

Introduction to Machine Learning:

Lecture 5 – Deep Generative Models

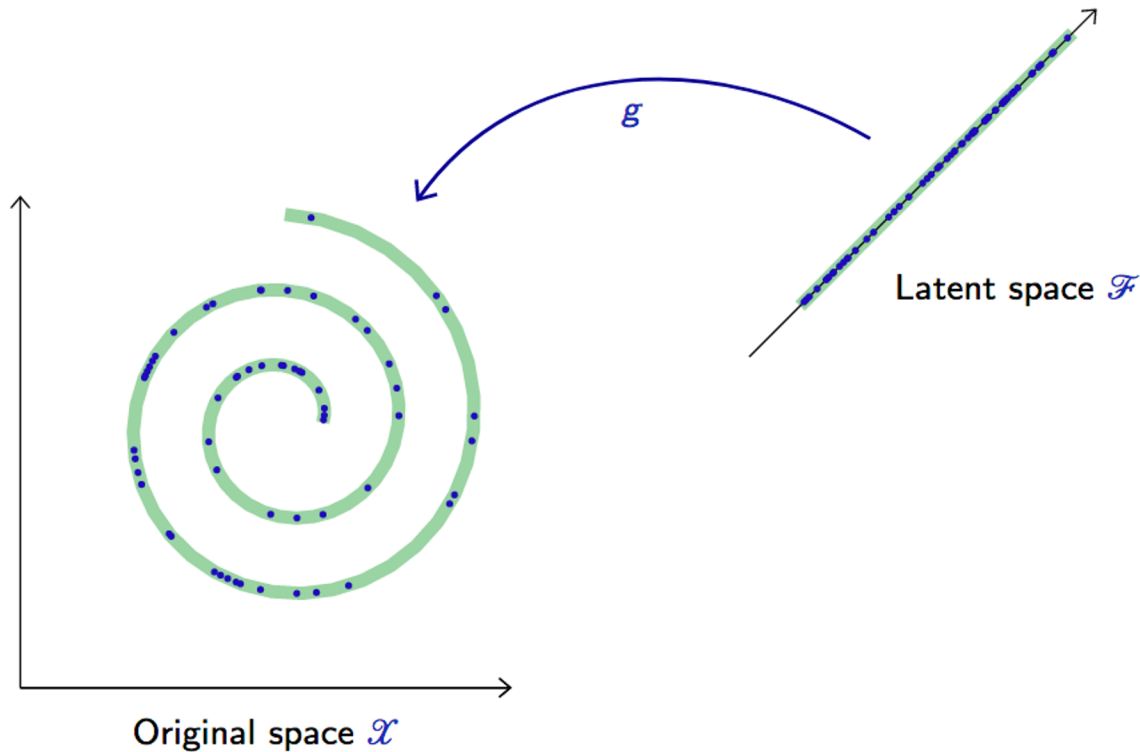
Michael Kagan



TRISEP Summer School
July 8-12, 2024

- Lecture 1 – Machine Learning Fundamentals
- Lecture 2 – Intro to Neural Networks
- Lecture 3 – Intro to Deep Learning
- Lecture 4 – Intro to Unsupervised Learning
- Lecture 5 – Intro to Deep Generative Models

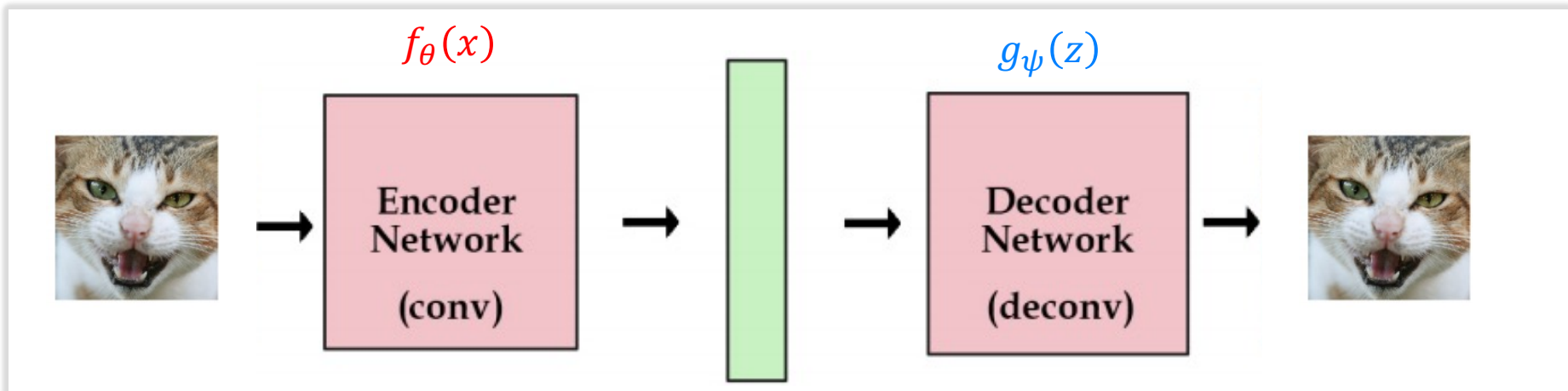
How can we find the “meaningful degrees of freedom” in the data?

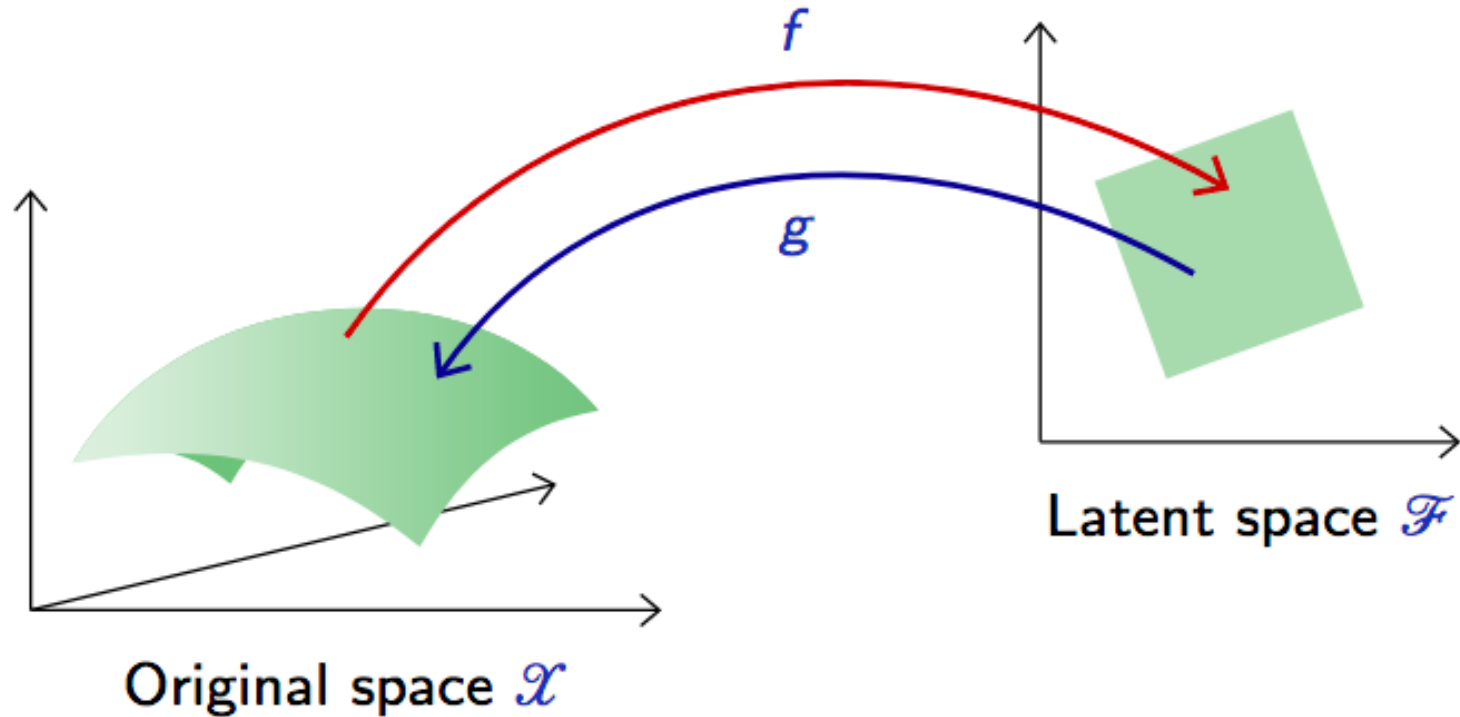


- Map a space to itself through a compression

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

- **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

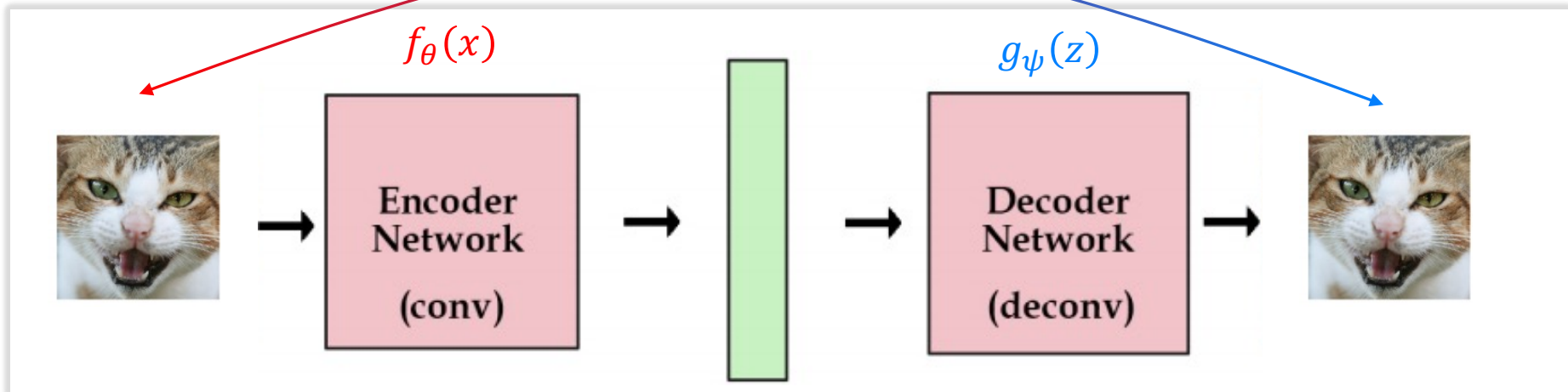




- Latent space is of lower dimension than data
- Model must learn a “good” parametrization and capture dependencies between components

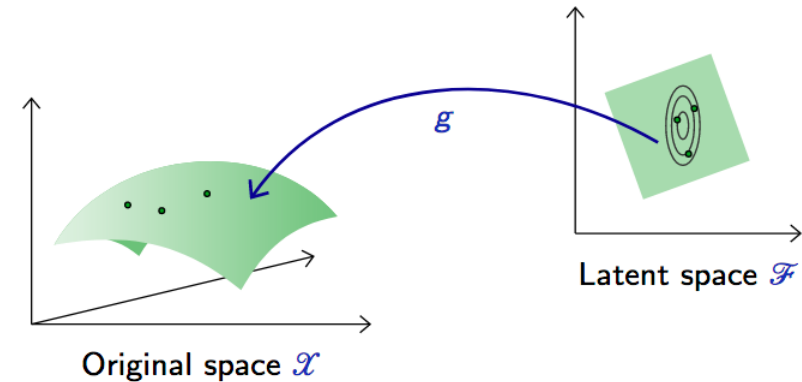
$$L(\theta, \psi) = \frac{1}{N} \sum_n \|x_n - g_\psi(f_\theta(x_n))\|^2$$

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



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 ($d = 16$)



- Don't know the right latent space density

- Autoencoders learn the latent space, but we don't know what is the latent space distribution
- Autoencoder prescribes a deterministic relationship between data space and latent space
- One set of “meaningful degrees of freedom” can only describe one data space point

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

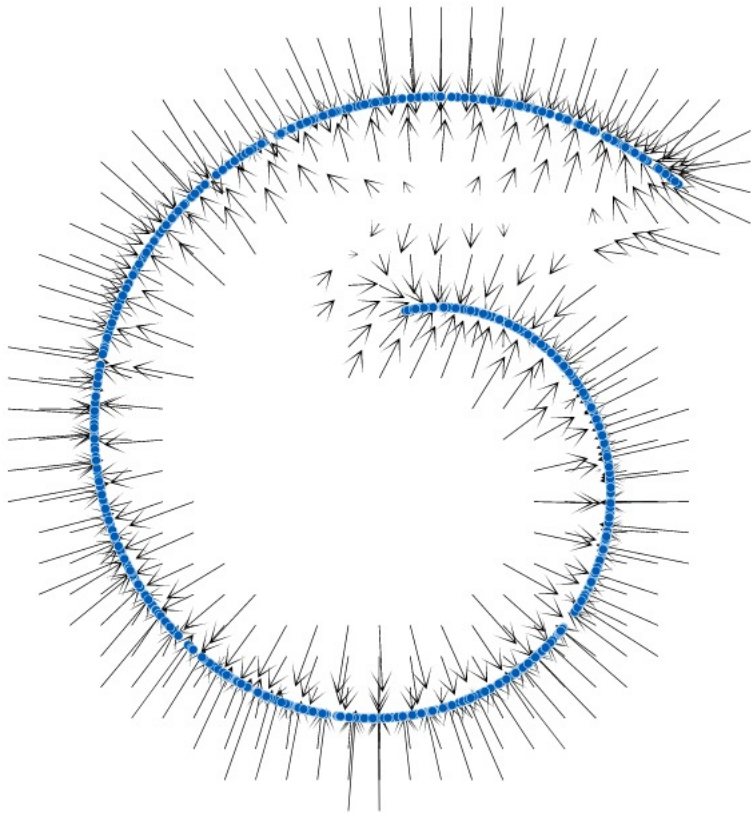
- Loss:

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



Perturbation, e.g. Gaussian noise

Denoising Autoencoders Examples



Original

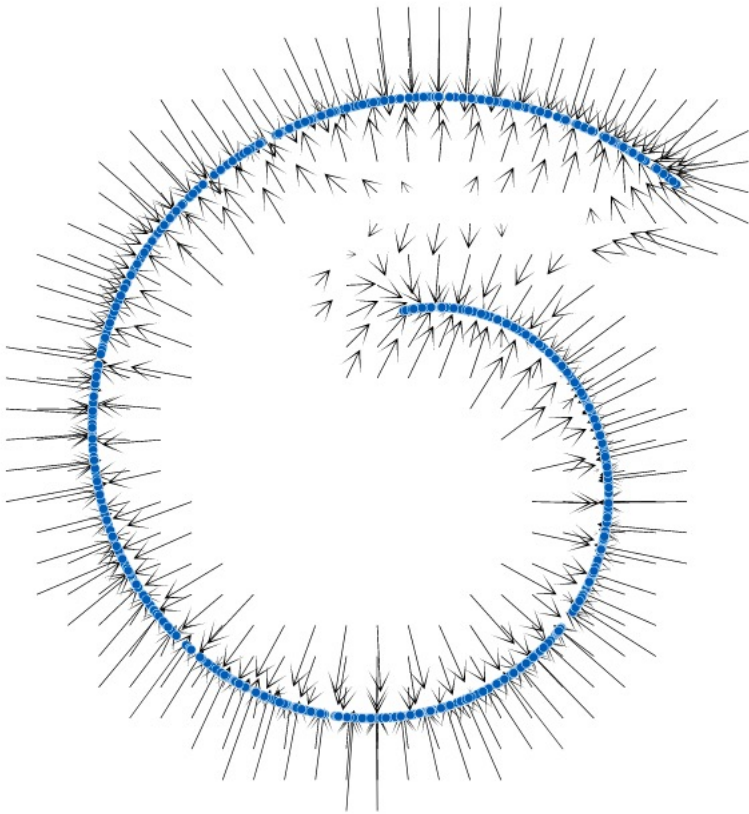
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ($\sigma = 4$)

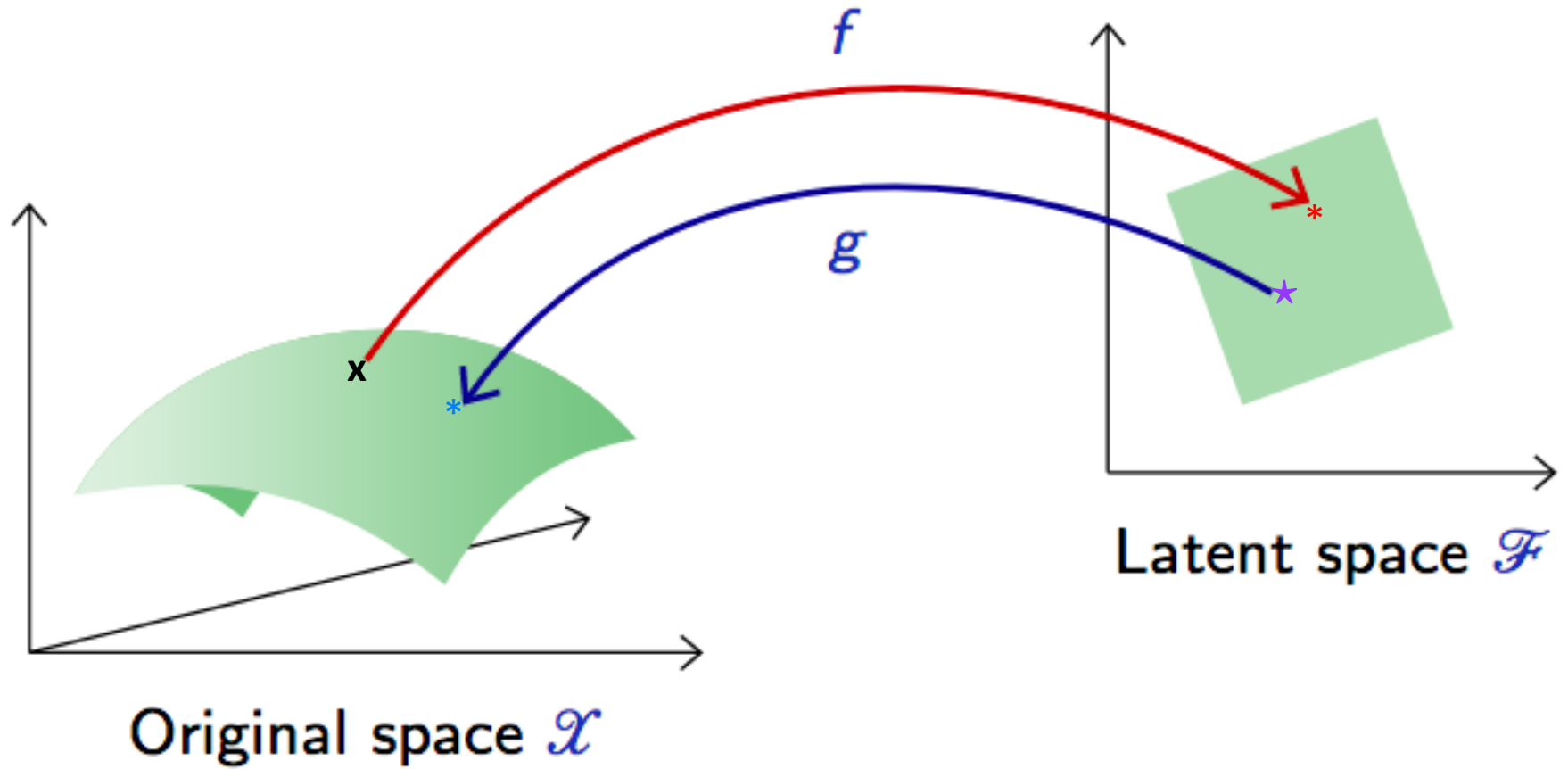
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

