

Introduction to Machine Learning and Artificial Intelligence: Lecture III

Michael Kagan



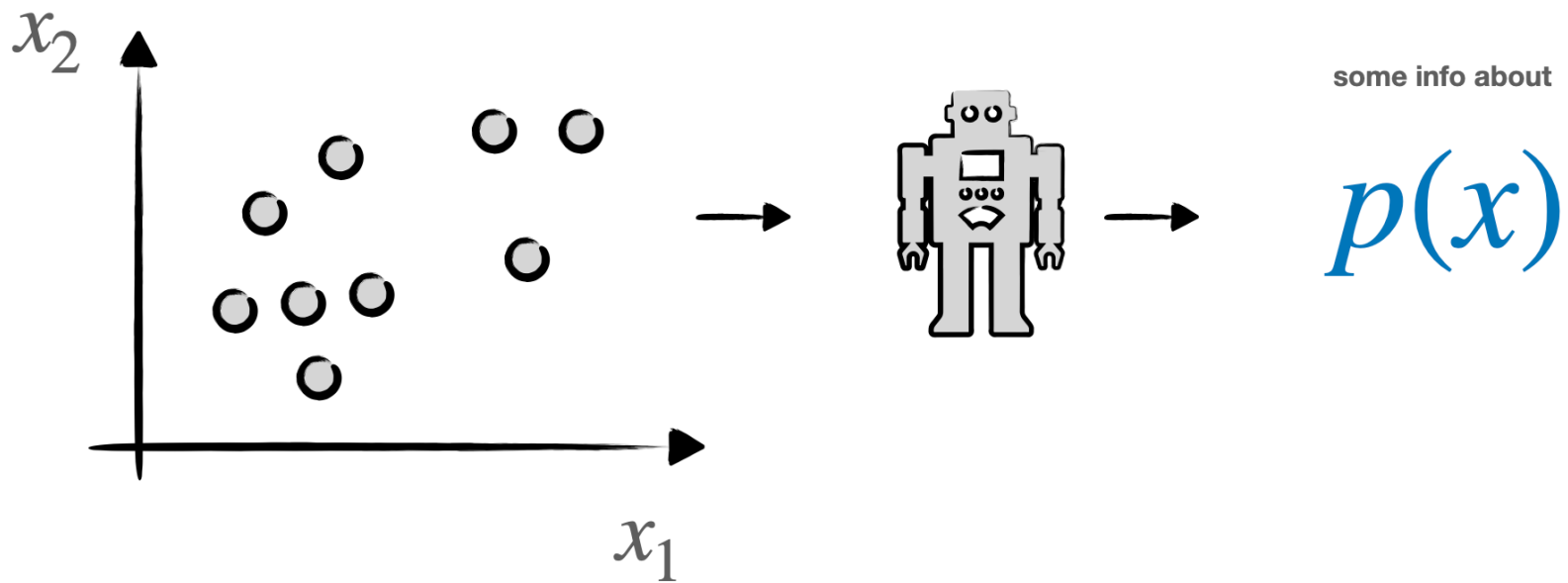
2nd COFI Advanced Instrumentation and Analysis Techniques School
December 10, 2023

- Lecture 1
 - Introduction to Machine Learning fundamentals
 - Linear Models
- Lecture 2
 - Neural Networks
 - Deep Neural Networks
 - Inductive Bias and Model Architectures
- Lecture 3
 - Unsupervised Learning
 - Autoencoders
 - Towards generative modeling: Variation Autoencoders

Beyond Regression and Classification

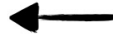
- Not all tasks are predicting a label from features, as in classification and regression
- May want to model a high-dim. signal
 - Data synthesis / simulation
 - Density estimation
 - Anomaly detection
 - Denoising, super resolution
 - Data compression
 - ...
- Often don't have labels → Unsupervised Learning

- Our goal is to study the data density $p(x)$
- Even w/o labels, aim to characterize the distribution



A process

$$\text{Die} \rightarrow \mathbb{R}^2$$



A formula

$$\mathbb{R}^2 \rightarrow \mathbb{R}$$

$$p_{\mu, \Sigma}(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma (x - \mu)\right)$$

**Evaluating the Probability
for a given sample**

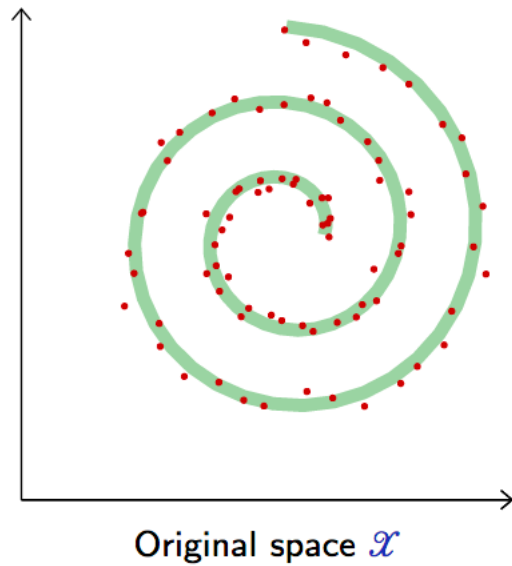
"Understanding $p(x)$ " – ability to do either or both of these

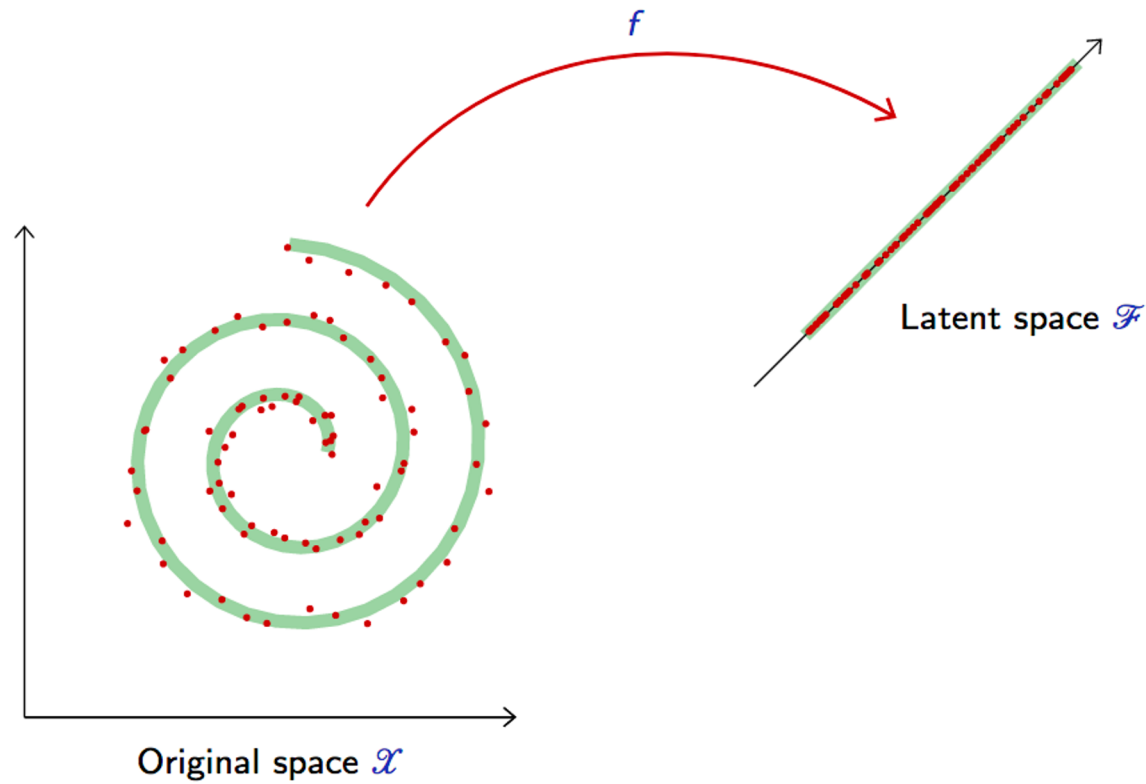
- In many cases, we don't have a theory of the underlying process → *Can still learn to sample*
- Deep learning can be very good at this!

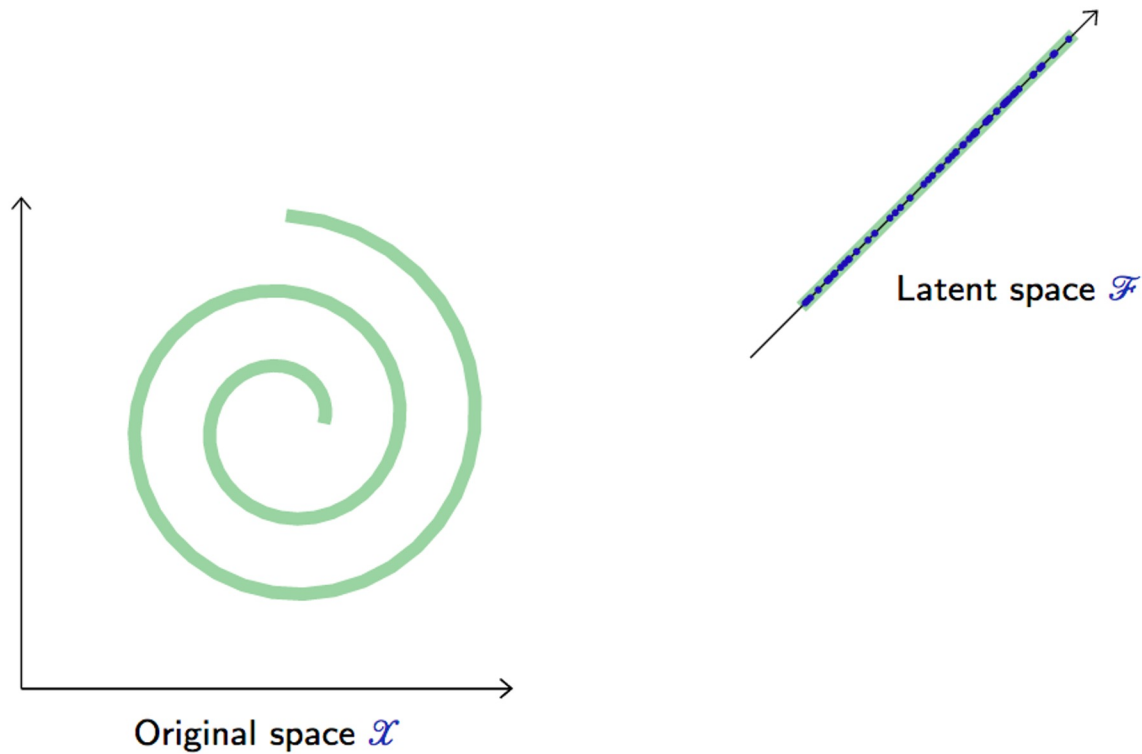
face $\sim p(\text{face})$

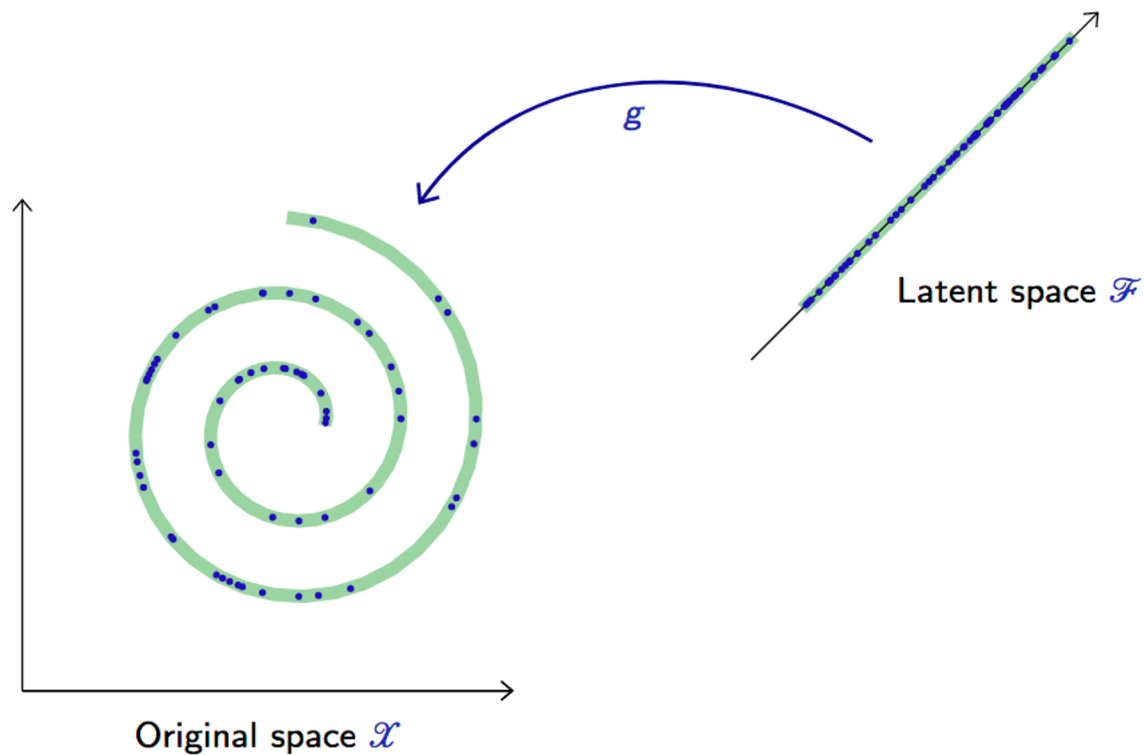


- Unsupervised learning is more heterogeneous than supervised learning
- Many architectures, losses, learning strategies
- Often constructed so model converges to $p(x)$
 - Variational inference, Adversarial learning, Self-supervision, ...
- Often framed as **modeling the lower dimensional “meaningful degrees of freedom”** that describe the data

