

# Hierarchical Variational Autoencoder and Denoise Diffusion Probabilistic Model

Duan zhibin

# Variational Autoencoder

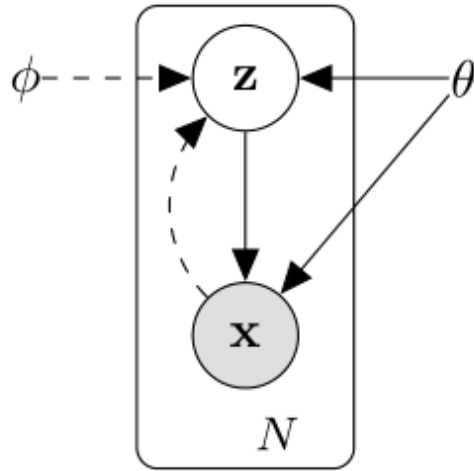


Figure 1: The type of directed graphical model under consideration.

## □ Generate model

Given  $N$  i.i.d samples:  $\mathbf{X} = \{x^{(i)}\}_i^N$

(i).  $z^{(i)} \sim p_{\theta^*}(z)$

(ii).  $x^{(i)} \sim p_{\theta^*}(x|z)$

## □ Intractable posterior distributions

$$p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$$

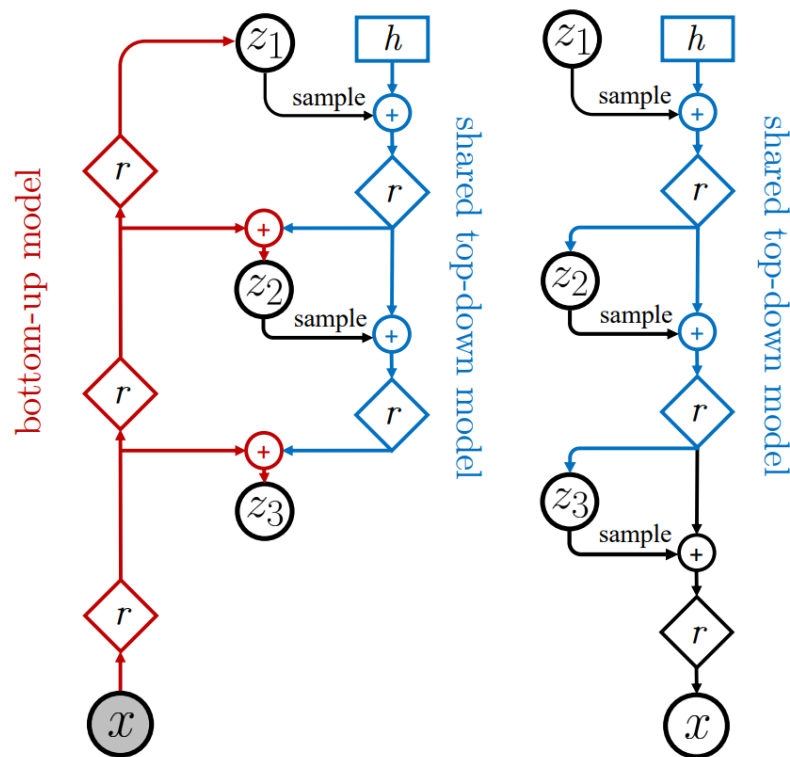
## □ Variational posterior

$$q_{\phi}(z^{(i)}|x^{(i)}) \sim N(f_{\mu}(x^{(i)}), f_{\sigma}(x^{(i)}))$$

## □ The variational bound

$$\log p_{\theta}(x^{(i)}) = D_{KL}(q_{\phi}(z^{(i)}|x^{(i)}) || p_{\theta}(z^{(i)}|x^{(i)})) + L(\theta, \phi; x^{(i)})$$

# Hierarchical variational autoencoder



(a) Bidirectional Encoder (b) Generative Model

## Generate model

Given  $N$  i.i.d samples:  $\mathbf{X} = \{x^{(i)}\}_i^N$

(i).  $z_1^{(i)} \sim p_{\theta^*}(z_1), \dots, z_l^{(i)} \sim p_{\theta^*}(z_l),$

(ii).  $x^{(i)} \sim p_{\theta^*}(x | z_L)$

## Variational conditional posterior

$$q_{\phi}(z_l^{(i)} | x_t^{(i)}, z_{l-1}^{(i)}) \sim N\left(f_{\mu}(x^{(i)}, z_{l-1}^{(i)}), f_{\sigma}(x^{(i)}, z_{l-1}^{(i)})\right)$$

