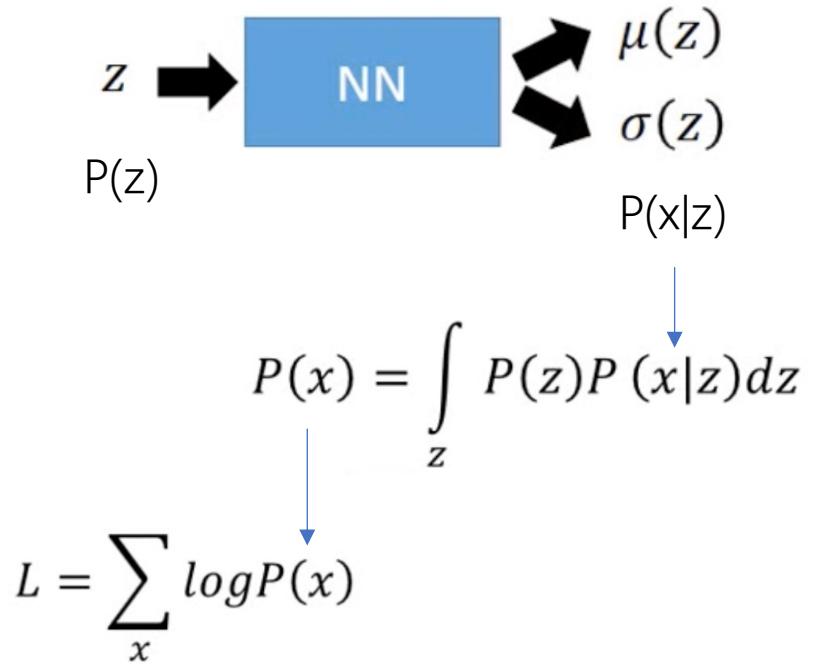
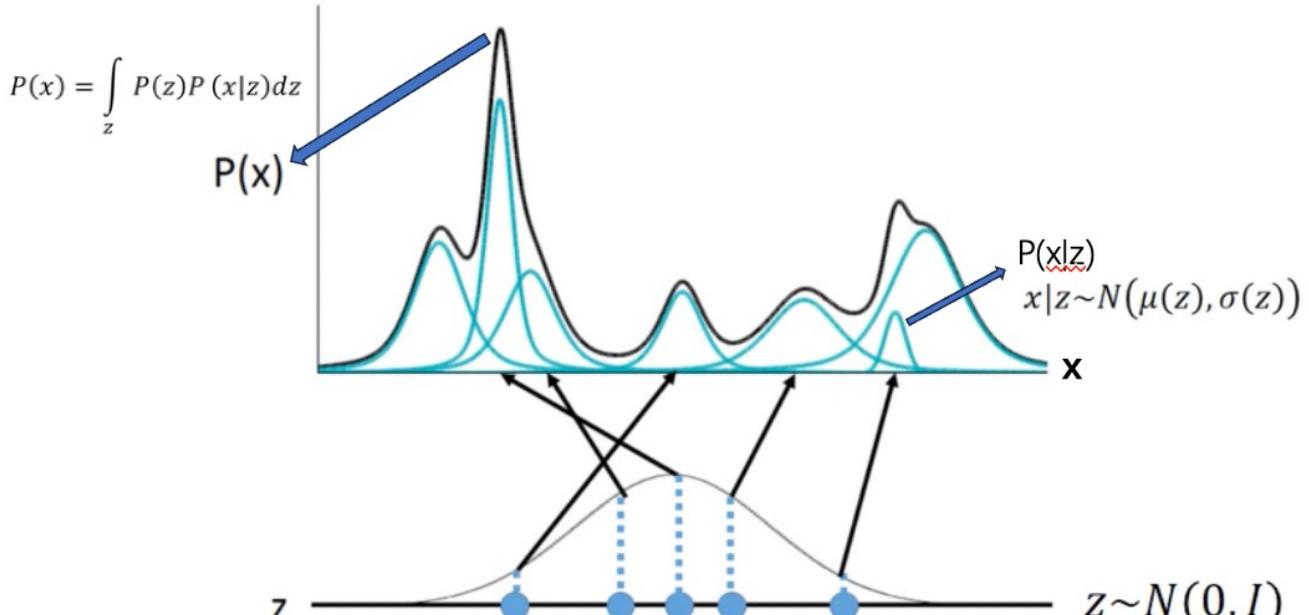
# VAE -Variational Autoencoder -变分自编码器



# VAE -Variational Autoencoder -变分自编码器



# VAE -Variational Autoencoder -变分自编码器



$$L = \sum_x \log P(x)$$

$$\log P(x) = \int_z q(z|x) \log P(x) dz \quad \text{q}(z|x) \text{ can be any distribution}$$

$$\int q(z|x) dz = 1$$

$$P(z, x) = P(z|x)P(x)$$

$$= \int_z q(z|x) \log \left( \frac{P(z,x)}{P(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x) \frac{P(z|x)}{P(z|x)}} \right) dz$$

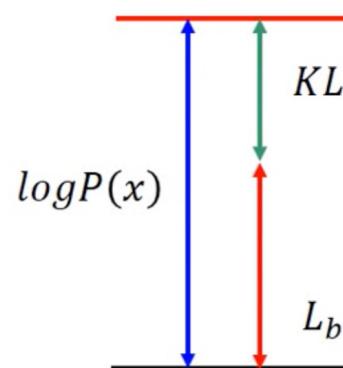
$$= \underbrace{\int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz}_{L_b} + \underbrace{\int_z q(z|x) \log \left( \frac{q(z|x)}{P(z|x)} \right) dz}_{KL(q(z|x)||P(z|x))} \geq \underbrace{\int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz}_{L_b}$$

$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

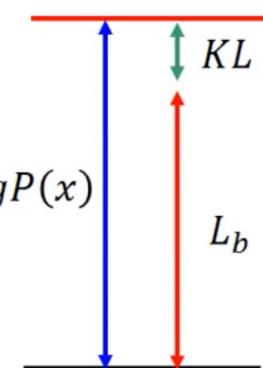
$$\downarrow$$

$$L_b = \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz$$

$$= \underbrace{\int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz}_{-KL(q(z|x)||P(z))} + \underbrace{\int_z q(z|x) \log P(x|z) dz}_{E_{q(z|x)}[\log P(x|z)]}$$



$\rightarrow$   
Maximize  $L_b$   
by  $q(z|x)$



$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

$$L_b = \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz$$

$$= \int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz + \int_z q(z|x) \log P(x|z) dz$$

$$-KL(q(z|x)||P(z))$$

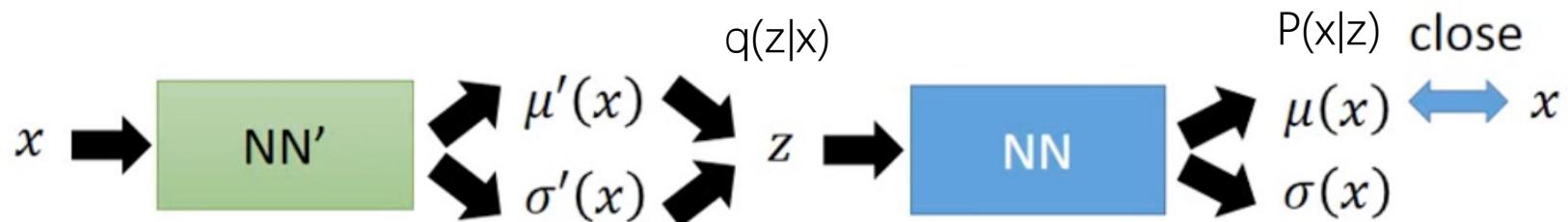
$$E_{q(z|x)}[\log P(x|z)]$$



$$p_\theta(x|z) = N(x; \mu_\theta(z), \sigma^2 I)$$

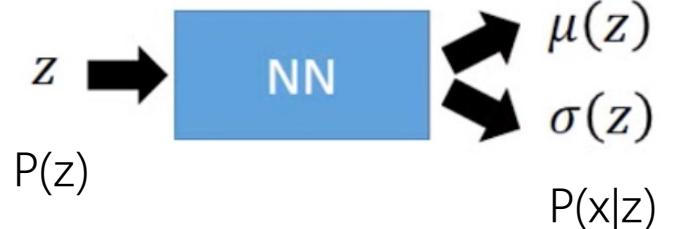
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\log p_\theta(x|z) = -c\|x - \mu_\theta(z)\|^2 + d$$



$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

$$\begin{aligned}
 L_b &= \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz \\
 &= \underbrace{\int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz}_{-KL(q(z|x)||P(z))} + \underbrace{\int_z q(z|x) \log P(x|z) dz}_{E_{q(z|x)}[\log P(x|z)]}
 \end{aligned}$$

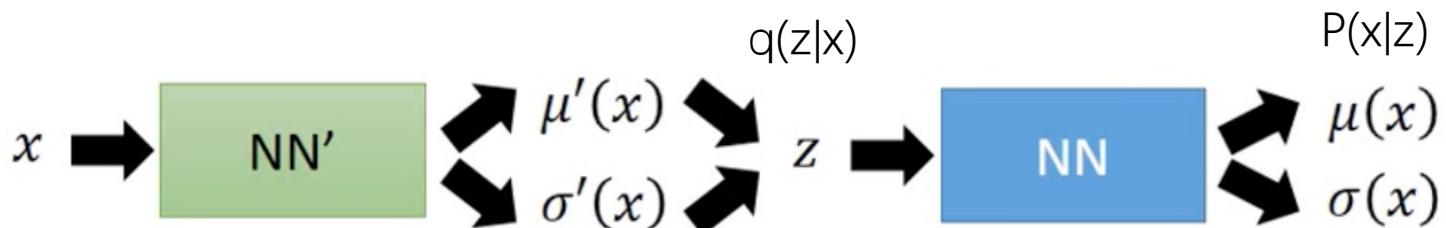


$$q_\varphi(z|x) = N(z; \mu_\varphi(x), \Sigma_\varphi(x))$$

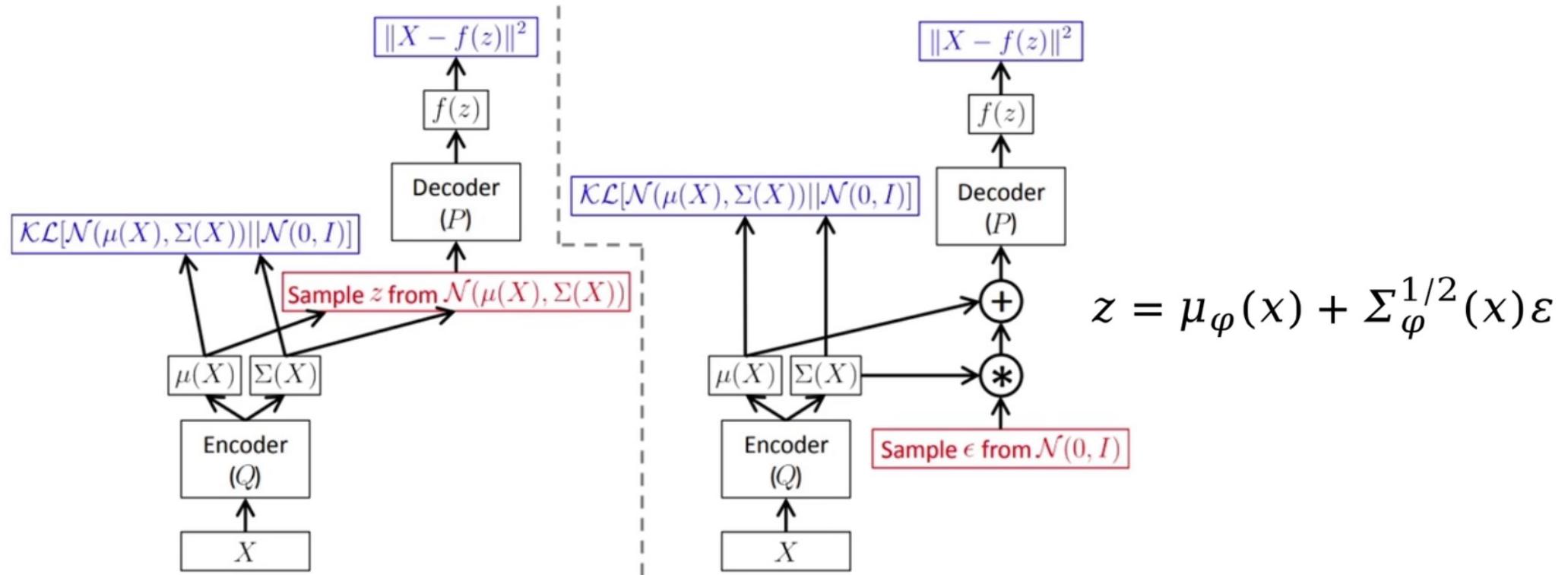
$$p(z) = N(z; 0, I)$$

$$D(N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)) = \frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - k + \log(\frac{\det \Sigma_1}{\det \Sigma_0}))$$

$$D(q_\varphi(z|x) \| p(z)) = \frac{1}{2} (\text{tr}(\Sigma_\varphi(x)) + \mu_\varphi(x)^T \mu_\varphi(x) - k - \log(\det \Sigma_\varphi(x))) \quad k \text{ is the dimension}$$

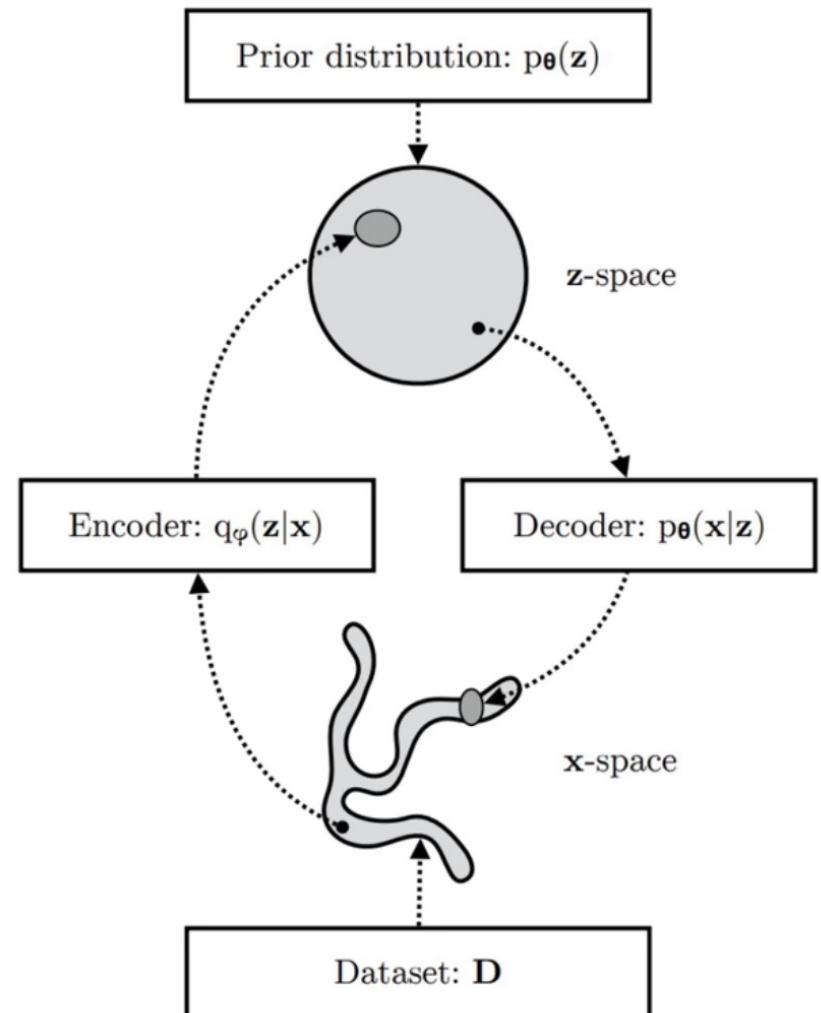
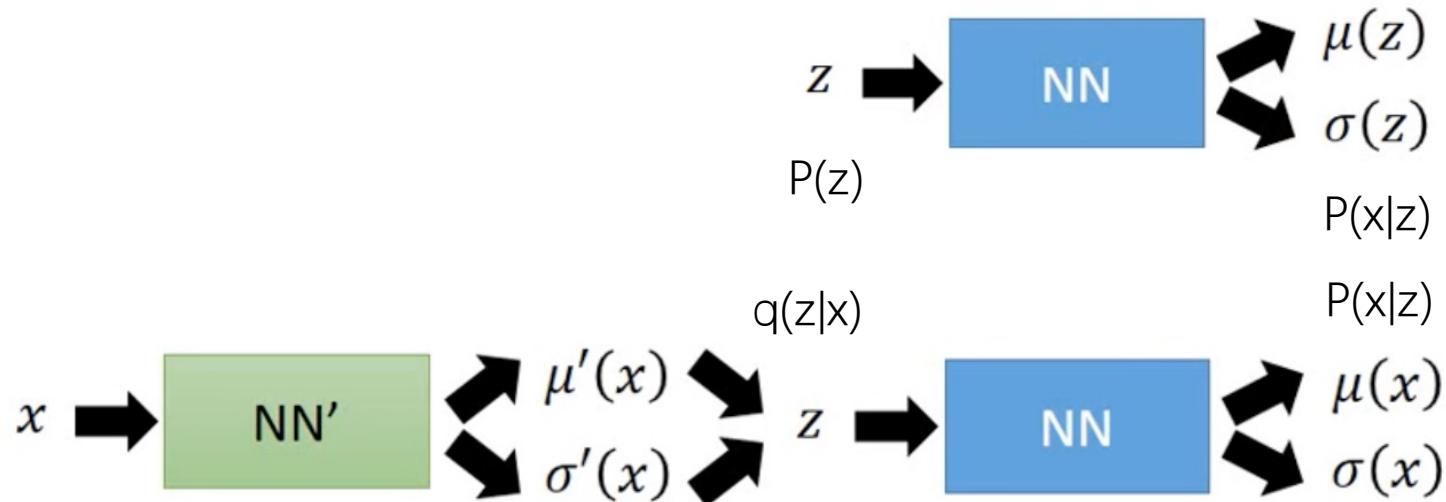


# VAE -Variational Autoencoder -变分自编码器



$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

$$\begin{aligned}
 L_b &= \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz \\
 &= \int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz + \int_z q(z|x) \log P(x|z) dz \\
 &\quad \underline{-KL(q(z|x)||P(z))} \quad \underline{E_{q(z|x)}[\log P(x|z)]} \\
 &= \underline{E_{q_\phi(z|x)}[\log(p_\theta(x|z))]} - \underline{D_{KL}(q_\phi(z|x)||p_\theta(z))} \\
 &\quad \text{Reconstruction Loss} \quad \quad \quad \text{Regularization Loss}
 \end{aligned}$$