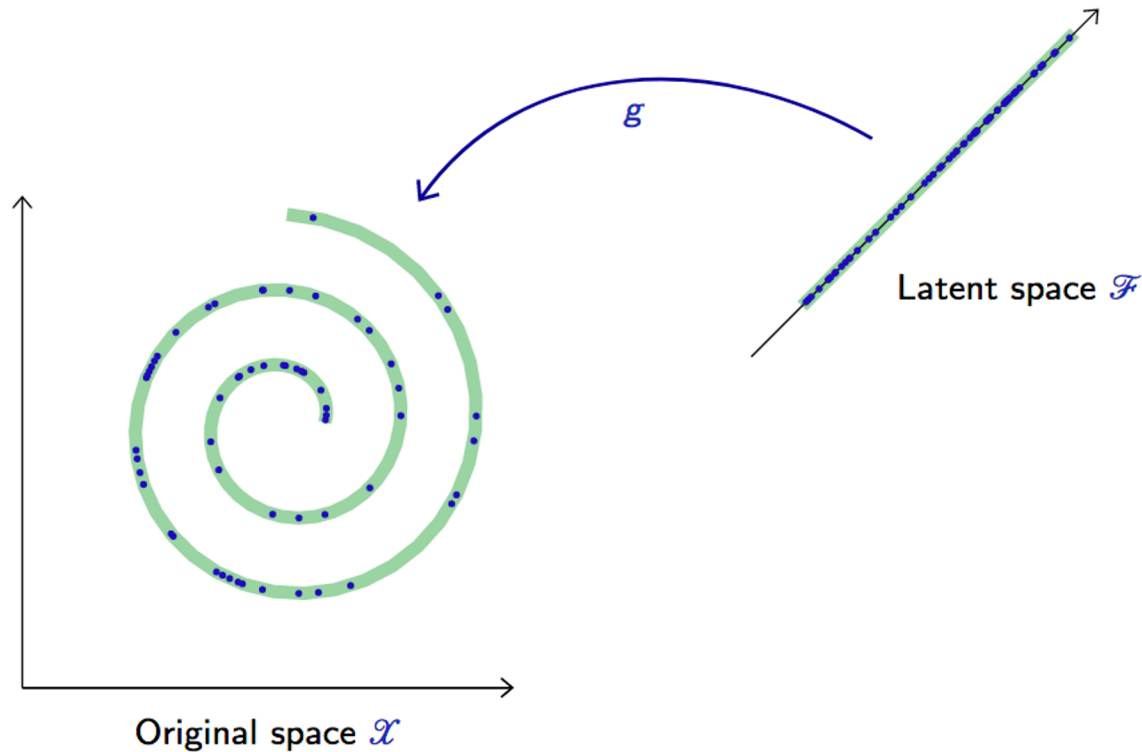








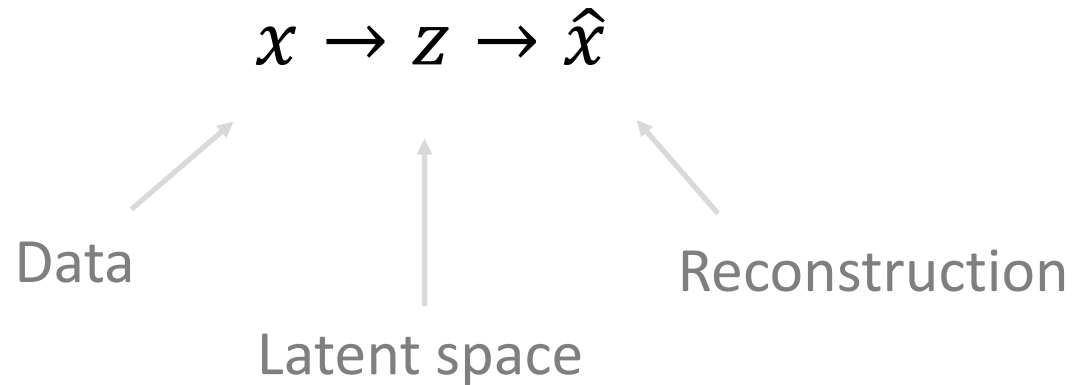
How can we find the “meaningful degrees of freedom” in the data?



- Dimensionality Reduction / Compression
- Can we learn to:
 1. Compress the data to a *latent space* with smaller number of dimensions
 2. Recover the original data from this latent space?
- Latent space must encode and retain the important information about the data

Autoencoders

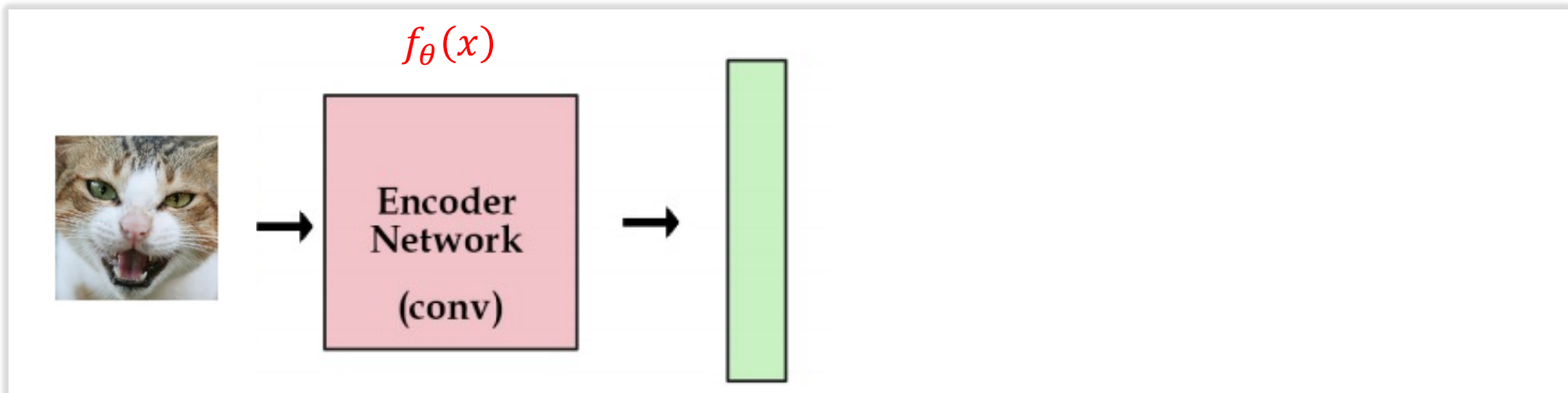
- Map a space to itself through a compression



- Map a space to itself through a compression

$$x \rightarrow z \rightarrow \hat{x}$$

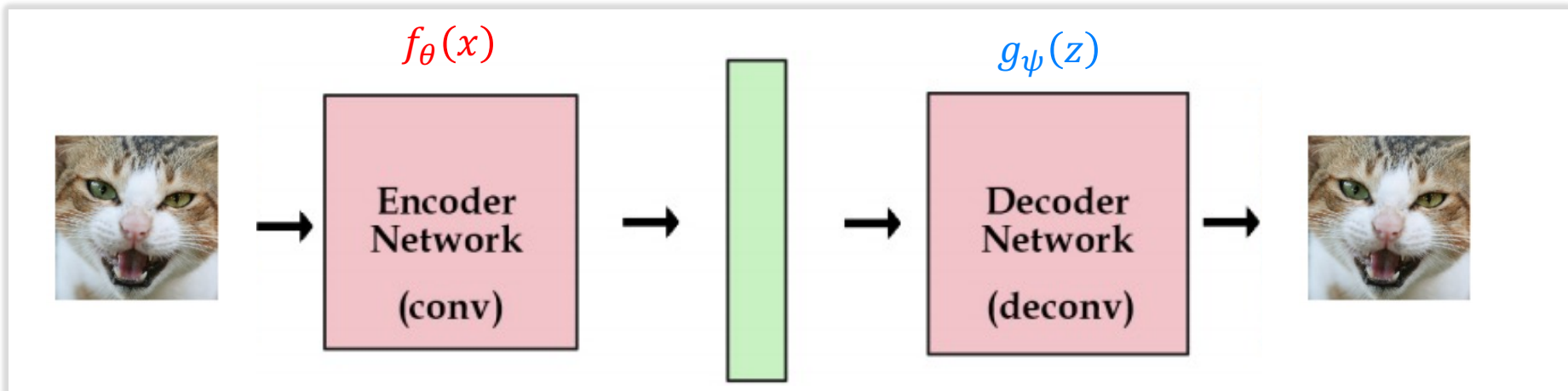
- **Encoder**: Map from data to a lower dim. latent space
 - Neural network $f_{\theta}(x)$ with parameters θ

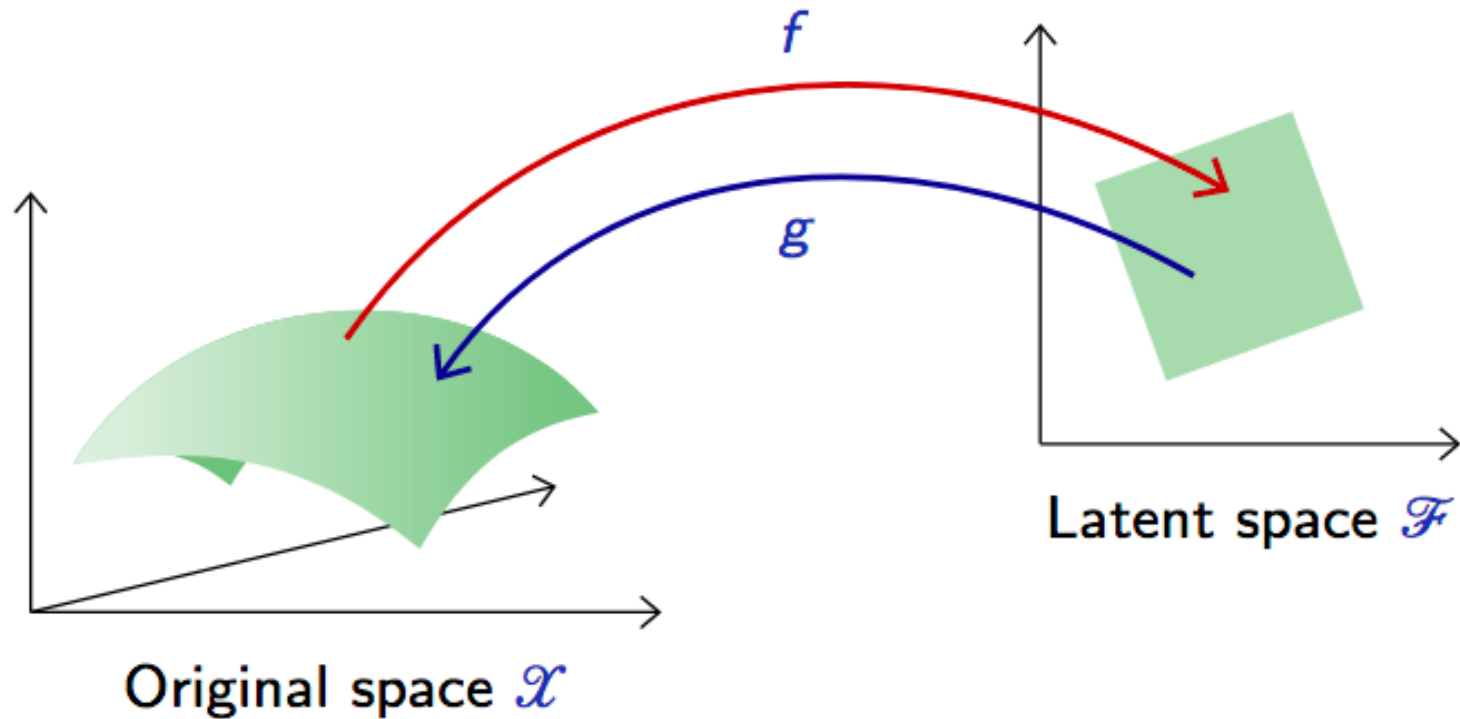


- Map a space to itself through a compression

$$x \rightarrow z \rightarrow \hat{x}$$

- **Encoder**: Map from data to a lower dim. latent space
 - Neural network $f_{\theta}(x)$ with parameters θ
- **Decoder**: Map from latent space back to data space
 - Neural network $g_{\psi}(z)$ with parameters ψ

