



भारतीय प्रौद्योगिकी संस्थान दिल्ली
Indian Institute of Technology Delhi

COV877

Special Module on Visual Computing

Generative AI for Visual Content Creation: Image, Video, and 3D

VAE & VQVAE

Instructor:

Dr. Lokender Tiwari

Research Scientist

Variational Autoencoder (VAE)

Variational Autoencoders

- Probabilistic Generative Model (like normalizing flows)
- Learn distribution $Pr(\mathbf{x})$ over data
- Post training we can sample from this distribution

Latent Variable Model

- An indirect approach to represent probability distribution of data $Pr(\mathbf{x})$

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

joint distribution of data & unobserved latent variable

marginalization

conditional probability

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}|\mathbf{z}) Pr(\mathbf{z}) d\mathbf{z}$$

prior

likelihood

complex distributions $Pr(\mathbf{x})$ can be represented using simple $Pr(\mathbf{x}|\mathbf{z})$ and $Pr(\mathbf{z})$

e.g. Mixture of Gaussians

Latent Variable Model - Mixture of Gaussian Example

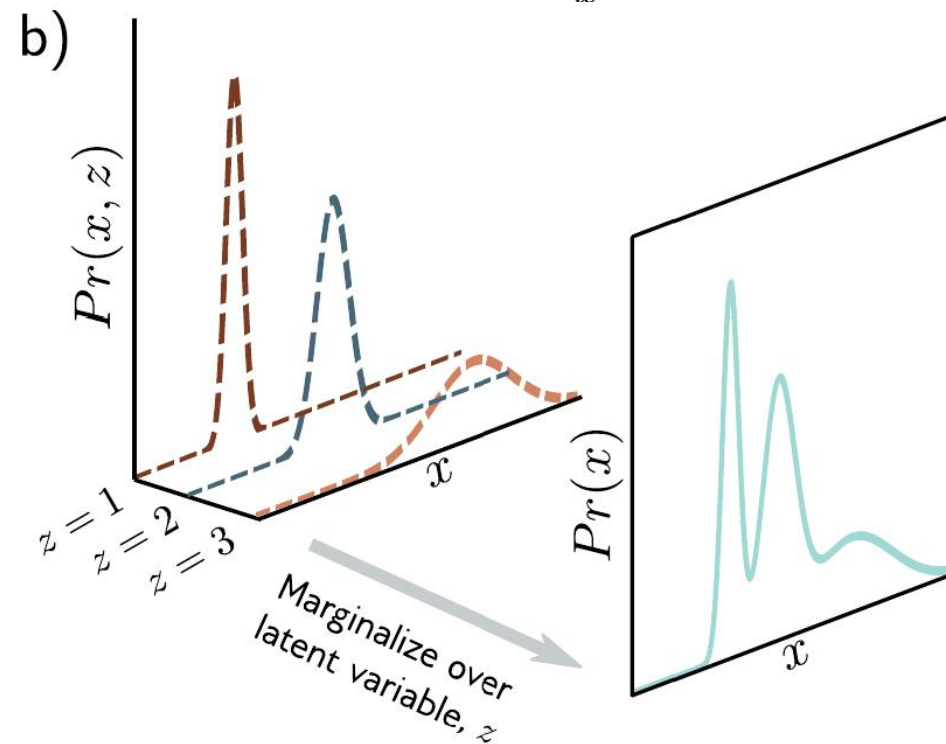
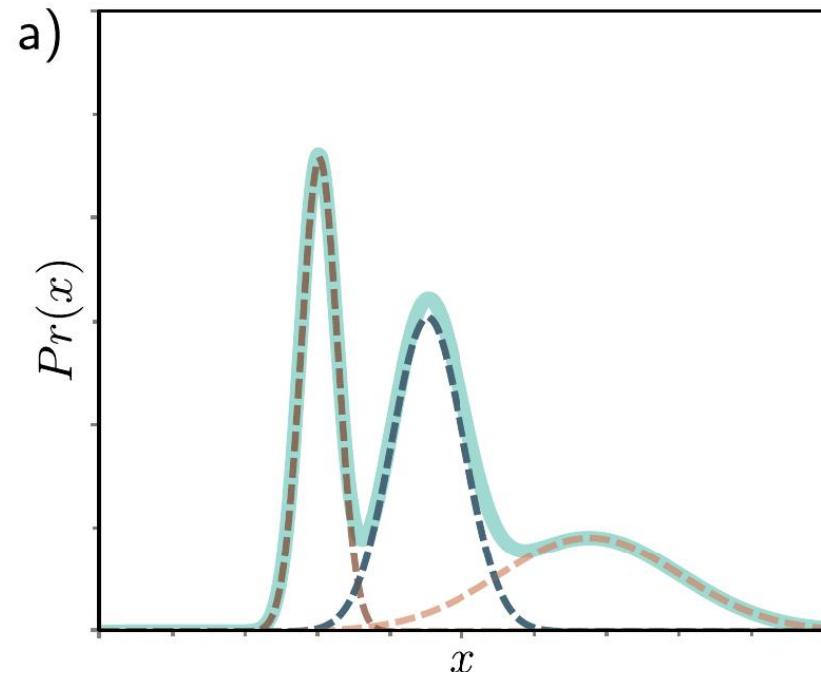
- \mathbf{z} is a discrete latent variable with categorical prior distribution $Pr(\mathbf{z})$

$$Pr(z = n) = \lambda_n$$

- Likelihood $Pr(\mathbf{x}|\mathbf{z} = \mathbf{n})$ of data \mathbf{x} given $\mathbf{z} = \mathbf{n}$ is normally distributed with mean μ_n and variance σ_n^2

$$Pr(x|z = n) = \text{Norm}_x[\mu_n, \sigma_n^2]$$

$$\begin{aligned} Pr(x) &= \sum_{n=1}^N Pr(x, z = n) && \text{Marginalization} \\ &= \sum_{n=1}^N Pr(x|z = n) \cdot Pr(z = n) \\ &= \sum_{n=1}^N \lambda_n \cdot \text{Norm}_x[\mu_n, \sigma_n^2]. \end{aligned}$$



Nonlinear Latent Variable Model

Both data \mathbf{x} and latent variable \mathbf{z} are continuous and multivariate

