# So, how does it work?

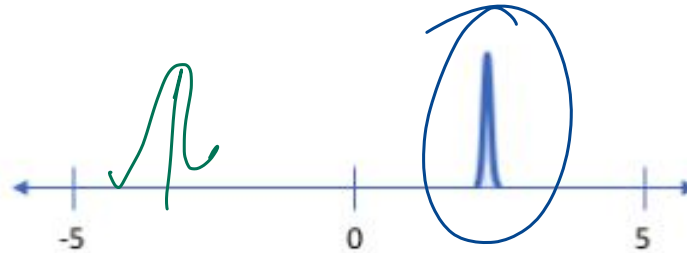*Encode* *Decode*

Penalizing reconstruction loss encourages the distribution to describe the input
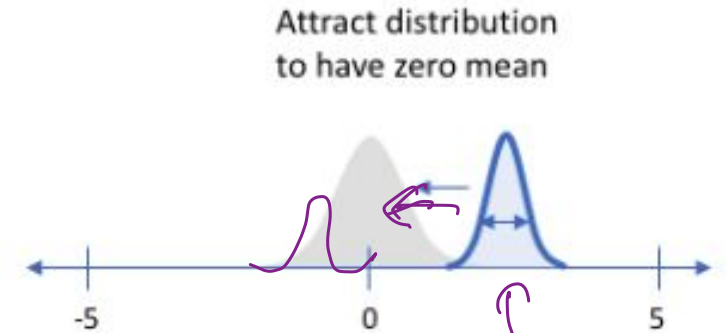
Without regularization, our network can "cheat" by learning narrow distributions

Penalizing KL divergence acts as a regularizing force

Attract distribution to have zero mean



Our distribution deviates from the prior to describe some characteristic of the data

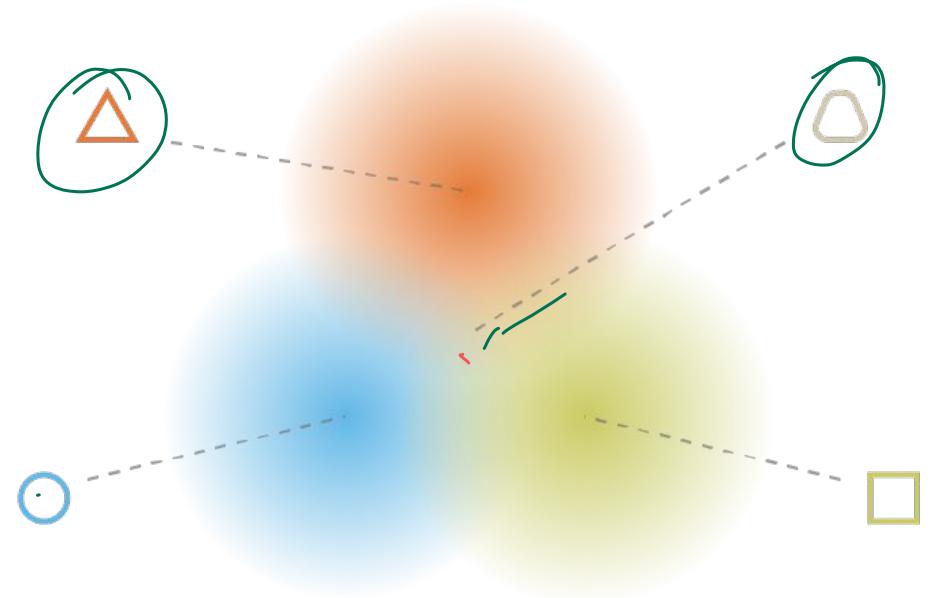With a small enough variance, this distribution is effectively only representing a single value

Ensure sufficient variance to yield a smooth latent space

# How does KL divergence help?



what can happen without regularisation ❌          ✔ what we want to obtain with regularisation