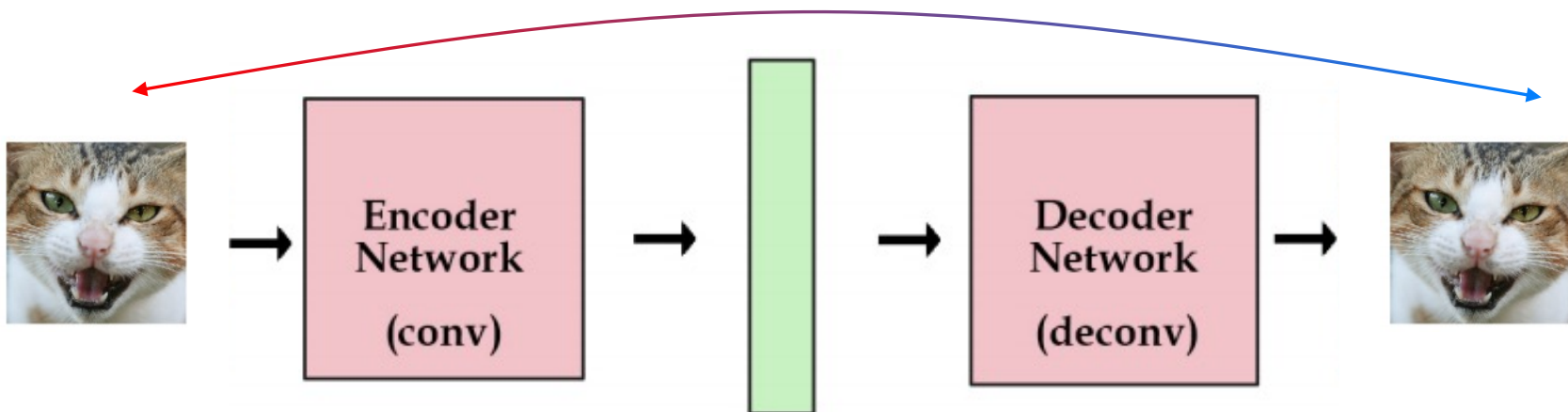




- Latent space is of lower dimension than data
- Model must learn a “good” parametrization and capture dependencies between components

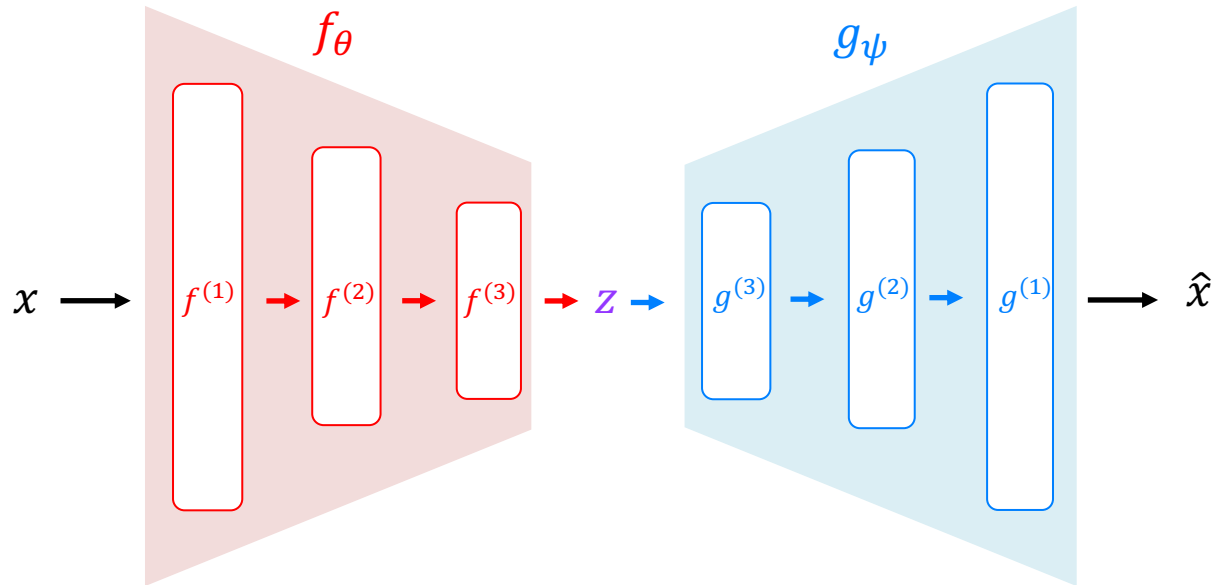
$$L(\theta, \psi) = \frac{1}{N} \sum_n \|x_n - g_\psi(f_\theta(x_n))\|^2$$

- **Loss:** mean *reconstruction loss* (MSE) between data and encoded-decoded data
- Min. over params. of encoder (θ) and decoder (ψ).



$$L(\theta, \psi) = \frac{1}{N} \sum_n \|x_n - g_\psi(f_\theta(x_n))\|^2$$

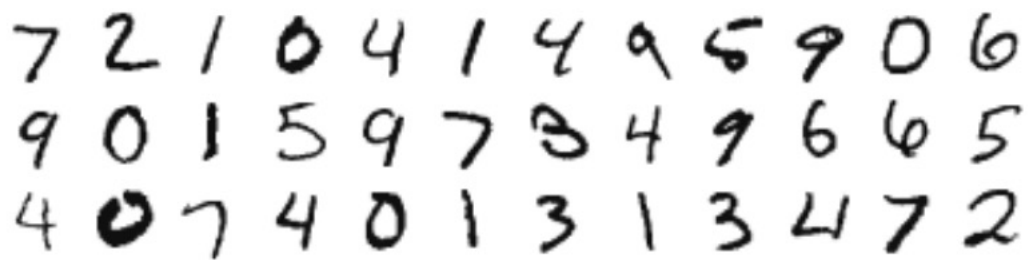
- Loss: mean *reconstruction loss* (MSE) between data and encoded-decoded data
- Min. over params. of encoder (θ) and decoder (ψ).
- NOTE: if $f_\theta(x)$ and $g_\psi(z)$ are linear, optimal solution given by Principle Components Analysis



- When f_{θ} and g_{ψ} are multiple neural network layers, can learn complex mappings
 - f_{θ} and g_{ψ} can be Fully Connected, CNNs, RNNs, etc.
 - Choice of network structure will depend on data

Deep Convolutional Autoencoder

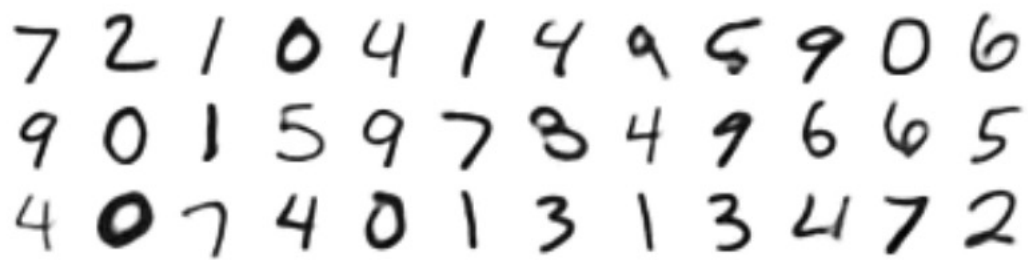
X (original samples)



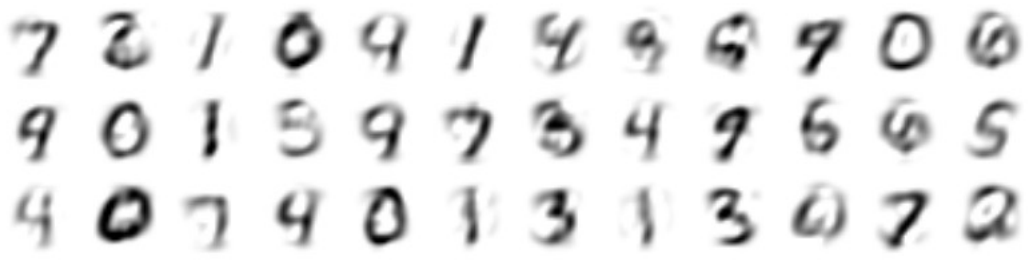
$g \circ f(X)$ (CNN, $d = 16$)



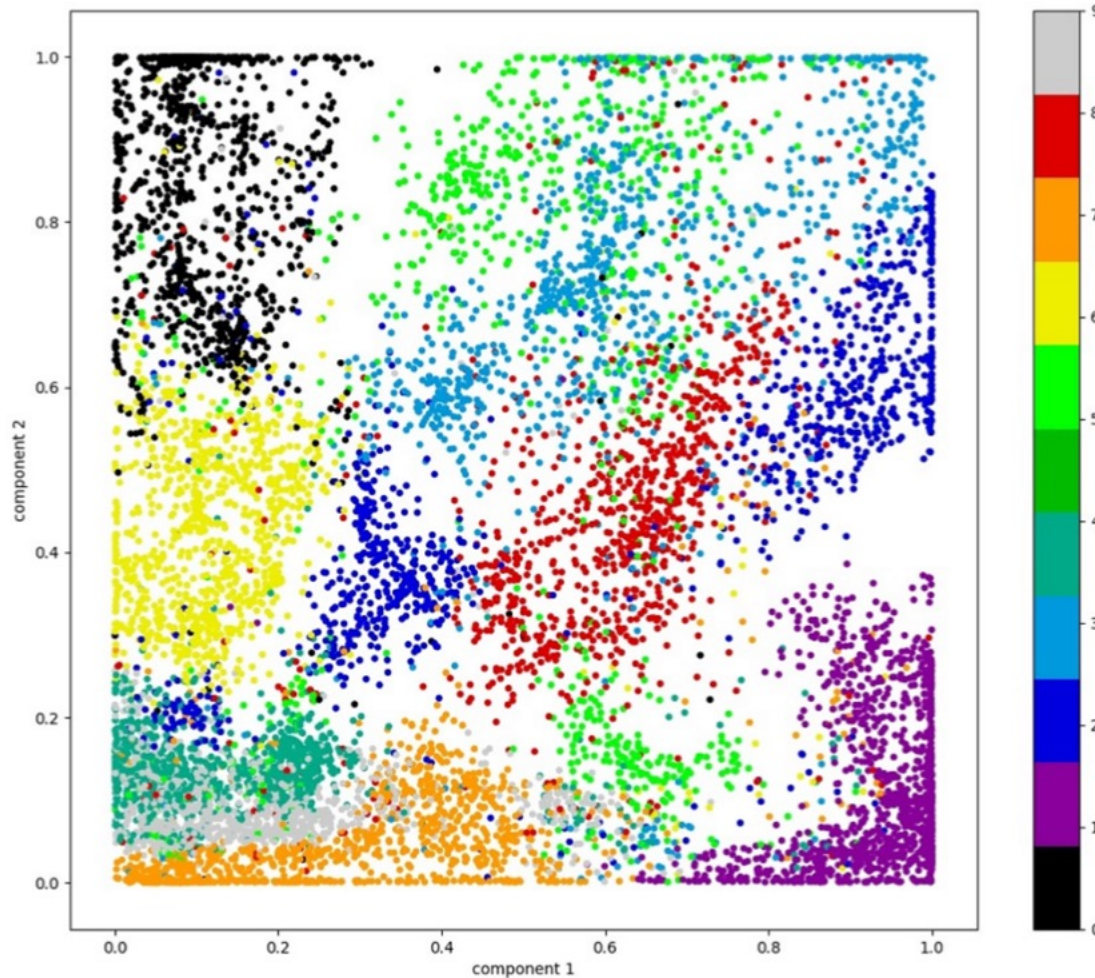
f_θ and g_ψ are five convolutional layers



$g \circ f(X)$ (PCA, $d = 16$)