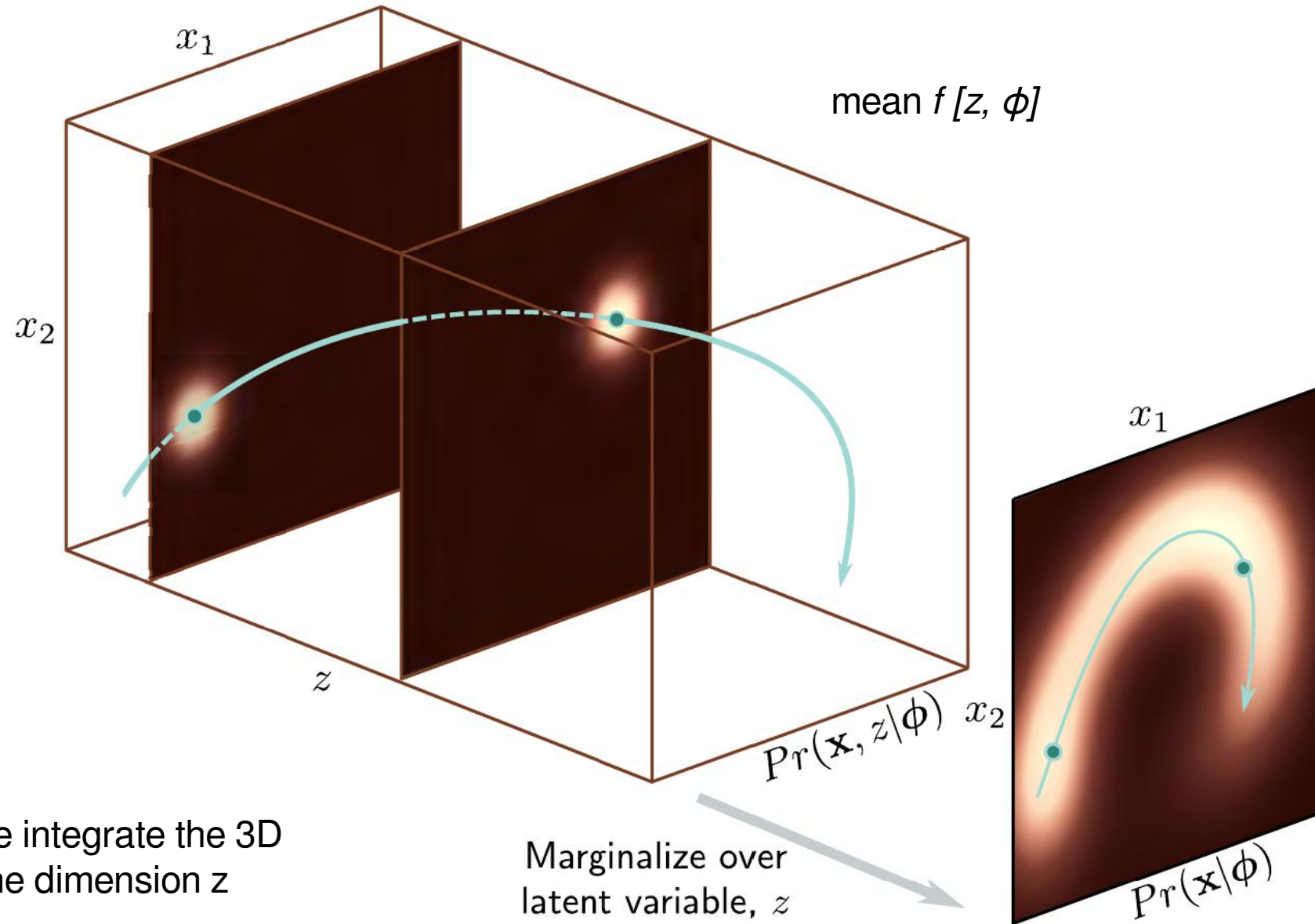
$$Pr(\mathbf{z}) = \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}]$$

likelihood $Pr(\mathbf{x}|\mathbf{z}, \phi) = \text{Norm}_{\mathbf{x}}[\mathbf{f}[\mathbf{z}, \phi], \sigma^2 \mathbf{I}]$ mean is a nonlinear function of the latent variable

$$\begin{aligned} Pr(\mathbf{x}|\phi) &= \int Pr(\mathbf{x}, \mathbf{z}|\phi) d\mathbf{z} \\ &= \int Pr(\mathbf{x}|\mathbf{z}, \phi) \cdot Pr(\mathbf{z}) d\mathbf{z} \\ &= \int \text{Norm}_{\mathbf{x}}[\mathbf{f}[\mathbf{z}, \phi], \sigma^2 \mathbf{I}] \cdot \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}] d\mathbf{z} \end{aligned}$$

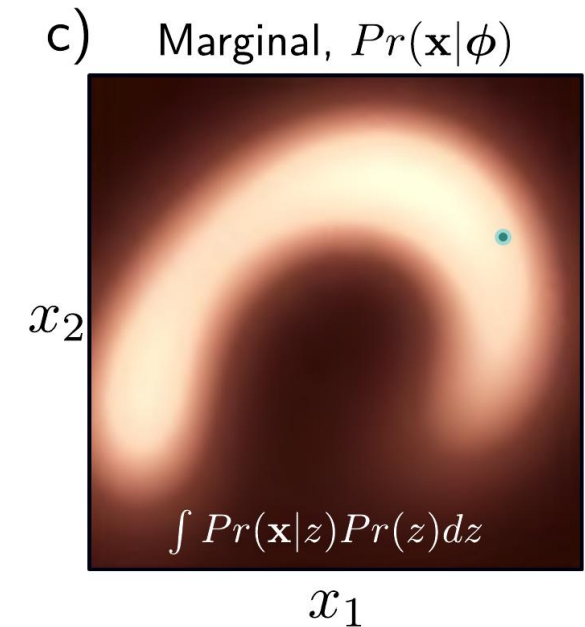
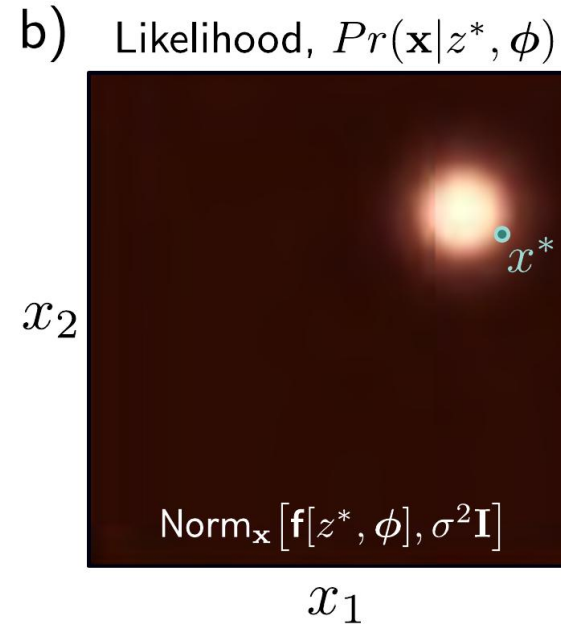
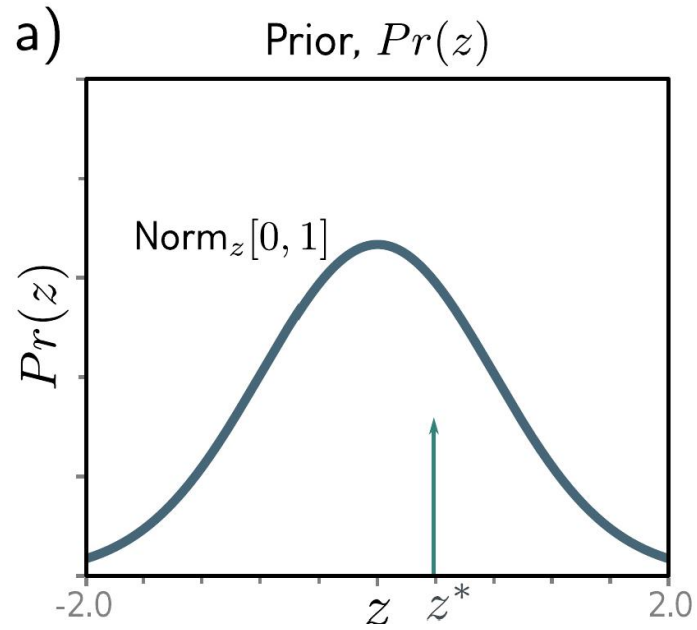
Nonlinear Latent Variable Model

