

# Deep Learning

Lecture 11: Auto-encoders and variational auto-encoders

Prof. Gilles Louppe  
[g.louppe@uliege.be](mailto:g.louppe@uliege.be)

# Today

Learn a model of the data.

- Auto-encoders
- Variational inference
- Variational auto-encoders



*"The brain has about  $10^{14}$  synapses and we only live for about  $10^9$  seconds. So we have a lot more parameters than data. This motivates the idea that we must do a lot of unsupervised learning since the perceptual input (including proprioception) is the only place we can get  $10^5$  dimensions of constraint per second."*

Geoffrey Hinton, 2014.



#### ■ "Pure" Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

#### ■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

#### ■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**

■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)



*"We need tremendous amount of information to build machines that have common sense and generalize."*

Yann LeCun, 2016.

## Deep unsupervised learning

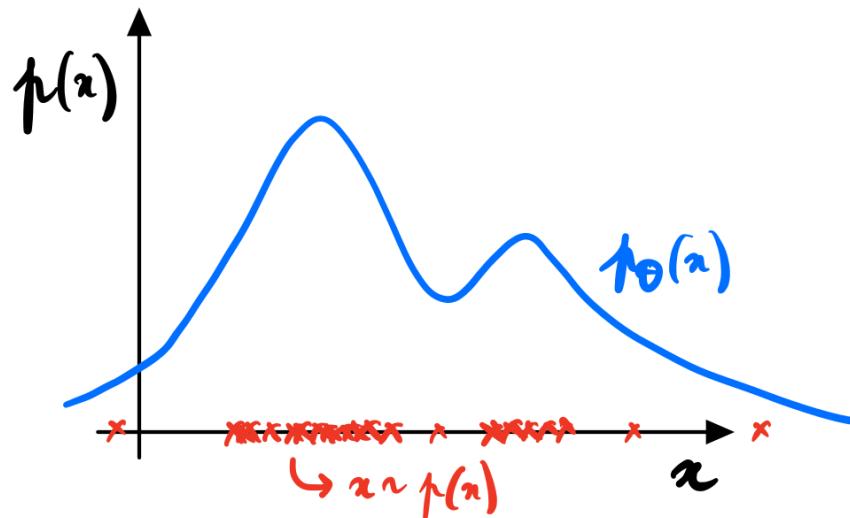
Deep unsupervised learning is about learning a model of the data, explicitly or implicitly, without requiring labels.

- **Generative models**: recreate the raw data distribution (e.g., the distribution of natural images).
- **Self-supervised learning**: solve puzzle tasks that require semantic understanding (e.g., predict a missing word in a sequence).

## Generative models

A (deep) **generative model** is a probabilistic model  $p_{\theta}$  that can be used as a simulator of the data.

Formally, a generative model defines a probability distribution  $p_{\theta}(\mathbf{x})$  over the data  $\mathbf{x} \in \mathcal{X}$ , where the parameters  $\theta$  are learned to match the (unknown) data distribution  $p(\mathbf{x})$ .





Variational auto-encoders  
(Kingma and Welling, 2013)



Diffusion models  
(Midjourney, 2023)



# What can we do with generative models?

Produce samples

$$\mathbf{x} \sim p(\mathbf{x}|\theta)$$

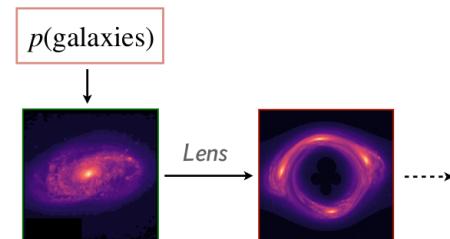
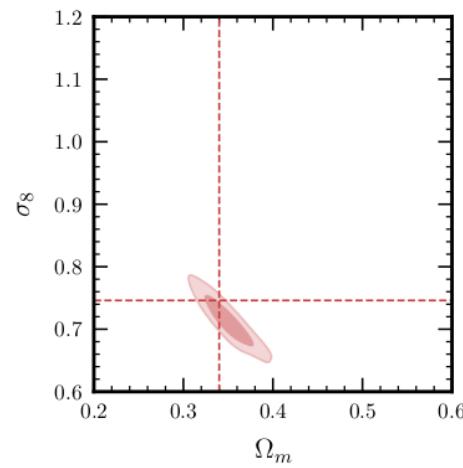
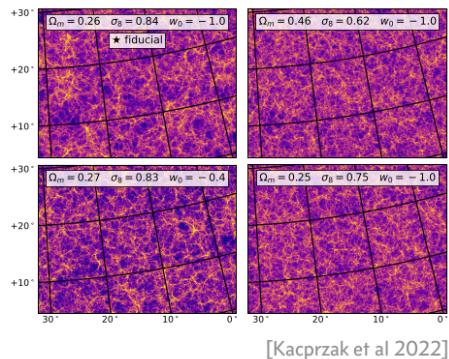
Evaluate densities