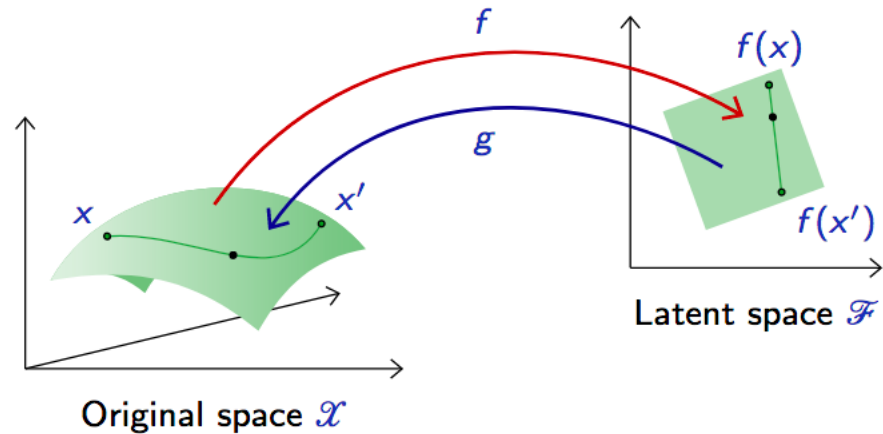


- Can look at latent space to see how the model arranges the data

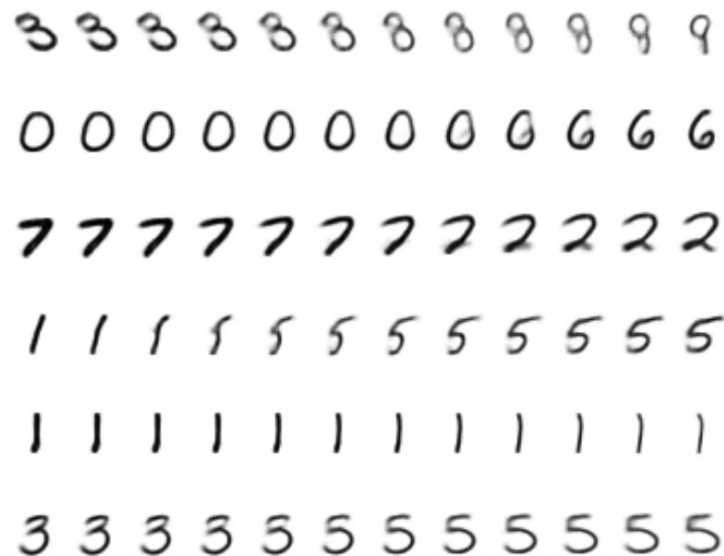


Interpolating in Latent Space

$$\alpha \in [0, 1], \quad \xi(x, x', \alpha) = g((1 - \alpha)f(x) + \alpha f(x')).$$

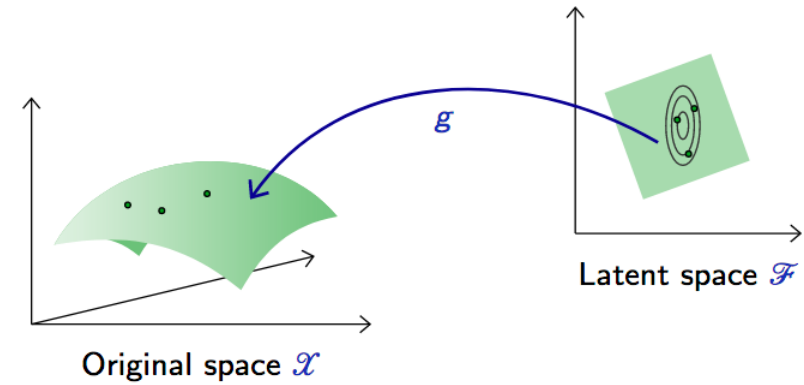


Autoencoder interpolation ($d = 8$)



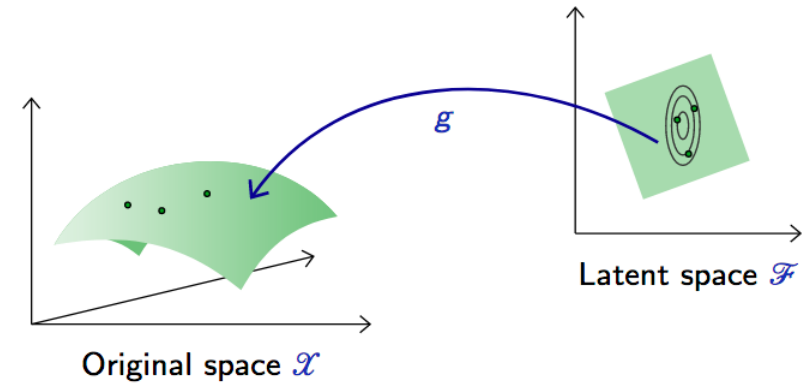
Can We Generate Data with Decoder?

- Can we sample in latent space and decode to generate data?



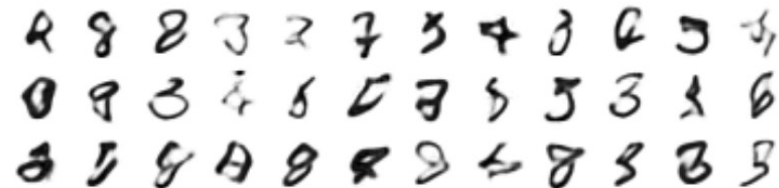
Can We Generate Data with Decoder?

- Can we sample in latent space and decode to generate data?

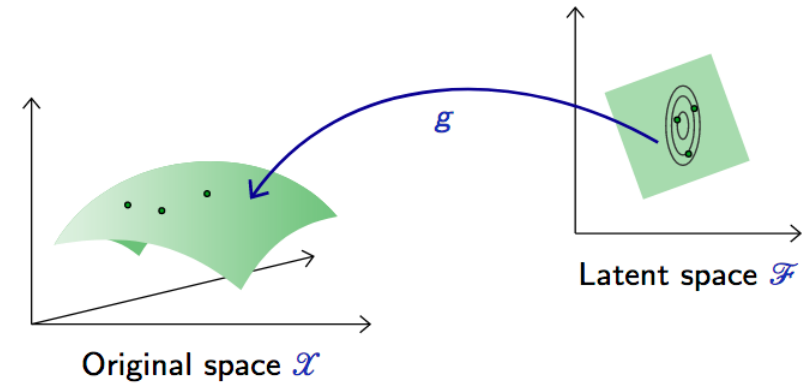


- What distribution to sample from in latent space?
 - Try Gaussian with mean and variance from data

Autoencoder sampling ($d = 16$)



- Can we sample in latent space and decode to generate data?



- What distribution to sample from in latent space?
 - Try Gaussian with mean and variance from data



Autoencoder sampling ($d = 16$)



- Don't know the right latent space density

Generative Models

A **generative model** is a probabilistic model q that can be used as a simulator of the data.

Goal: generate synthetic, realistic high-dimension data

$$x \sim q(x; \theta)$$

that is as close as possible to the unknown data distribution $p(x)$ for which we have empirical samples.

i.e. want to recreate the raw data distribution (such as the distribution of natural images).

- Generative models aim to:
 - Learn a distribution $p(x)$ that explains the data
 - Draw samples of plausible data points
- Explicit Models
 - Can evaluate the density $p(x)$ of a data point x
- Implicit Models
 - Can only sample $p(x)$, but not evaluate density

Variational Autoencoders

- Learn a mapping from corrupted data space $\tilde{\mathcal{X}}$ back to original data space

- Mapping $\phi_w(\tilde{x}) = x$

- ϕ_w will be a neural network with parameters w

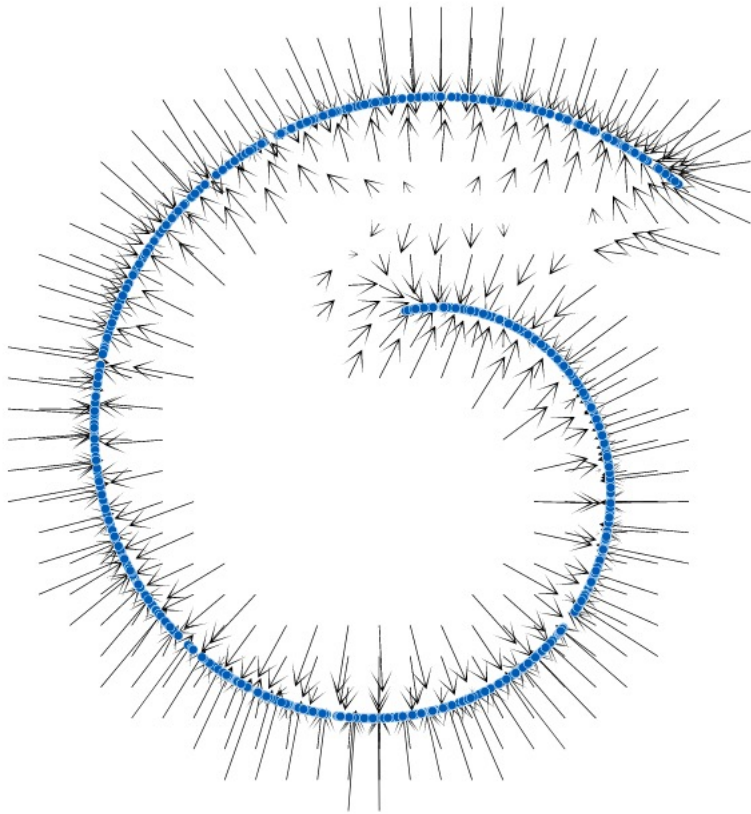
- Loss:

$$L = \frac{1}{N} \sum_n \|x_n - \phi_w(x_n + \epsilon_n)\|$$



Perturbation, e.g. Gaussian noise

Denoising Autoencoders Examples



Original

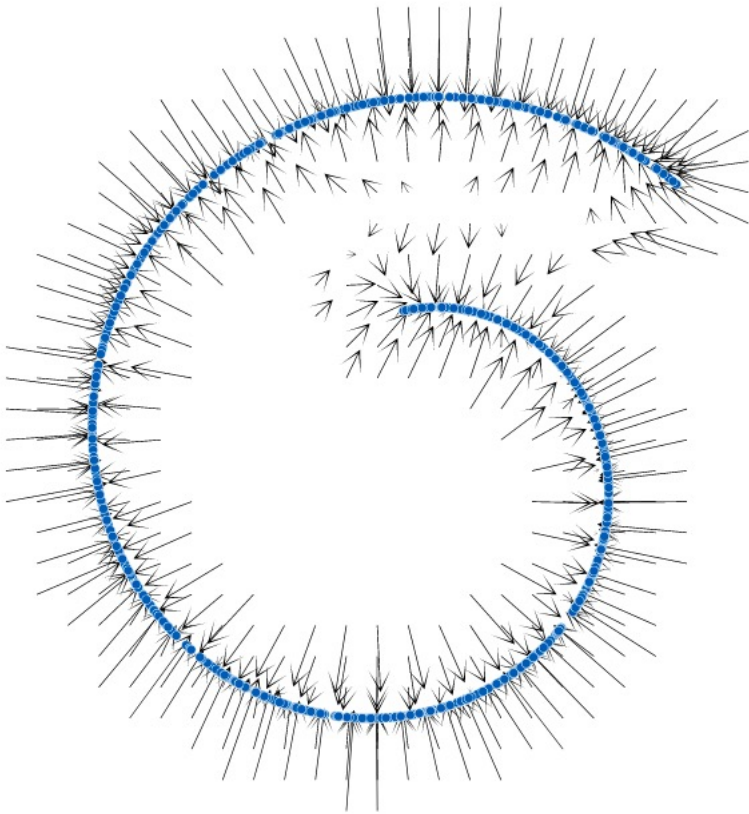
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ($\sigma = 4$)

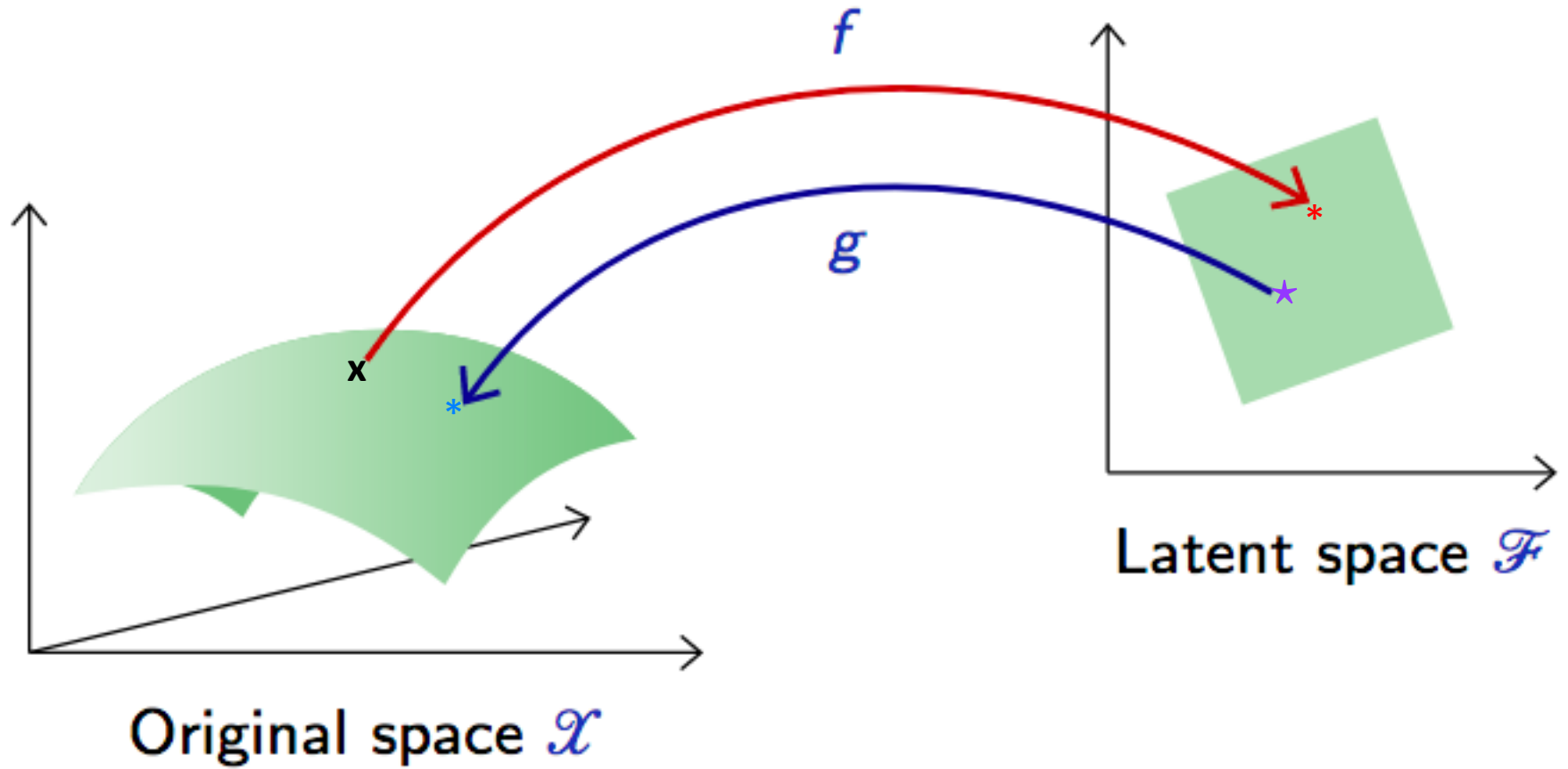
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