- Both data \mathbf{x} and latent variable \mathbf{z} are continuous and multivariate



Nonlinear Latent Variable Model - Generation from nonlinear latent variable model

- Both data \mathbf{x} and latent variable \mathbf{z} are continuous and multivariate



Nonlinear Latent Variable Model - Training

- Maximize the log-likelihood w.r.t the model parameters over training set
- For simplicity assume the variance term is known
- Goal is to learn ϕ

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log \left[\operatorname{Pr}(\mathbf{x}_i | \phi) \right] \right]$$

$$\operatorname{Pr}(\mathbf{x}_i | \phi) = \int \operatorname{Norm}_{\mathbf{x}_i}[\mathbf{f}[\mathbf{z}, \phi], \sigma^2 \mathbf{I}] \cdot \operatorname{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}] d\mathbf{z}$$

- Unfortunately, this is intractable.
- There is no closed-form expression for the integral and no easy way to evaluate it for a particular value of x

Nonlinear Latent Variable Model - ELBO

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log \left[Pr(\mathbf{x}_i | \phi) \right] \right]$$

- Define lower bound on the log-likelihood
- a function which is always less than or equal to the log-likelihood for a given ϕ and other parameters
- ELBO (Evidence Lower Bound)

Nonlinear Latent Variable Model - ELBO

- ELBO (Evidence Lower Bound) - Derivation

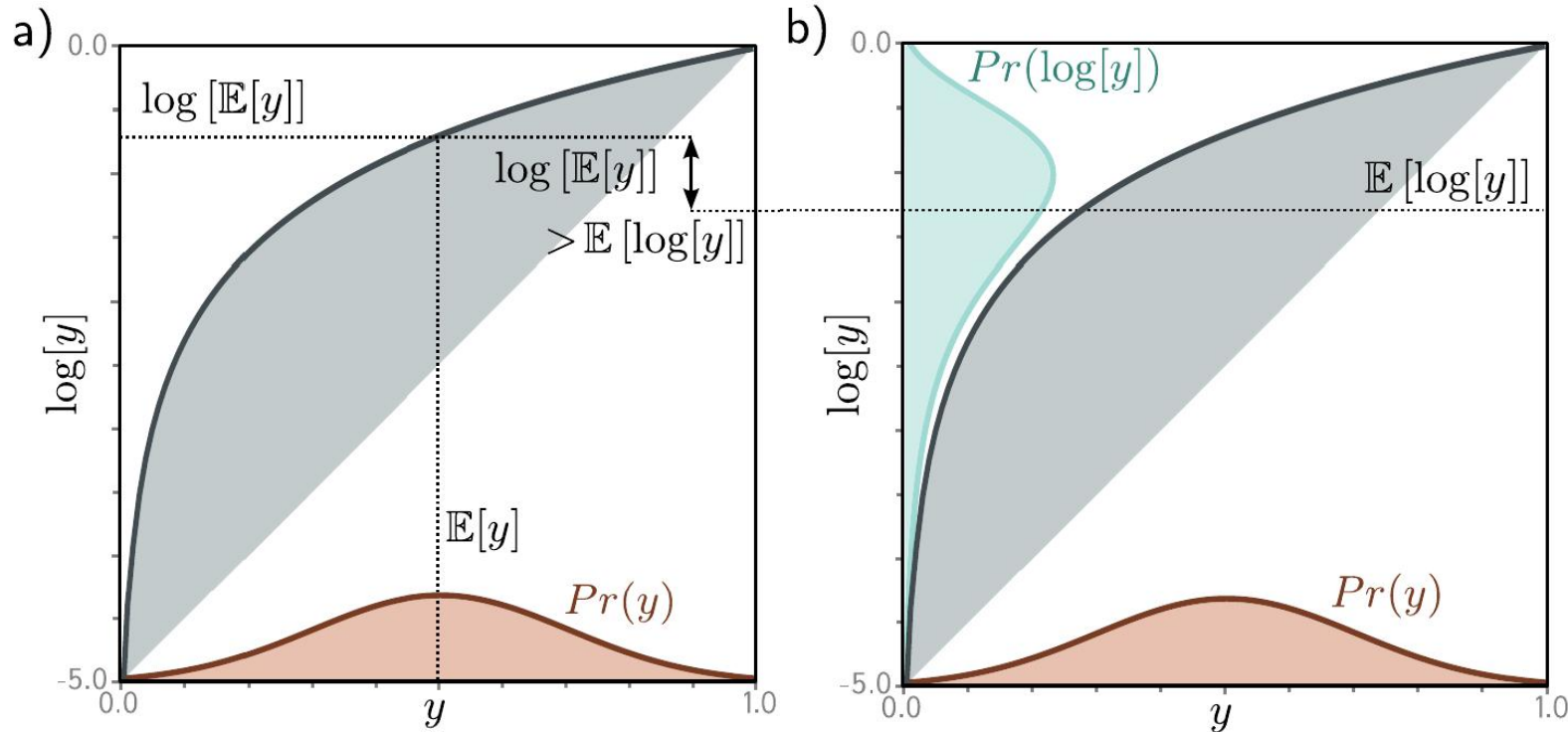
$$g[\mathbb{E}[y]] \geq \mathbb{E}[g[y]]$$

