

# Introduction to Machine Learning: Lecture 3

Michael Kagan

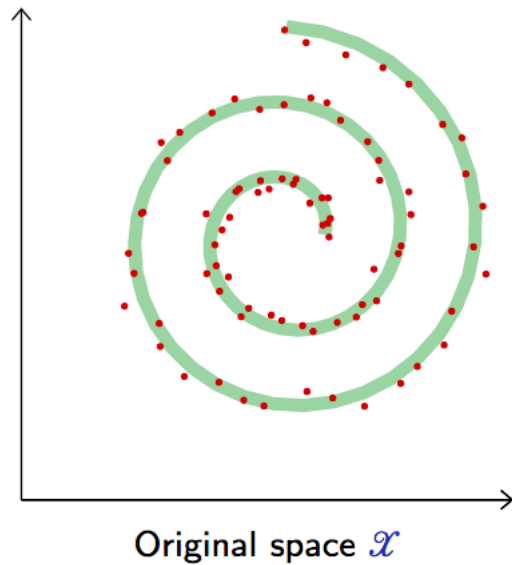
SLAC

Hadron Collider Physics Summer School

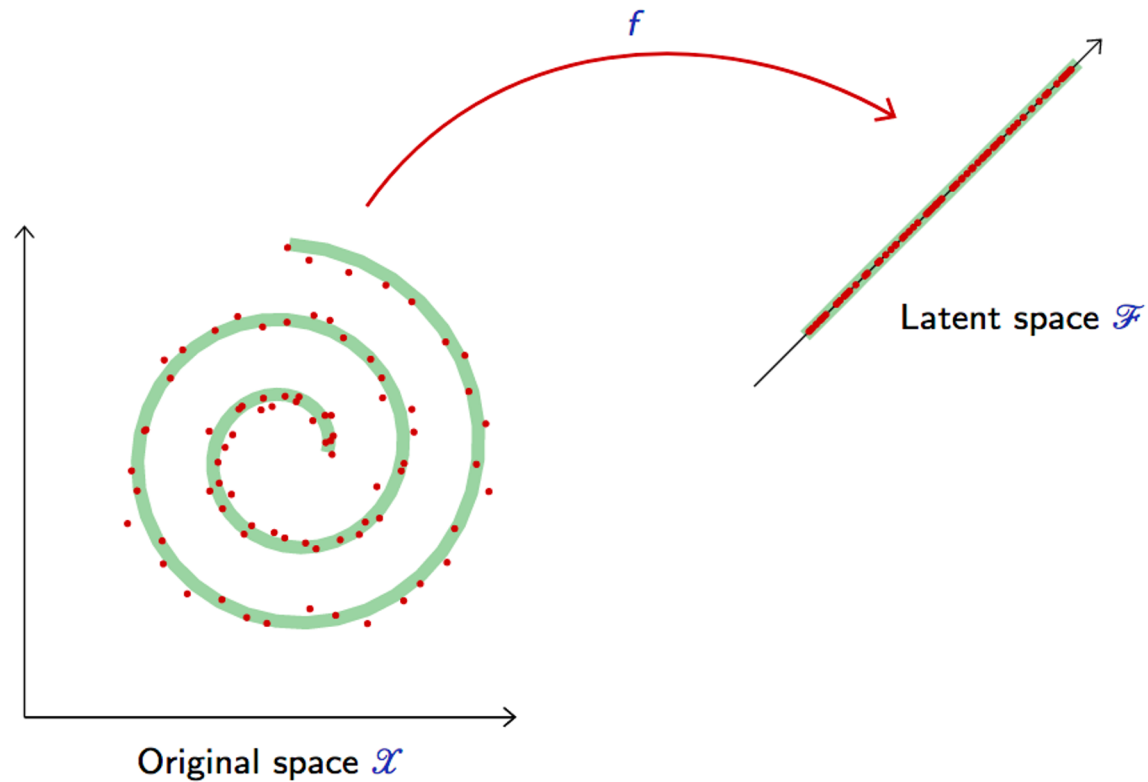
August 25, 2021

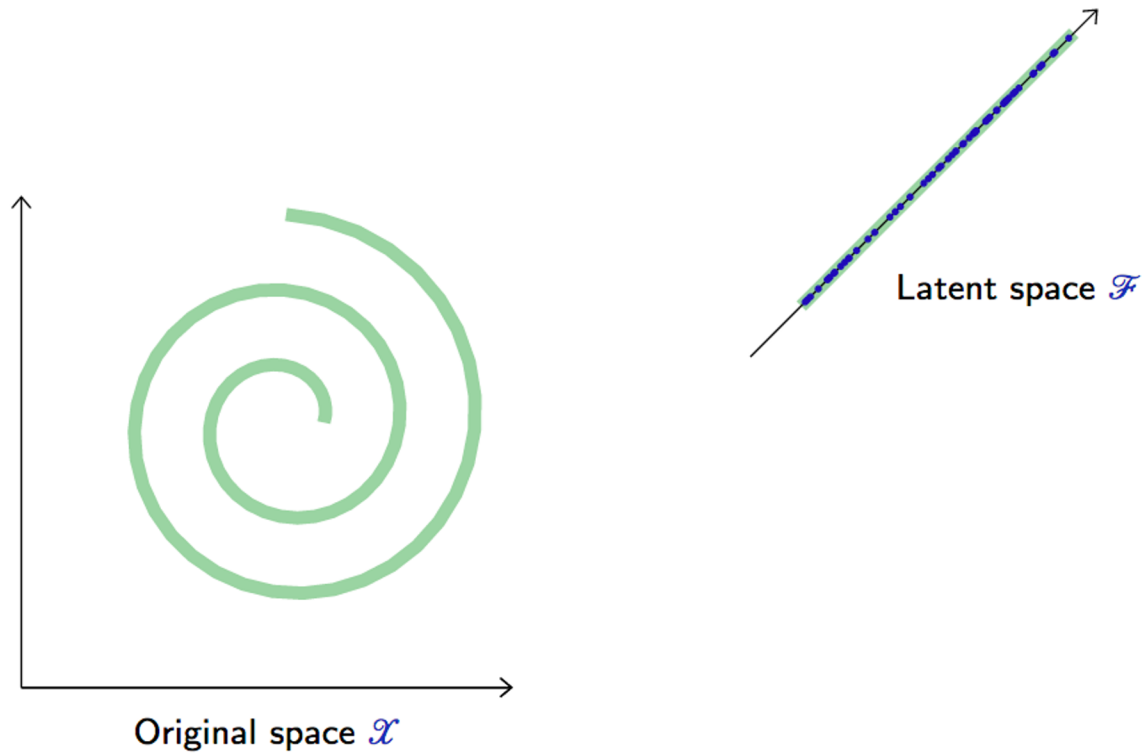
- Lecture 1
  - Brief introduction to probability and statistics
  - Introduction to Machine Learning fundamentals
  - Linear Models
- Lecture 2
  - Neural Networks
  - Deep Neural Networks
  - Convolutional, Recurrent, and Graph Neural Networks
- Lecture 3
  - Unsupervised Learning
  - Autoencoders
  - Generative Adversarial Networks and Normalizing Flows

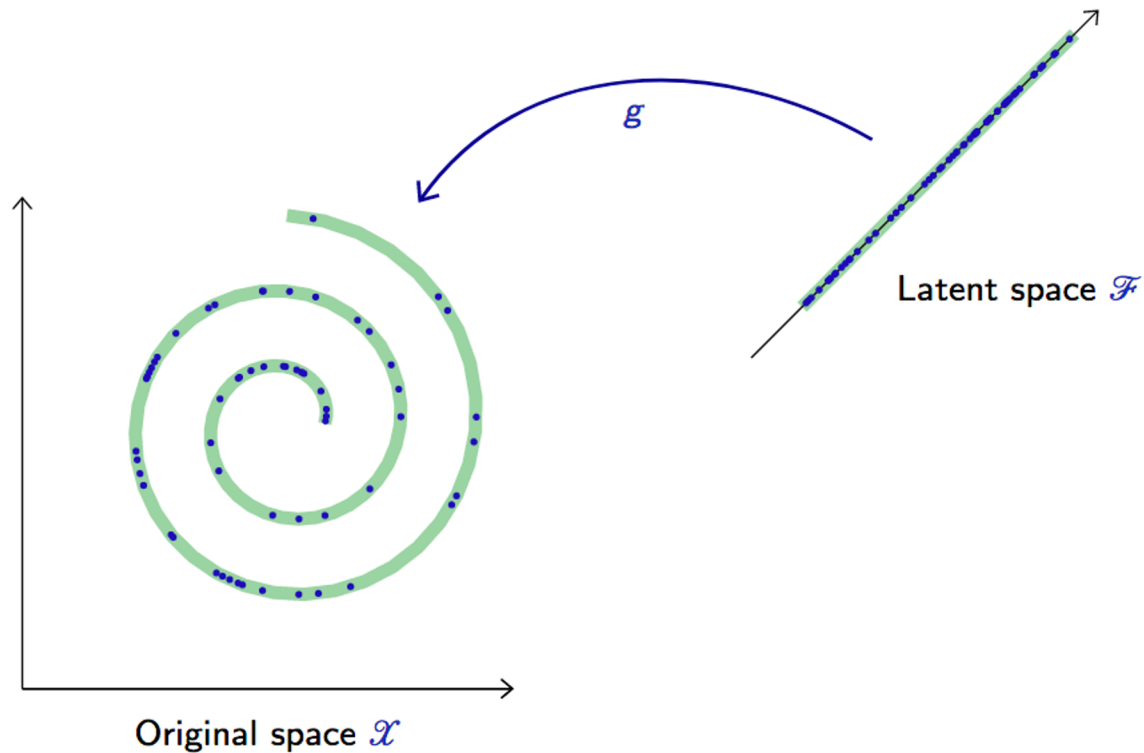
- Not all tasks are predicting a label from features
  - Data synthesis / simulation
  - Density estimation
  - Anomaly detection
  - Denoising, super resolution
  - Data compression
  - ...
- Requires **Unsupervised Learning**
- Often framed as **modeling the lower dimensional “meaningful degrees of freedom”** that describe the data











# Autoencoders

- Dimensionality Reduction / Compression
  - Compress the data to a *latent space* with smaller number of dimensions
  - Latent space must **encode and retain the important information about the data**
  - One way to frame “retaining important information”:  
**We can reconstruct original data from latent space**
  - Can we learn this compression and latent space?

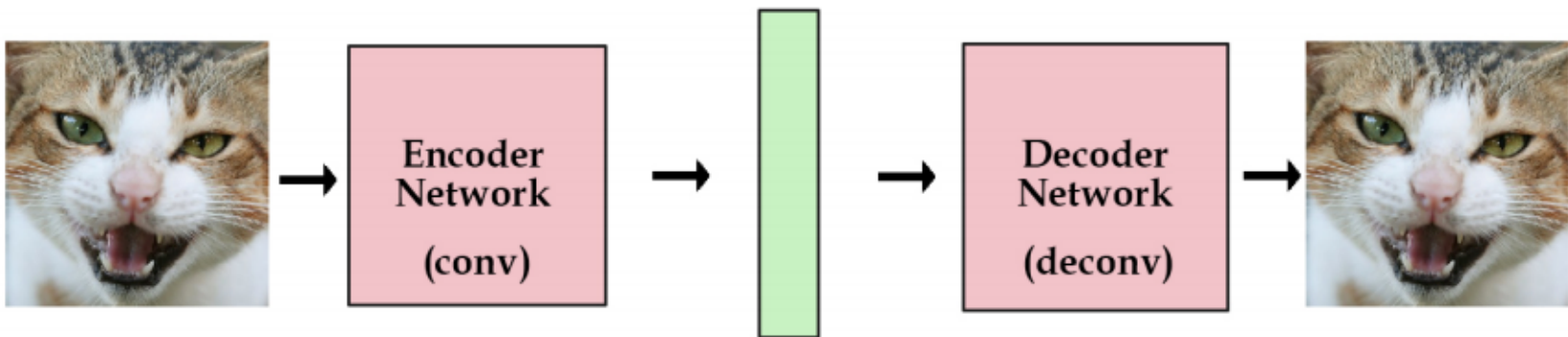
- Autoencoders map a space to itself through a **compression**

$x = \text{Data}$

$z = \text{Latent Space}$

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

Full transformation **should be close to the identity** on the data



- Autoencoders map a space to itself through a **compression**

$x = \text{Data}$

$z = \text{Latent Space}$

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

Full transformation **should be close to the identity** on the data

$$x \rightarrow f(x) = z$$

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

- Autoencoders map a space to itself through a **compression**

$$x = \text{Data} \qquad z = \text{Latent Space}$$

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

Full transformation **should be close to the identity** on the data

$$x \rightarrow f(x) = z$$

$$z \rightarrow g(z) = \hat{x}$$

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



- Autoencoders map a space to itself through a **compression**

$x = \text{Data}$

$z = \text{Latent Space}$

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

Full transformation **should be close to the identity** on the data

$$x \rightarrow f(x) = z$$

$$z \rightarrow g(z) = \hat{x}$$

We must:

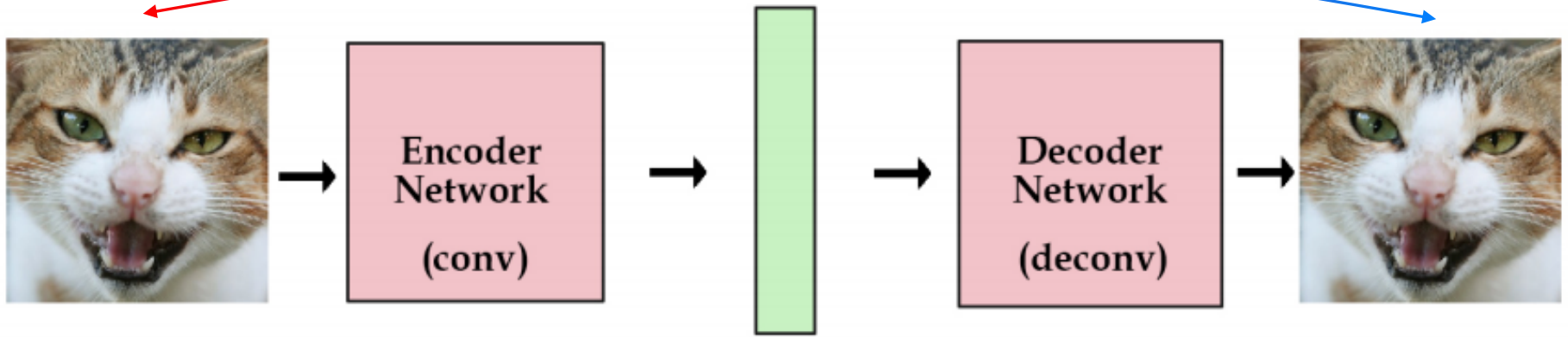
- Choose latent dimension  $D$
- Learn mapping  $f(\cdot)$  and  $g(\cdot)$

– **Encoder**: Map from to a lower dimensional latent space

- Neural network  $f_{\theta}(x)$  with parameters  $\theta$

– **Decoder**: Map from latent space back to data space

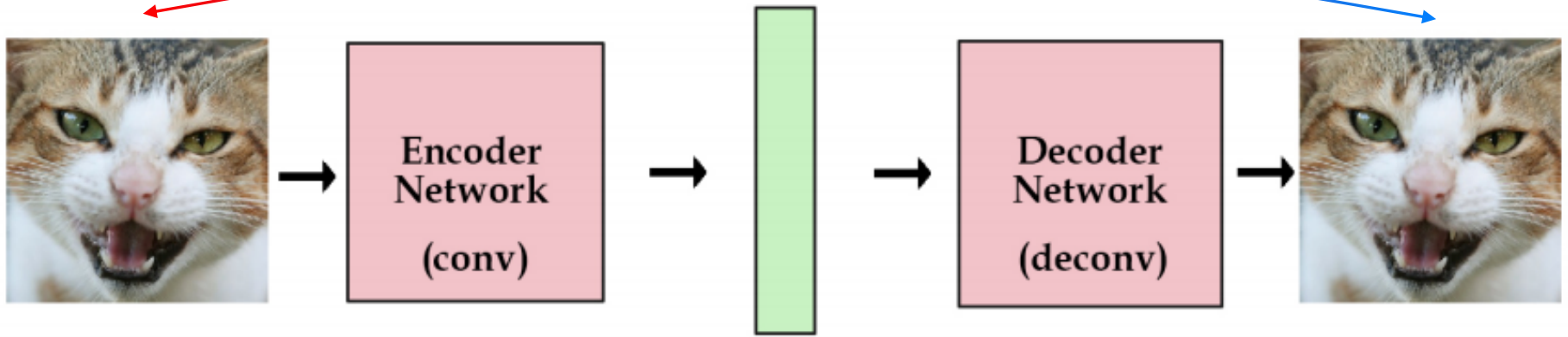
- Neural network  $g_{\psi}(z)$  with parameters  $\psi$



- Loss: mean *reconstruction loss* (MSE) between data and encoded-decoded data

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

- Minimize this loss over parameters of encoder ( $\theta$ ) and decoder ( $\psi$ ).