## Training object

$$L_{\text{VAE}}(x) := \mathbb{E}_{q(z|x)} \left[ \log p(x | z) - \text{KL}(q(z_1 | x) \| p(z_1)) \right] \\ - \sum_{l=2}^L \mathbb{E}_{q(z_{<l}|x)} \left[ \text{KL}(q(z_l | x, z_{<l}) \| p(z_l | z_{<l})) \right]$$

# Hierarchical variational autoencoder



Fig 2: unconditional CIFAR10 generative  
From NVAE

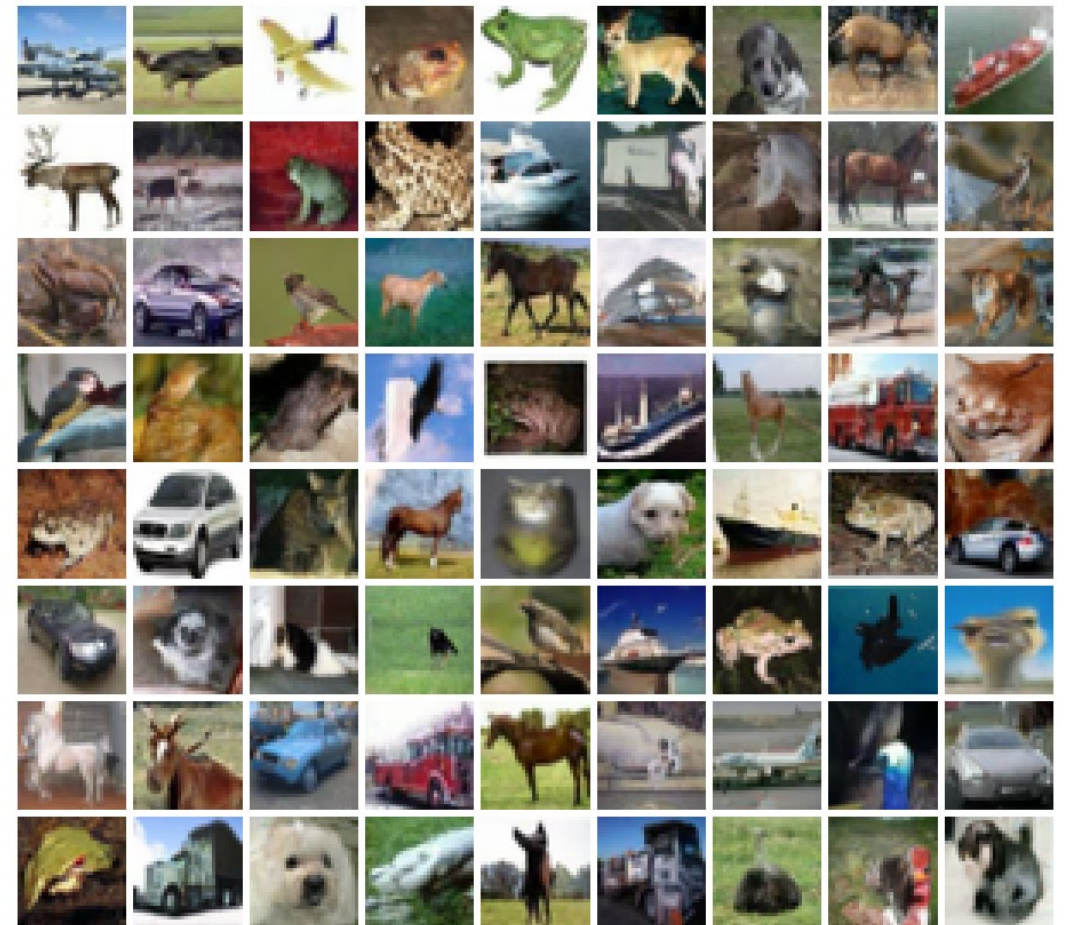


Fig 2: unconditional CIFAR10 generative  
from denoise diffusion model

# Denoise diffusion probabilistic model

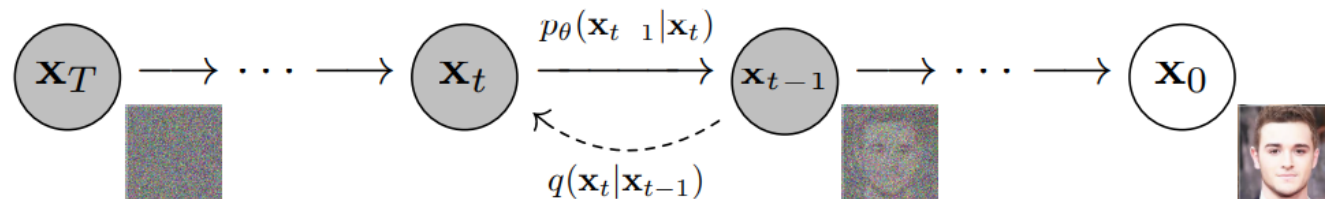


Figure 2: The directed graphical model in this work.