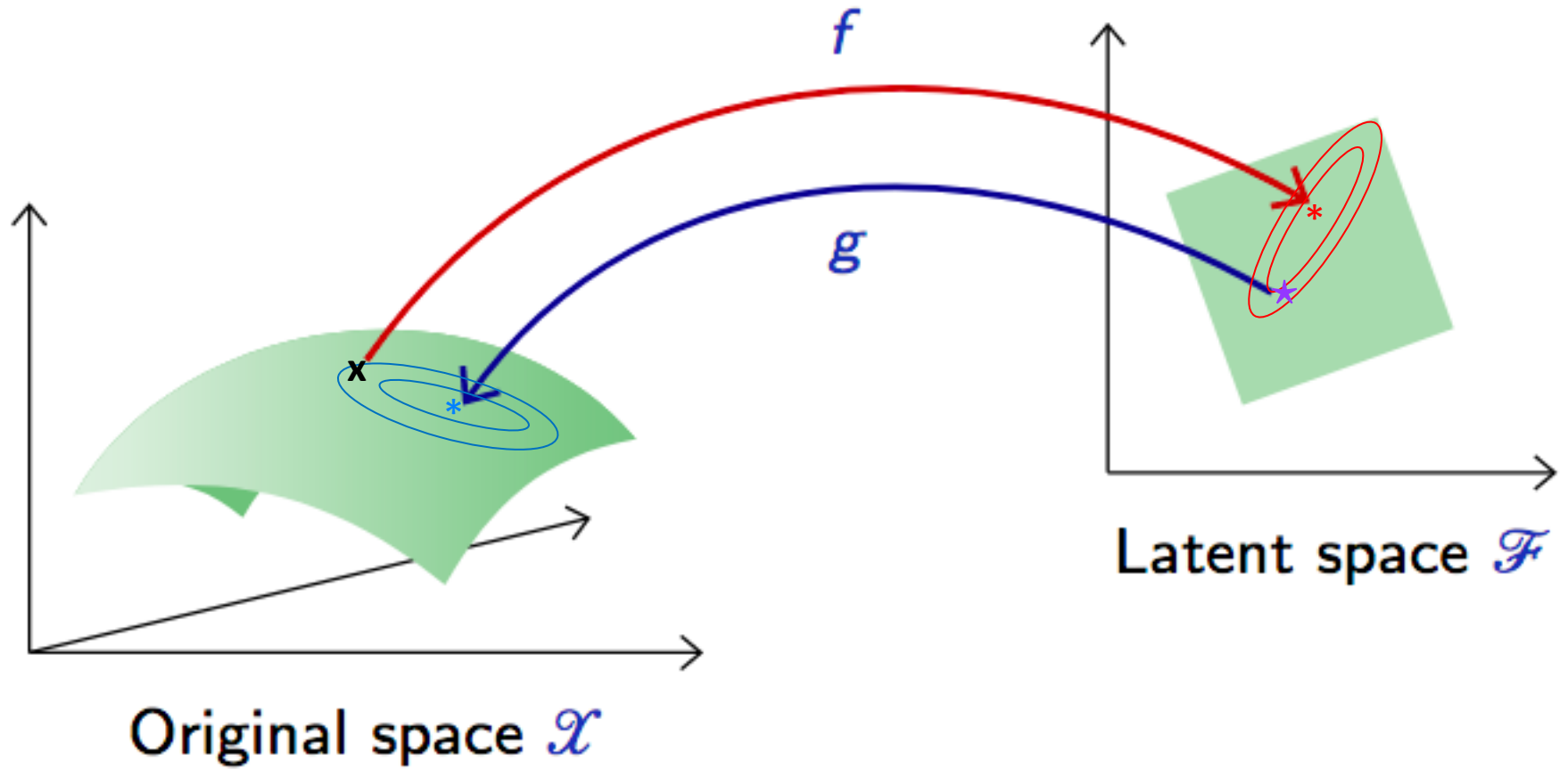
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2



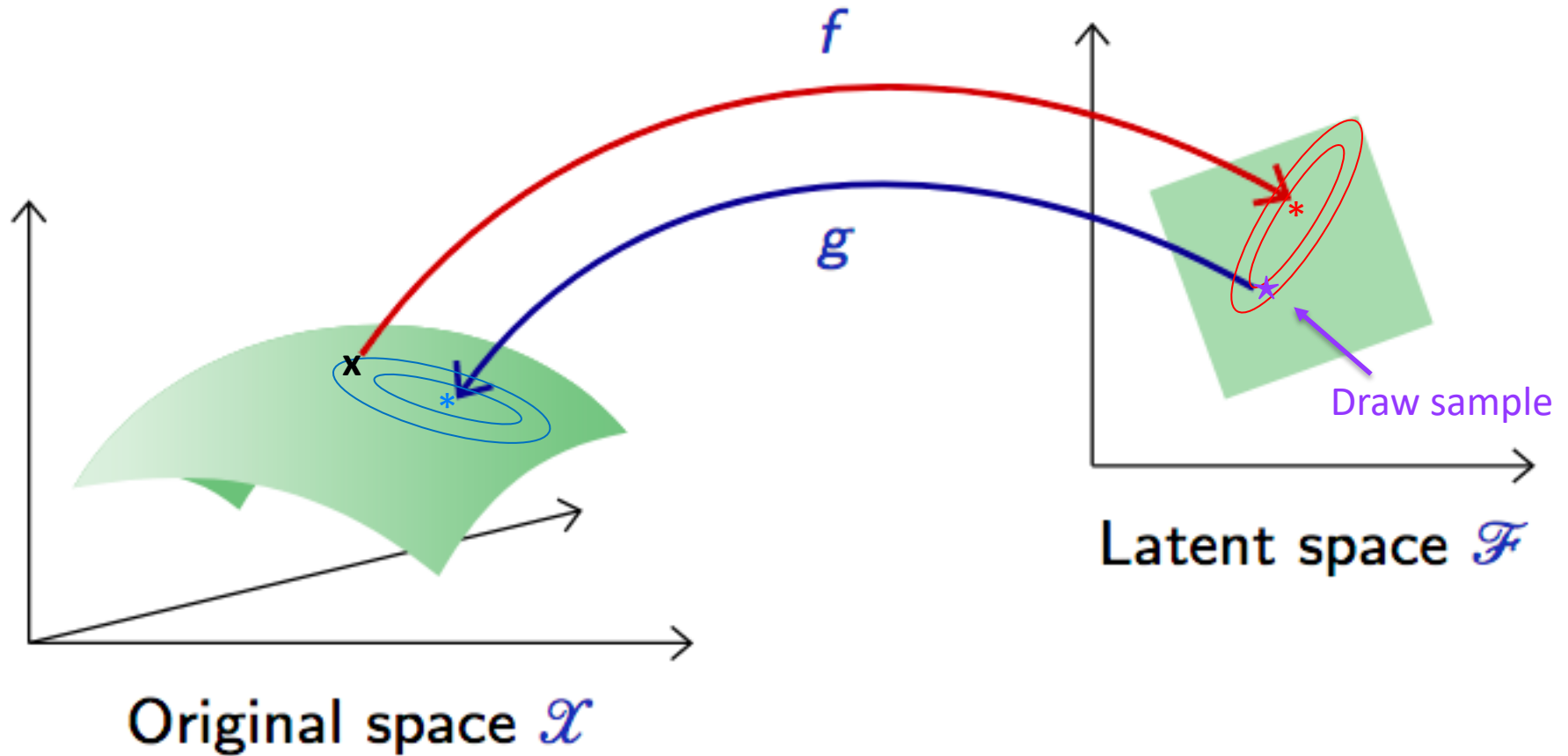
- **Autoencoder learns the average behavior**
- What if we care about these variations?
- Can we add a notion of variation in the autoencoder?

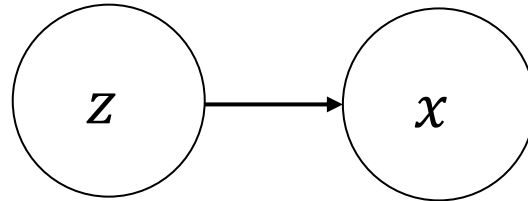


Variational Autoencoder



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 $z \rightarrow x$

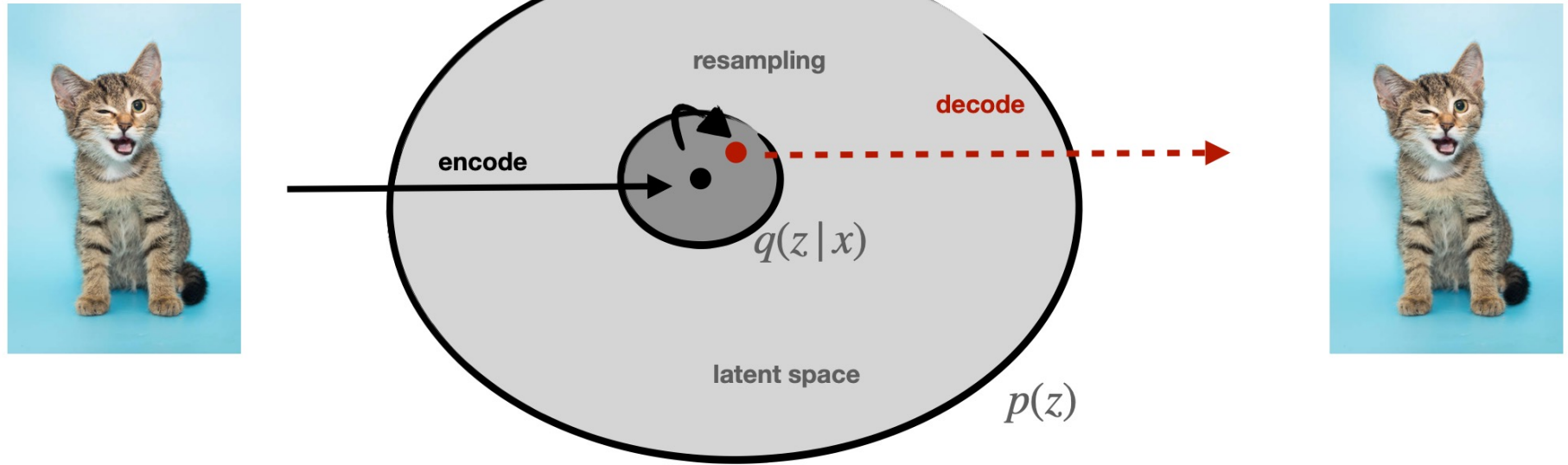
- Probabilistic relationship between data and latents

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

- 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

- Classification / regression models make single prediction

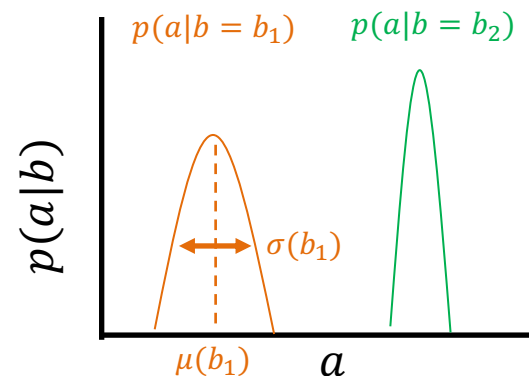
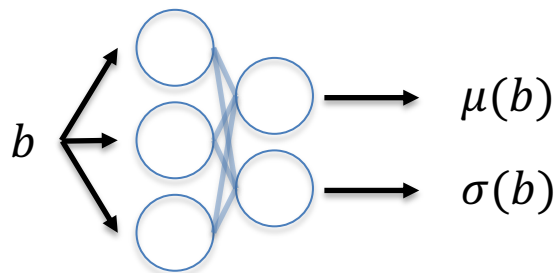
How to model a conditional density $p(a|b)$?

- Assume a known form of density, e.g. normal

$$p(a|b) = \mathcal{N}(a; \mu(b), \sigma(b))$$

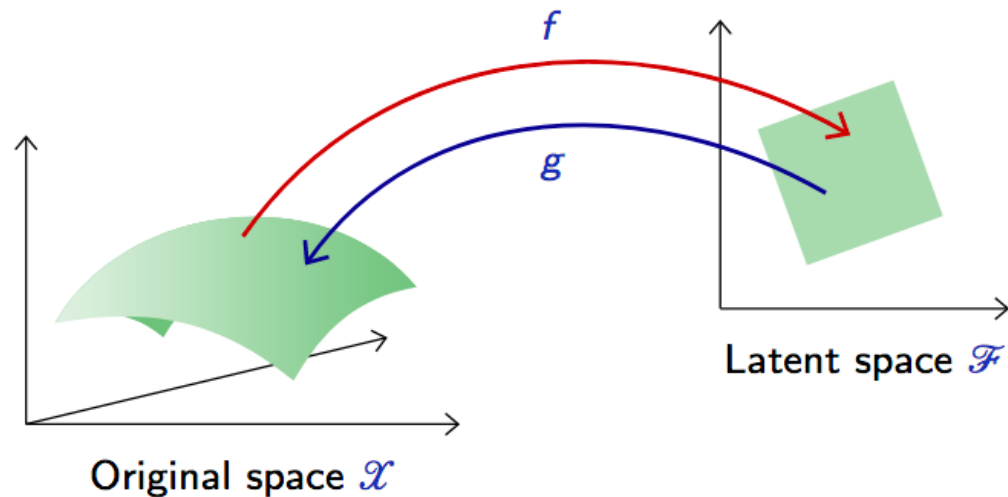
– Parameters of density depend on conditioned variable

- Use neural network to model density parameters



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

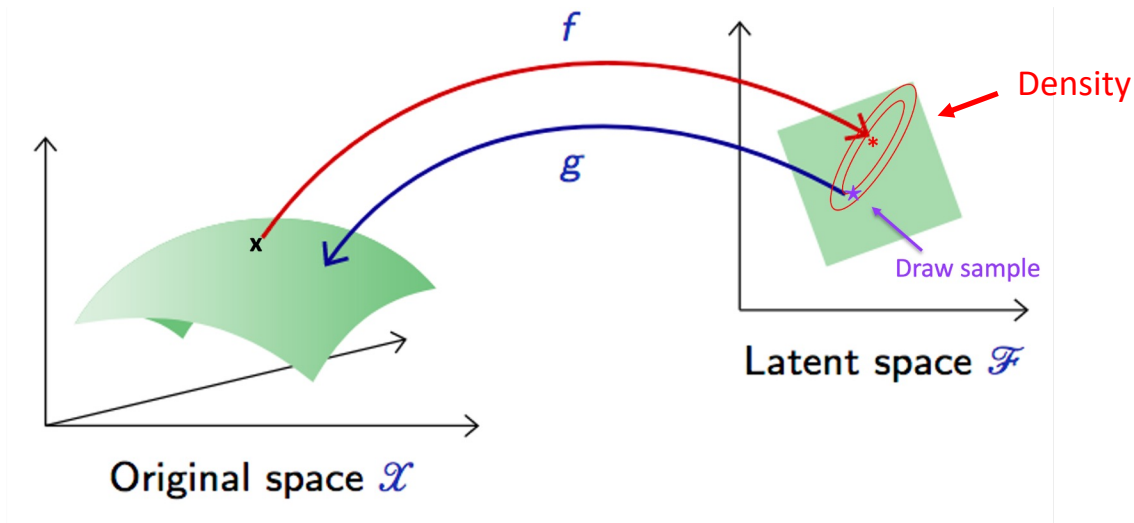
$$f(x) = \mu(x)$$



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

$$f_{\psi}(x) = \{\mu_{\psi}(x), \sigma_{\psi}^2(x)\}$$

ψ are parameters of the NN

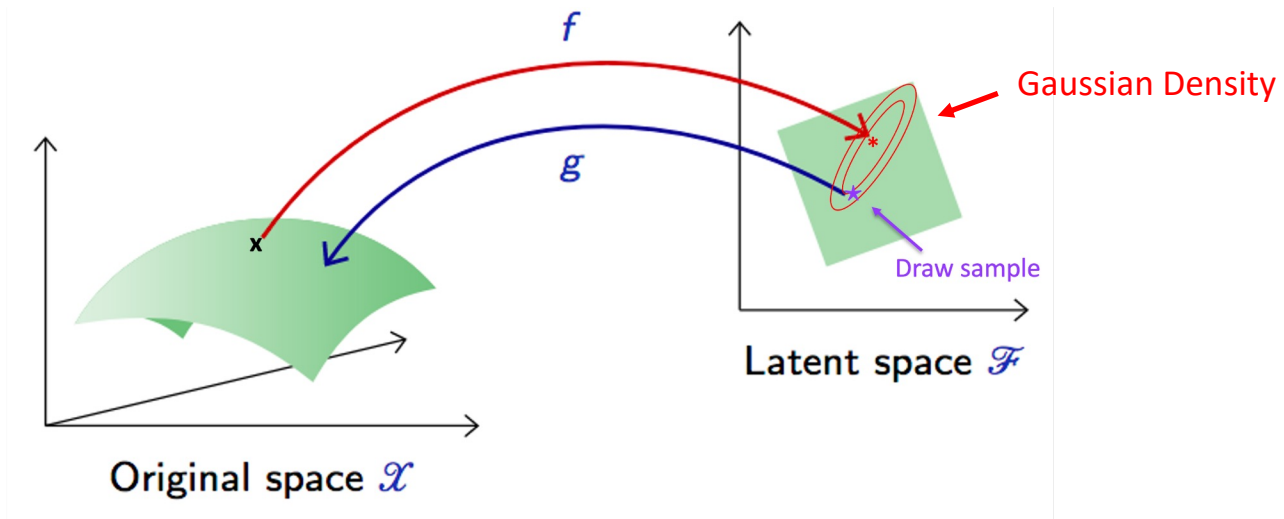


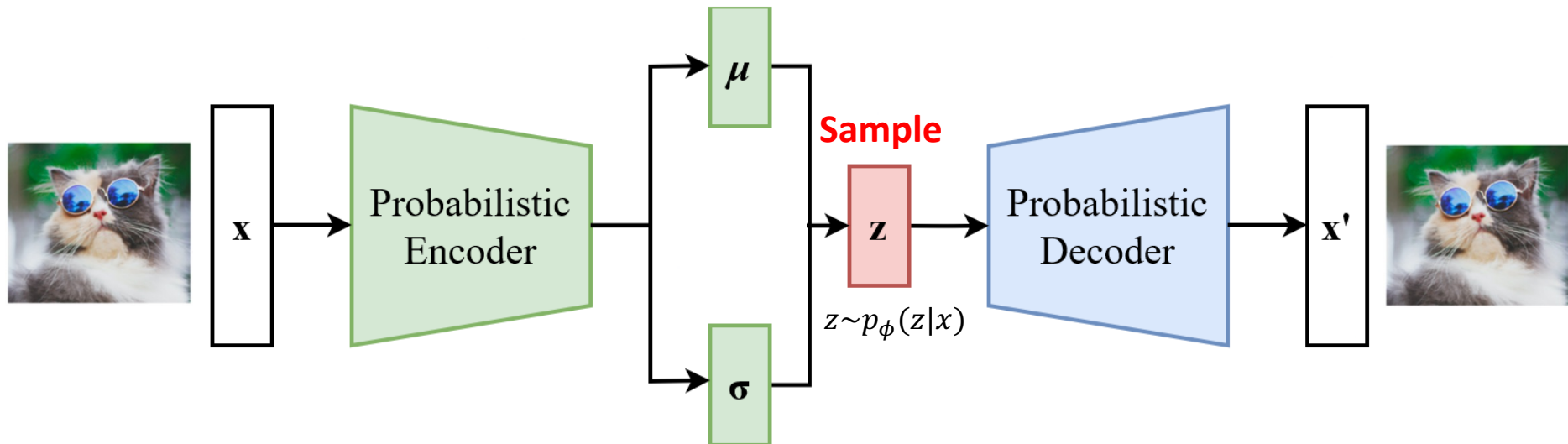
- A VAE Encoder has two outputs: mean & variance function

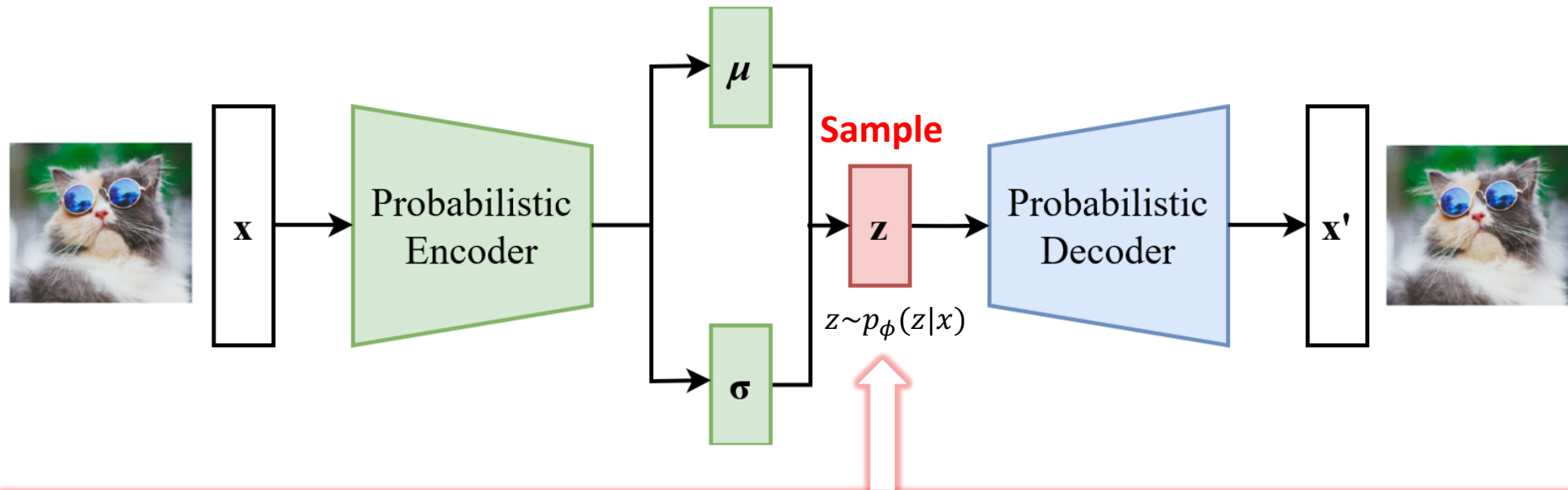
$$f_{\psi}(x) = \{\mu_{\psi}(x), \sigma_{\psi}^2(x)\} \quad \psi \text{ are parameters of the NN}$$

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

$$p_{\psi}(z|x) = N(z | \mu_{\psi}(x), \sigma_{\psi}^2(x)) \quad \text{Could choose different density Gaussian is easiest}$$







But for training:

How do we take a derivative through a randomly sampled number?