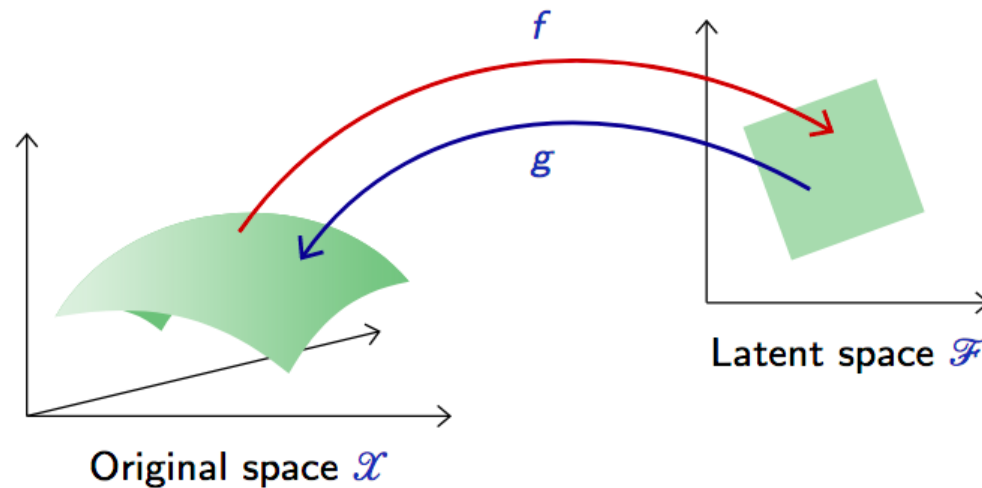


- Loss: mean *reconstruction loss* (MSE) between data and encoded-decoded data

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

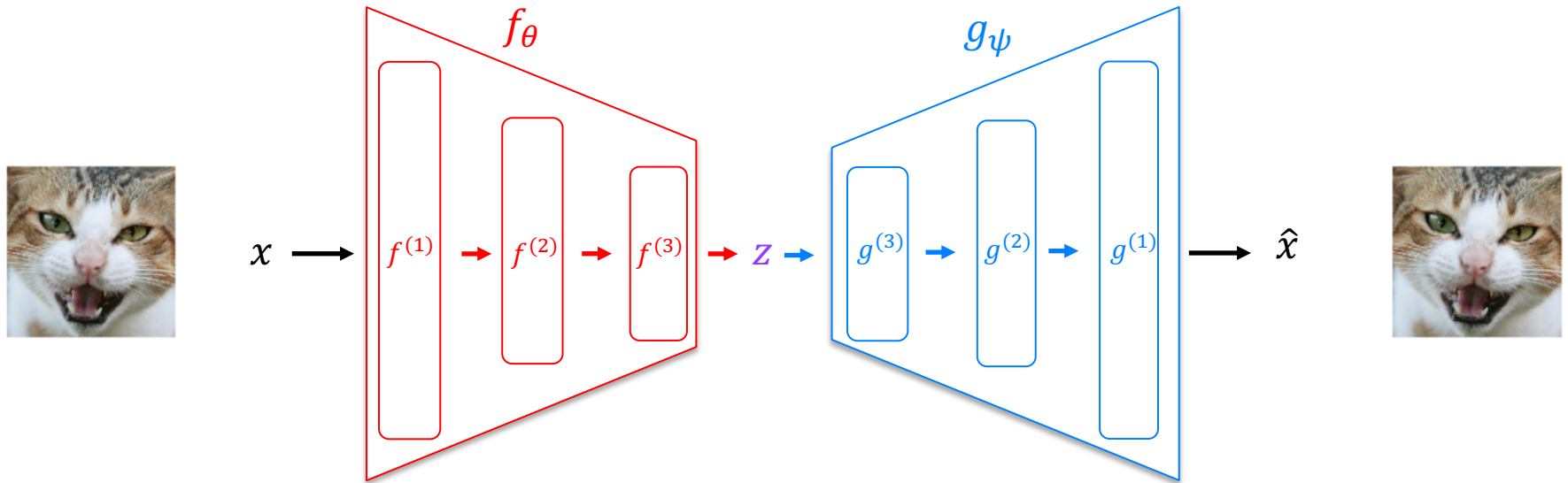
- Minimize this loss over parameters of encoder ( $\theta$ ) and decoder ( $\psi$ ).

NOTE: if  $f_\theta(x)$  and  $g_\psi(z)$  are linear, optimal solution given by Principle Components Analysis



If the latent space is of lower dimension:

autoencoder must capture a “good” parametrization, and in particular dependencies between components



- When  $f_\theta$  and  $g_\psi$  are multiple neural network layers, can learn complex mappings between  $x$  and  $z$ 
  - $f_\theta$  and  $g_\psi$  can be Fully Connected, CNNs, RNNs, etc.
  - Choice of network structure will depend on data

$X$  (original samples)

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

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



$f_\theta$  and  $g_\psi$  are each  
5 convolutional layers

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

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

7 2 1 0 9 1 4 9 9 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

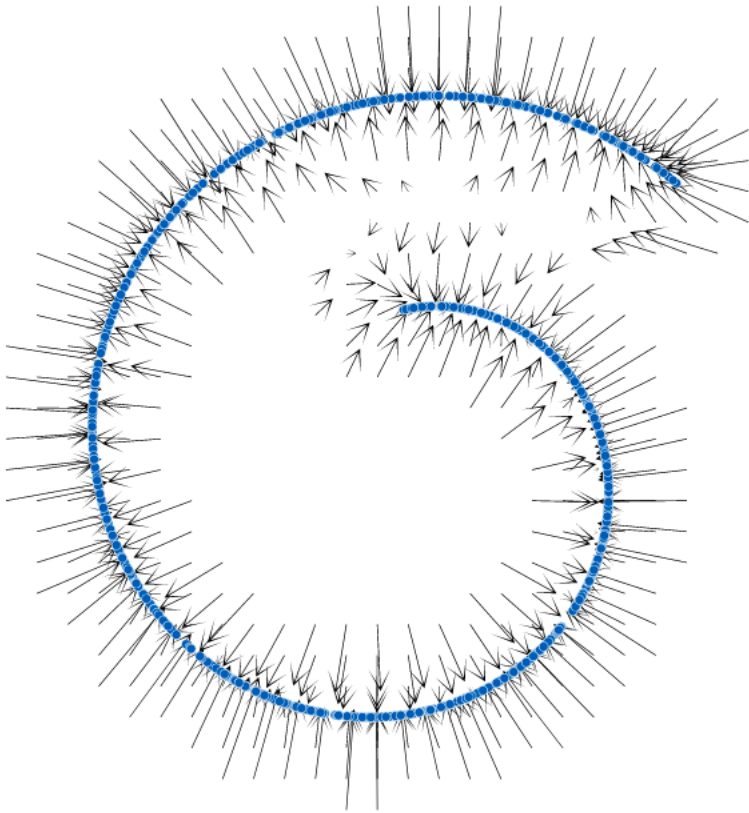
- Learn a mapping from corrupted data space  $\tilde{\mathcal{X}}$  back to original data space
  - Mapping  $\phi_w(\tilde{x}) = x$
  - $\phi_w$  will be a neural network with parameters  $w$
- Loss:

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



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

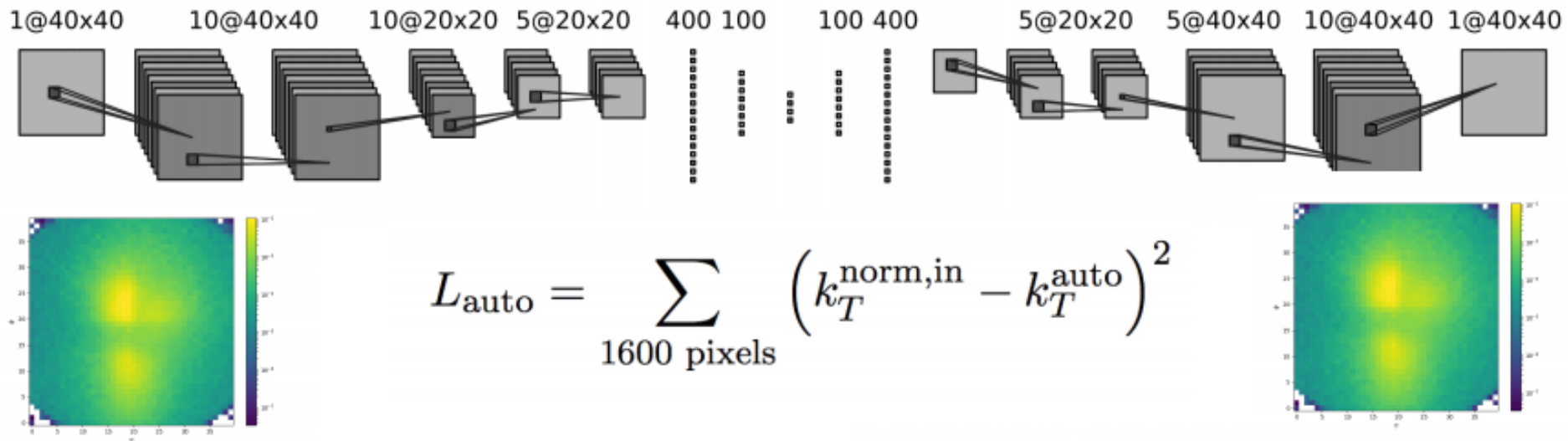
Corrupted ( $\sigma = 4$ )

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

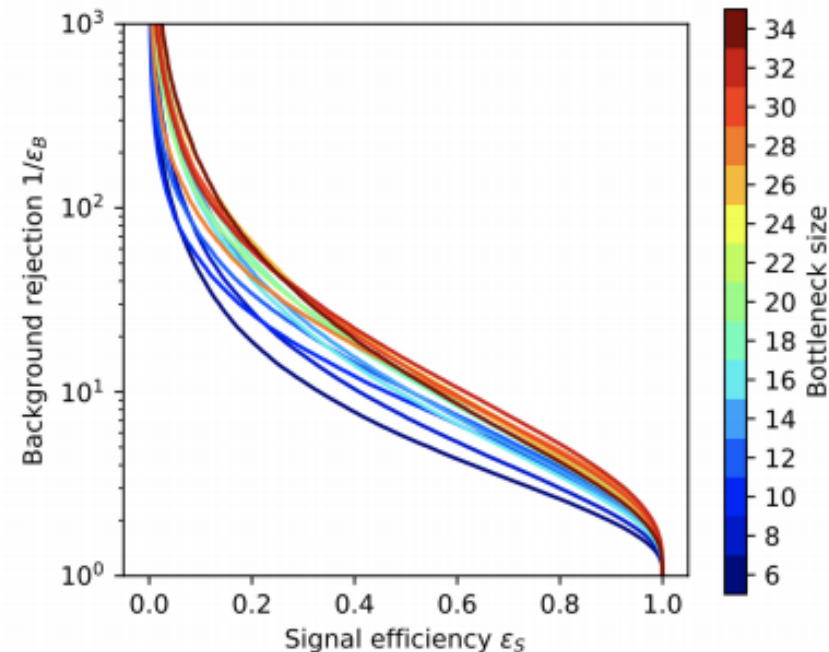
Reconstructed

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





- Train on pure QCD light quark/ gluon jets and apply to top tagging
- Top quarks/ new physics identified as anomaly if poorly reconstructed by autoencoder



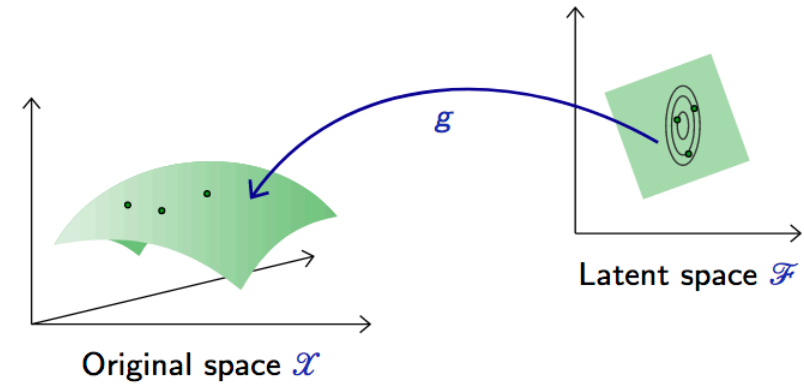
*QCD or What?*

T Heime1, GK, T Plehn, JM Thompson, 1808.08979

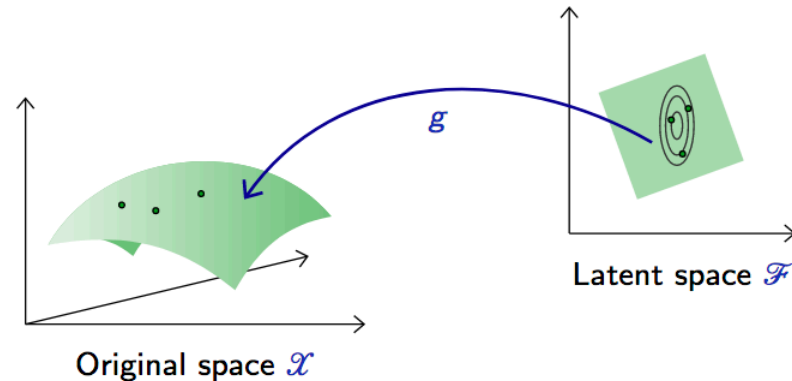
*Searching for New Physics with Deep Autoencoders*

M Farina, Y Nakai, D Shih, 1808.08992

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



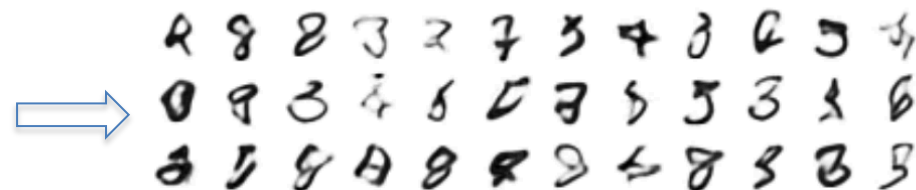
- Can we sample in latent space and decode to generate data?



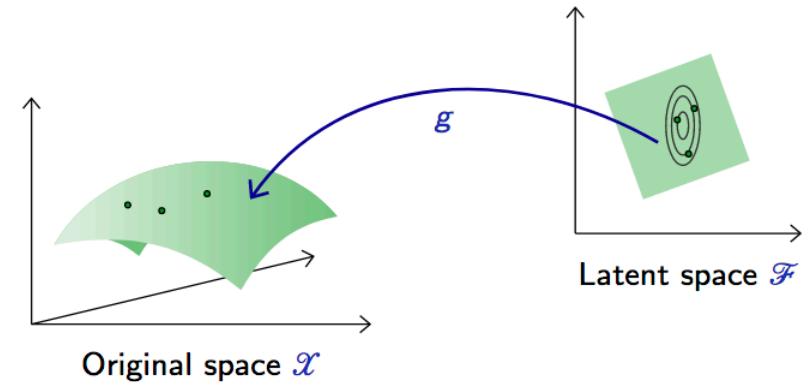
- What distribution to sample from in latent space?