Jensen Inequality
(Concave Function)

$$\log[\mathbb{E}[y]] \geq \mathbb{E}[\log[y]]$$

$$\log \left[\int Pr(y)h[y]dy \right] \geq \int Pr(y) \log[h[y]]dy$$

Expansion



Nonlinear Latent Variable Model - ELBO

- ELBO (Evidence Lower Bound) - Derivation

$$\begin{aligned}\log[Pr(\mathbf{x}|\phi)] &= \log \left[\int Pr(\mathbf{x}, \mathbf{z}|\phi) d\mathbf{z} \right] \\ &= \log \left[\int q(\mathbf{z}) \frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} d\mathbf{z} \right]\end{aligned}$$

multiply & divide by arbitrary probability distribution over latent i.e., $q(\mathbf{z})$

$$\log \left[\int q(\mathbf{z}) \frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} d\mathbf{z} \right] \geq \underbrace{\int q(\mathbf{z}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} \right] d\mathbf{z}}_{\text{Evidence lower bound}}$$

using Jensen's inequality

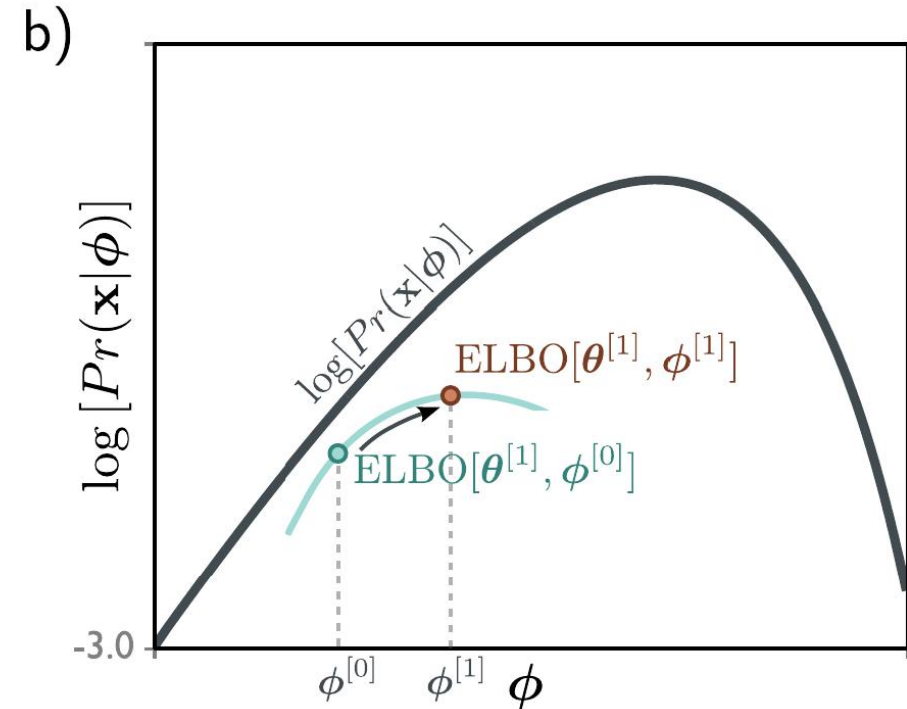
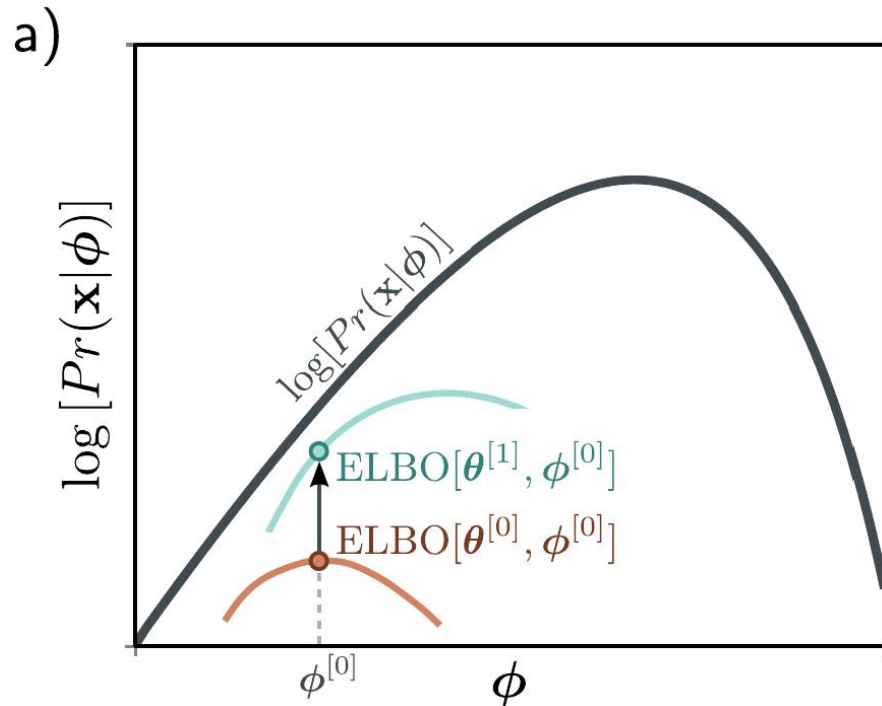
- In practice $q(\mathbf{z})$ is parameterized by θ

$$ELBO[\theta, \phi] = \int q(\mathbf{z}|\theta) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z}$$

maximize w.r.t both θ and ϕ

Nonlinear Latent Variable Model - ELBO

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log [Pr(\mathbf{x}_i | \phi)] \right]$$



- **Objective** : maximize the log-likelihood (black curve) with respect to the parameters ϕ
- For fixed θ , we get a function of ϕ (two colored curves for different values of θ)
- The log-likelihood can be increased by either ways