How do we know the dependence on the parameters?

- Given $x \sim p(x|\theta)$
- Sometimes, we can rewrite x as a function of the parameters and a simpler distribution without parameter dependence

$$x = g(\epsilon, \theta) \quad \epsilon \sim p(\epsilon)$$

- Example:

$$x \sim N(x|\mu, \sigma) \rightarrow x = \sigma * \epsilon + \mu \quad \text{with } \epsilon \sim N(0,1)$$

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

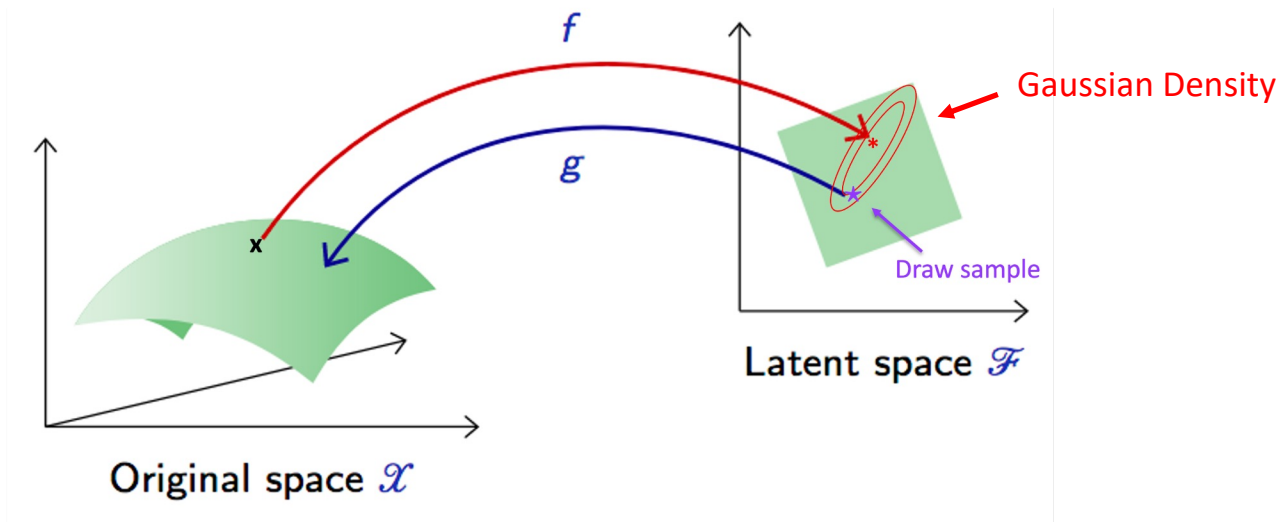
$$f_{\psi}(x) = \{\mu_{\psi}(x), \sigma_{\psi}^2(x)\} \quad \psi \text{ are parameters of the NN}$$

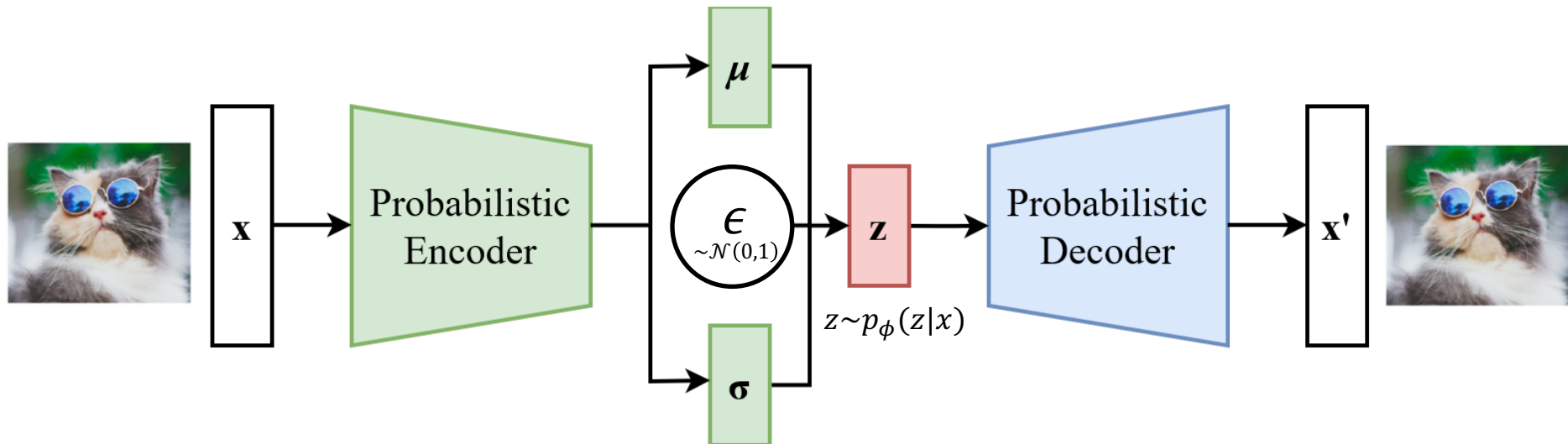
- What is the probability of a point in latent space?

$$p_{\psi}(z|x) = N(z | \mu_{\psi}(x), \sigma_{\psi}^2(x)) \quad \begin{array}{l} \text{Could choose different density} \\ \text{Gaussian is easiest} \end{array}$$

- How do we draw a sample in latent space?

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





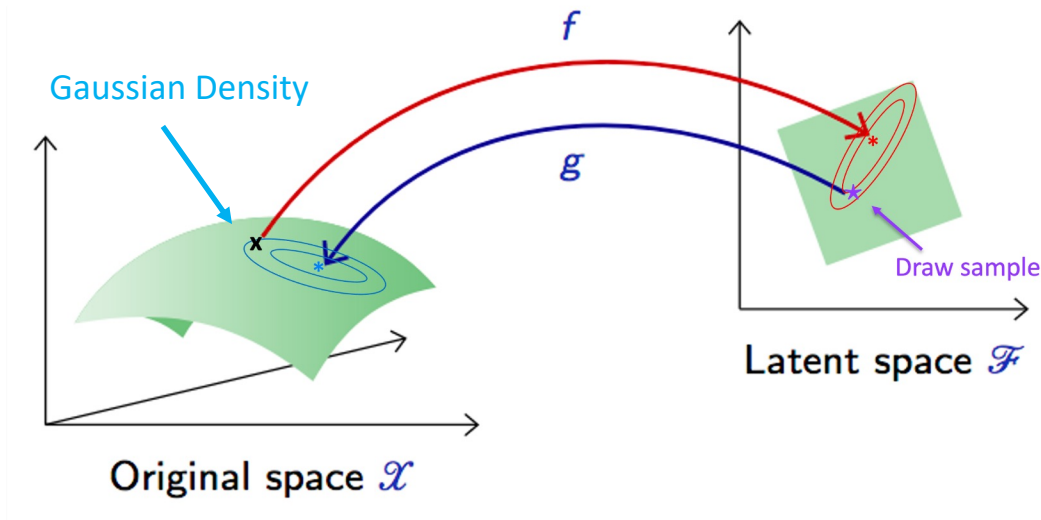
- Same as autoencoder

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

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

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

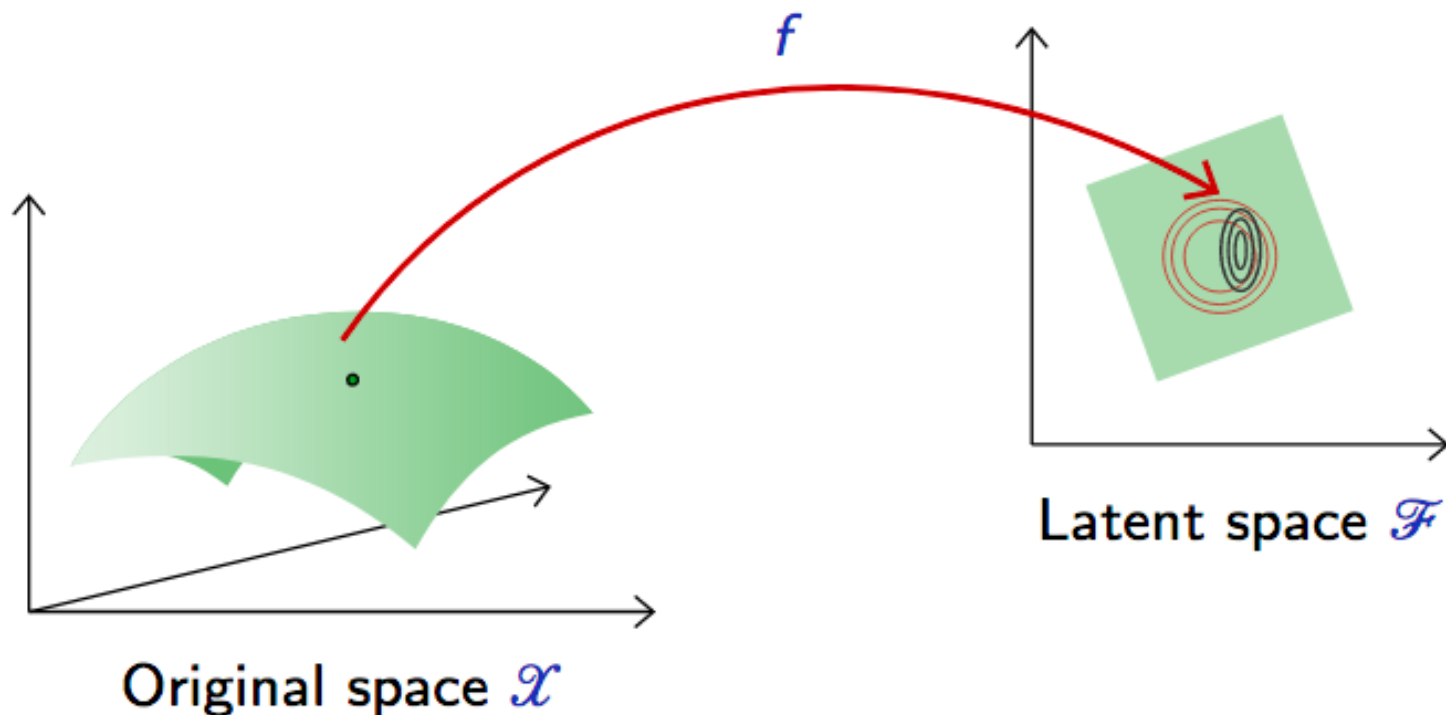
- **“Reconstruction Loss”**: Maximum likelihood

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

Same as the autoencoder loss

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

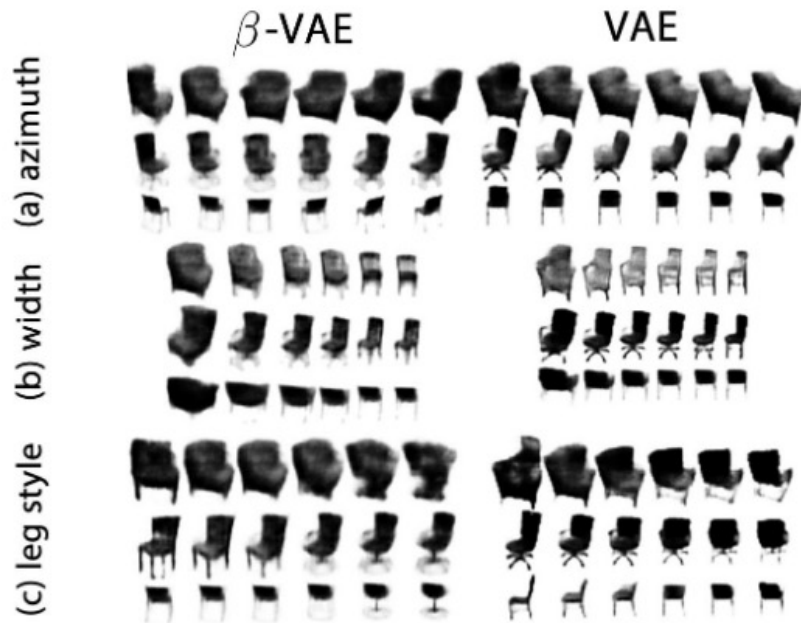
- How do we make sure system doesn't collapse to an autoencoder (i.e. VAE encoder only predicts mean)?
- Use prior $p(z)$ for the latent space distribution, **need to ensure the encoder is consistent with prior**



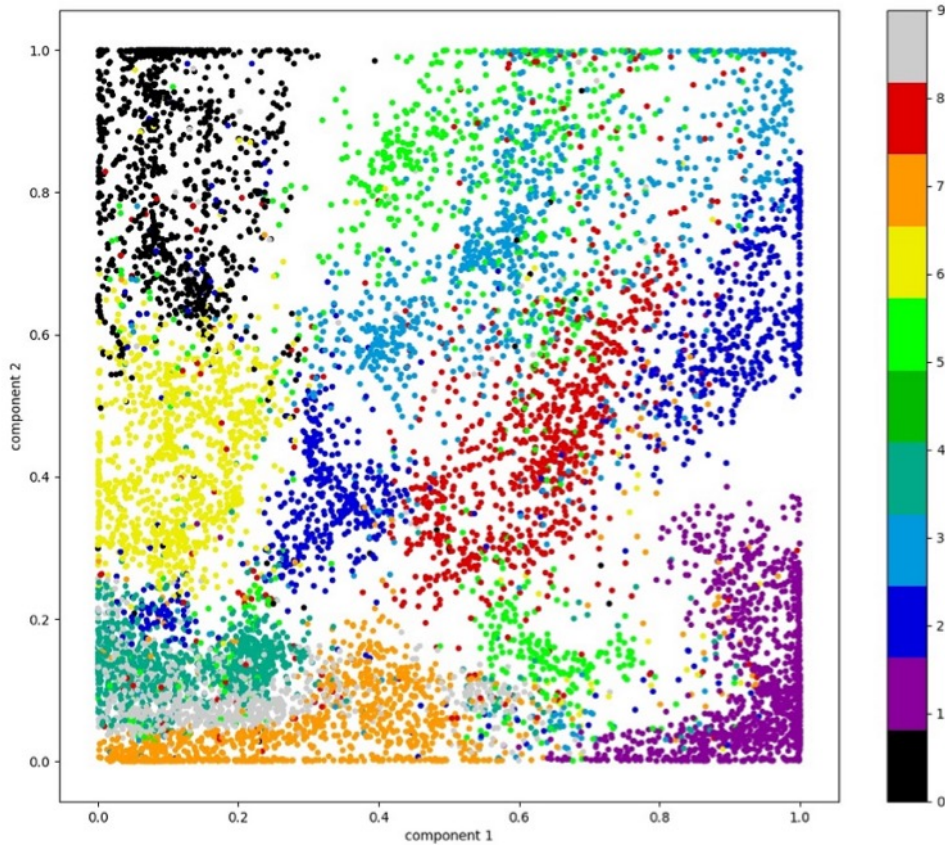
- Constrain difference between distributions with **Kullback–Leibler divergence**

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

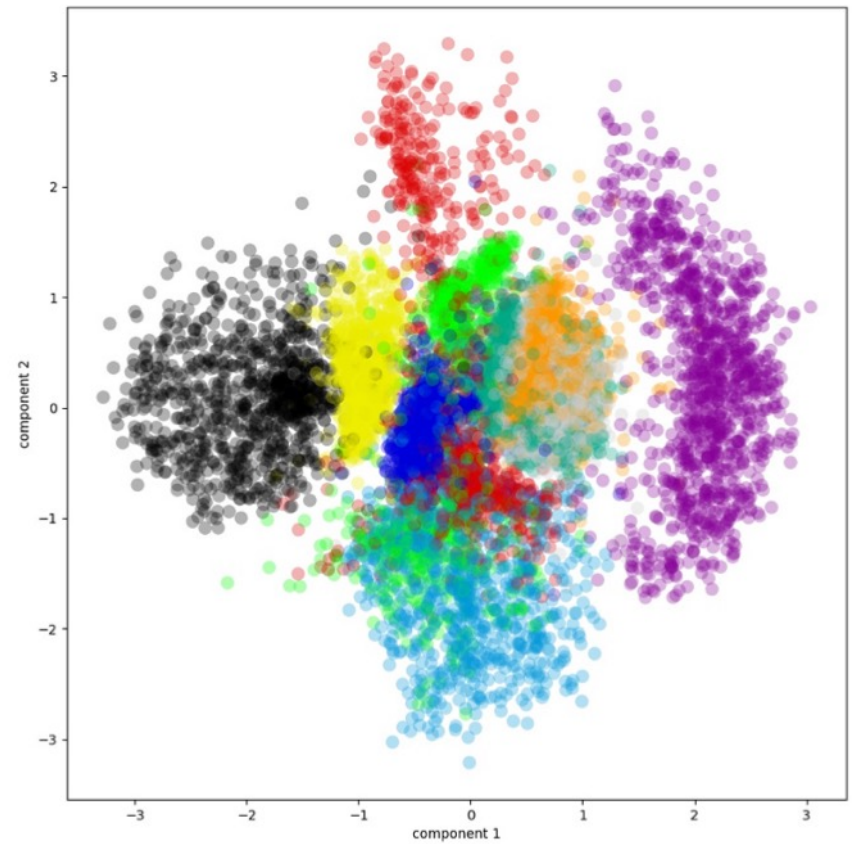
– $D_{KL}[q|p] \geq 0$ and is only 0 when $q = p$