- Try Gaussian with mean and variance from data

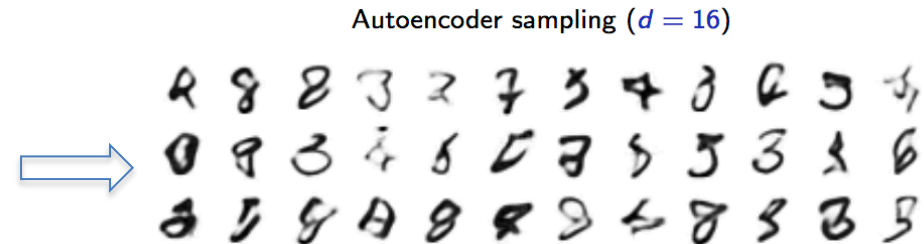
Autoencoder sampling ( $d = 16$ )



- 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



- Doesn't work! Don't know the right latent space density
  - This can be done with a Variational Autoencoder (See Backup)

- Generative models aim to:
  - Learn distribution  $p(x)$  that models PDF of 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 from  $p(x)$ , but not evaluate density



- Formulate generative modeling task 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:
  - *A generator that can produce samples*
  - *A measure of not too far from the real data*

- **Generator network  $g_{\theta}(z)$**  with parameters  $\theta$ 
  - 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 learned distribution  $p_{\theta}(x)$  is, but we can sample from it  $\rightarrow$  *Implicit Model*

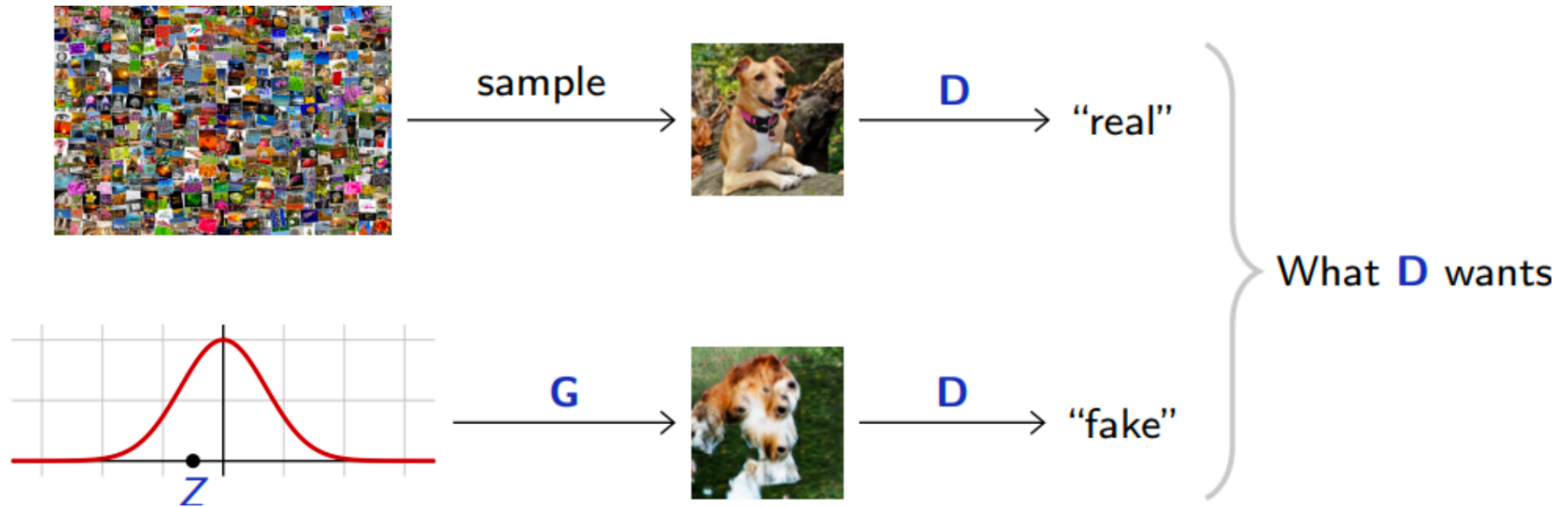


- **Generator network  $g_{\theta}(z)$**  with parameters  $\theta$ 
  - Map sample from known  $p(z)$  to sample in data space

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

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

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



- **Generator goal:**
  - Produce *fake* data to trick discriminator to classify as *real*
- **Discriminator goal:**
  - Minimizes miss-classification of data as real or fake
- **Adversarial setup:** two networks w/ 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:  $z_i \sim p(z)$

↑  
Usually Gaussian  $\mathcal{N}(0,1)$

- 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 binary cross entropy

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

- 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 binary cross entropy

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

- 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 binary cross entropy

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

- Generator isn't fixed  $\rightarrow$  Must be trained

- Consider objective as a *value function* of  $\phi$  and  $\theta$

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

- Consider objective as a *value function* of  $\phi$  and  $\theta$

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

- For fixed generator,  $V(\phi, \theta)$  is high when discriminator is good, i.e. when generator is not producing good fakes



- Consider objective as a *value function* of  $\phi$  and  $\theta$

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

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

- Consider objective as a *value function* of  $\phi$  and  $\theta$

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

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

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

- Consider objective as a *value function* of  $\phi$  and  $\theta$

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

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

- So our optimization goal becomes:

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

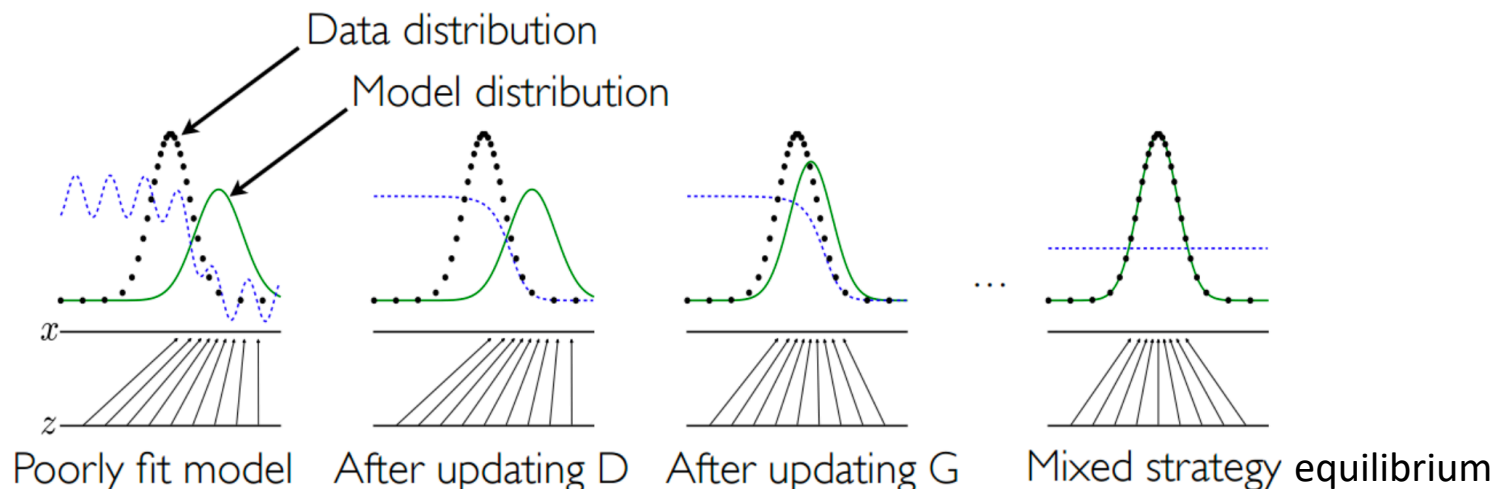
NOTE: can prove that minimax solution corresponds to generator that perfectly reproduces data distribution

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

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

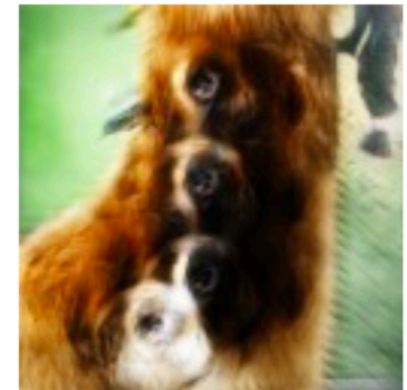
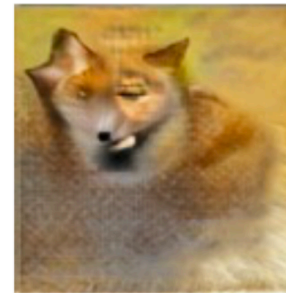
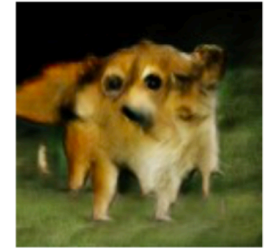
$$\phi \leftarrow \phi - \gamma \nabla_{\phi} V(\phi, \theta) = \phi - \gamma \frac{\partial V}{\partial d} \frac{d(d_{\phi})}{d\phi}$$

- For each  $\theta$  step, take  $k$  steps in  $\phi$  to keep discriminator near optimal



# Examples

Goodfellow et. al., 2014

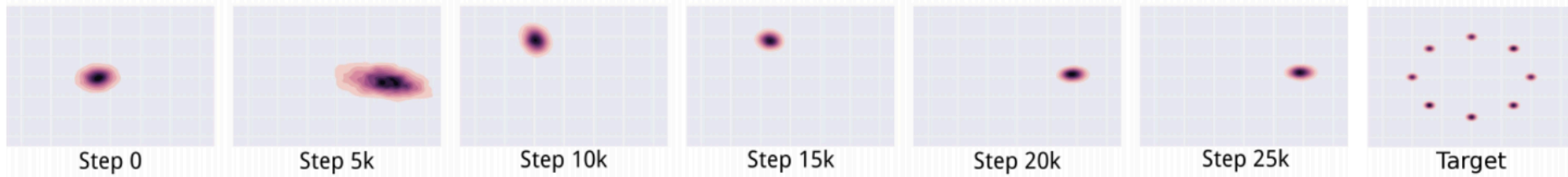


Not so good

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:** is generated data good enough?
- **Improvements** in training objective (WGAN) and model design have significantly helped with these challenges



## StyleGAN v2



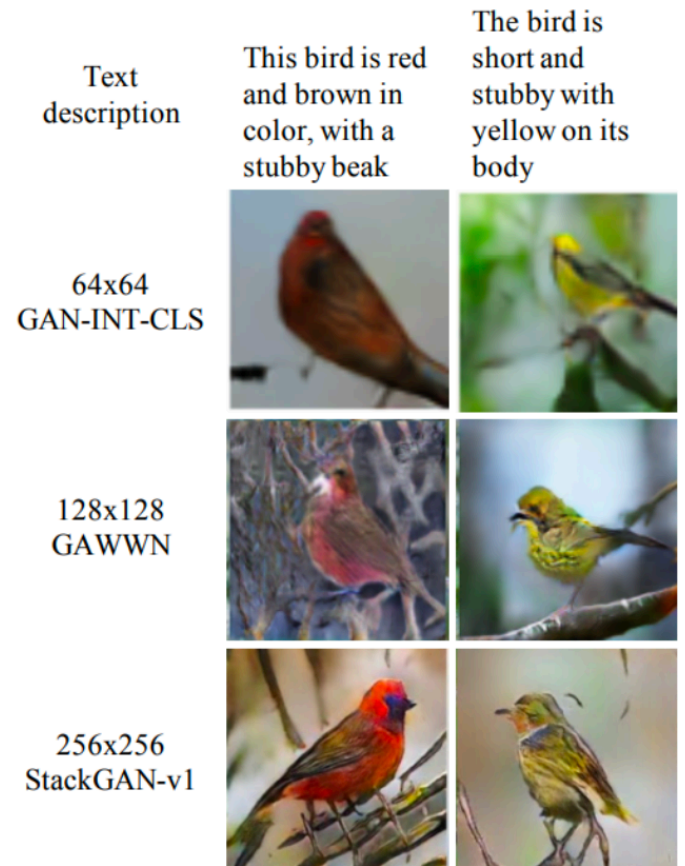
(Karras et al, 2019)

## Image-to-Image Translation with CycleGAN

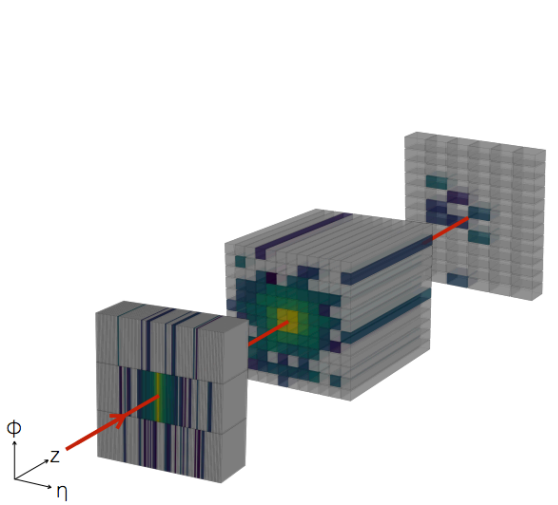
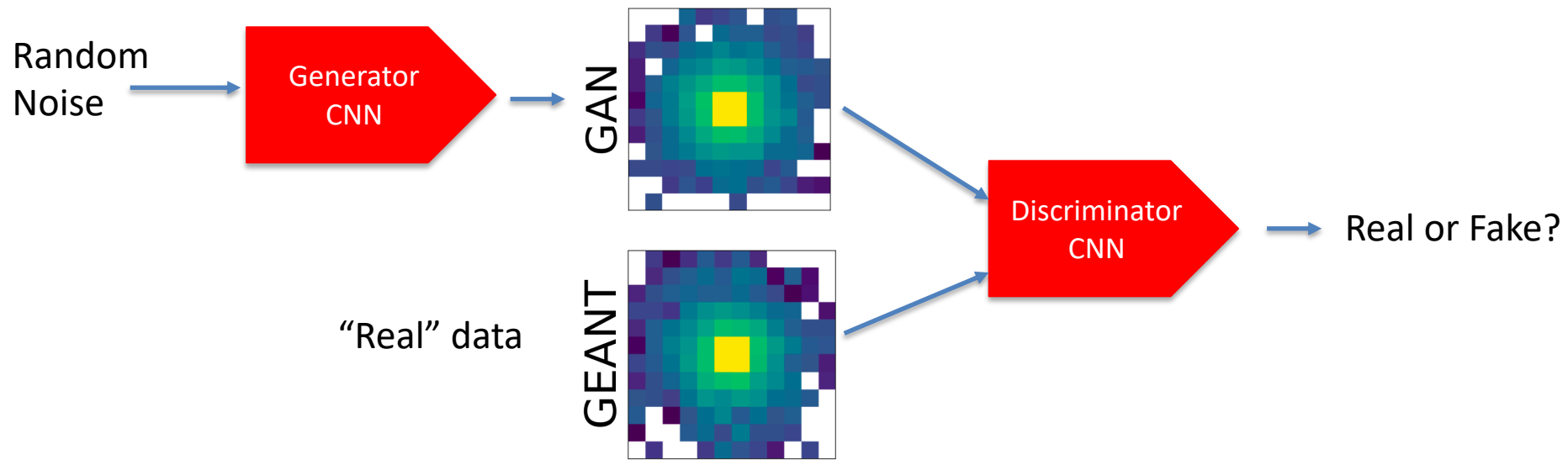


[Zhu et. al. 2017](#)

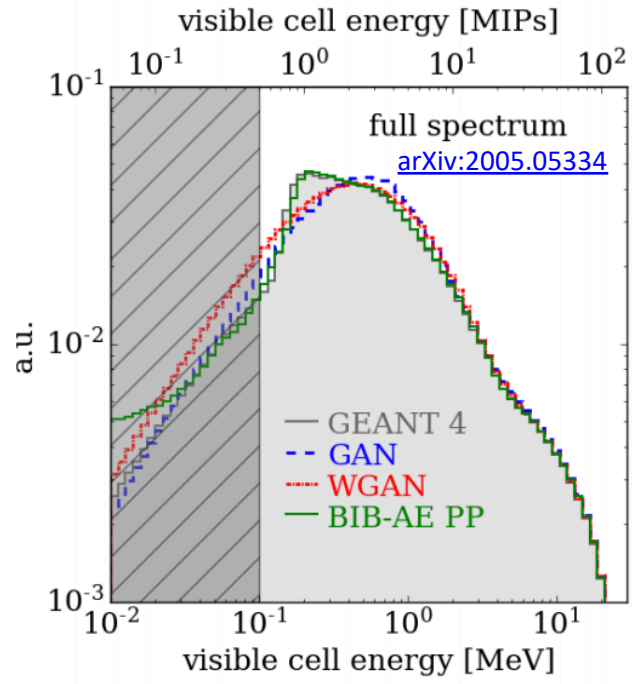
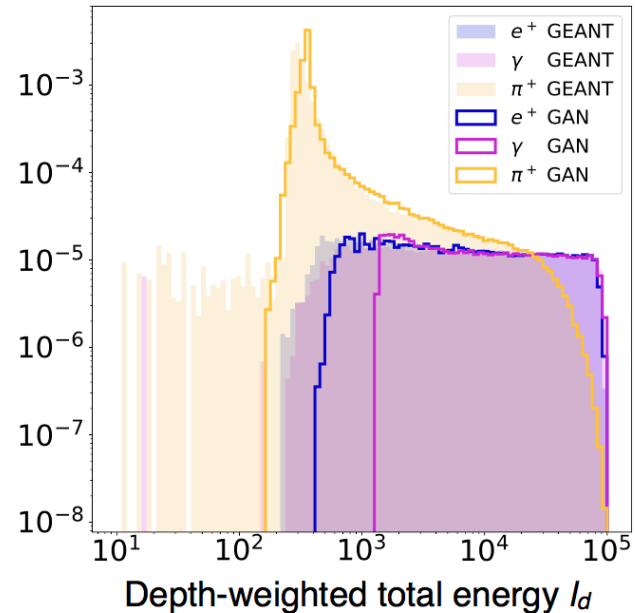
## Text-to-Image Synthesis with StackGAN



Zhang et. al. 2017



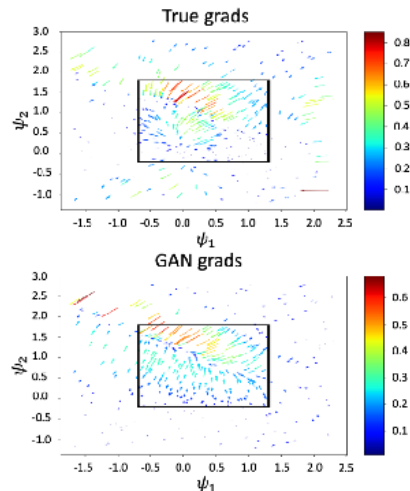
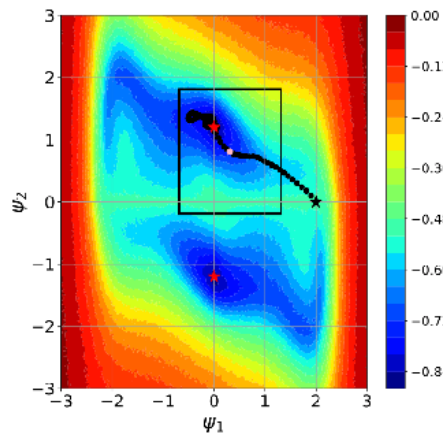
[PRD97, 014021 \(2018\)](#)  
[arXiv:1705.02355](#)  
[arXiv:1701.05927](#)