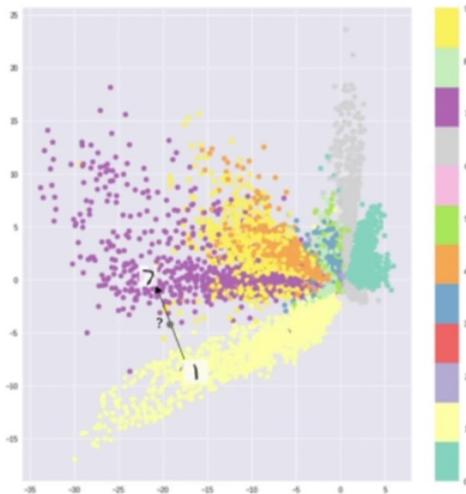


$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

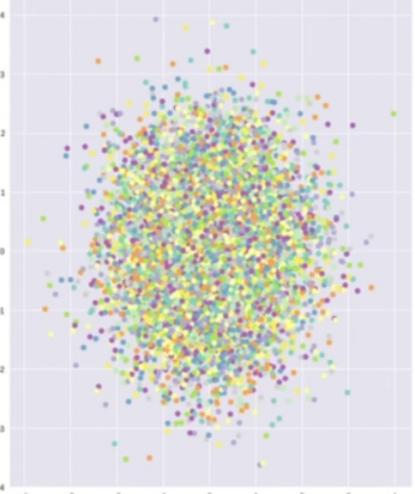
$$\begin{aligned}
 L_b &= \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz \\
 &= \int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz + \int_z q(z|x) \log P(x|z) dz \\
 &\quad \underline{-KL(q(z|x)||P(z))} \quad \underline{E_{q(z|x)}[\log P(x|z)]} \\
 &= \underline{E_{q_\phi(z|x)}[\log(p_\theta(x|z))]} - \underline{D_{KL}(q_\phi(z|x)||p_\theta(z))}
 \end{aligned}$$

Reconstruction Loss                      Regularization Loss

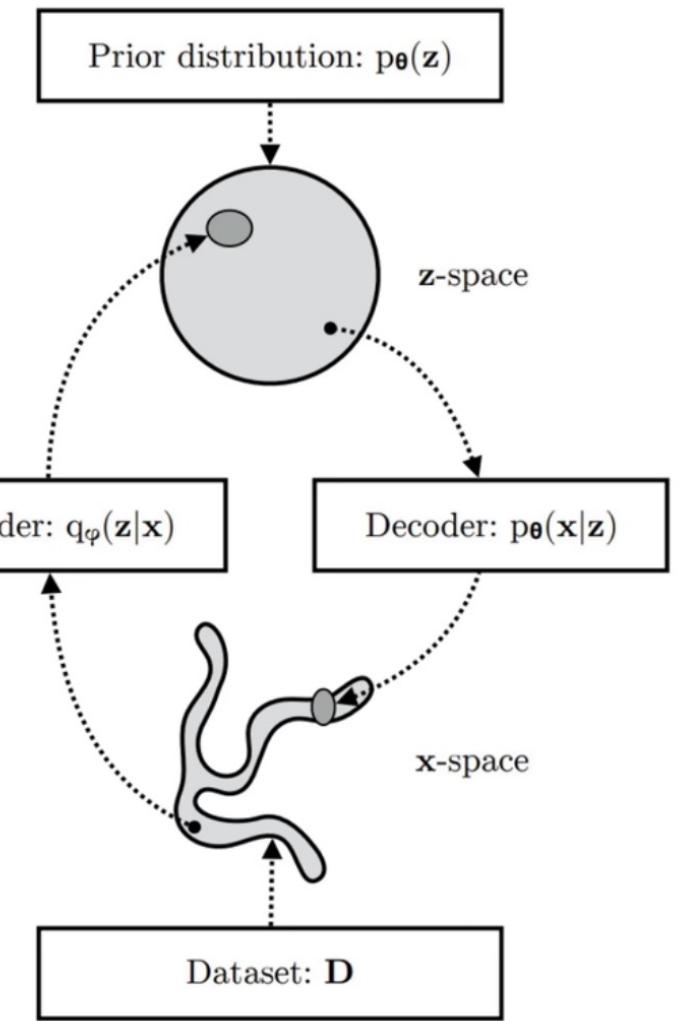
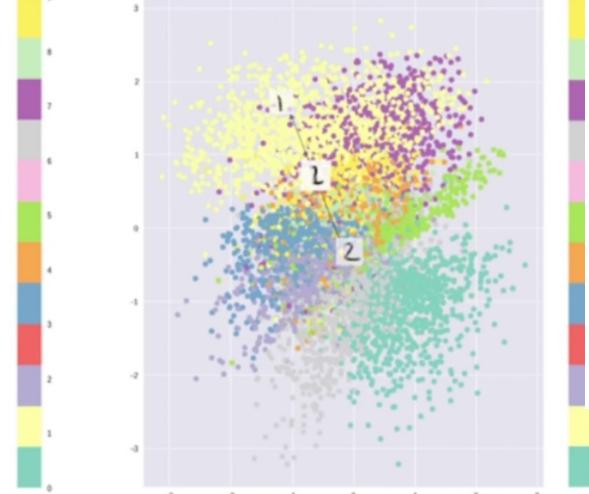
Only reconstruction loss



Only KL divergence



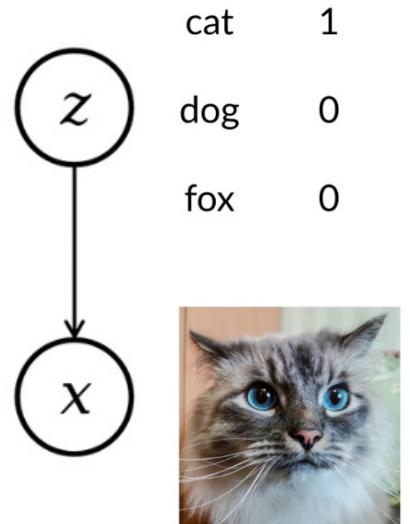
Combination



# VAE

-Variational Autoencoder -变分自编码器

## 变分推断



$$z \sim p(z) = \begin{cases} e^{-z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = e^{-z} I(z \geq 0)$$

$$x \sim p(x|z) = N(x, \mu = z, \sigma = 1) = \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)}$$

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)} = \frac{p(x|z)p(z)}{\int_z p(x, z) dz}$$

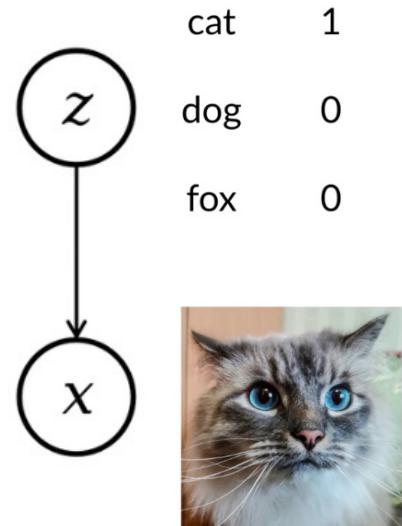
$$q_{\theta}(z) \approx p(z|x)$$

$$p(x) = \int_0^{\infty} p(x, z) dz = \underbrace{\int_0^{\infty} e^{-z} \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)} dz}_{\text{the integral has no closed-form solution}}$$

thus the posterior has  
no closed-form solution

the integral has no closed-form solution

# VAE -Variational Autoencoder -变分自编码器



$$z \sim p(z) = \begin{cases} e^{-z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = e^{-z} I(z \geq 0)$$

$$x \sim p(x|z) = N(x, \mu = z, \sigma = 1) = \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)}$$

$$q_\theta(z) \approx p(z|x)$$

$$\begin{aligned} D(q_\theta(z)||p(z|x)) &= E_{z \sim q}[\log \frac{q_\theta(z)}{p(z|x)}] = E_{z \sim q}[\log q_\theta(z) - \log p(z|x)] \\ &= E_{z \sim q}[\log q_\theta(z) - \log \frac{p(z, x)}{p(x)}] \\ &= E_{z \sim q}[\log q_\theta(z) - \log p(z, x)] + \log p(x) \end{aligned}$$

$$\log p(x) = E_{z \sim q}[\log p(z, x) - \log q_\theta(z)] + D(q_\theta(z)||p(z|x))$$

$$\log p(x) \geq E_{z \sim q}[\log p(z, x) - \log q_\theta(z)] \equiv \mathcal{L}_q$$

$$q_\theta(z) \approx p(z|x)$$

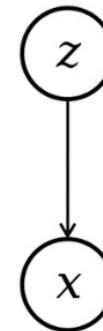
$$\begin{aligned} D(q_\theta(z)\|p(z|x)) &= E_{z\sim q}[\log \frac{q_\theta(z)}{p(z|x)}] = E_{z\sim q}[\log q_\theta(z) - \log p(z|x)] \\ &= E_{z\sim q}[\log q_\theta(z) - \log \frac{p(z,x)}{p(x)}] \\ &= E_{z\sim q}[\log q_\theta(z) - \log p(z,x)] + \log p(x) \end{aligned}$$

$$\log p(x) = E_{z\sim q}[\log p(z,x) - \log q_\theta(z)] + D(q_\theta(z)\|p(z|x))$$

$$\log p(x) \geq E_{z\sim q}[\log p(z,x) - \log q_\theta(z)] \equiv \mathcal{L}_q$$

$$q_\theta(z) \approx p(z|x)$$

$$q_\theta(z) = \begin{cases} \theta e^{-\theta z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = \theta e^{-\theta z} I(z \geq 0)$$



$$z \sim p(z) = \begin{cases} e^{-z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = e^{-z} I(z \geq 0)$$

$$x \sim p(x|z) = N(x, \mu = z, \sigma = 1) = \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)}$$

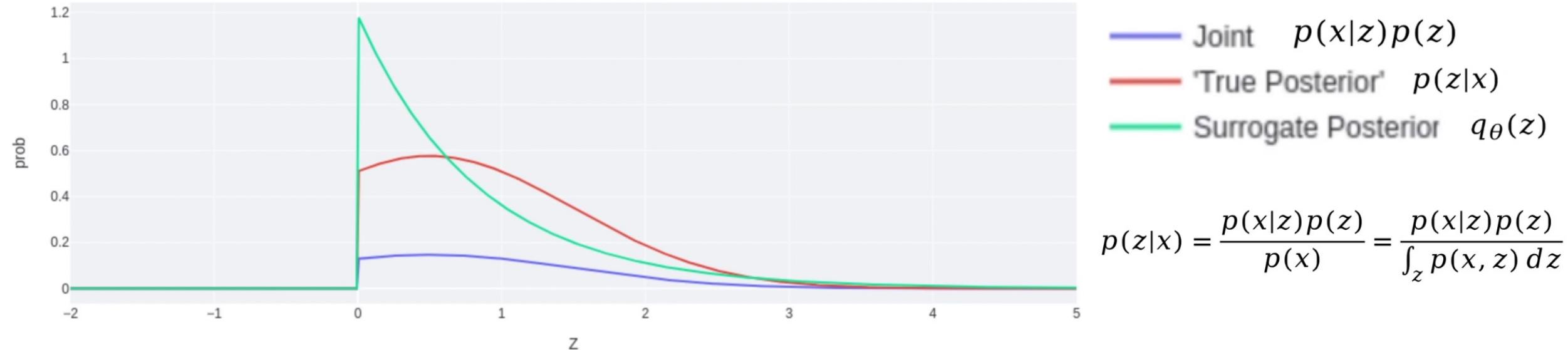
$$\mathcal{L}_q = E_{z\sim q}[\log p(z,x) - \log q_\theta(z)]$$

$$= E_{z\sim q}[\log \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)} e^{-z} I(z \geq 0) - \log \theta e^{-\theta z} I(z \geq 0)]$$

$$= E_{z\sim q}[-\frac{1}{2}z^2 + (x-1+\theta)z + C] = -\frac{1}{\theta^2} + \frac{x-1+\theta}{\theta} - \log \theta + C$$

$$\frac{\partial \mathcal{L}_q}{\partial \theta} = \frac{2}{\theta^3} - \frac{0.5}{\theta^2} - \frac{1}{\theta} = 0$$

↑  
expectation of exponential distribution  $E_\theta[z^n] = \frac{n!}{\theta^n}$   
 $x = 1.5$



$$p(z|x) = \frac{p(x|z)p(z)}{p(x)} = \frac{p(x|z)p(z)}{\int_z p(x, z) dz}$$

$$q_\theta(z) \approx p(z|x)$$

$$q_\theta(z) = \begin{cases} \theta e^{-\theta z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = \theta e^{-\theta z} I(z \geq 0)$$



$$z \sim p(z) = \begin{cases} e^{-z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = e^{-z} I(z \geq 0)$$

$$x \sim p(x|z) = N(x, \mu = z, \sigma = 1) = \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)}$$

$$\mathcal{L}_q = E_{z \sim q}[\log p(z, x) - \log q_\theta(z)]$$

$$= E_{z \sim q}[\log \frac{1}{\sqrt{2\pi}} e^{(-\frac{1}{2}(x-z)^2)} e^{-z} I(z \geq 0) - \log \theta e^{-\theta z} I(z \geq 0)]$$

$$= E_{z \sim q}[-\frac{1}{2}z^2 + (x-1+\theta)z + C] = -\frac{1}{\theta^2} + \frac{x-1+\theta}{\theta} - \log \theta + C$$

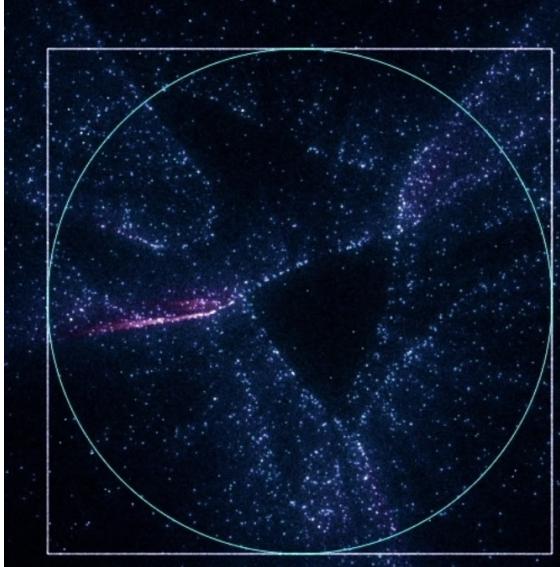
$$\frac{\partial \mathcal{L}_q}{\partial \theta} = \frac{2}{\theta^3} - \frac{0.5}{\theta^2} - \frac{1}{\theta} = 0$$

↑

expectation of exponential distribution  $E_\theta[z^n] = \frac{n!}{\theta^n}$

$x = 1.5$

# VAE -Variational Autoencoder -变分自编码器

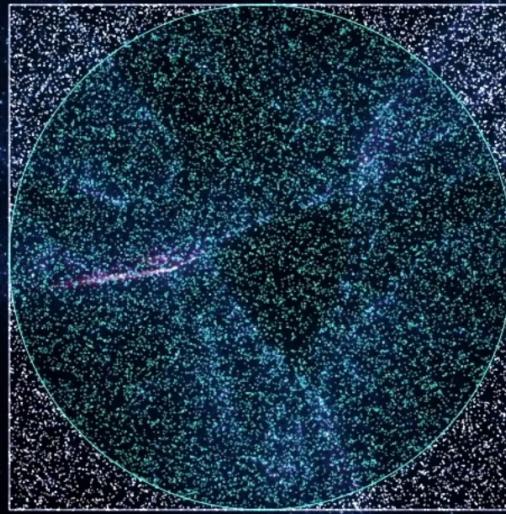


$$S_{\text{方}} = 4 \times r^2$$

$$S_{\text{圆}} = \pi \times r^2$$

T  
(总打点数)

N  
(圆中点数)



$$\pi = 4 \times \frac{N}{T} \begin{matrix} (\text{圆中点数}) \\ (\text{总打点数}) \end{matrix}$$

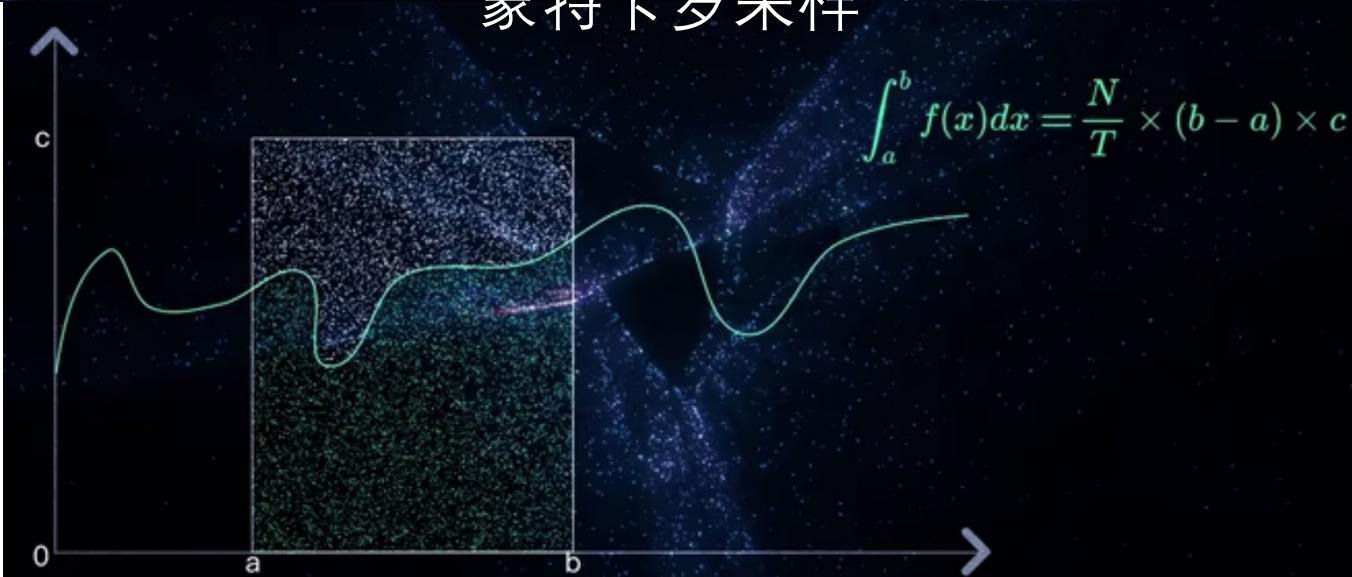
$$T=28, \quad \pi=2.857143$$

$$T=45, \quad \pi=3.111111$$

$$T=888, \quad \pi=3.144144$$

$$T=30507, \pi=3.141574$$

## 蒙特卡罗采样

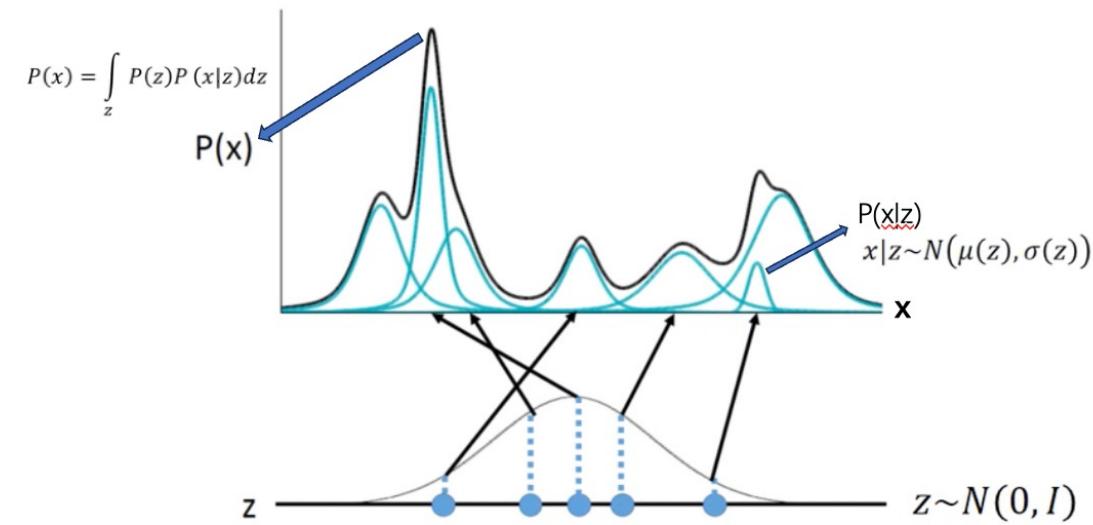


# VAE -Variational Autoencoder -变分自编码器



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

$$p_{\theta}(X | z) = \mathcal{N}(X | \mu'_i(z; \theta), \sigma'^2_i(z; \theta) * I)$$



$$p_{\theta}(X) = \int_z p_{\theta}(X | z)p(z)dz \approx \frac{1}{m} \sum_{j=1}^m p_{\theta}(X | z_j)$$