Nonlinear Latent Variable Model - ELBO

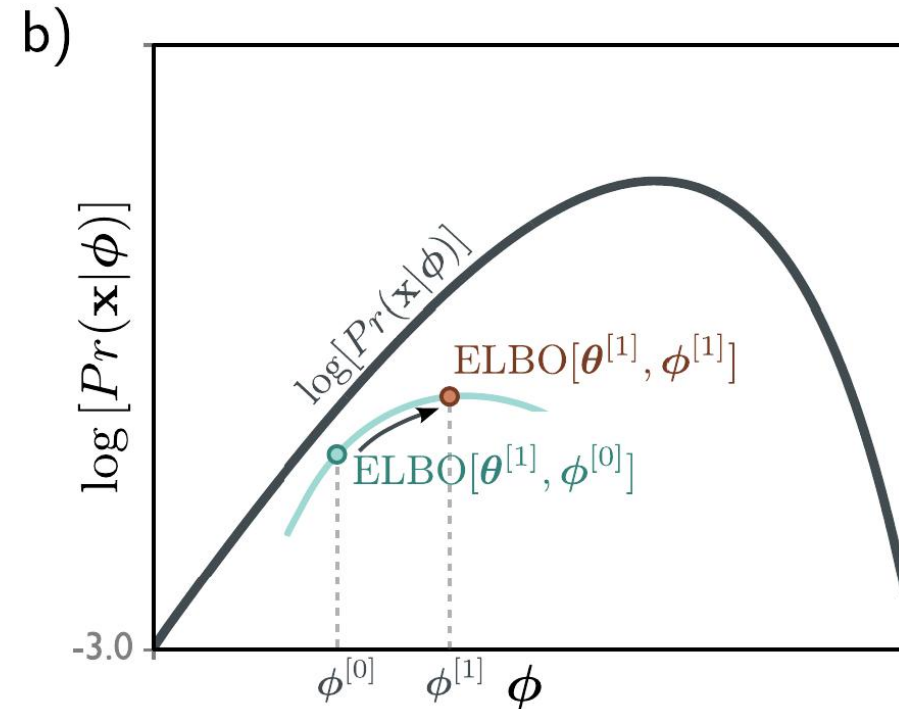
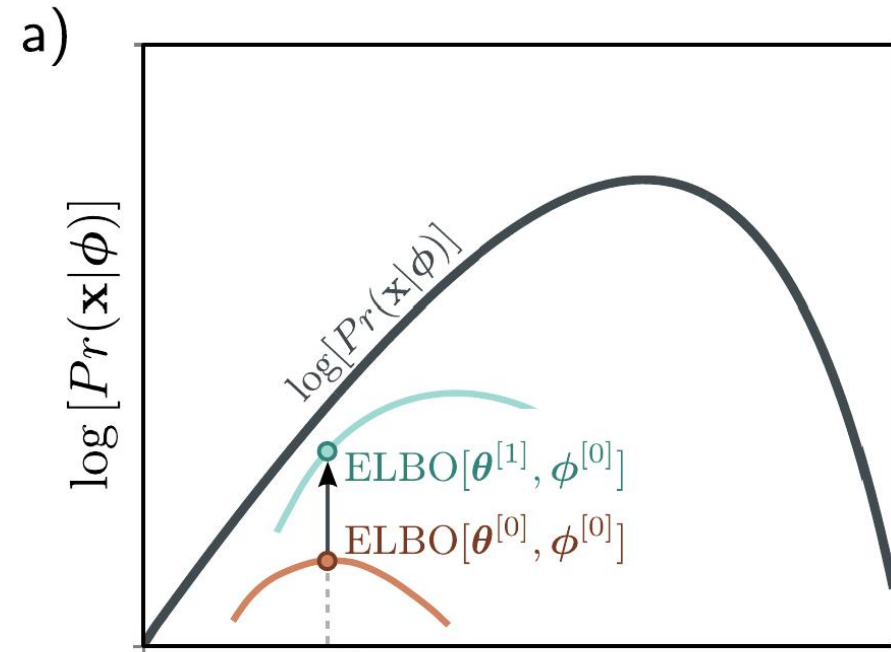
- **Tight ELBO** - ELBO coincide with the likelihood function
 - for a fixed ϕ

$$\begin{aligned}
 \text{ELBO}[\theta, \phi] &= \int q(\mathbf{z}|\theta) \log \left[\frac{\text{Pr}(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\
 &= \int q(\mathbf{z}|\theta) \log \left[\frac{\text{Pr}(\mathbf{z}|\mathbf{x}, \phi) \text{Pr}(\mathbf{x}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \quad \text{conditional probability} \\
 &= \int q(\mathbf{z}|\theta) \log[\text{Pr}(\mathbf{x}|\phi)] d\mathbf{z} + \int q(\mathbf{z}|\theta) \log \left[\frac{\text{Pr}(\mathbf{z}|\mathbf{x}, \phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\
 &= \log[\text{Pr}(\mathbf{x}|\phi)] + \int q(\mathbf{z}|\theta) \log \left[\frac{\text{Pr}(\mathbf{z}|\mathbf{x}, \phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\
 &= \log[\text{Pr}(\mathbf{x}|\phi)] - D_{KL} \left[q(\mathbf{z}|\theta) \parallel \boxed{\text{Pr}(\mathbf{z}|\mathbf{x}, \phi)} \right]. \quad \text{posterior distribution}
 \end{aligned}$$

the bound is tight, when
KL distance is zero i.e.,

$$q(\mathbf{z}|\theta) = \text{Pr}(\mathbf{z}|\mathbf{x}, \phi)$$

Kullback-Leibler (KL) Divergence
(Measures the distance between
two distributions)



Nonlinear Latent Variable Model - ELBO

- ELBO as reconstruction loss - KL distance to prior

$$\begin{aligned} \text{ELBO}[\boldsymbol{\theta}, \boldsymbol{\phi}] &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{\text{Pr}(\mathbf{x}, \mathbf{z}|\boldsymbol{\phi})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\ &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{\text{Pr}(\mathbf{x}|\mathbf{z}, \boldsymbol{\phi}) \text{Pr}(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \quad \text{conditional probability} \\ &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log [\text{Pr}(\mathbf{x}|\mathbf{z}, \boldsymbol{\phi})] d\mathbf{z} + \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{\text{Pr}(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\ &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log [\text{Pr}(\mathbf{x}|\mathbf{z}, \boldsymbol{\phi})] d\mathbf{z} - \underbrace{D_{KL} [q(\mathbf{z}|\boldsymbol{\theta}) || \text{Pr}(\mathbf{z})]}_{\text{degree of agreement between auxiliary distribution and prior}}, \end{aligned}$$

reconstruction accuracy

This formulation is generally used in the **variational autoencoder**

Nonlinear Latent Variable Model - ELBO

- Tight ELBO - ELBO coincide with the likelihood function
 - for a fixed ϕ

$$q(\mathbf{z}|\boldsymbol{\theta}) = Pr(\mathbf{z}|\mathbf{x}, \phi)$$

$$Pr(\mathbf{z}|\mathbf{x}, \phi) = \frac{Pr(\mathbf{x}|\mathbf{z}, \phi)Pr(\mathbf{z})}{Pr(\mathbf{x}|\phi)}$$

Baye's rule

intractable, we can't evaluate this

multivariate normal distribution
mean μ and diagonal variance Σ

Solution ?