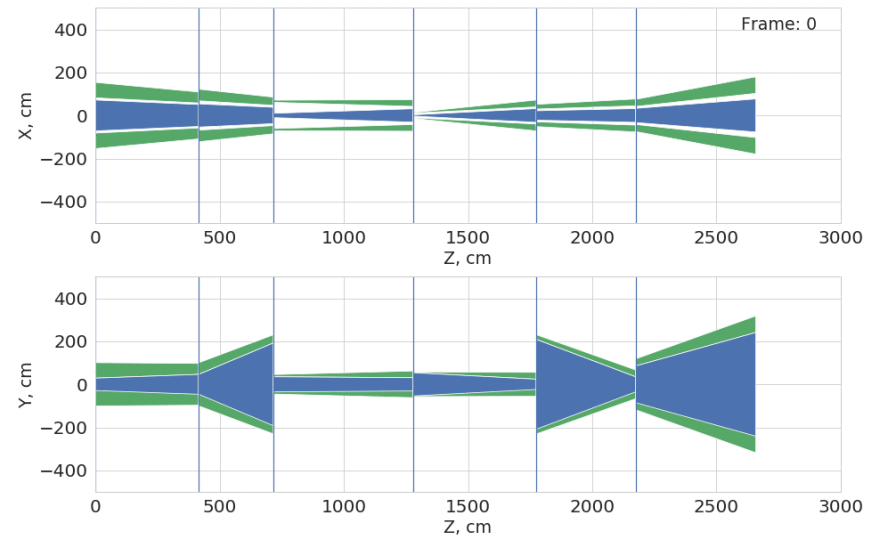
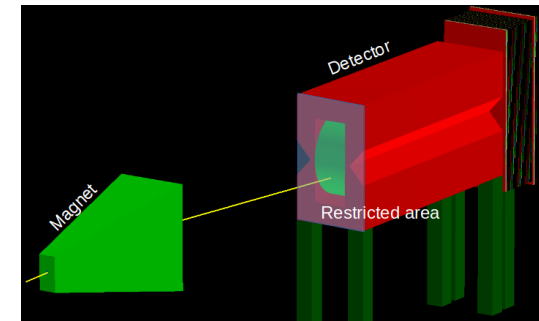


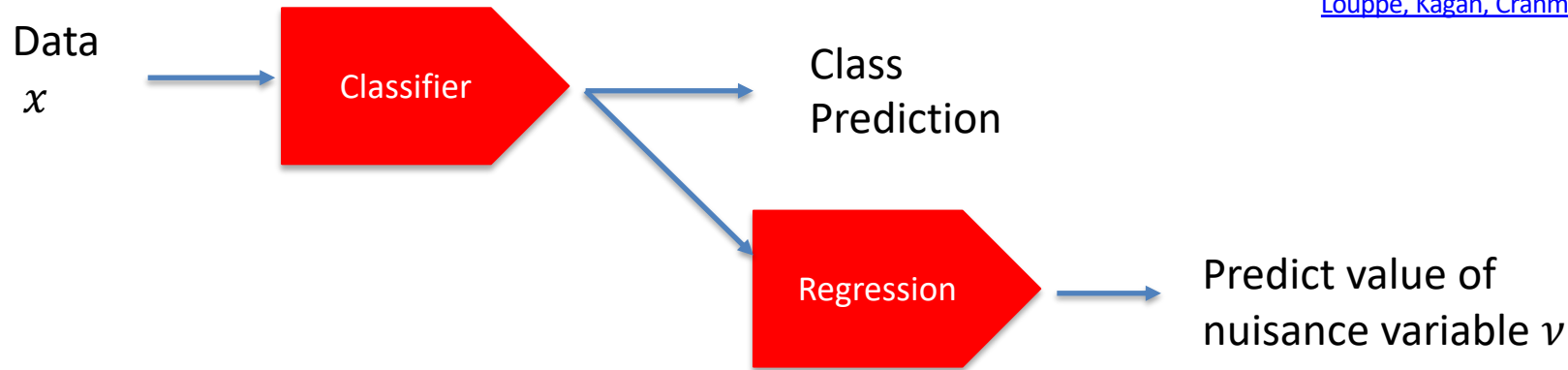
- Train GAN to emulate data from simulated detector,  $\tilde{x} = g(z|\psi)$  conditioned on detector parameters  $\psi$  (e.g. magnet shape below)
- Define objective  $\mathcal{C}$  to minimize:  $\min_{\psi} \mathbb{E}_{\tilde{x}}[\mathcal{C}(\tilde{x} = g(z|\psi))]$
- GAN is differentiable
  - Minimize with gradient descent

## Toy Example

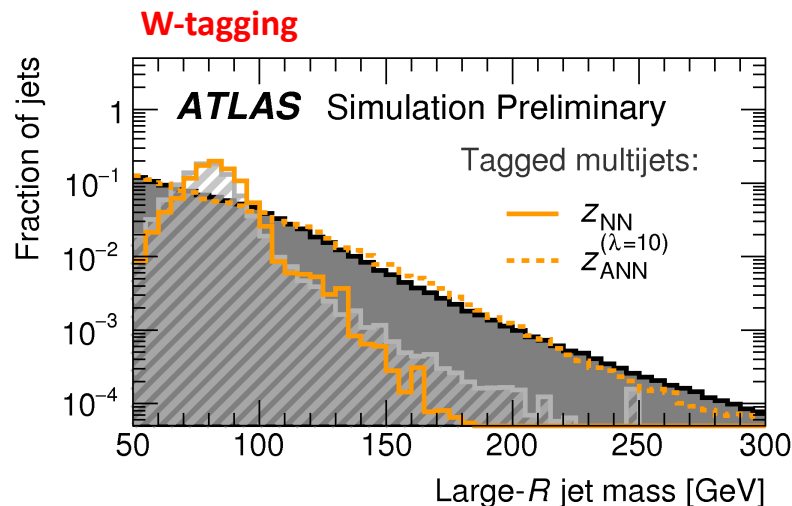
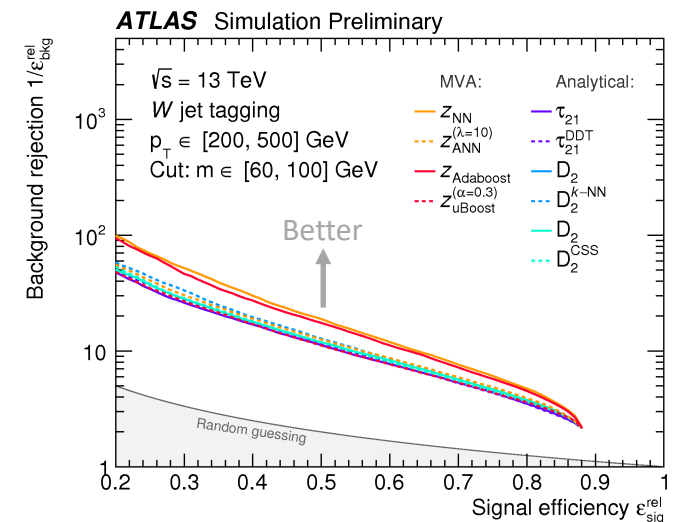


## Magnet Optimization





- Want to remove dependence of classifier on “nuisance” variable  $\nu$ , e.g. a systematic, mass, etc.



[ATL-PHYS-PUB-2018-014](#)


- GAN can only learn to sample from a distribution
- Is there a way to learn the explicit density  $p(x)$  ?

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

Multivariate:

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



Change of volume:  
Determinant of Jacobian  
of the transformation

- If  $f(x)$  is the pdf of  $x$  and  $y(x)$  is a change of variables, since probability should not change:

$$P(x_a < x < x_b) = P(y(x_a) < y < y(x_b))$$

$$\int_{x_a}^{x_b} f(x) dx = \int_{y(x_a)}^{y(x_b)} g(y) dy$$

$$= \int_{x_a}^{x_b} g(y(x)) \left| \frac{dy}{dx} \right| dx$$

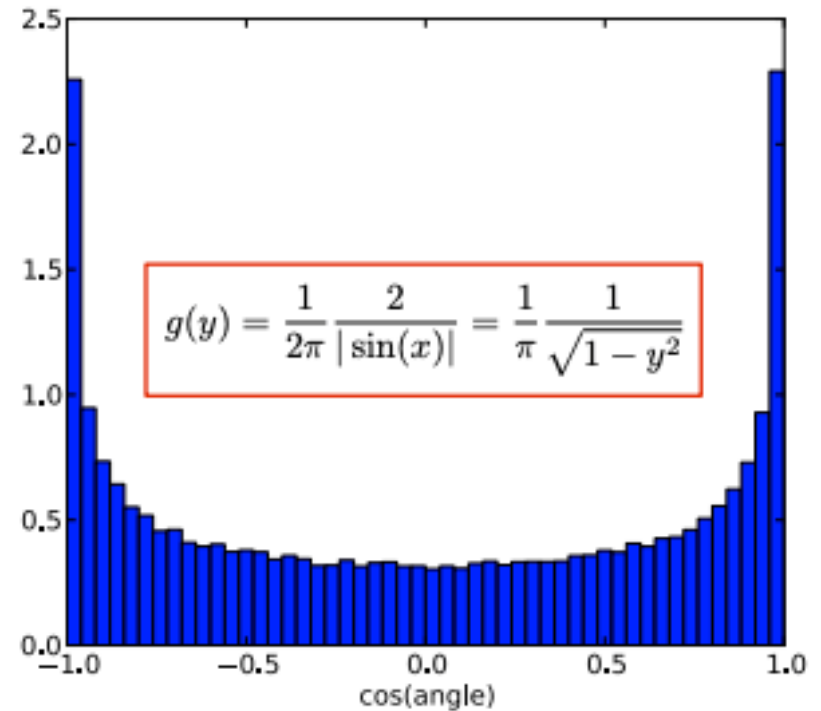
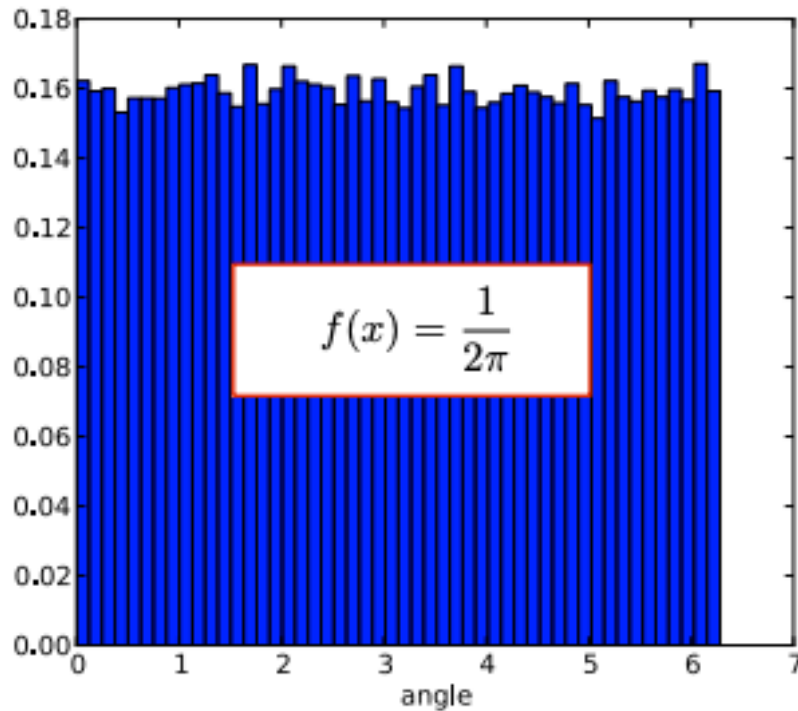
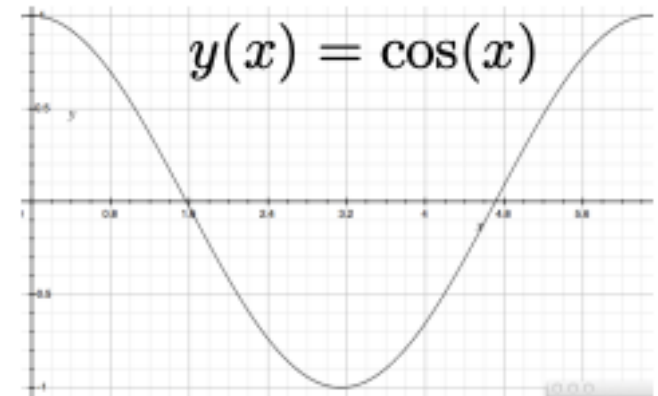
Rewrite of r.h.s.

Thus:

$$f(x) = g(y) \left| \frac{dy}{dx} \right|$$

Distributions are related through the Jacobian

$$f(x) = g(y) \left| \frac{dy}{dx} \right|$$



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

- $\mathbf{x} \equiv$  data we want to model
- $\mathbf{z} \sim p(\mathbf{z})$  is a chosen noise distribution, usually Gaussian
- $\phi$  is continuous, invertible, differentiable,  $\mathbf{z} = \phi^{-1}(\mathbf{x})$

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

- $\mathbf{x} \equiv$  data we want to model
- $\mathbf{z} \sim p(\mathbf{z})$  is a chosen noise distribution, usually Gaussian
- $\phi$  is continuous, invertible, differentiable,  $\mathbf{z} = \phi^{-1}(\mathbf{x})$
- Want to find  $\phi(\mathbf{z})$  that transforms data  $\mathbf{x} \Leftrightarrow$  noise  $\mathbf{z} \sim p(\mathbf{z})$



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

- $\mathbf{x} \equiv$  data we want to model
- $\mathbf{z} \sim p(\mathbf{z})$  is a chosen noise distribution, usually Gaussian
- $\phi$  is continuous, invertible, differentiable,  $\mathbf{z} = \phi^{-1}(\mathbf{x})$
- Want to find  $\phi(\mathbf{z})$  that transforms data  $\mathbf{x} \Leftrightarrow$  noise  $\mathbf{z} \sim p(\mathbf{z})$