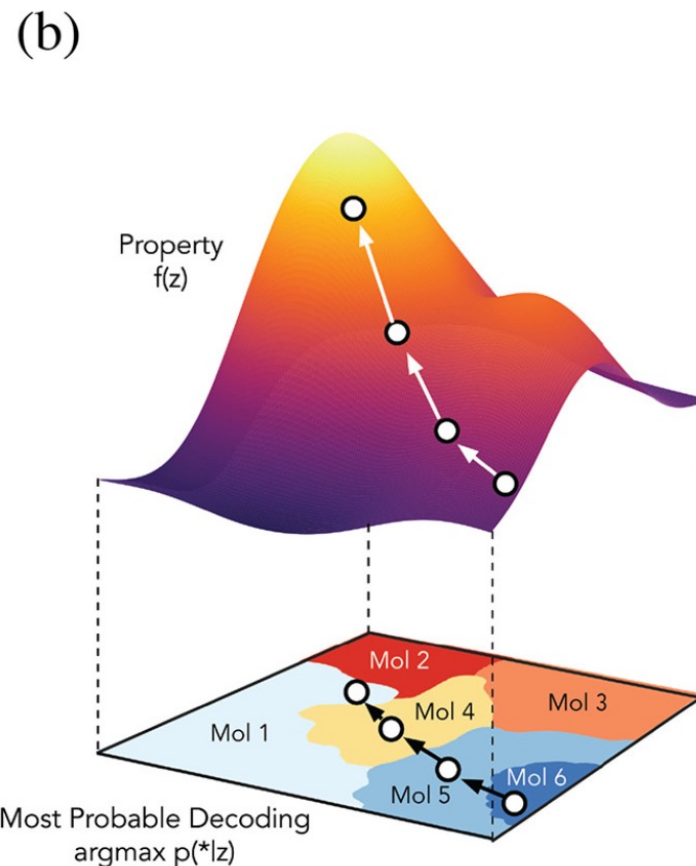
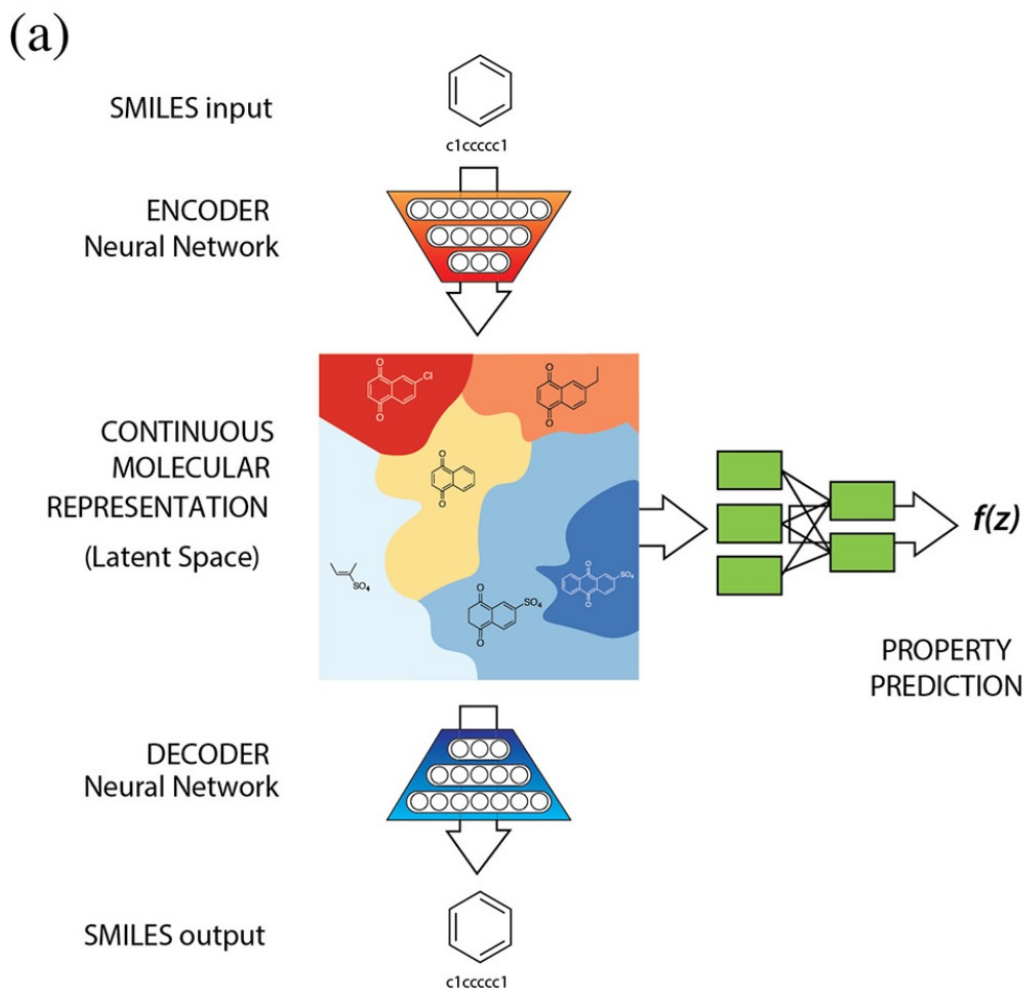
Autoencoder



Variational Autoencoder



Data: MNIST data set of hand-written digits

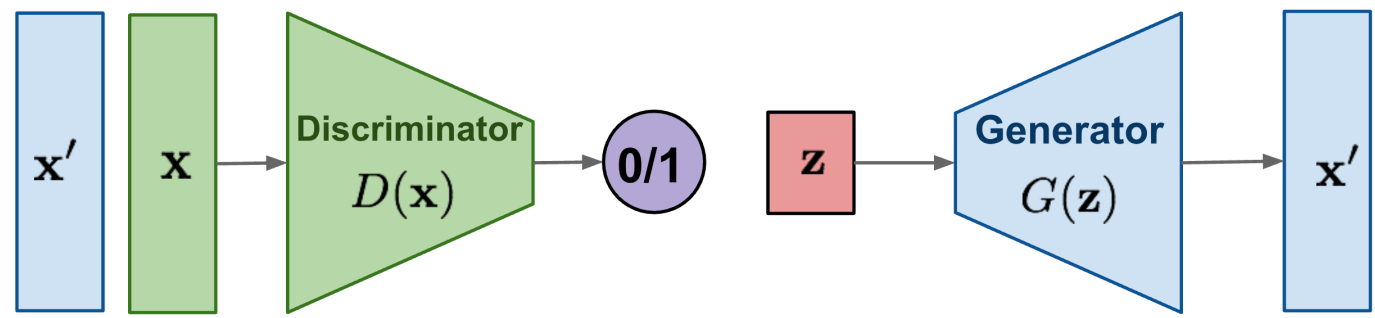


Design of new molecules with desired chemical properties.
(Gomez-Bombarelli et al, 2016)

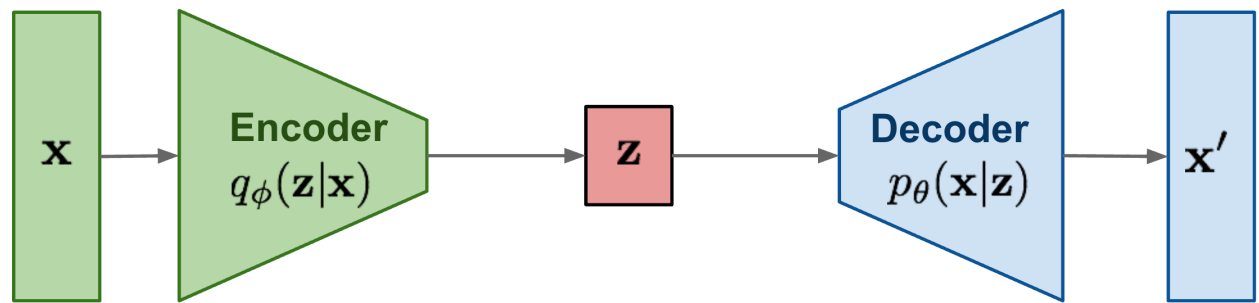
- In generative modeling, want to learn the lower dimensional degrees of freedom that describe the features of the data
- “Degrees of freedom” are modeled with a latent distribution (kept simple for convenience) and complex neural network mappings
- Need to think about **probabilistic systems**
- Design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

The Zoo of Generative Models...

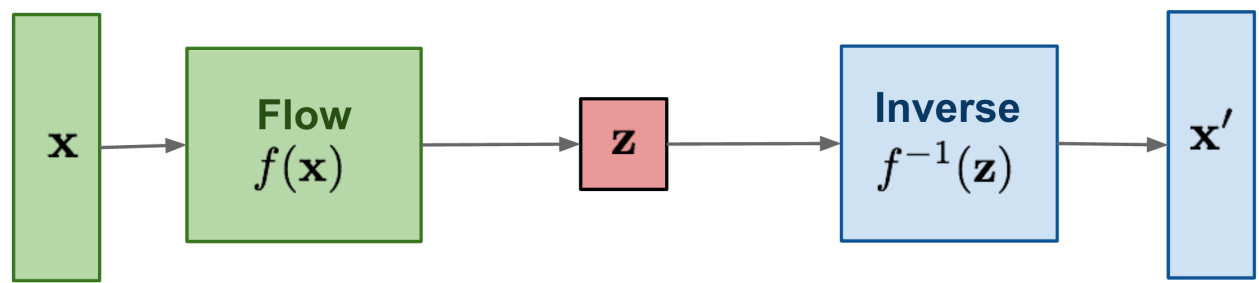
GAN: Adversarial training



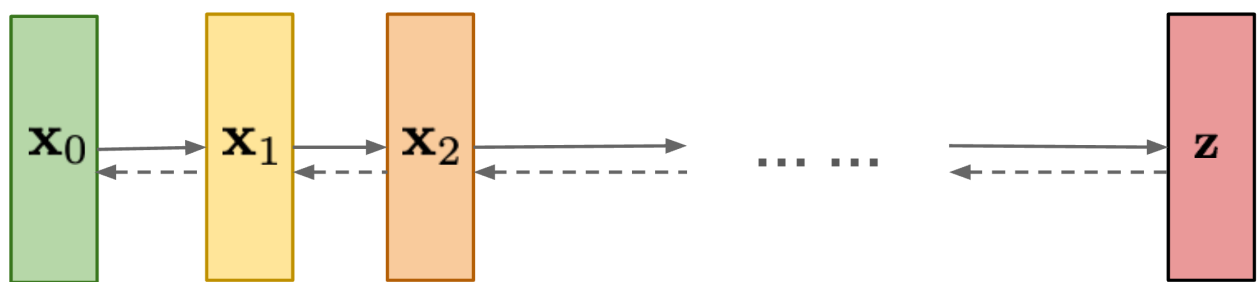
VAE: maximize variational lower bound



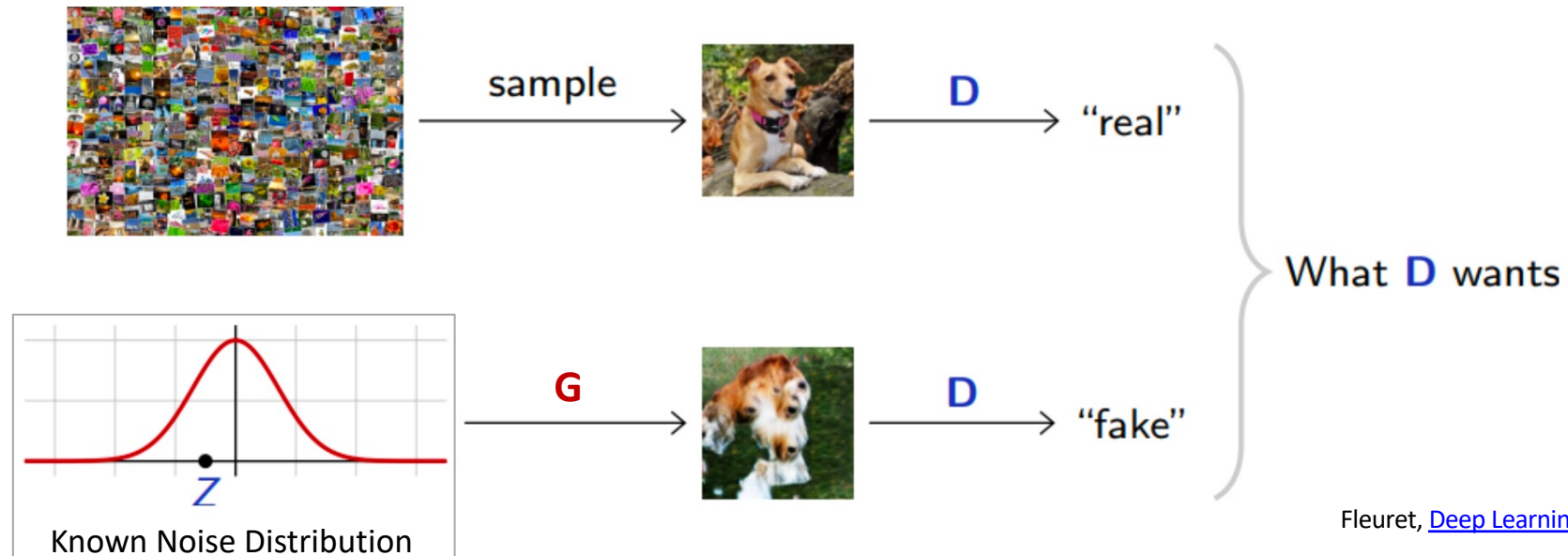
Flow-based models:
Invertible transform of distributions



Diffusion models:
Gradually add Gaussian noise and then reverse



Generative Adversarial Networks (GAN)



Fleuret, [Deep Learning Course](#)

- **Generator** creates data from noise, trained to trick **Discriminator** that classifies data as real or fake

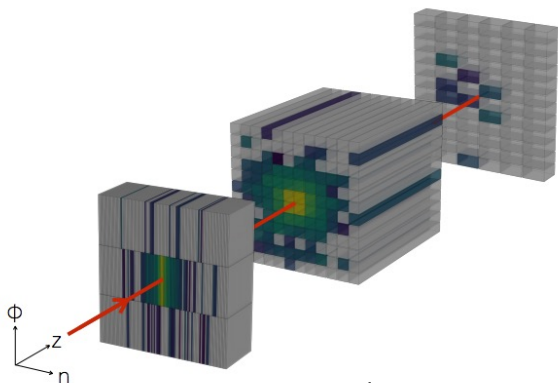
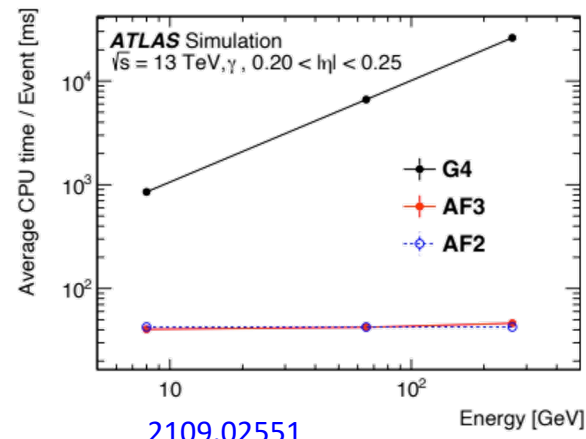
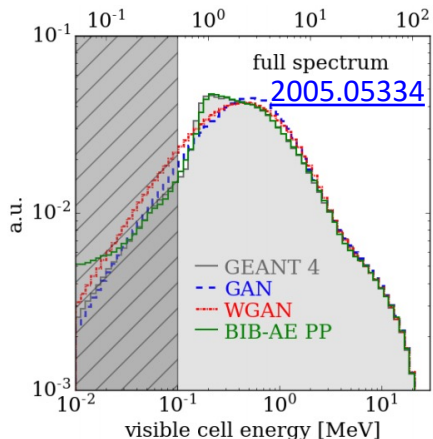
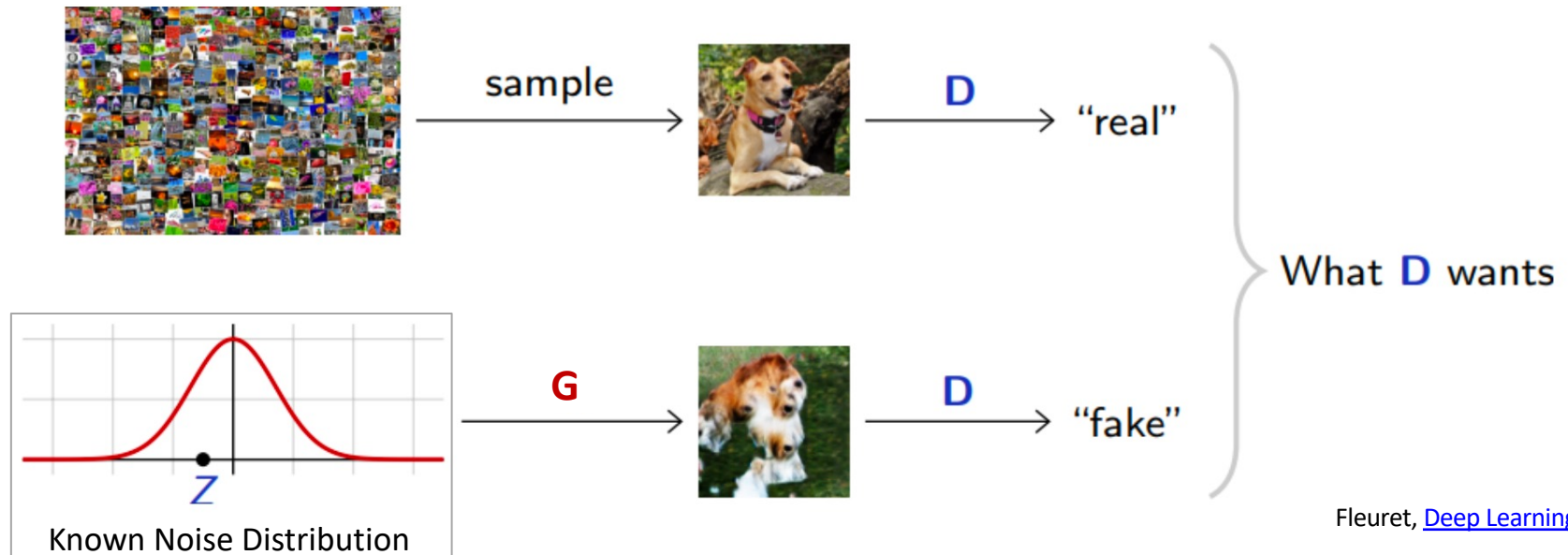


Image credit: [1705.02355](#)



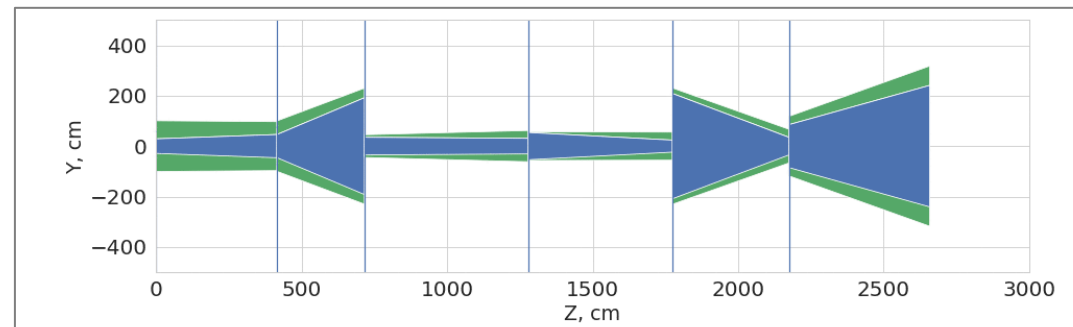
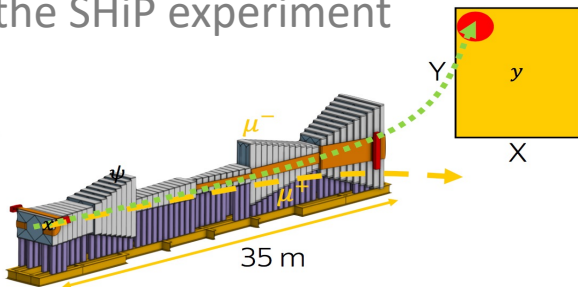
Generative Adversarial Networks (GAN)