- **Solution** : make a variational approximation
- Choose a **simple auxiliary distribution** $q(\mathbf{z} | \boldsymbol{\theta})$ to approximate **true posterior**

during training the goal is to find values of μ and Σ

such that the normal distribution is close to the true posterior $Pr(\mathbf{z} | \mathbf{x})$

} = minimize the KL divergence

$$q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) = \text{Norm}_{\mathbf{z}} \left[\mathbf{g}_{\mu}[\mathbf{x}, \boldsymbol{\theta}], \mathbf{g}_{\Sigma}[\mathbf{x}, \boldsymbol{\theta}] \right]$$

Neural network with params $\boldsymbol{\theta}$

The Variational Autoencoder

$$\text{ELBO}[\boldsymbol{\theta}, \phi] = \underbrace{\int q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \log [Pr(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z}}_{\text{Intractable Integral}} - D_{KL} [q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})]$$

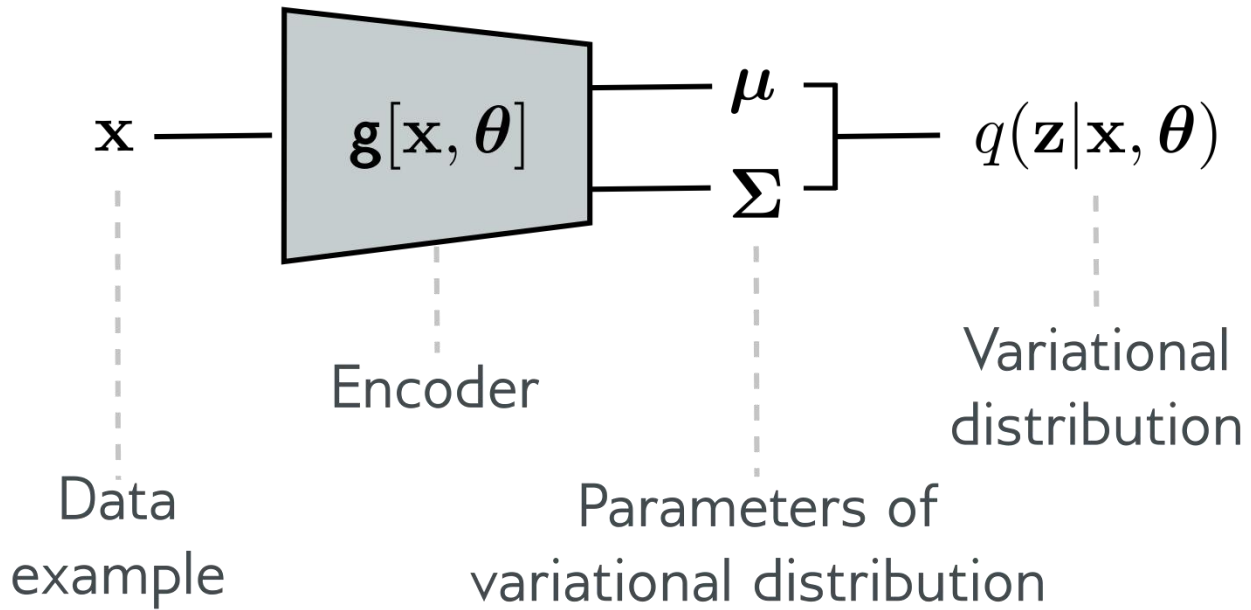
Solution : approximate it by sampling i.e., take a Monte Carlo estimate

$$\text{ELBO}[\boldsymbol{\theta}, \phi] \approx \log [Pr(\mathbf{x}|\mathbf{z}^*, \phi)] - D_{KL} [q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})]$$

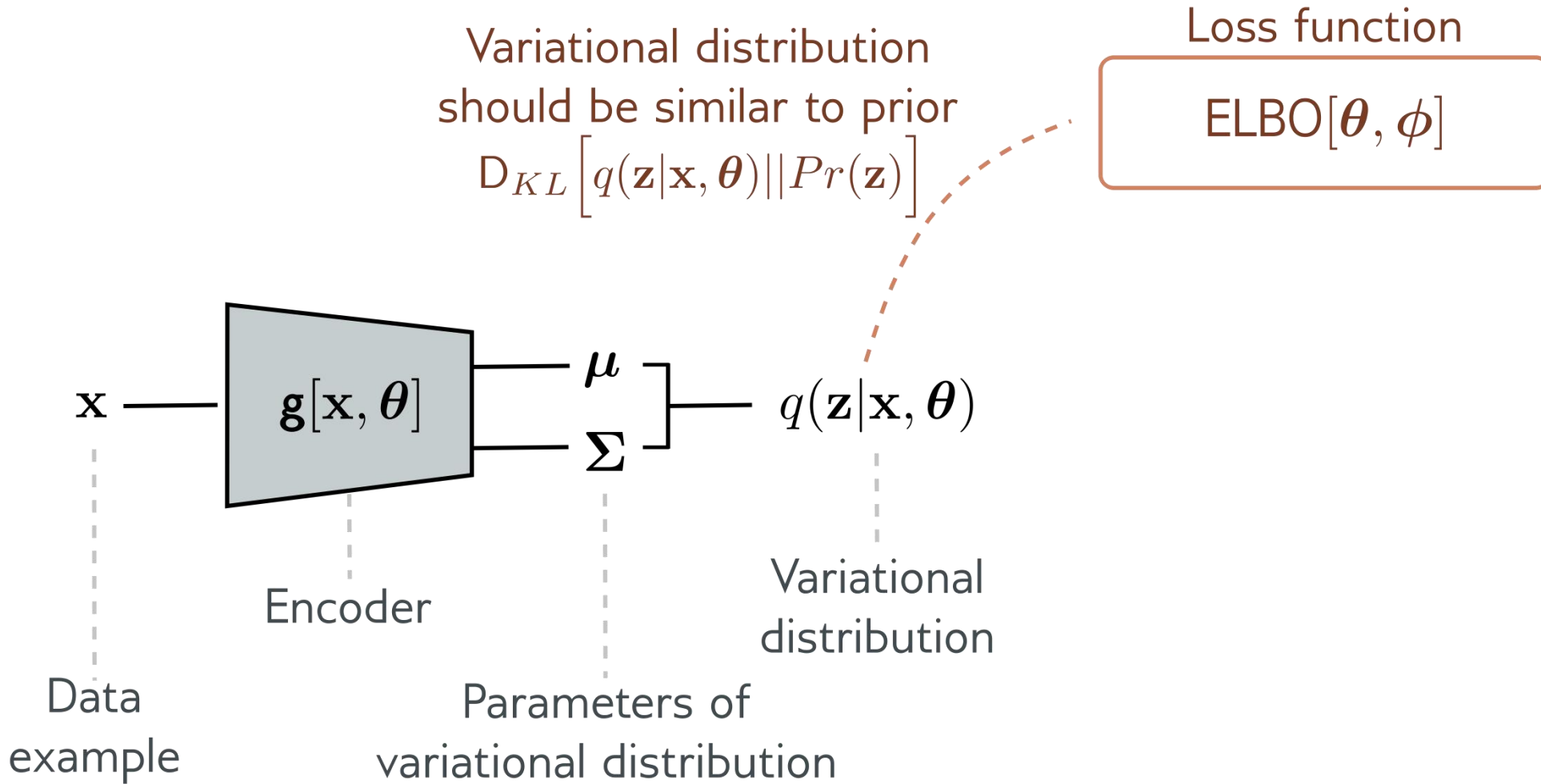
$$D_{KL} [q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})] = \frac{1}{2} \left(\text{Tr}[\boldsymbol{\Sigma}] + \boldsymbol{\mu}^T \boldsymbol{\mu} - \boxed{D_{\mathbf{z}}} - \log [\det[\boldsymbol{\Sigma}]] \right)$$