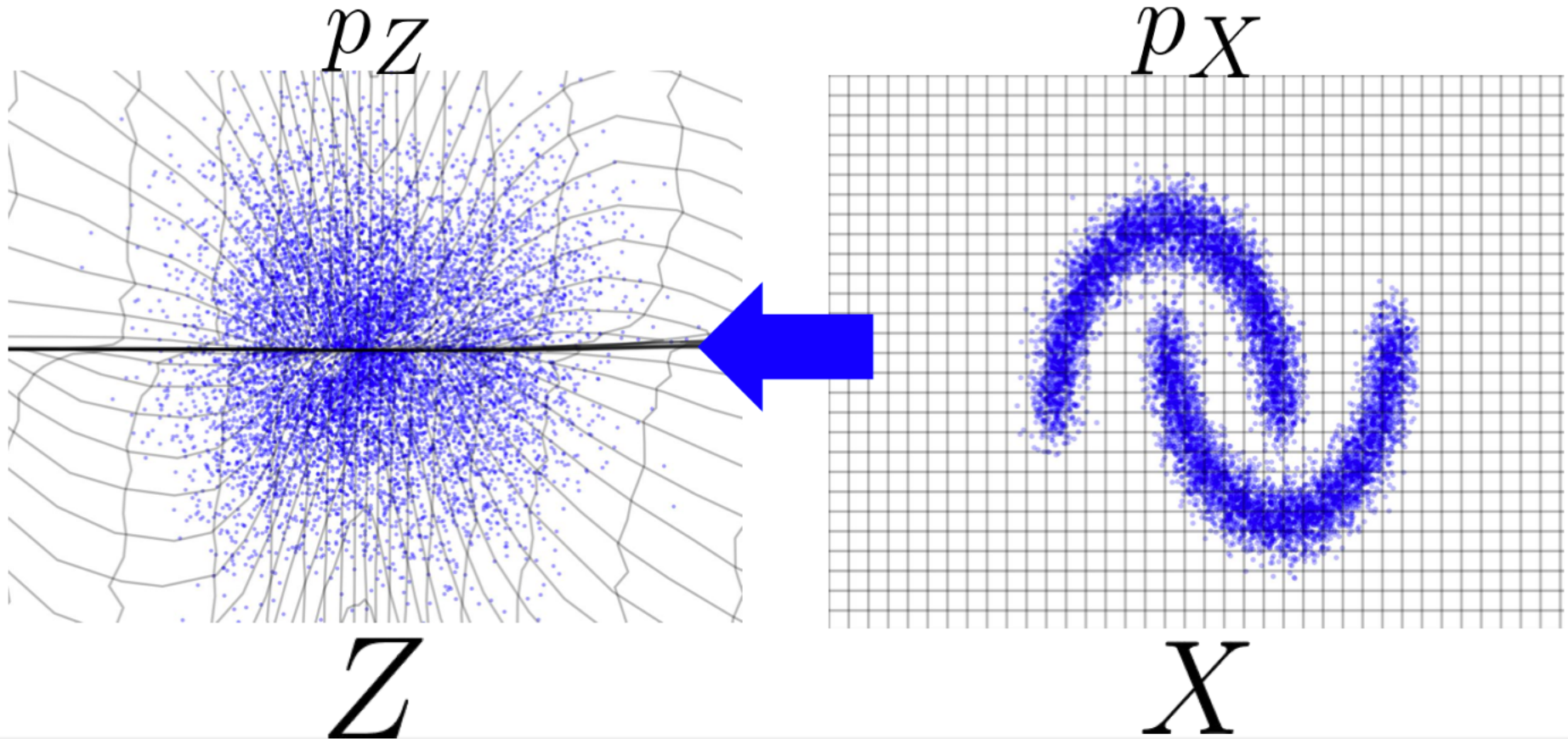
$\phi(\mathbf{z})$  neural network

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

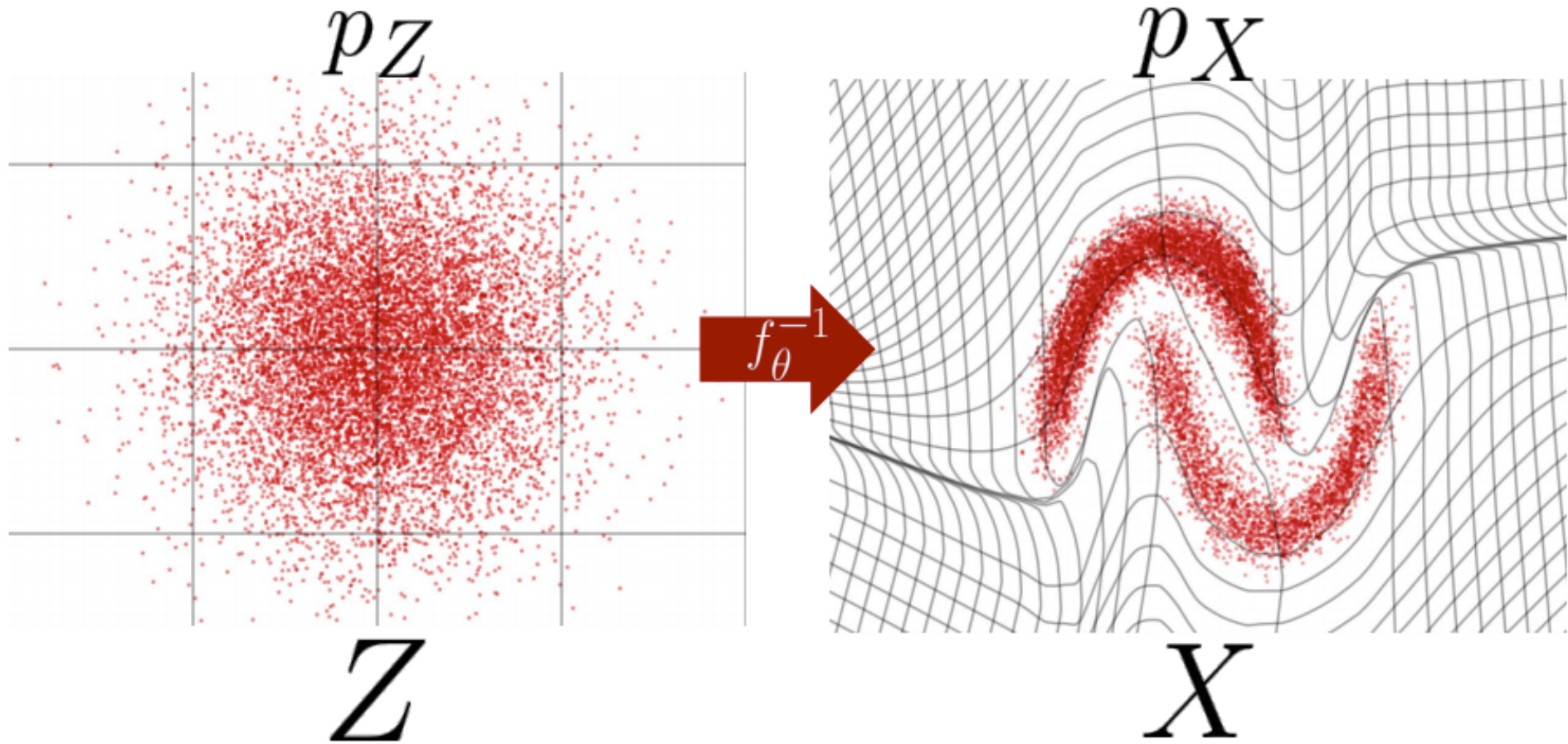


$\phi^{-1}(\mathbf{x})$  inverse

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



## Generation through process reversion



- Data vector  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$
- **Transformation:** where  $f(\cdot)$  and  $g(\cdot)$  are neural networks

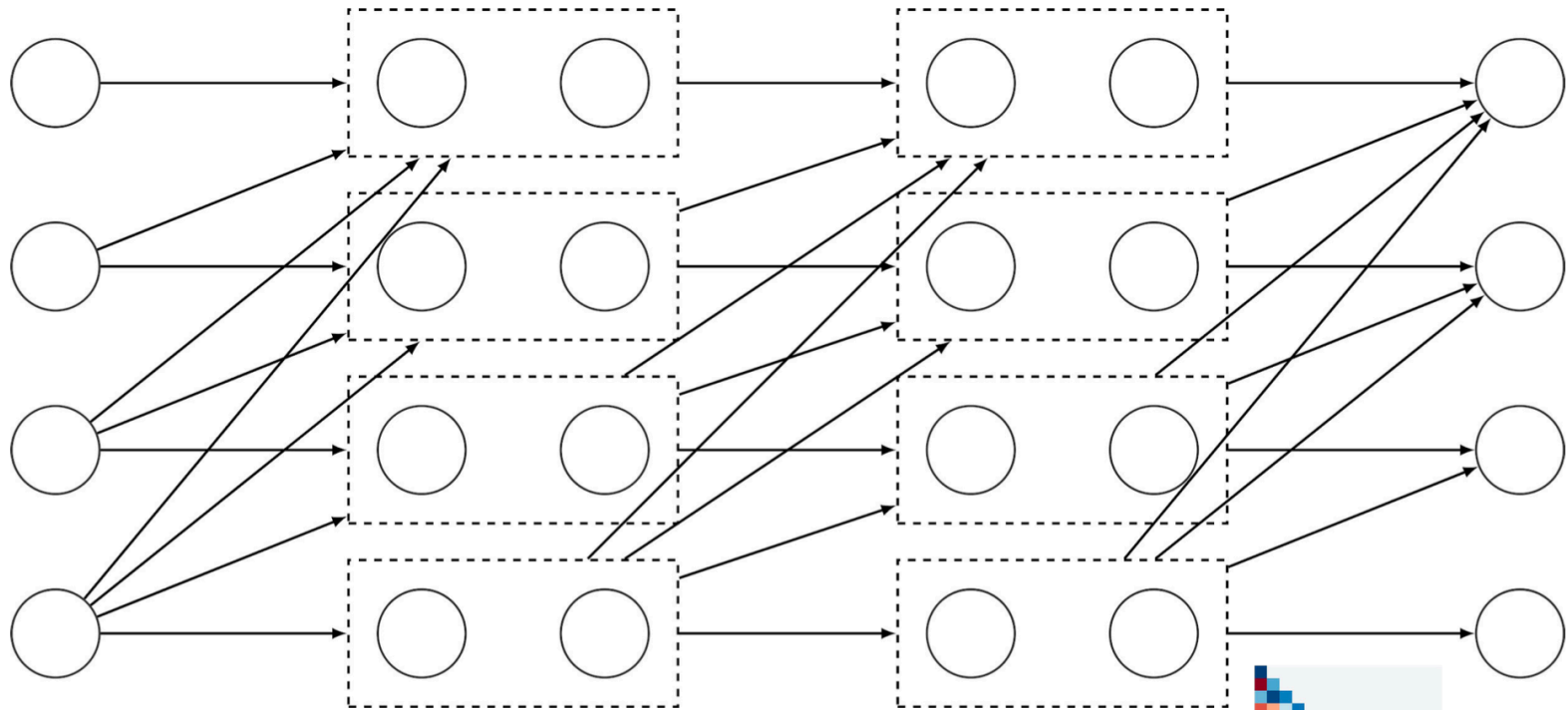
$$\phi(\mathbf{z}): \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \phi_1(\mathbf{z}) = z_1 \\ \phi_2(\mathbf{z}) = z_2 f(z_1) + g(z_1) \end{pmatrix}$$

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

- **Determinant:** Use fact that Jacobian is lower triangular

$$\det\left(\frac{\partial\phi(\mathbf{z})}{d\mathbf{z}}\right) = \det\left(\begin{pmatrix} 1 & 0 \\ \frac{\partial\phi_2(\mathbf{z})}{dz_1} & f(z_1) \end{pmatrix}\right) = \prod \text{Diag}\left(\frac{\partial\phi(\mathbf{z})}{d\mathbf{z}}\right) = f(z_2)$$

(Bengio<sup>2</sup>, 1999; Larochelle & Murray, 2011; van den Oord et al., 2015; Uria et al., 2016)



$$z_d = f_d(x_d; x_{<d})$$

$$x_d = f_d^{-1}(z_d; x_{<d})$$

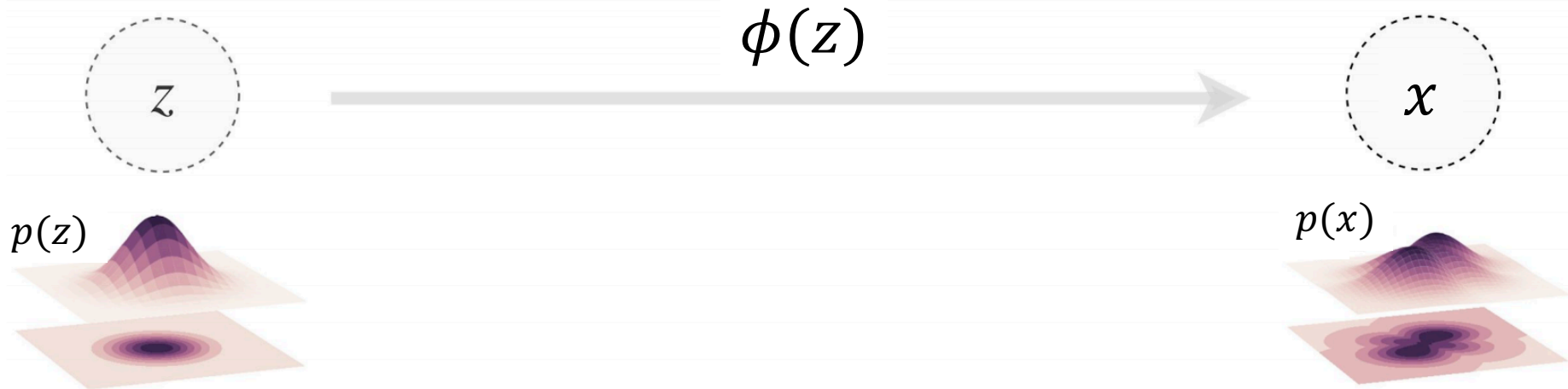


$$f_3 \circ f_2 \circ f_1$$

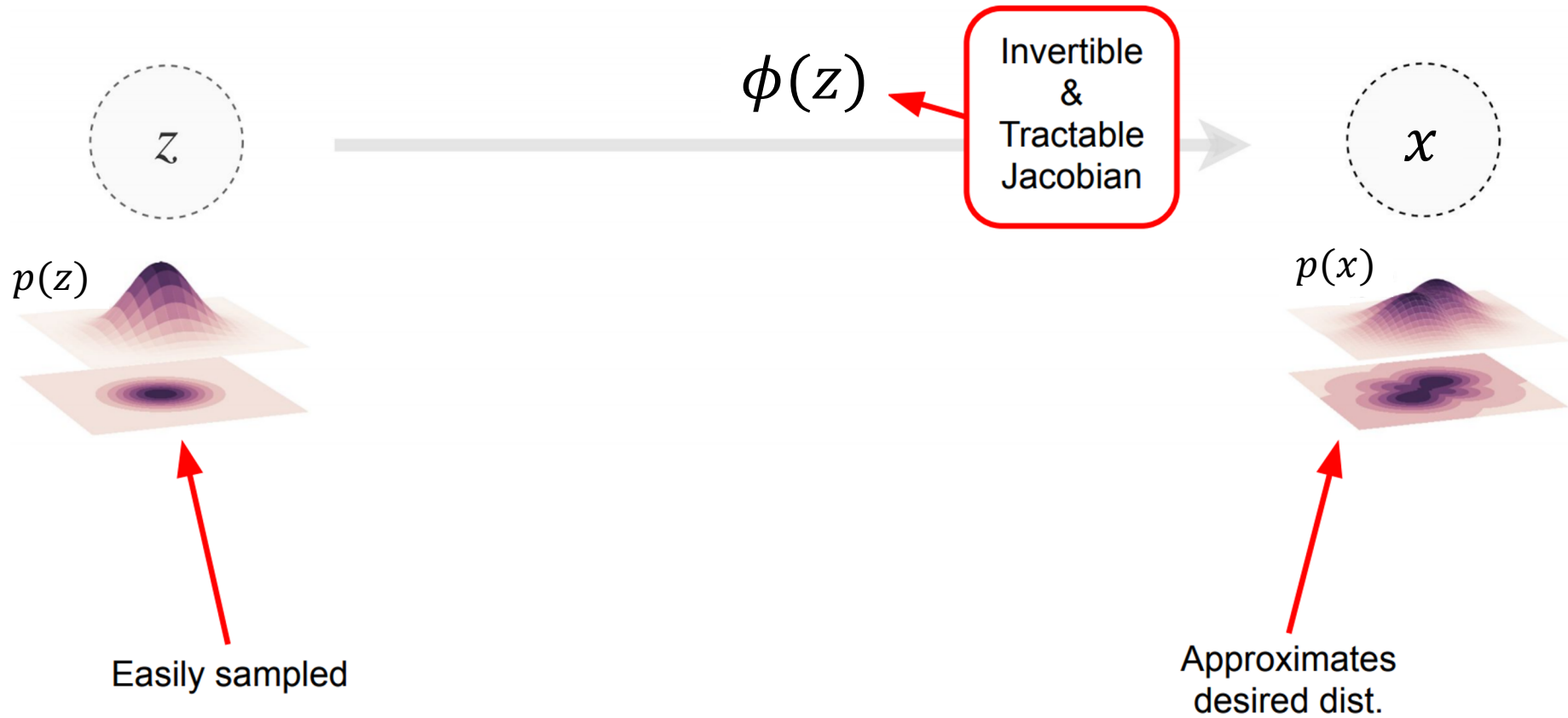
Inverse:  $(f_2 \circ f_1)^{-1} = f_1^{-1} \circ f_2^{-1}$

Jacobian:  $\det(A \cdot B) = \det(A) \cdot \det(B)$

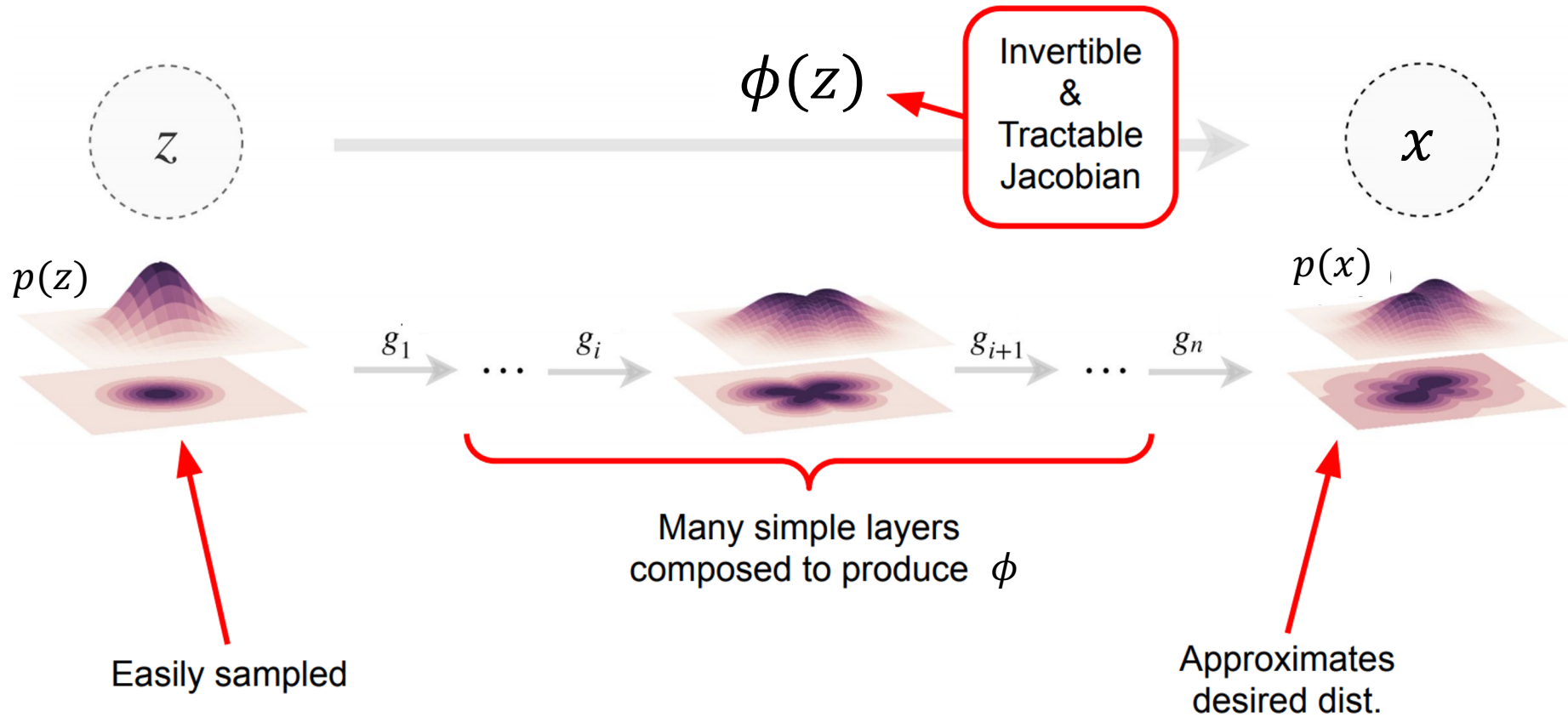
$$p_x(\mathbf{x}) = p_z(\mathbf{z}) \left| \det \left( \frac{\partial \phi(\mathbf{z})}{d\mathbf{z}} \right)^{-1} \right|$$