Fleuret, [Deep Learning Course](#)

- **Generator** creates data from noise, trained to trick **Discriminator** that classifies data as real or fake

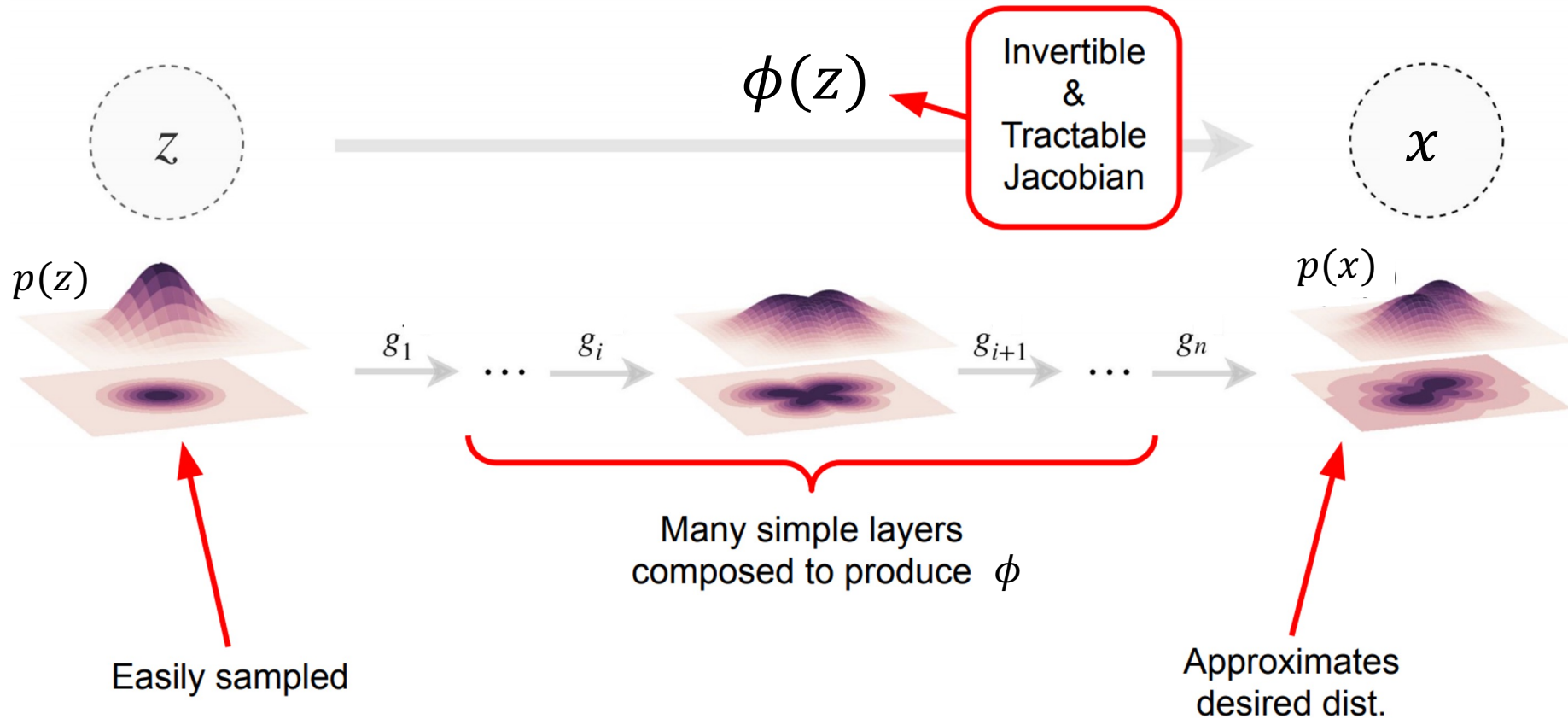
Optimization of the magnet system
For the SHiP experiment



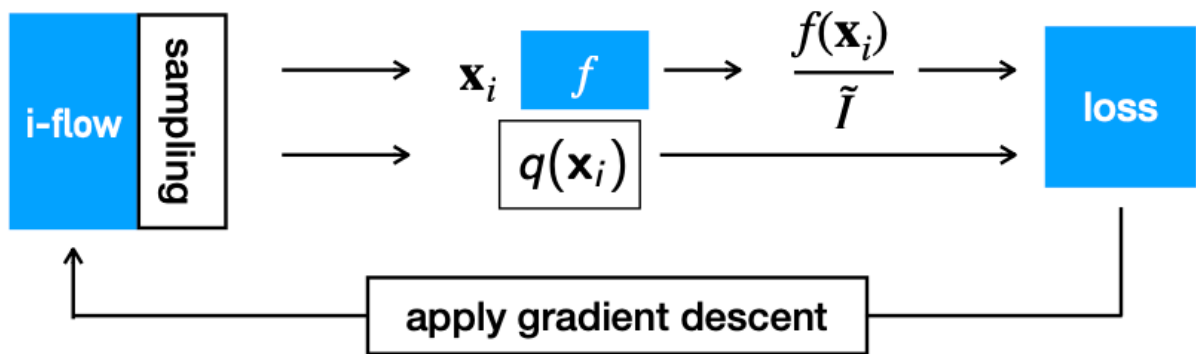
Explicit density estimation

We can evaluate density $p(x)$

$$p_x(\mathbf{x}) = p_z(\mathbf{z}) \left| \det \left(\frac{\partial \phi(\mathbf{z})}{d\mathbf{z}} \right)^{-1} \right|$$

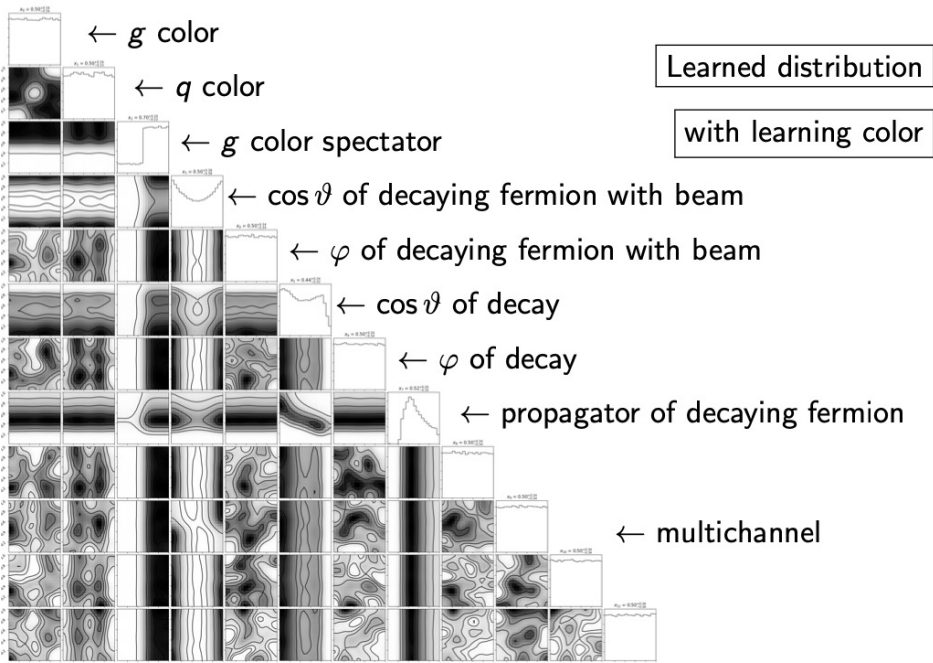
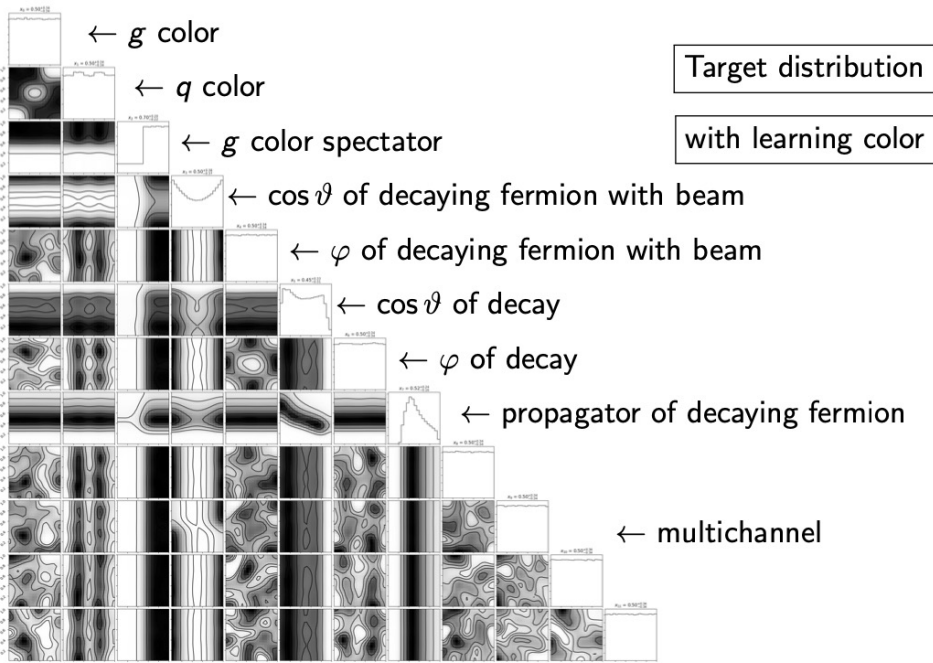


Event Generation with Normalizing Flows



arXiv: 2001.05486, ML:ST
 arXiv: 2001.10028, PRD
 Slide credit: [C. Krause](#)

Example: Learning $e^+ e^- \rightarrow 3j$



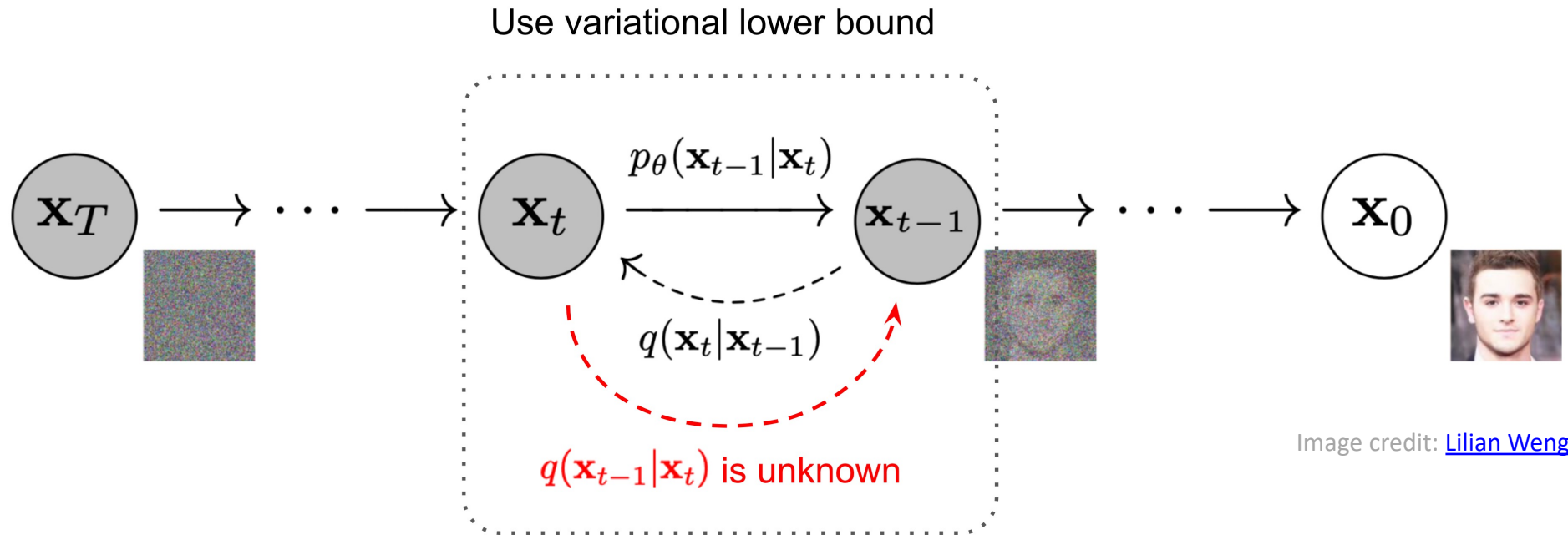
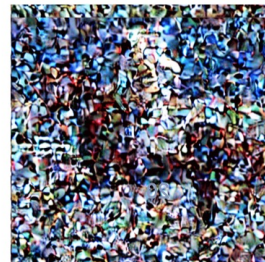


Image credit: [Lilian Weng](#)

- Iteratively add noise to data,
Train model to learn how to denoise step by step



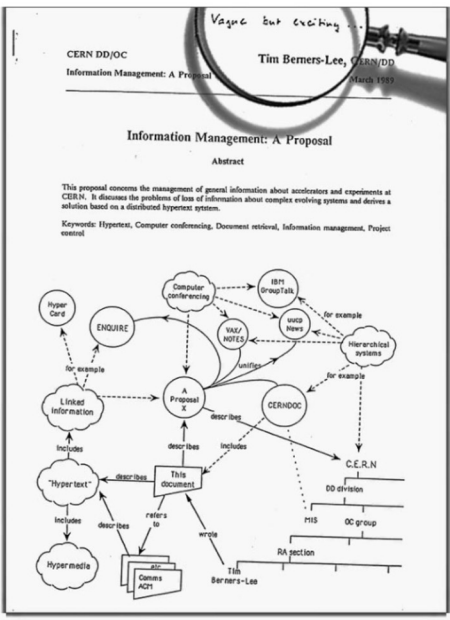
noise



data

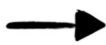
Some Final Thoughts

Since Tim Berners-Lee Invented the World Wide Web...



1989

~10 years



~25 years



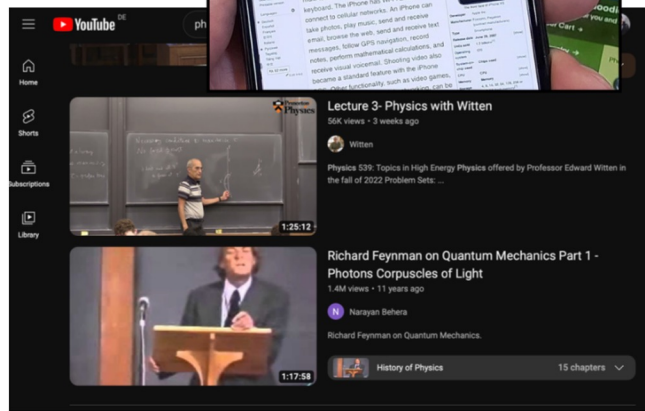
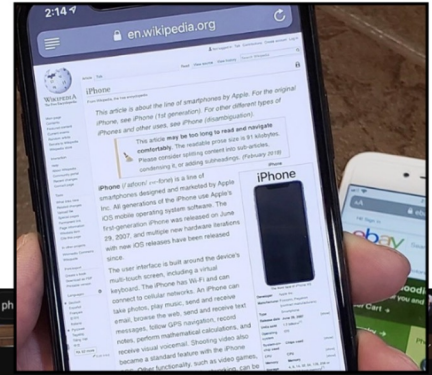
Search the web using Google

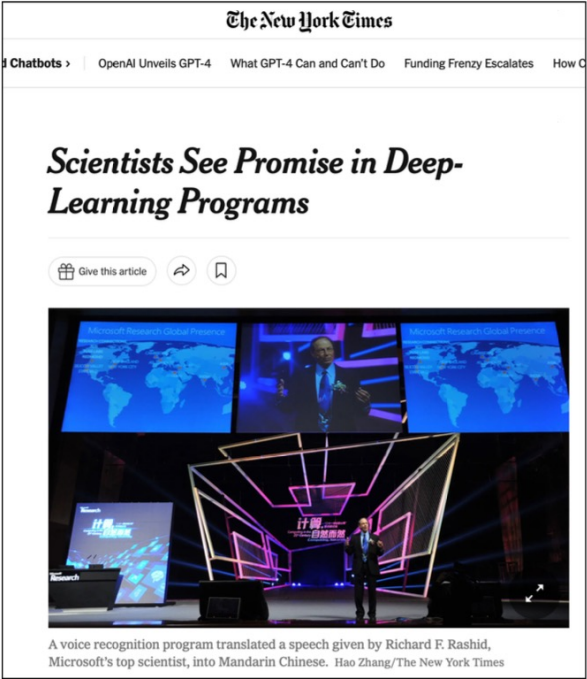


49

[More Google!](#)

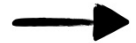
Copyright ©1999 Google Inc.





2012

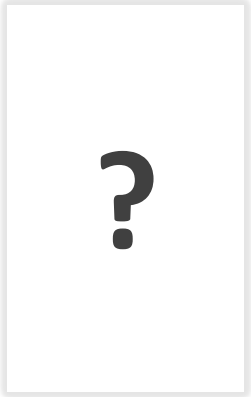
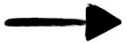
~10 years



Prompt: Several giant woolly mammoths approach treading through a snowy meadow [...]

[OpenAI Sora](#)

~25 years



Do These Models Know Physics?... Maybe Not Yet

52



- 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!