$$p(\mathbf{x}|\theta)$$

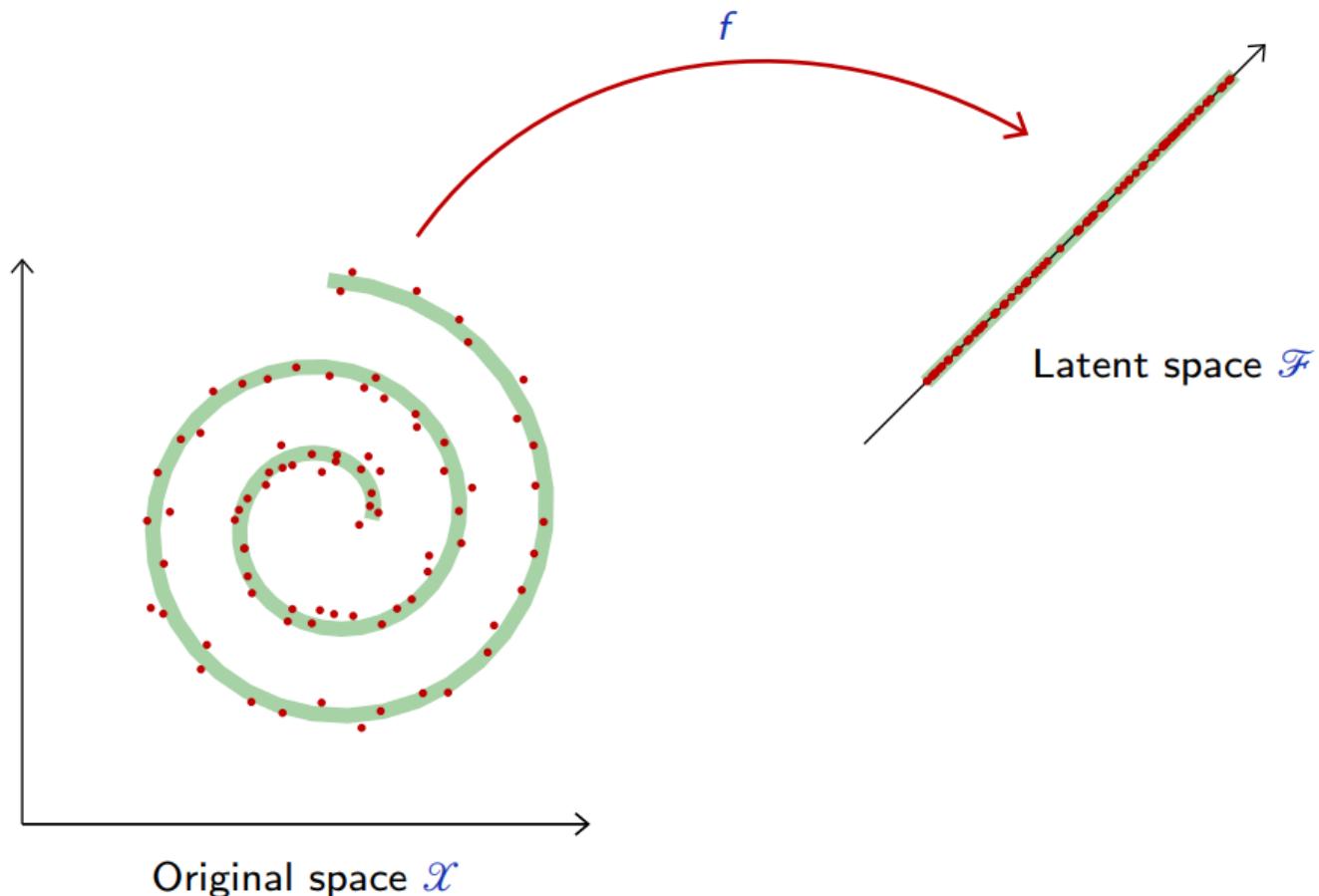
Encode complex priors

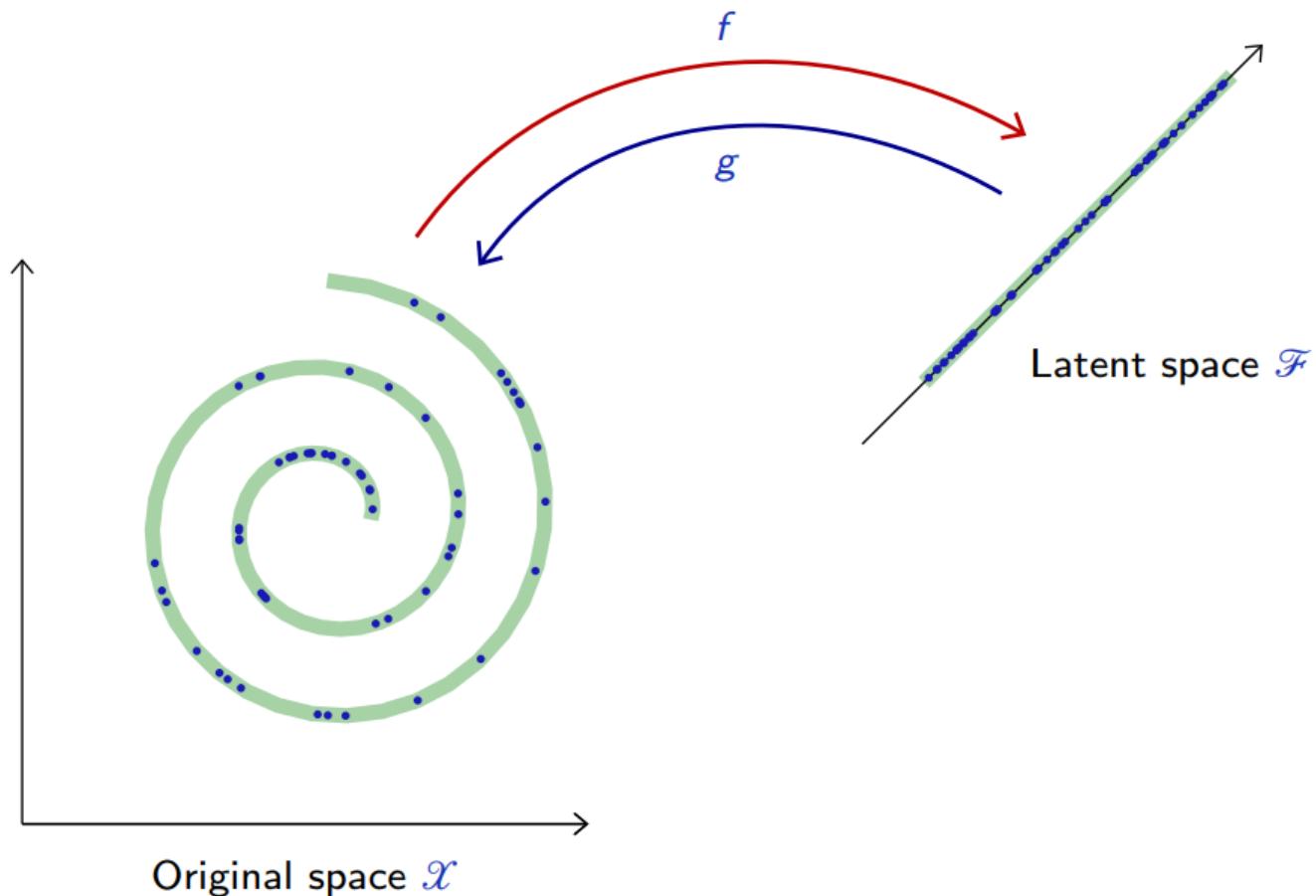
$$p(\mathbf{x})$$

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$



# Auto-encoders



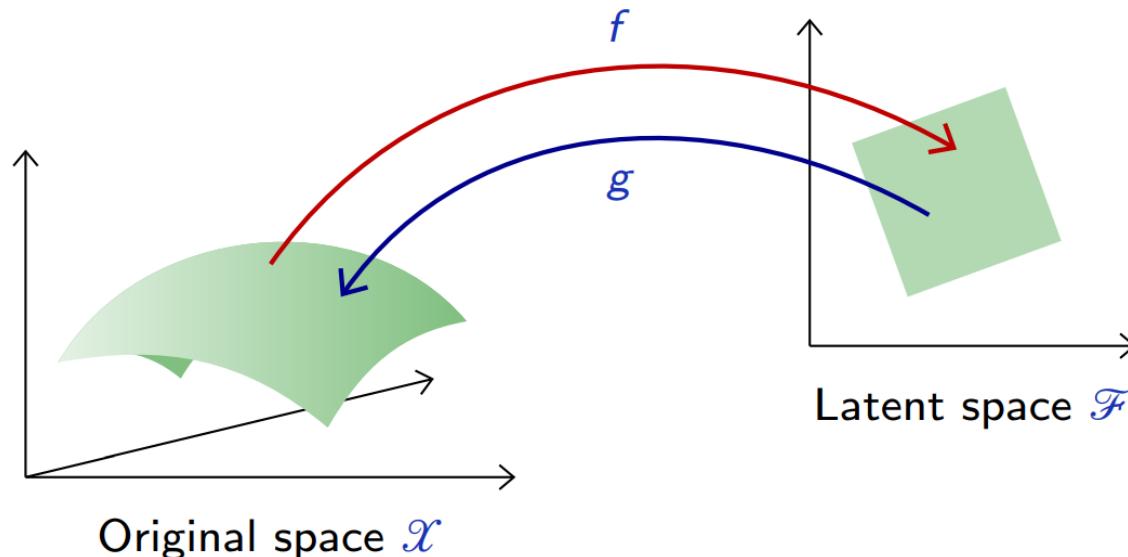


# Auto-encoders

An auto-encoder is a composite function made of

- an **encoder**  $f$  from the original space  $\mathcal{X}$  to a latent space  $\mathcal{Z}$ ,
- a **decoder**  $g$  to map back to  $\mathcal{X}$ ,

such that  $g \circ f$  is close to the identity on the data.



Let  $p(\mathbf{x})$  be the data distribution over  $\mathcal{X}$ . A good auto-encoder could be characterized with the reconstruction loss

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [||\mathbf{x} - g \circ f(\mathbf{x})||^2] \approx 0.$$

Given two parameterized mappings  $f(\cdot; \theta_f)$  and  $g(\cdot; \theta_g)$ , training consists of minimizing an empirical estimate of that loss,

$$\theta_f, \theta_g = \arg \min_{\theta_f, \theta_g} \frac{1}{N} \sum_{i=1}^N ||\mathbf{x}_i - g(f(\mathbf{x}_i, \theta_f), \theta_g)||^2.$$

For example, when the auto-encoder is linear,

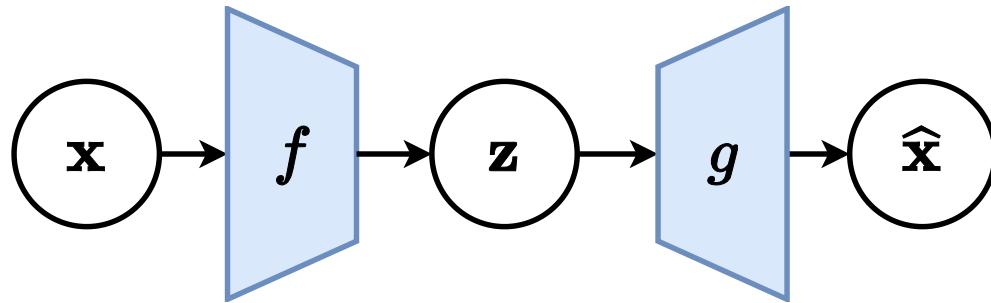
$$\begin{aligned}f &: \mathbf{z} = \mathbf{U}^T \mathbf{x} \\g &: \hat{\mathbf{x}} = \mathbf{U} \mathbf{z},\end{aligned}$$

with  $\mathbf{U} \in \mathbb{R}^{p \times d}$ , the reconstruction error reduces to

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [||\mathbf{x} - \mathbf{U}\mathbf{U}^T \mathbf{x}||^2].$$

In this case, an optimal solution is given by PCA.

## Deep auto-encoders



Better results can be achieved with more sophisticated classes of mappings than linear projections: use deep neural networks for  $f$  and  $g$ .

For instance,

- by combining a multi-layer perceptron encoder  $f : \mathbb{R}^p \rightarrow \mathbb{R}^d$  with a multi-layer perceptron decoder  $g : \mathbb{R}^d \rightarrow \mathbb{R}^p$ .
- by combining a convolutional network encoder  $f : \mathbb{R}^{w \times h \times c} \rightarrow \mathbb{R}^d$  with a decoder  $g : \mathbb{R}^d \rightarrow \mathbb{R}^{w \times h \times c}$  composed of the reciprocal transposed convolutional layers.