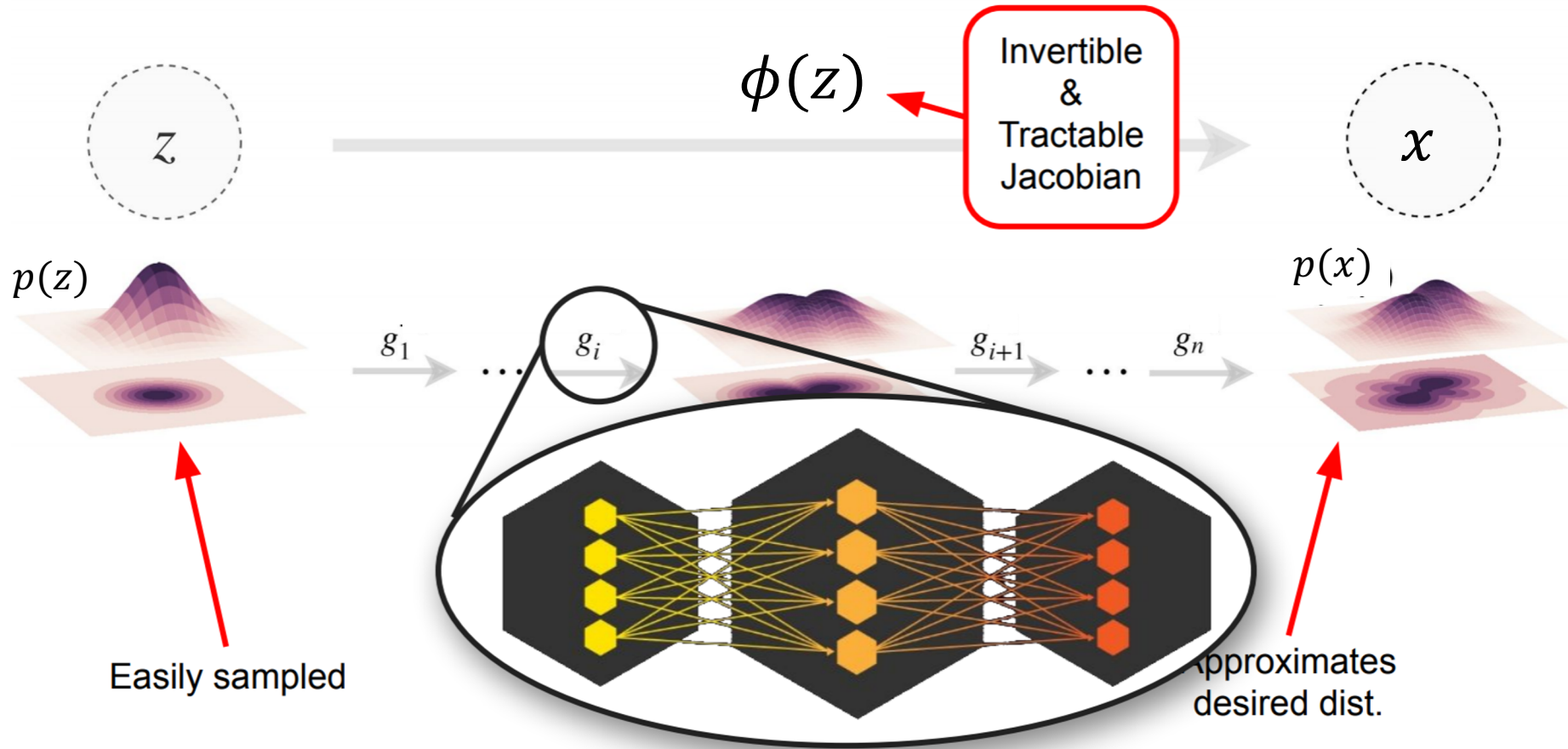








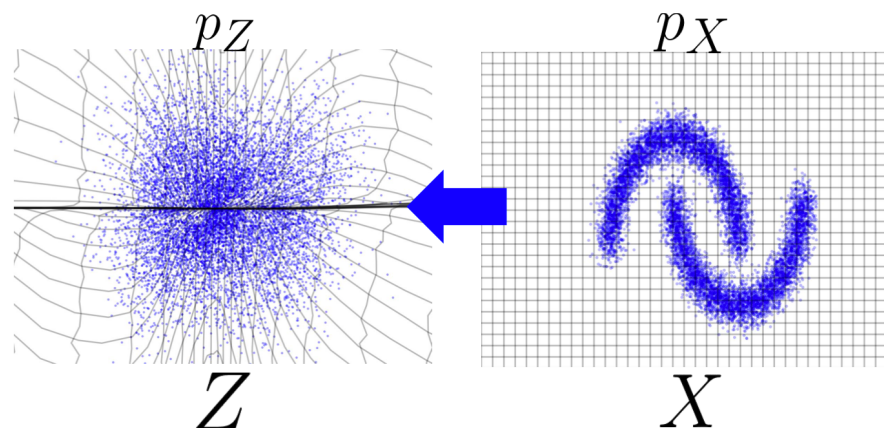




- Learn  $\theta$  with maximum likelihood

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

Where  $z = \phi^{-1}(x)$

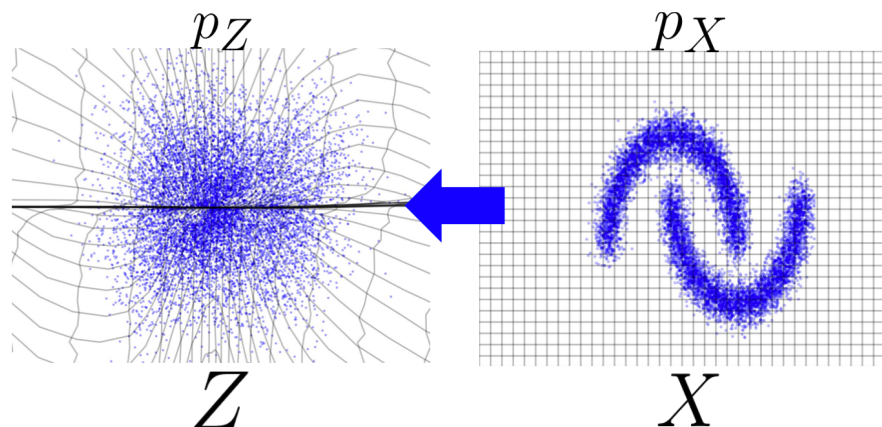


- Learn  $\theta$  with maximum likelihood

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

For each data point  $x$

- Map back to point in  $z$ -space with  $\phi^{-1}(x)$
- Evaluate probability in  $z$ -space with  $p_z(\cdot)$



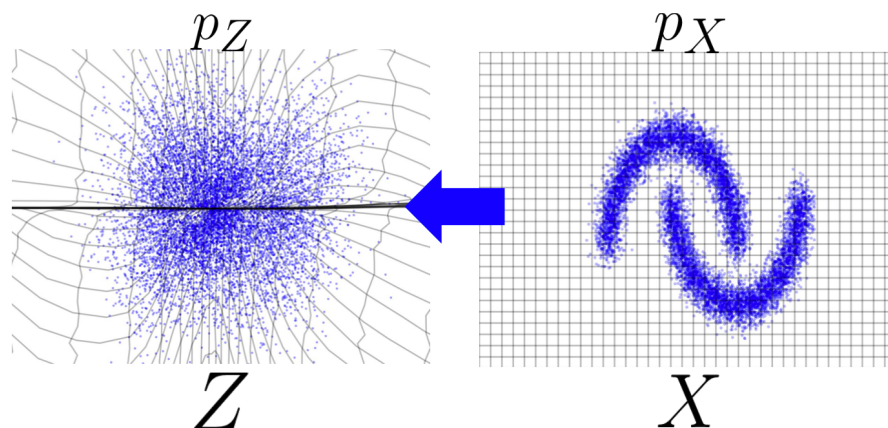
- Learn  $\theta$  with maximum likelihood

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

For each data point  $x$

- Map back to point in  $z$ -space with  $\phi^{-1}(x)$
- Evaluate probability in  $z$ -space with  $p_z(\cdot)$

Account for volume change  
due to transformation  $\phi^{-1}(x)$

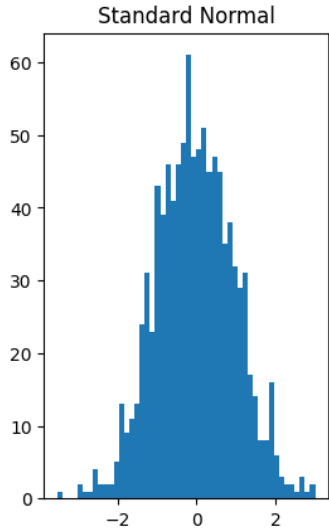
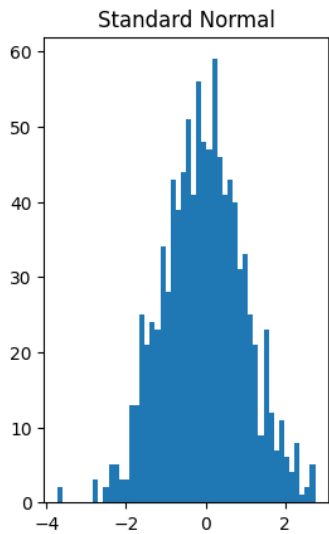
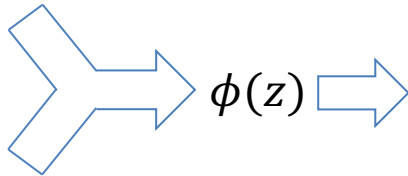
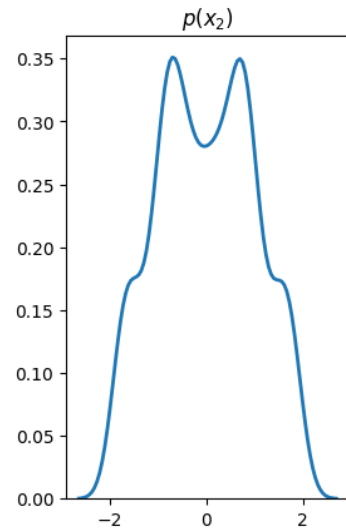
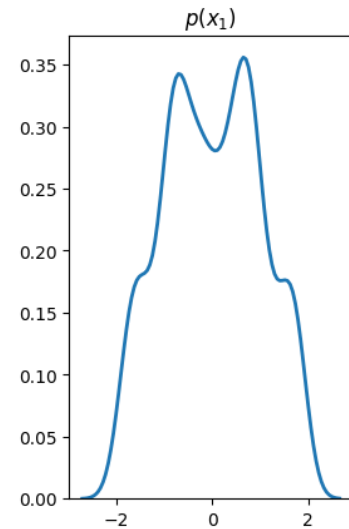
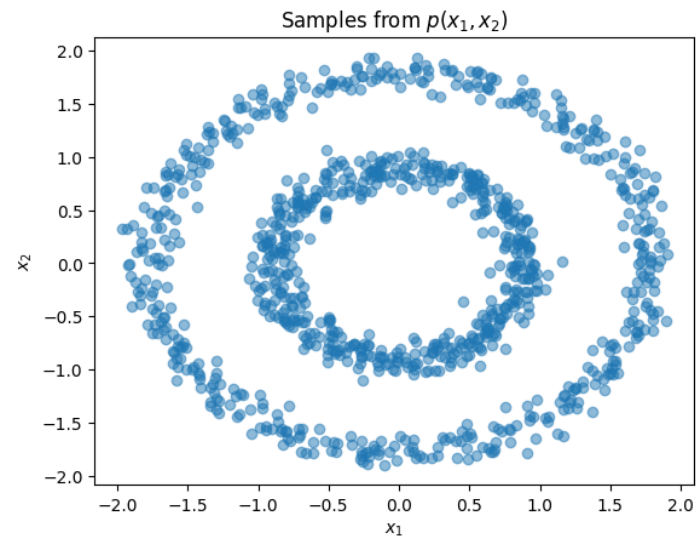


- **Learn  $\theta$**  with maximum likelihood

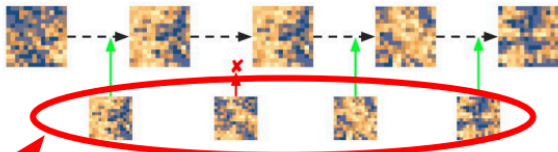
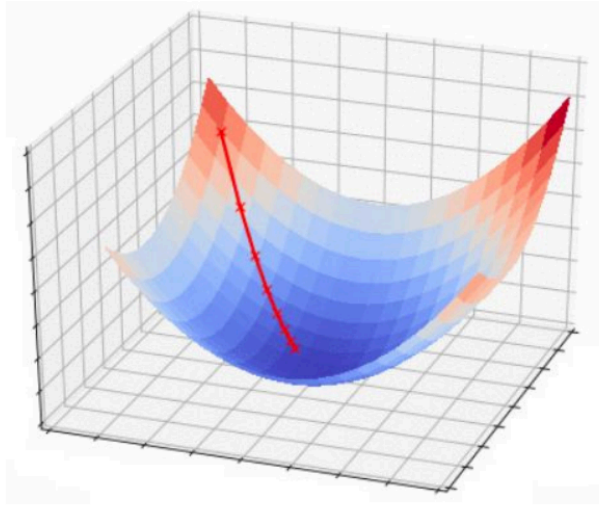
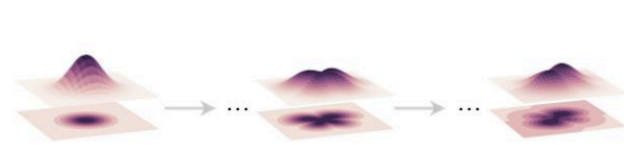
$$\max_{\theta} p(x) = \max_{\theta} p_z(\phi_{\theta}^{-1}(x)) \left| \det \left( \frac{\partial \phi_{\theta}^{-1}(x)}{dx} \right) \right|$$

- Gradient descent on  $\theta$
- 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:  $\phi(z) = x$