Backup

Explicit Density Estimation with Normalizing Flows

$$\int f(g(x)) \frac{\partial g(x)}{dx} dx = \int f(u) du \quad \text{where } u = g(x)$$

Multivariate:

$$\int f(g(\mathbf{x})) \left| \det \frac{\partial g(\mathbf{x})}{d\mathbf{x}} \right| d\mathbf{x} = \int f(\mathbf{u}) d\mathbf{u} \quad \text{where } \mathbf{u} = g(\mathbf{x})$$

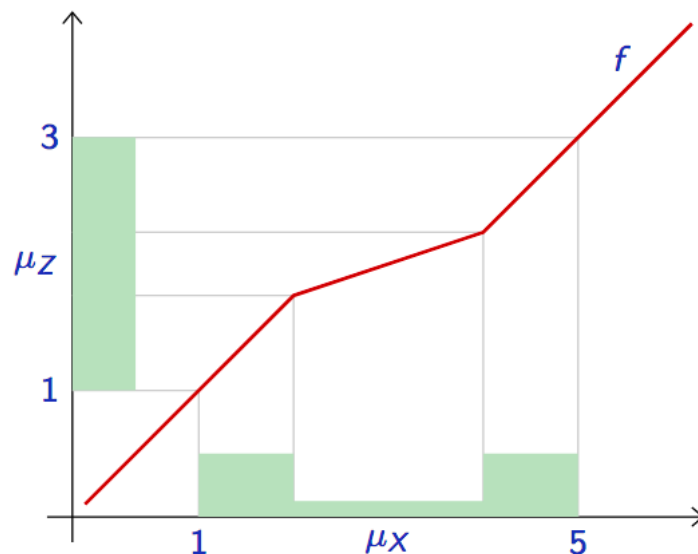
Determinant of Jacobian
of the transformation

→ Change of volume

Change of Variables in Probability

- If f is continuous, invertible, differentiable, and $\mathbf{x} = f^{-1}(\mathbf{z}) \equiv \phi(\mathbf{z})$ then

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{z}) \left| \det \left(\frac{\partial \phi(\mathbf{z})}{d\mathbf{z}} \right)^{-1} \right| \quad \text{where } \mathbf{x} = \phi(\mathbf{z})$$



The term $\left| \det \left(\frac{\partial \phi(\mathbf{z})}{d\mathbf{z}} \right)^{-1} \right|$ accounts for the local stretching of space

- If f is continuous, invertible, differentiable, and $\mathbf{x} = f^{-1}(\mathbf{z}) \equiv \phi(\mathbf{z})$ then

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{z}) \left| \det \left(\frac{\partial \phi(\mathbf{z})}{\partial \mathbf{z}} \right)^{-1} \right| \quad \text{where } \mathbf{x} = \phi(\mathbf{z})$$

- \mathbf{x} = data we want to model, \mathbf{z} = known noise
- $\phi_{\theta}(\mathbf{z})$ will be a neural network with parameters θ
 - Must be continuous, invertible, differentiable
- Output of ϕ is a potential sample \mathbf{x}
 - **Learn the right ϕ** : adjust weights θ to maximize data probability (formula above)

- If f is continuous, invertible, differentiable, and $\mathbf{x} = f^{-1}(\mathbf{z}) \equiv \phi(\mathbf{z})$ then

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{z}) \left| \det \left(\frac{\partial \phi(\mathbf{z})}{\partial \mathbf{z}} \right)^{-1} \right| \quad \text{where } \mathbf{x} = \phi(\mathbf{z})$$

- \mathbf{x} = data we want to model, \mathbf{z} = known noise

$\phi(\mathbf{z})$ neural network

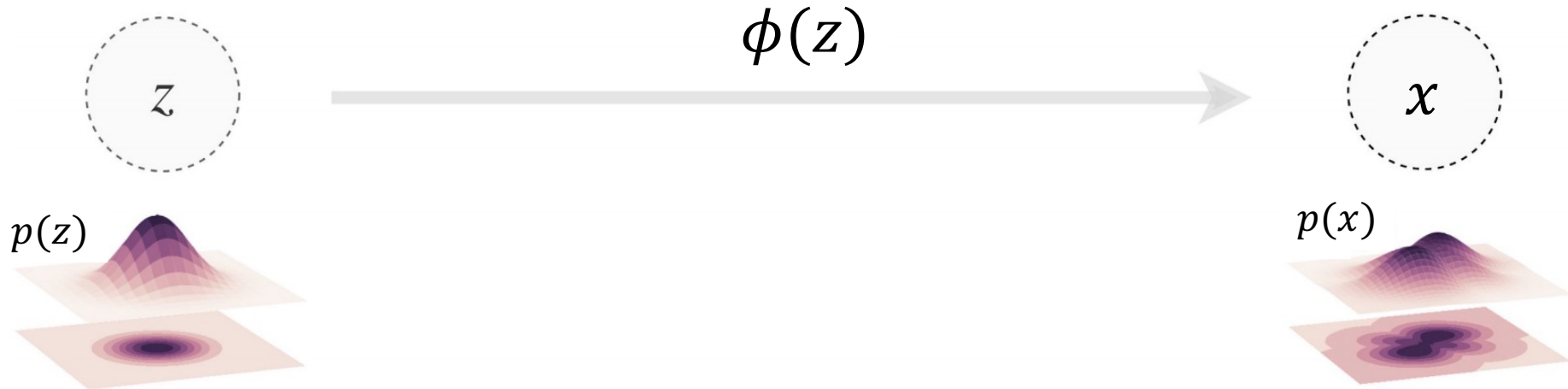
- Input = a sample of noise
- Output = a sample of \mathbf{X}

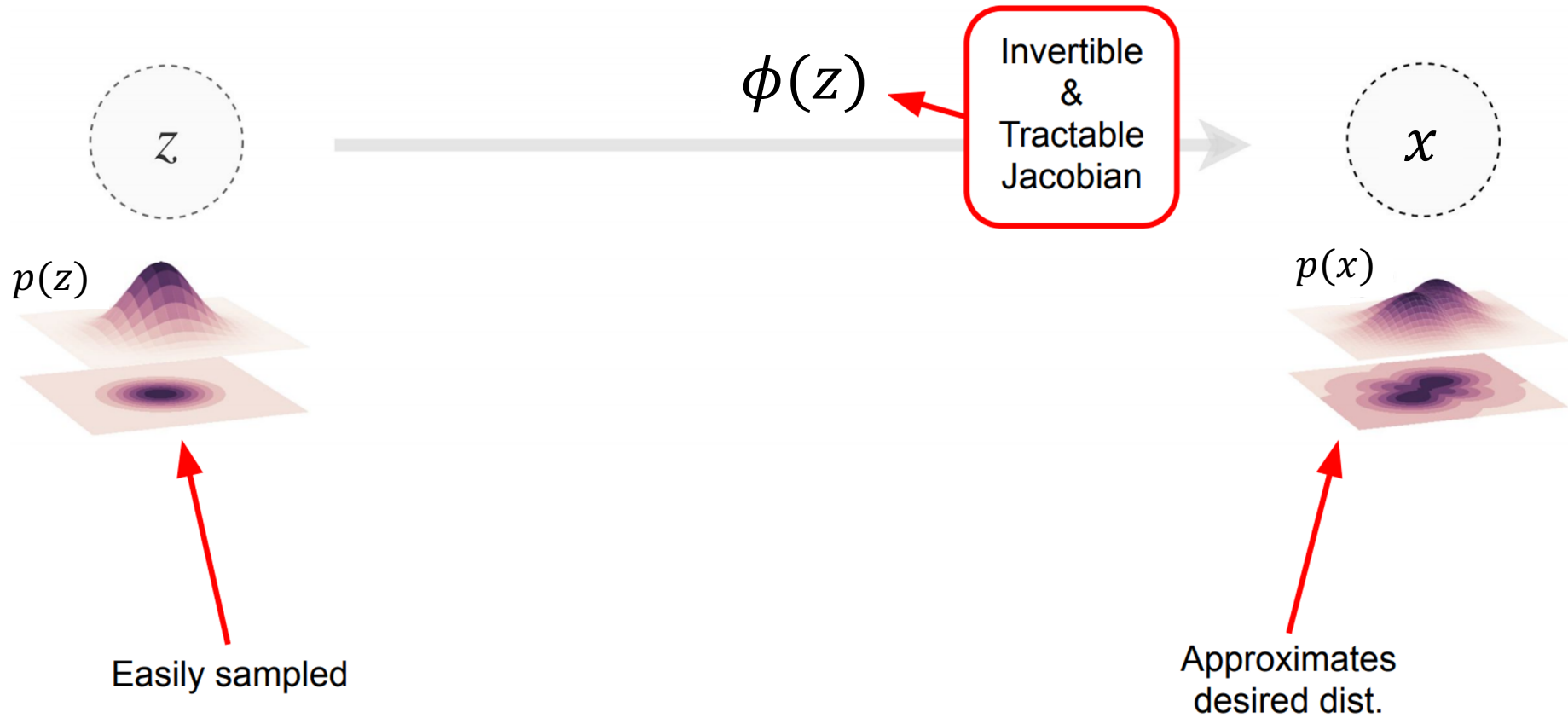
$\phi^{-1}(\mathbf{x})$ inverse

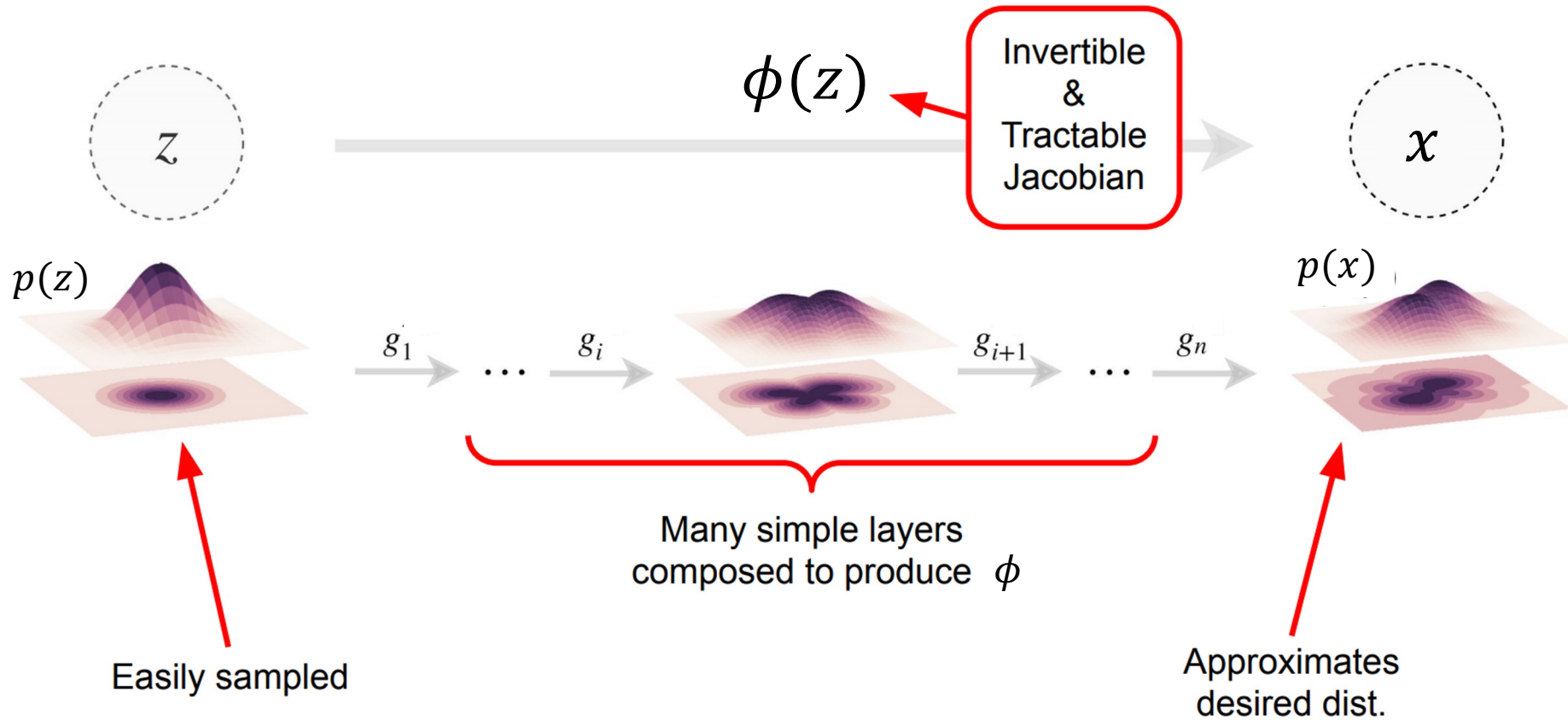
- Input = a sample \mathbf{X}
- Output = a sample of noise

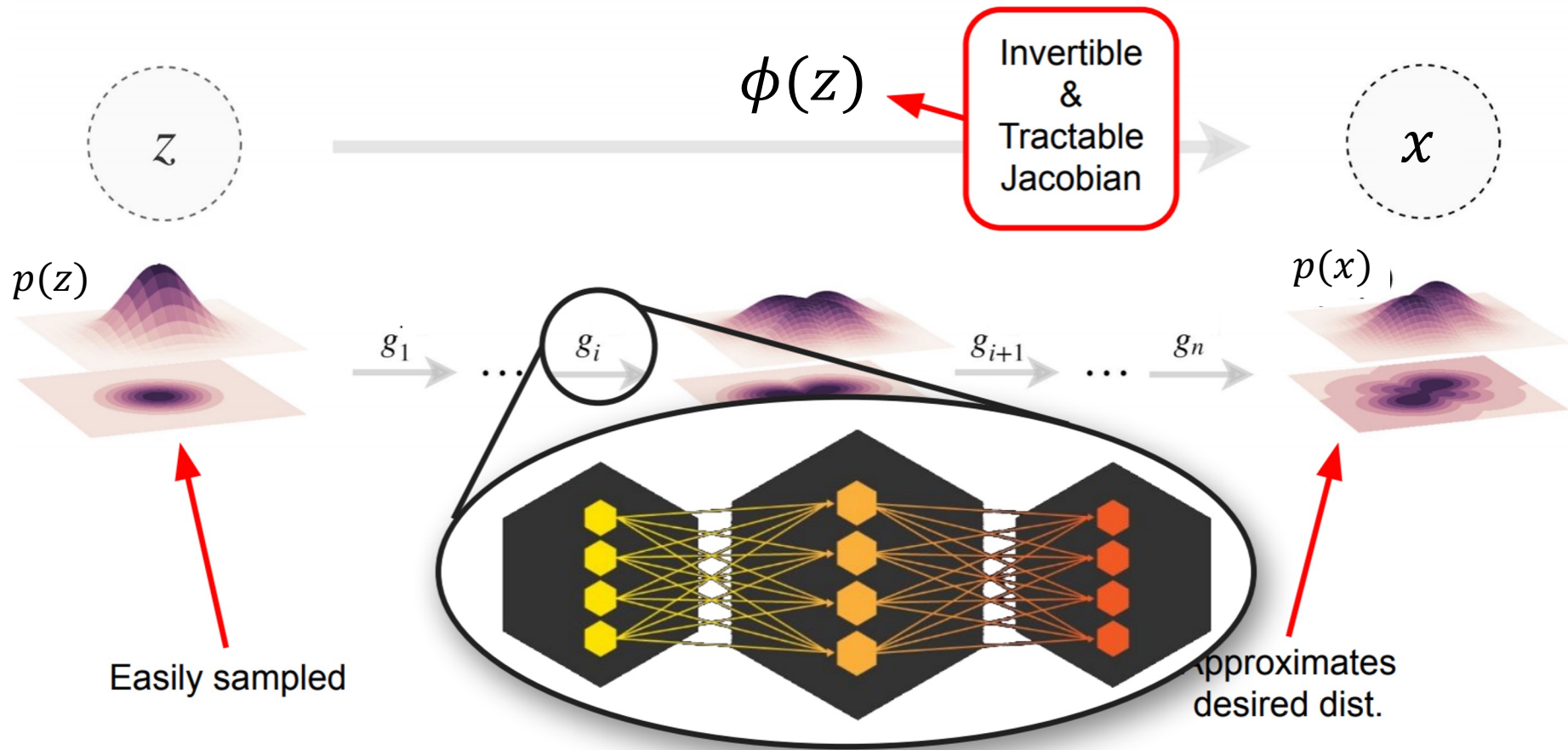
- Calculate the probability of a sample using the formula above

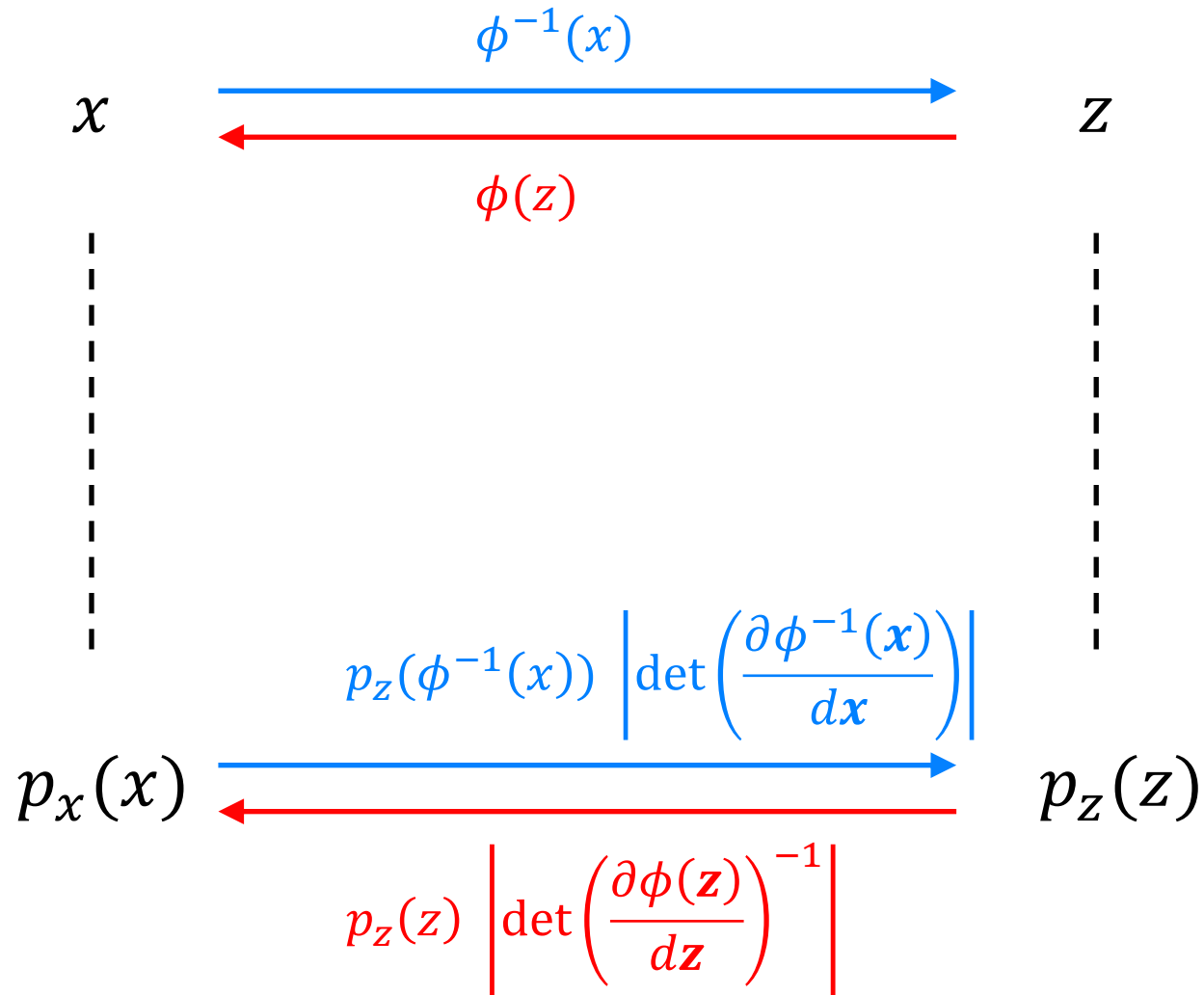
$$p_x(\mathbf{x}) = p_z(\mathbf{z}) \left| \det \left(\frac{\partial \phi(\mathbf{z})}{d\mathbf{z}} \right)^{-1} \right|$$











- **Learn** θ with maximum likelihood

$$\max_{\theta} p(x) = \max_{\theta} p_z(\phi_{\theta}^{-1}(x)) \left| \det \left(\frac{\partial \phi_{\theta}^{-1}(x)}{dx} \right) \right|$$

- Gradient descent on θ
- Find transformation s.t. data is most likely

- **Benefits** once trained

- Can evaluate $p(x)$ for any point X
- Can generate “new” data points
 - Sample noise: $z \sim p(z)$
 - Transform: $\phi(z) = x$

Example Normalizing Flow: Real NVP

- Data vector $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$
- Transformation

Functions $f()$ and $g()$
are neural networks

$$\phi(z): \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \phi_1(z) \\ \phi_2(z) \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 * f(z_1) + g(z_1) \end{pmatrix}$$