Dimensionality of the latent space

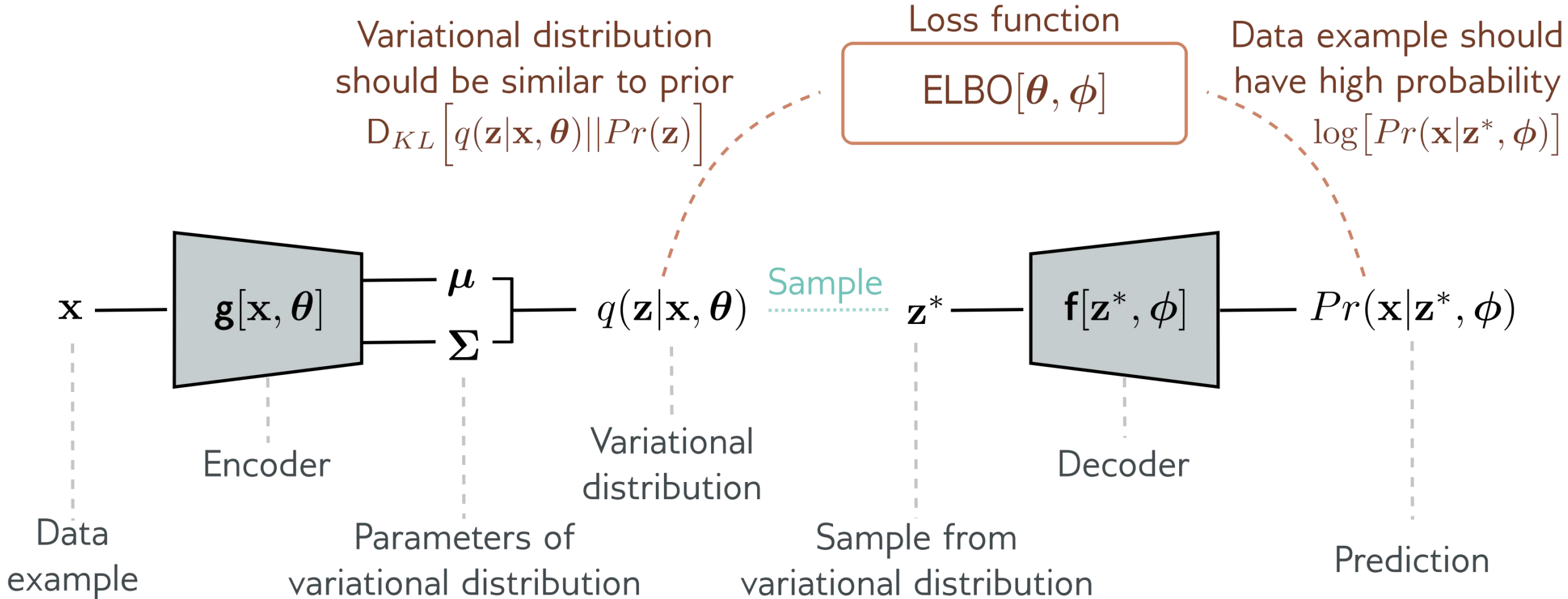
The Variational Autoencoder



The Variational Autoencoder



The Variational Autoencoder



Vector Quantised Variational Autoencoder (VAE)

The VQVAE - Vector Quantized VAE

The Vector Quantised- Variational AutoEncoder (VQ-VAE), differs from VAEs in two key ways:

- the encoder network outputs discrete, rather than continuous codes
- the prior is learnt rather than static

Neural Discrete Representation Learning

Aaron van den Oord
DeepMind
avdnoord@google.com

Oriol Vinyals
DeepMind
vinyals@google.com

Koray Kavukcuoglu
DeepMind
korayk@google.com

The VQVAE - Vector Quantized VAE

Key Idea :

- Incorporate ideas from vector quantisation (VQ)
- VQ method allows model to circumvent issues of “posterior collapse”
- where the latents are ignored when they are paired with a powerful autoregressive decoder =
- typically observed in the VAE framework.