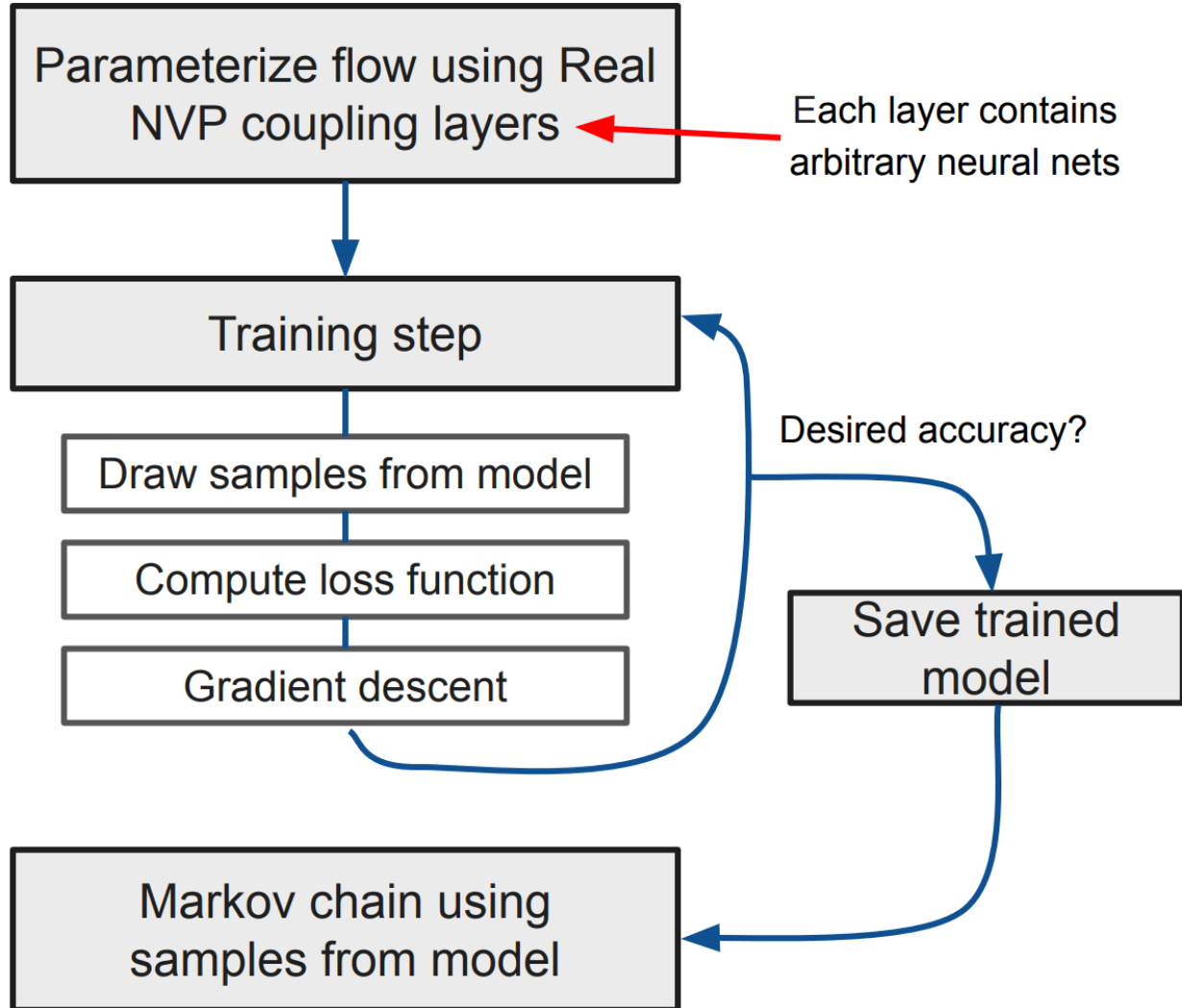
# Example Normalizing flow

 $z_1$  $z_2$ 

# Applications: Sampling in Lattice QCD

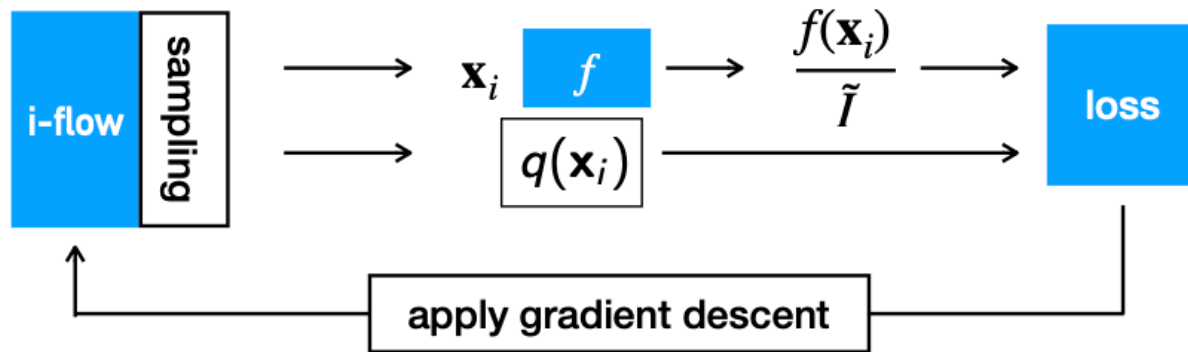


generating samples is  
"embarrassingly parallel"



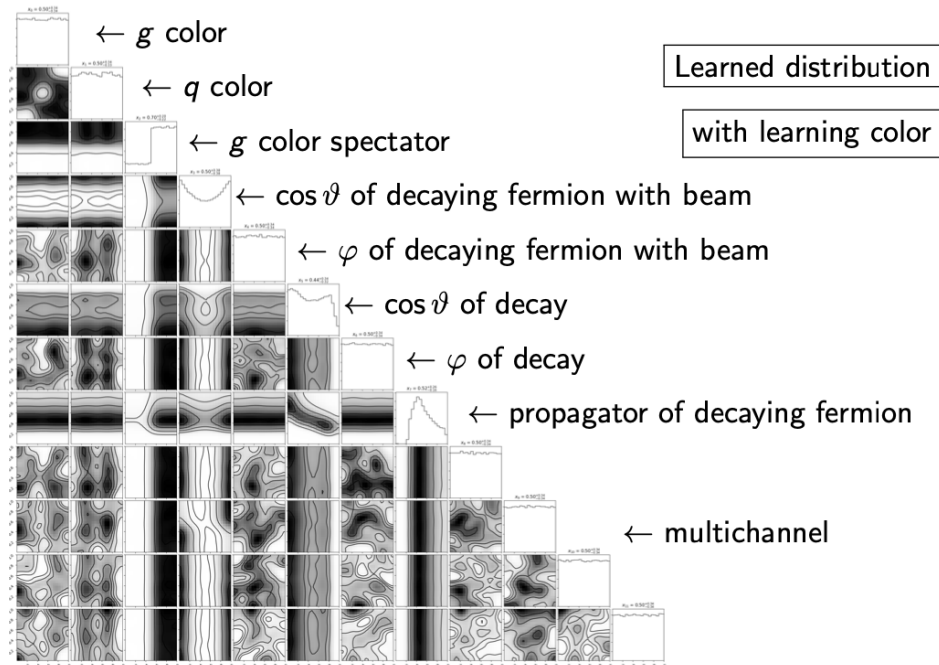
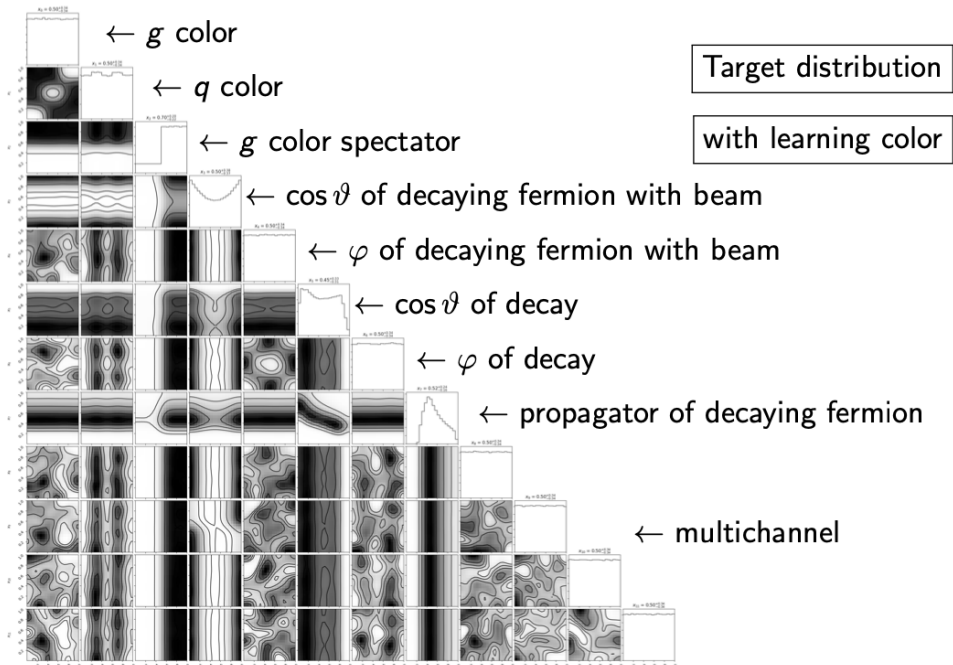


# Event Generation with Normalizing Flows

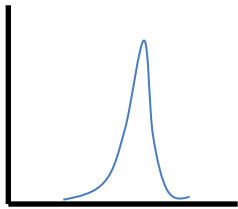


arXiv: 2001.05486, ML:ST  
arXiv: 2001.10028, PRD  
Slide credit: [C. Krause](#)

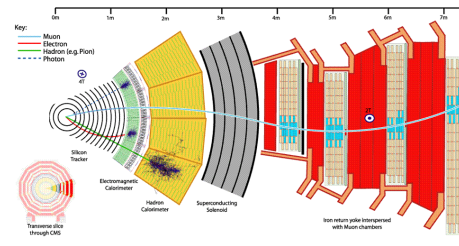
## Example: Learning $e^+e^- \rightarrow 3j$



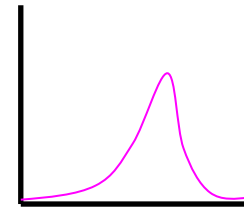
# Unfolding with Normalizing Flows



$x \sim p(x)$  = input / true distribution



$p(y|x)$  = Detector smearing

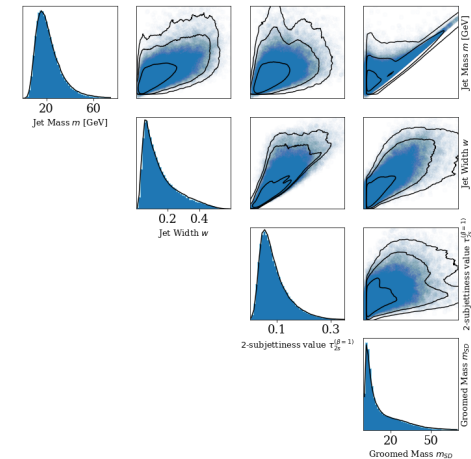
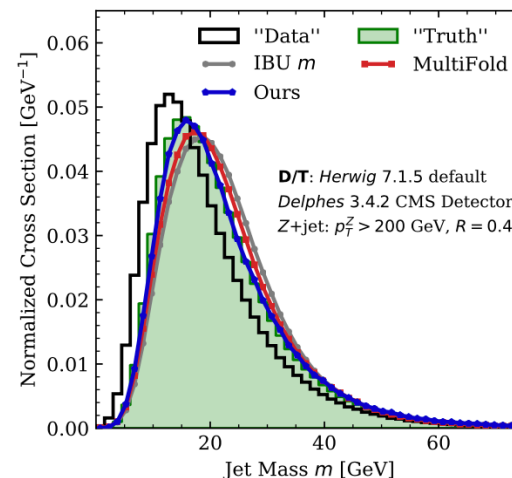
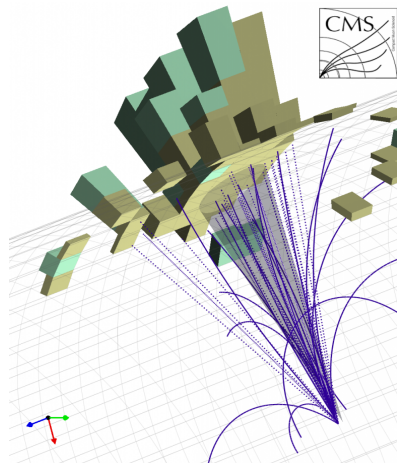


$y \sim p(y)$  = output / observed distribution

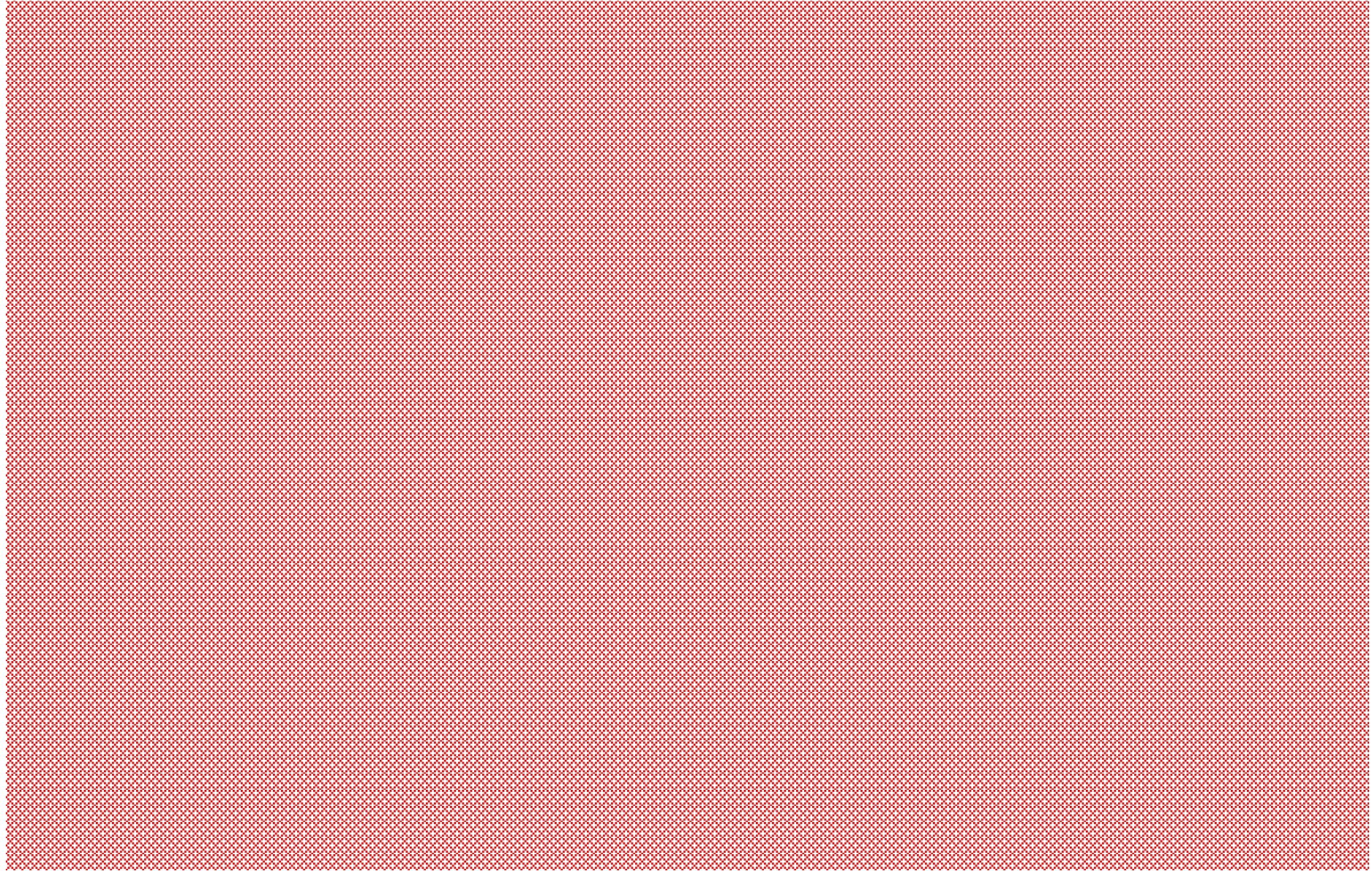
$$\text{Observed distribution: } p(y) = \int p(y|x)p(x)dx \approx \sum_{x \sim p(x)} p(y|x)$$

- Normalizing flows to model detector  $p(y|x)$  trained with simulation
- Normalizing flow to model unknown truth  $p_\theta(x)$
- Maximize data likelihood  $p(y) \rightarrow$  Gradient descent to learn parameters  $\theta$

Unfolding  
Jet variables  
in Z+jet events



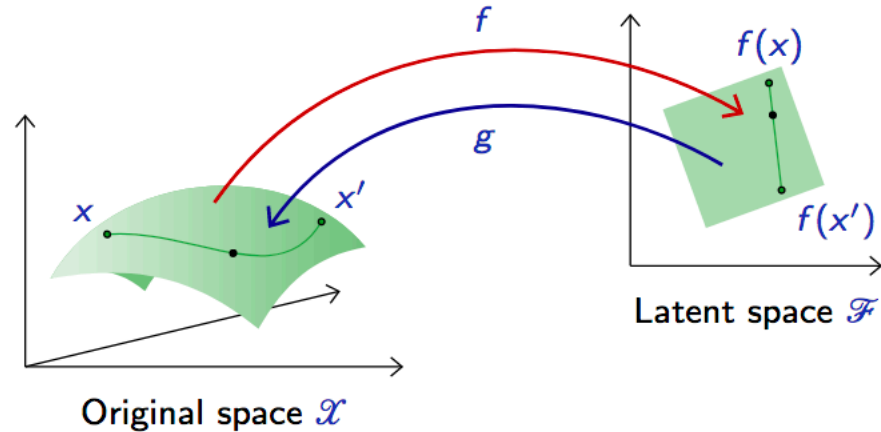
- Deep neural networks are an extremely powerful class of models
- We can express our inductive bias about a system in terms of model design, and can be adapted to a many types of data
- Even beyond classification and regression, deep neural networks allow for powerful model schemes such as Generative adversarial Networks and normalizing flows that open many new possible tasks where Machine Learning can be applied in HEP