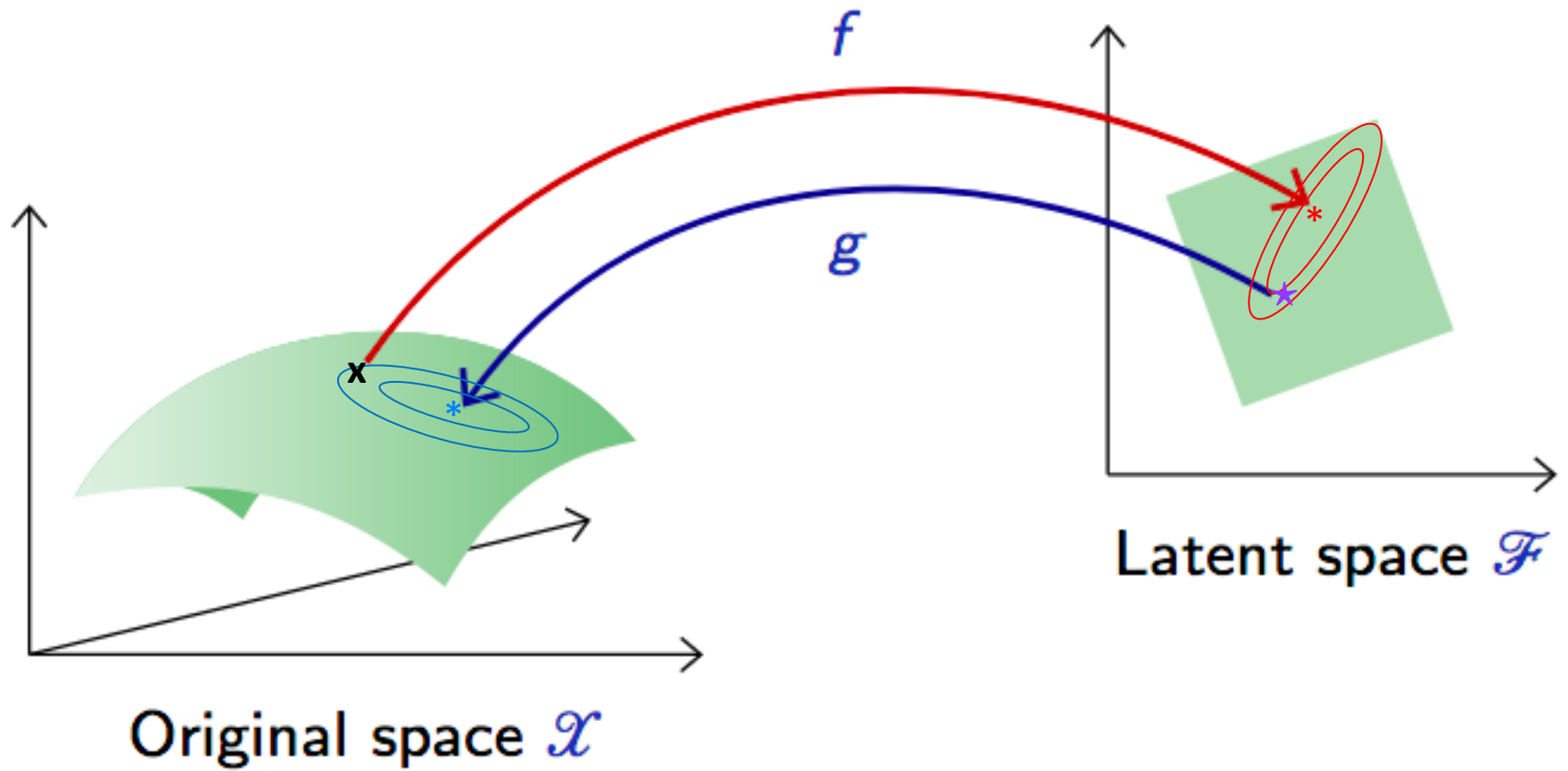
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2



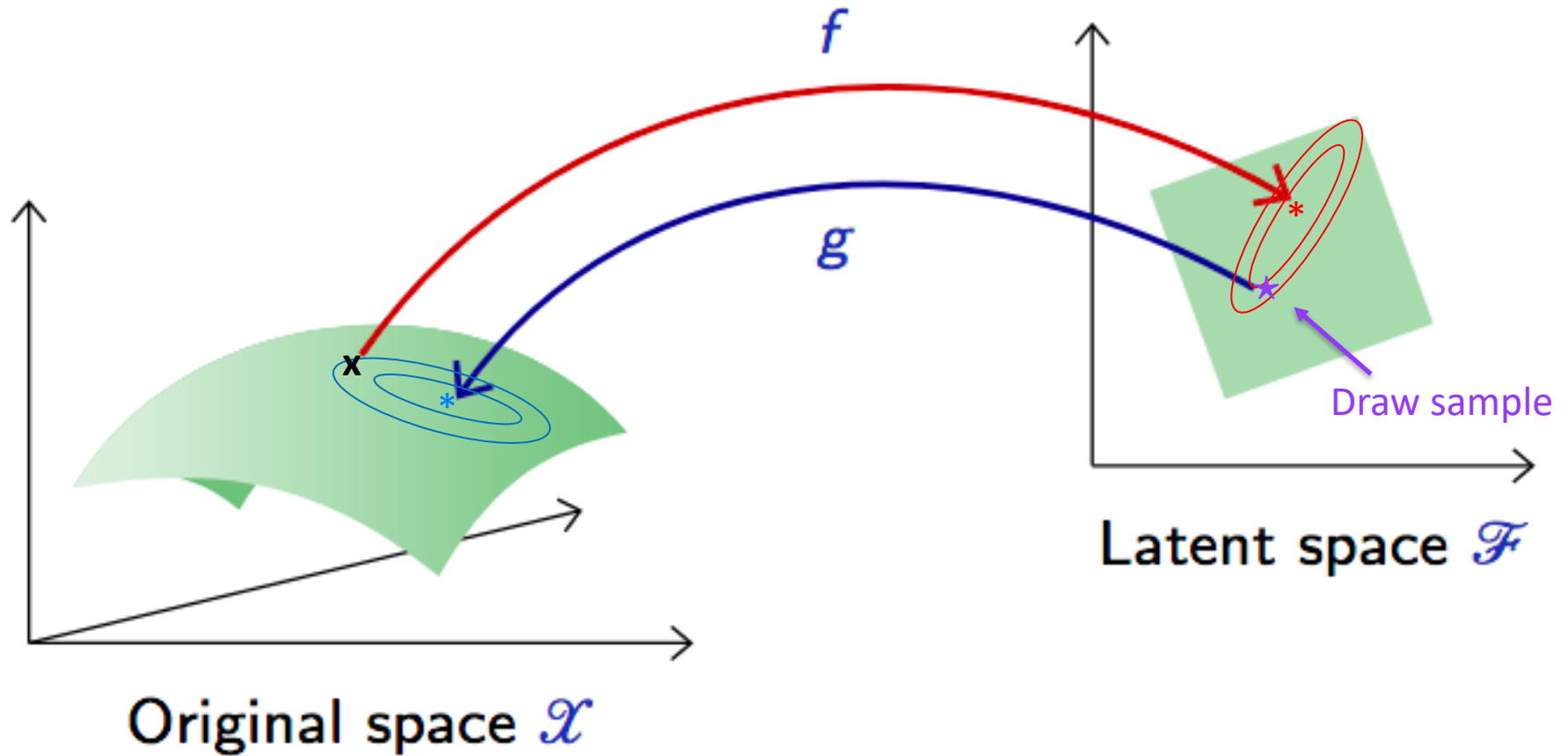
- **Autoencoder learns the average behavior**
- What if we care about these variations?
- Can we add a notion of variation in the autoencoder?

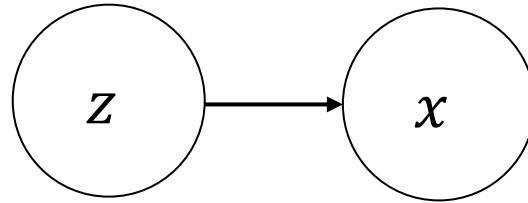


Variational Autoencoder



Variational Autoencoder





- Observed random variable x depends on unobserved latent random variable z
- Joint probability: $p(x, z) = p(x|z)p(z)$
- $p(x|z)$ is stochastic generation process from $z \rightarrow x$

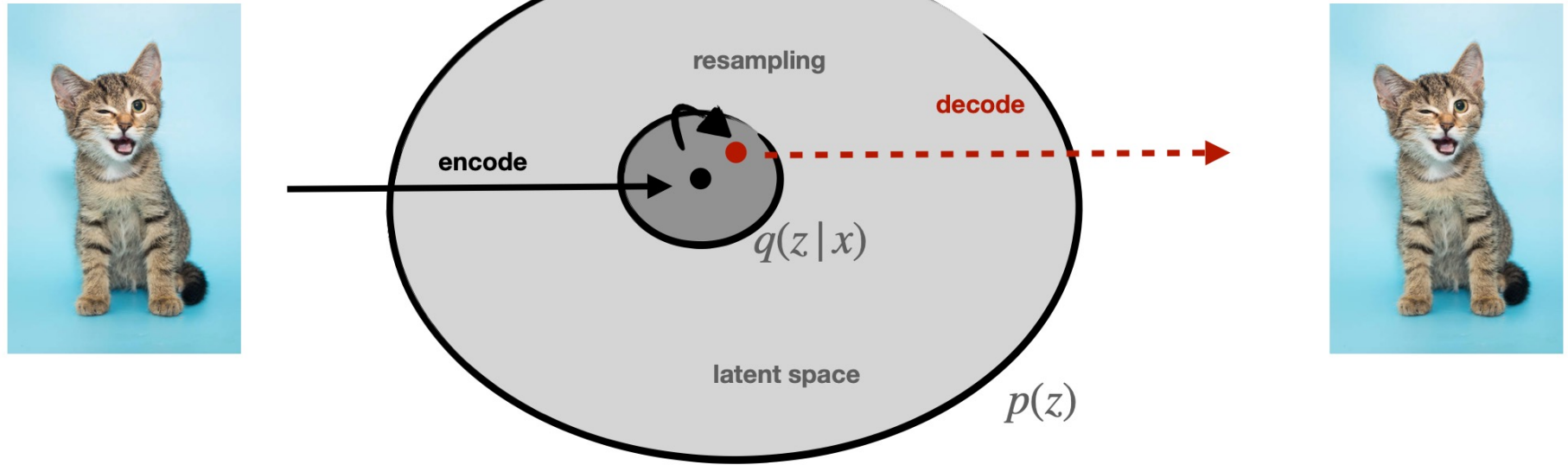
- Probabilistic relationship between data and latents

$$x, z \sim p(x, z) = p(x|z)p(z)$$

- Autoencoding

$$x \rightarrow q(z|x) \xrightarrow[\text{sample}]{} z \rightarrow p(x|z)$$

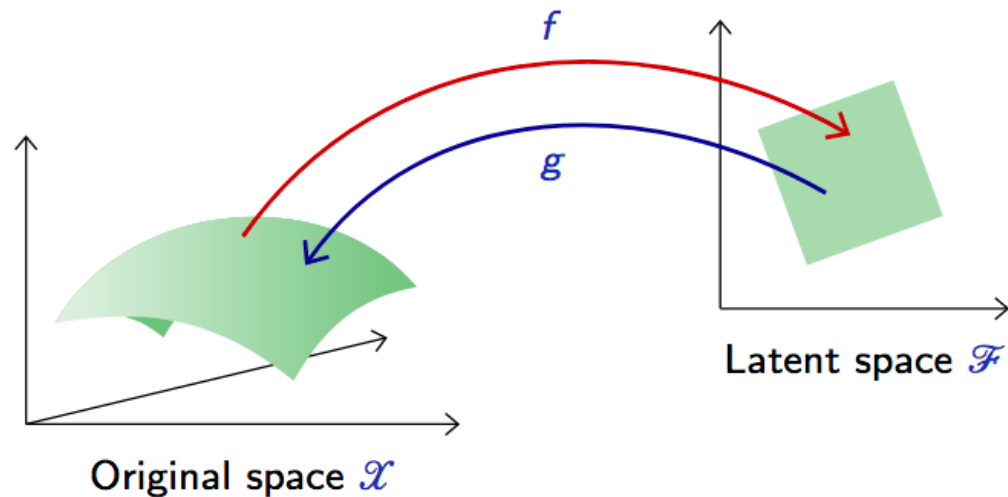
- **Encoder:** Learn what latents can produced data: $q(z|x)$
- **Decoder:** Learn what data is produced by latent: $p(x|z)$