# Continuity and completeness



A point of the latent space that would be halfway between the means of two encoded distributions coming from different training data should be decoded in something that is somewhere between the data that gave the first distribution and the data that gave the second distribution as it may be sampled by the autoencoder in both cases.
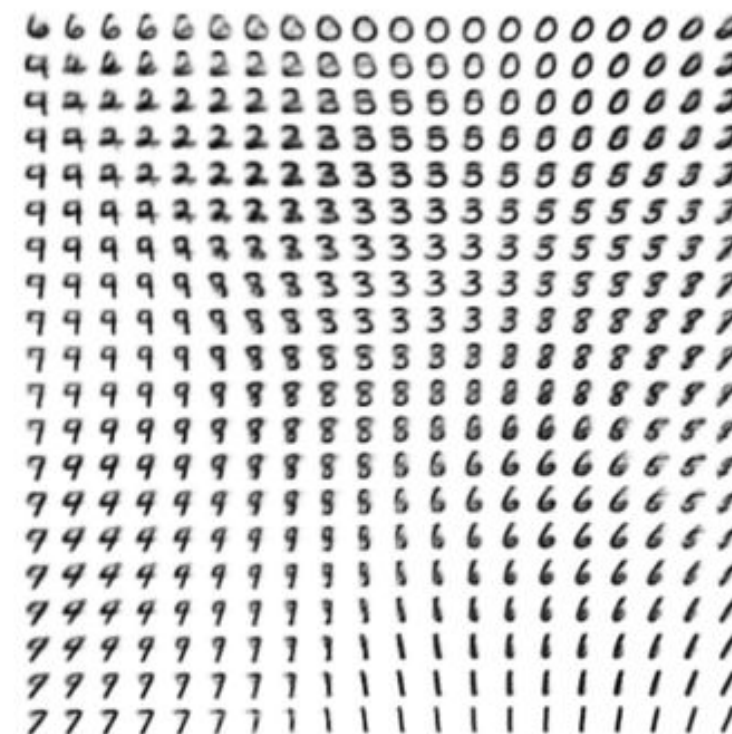
# Visualization of the learned data manifold

**Auto-Encoding Variational Bayes**

**Diederik P. Kingma**
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

**Max Welling**
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

(a) Learned Frey Face manifold

(b) Learned MNIST manifold

Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables $z$. For each of these values $z$, we plotted the corresponding generative $p_\theta(x|z)$ with the learned parameters $\theta$.
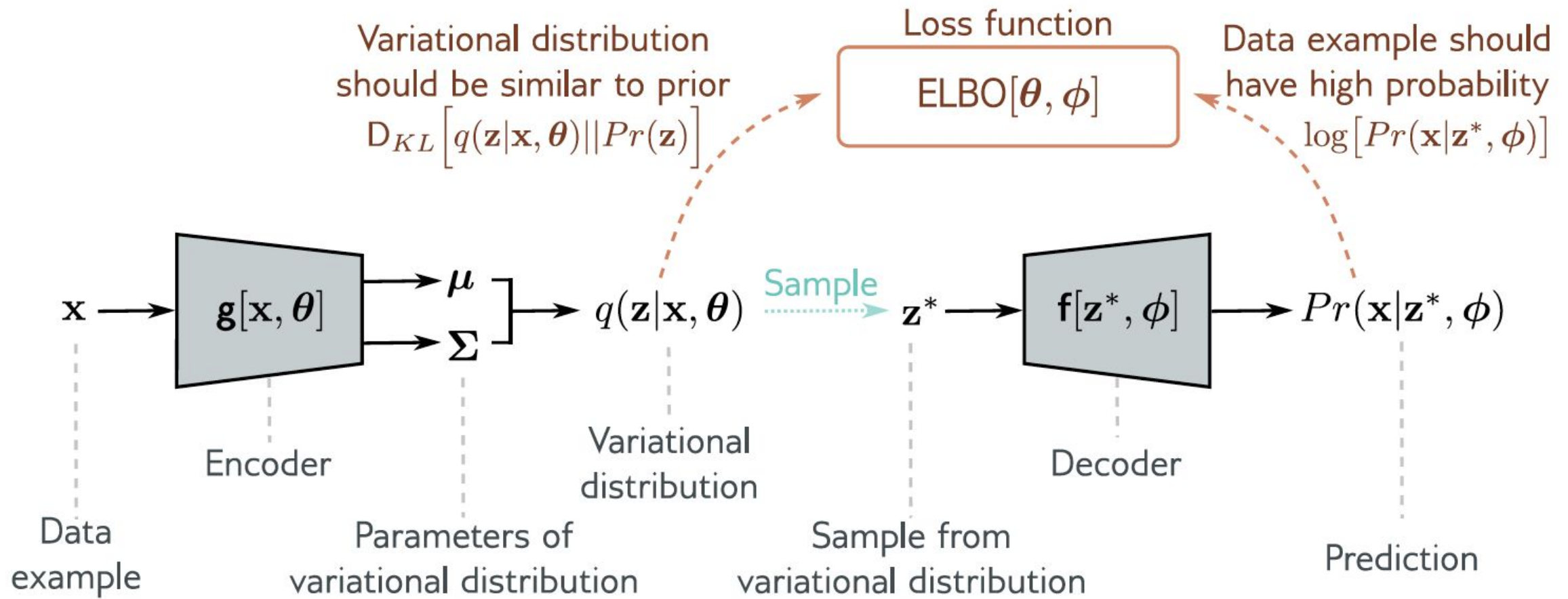
# Maths behind VAEs