- 1) 从  $p(z)$  中多次采样  $z_1, z_2, \dots, z_m$
- 2) 根据  $p(x|z; \theta)$  计算  $x_1, x_2, \dots, x_m$
- 3) 求  $x$  的均值  $\frac{1}{m} \sum_{j=1}^m x_j$

# VAE -Variational Autoencoder -变分自编码器



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

$$p_{\theta}(X | z) = \mathcal{N}(X | \mu'_i(z; \theta), \sigma'^2_i(z; \theta) * I)$$

$$p_{\theta}(X) = \int_z p_{\theta}(X | z)p(z)dz$$

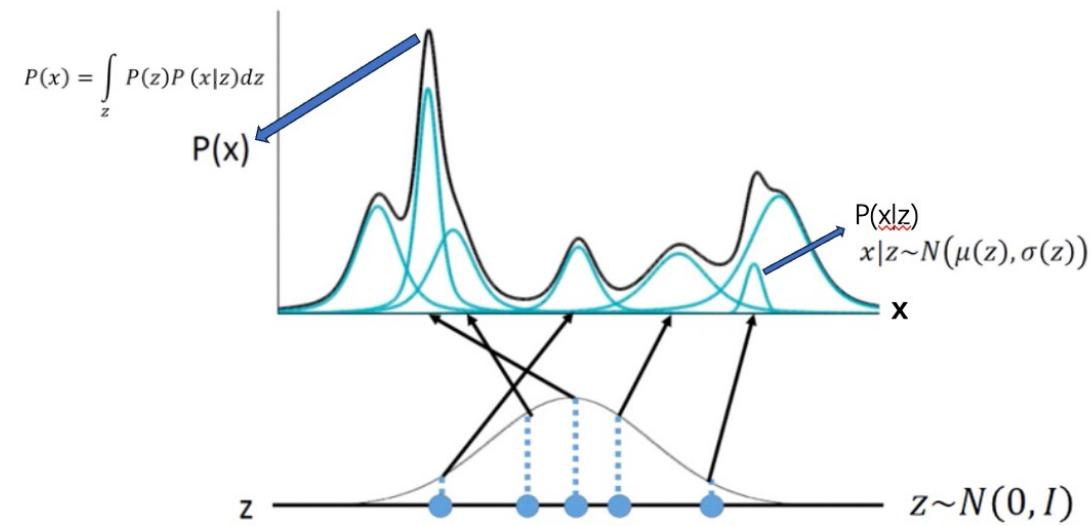
$$\approx \frac{1}{m} \sum_{j=1}^m p_{\theta}(X | z_j).$$



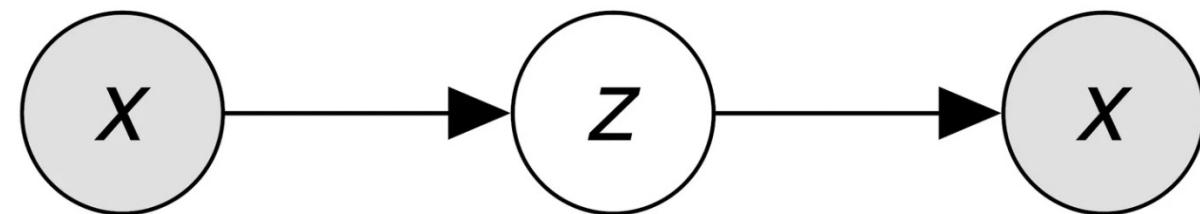
(a)



(b)



$$q_{\theta}(z | X) \rightarrow p_{\theta}(z | X) \quad p_{\theta}(X | z)$$



$$\begin{aligned} p_{\theta}(z | x_i) &= \frac{p_{\theta}(x_i | z)p(z)}{p_{\theta}(x_i)} \\ &= \frac{p_{\theta}(x_i | z)p(z)}{\int_{\hat{z}} p_{\theta}(x_i | \hat{z})p(\hat{z})d\hat{z}} \end{aligned}$$

# VAE -Variational Autoencoder -变分自编码器

$$L = \sum_x \log P(x)$$

$$\log P(x) = \int_z q(z|x) \log P(x) dz \quad q(z|x) \text{ can be any distribution}$$

$$= \int_z q(z|x) \log \left( \frac{P(z,x)}{P(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x) P(z|x)} \right) dz$$

$$= \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz + \int_z q(z|x) \log \left( \frac{q(z|x)}{P(z|x)} \right) dz \geq \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz$$

$L_b$

$KL(q(z|x)||P(z|x))$

$L_b$

$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

$$\downarrow \\ L_b = \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz$$

$$= \int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz + \int_z q(z|x) \log P(x|z) dz$$

$-KL(q(z|x)||P(z))$

$E_{q(z|x)}[\log P(x|z)]$

$$KL(q_\theta(z|x)||p(z|x))$$

$$= \int q_\theta(z|x) \ln \frac{q_\theta(z|x)}{p(z|x)} dz$$

$$= \mathbb{E}_{z \sim q_\theta(z|x)} [\ln \frac{q_\theta(z|x)}{p(z|x)}]$$

$$= \mathbb{E}_{z \sim q_\theta(z|x)} [\ln q_\theta(z|x) - \ln p(z|x)]$$

$$= \mathbb{E}_{z \sim q_\theta(z|x)} [\ln q_\theta(z|x) - \ln \frac{p(x|z)p(z)}{p(x)}]$$

$$= \mathbb{E}_{z \sim q_\theta(z|x)} [\ln q_\theta(z|x) - \ln p(z) - \ln p(x|z)] + \ln p(x)$$

$$= KL(q_\theta(z|x)||p(z)) - \mathbb{E}_{z \sim q_\theta(z|x)} [\ln p(x|z)] + \ln p(x)$$

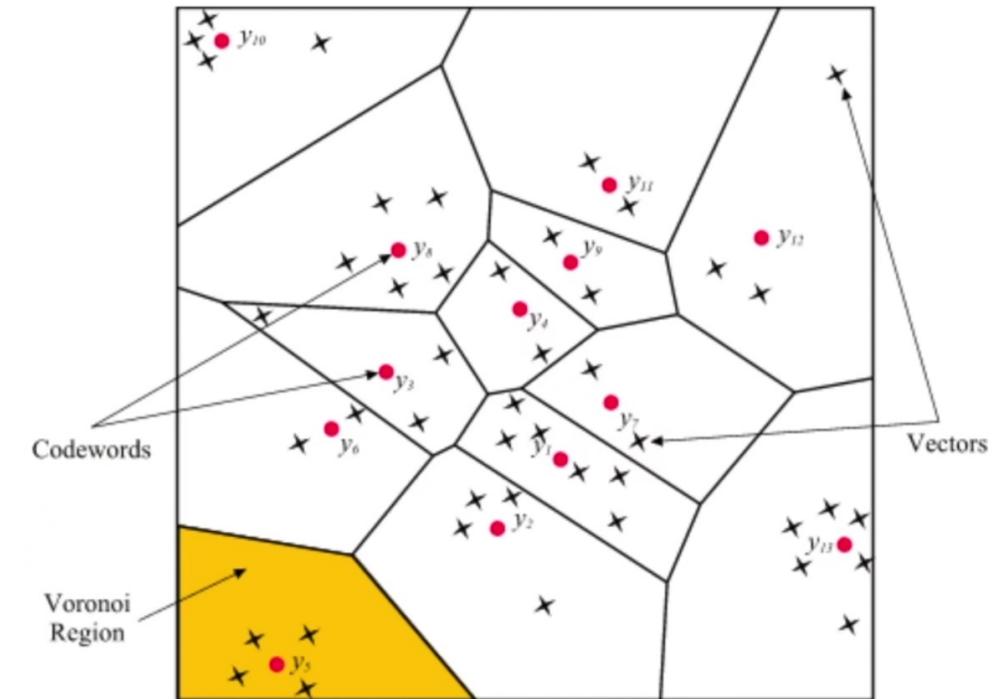
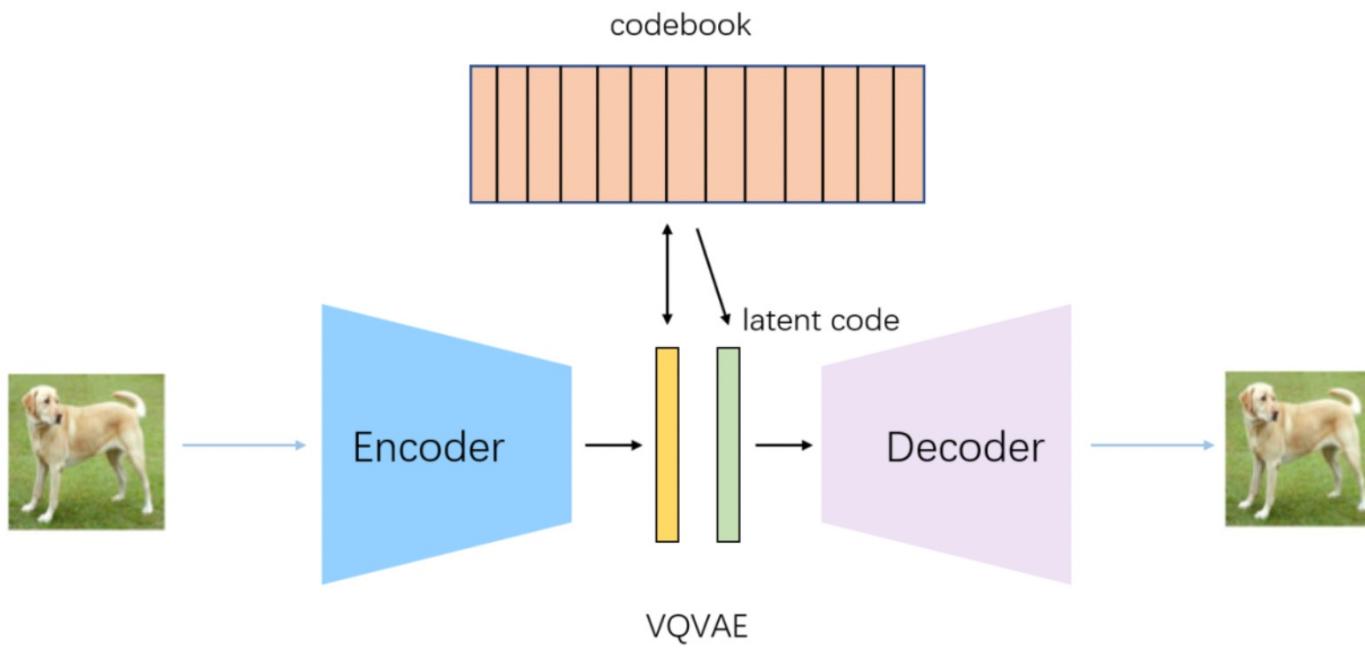
$$\ln p(x) - KL(q_\theta(z|x)||p(z|x)) = \mathbb{E}_{z \sim q_\theta(z|x)} [\ln p(x|z)] - KL(q_\theta(z|x)||p(z))$$

# VQ-VAE

-Vector Quantised Variational AutoEncoder –向量量化变分自动编码器

## Posterior Collapse

Have a powerful decoder, but latents are ignored.

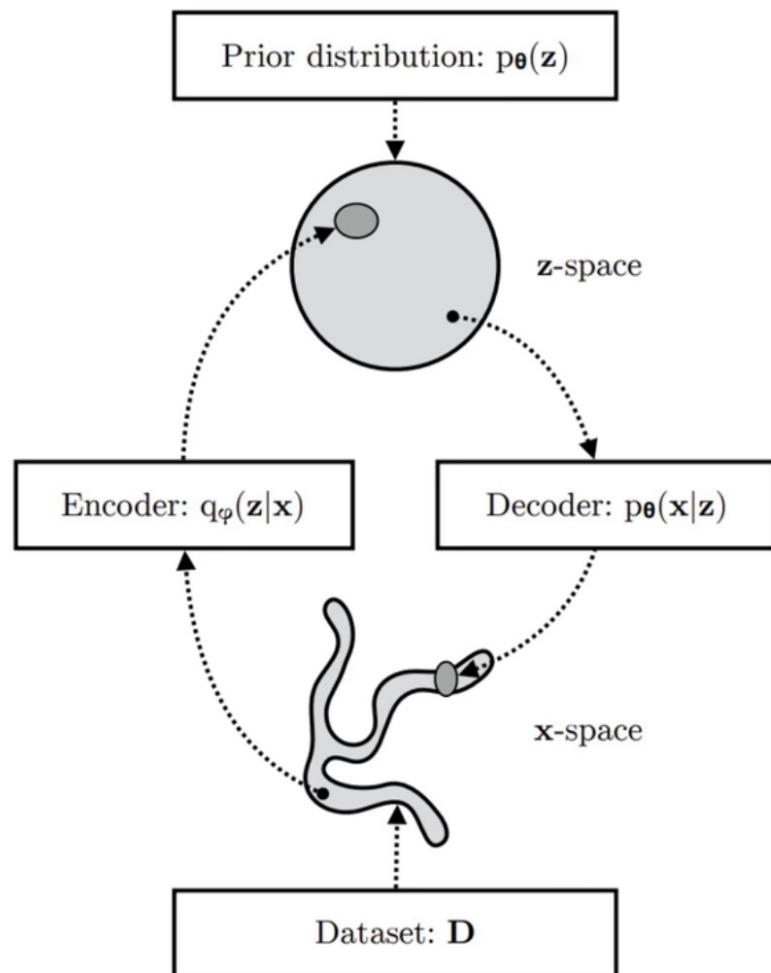


$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

$$\begin{aligned}
 L_b &= \int_z q(z|x) \log \left( \frac{P(z,x)}{q(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(x|z)P(z)}{q(z|x)} \right) dz \\
 &= \underbrace{\int_z q(z|x) \log \left( \frac{P(z)}{q(z|x)} \right) dz}_{-KL(q(z|x)||P(z))} + \underbrace{\int_z q(z|x) \log P(x|z) dz}_{E_{q(z|x)}[\log P(x|z)]} \\
 &= \underbrace{E_{q_\phi(z|x)}[\log(p_\theta(x|z))]}_{\text{Reconstruction Loss}} - \underbrace{D_{KL}(q_\phi(z|x)||p_\theta(z))}_{\text{Regularization Loss}}
 \end{aligned}$$

$$q_\phi(z|x) \simeq q_\phi(z) = \mathcal{N}(a, b)$$

$$\mathcal{L} = \mathcal{L}_1 + \alpha \mathcal{L}_2$$



$$\log P(x) = L_b + KL(q(z|x)||P(z|x))$$

$$= \underline{E_{q_\phi(z|x)}[\log(p_\theta(x|z))]} - \underline{D_{KL}(q_\phi(z|x)||p_\theta(z))}$$