VAE setup

- Posterior Distribution $q(z|x)$
- Prior Distribution $p(z)$
- Decoder with distribution $p(x|z)$

Posterior and prior are assumed normally distributed

The VQVAE - Vector Quantized VAE

VQ-VAE:

- Discrete latent variables
- New way of training leveraging the concepts from vector quantisation (VQ)
- Posterior and prior distributions are categorical

The VQVAE - Formulation

Discrete Latent Variables

Define a latent embedding space $e \in R^{K \times D}$
 K : Size of the discrete latent space (K-way categorical)

D is the dimensionality of each latent embedding vector

$$e_i \in R^D, i \in 1, 2, \dots, K$$

The encoder takes an input and produce output $z_e(x)$

Discrete latent variable is then calculated by a nearest neighbour look-up using a shared embedding space

$$q(z = k|x) = \begin{cases} 1 & \text{for } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2 \\ 0 & \text{otherwise} \end{cases}$$

The VQVAE - Formulation

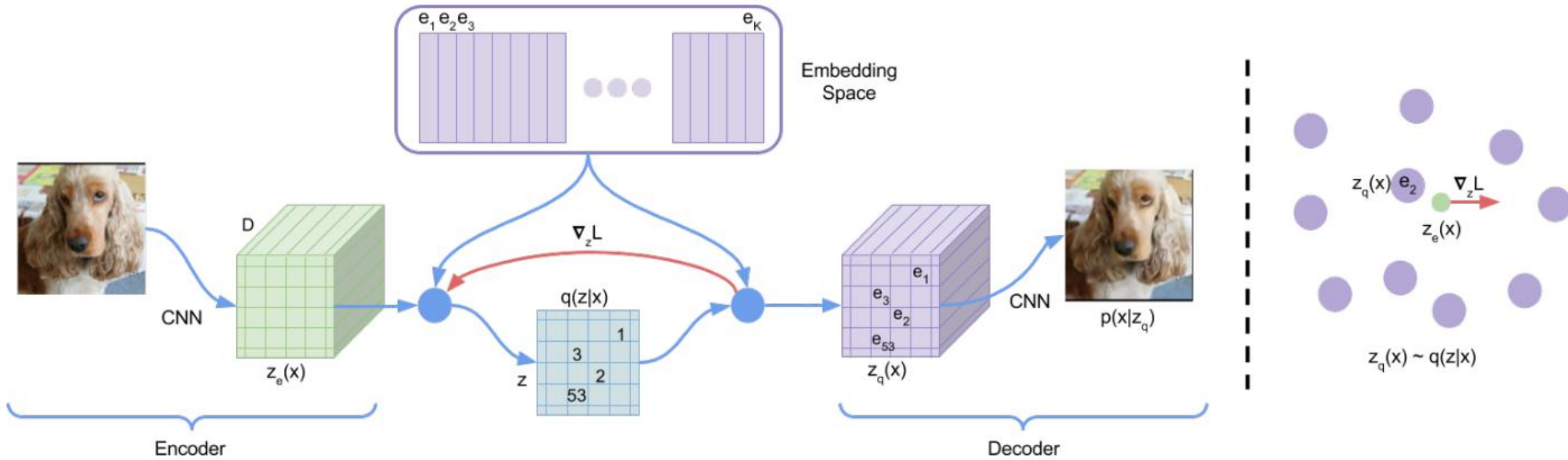
the posterior is a categorical distribution as one-hot

$$q(z = k|x) = \begin{cases} 1 & \text{for } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2 \\ 0 & \text{otherwise} \end{cases}$$

Decoder input is the embedding vector calculated from NN search

$$z_q(x) = e_k, \quad \text{where } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2$$

The VQVAE - Training



The VQVAE - Training

There is no real gradient defined for the below equation,

$$z_q(x) = e_k, \quad \text{where } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2$$

however we approximate the gradient similar to the straight-through estimator

- copy gradients from decoder input to encoder output .

$$z_q(x) \quad z_e(x)$$

- During forward computation the nearest embedding is passed to the decoder,
- During the backwards pass the gradient is passed unaltered to the encoder

Will this work ? Why?

- the output representation of the encoder and the input to the decoder share the same D dimensional space,
- the gradients contain useful information for how the encoder has to change its output to lower the reconstruction loss.

The VQVAE - Training

- the gradient must push the encoder's output to be discretised differently in the next forward pass, because the assignment will be different
- Total Loss

$$L = \log p(x|z_q(x)) + \|\text{sg}[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - \text{sg}[e]\|_2^2$$

The VQVAE - Training

- the gradient must push the encoder's output to be discretised differently in the next forward pass, because the assignment will be different

- Total Loss

reconstruction loss

Due to the straight-through gradient estimation of mapping from $z_e(x)$ to $z_q(x)$, the embeddings e_i receive no gradients from the reconstruction loss

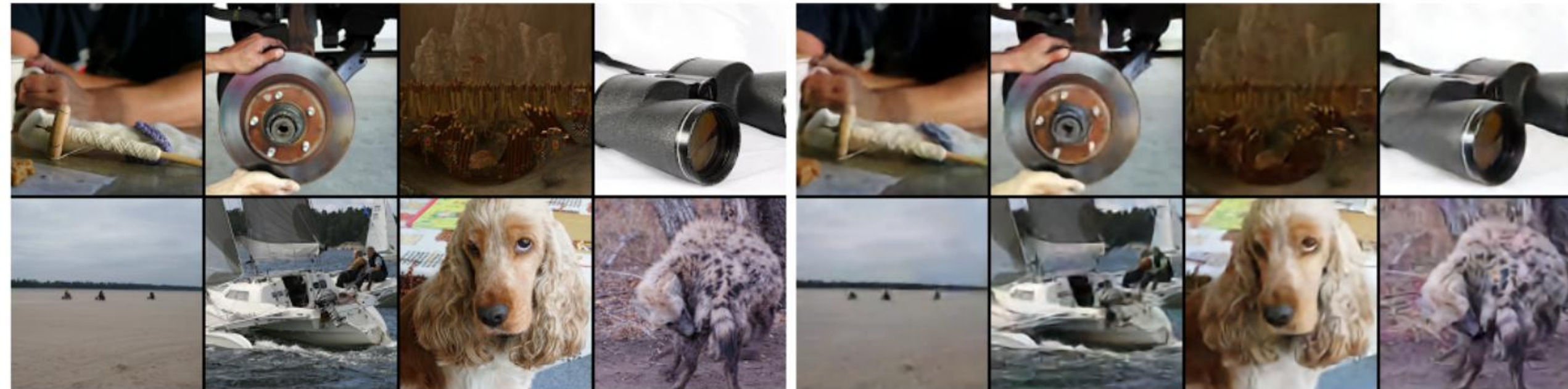
$$L = \log p(x|z_q(x)) + \|\text{sg}[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - \text{sg}[e]\|_2^2$$

VQ objective uses the l_2 error to move the embedding vectors e_i towards the encoder outputs $z_e(x)$

Commitment Loss

Volume of the embedding space can grow arbitrarily. To make sure the encoder commits to an embedding and its output does not grow.

The VQVAE - Results



Left: ImageNet 128x128x3 images, right: reconstructions from a VQ-VAE with a 32x32x1 latent space, with K=512.