**Figure 17.9** Variational autoencoder. The encoder $\mathbf{g}[\mathbf{x}, \boldsymbol{\theta}]$ takes a training example $\mathbf{x}$ and predicts the parameters $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ of the variational distribution $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$. We sample from this distribution and then use the decoder $\mathbf{f}[\mathbf{z}, \boldsymbol{\phi}]$ to predict the data $\mathbf{x}$. The loss function is the negative ELBO, which depends on how accurate this prediction is and how similar the variational distribution $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$ is to the prior $Pr(\mathbf{z})$ (equation 17.21).

# Using Deep Networks

# Deriving the bound

**Objective:** maximize log [*Pr(x|ϕ)*]

Which can be written by marginalizing over z

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

We now multiply and divide the log-likelihood by an arbitrary probability distribution q(z) over the latent variables:

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

Jensen's inequality says that a concave function g[•] of the expectation of data y is greater than or equal to the expectation of the function of the data

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

# Graphical Intuition: Jensen's inequality

**Figure 17.4** Jensen's inequality (discrete case). The logarithm (black curve) is a concave function; you can draw a straight line between any two points on the curve, and this line will always lie underneath it. It follows that any convex combination (weighted sum with positive weights that sum to one) of the six points on the log function must lie in the gray region under the curve. Here, we have weighted the points equally (i.e., taken the mean) to yield the cyan point. Since this point lies below the curve, $\log[\mathbb{E}[y]] > \mathbb{E}[\log[y]]$.

# Graphical Intuition: Jensen's inequality

# Deriving the bound

Jensen's inequality says that a concave function g[•] of the expectation of data y is greater than or equal to the expectation of the function of the data
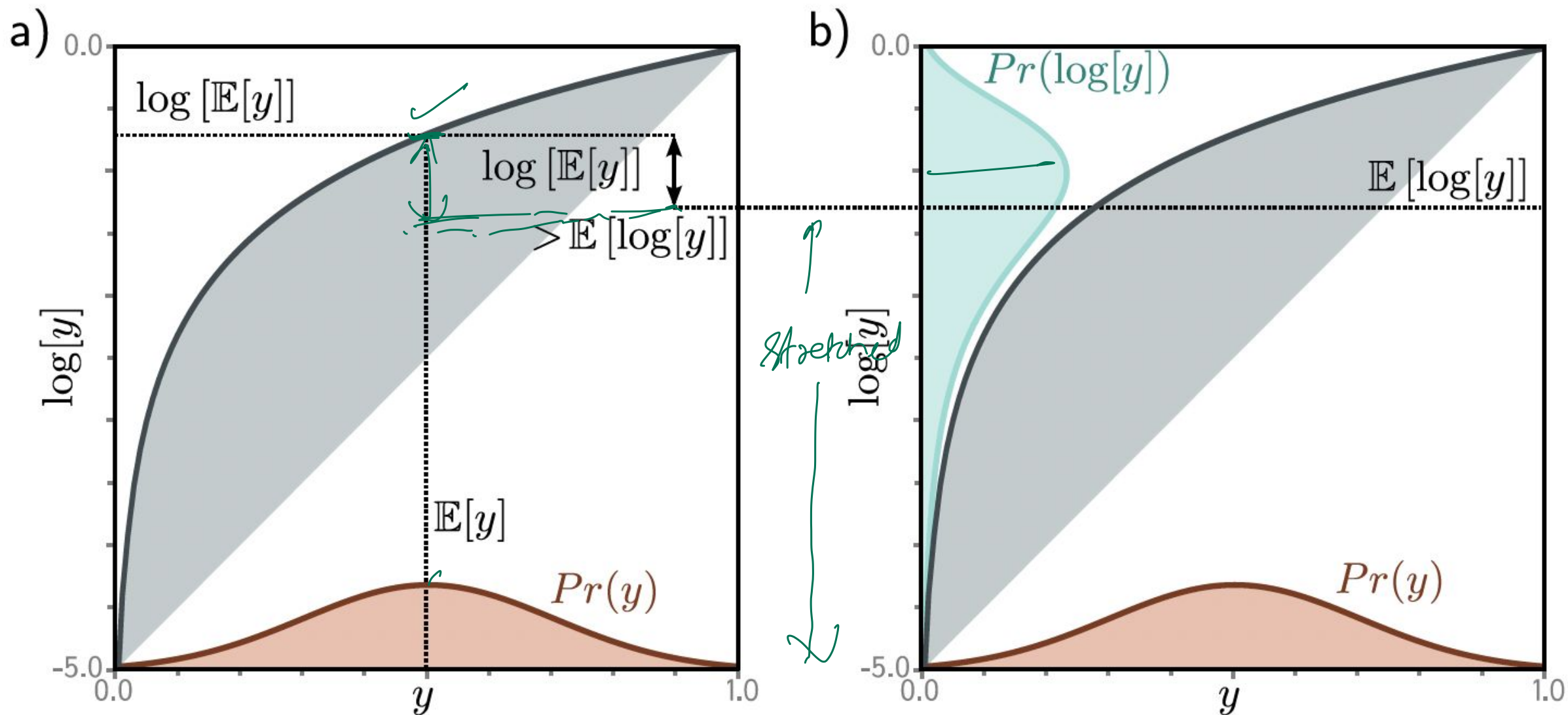
$$\log\left[\int Pr(y)y\,dy\right] \geq \int Pr(y)\log[y]\,dy.$$