$\textcolor{blue}{X}$  (original samples)

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(\textcolor{blue}{X})$  (CNN,  $d = 2$ )

7 2 1 0 9 1 9 9 6 9 0 6  
9 0 1 5 9 7 5 9 9 6 6 5  
9 0 7 9 0 1 3 1 3 6 7 2

$g \circ f(\textcolor{blue}{X})$  (PCA,  $d = 2$ )

9 3 1 0 9 1 9 9 0 9 0 0  
9 0 1 8 9 9 8 9 9 0 9 0  
9 0 9 9 0 1 8 1 3 0 9 0

$\textcolor{blue}{X}$  (original samples)

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(\textcolor{blue}{X})$  (CNN,  $d = 8$ )

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(\textcolor{blue}{X})$  (PCA,  $d = 8$ )

7 3 1 0 4 1 3 9 0 7 0 0  
9 0 1 0 9 7 3 4 7 6 0 5  
4 0 7 4 0 1 3 1 3 0 7 0

$\textcolor{blue}{X}$  (original samples)

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(\textcolor{blue}{X})$  (CNN,  $d = 32$ )

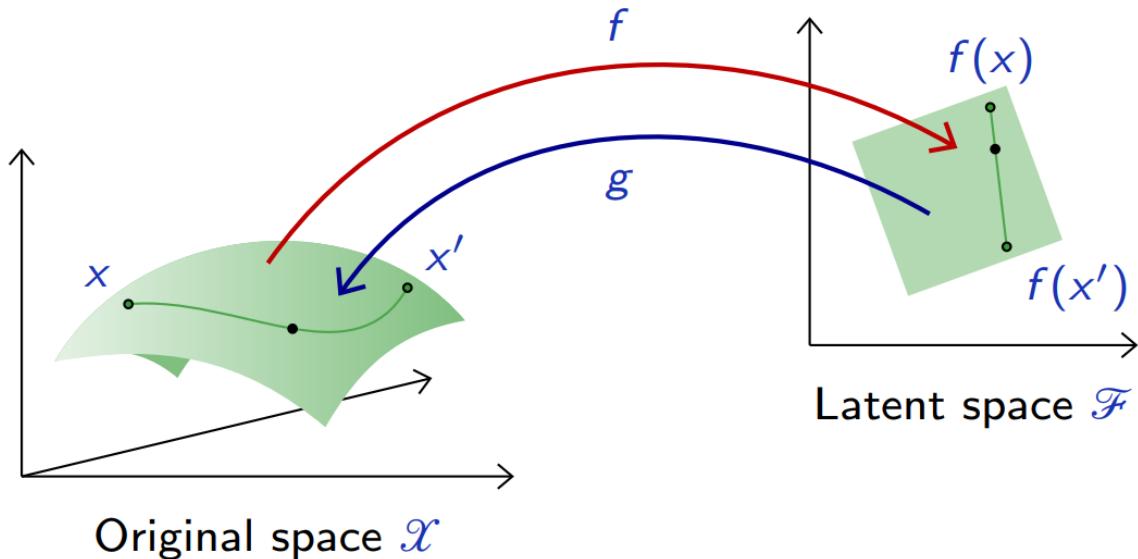
7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(\textcolor{blue}{X})$  (PCA,  $d = 32$ )

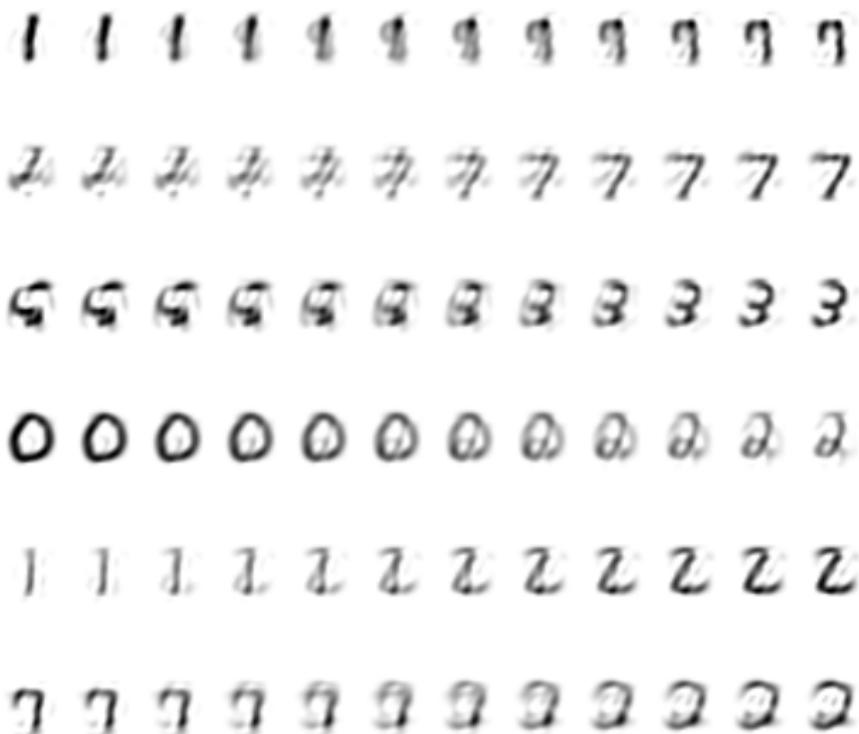
7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

# Interpolation

To get an intuition of the learned latent representation, we can pick two samples  $\mathbf{x}$  and  $\mathbf{x}'$  at random and interpolate samples along the line in the latent space.



PCA interpolation ( $d = 32$ )



Autoencoder interpolation ( $d = 32$ )

0 0 0 0 0 0 3 3 3 3 3 3

7 7 7 7 7 7 5 5 5 5 5 5

4 4 4 4 4 4 2 2 2 2 2 2

7 7 7 7 7 7 7 7 7 7 7 7

2 2 2 2 2 4 4 4 4 4 4 4

4 4 4 4 4 7 7 7 7 7 7 7

# Denoising auto-encoders

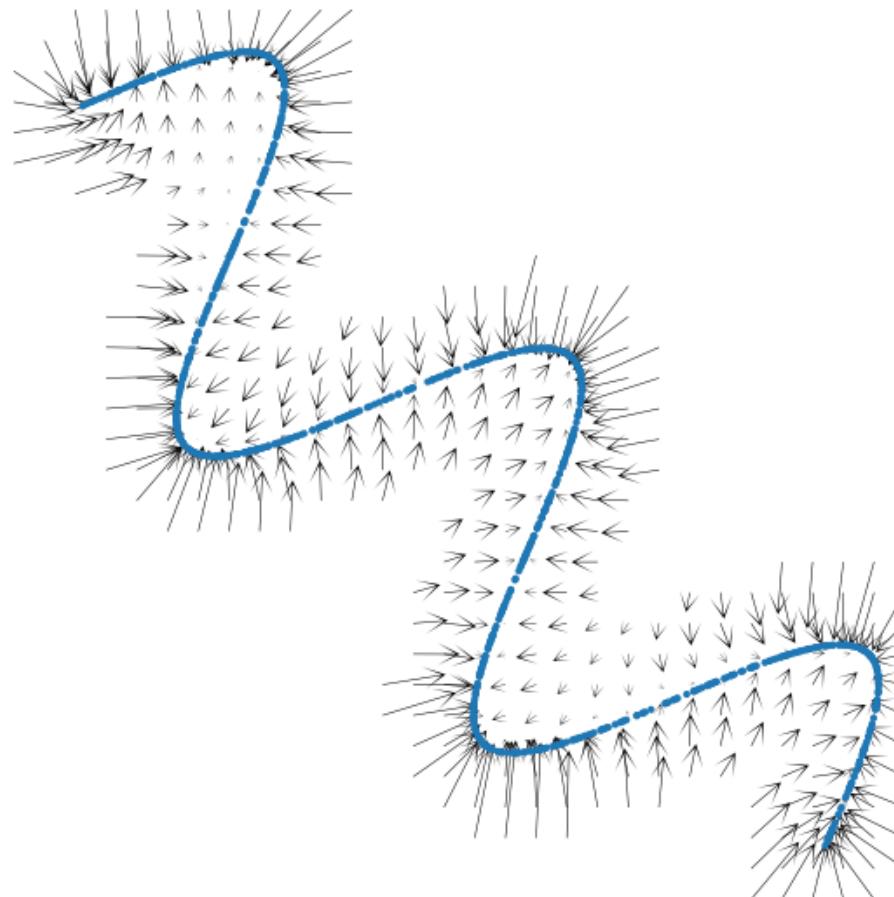
Besides dimension reduction, auto-encoders can capture dependencies between signal components to restore degraded or noisy signals. In this case, the composition

$$h = g \circ f : \mathcal{X} \rightarrow \mathcal{X}$$

is a **denoising** auto-encoder.