## Encoder

$$\begin{cases} q(x_t | x_{t-1}) = N(x_t; \alpha_t x_{t-1}, \beta_t^2 I) \\ q(x_{1:T} | x_0) := \prod_{t=1}^T q(x_t | x_{t-1}) \end{cases}$$
$$\Rightarrow q(x_{1:T} | x_0) := N(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

## Analytic conditional posterior

$$p(x_{t-1} | x_t, x_0) = N\left(x_{t-1}; \frac{\alpha_t \bar{\beta}_{t-1}}{\bar{\beta}_t^2} x_t + \frac{\bar{\alpha}_{t-1} \bar{\beta}_t^2}{\bar{\beta}_t^2} x_0, \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} I\right),$$

[1] Sohl-Dickstein, Jascha, et al. "Deep unsupervised learning using nonequilibrium thermodynamics." ICML2015

[2] Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." NeurIPS 2020.

# Denoise diffusion probabilistic model

## □ Generative model

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

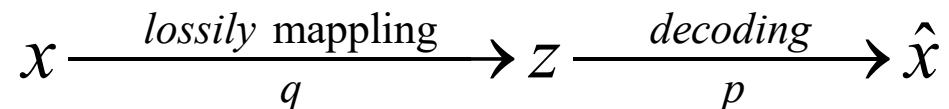
## □ Object function

$$L_{\text{vb}} = \mathbb{E}_{q(\mathbf{x}_0)} \left[ \underbrace{D_{\text{KL}}[q(\mathbf{x}_T | \mathbf{x}_0) || p(\mathbf{x}_T)]}_{L_T} + \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [D_{\text{KL}}[q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)]]}_{L_{t-1}} \right. \\ \left. \underbrace{- \mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)]}_{L_0} \right]. \quad (1)$$

[1] Sohl-Dickstein, Jascha, et al. "Deep unsupervised learning using nonequilibrium thermodynamics." ICML2015

[2] Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." NeurIPS 2020.

# Lossy compressor



## □ Rate

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

$$R = E_{p_{data}(x)} D_{KL}(q(z|x) \| p(z))$$

## □ Distortion

$$D = E_{p_{data}(x)} E_{q(z|x)} f(x, d(z))$$

## □ A rate-distortion trade-off

$$\begin{aligned} L_{vae}(x^{(i)}) &= D_{KL}(q_{\phi}(z^{(i)} | x^{(i)}) \| p_{\theta}(z^{(i)} | x^{(i)})) + L(\theta, \phi; x^{(i)}) \\ &= R(x^{(i)}) + D(x^{(i)}) \end{aligned}$$

## □ $\gamma$ – optimal

$$\lambda^* = \min_{d, p, q} E_{p_{data}(x)} E_{q(z|x)} (x, d(z))$$

$$s.t. \quad \min_{d, p, q} E_{p_{data}(x)} D_{KL}(q(z|x) \| p(z)) \leq \gamma$$

## □ $\lambda$ – optimal

$$\gamma^* = \min_{d, p, q} E_{p_{data}(x)} D_{KL}(q(z|x) \| p(z))$$

$$s.t. \quad E_{p_{data}(x)} E_{q(z|x)} (x, d(z)) \leq \lambda$$

# Progressively lossier compressors

## □ Rate

*Given :  $p(z_{1:k})$*

$$R_k = E_{p_{data}(x)} D_{KL} \left( q(z_{1:k} | x) \| p(z_{1:k}) \right)$$

## □ Distortion

$$D_k = E_{p_{data}(x)} E_{q(z_{1:k}|x)} f(x, d_k(z_{1:k}))$$

*A sequence of  $(R_k, D_k)$ , and  $d_k$*

## □ Pareto Optimality

帕累托最优 (Pareto Optimality)，是指资源分配的一种理想状态，假定固有的一群人和可分配的资源，从一种分配状态到另一种状态的变化中，在没有使任何人境况变坏的前提下，使得至少一个人变得更好。

## □ VAE are pareto optimality

$$L_{vae} \left( x^{(i)} \right) = D_{KL} \left( q_{\phi} \left( z^{(i)} | x^{(i)} \right) \| p_{\theta} \left( z^{(i)} | x^{(i)} \right) \right) + L \left( \theta, \phi; x^{(i)} \right)$$