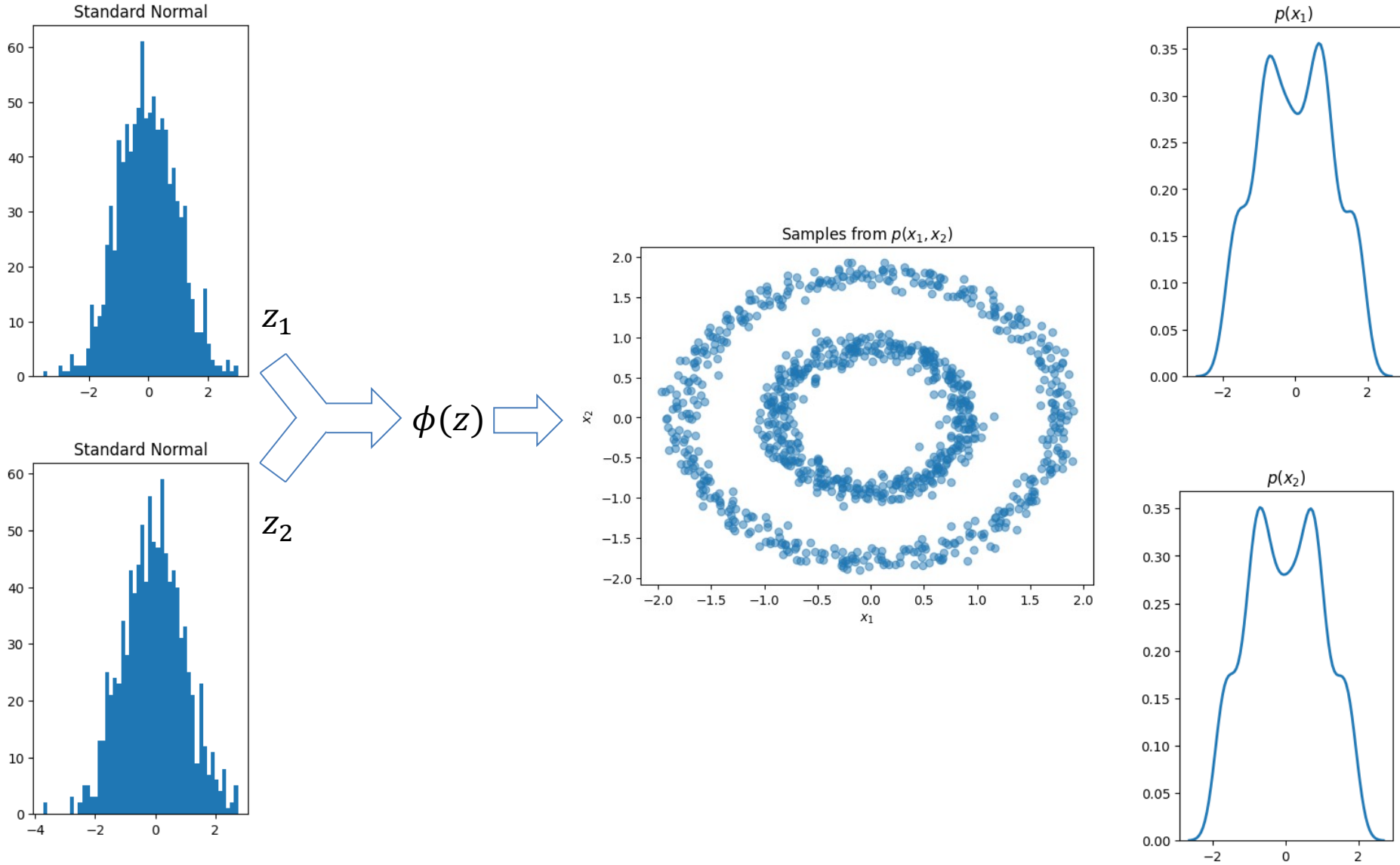
$$\phi^{-1}(x): \quad \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} \phi_1^{-1}(x) \\ \phi_2^{-1}(x) \end{pmatrix} = \begin{pmatrix} x_1 \\ (x_2 - g(x_1))/f(x_1) \end{pmatrix}$$

- Determinant:

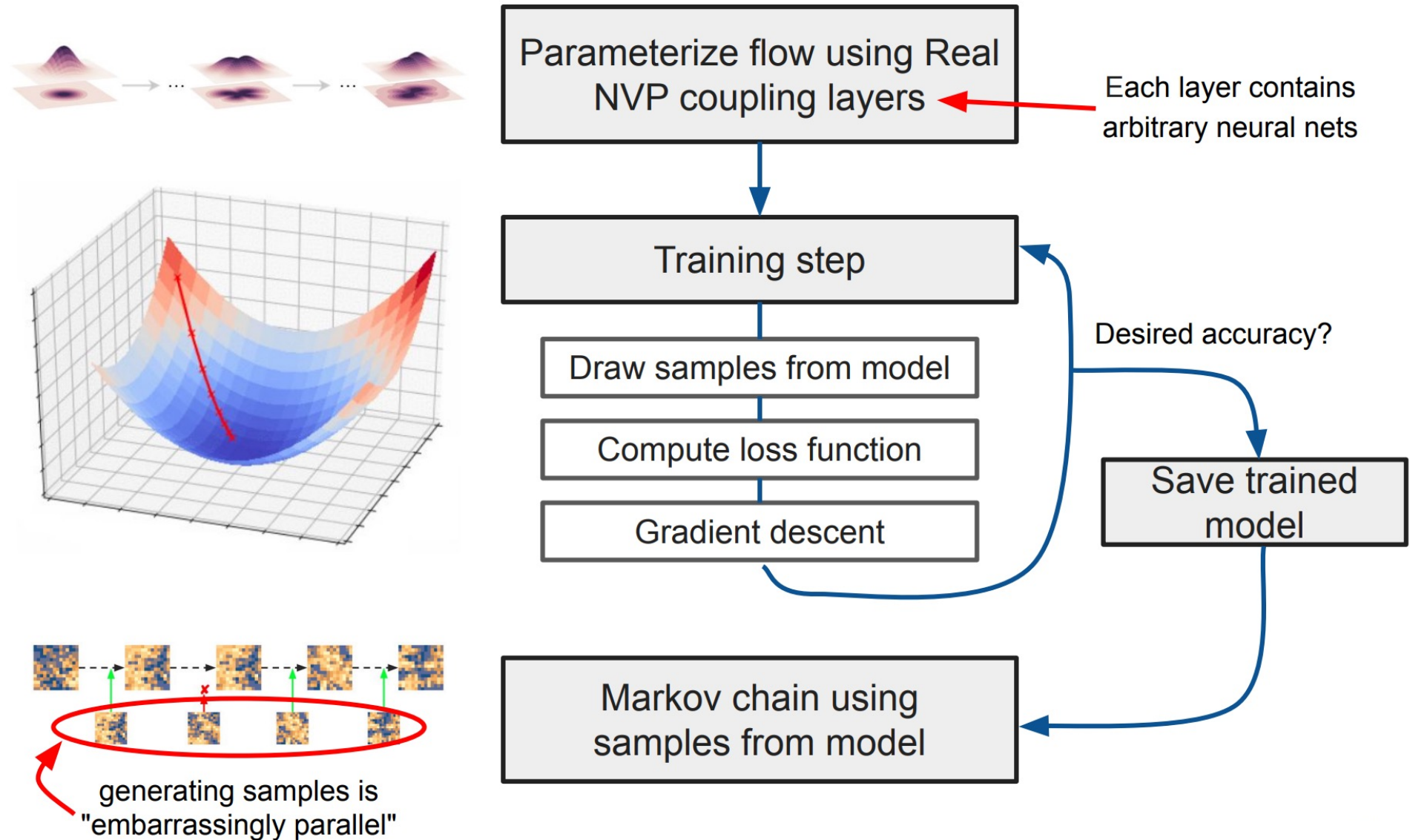
Jacobian is
lower triangular

$$\det\left(\frac{\partial\phi(z)}{\partial z}\right) = \det\left(\begin{pmatrix} 1 & 0 \\ \left(\frac{\partial\phi_2(z)}{\partial z_1}\right) & f(z_1) \end{pmatrix}\right) = f(z_2)$$

Example Normalizing flow



Applications: Sampling in Lattice QCD



GANS

- Formulate as a two player game
- One player tries to output data that looks as real as possible
- Another player tries to compare real and fake data
- In this case we need:
 1. A *generator* that can produce samples
 2. A measure of *not too far from the real data*

- **Generator network $g_{\theta}(z)$** with parameters θ
 - Map sample from known $p(z)$ to sample in data space

$$x = g_{\theta}(z) \quad z \sim p(z)$$

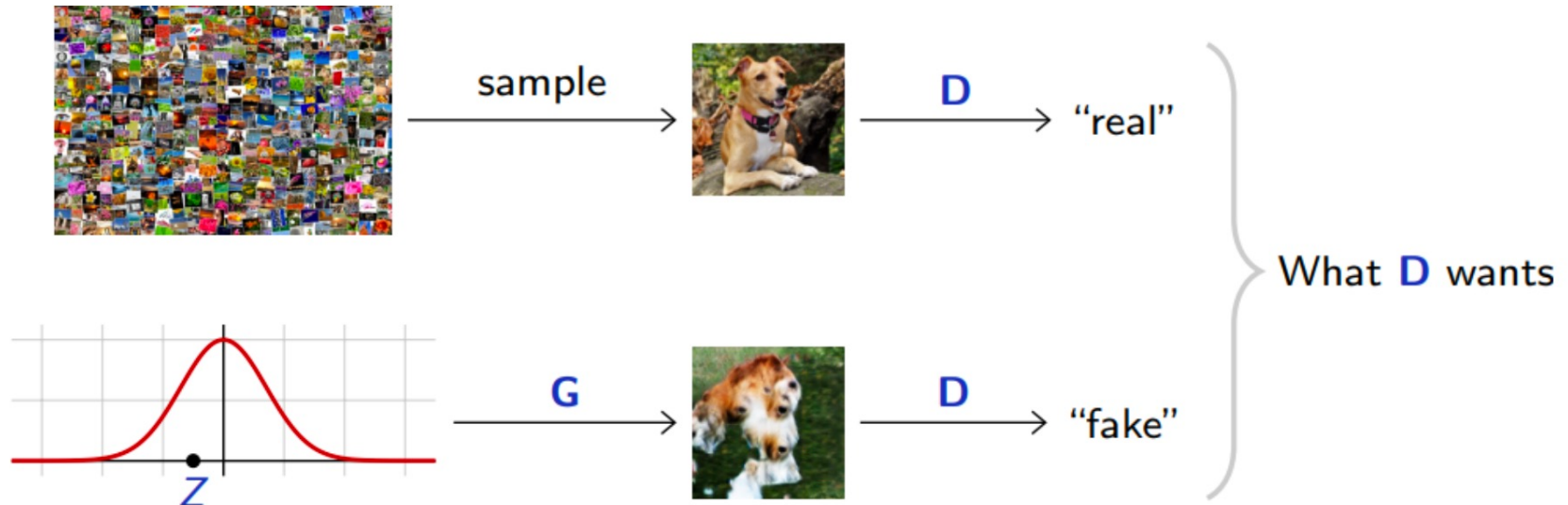
- We don't know what the generated distribution $p_{\theta}(x)$ is, but we can sample from it \rightarrow *Implicit Model*

- **Generator network $g_{\theta}(z)$** with parameters θ
 - Map sample from known $p(z)$ to sample in data space

$$x = g_{\theta}(z) \quad z \sim p(z)$$

- We don't know what the generated distribution $p_{\theta}(x)$ is, but we can sample from it \rightarrow *Implicit Model*

- **Discriminator Network $d_{\phi}(x)$** with parameters ϕ
 - Classifier trained to distinguish between real and fake data
 - Classifier is learning to predict $p(y = \text{real} \mid x)$
 - This classifier is our measure of *not too far from the real data*



- Generator's goal is to produce *fake* data that tricks the discriminator to think it is *real* data
- Discriminator wants to miss-classify data as real or fake as little as possible
- The setup is *adversarial* because the two networks have opposing objectives

- Data
 - Real data samples: $\{x_i, y_i = 1\}$
 - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$ with: $z_i \sim p(z)$

- Data
 - Real data samples: $\{x_i, y_i = 1\}$
 - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$ with: $z_i \sim p(z)$
- For a fixed generator, can train discriminator by minimizing the cross entropy

$$L(\phi) = -\frac{1}{2N} \sum_{i=1}^N \left[y_i \log d_\phi(x_i) + (1 - \tilde{y}_i) \log(1 - d_\phi(\tilde{x}_i)) \right]$$

- Data
 - Real data samples: $\{x_i, y_i = 1\}$
 - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$ with: $z_i \sim p(z)$
- For a fixed generator, can train discriminator by minimizing the cross entropy

$$\begin{aligned} L(\phi) &= -\frac{1}{2N} \sum_{i=1}^N \left[y_i \log d_\phi(x_i) + (1 - \tilde{y}_i) \log(1 - d_\phi(\tilde{x}_i)) \right] \\ &= -\frac{1}{2N} \sum_{i=1}^N \left[\log d_\phi(x_i) + \log(1 - d_\phi(g_\theta(z_i))) \right] \end{aligned}$$

- Data
 - Real data samples: $\{x_i, y_i = 1\}$
 - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$ with: $z_i \sim p(z)$
- For a fixed generator, can train discriminator by minimizing the cross entropy

$$\begin{aligned} L(\phi) &= -\frac{1}{2N} \sum_{i=1}^N \left[y_i \log d_\phi(x_i) + (1 - \tilde{y}_i) \log(1 - d_\phi(\tilde{x}_i)) \right] \\ &= -\frac{1}{2N} \sum_{i=1}^N \left[\log d_\phi(x_i) + \log(1 - d_\phi(g_\theta(z_i))) \right] \\ &= -\mathbb{E}_{x \sim p_{\text{data}}(x)} \left[\log d_\phi(x) \right] - \mathbb{E}_{z \sim p(z)} \left[\log(1 - d_\phi(g_\theta(z))) \right] \end{aligned}$$

- However, generator isn't fixed... have to train it!

- However, generator isn't fixed... have to train it!
- Consider objective as a *value function* of ϕ and θ

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[\log d_{\phi}(x) \right] + \mathbb{E}_{z \sim p(z)} \left[\log(1 - d_{\phi}(g_{\theta}(z))) \right]$$

- However, generator isn't fixed... have to train it!
- Consider objective as a *value function* of ϕ and θ

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[\log d_{\phi}(x) \right] + \mathbb{E}_{z \sim p(z)} \left[\log(1 - d_{\phi}(g_{\theta}(z))) \right]$$

- For fixed generator, $V(\phi, \theta)$ is high when discriminator is good, i.e. when generator is not producing good fakes
- For a perfect discriminator, a good generator will confuse discriminator and $V(\phi, \theta)$ will be low

- However, generator isn't fixed... have to train it!
- Consider objective as a *value function* of ϕ and θ

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[\log d_{\phi}(x) \right] + \mathbb{E}_{z \sim p(z)} \left[\log(1 - d_{\phi}(g_{\theta}(z))) \right]$$

- For fixed generator, $V(\phi, \theta)$ is high when discriminator is good, i.e. when generator is not producing good fakes
 - For a perfect discriminator, a good generator will confuse discriminator and $V(\phi, \theta)$ will be low
- So our optimization goal becomes:

$$\theta^* = \arg \min_{\theta} \max_{\phi} V(\phi, \theta)$$

- However, generator isn't fixed... have to train it!
- Consider objective as a *value function* of ϕ and θ

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[\log d_{\phi}(x) \right] + \mathbb{E}_{z \sim p(z)} \left[\log(1 - d_{\phi}(g_{\theta}(z))) \right]$$

- For fixed generator, $V(\phi, \theta)$ is high when discriminator is good, i.e. when generator is not producing good fakes
- For a perfect discriminator, a good generator will confuse discriminator and $V(\phi, \theta)$ will be low
- So our optimization goal becomes:

$$\theta^* = \arg \min_{\theta} \max_{\phi} V(\phi, \theta)$$

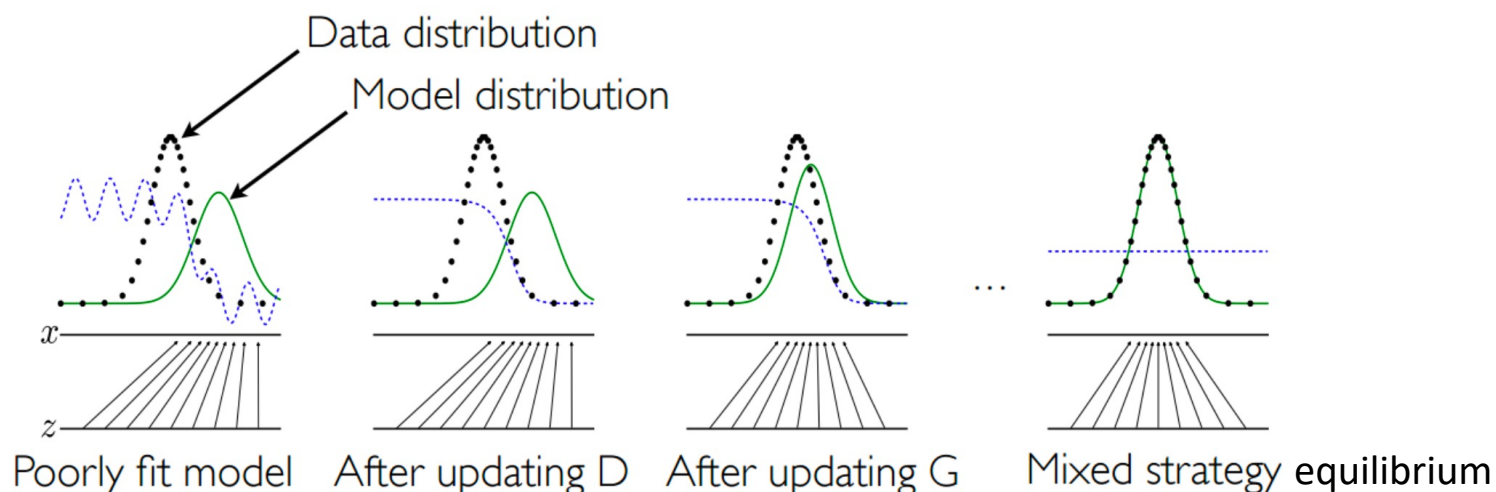
NOTE: can prove that
minimax solution
corresponds to generator
that perfectly reproduces
data distribution
 $q_{\theta^*}(x) = p_{\text{data}}(x)$

- Alternating Gradient descent to solve the min-max problem:

$$\theta \leftarrow \theta - \gamma \nabla_{\theta} V(\phi, \theta) = \theta - \gamma \frac{\partial V}{\partial d} \frac{\partial (d_{\phi})}{\partial g} \frac{\partial g_{\theta}}{\partial \theta}$$

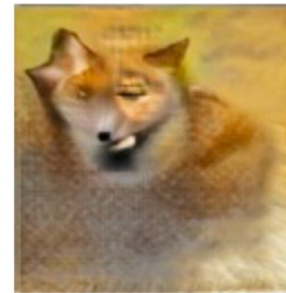
$$\phi \leftarrow \phi - \gamma \nabla_{\phi} V(\phi, \theta) = \phi - \gamma \frac{\partial V}{\partial d} \frac{d(d_{\phi})}{d\phi}$$

- For each θ step, take k steps in ϕ to keep discriminator near optimal



Examples

Goodfellow et. al., 2014



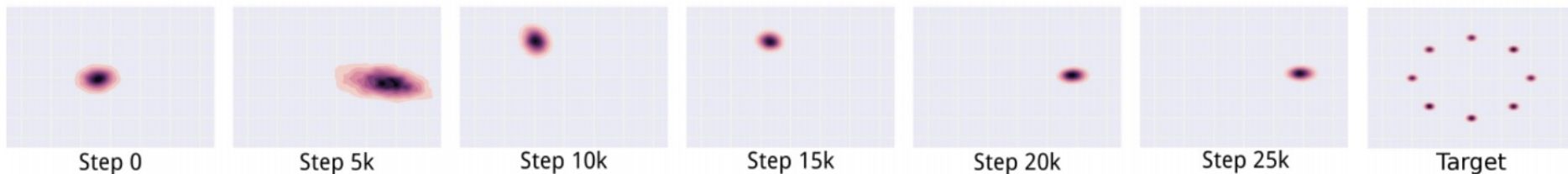
Not so good

Goodfellow 2016



Radford et al, 2015

- **Oscillations without convergence:** unlike standard loss minimization, alternating stochastic gradient descent has no guarantee of convergence.
- **Vanishing gradients:** if classifier is too good, value function saturates \rightarrow no gradient to update generator
- **Mode collapse:** generator models only a small sub-population, concentrating on a few data distribution modes.
- **Difficult to assess performance,** when are generated data good enough?