The goal is to optimize  $h$  such that a perturbation  $\tilde{\mathbf{x}}$  of the signal  $\mathbf{x}$  is restored to  $\mathbf{x}$ , hence

$$h(\tilde{\mathbf{x}}) \approx \mathbf{x}.$$



Original

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 8 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ( $p = 0.5$ )

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 8 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 8 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

Original

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 4 5  
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ( $p = 0.9$ )

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 4 5  
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

7 2 1 0 4 1 4 3 4 9 0 6  
9 0 1 5 9 7 3 4 9 6 4 5  
4 0 7 4 0 1 3 1 3 4 7 2

Original

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ( $\sigma = 4$ )

7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

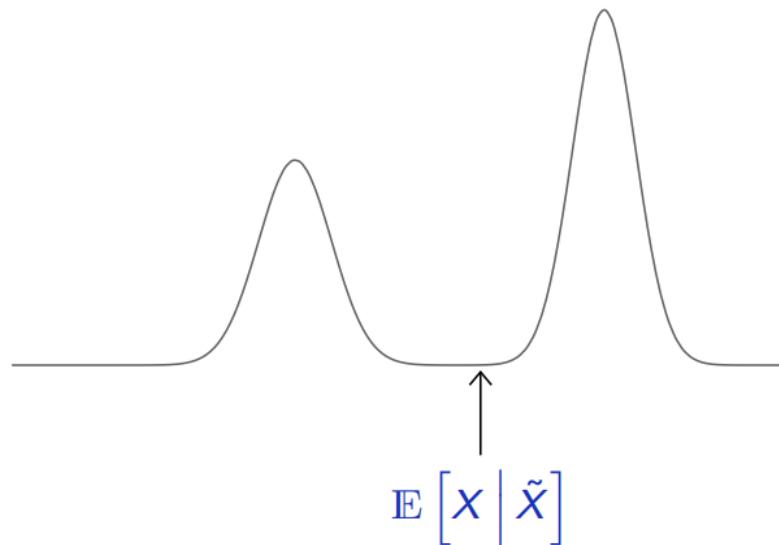
7 2 1 0 4 1 4 9 5 9 0 6  
9 0 1 5 9 7 3 4 9 6 6 5  
4 0 7 4 0 1 3 1 3 4 7 2

A fundamental weakness of denoising auto-encoders is that the posterior  $p(\mathbf{x}|\tilde{\mathbf{x}})$  is possibly multi-modal.

If we train an auto-encoder with the quadratic loss (i.e., implicitly assuming a Gaussian likelihood), then the best reconstruction is

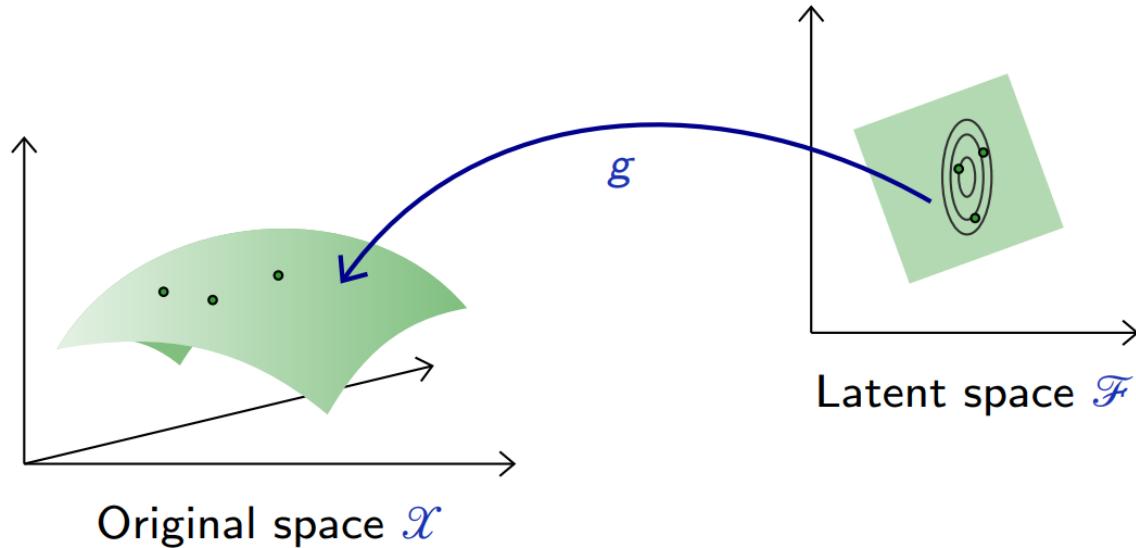
$$h(\tilde{\mathbf{x}}) = \mathbb{E}[\mathbf{x}|\tilde{\mathbf{x}}],$$

which may be very unlikely under  $p(\mathbf{x}|\tilde{\mathbf{x}})$ .



# Sampling from an AE's latent space

The generative capability of the decoder  $g$  in an auto-encoder can be assessed by introducing a (simple) density model  $q$  over the latent space  $\mathcal{Z}$ , sample there, and map the samples into the data space  $\mathcal{X}$  with  $g$ .

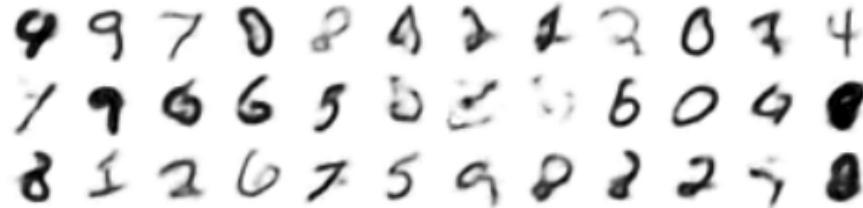


For instance, a factored Gaussian model with diagonal covariance matrix,

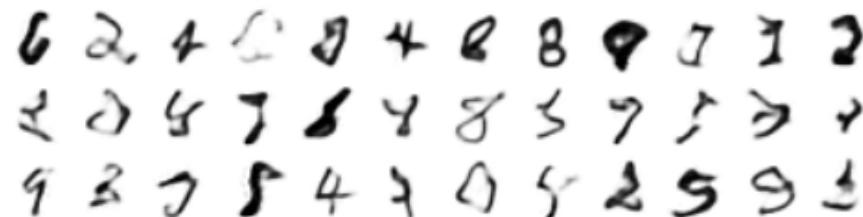
$$q(\mathbf{z}) = \mathcal{N}(\hat{\mu}, \hat{\Sigma}),$$

where both  $\hat{\mu}$  and  $\hat{\Sigma}$  are estimated on training data.

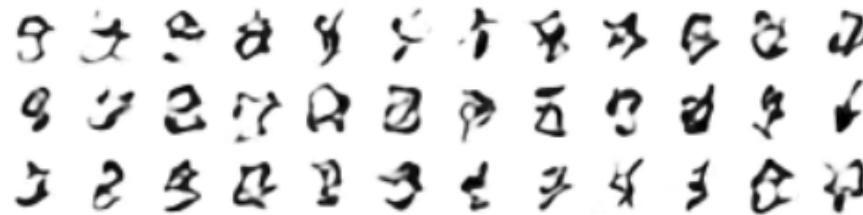
Autoencoder sampling ( $d = 8$ )



Autoencoder sampling ( $d = 16$ )



Autoencoder sampling ( $d = 32$ )

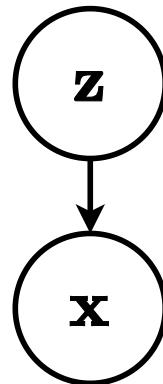


These results are not satisfactory because the density model on the latent space is **too simple and inadequate**.

Building a good model in latent space amounts to our original problem of modeling an empirical distribution, although it may now be in a lower dimension space.