*In fact, the slightly more general statement is true:*

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

where h[y] is a function of y. This follows because h[y] is another random variable with a new distribution. Since we never specified Pr(y), the relation remains true.

# Deriving the bound

Using Jensen's inequality

$$\log\left[\int q(\mathbf{z})\frac{Pr(\mathbf{x},\mathbf{z}|\phi)}{q(\mathbf{z})}d\mathbf{z}\right] \geq \int q(\mathbf{z})\log\left[\frac{Pr(\mathbf{x},\mathbf{z}|\phi)}{q(\mathbf{z})}\right]d\mathbf{z}$$ 
*Evidence Lower Bound (ELBO)*

In practice, the distribution $q(z)$ has parameters $\theta$, so the ELBO can be written as:

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

# ELBO as reconstruction loss - KL divergence

$$
\begin{aligned}
\text{ELBO}[\boldsymbol{\theta}, \phi] &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[ \frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[ \frac{Pr(\mathbf{x}|\mathbf{z}, \phi) Pr(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[ Pr(\mathbf{x}|\mathbf{z}, \phi) \right] d\mathbf{z} + \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[ \frac{Pr(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[ Pr(\mathbf{x}|\mathbf{z}, \phi) \right] d\mathbf{z} - D_{KL} \left[ q(\mathbf{z}|\boldsymbol{\theta}) \middle\| Pr(\mathbf{z}) \right],
\end{aligned}
$$

For our network: *q(z|θ) = q(z|x, θ)*

$$
\text{ELBO}[\boldsymbol{\theta}, \phi] = \int q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \log \left[ Pr(\mathbf{x}|\mathbf{z}, \phi) \right] d\mathbf{z} - D_{KL} \left[ q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \middle\| Pr(\mathbf{z}) \right]
$$

# Approximating ELBO

$$\text{ELBO}[\boldsymbol{\theta}, \boldsymbol{\phi}] = \int q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \log\left[Pr(\mathbf{x}|\mathbf{z}, \boldsymbol{\phi})\right] d\mathbf{z} - \text{D}_{KL}\left[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \big|\big| Pr(\mathbf{z})\right]$$

*Still intractable*

*Can you think of a trick to approximate the first term?*

$$\mathbb{E}_{\mathbf{z}}\left[\text{a}[\mathbf{z}]\right] = \int \text{a}[\mathbf{z}] q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) d\mathbf{z} \approx \frac{1}{N}\sum_{n=1}^{N} \text{a}[\mathbf{z}_n^*]$$

*Monte Carlo Markov Sampling*

For a very approximate estimate, we can just use a single sample $\mathbf{z}^*$ from $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$

$$\text{ELBO}[\boldsymbol{\theta}, \boldsymbol{\phi}] \approx \log\left[Pr(\mathbf{x}|\mathbf{z}^*, \boldsymbol{\phi})\right] - \text{D}_{KL}\left[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \big|\big| Pr(\mathbf{z})\right]$$