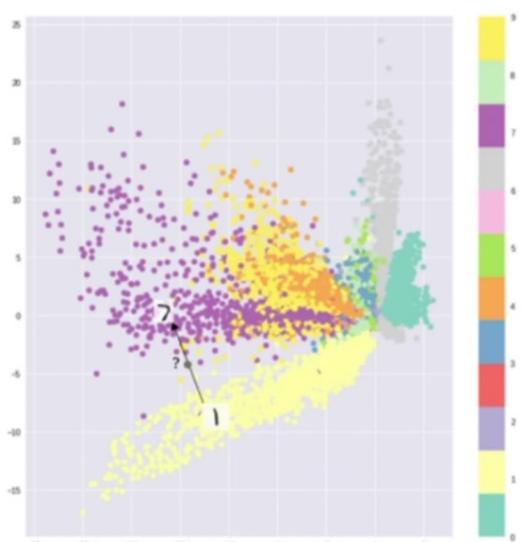
Reconstruction Loss

Regularization Loss

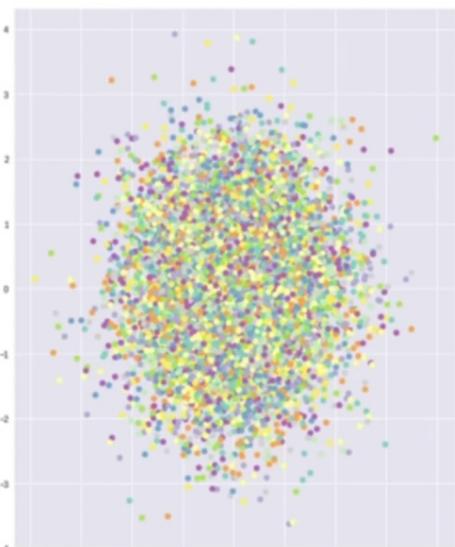
$$\mathcal{L} = \mathcal{L}_1 + \alpha \mathcal{L}_2$$

$$q_\phi(z|x) \simeq q_\phi(z) = \mathcal{N}(a, b)$$

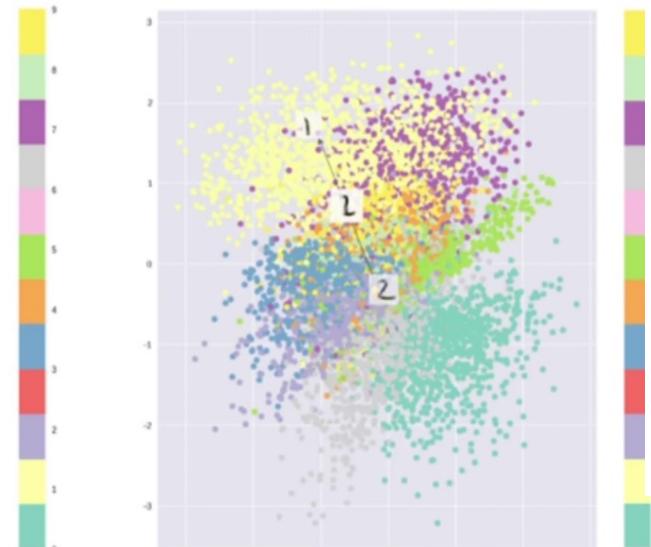
Only reconstruction loss



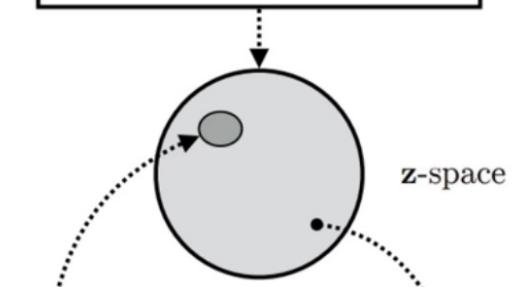
Only KL divergence



Combination



Prior distribution:  $p_\theta(z)$



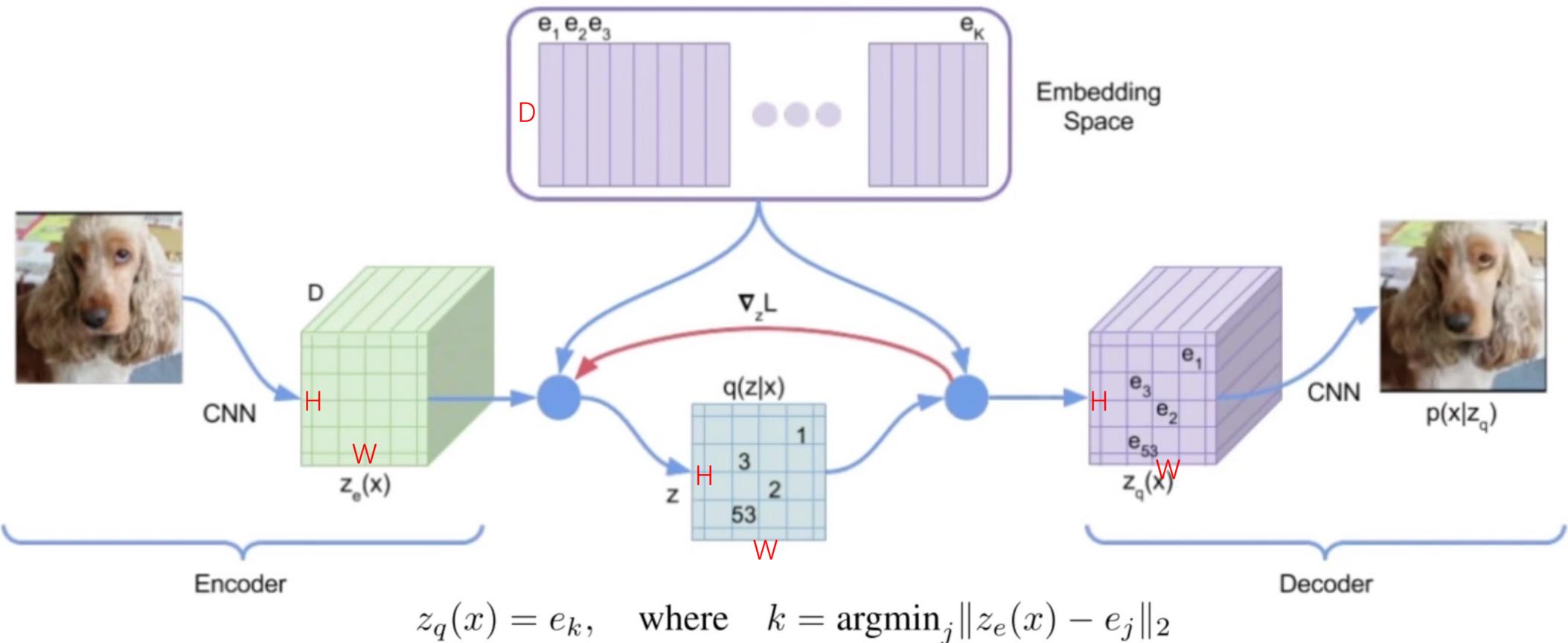
Encoder:  $q_\phi(z|x)$

Decoder:  $p_\theta(x|z)$

Dataset:  $\mathbf{D}$

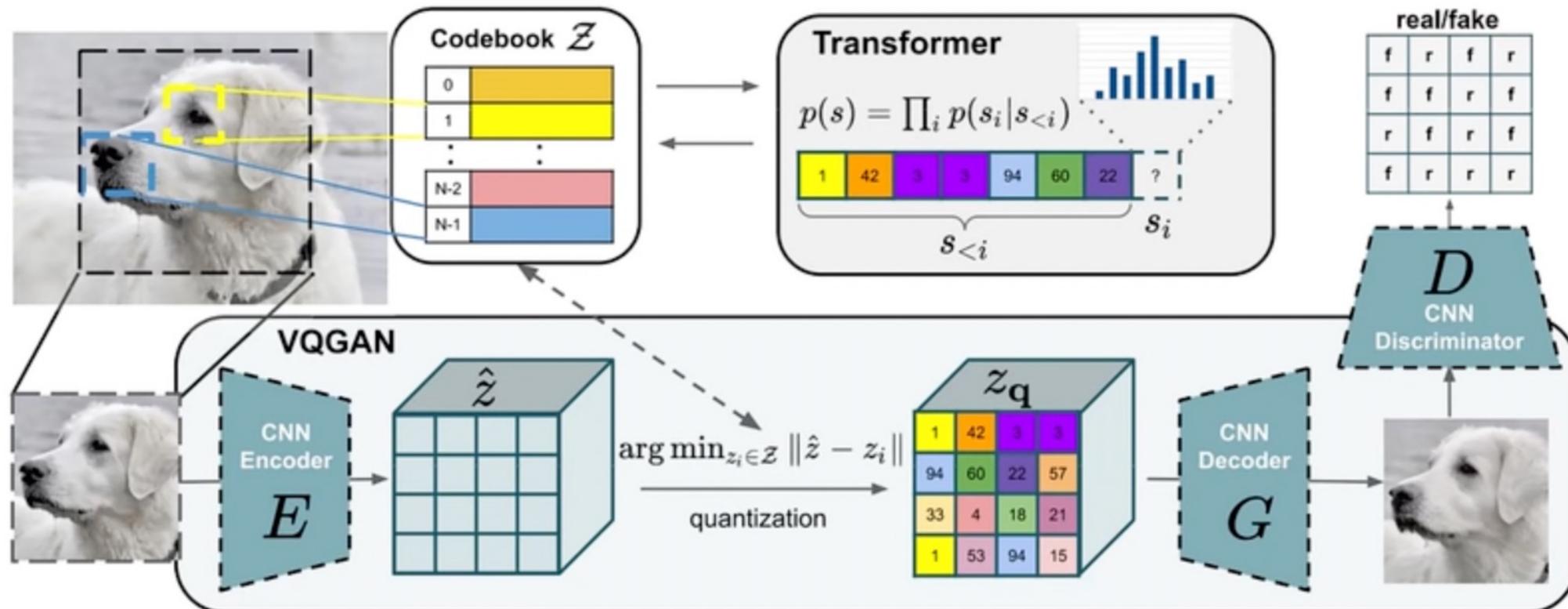
# VQ-VAE

-Vector Quantised Variational AutoEncoder – 向量量化变分自动编码器

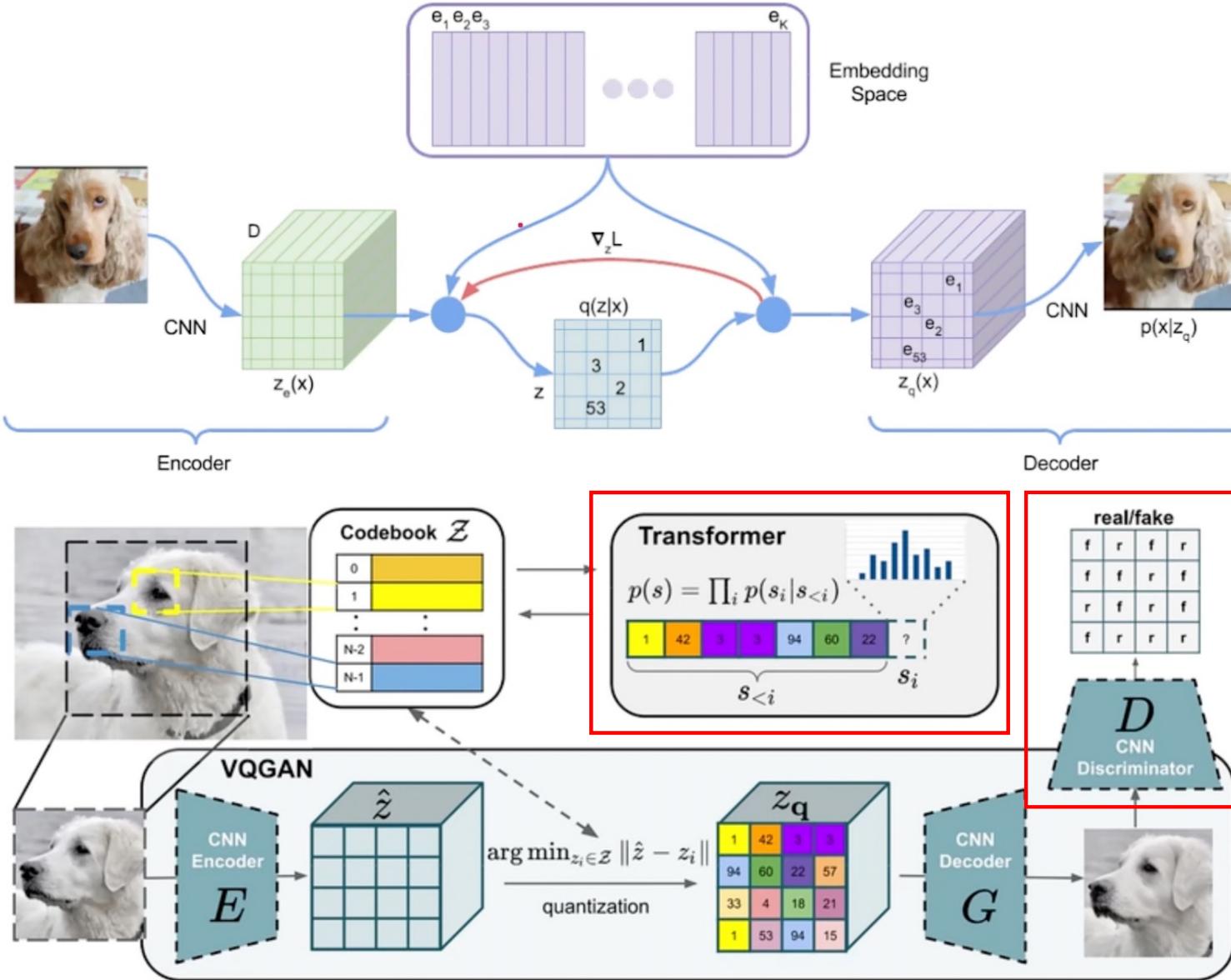


$$L = \underbrace{\log p(x|z_q(x))}_{\text{reconstruction loss}} + \underbrace{\|\operatorname{sg}[z_e(x)] - e\|_2^2}_{\text{VQ loss}} + \underbrace{\beta \|z_e(x) - \operatorname{sg}[e]\|_2^2}_{\text{commitment loss}}$$

# VQ-GAN -Vector Quantized Generative Adversarial Network -向量量化生成对抗网络



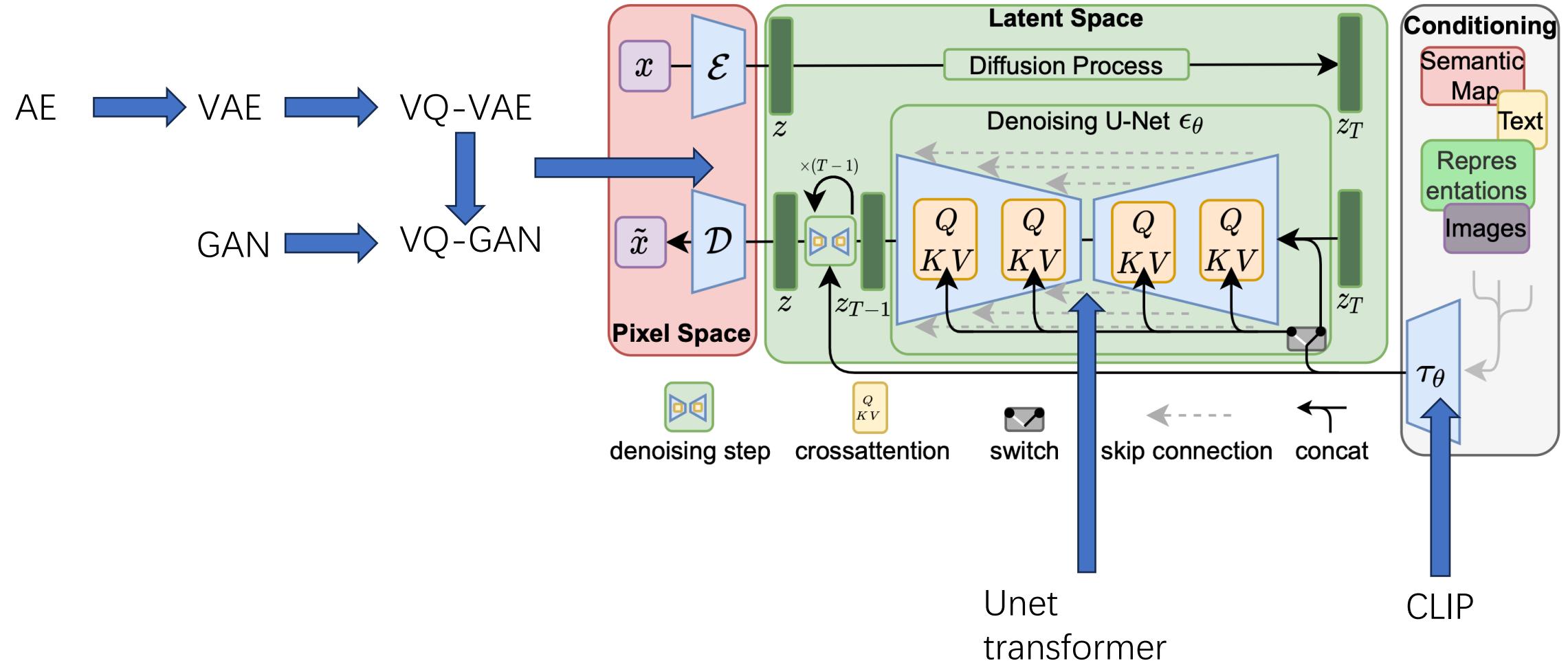
# VQ-GAN -Vector Quantized Generative Adversarial Network -向量量化生成对抗网络



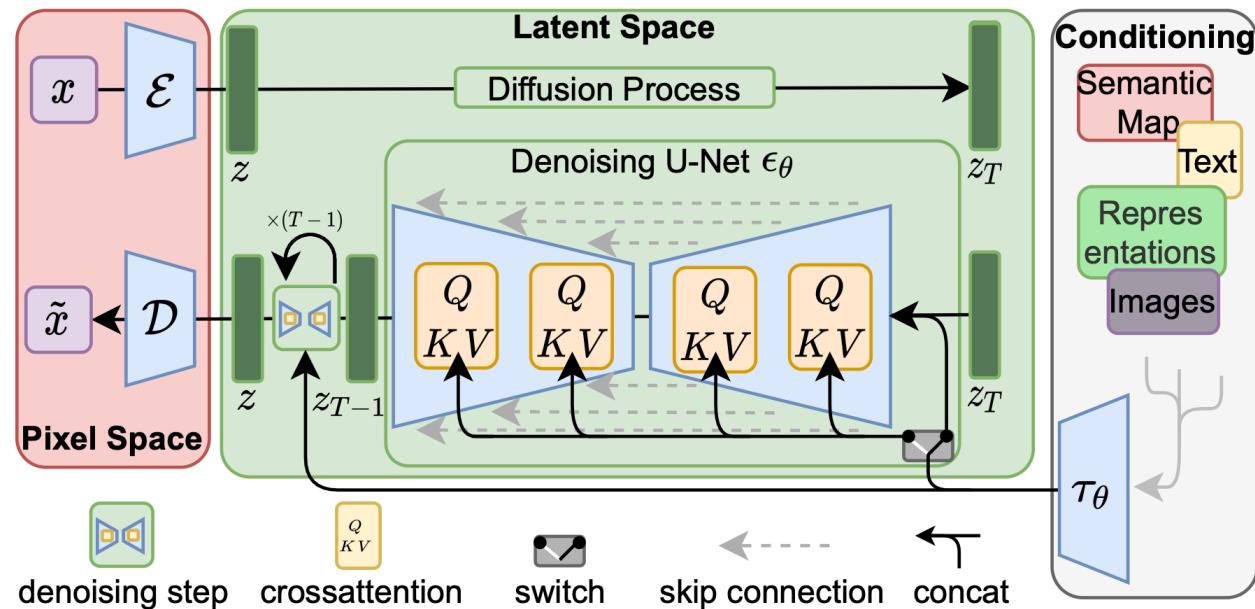
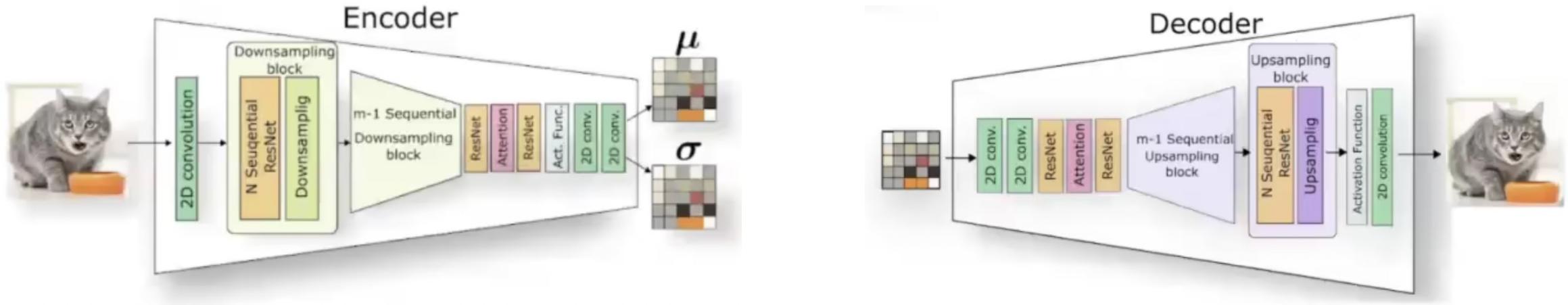
$$\mathcal{L}_{\text{VQ}}(E, G, \mathcal{Z}) = \|x - \hat{x}\|^2 + \|\text{sg}[E(x)] - z_q\|_2^2 + \|\text{sg}[z_q] - E(x)\|_2^2.$$

$$\mathcal{Q}^* = \arg \min_{E, G, \mathcal{Z}} \max_D \mathbb{E}_{x \sim p(x)} [\mathcal{L}_{\text{VQ}}(E, G, \mathcal{Z}) + \lambda \mathcal{L}_{\text{GAN}}(\{E, G, \mathcal{Z}\}, D)]$$