# Variational inference

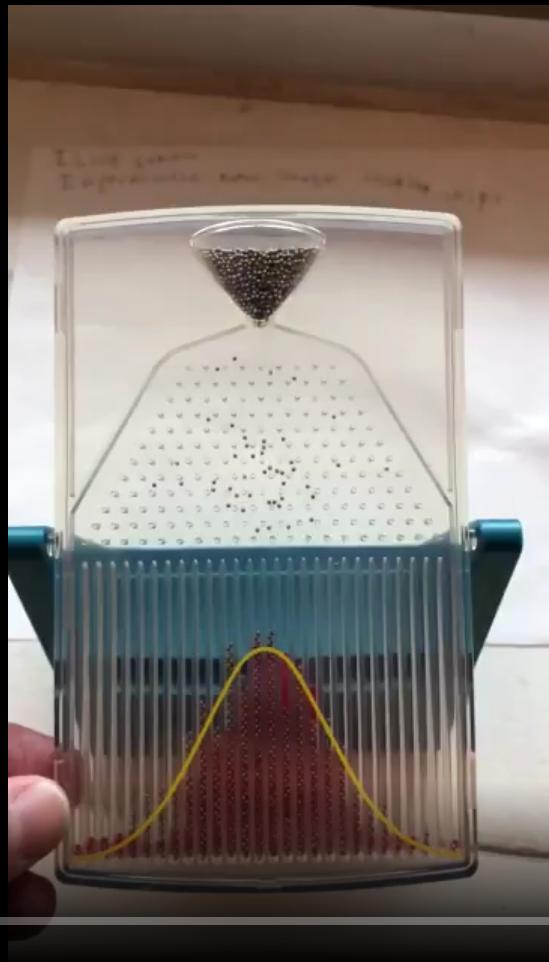
## Latent variable model



Consider for now a **prescribed latent variable model** that relates a set of observable variables  $\mathbf{x} \in \mathcal{X}$  to a set of unobserved variables  $\mathbf{z} \in \mathcal{Z}$ .

The probabilistic model defines a joint probability distribution  $p_{\theta}(\mathbf{x}, \mathbf{z})$ , which decomposes as

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z}).$$



II 0:00 / 0:45

🔇 🔍 ⏮

## How to fit a latent variable model?

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p_{\theta}(\mathbf{x}) \\&= \arg \max_{\theta} \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\&= \arg \max_{\theta} \mathbb{E}_{p(\mathbf{z})} [p_{\theta}(\mathbf{x}|\mathbf{z})] d\mathbf{z} \\&\approx \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N p_{\theta}(\mathbf{x}|\mathbf{z}_i)\end{aligned}$$

## How to fit a latent variable model?

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p_{\theta}(\mathbf{x}) \\&= \arg \max_{\theta} \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\&= \arg \max_{\theta} \mathbb{E}_{p(\mathbf{z})} [p_{\theta}(\mathbf{x}|\mathbf{z})] d\mathbf{z} \\&\approx \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N p_{\theta}(\mathbf{x}|\mathbf{z}_i)\end{aligned}$$

The curse of dimensionality will lead to poor estimates of the expectation.

## Variational inference

Let us instead consider a variational approach to fit the model parameters  $\theta$ .

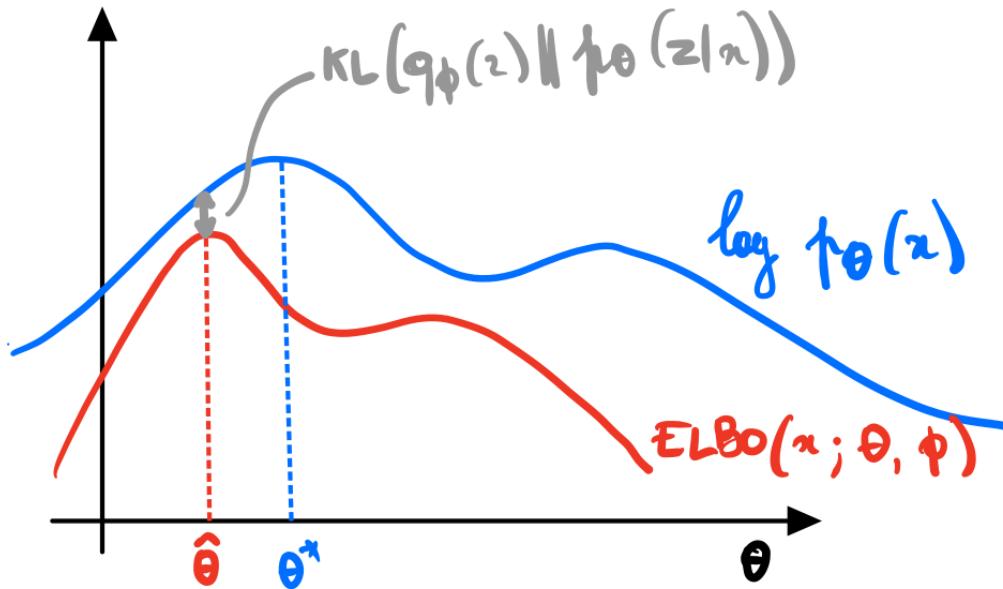
Using a **variational distribution**  $q_\phi(\mathbf{z})$  over the latent variables  $\mathbf{z}$ , we have

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

Using the Bayes rule, we can also write

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

Therefore,  $\log p_\theta(\mathbf{x}) = \text{ELBO}(\mathbf{x}; \theta, \phi) + \text{KL}(q_\phi(\mathbf{z}) || p_\theta(\mathbf{z}|\mathbf{x}))$ .



Provided the KL gap remains small, the model parameters can now be optimized by maximizing the ELBO,

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \text{ELBO}(x; \theta, \phi).$$

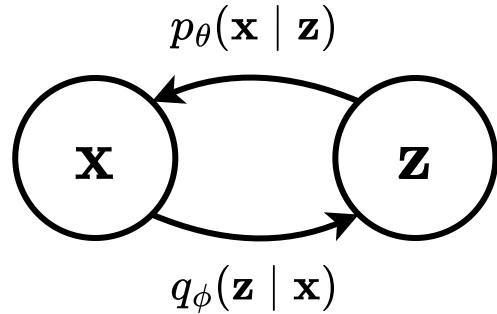
## Optimization

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \text{ELBO}(\mathbf{x}; \theta, \phi).$$

We can proceed by gradient ascent, provided we can evaluate  $\nabla_\theta \text{ELBO}(\mathbf{x}; \theta, \phi)$  and  $\nabla_\phi \text{ELBO}(\mathbf{x}; \theta, \phi)$ .

In general, the gradient of the ELBO is intractable to compute, but we can estimate it with Monte Carlo integration.

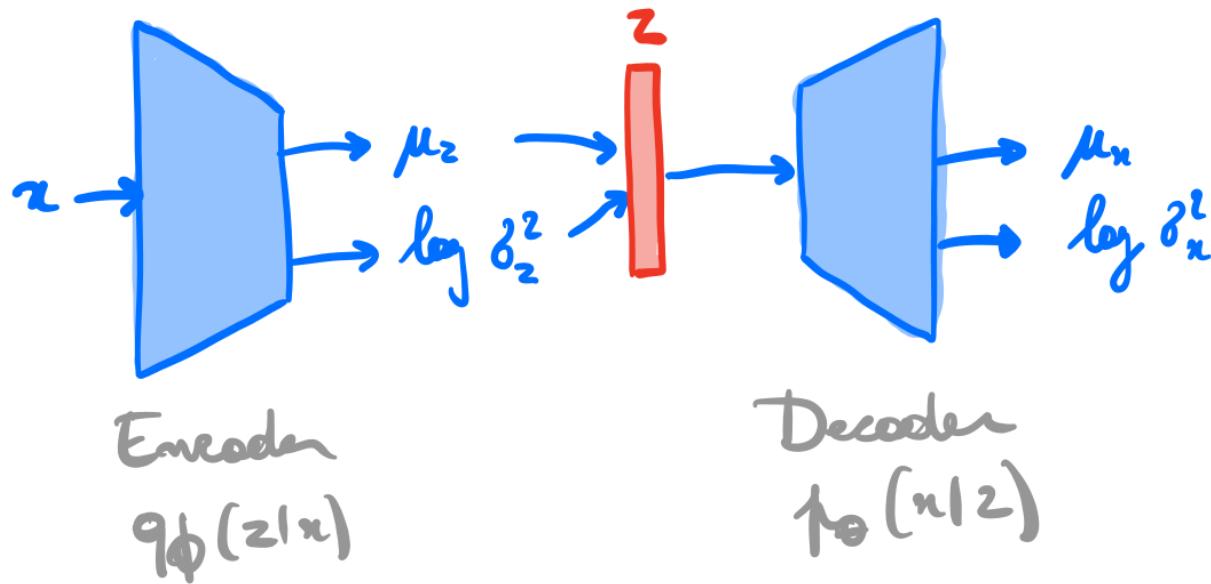
# Variational auto-encoders



So far we assumed a prescribed probabilistic model motivated by domain knowledge. We will now directly learn a stochastic generating process  $p_\theta(\mathbf{x}|\mathbf{z})$  with a neural network.

We will also amortize the inference process by learning a second neural network  $q_\phi(\mathbf{z}|\mathbf{x})$  approximating the posterior, conditionally on the observed data  $\mathbf{x}$ .