## □ HVAE are not pareto optimality

$$L_{VAE} (x) := E_{q(z|x)} \left[ \log p(x|z) - \text{KL} \left( q(z_1|x) \| p(z_1) \right) \right] \\ - \sum_{l=2}^L E_{q(z_{<l}|x)} \left[ \text{KL} \left( q(z_l|x, z_{<l}) \| p(z_l|z_{<l}) \right) \right]$$



# DDPM are good progressively lossier compressors

## ▣ Progressive coding objective

$$\begin{aligned} \min_{d_{1:T}, p, q} \sum_{k=1}^T \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z|x)} f(x, d_k(z_{1:k})) \\ \text{s. t. } \mathbb{E}_{p_{\text{data}}(x)} D_{\text{KL}}(q(z_{1:k} | x) \parallel p(z_{1:k})) \leq \gamma_k \\ \forall k \in \{1, \dots, T\}. \end{aligned}$$

## ▣ Optimizing

$$\begin{aligned} d_k^* &= \arg \min_{d_k} \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z_{1:k}|x)} f(x, d_k(z_{1:k})) \\ p^* &= \arg \min_p \mathbb{E}_{p_{\text{data}}(x)} D_{\text{KL}}(q(z_{1:T} | x) \parallel p(z_{1:T})) \end{aligned}$$

## ▣ Rate $\{\gamma_1, \dots, \gamma_T\}$

$$\gamma_k = \mathbb{E}_{p_{\text{data}}(x)} D_{\text{KL}}(q(z_{1:k} | x) \parallel p^*(z_{1:k})).$$

## ▣ Parameter-sharing setting

$$p(z_k | z_{k-1}) = q(z_k | x = d_{k-1}(z_{k-1})).$$

## ▣ The progressive coding objective

$$\min_{d_{1:T}} \sum_{k=1}^T \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z_k|x)} f(x, d_k(z_k))$$

# HAVE are not good progressively lossier compressors

## □ Rate

$$R_K = \mathbb{E}_{p_{\text{data}}(x)} D_{\text{KL}}(q(z_{1:K} \mid x) \parallel p(z_{1:K}))$$

## □ Distortion

$$D_K = \min_d \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z_{1:K} \mid x)} \|x - d(z_{1:K})\|^2$$

## □ $\gamma$ – optimal using a Lagrange multiplier

$$\begin{aligned} \lambda^* &= \min_{d, p, q} \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z_{1:K} \mid x)} \|x - d(z_{1:K})\|^2 \\ &\text{s. t. } \mathbb{E}_{p_{\text{data}}(x)} D_{\text{KL}}(q(z_{1:K} \mid x) \parallel p(z_{1:K})) \leq R_K \end{aligned}$$

# HAVE are not good progressively lossier compressors

## ▣ Rate

Dataset	Conventional Training		$\gamma$ -Optimal Training	
	$R_K$	$D_K$	$R'_K$	$\geq \lambda^*$
SVHN	0.0168	0.0200	0.0162	<b>0.0063</b>
CelebA	0.0142	0.0274	0.0135	<b>0.0128</b>
LSUN	0.0187	0.0611	0.0181	<b>0.0339</b>

# Two-stage vae

To enforce pareto-optimality  $(R_K, D_K)$

## □ First stage

$$\begin{aligned} \min_{d, p, q} \mathbb{E}_{p_{\text{data}}(x)} D_{\text{KL}}(q(z_{1:K} \mid x) \parallel p(z_{1:K})) \\ \text{s. t. } \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z_{1:K} \mid x)} f(x, d(z_{1:K})) \leq \lambda \end{aligned}$$

## □ Second stage

$$\begin{aligned} \min_{p', q'} \mathbb{E}_{p_{\text{data}}(x)} \left[ \mathbb{E}_{q(z_{1:T})} \ln p(x \mid z_{1:T}) \right. \\ \left. + D_{\text{KL}}(q(z_{1:T} \mid x) \parallel p(z_{1:T})) \right], \end{aligned}$$

# Experiments

Table 2. Quantitative evaluation of conventionally-trained versus our progressively-coded HVAEs. ELBOs are per-dimension (i.e., normalized by the data dimensionality).

Dataset	Conventional Training		Progressive Coding	
	ELBO	FID	ELBO	FID
SVHN	−1.31	18.92	−1.46	<b>10.07</b>
CelebA	−1.44	13.33	−1.58	<b>8.54</b>
LSUN	−1.72	40.71	−1.78	<b>36.87</b>

Conventional  
training



Progressive  
coding