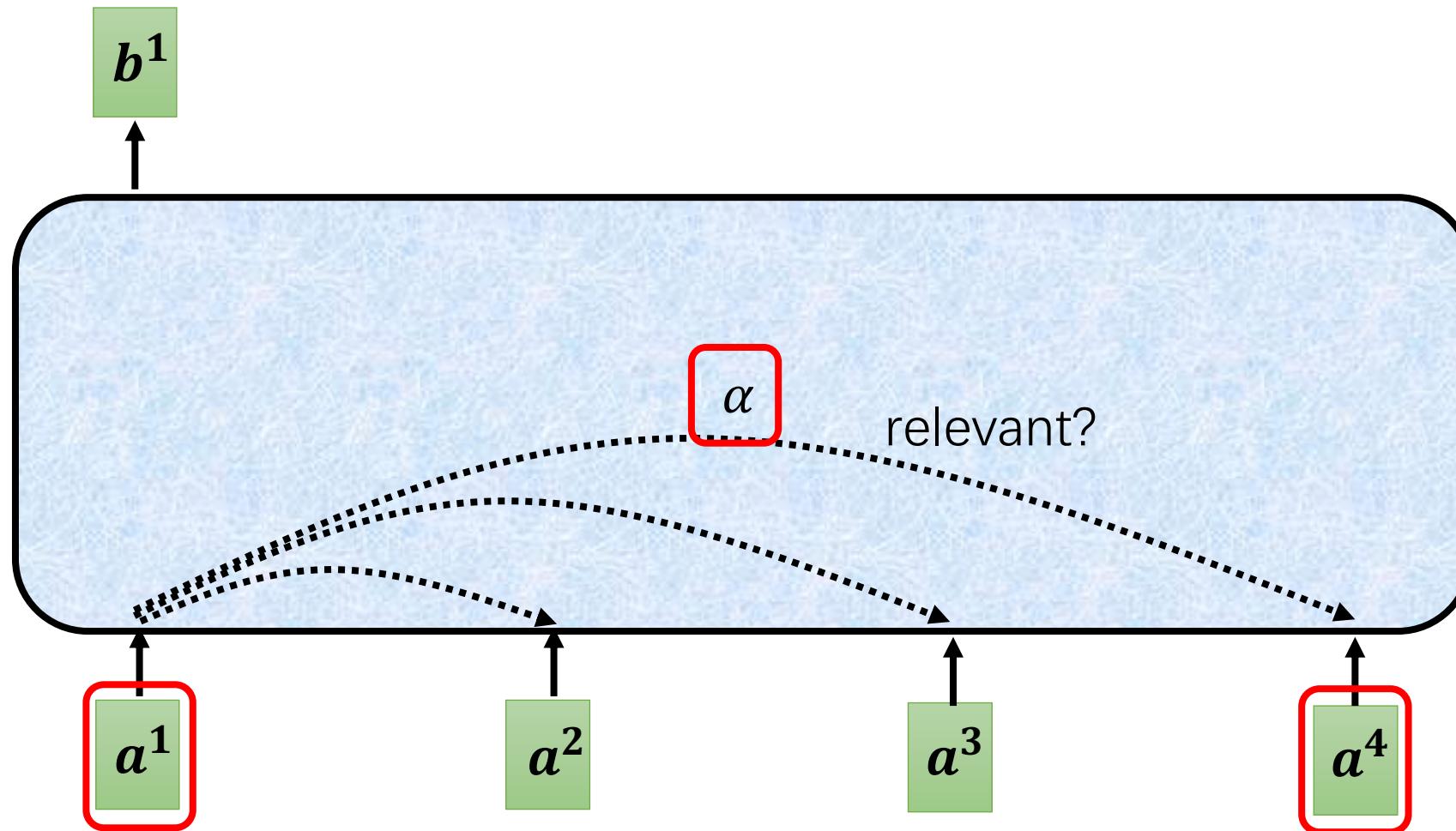
# LDM -Latent Diffusion Models -潜在扩散模型



# LDM -Latent Diffusion Models -潜在扩散模型

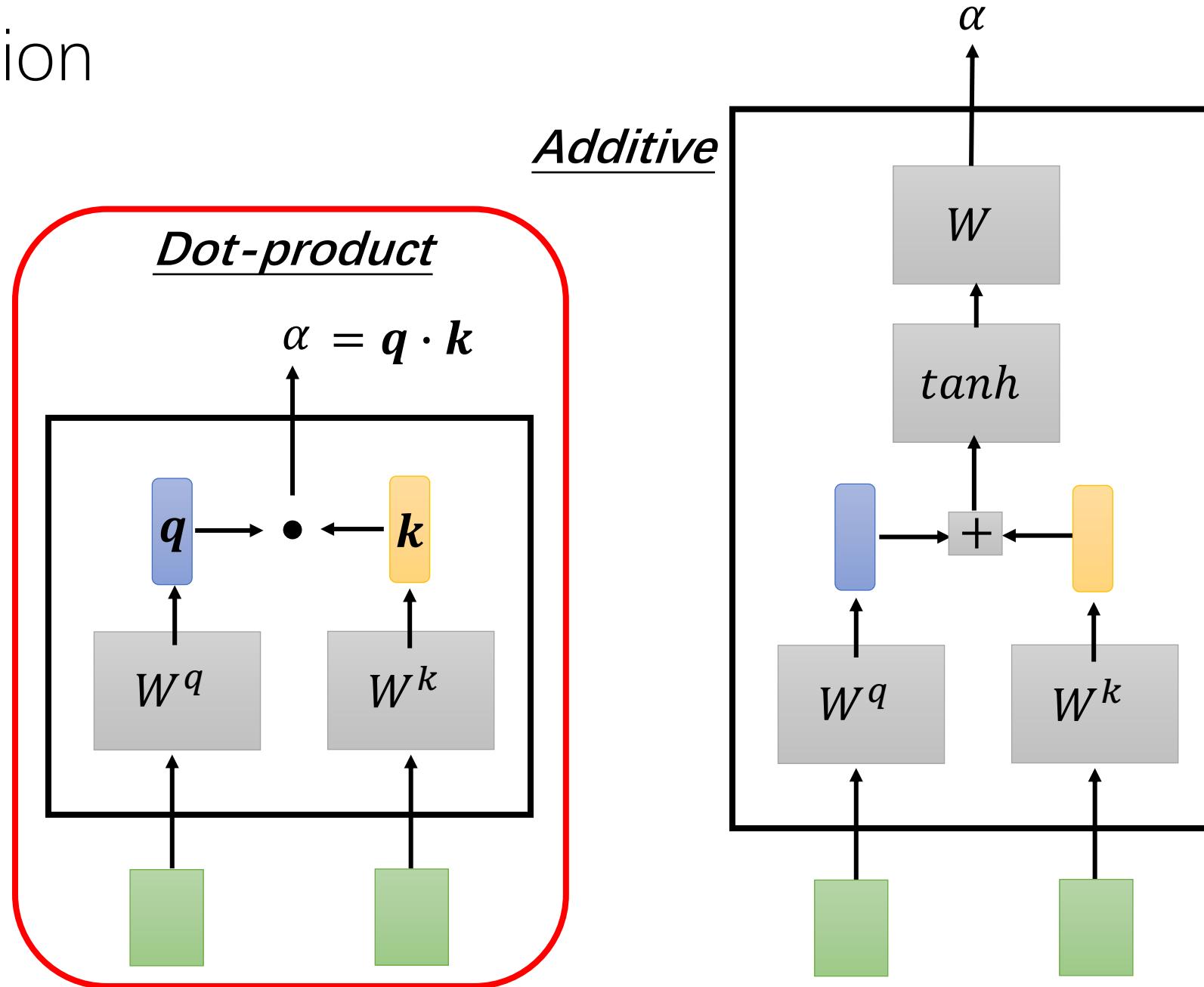


# Self-attention



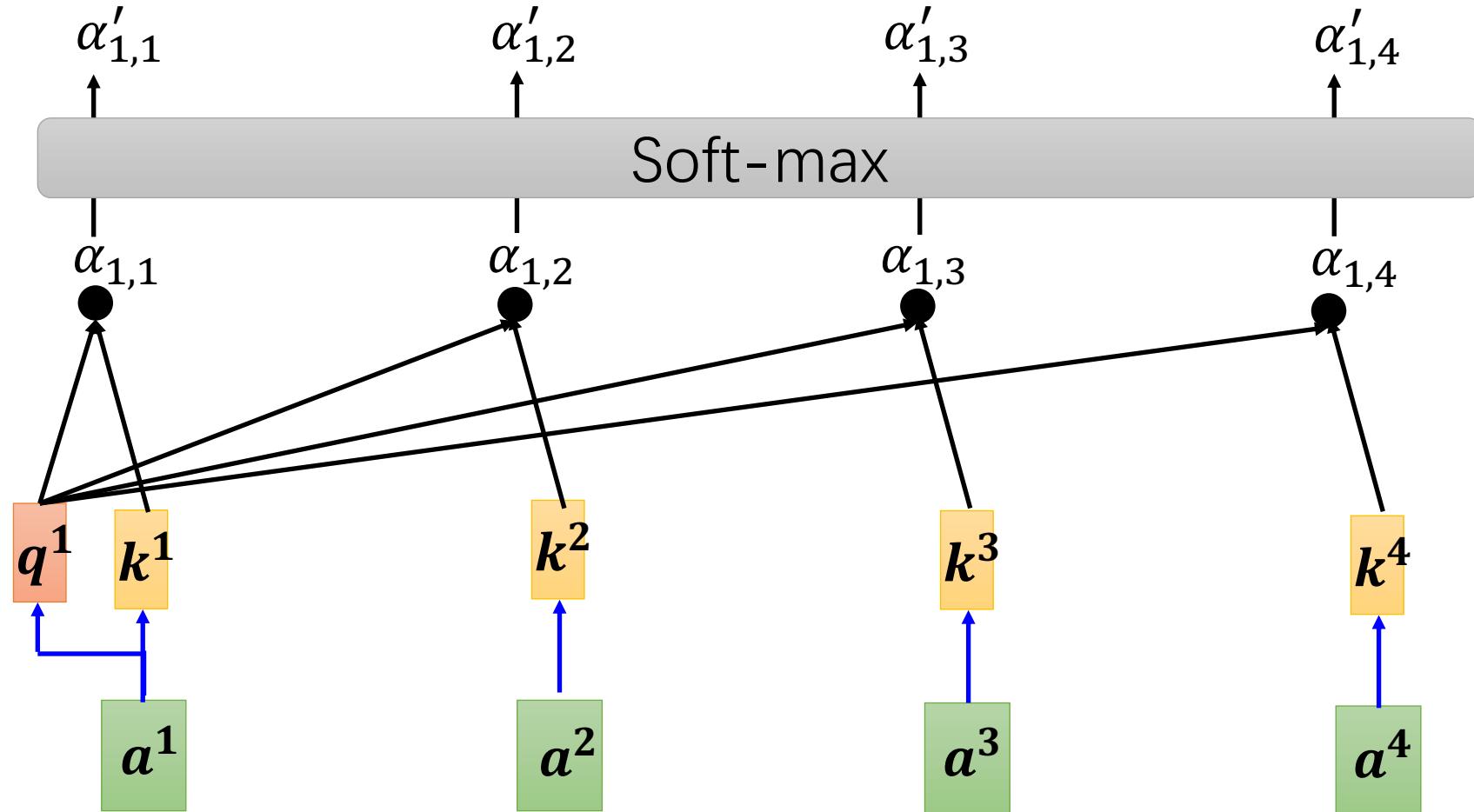
Find the relevant vectors in a sequence

# Self-attention



## Self-attention

$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



$$q^1 = W^q a^1$$

$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

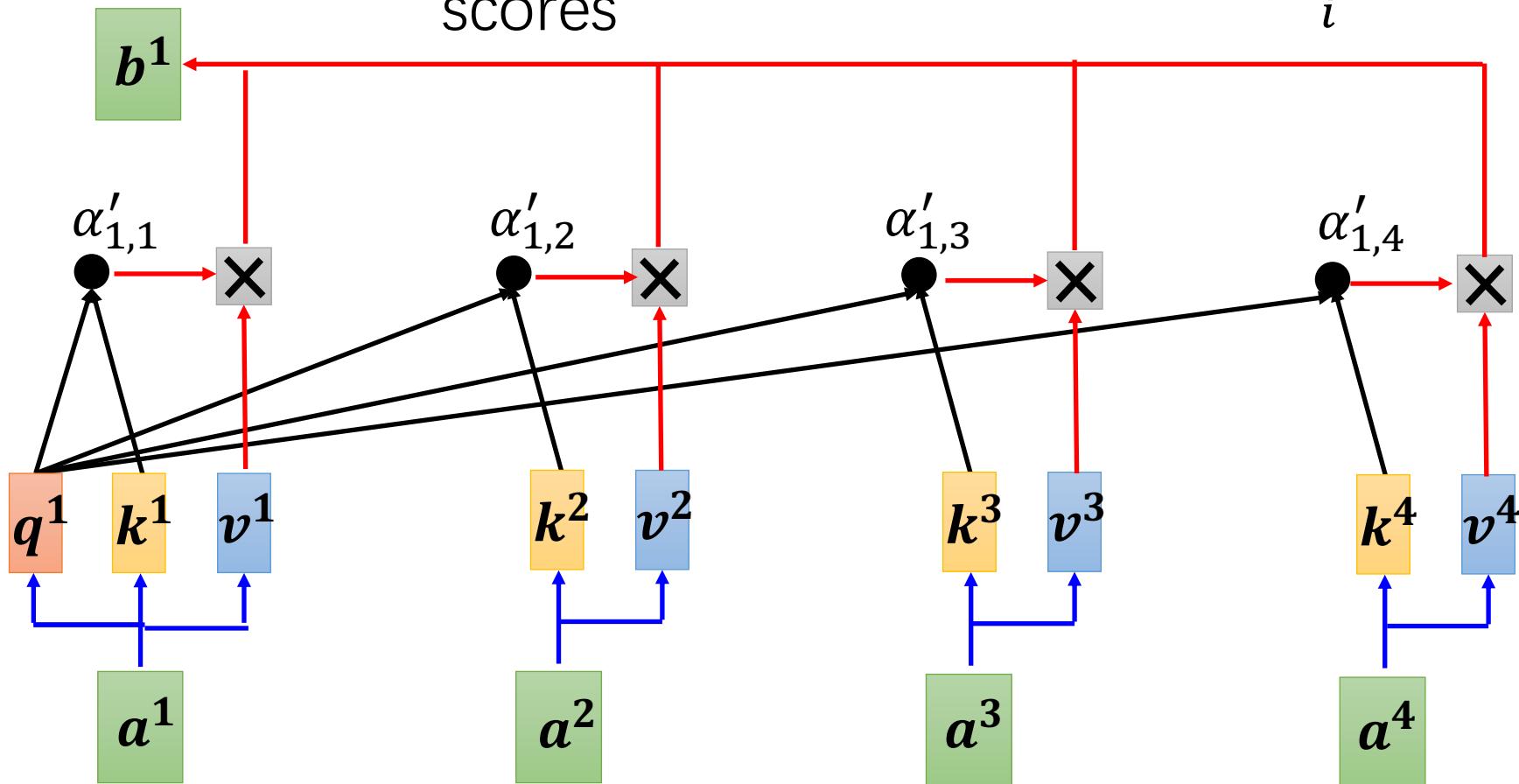
$$k^4 = W^k a^4$$

$$k^1 = W^k a^1$$

# Self-attention

Extract information  
based on attention  
scores

$$\mathbf{b}^1 = \sum_i \alpha'_{1,i} \mathbf{v}^i$$



$$\mathbf{v}^1 = W^v \mathbf{a}^1$$

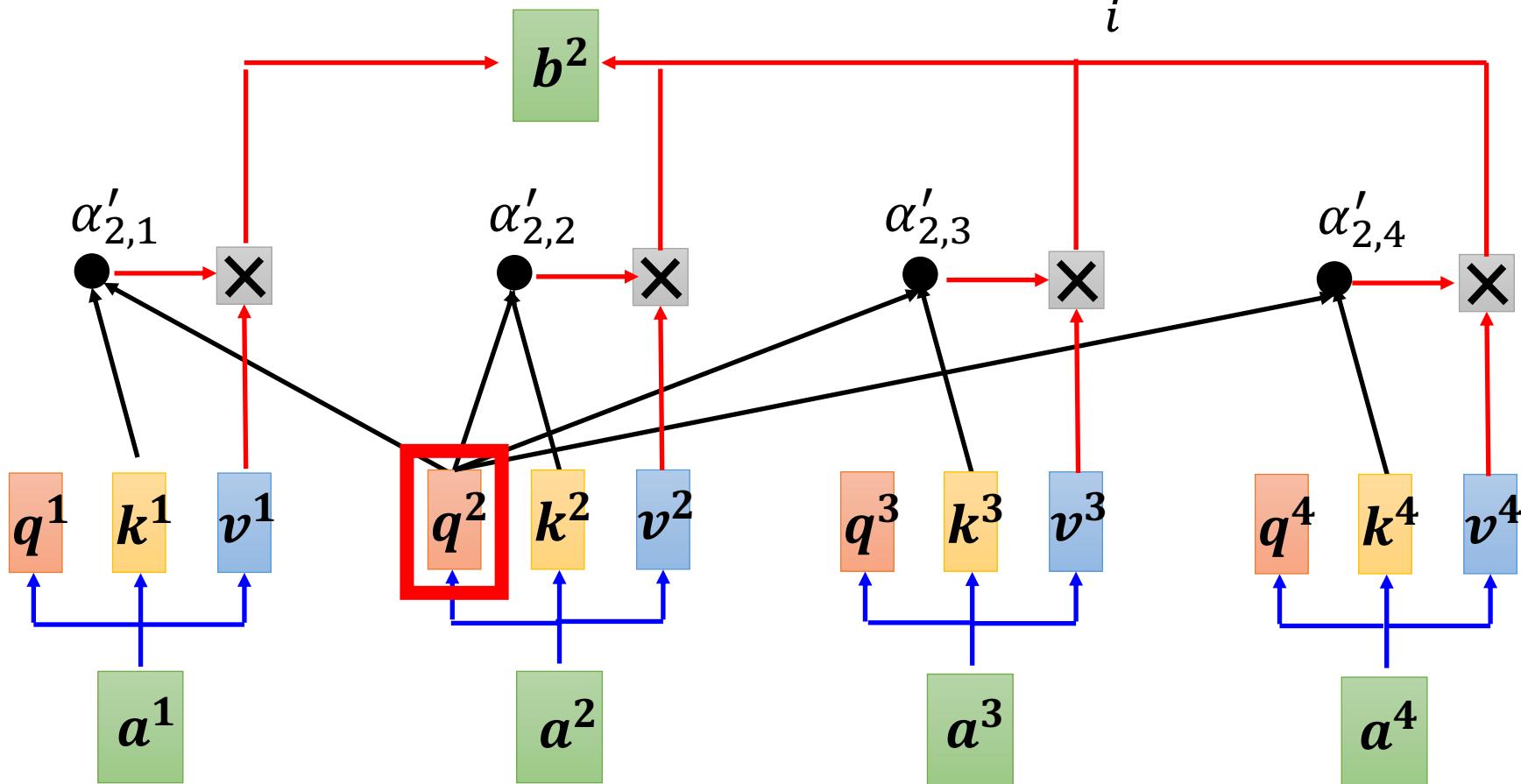
$$\mathbf{v}^2 = W^v \mathbf{a}^2$$

$$\mathbf{v}^3 = W^v \mathbf{a}^3$$

$$\mathbf{v}^4 = W^v \mathbf{a}^4$$