## Variational auto-encoders



A variational auto-encoder is a deep latent variable model where:

- The prior  $p(\mathbf{z})$  is prescribed, and usually chosen to be Gaussian.
- The likelihood  $p_\theta(\mathbf{x}|\mathbf{z})$  is parameterized with a generative network  $\text{NN}_\theta$  (or decoder) that takes as input  $\mathbf{z}$  and outputs parameters  $\varphi = \text{NN}_\theta(\mathbf{z})$  to the data distribution. E.g.,

$$\mu, \sigma = \text{NN}_\theta(\mathbf{z})$$

$$p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mu, \sigma^2 \mathbf{I})$$

- The approximate posterior  $q_\phi(\mathbf{z}|\mathbf{x})$  is parameterized with an inference network  $\text{NN}_\phi$  (or encoder) that takes as input  $\mathbf{x}$  and outputs parameters  $\nu = \text{NN}_\phi(\mathbf{x})$  to the approximate posterior. E.g.,

$$\mu, \sigma = \text{NN}_\phi(\mathbf{x})$$

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu, \sigma^2 \mathbf{I})$$

As before, we can use variational inference to jointly optimize the encoder and decoder networks parameters  $\phi$  and  $\theta$ , but now in expectation over the data distribution  $p(\mathbf{x})$ :

$$\begin{aligned}\theta^*, \phi^* &= \arg \max_{\theta, \phi} \mathbb{E}_{p(\mathbf{x})} [\text{ELBO}(\mathbf{x}; \theta, \phi)] \\ &= \arg \max_{\theta, \phi} \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \right] \\ &= \arg \max_{\theta, \phi} \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \right].\end{aligned}$$

Interpretation:

- Given some decoder network set at  $\theta$ , we want to put the mass of the latent variables, by adjusting  $\phi$ , such that they explain the observed data, while remaining close to the prior.
- Given some encoder network set at  $\phi$ , we want to put the mass of the observed variables, by adjusting  $\theta$ , such that they are well explained by the latent variables.

Unbiased gradients of the ELBO with respect to the generative model parameters  $\theta$  are simple to obtain, as

$$\begin{aligned}\nabla_{\theta} \text{ELBO}(\mathbf{x}; \theta, \phi) &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z})],\end{aligned}$$

which can be estimated with Monte Carlo integration.

However, gradients with respect to the inference model parameters  $\phi$  are more difficult to obtain since

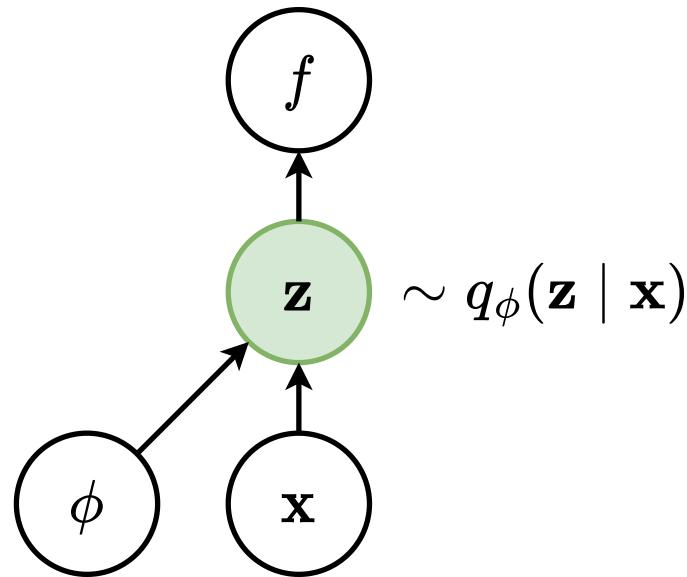
$$\begin{aligned}\nabla_{\phi} \text{ELBO}(\mathbf{x}; \theta, \phi) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &\neq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\nabla_{\phi} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))].\end{aligned}$$

## Reparameterization trick

Let us abbreviate

$$\begin{aligned}\text{ELBO}(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [f(\mathbf{x}, \mathbf{z}; \phi)].\end{aligned}$$

The computational graph of a Monte Carlo estimate of the ELBO would look like



Issue: We cannot backpropagate through the stochastic node  $\mathbf{z}$  to compute  $\nabla_\phi f$ !

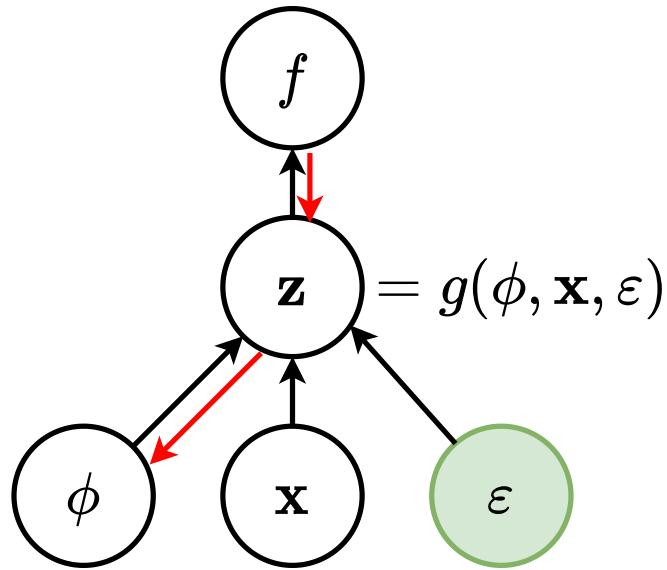
The reparameterization trick consists in re-expressing the variable

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$$

as some differentiable and invertible transformation of another random variable  $\epsilon$  given  $\mathbf{x}$  and  $\phi$ ,

$$\mathbf{z} = g(\phi, \mathbf{x}, \epsilon),$$

such that the distribution of  $\epsilon$  is independent of  $\mathbf{x}$  or  $\phi$ .



If  $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu(\mathbf{x}; \phi), \sigma^2(\mathbf{x}; \phi))$ , where  $\mu(\mathbf{x}; \phi)$  and  $\sigma^2(\mathbf{x}; \phi)$  are the outputs of the inference network  $NN_\phi$ , then a common reparameterization is

$$p(\epsilon) = \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$

$$\mathbf{z} = \mu(\mathbf{x}; \phi) + \sigma(\mathbf{x}; \phi) \odot \epsilon.$$

Given this change of variable, the ELBO can be rewritten as

$$\begin{aligned}\text{ELBO}(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [f(\mathbf{x}, \mathbf{z}; \phi)] \\ &= \mathbb{E}_{p(\epsilon)} [f(\mathbf{x}, g(\phi, \mathbf{x}, \epsilon); \phi)].\end{aligned}$$

Therefore estimating the gradient of the ELBO with respect to  $\phi$  is now easy, as

$$\begin{aligned}\nabla_\phi \text{ELBO}(\mathbf{x}; \theta, \phi) &= \nabla_\phi \mathbb{E}_{p(\epsilon)} [f(\mathbf{x}, g(\phi, \mathbf{x}, \epsilon); \phi)] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_\phi f(\mathbf{x}, g(\phi, \mathbf{x}, \epsilon); \phi)],\end{aligned}$$

which we can now estimate with Monte Carlo integration.

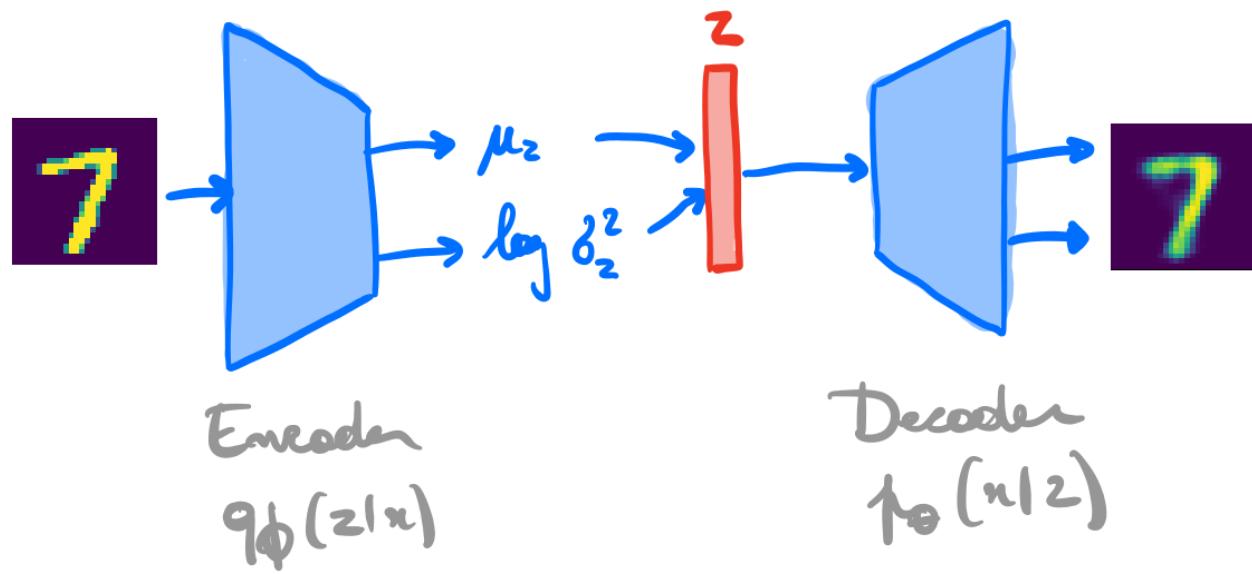
The last required ingredient is the evaluation of the approximate posterior  $q_\phi(\mathbf{z}|\mathbf{x})$  given the change of variable  $\mathbf{g}$ . As long as  $\mathbf{g}$  is invertible, we have

$$\log q_\phi(\mathbf{z}|\mathbf{x}) = \log p(\epsilon) - \log \left| \det \left( \frac{\partial \mathbf{z}}{\partial \epsilon} \right) \right|.$$