- Standard GANS compare real and fake distributions with Jensen-Shannon Divergence, “vertically”
- Wasserstein-GAN (Arjovsky et al, [2017](#)) compares “horizontally” with Wasserstein-1 distance (a.k.a. Earth Movers distance)
- Substantially improves *vanishing gradient* and *mode collapse* problems!

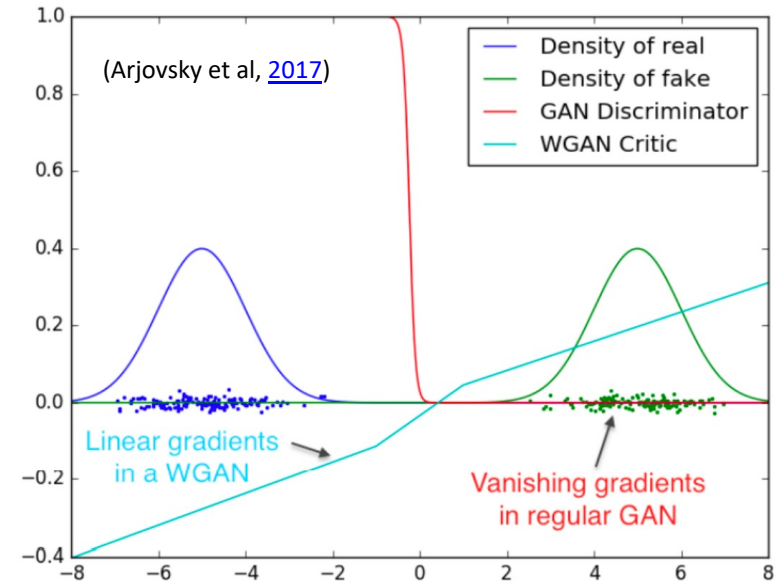
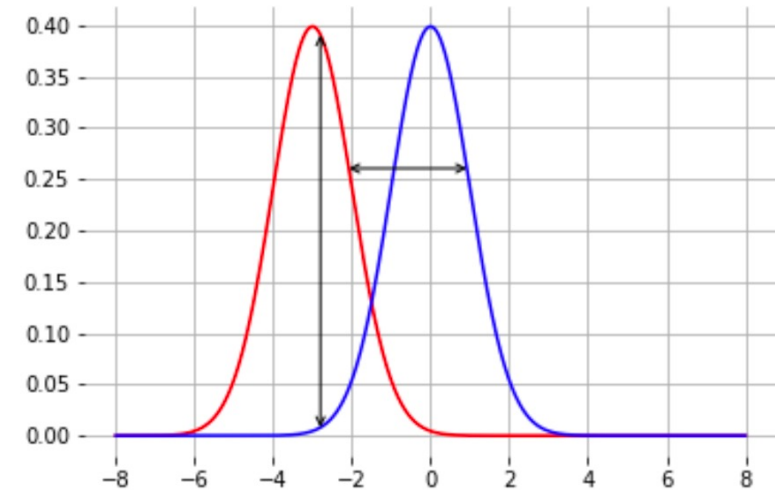
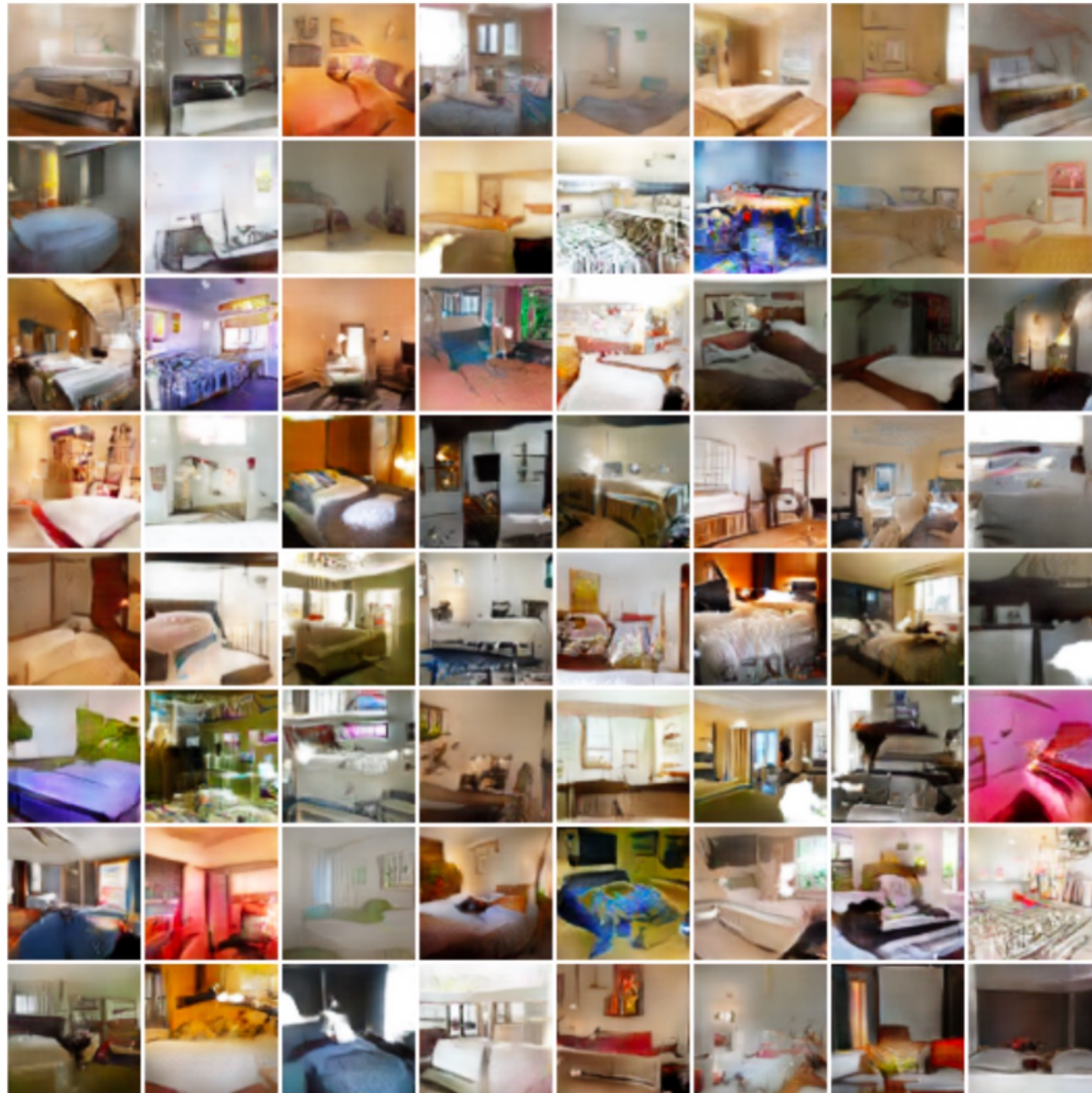
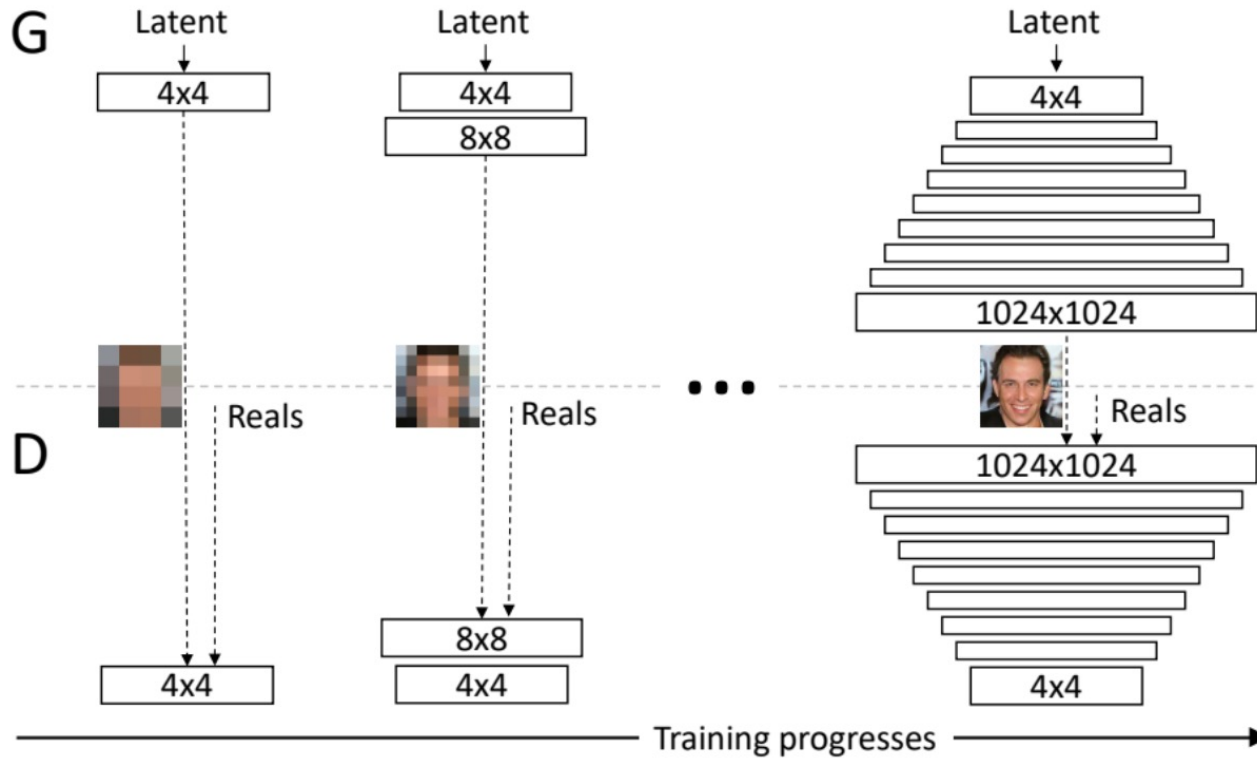


Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.

WGAN Examples



Progressive GAN



(Karras et al, 2017)

StyleGAN v2



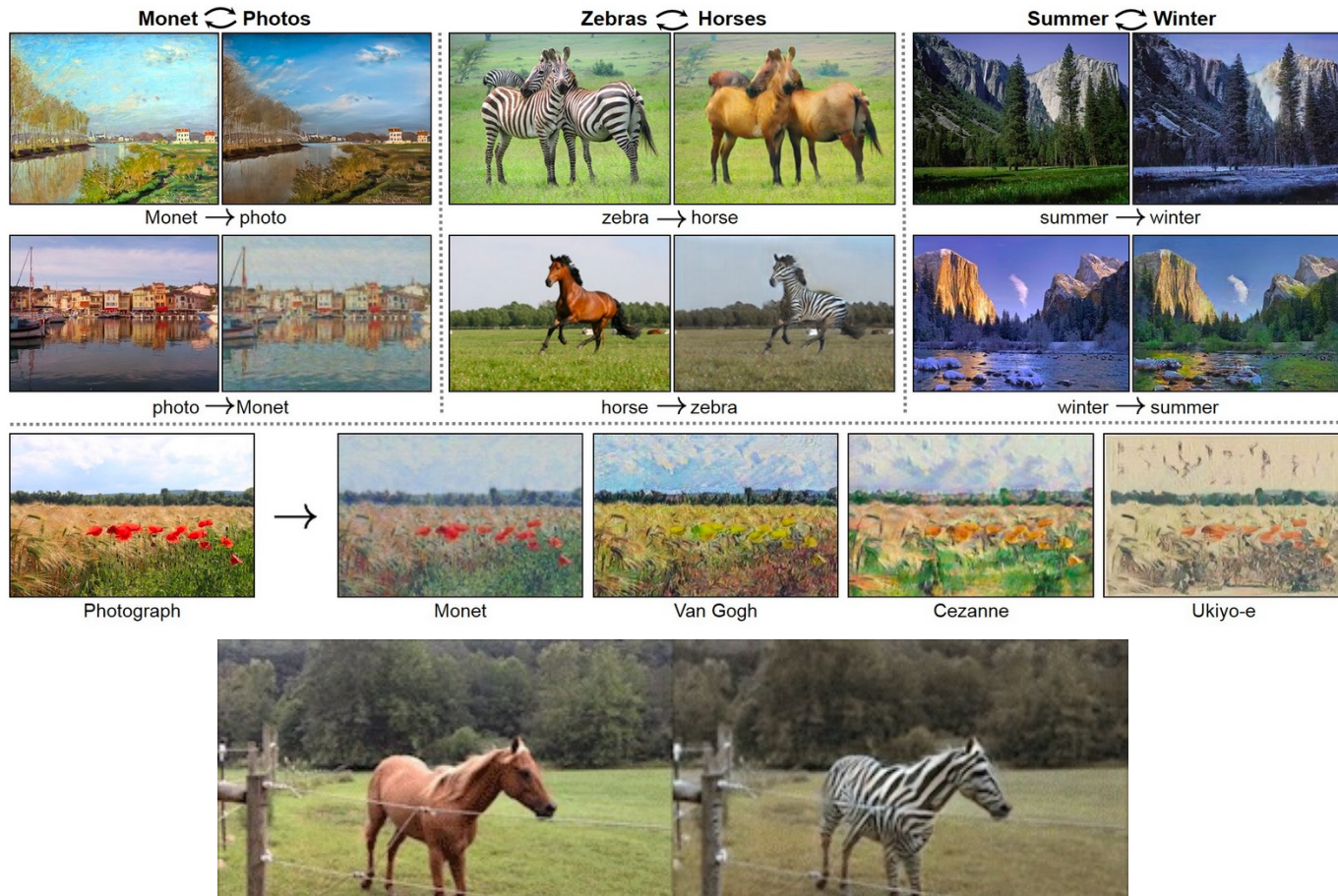
BigGAN

(Karras et al, 2019)



(Brock et al, 2018)

- $p(z)$ doesn't have to be random noise
- CycleGAN uses *cycle-consistency loss* in addition to GAN loss
 - Translating from $A \rightarrow B \rightarrow A$ should be consistent with original A



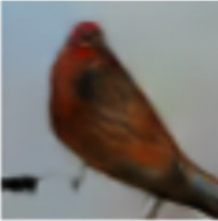

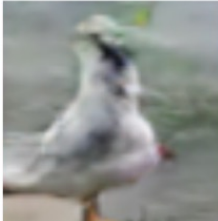

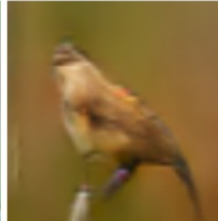
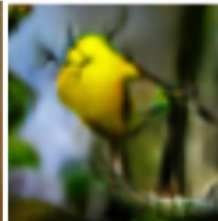
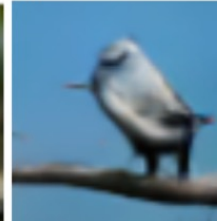

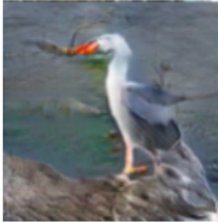

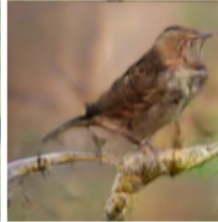
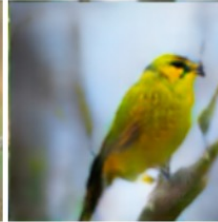




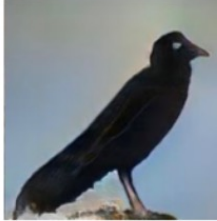


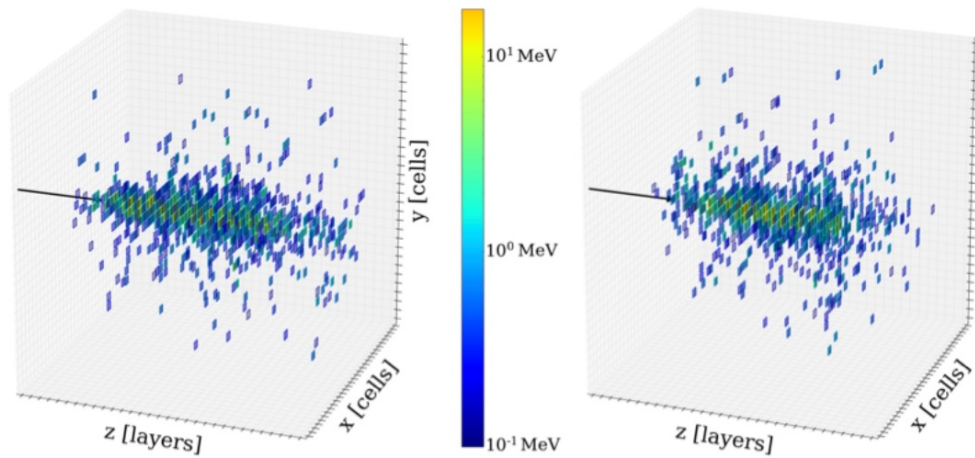
Text description	This bird is red and brown in color, with a stubby beak	The bird is short and stubby with yellow on its body	A bird with a medium orange bill white body gray wings and webbed feet	This small black bird has a short, slightly curved bill and long legs	A small bird with varying shades of brown with white under the eyes	A small yellow bird with a black crown and a short black pointed beak	This small bird has a white breast, light grey head, and black wings and tail
64x64 GAN-INT-CLS							
128x128 GAWWN							
256x256 StackGAN-v1							

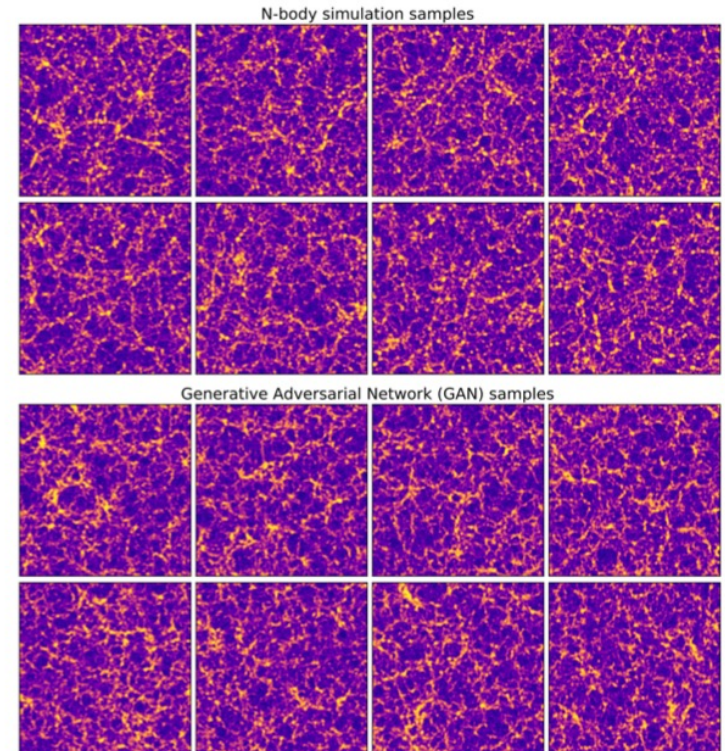
Fig. 3: Example results by our StackGAN-v1, GAWWN [29], and GAN-INT-CLS [31] conditioned on text descriptions from CUB test set.

(Zhang et al, 2017)

- Often studied for fast approximate simulation, simulation-based inference, optimization, ...



[2005.05334](#)



[1801.09070](#)