- Close-by points must decode to similar images

- Typical encoder maps input x to “average” point in latent space

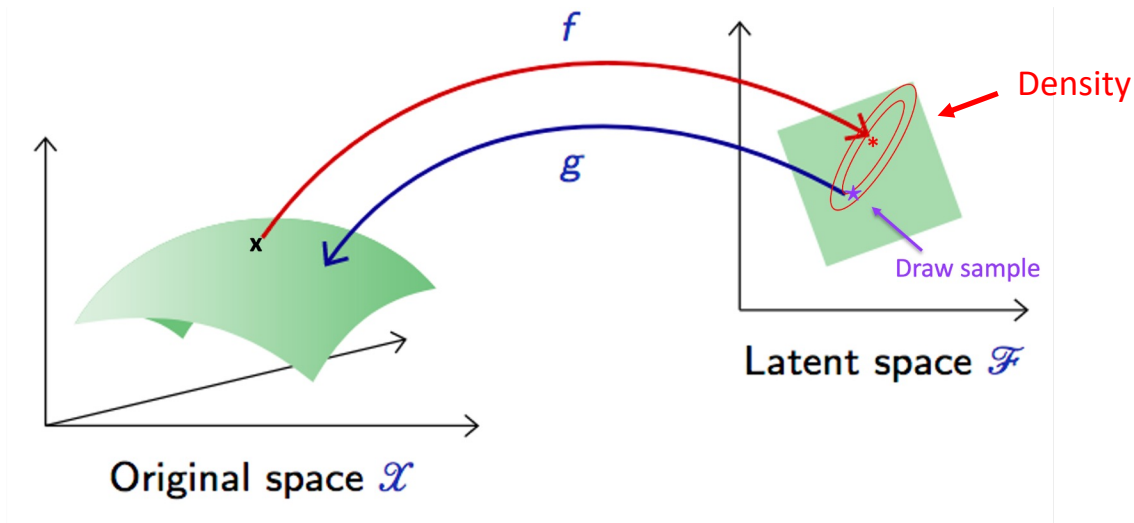
$$f(x) = \mu(x)$$



- A VAE Encoder has two outputs: mean & variance function

$$f_{\psi}(x) = \{\mu_{\psi}(x), \sigma_{\psi}(x)\}$$

ψ are parameters of the NN

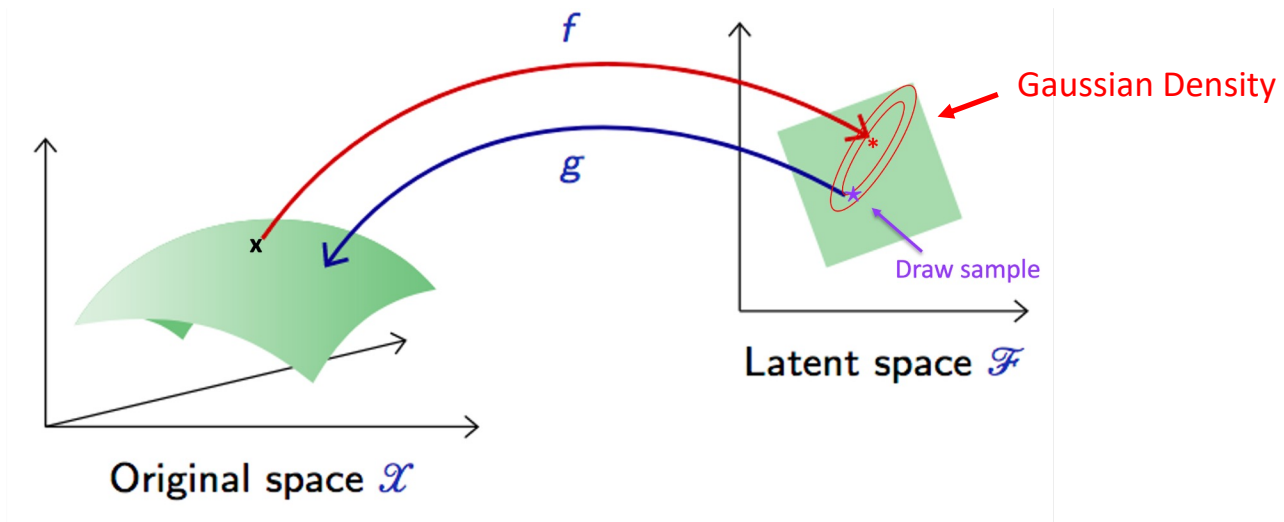


- A VAE Encoder has two outputs: mean & variance function

$$f_{\psi}(x) = \{\mu_{\psi}(x), \sigma_{\psi}(x)\} \quad \psi \text{ are parameters of the NN}$$

- What is the probability of a point in latent space?

$$p_{\psi}(z|x) = N(z | \mu_{\psi}(x), \sigma_{\psi}(x)) \quad \text{Could choose different density Gaussian is easiest}$$



- A VAE Encoder has two outputs: mean & variance function

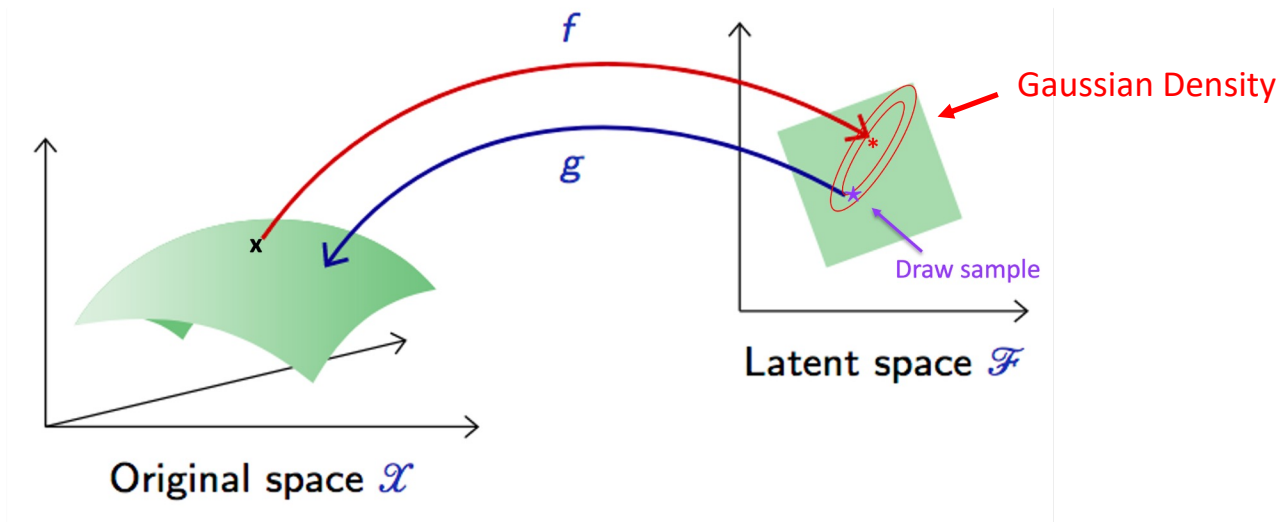
$$f_{\psi}(x) = \{\mu_{\psi}(x), \sigma_{\psi}(x)\} \quad \psi \text{ are parameters of the NN}$$

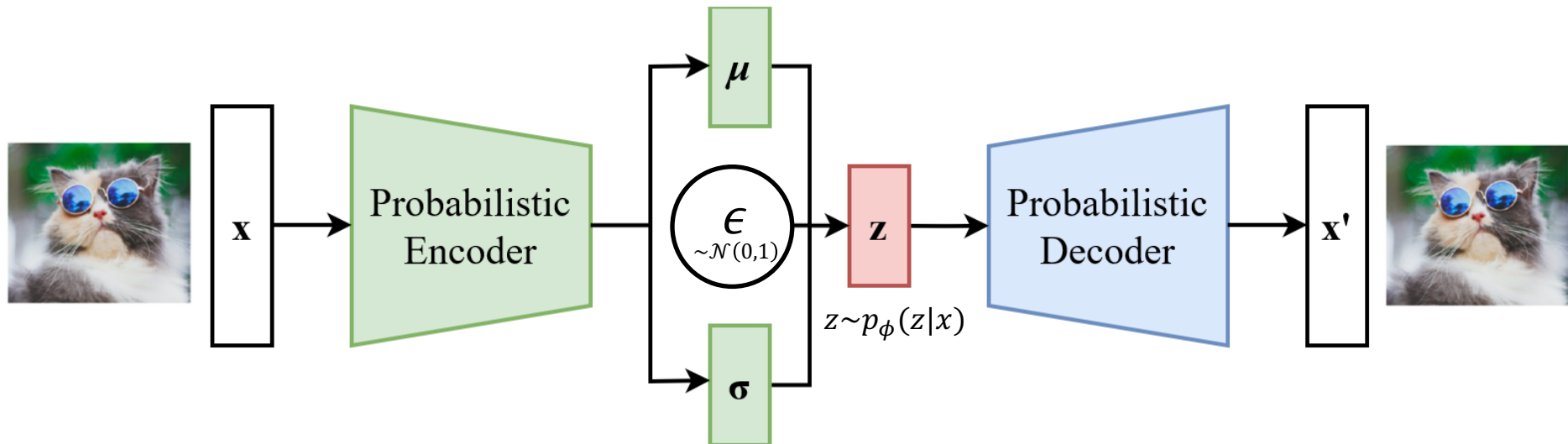
- What is the probability of a point in latent space?

$$p_{\psi}(z|x) = N(z | \mu_{\psi}(x), \sigma_{\psi}(x)) \quad \text{Could choose different density Gaussian is easiest}$$

- How do we draw a sample in latent space?

$$z = \sigma_{\psi}(x) * \epsilon + \mu_{\psi}(x) \quad \epsilon \sim N(0, I) \quad \text{Re-parameterization trick}$$





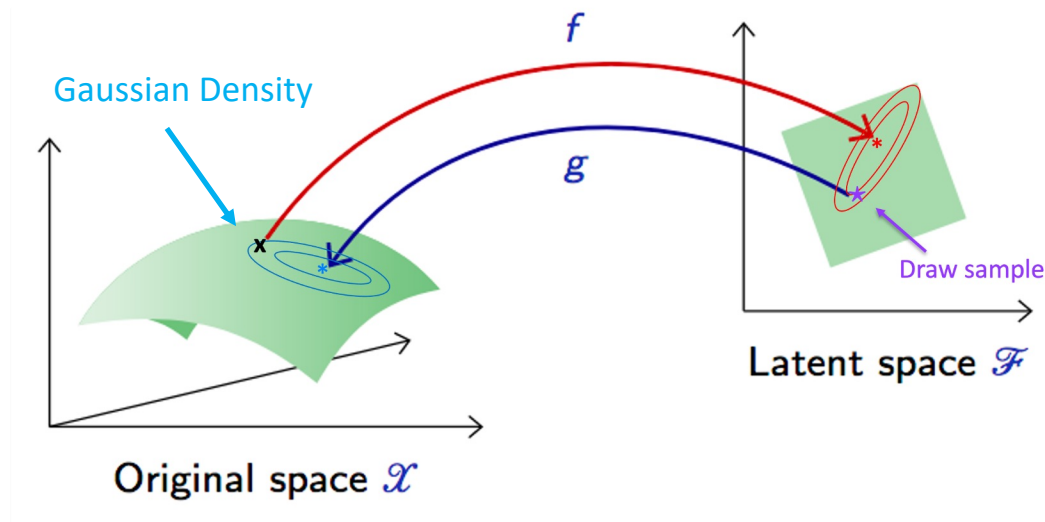
- Same as autoencoder

$$g_{\theta}(z) \equiv \mu_{\theta}(z)$$

θ are parameters of the NN

- Likelihood of an observation x

$$p_{\theta}(x|z) = N(x | \mu_{\theta}(z), I)$$



- Same as autoencoder

$$g_{\theta}(z) \equiv \mu_{\theta}(z)$$

θ are parameters of the NN

- Likelihood of an observation x

$$p_{\theta}(x|z) = N(x | \mu_{\theta}(z), I)$$

- **“Reconstruction Loss”**: Maximum likelihood

$$L_{reco} = \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)]$$

- Same as autoencoder

$$g_{\theta}(z) \equiv \mu_{\theta}(z)$$

θ are parameters of the NN

- Likelihood of an observation x

$$p_{\theta}(x|z) = N(x | \mu_{\theta}(z), I)$$

- **“Reconstruction Loss”**: Maximum likelihood

$$L_{reco} = \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] \approx \frac{1}{N} \sum_{z_i \sim q(z|x)} \log N(x | g_{\theta}(z_i), I)$$