## Step-by-step example

Consider as data **d** the MNIST digit dataset:

0  
1  
2  
3  
4  
5  
6  
7  
8  
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9



- Decoder  $p_\theta(\mathbf{x}|\mathbf{z})$ :

$$\mathbf{z} \in \mathbb{R}^d$$

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

$$p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mu(\mathbf{z}; \theta), \sigma^2(\mathbf{z}; \theta)\mathbf{I})$$

$$\mu(\mathbf{z}; \theta) = \mathbf{W}_2^T \mathbf{h} + \mathbf{b}_2$$

$$\log \sigma^2(\mathbf{z}; \theta) = \mathbf{W}_3^T \mathbf{h} + \mathbf{b}_3$$

$$\mathbf{h} = \text{ReLU}(\mathbf{W}_1^T \mathbf{z} + \mathbf{b}_1)$$

$$\theta = \{\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \mathbf{W}_3, \mathbf{b}_3\}$$

- Encoder  $q_\phi(\mathbf{z}|\mathbf{x})$ :

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu(\mathbf{x}; \phi), \sigma^2(\mathbf{x}; \phi)\mathbf{I})$$

$$p(\epsilon) = \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$

$$\mathbf{z} = \mu(\mathbf{x}; \phi) + \sigma(\mathbf{x}; \phi) \odot \epsilon$$

$$\mu(\mathbf{x}; \phi) = \mathbf{W}_5^T \mathbf{h} + \mathbf{b}_5$$

$$\log \sigma^2(\mathbf{x}; \phi) = \mathbf{W}_6^T \mathbf{h} + \mathbf{b}_6$$

$$\mathbf{h} = \text{ReLU}(\mathbf{W}_4^T \mathbf{x} + \mathbf{b}_4)$$

$$\phi = \{\mathbf{W}_4, \mathbf{b}_4, \mathbf{W}_5, \mathbf{b}_5, \mathbf{W}_6, \mathbf{b}_6\}$$

Note that there is no restriction on the encoder and decoder network architectures. They could as well be arbitrarily complex convolutional networks.

Plugging everything together, the objective can be expressed as

$$\begin{aligned}\text{ELBO}(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \\ &= \mathbb{E}_{p(\epsilon)} [\log p(\mathbf{x}|\mathbf{z} = g(\phi, \mathbf{x}, \epsilon); \theta)] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})),\end{aligned}$$

where the negative KL divergence can be expressed analytically as

$$-\text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = \frac{1}{2} \sum_{j=1}^d (1 + \log(\sigma_j^2(\mathbf{x}; \phi)) - \mu_j^2(\mathbf{x}; \phi) - \sigma_j^2(\mathbf{x}; \phi)),$$

which allows to evaluate its derivative without approximation.

|                     |                     |                     |                     |
|---------------------|---------------------|---------------------|---------------------|
| 8 6 / 7 8 1 4 8 2 8 | 5 1 6 5 1 6 7 6 7 2 | 2 8 3 8 3 8 5 7 3 8 | 8 2 0 8 9 2 3 9 0 0 |
| 9 6 8 3 9 6 0 3 1 9 | 8 5 9 4 6 8 2 1 6 8 | 8 3 8 2 7 9 3 3 3 8 | 7 5 1 9 1 1 7 1 4 4 |
| 3 3 9 1 3 6 8 1 7 9 | 6 1 5 3 2 8 8 1 3 8 | 2 5 9 9 4 3 9 5 1 6 | 8 9 6 2 0 8 2 8 2 9 |
| 8 9 0 8 6 9 1 9 6 3 | 2 8 6 8 9 1 0 0 4 1 | 1 9 1 8 9 3 3 4 9 2 | 2 9 8 6 3 1 7 0 6 1 |
| 9 2 3 3 3 3 1 3 8 6 | 5 1 9 3 0 1 5 3 5 9 | 2 7 3 6 4 2 0 2 0 3 | 5 9 7 9 8 9 9 9 1 0 |
| 6 9 9 8 6 1 6 6 6 6 | 6 5 6 1 4 9 1 7 5 8 | 5 9 7 0 5 9 3 8 4 5 | 6 8 8 6 2 4 8 2 8 1 |
| 9 5 2 6 6 5 1 8 9 9 | 1 3 4 3 9 8 3 2 7 0 | 6 9 4 3 6 2 8 5 5 2 | 7 5 8 2 4 6 1 3 8 8 |
| 9 9 7 7 3 1 2 8 2 3 | 4 5 8 2 9 7 0 4 5 3 | 8 4 9 0 8 0 7 0 6 6 | 9 9 3 9 2 2 9 3 9 0 |
| 0 4 6 1 2 3 2 0 8 8 | 6 9 4 4 9 7 2 3 9 8 | 7 4 3 6 3 0 3 6 0 1 | 4 5 2 4 3 9 0 1 8 4 |
| 9 7 5 4 9 3 4 8 5 1 | 2 6 4 5 6 0 9 7 9 8 | 2 1 2 0 9 7 1 0 0 0 | 8 8 7 2 3 1 6 2 3 6 |