# Self-attention

$$\mathbf{b}^2 = \sum_i \alpha'_{2,i} \mathbf{v}^i$$



# Self-attention

$$\alpha_{1,1} = \begin{matrix} k^1 \\ q^1 \end{matrix}$$

$$\alpha_{1,2} = \begin{matrix} k^2 \\ q^1 \end{matrix}$$

$$\alpha_{1,1}$$

$$k^1$$

$$\alpha_{1,2}$$

$$k^2$$

$$\alpha_{1,3}$$

$$k^3$$

$$\alpha_{1,4}$$

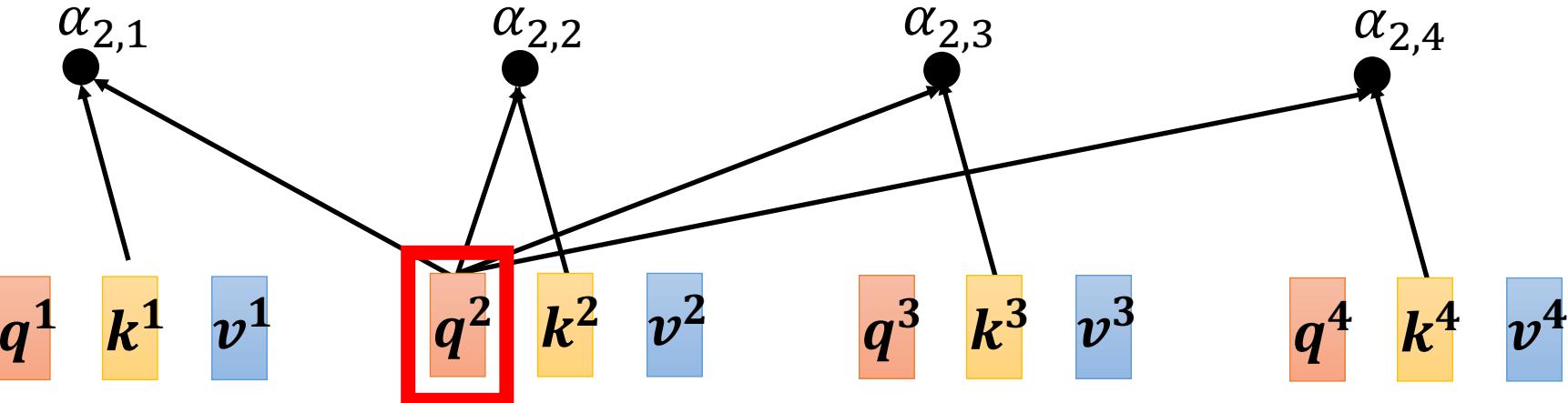
$$k^4$$

$$\alpha_{1,3} = \begin{matrix} k^3 \\ q^1 \end{matrix}$$

$$\alpha_{1,4} = \begin{matrix} k^4 \\ q^1 \end{matrix}$$

$$\alpha_{1,4} = \begin{matrix} k^4 \\ q^1 \end{matrix}$$

$$q^1$$



$$\begin{matrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{matrix}$$

$$A'$$

$$\text{softmax}$$

$$\begin{matrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{matrix}$$

$$A$$

$$k^1$$

$$k^2$$

$$k^3$$

$$k^4$$

$$K^T$$

$$Q$$

$$k^1$$

$$k^2$$

$$k^3$$

$$k^4$$