- Same as autoencoder

$$g_{\theta}(z) \equiv \mu_{\theta}(z)$$

θ are parameters of the NN

- Likelihood of an observation x

$$p_{\theta}(x|z) = N(x | \mu_{\theta}(z), I)$$

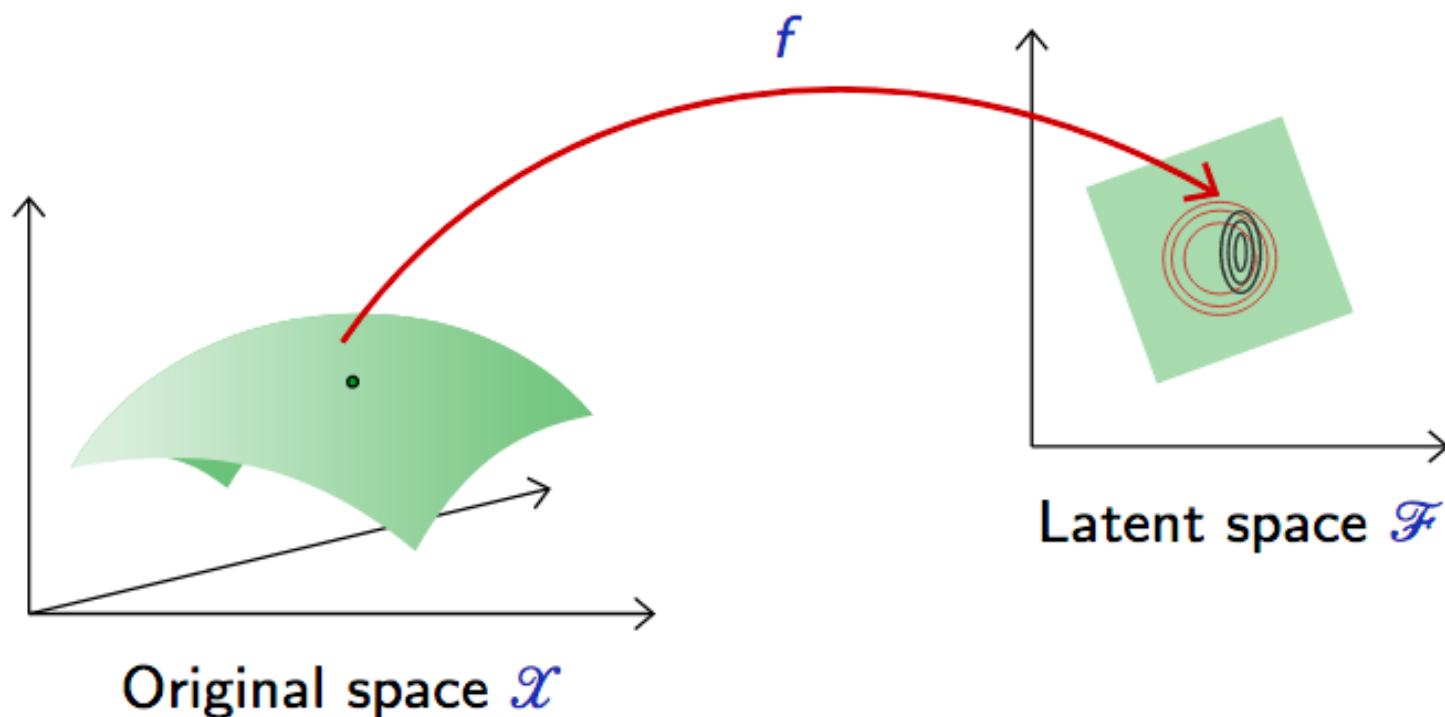
- **“Reconstruction Loss”**: Maximum likelihood

$$L_{reco} = \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] \approx -\frac{1}{N} \sum_{z_i \sim q(z|x)} (x - g_{\theta}(z_i))^2$$

Same as the autoencoder loss

- How do we make sure system doesn't collapse to an autoencoder (i.e. VAE encoder only predicts mean)?

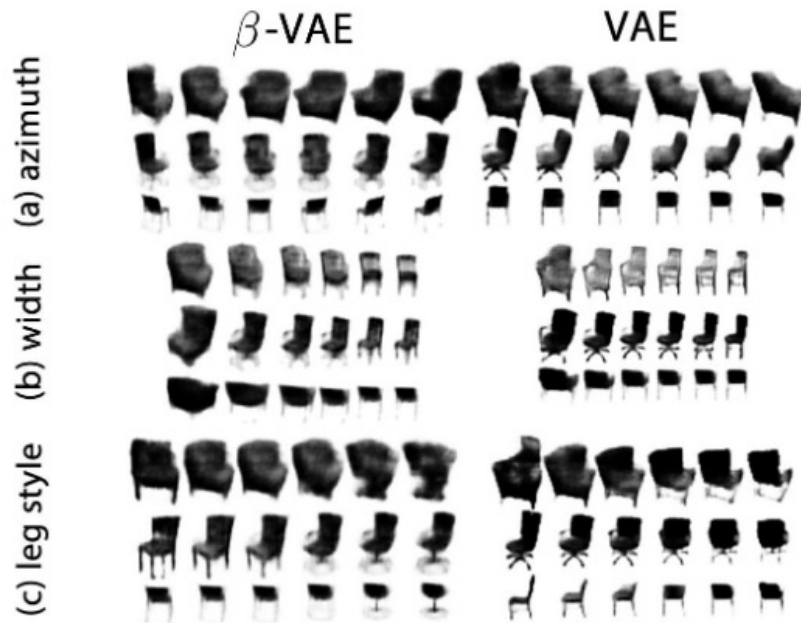
- How do we make sure system doesn't collapse to an autoencoder (i.e. VAE encoder only predicts mean)?
- Use prior $p(z)$ for the latent space distribution, **need to ensure the encoder is consistent with prior**



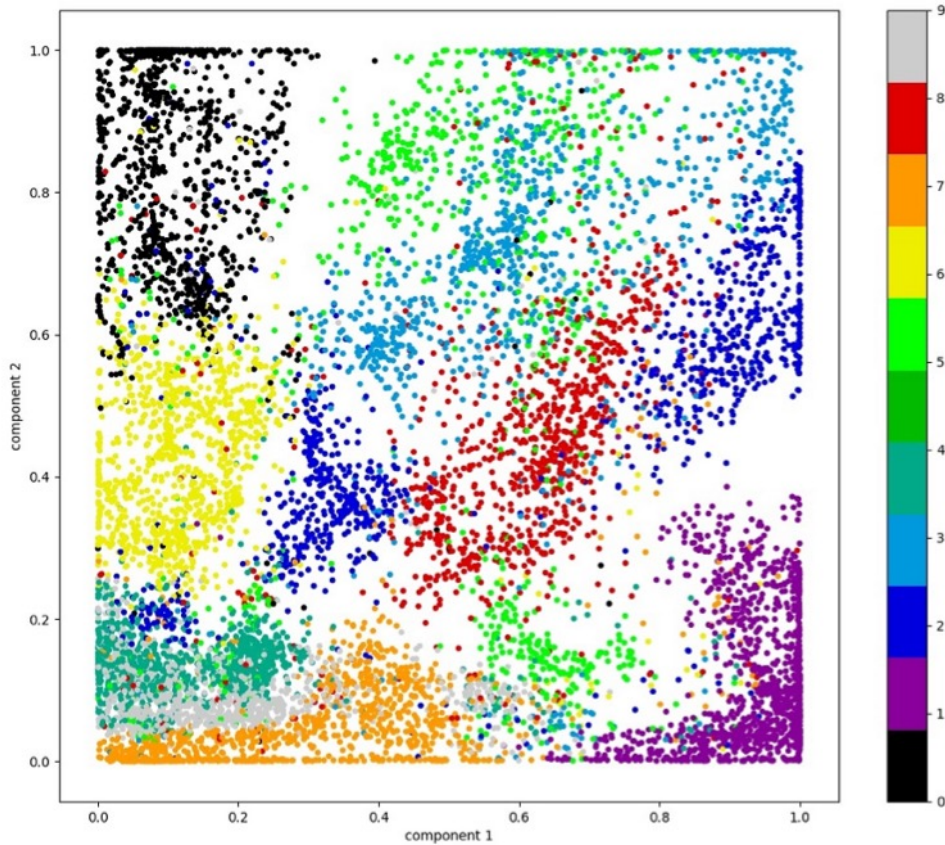
- Constrain difference between distributions with **Kullback–Leibler divergence**

$$D_{KL}[q(z|x)|p(z)] = \mathbb{E}_{q(z|x)} \left[\log \frac{q(z|x)}{p(z)} \right] = \int q(z|x) \log \frac{q(z|x)}{p(z)} dz$$

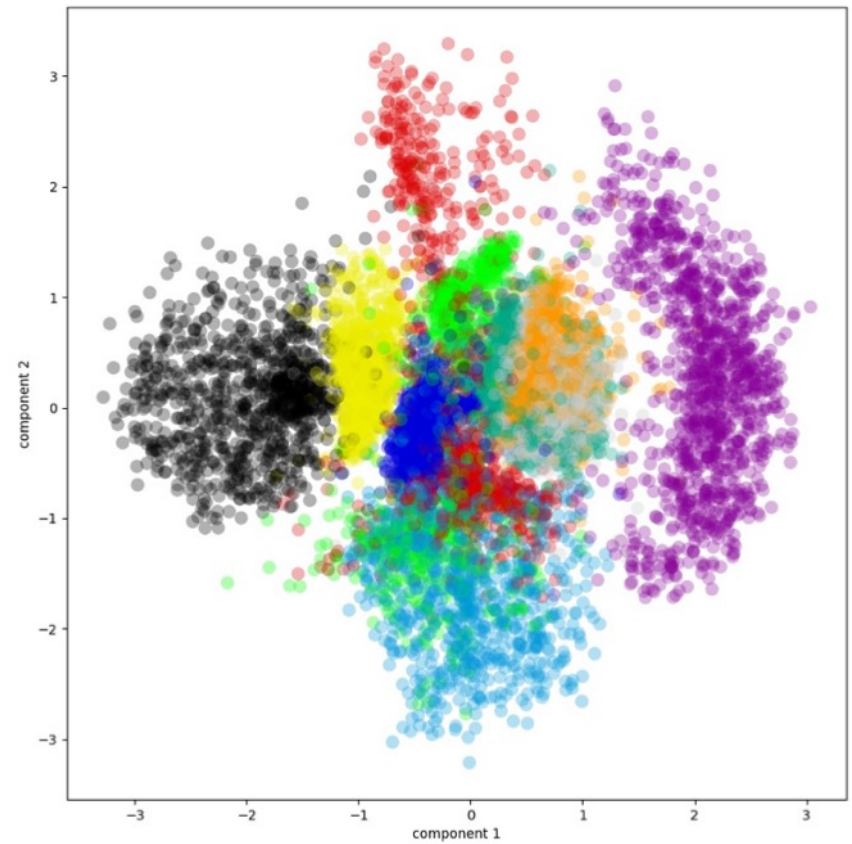
– $D_{KL}[q|p] \geq 0$ and is only 0 when $q = p$