*KL divergence can be calculated in closed form*

# Computing KL divergence

$$\text{ELBO}[\boldsymbol{\theta}, \boldsymbol{\phi}] \approx \log\left[Pr(\mathbf{x}|\mathbf{z}^*, \boldsymbol{\phi})\right] - \text{D}_{KL}\left[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})\middle|\middle|Pr(\mathbf{z})\right]$$
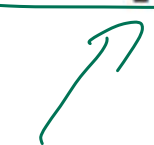
MSE

The second term is the KL divergence between the variational distribution $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) = \text{Norm}_{\mathbf{z}}[\boldsymbol{\mu}, \boldsymbol{\Sigma}]$ and the prior $Pr(\mathbf{z}) = \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}]$. The KL divergence between two normal distributions can be calculated in closed form. For the special case where one distribution has parameters $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ and the other is a standard normal, it is given by:

$$\text{D}_{KL}\left[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})\middle|\middle|Pr(\mathbf{z})\right] = \frac{1}{2}\left(\text{Tr}[\boldsymbol{\Sigma}] + \boldsymbol{\mu}^T\boldsymbol{\mu} - D_{\mathbf{z}} - \log\left[\det[\boldsymbol{\Sigma}]\right]\right). \qquad (17.24)$$

where $D_{\mathbf{z}}$ is the dimensionality of the latent space.

# VAE Algorithm

We aim to build a model that computes the evidence lower bound for a point $x$.
Then we use an optimization algorithm to maximize this lower bound over the dataset and hence improve the log-likelihood.
To compute the ELBO, we:

- compute the mean $\mu$ and variance $\Sigma$ of the variational posterior distribution $q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{x})$ for this data point $\mathbf{x}$ using the network $\mathbf{g}[\mathbf{x}, \boldsymbol{\theta}]$, *Encoder*

- draw a sample $\mathbf{z}^*$ from this distribution, and

- compute the ELBO using equation 17.23.

$$\text{ELBO}[\boldsymbol{\theta}, \phi] \approx \log\left[Pr(\mathbf{x}|\mathbf{z}^*, \phi)\right] - \text{D}_{KL}\left[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})\,\big|\big|\,Pr(\mathbf{z})\right]$$

To maximize this bound, we run mini-batches of samples through the network and update these parameters with an optimization algorithm such as SGD or Adam
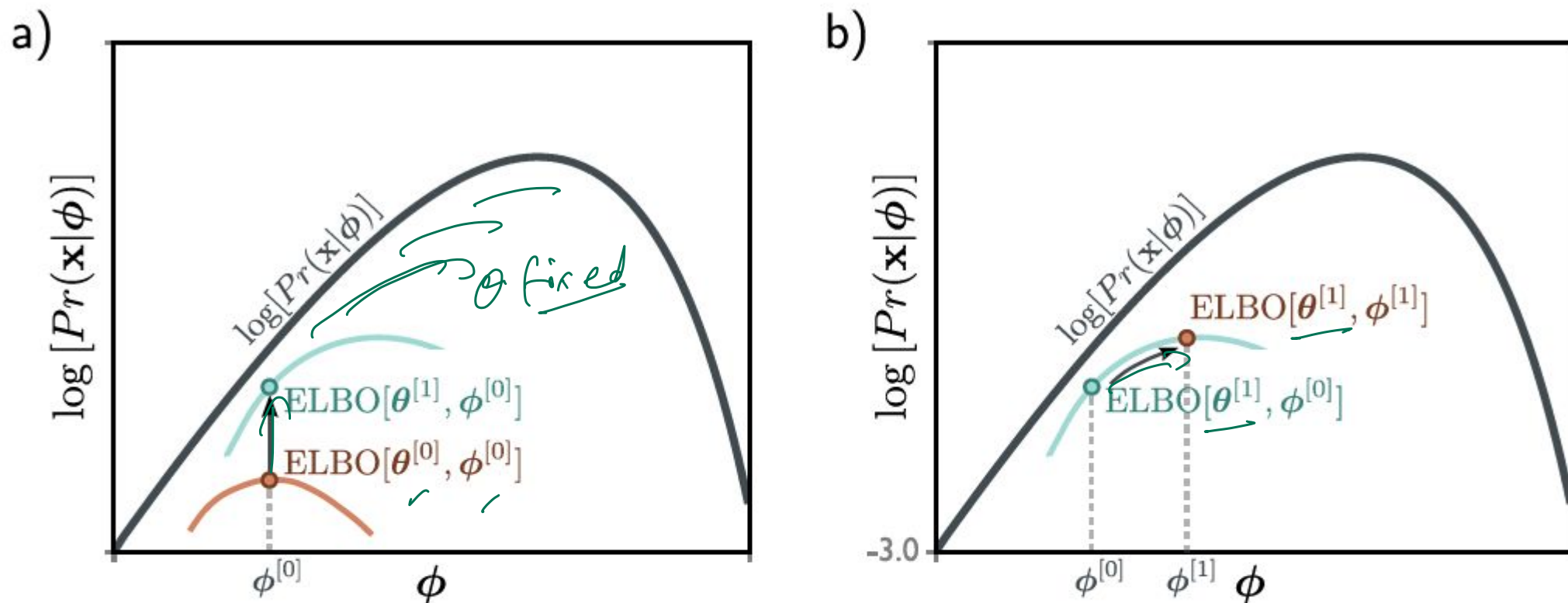
**Figure 17.6** Evidence lower bound (ELBO). The goal is to maximize the log-likelihood $\log[Pr(\mathbf{x}|\phi)]$ (black curve) with respect to the parameters $\phi$. The ELBO is a function that lies everywhere below the log-likelihood. It is a function of both $\phi$ and a second set of parameters $\theta$. For fixed $\theta$, we get a function of $\phi$ (two colored curves for different values of $\theta$). Consequently, we can increase the log-likelihood by either improving the ELBO with respect to a) the new parameters $\theta$ (moving from colored curve to colored curve) or b) the original parameters $\phi$ (moving along the current colored curve).