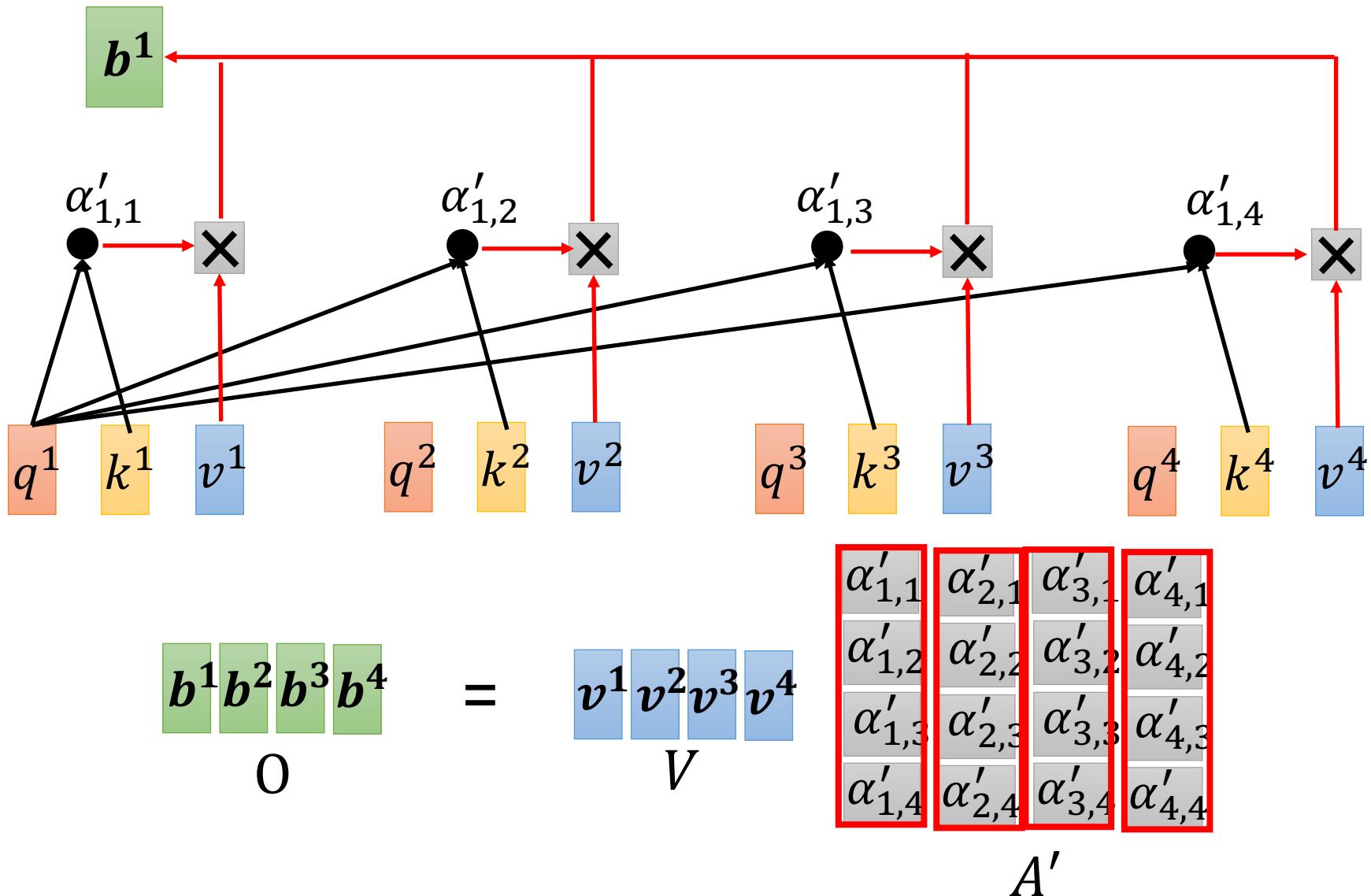
$$q^1$$

$$q^2$$

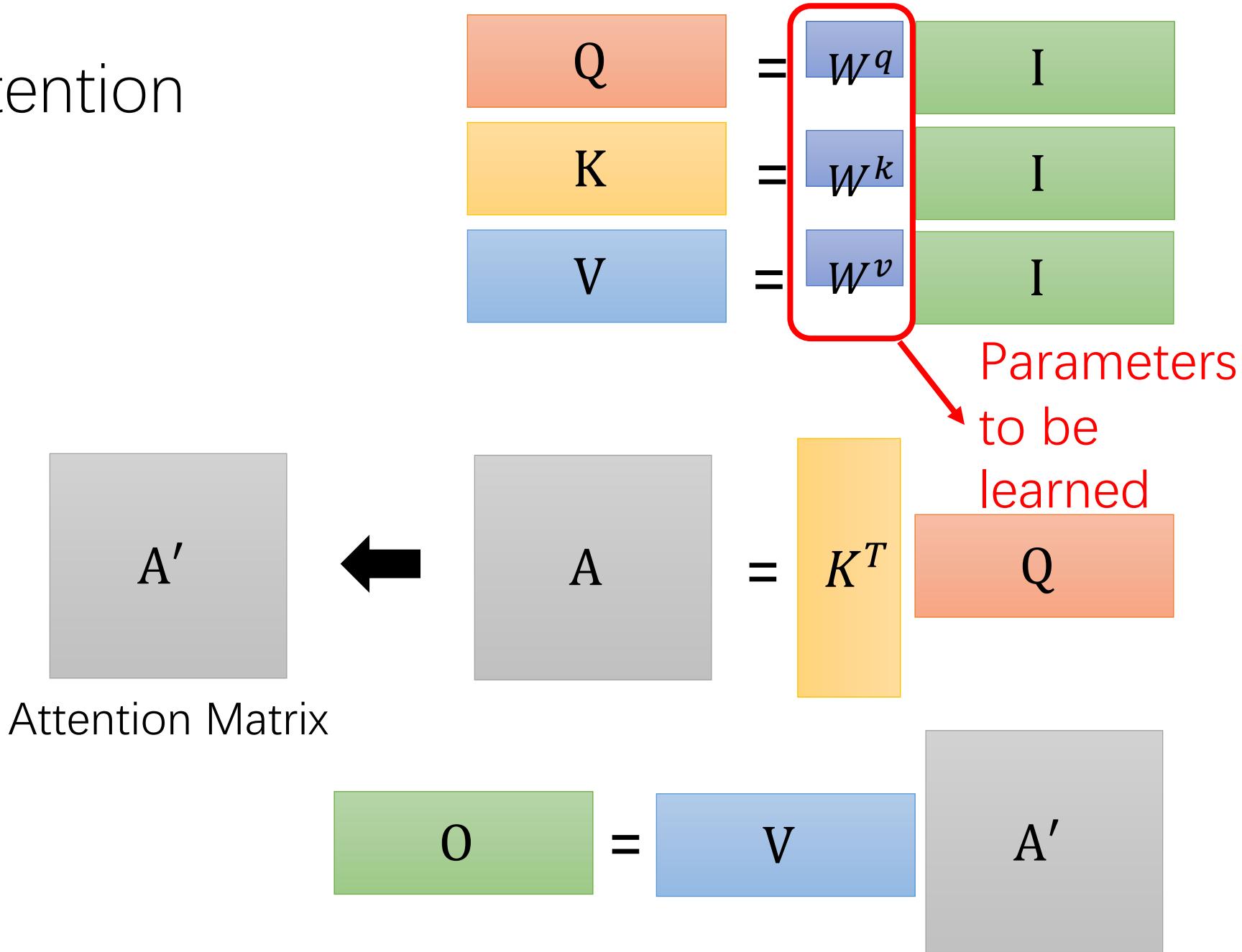
$$q^3$$

$$q^4$$

# Self-attention



# Self-attention

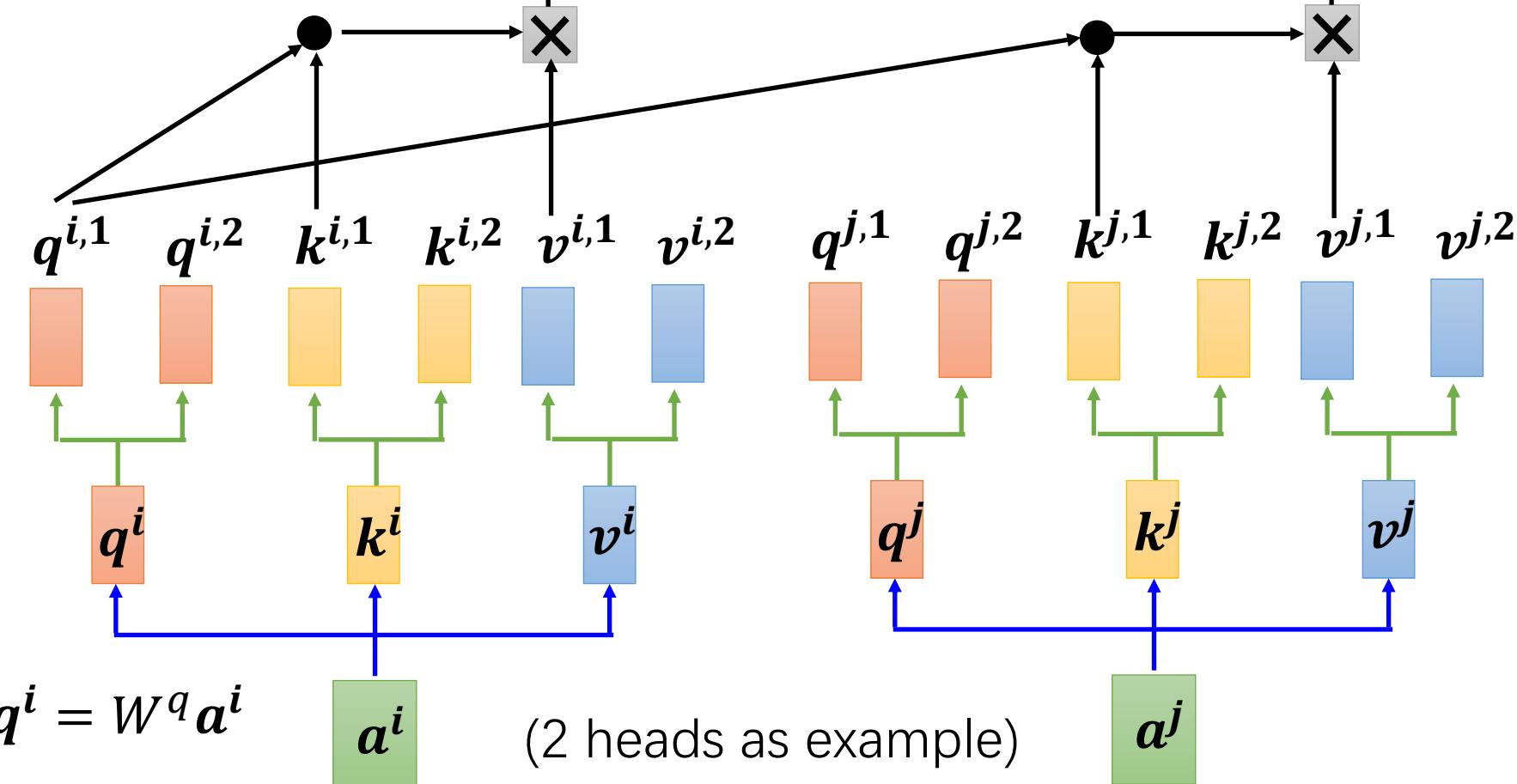


# Multi-head Self-attention

Different types of relevance

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$

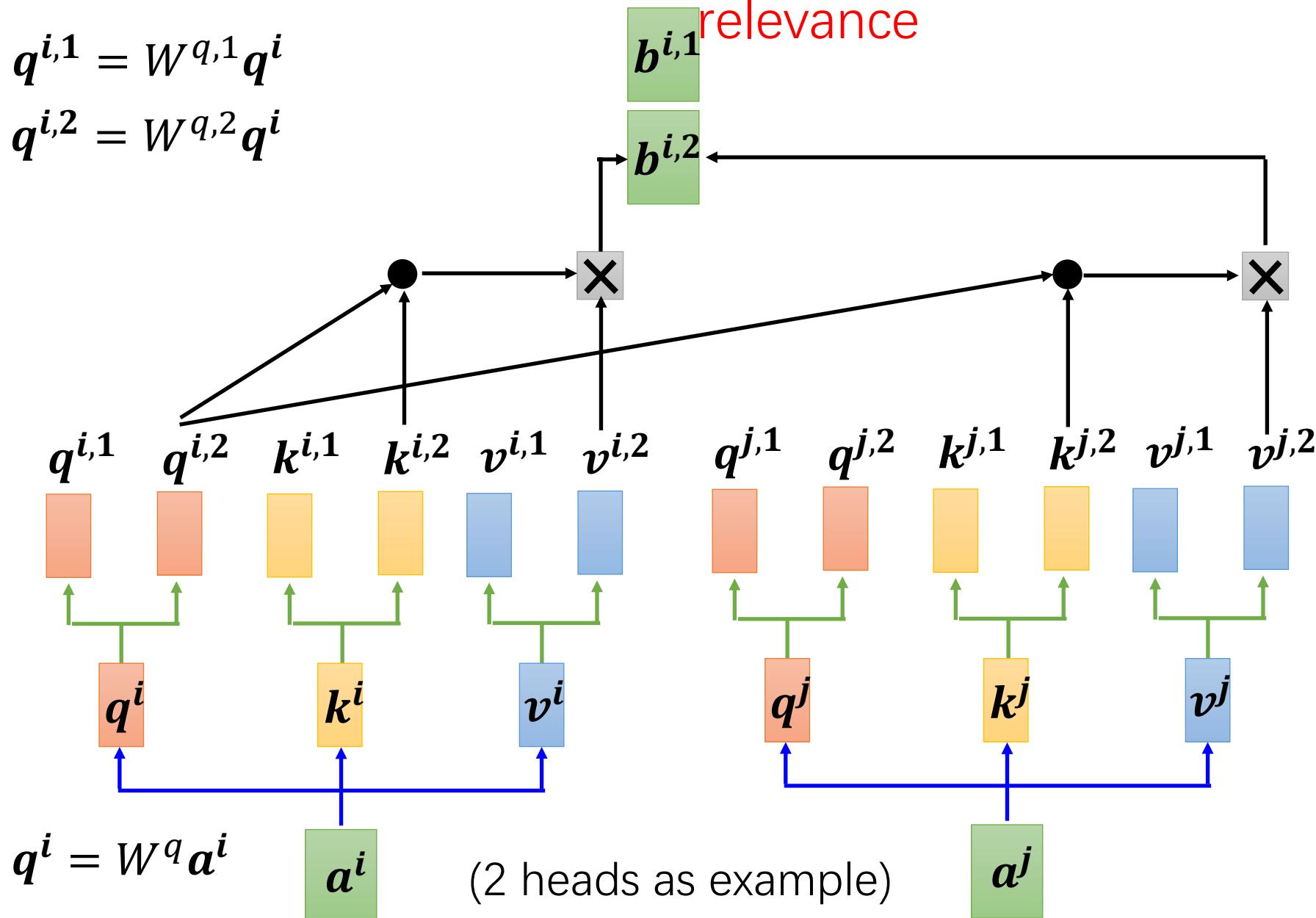


# Multi-head Self-attention

$$q^{i,1} = W^{q,1} q^i$$

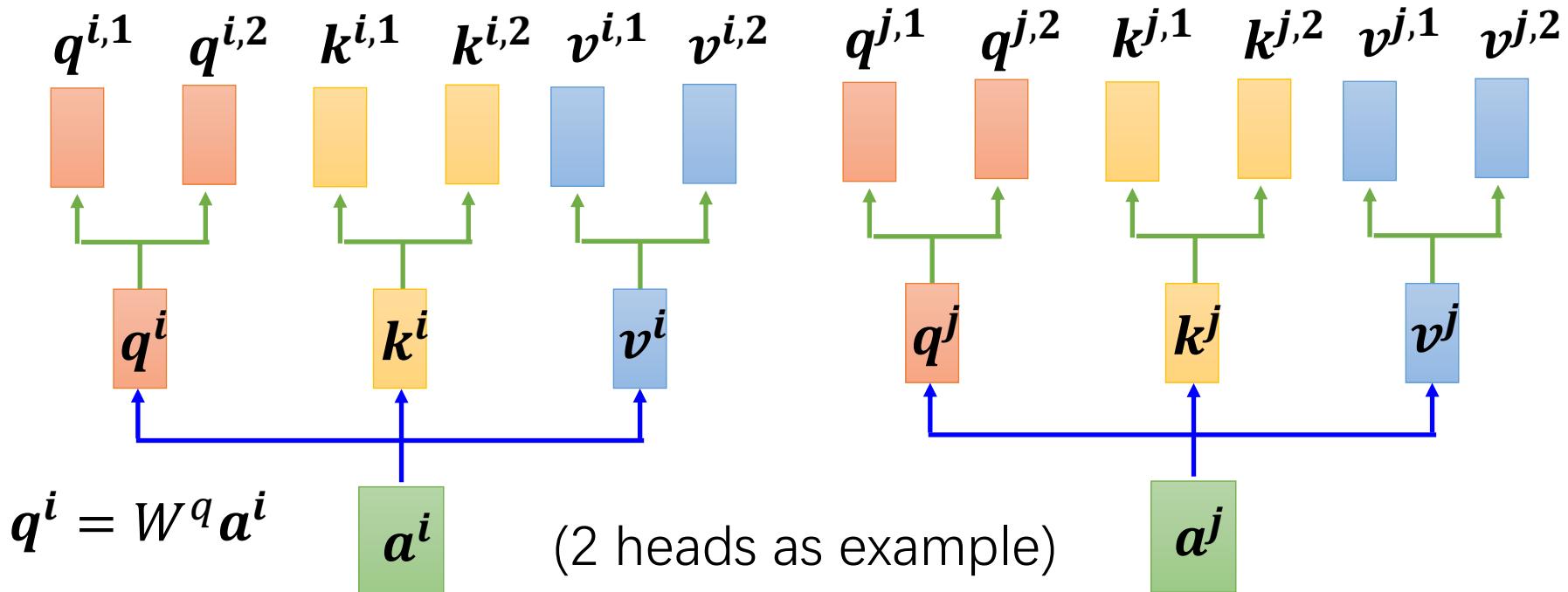
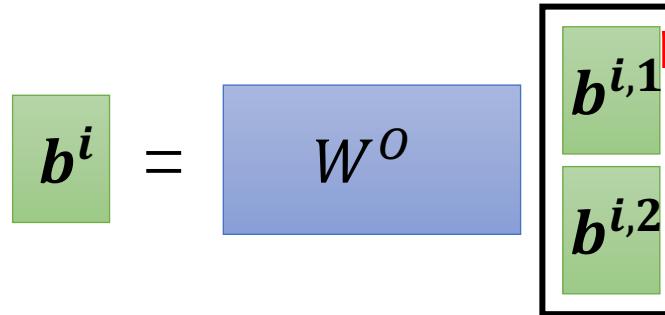
$$q^{i,2} = W^{q,2} q^i$$

Different types of relevance



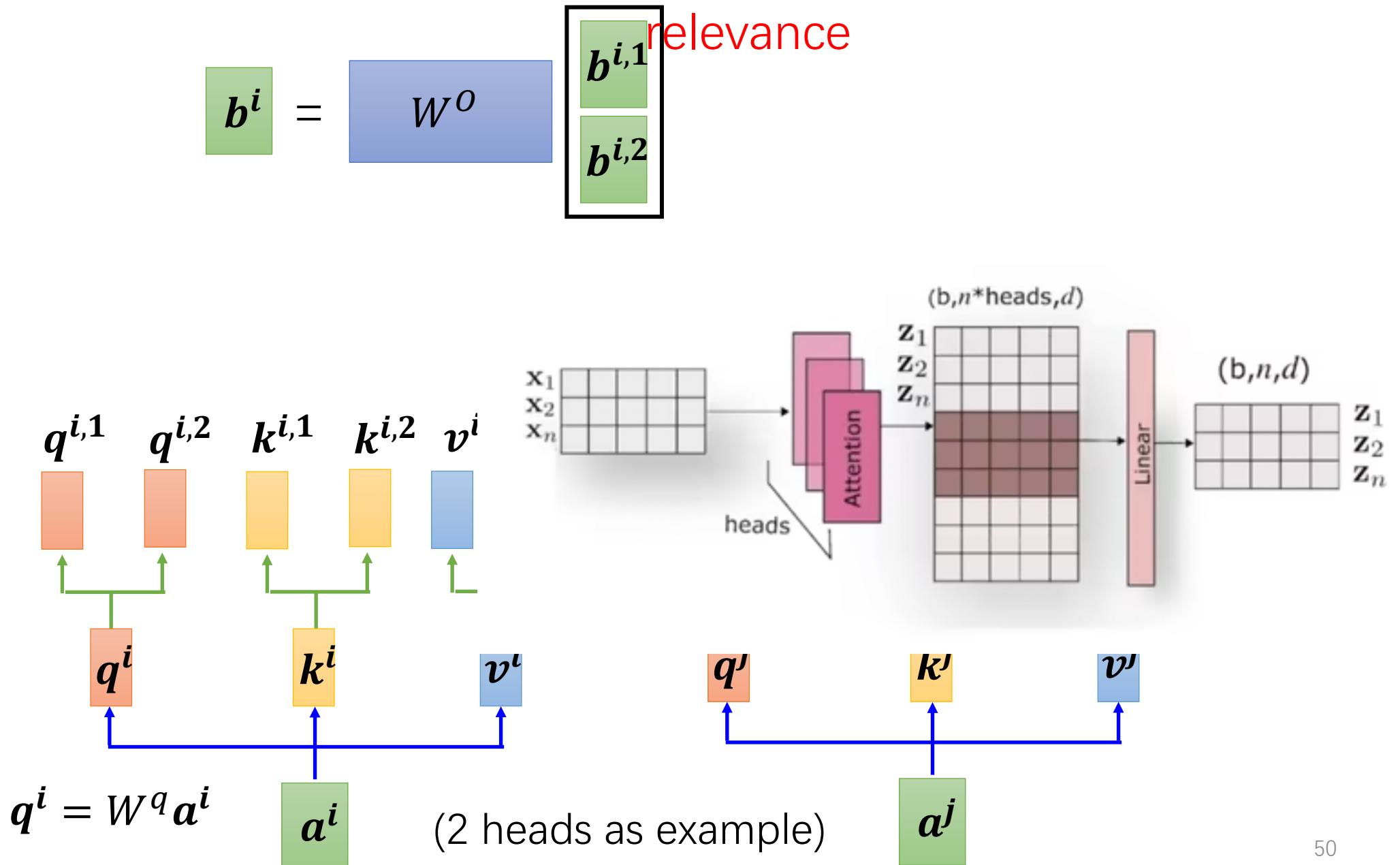
# Multi-head Self-attention

Different types of relevance

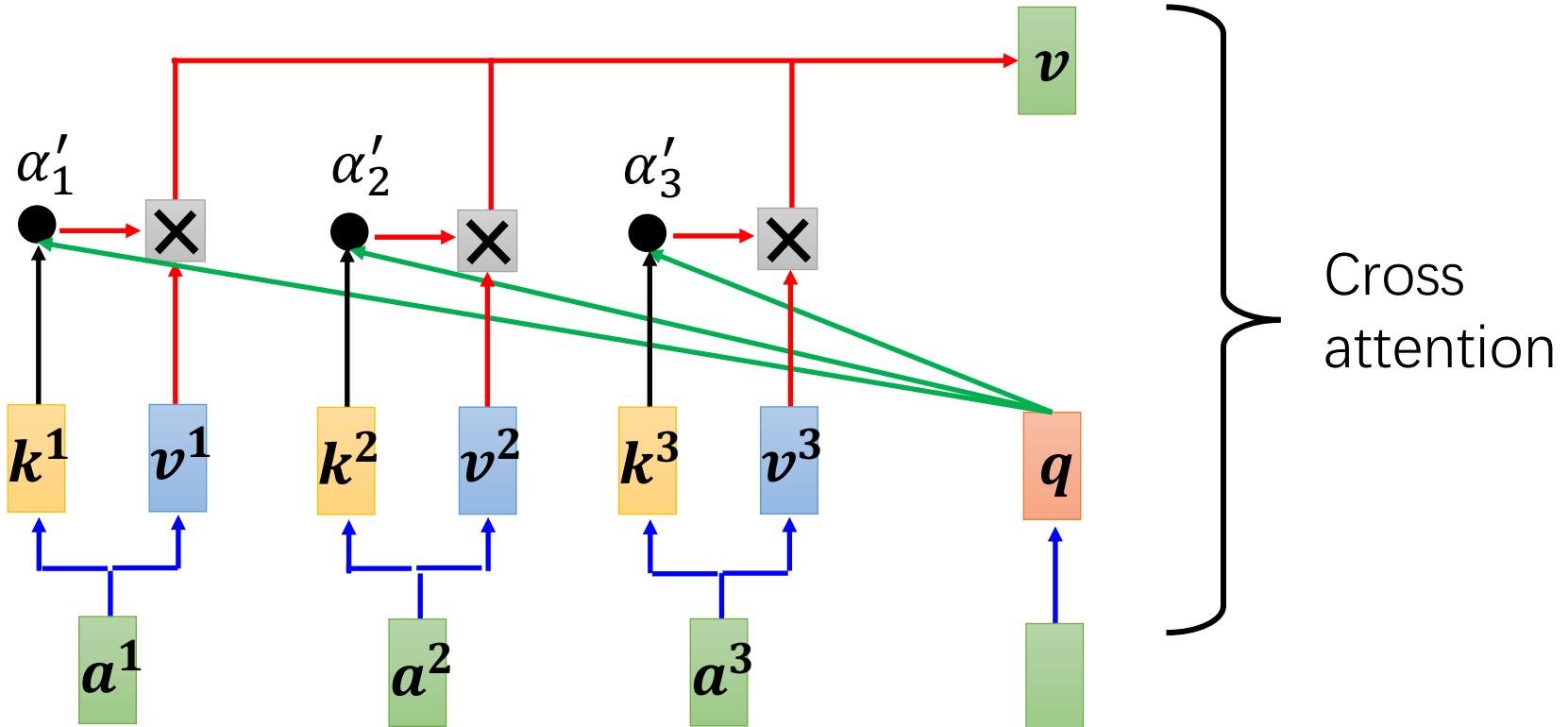


# Multi-head Self-attention

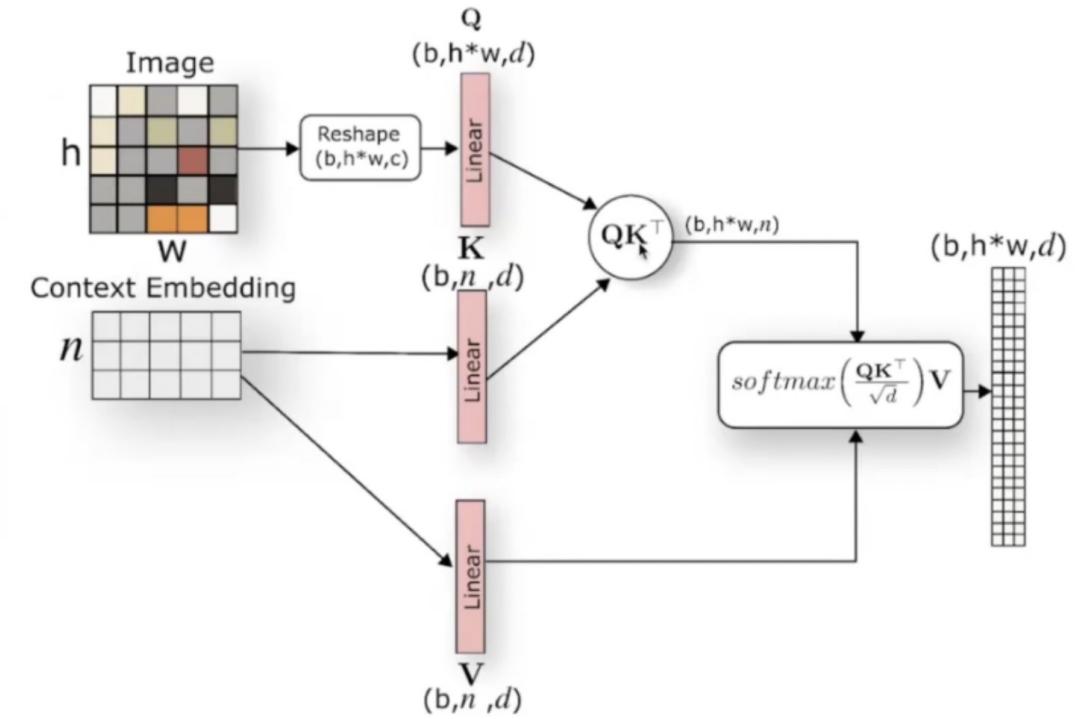
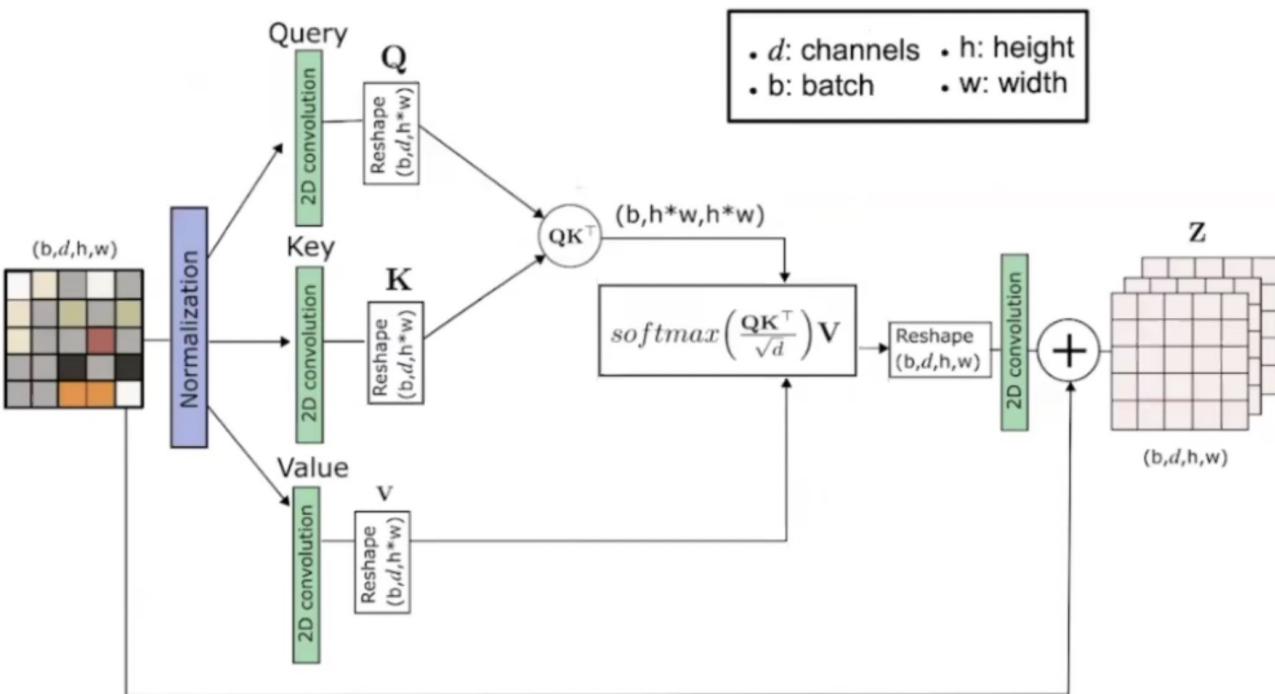
Different types of relevance