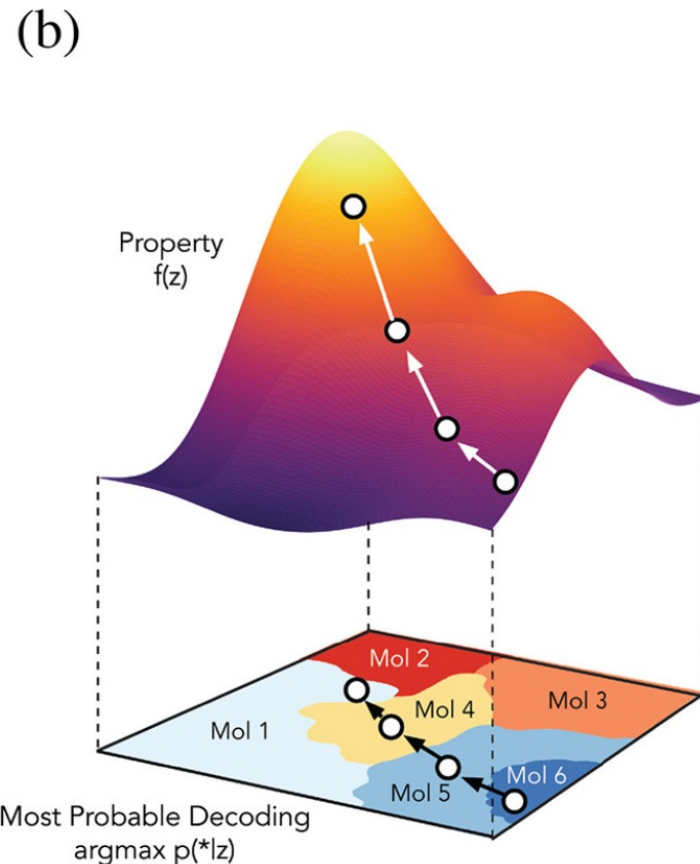
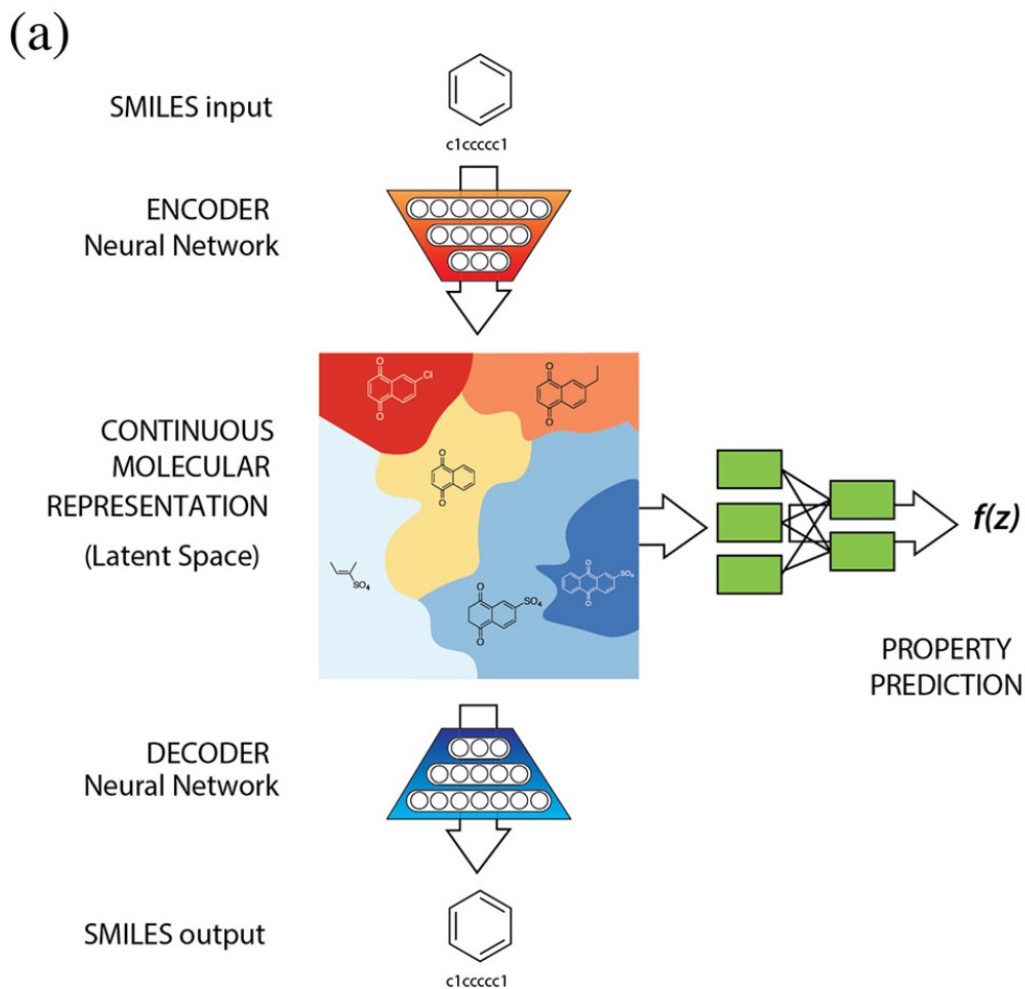
Autoencoder



Variational Autoencoder



Data: MNIST data set of hand-written digits

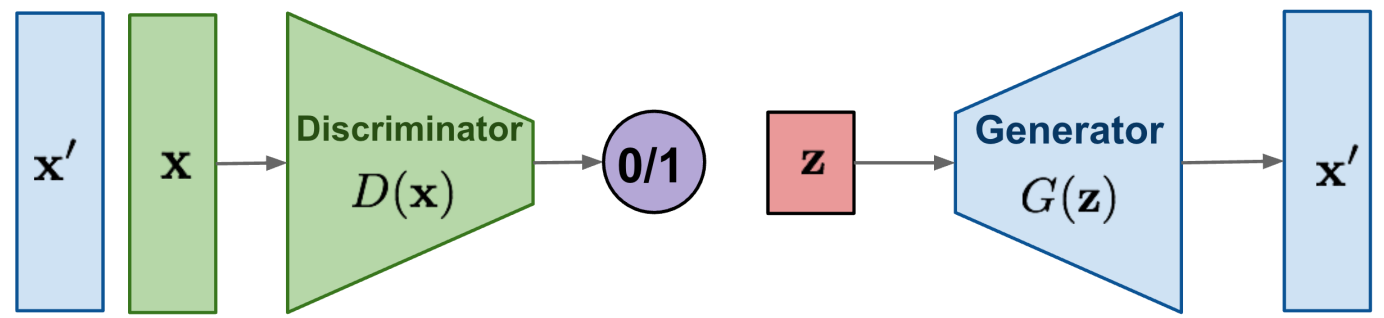


Design of new molecules with desired chemical properties.
(Gomez-Bombarelli et al, 2016)

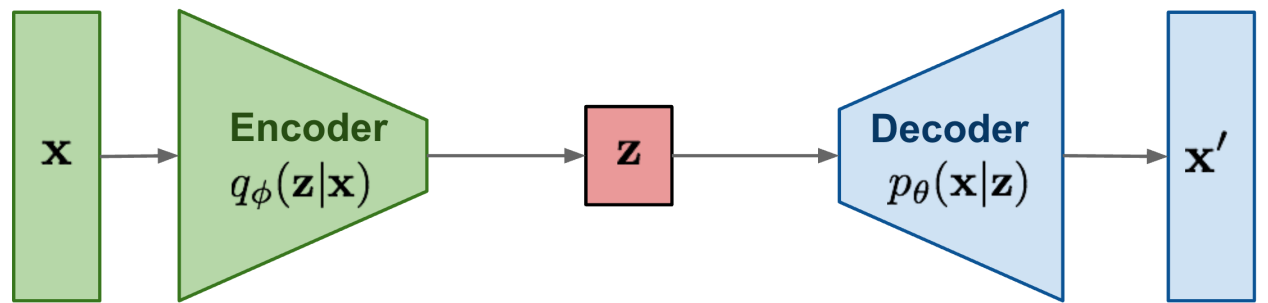
- In generative modeling, want to learn the lower dimensional degrees of freedom that describe the features of the data
- “Degrees of freedom” are modeled with a latent distribution (kept simple for convenience) and complex neural network mappings
- Need to think about **probabilistic systems**
- Design loss around this probabilistic model

The Zoo of Generative Models...

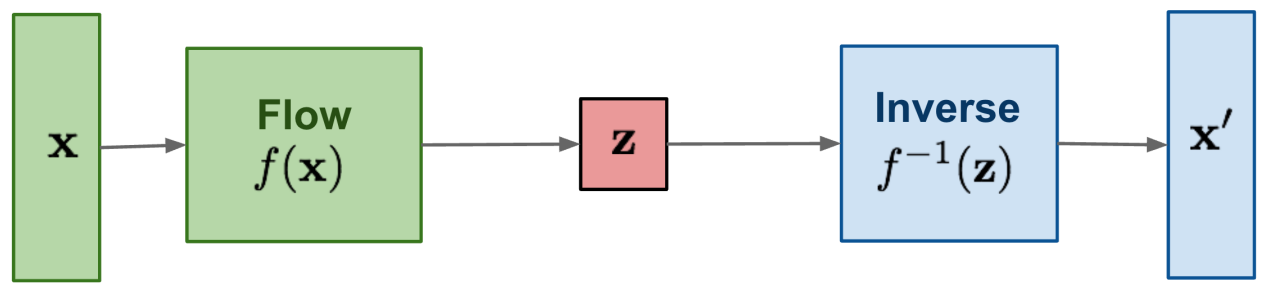
GAN: Adversarial training



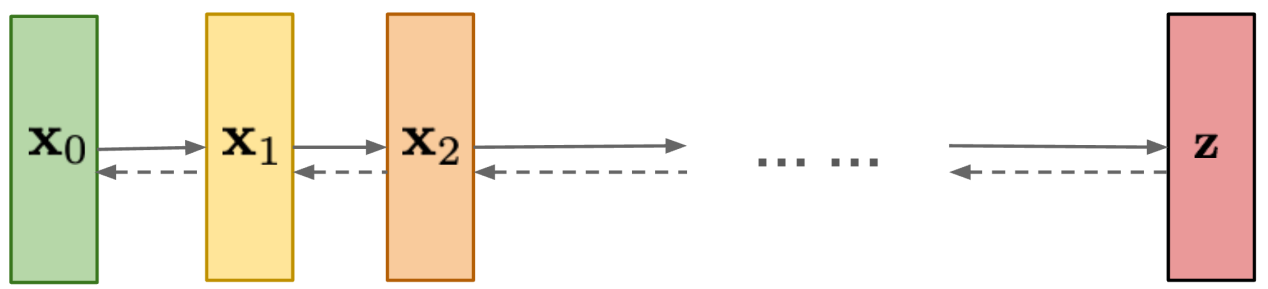
VAE: maximize variational lower bound



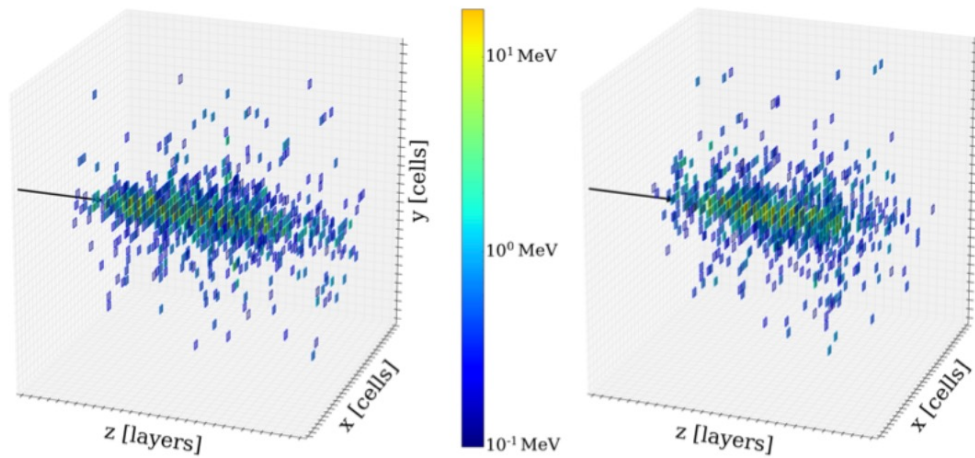
Flow-based models: Invertible transform of distributions



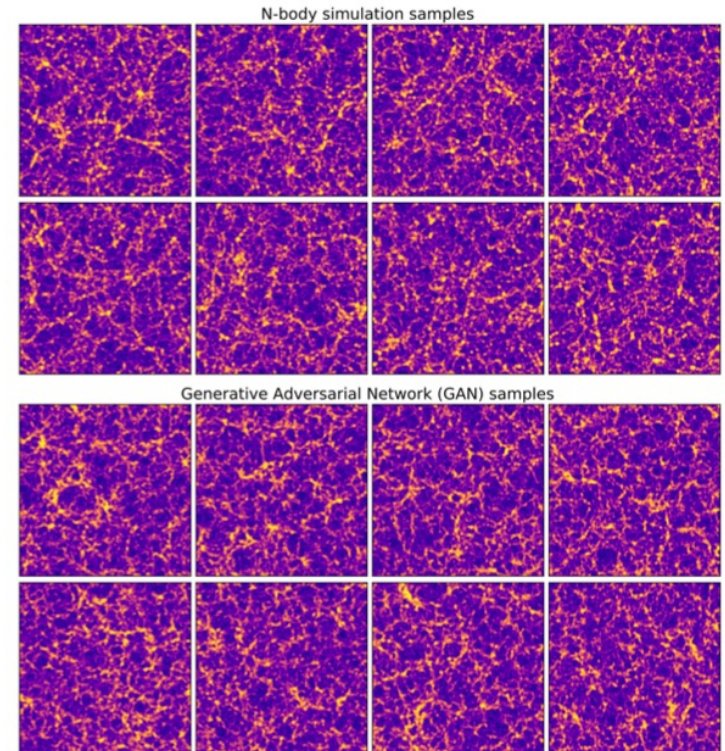
Diffusion models: Gradually add Gaussian noise and then reverse



- Often studied for fast approximate simulation, simulation-based inference, optimization, anomaly detection, ...
 - See talks by [G. Kasieczka](#), [D. Shih](#), [B. Nevin](#), [A. Edelen](#)



[2005.05334](#)



[1801.09070](#)

- Deep neural networks are an extremely powerful class of models
- We can express our inductive bias about a system in terms of model design, and can be adapted to a many types of data
- Even beyond classification and regression, deep neural networks allow powerful unsupervised learning and Generative modeling!