Variational distribution should be similar to prior
$$D_{KL}\left[q(\mathbf{z}|\mathbf{x},\boldsymbol{\theta})||Pr(\mathbf{z})\right]$$

Loss function
$$\text{ELBO}[\boldsymbol{\theta},\boldsymbol{\phi}]$$

Data example should have high probability
$$\log\left[Pr(\mathbf{x}|\mathbf{z}^*,\boldsymbol{\phi})\right]$$

$\mathbf{x} \longrightarrow \mathbf{g}[\mathbf{x},\boldsymbol{\theta}] \longrightarrow \begin{matrix}\boldsymbol{\mu}\\\boldsymbol{\Sigma}\end{matrix} \longrightarrow q(\mathbf{z}|\mathbf{x},\boldsymbol{\theta}) \xrightarrow{\text{Sample}} \mathbf{z}^* \longrightarrow \mathbf{f}[\mathbf{z}^*,\boldsymbol{\phi}] \longrightarrow Pr(\mathbf{x}|\mathbf{z}^*,\boldsymbol{\phi})$

Encoder

Variational distribution

Decoder

Data example

Parameters of variational distribution

Sample from variational distribution

Prediction

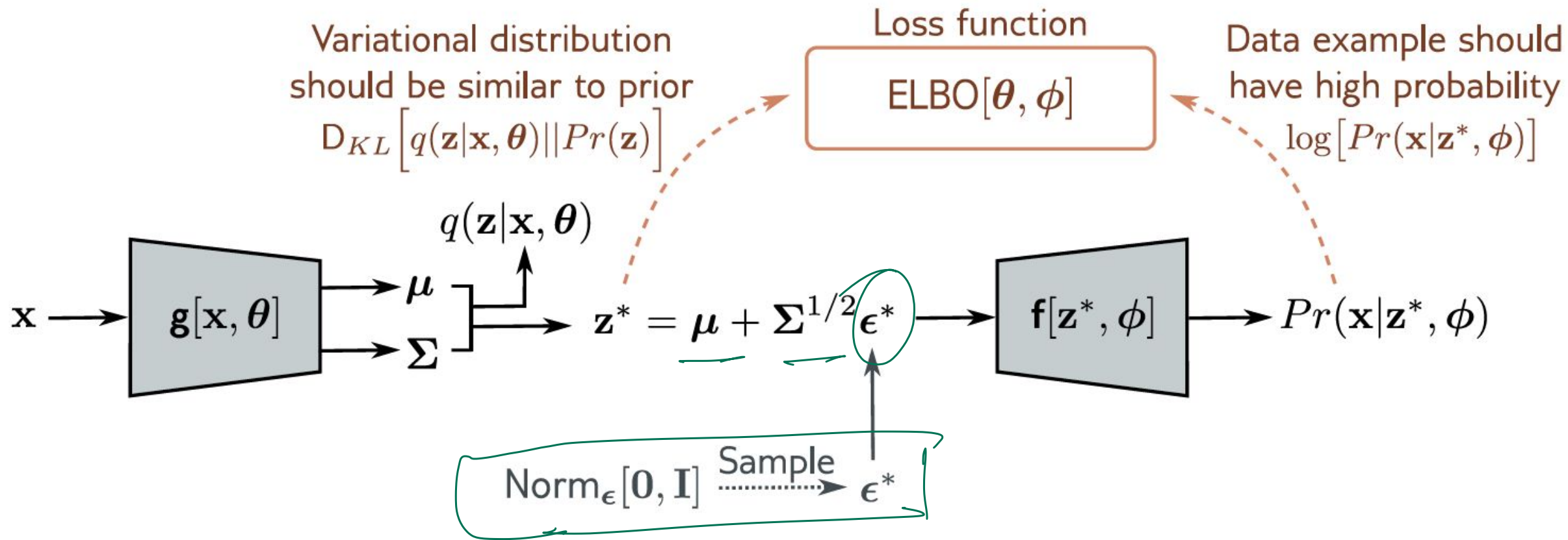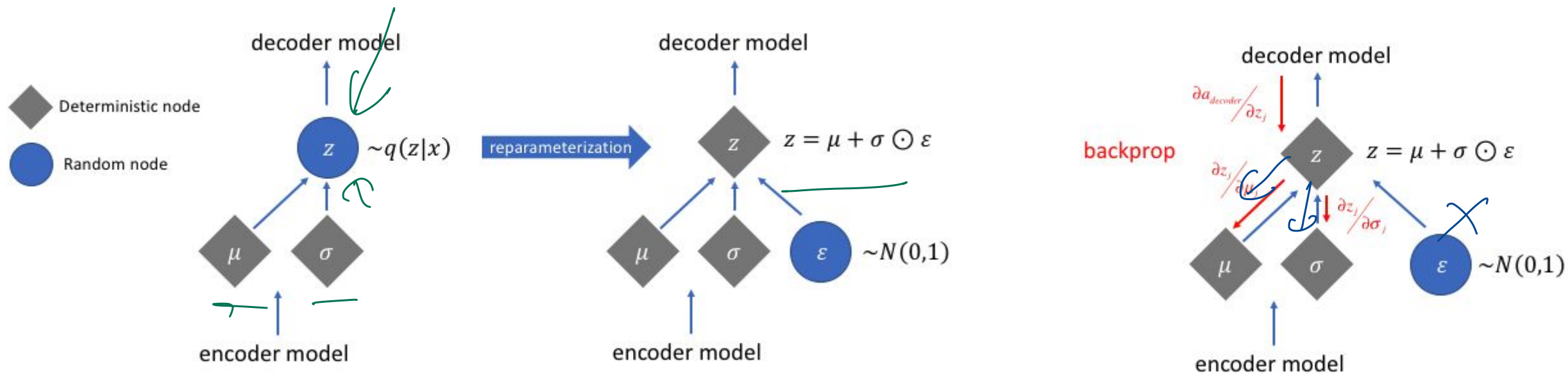*Can you see any problem with the optimization?*

**Figure 17.11** Reparameterization trick. With the original architecture (figure 17.9), we cannot easily backpropagate through the sampling step. The reparameterization trick removes the sampling step from the main pipeline; we draw from a standard normal and combine this with the predicted mean and covariance to get a sample from the variational distribution.

# The reparameterization trick

# Application

A data point x can be projected into the latent space by taking the mean of the distribution predicted by the encoder

Multiple images labeled as "neutral" or "smiling" are projected into latent space. The vector representing this change is estimated by taking the difference in latent space between the means of these two groups. A second vector is estimated to represent "mouth closed" versus "mouth open."



**Figure 17.13** Resynthesis. The original image on the left is projected into the latent space using the encoder, and the mean of the predicted Gaussian is chosen to represent the image. The center-left image in the grid is the reconstruction of the input. The other images are reconstructions after manipulating the latent space in directions representing smiling/neutral (horizontal) and mouth open/closed (vertical). Adapted from White (2016).