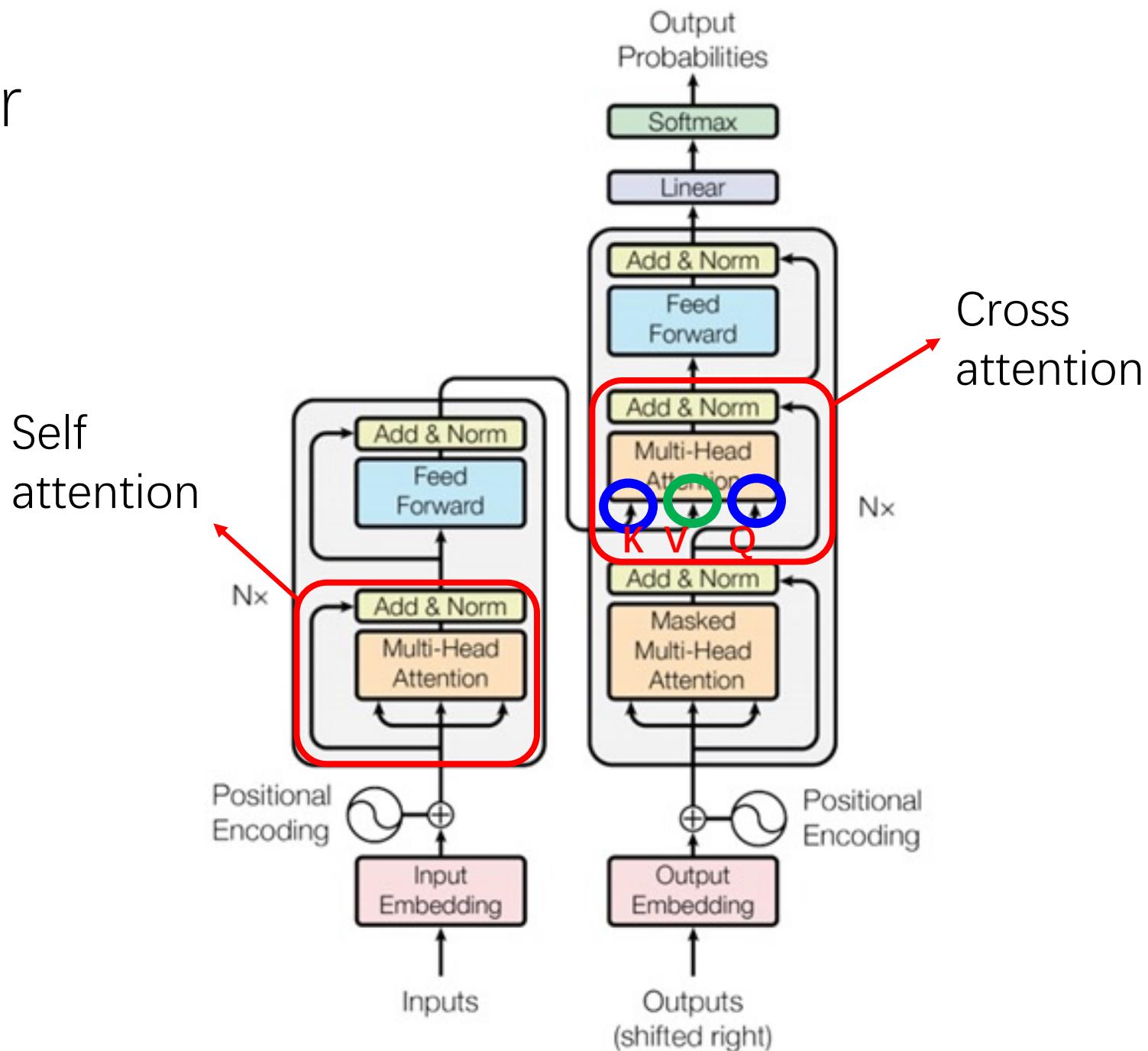
# Cross Attention



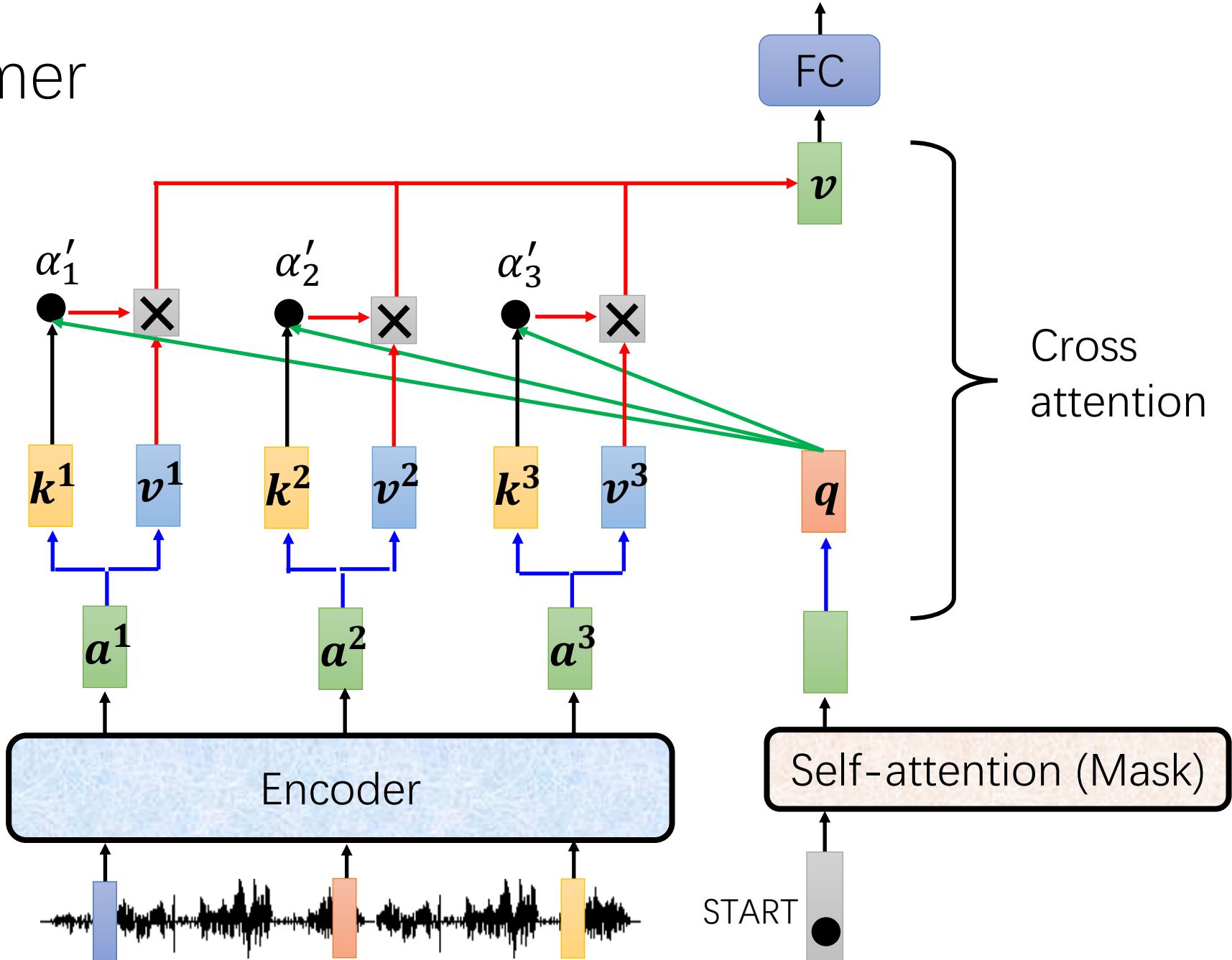
# Cross Attention



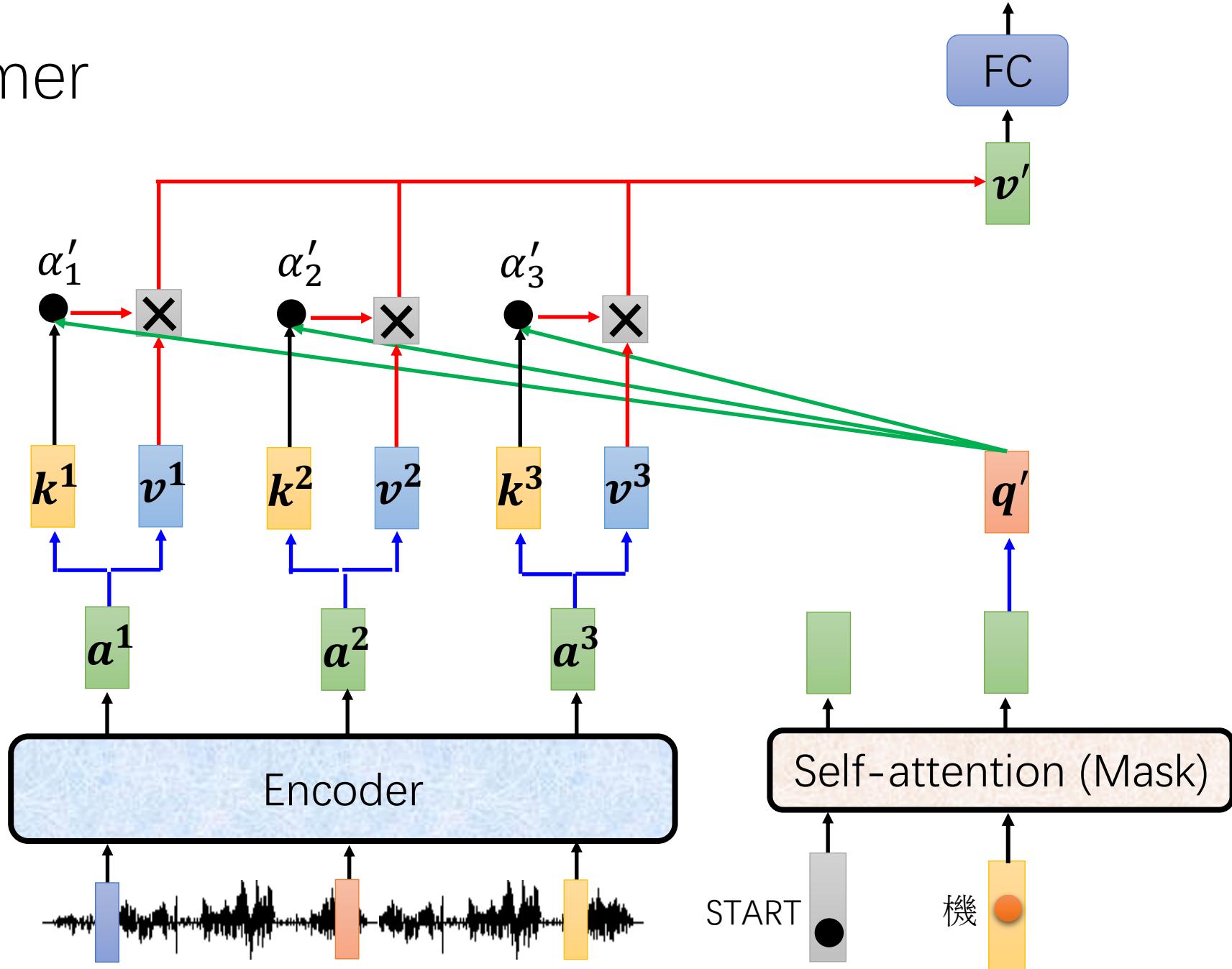
# Transformer



# Transformer

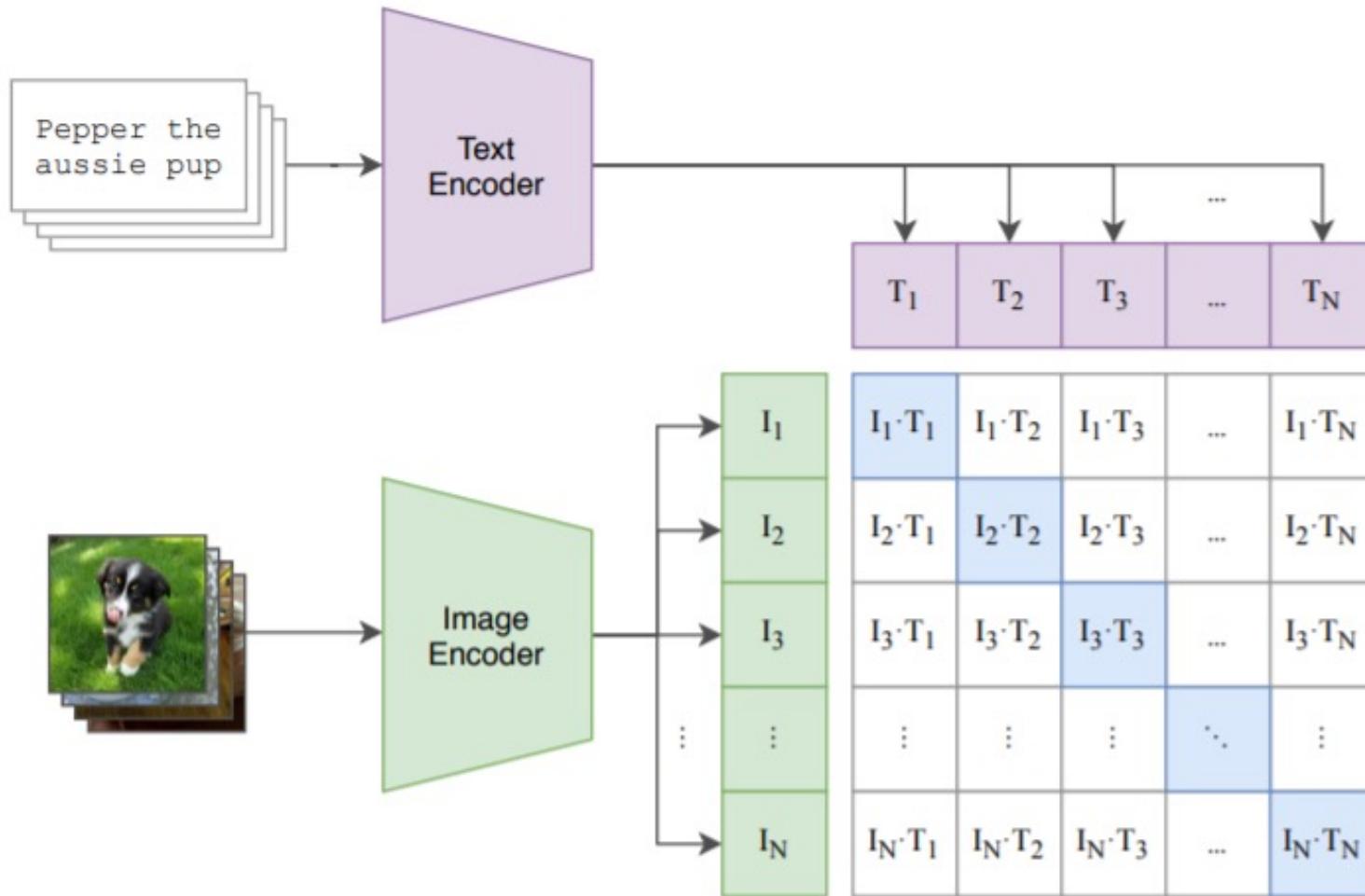


# Transformer



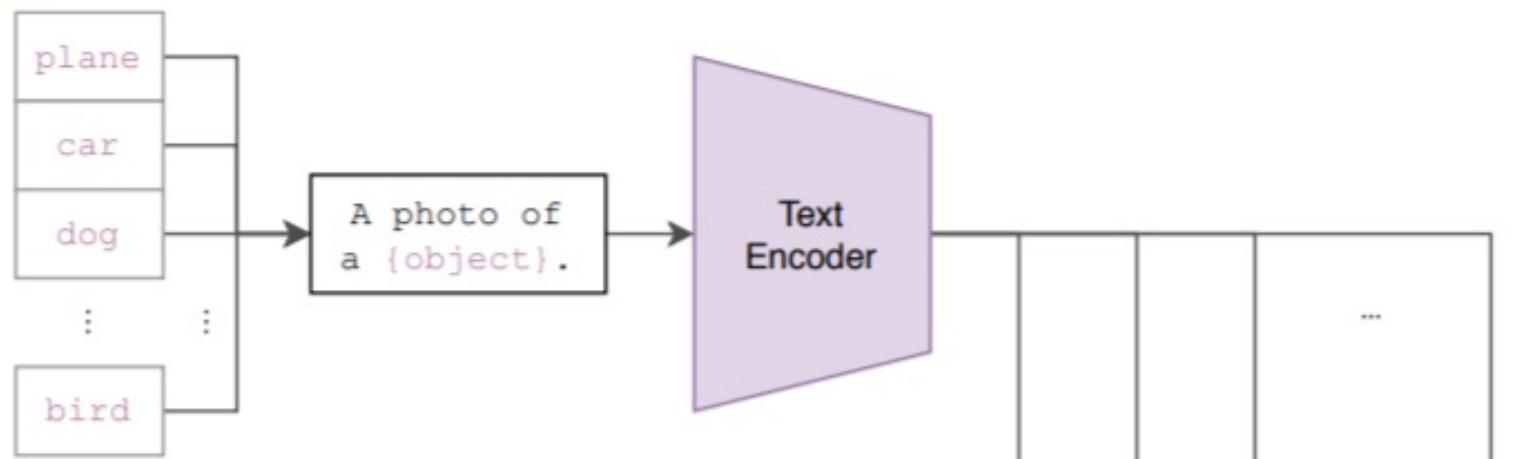
# CLIP Contrastive Language-Image Pre-Training

## (1) Contrastive pre-training

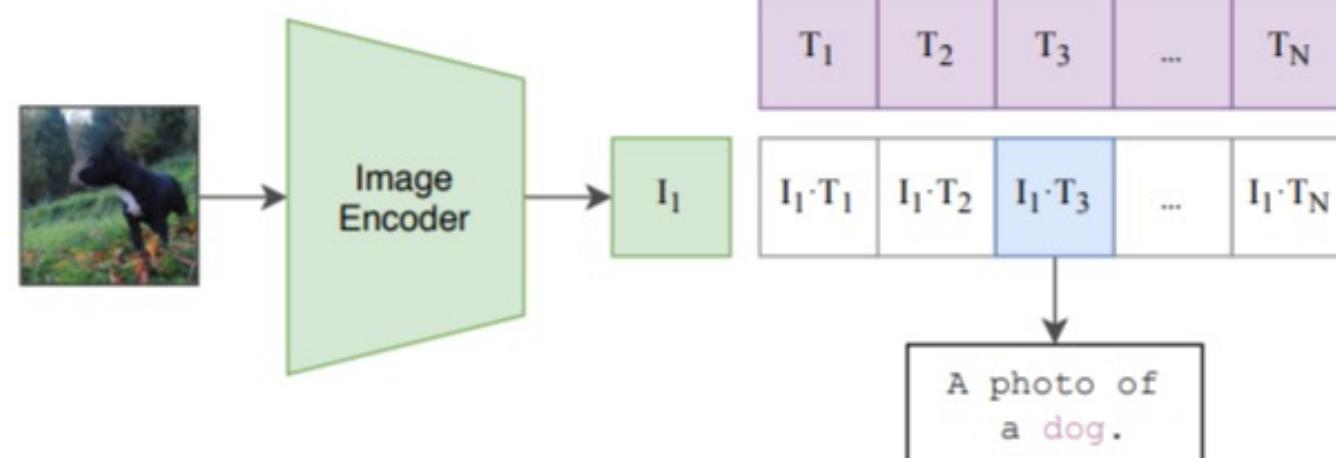


# CLIP Contrastive Language-Image Pre-Training

## (2) Create dataset classifier from label text



## (3) Use for zero-shot prediction



# LDM -Latent Diffusion Models -潜在扩散模型