(a) 2-D latent space

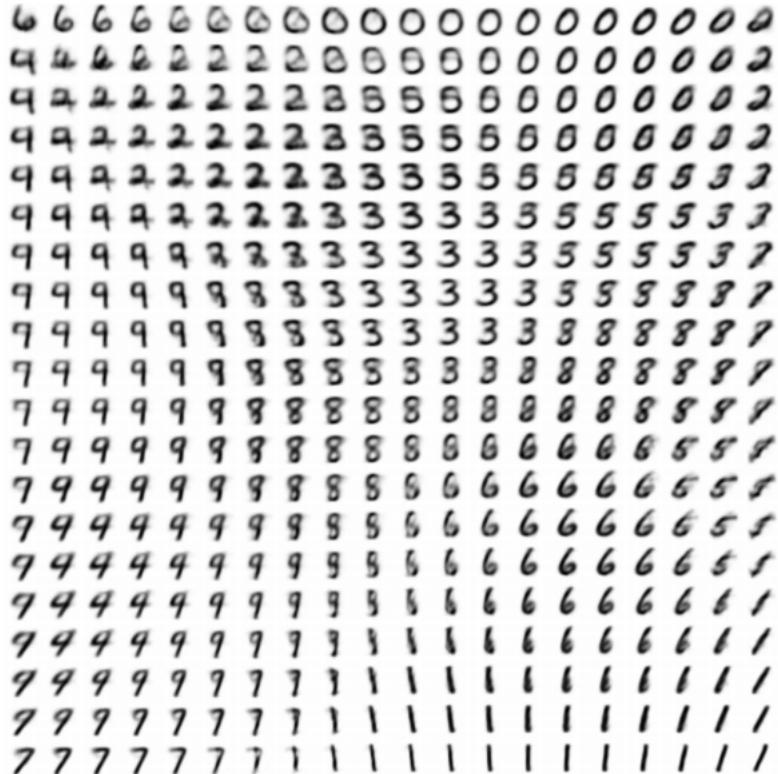
(b) 5-D latent space

(c) 10-D latent space

(d) 20-D latent space



(a) Learned Frey Face manifold

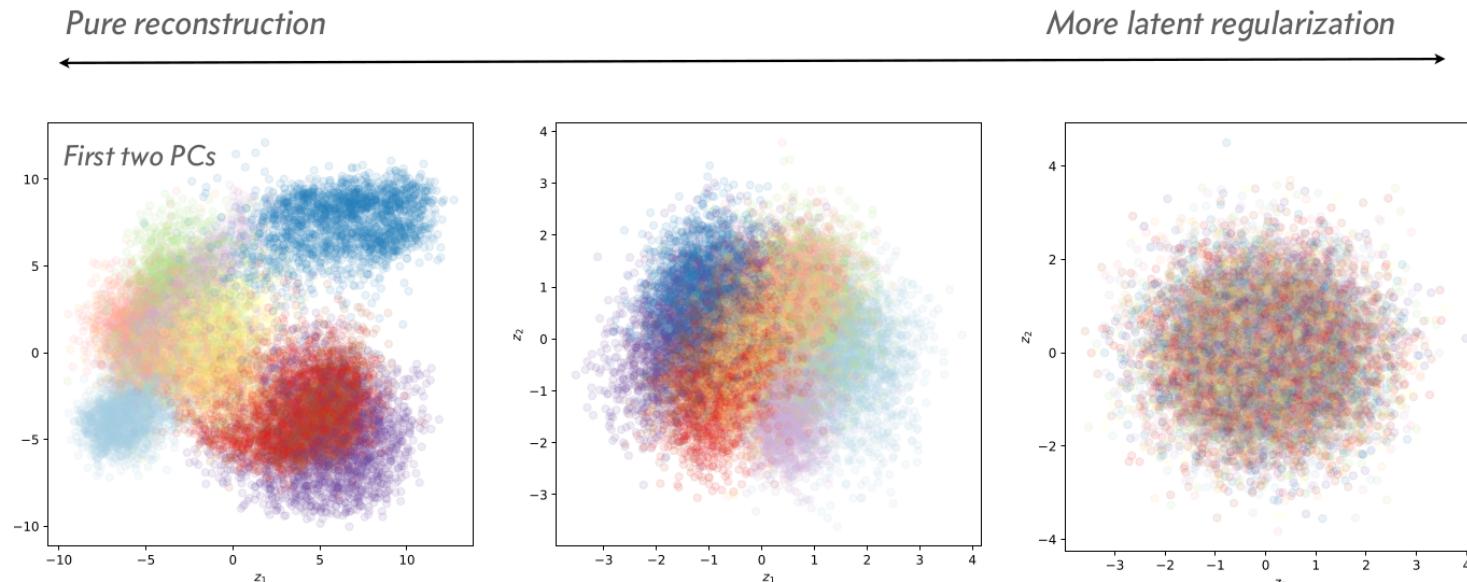


(b) Learned MNIST manifold

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  $\mathbf{z}$ . For each of these values  $\mathbf{z}$ , we plotted the corresponding generative  $p_{\theta}(\mathbf{x}|\mathbf{z})$  with the learned parameters  $\theta$ .

## A semantically meaningful latent space

The prior-matching term  $\text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$  enforces simplicity in the latent space, encouraging learned semantic structure and disentanglement.



# Some selected applications

Original images



Compression rate: 0.2bits/dimension

JPEG



JPEG-2000



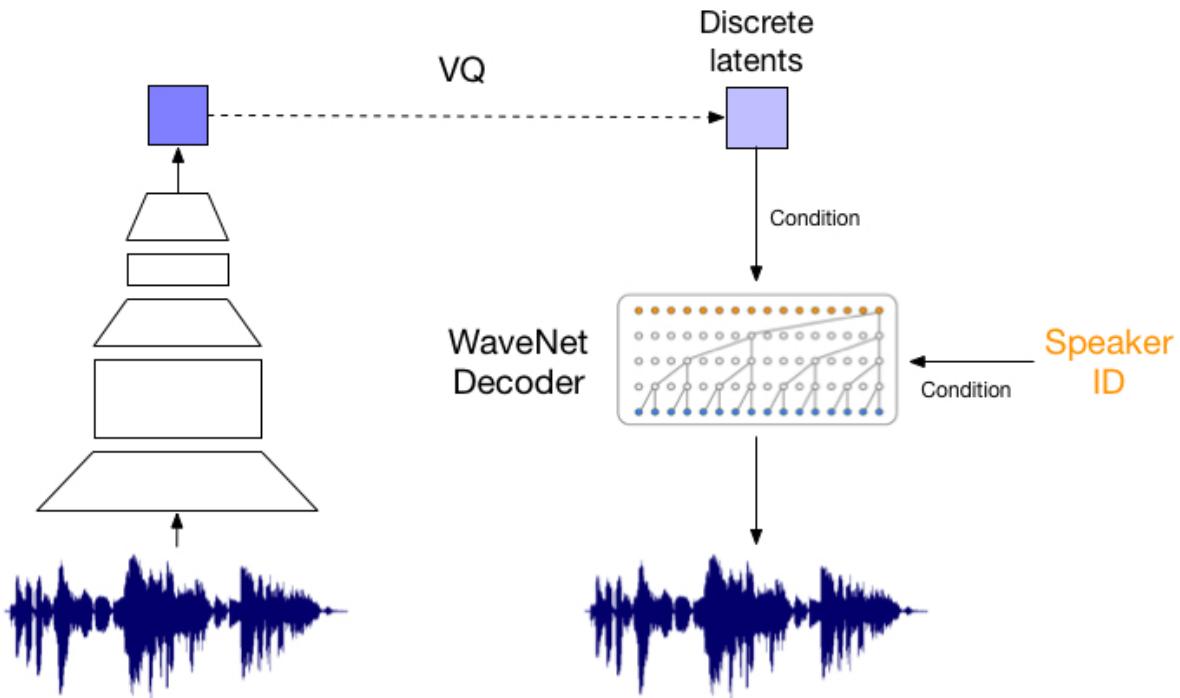
RVAE v1



RVAE v2

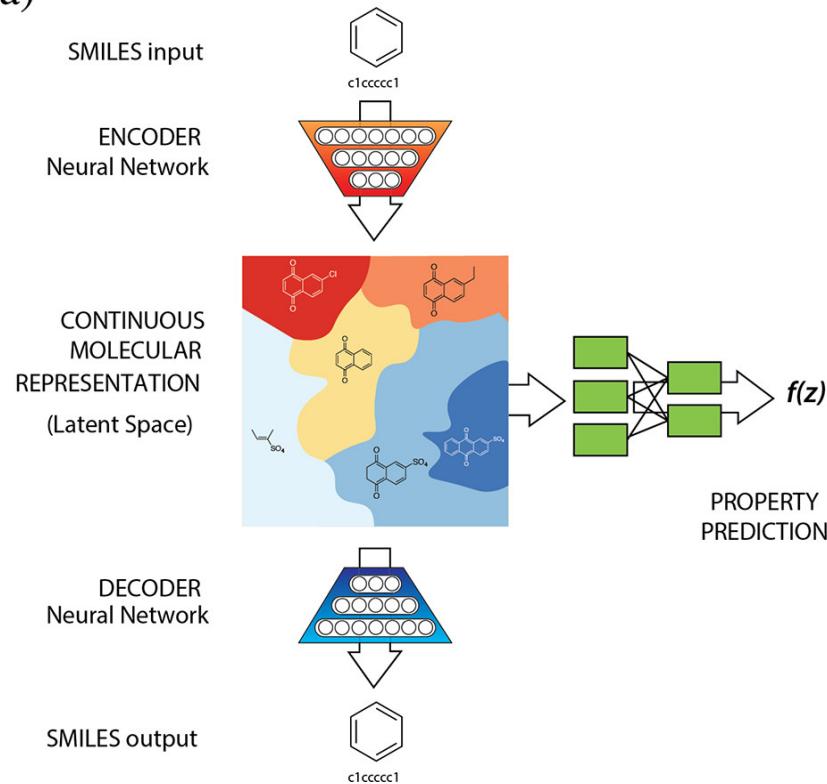


Hierarchical **compression of images and other data**,  
e.g., in video conferencing systems (Gregor et al, 2016).

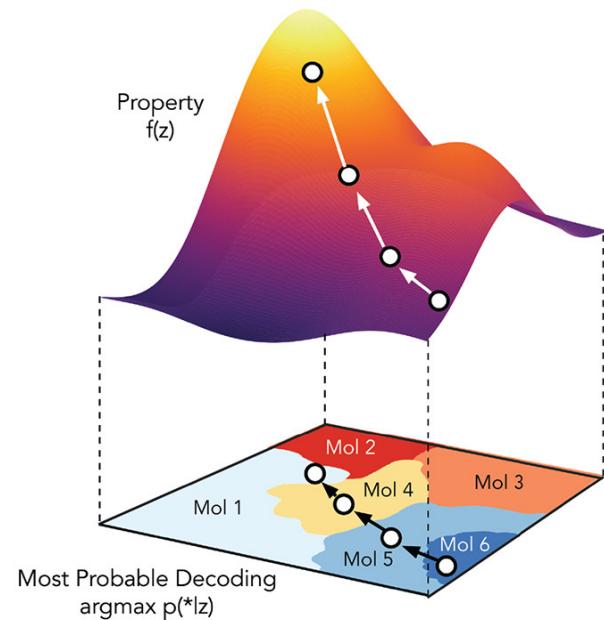


**Voice style transfer** [demo] (van den Oord et al, 2017).

(a)



(b)



**Design of new molecules** with desired chemical properties  
 (Gomez-Bombarelli et al, 2016).