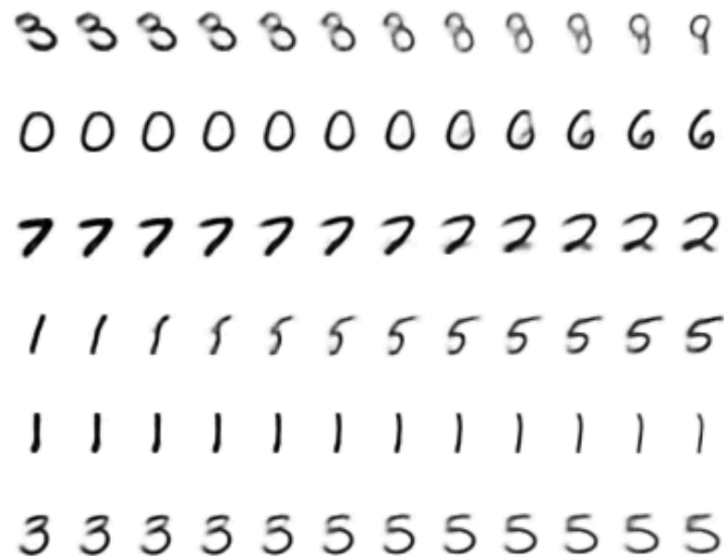
- Autoencoders learn the latent space, but we don't know what is the latent space distribution
- Autoencoder prescribes a deterministic relationship between data space and latent space
- One set of “meaningful degrees of freedom” can only describe one data space point

# Interpolating in Latent Space

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



Autoencoder interpolation ( $d = 8$ )





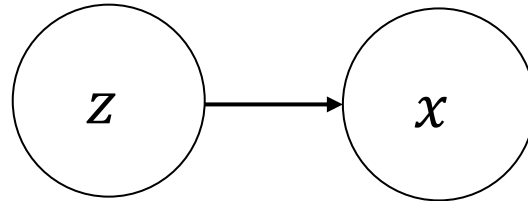
- For  $z \sim p_\theta(z)$ , rewrite  $z$  as a function of a random variable  $\epsilon$  whose distributions  $p(\epsilon)$  does not depend on  $\theta$ 
  - Gaussian Example:

$$z \sim \mathcal{N}(\mu, \sigma) \rightarrow z = \sigma * \epsilon + \mu \quad \text{where } \epsilon \sim \mathcal{N}(0,1)$$

- VAE Loss

$$\max_{\theta, \psi} L(\theta, \psi) = \max_{\theta, \psi} \sum_{\epsilon \sim p(\epsilon)} \log p_\theta(x | z_i = \epsilon * \sigma_\psi(x) + \mu_\psi(x)) - \log \left[ \frac{q_\psi(z_i | x)}{p(z_i)} \right]$$





- Observed random variable  $x$  depends on unobserved latent random variable  $z$ 
  - Interpret  $z$  as the causal factors for  $x$
- Joint probability:  $p(x, z) = p(x|z)p(z)$
- $p(x|z)$  is a stochastic generation process from  $z \rightarrow x$
- Inference from posterior: 
$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$
  - Usually can't compute marginal  $p(x) = \int p(x|z)p(z)dz$

- Consider probabilistic relationship between data and latent variables

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

Decoding data  $x$   
from latent  $z$

Prior over latent space

- Consider probabilistic relationship between data and latent variables

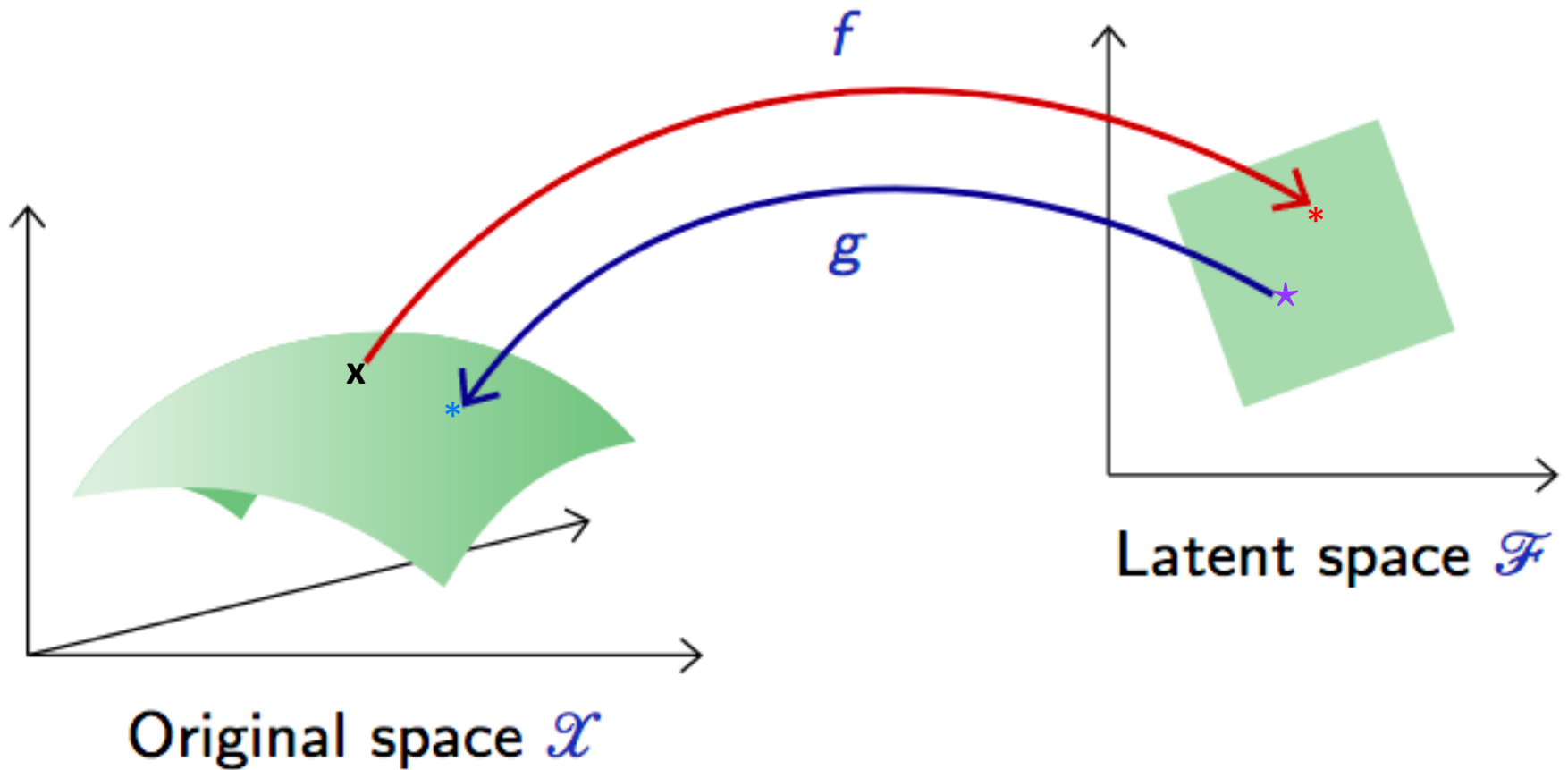
$$x, z \sim p(x, z) = p(x|z)p(z)$$

- Autoencoding

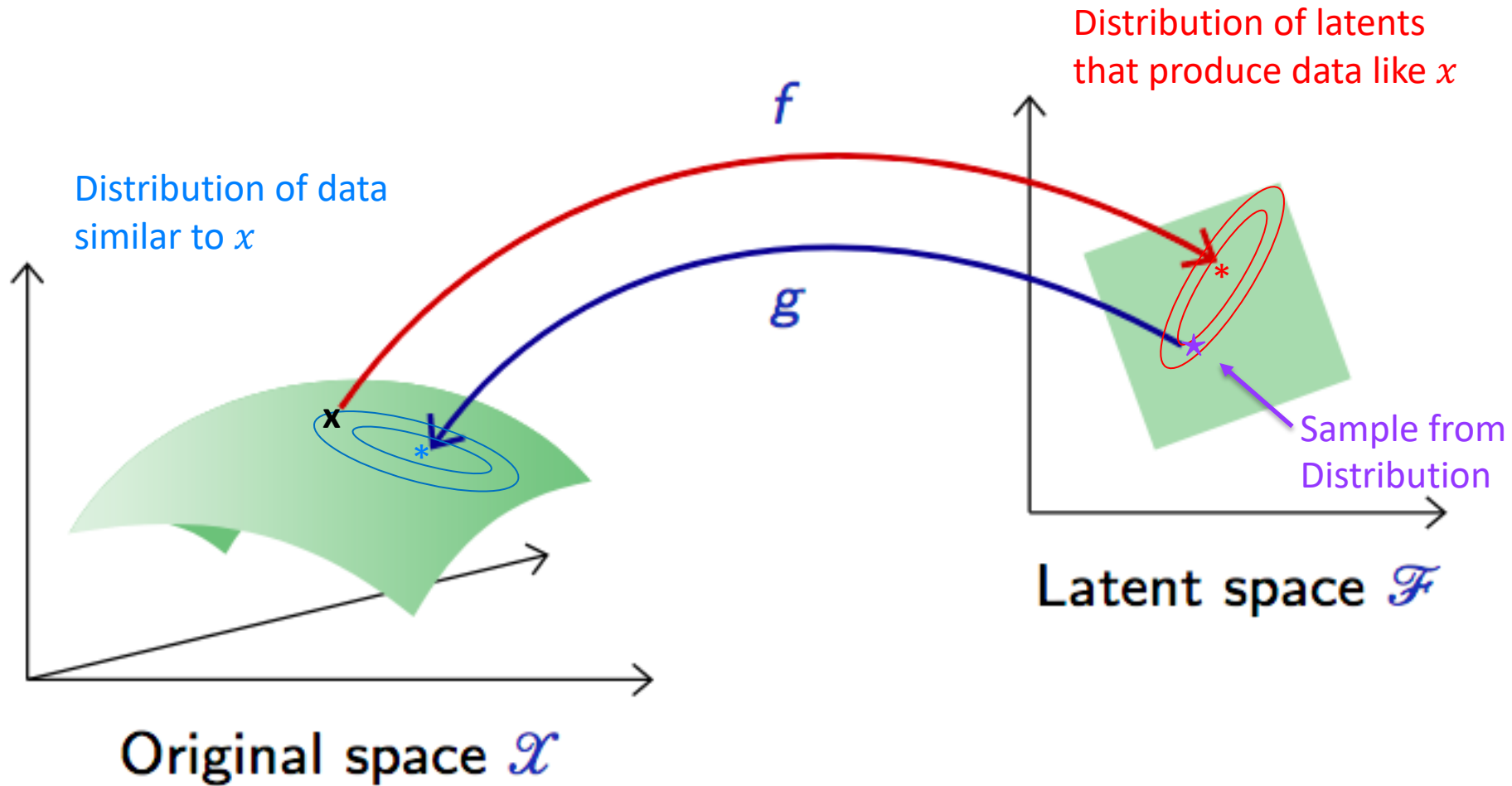
$$x \rightarrow q(z|x) \xrightarrow{\text{sample}} z \rightarrow p(x|z)$$

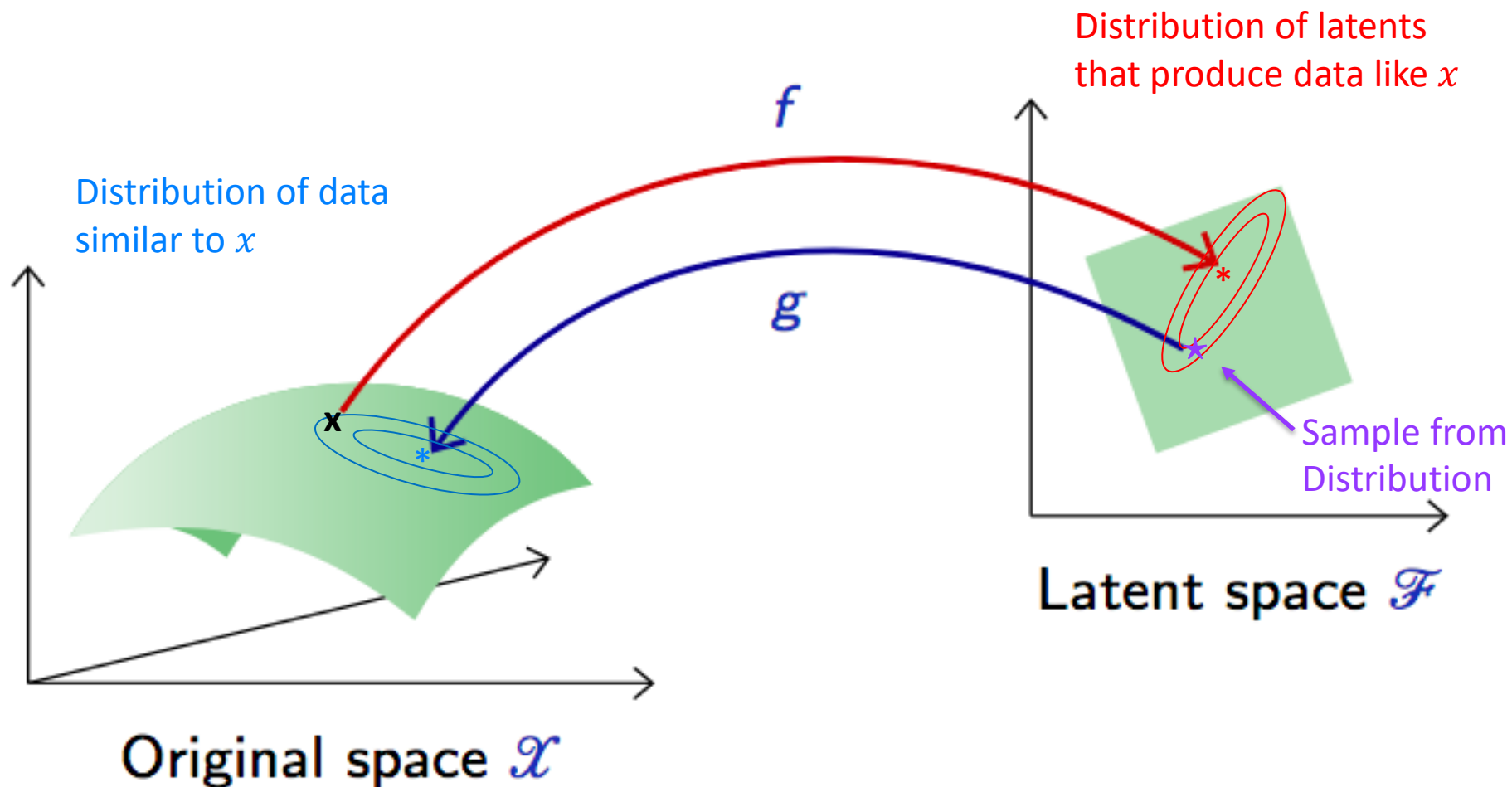
- Choose simple prior distribution
- **Encoder:** Learn what latents can produced data:  $q(z|x)$
- **Decoder:** Learn what data is produced by latent:  $p(x|z)$

# Autoencoder $\rightarrow$ Variational Autoencoder (VAE)



# Autoencoder $\rightarrow$ Variational Autoencoder (VAE)





- **Encode** data  $x$  into distribution over latents  $q(z|x)$
- For any sample  $z$  **decode** into distribution over data  $p(x|z)$

- PDF often depends on parameters  $\theta$  we are interested in
  - Write the density as  $f(x|\theta)$  or  $f(x; \theta)$
- Choose a family we know: Gaussian
- Model the parameters  $\theta$  as output of Neural Net

$$\mu(x) \equiv NN(x)$$

$$\sigma(x) \equiv NN(x)$$

$$p(z|x) = \mathcal{N}(z; \mu(x), \sigma(x))$$

Encoding Distribution

$$\log p(z|x) = \frac{(z - \mu(x))^2}{2\sigma^2(x)} - \frac{1}{2} \log \sigma(x) - \log \sqrt{\pi}$$

- PDF often depends on parameters  $\theta$  we are interested in
  - Write the density as  $f(x|\theta)$  or  $f(x; \theta)$
- Choose a family we know: Gaussian
- Model the parameters  $\theta$  as output of Neural Net

$$\mu(z) \equiv NN(z)$$

$$\sigma(z) \equiv NN(z)$$

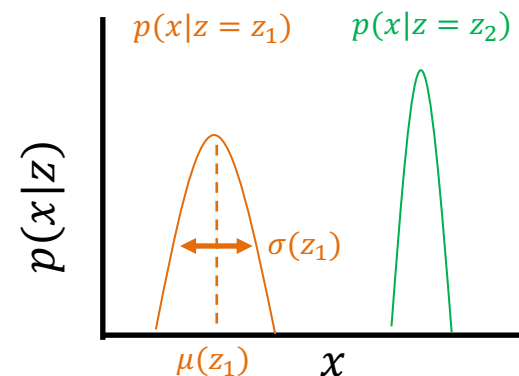
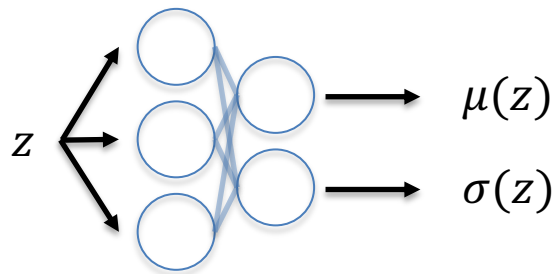
$$p(x|z) = \mathcal{N}(x; \mu(z), \sigma(z))$$

Decoding Distribution

$$\log p(x|z) = \frac{(x - \mu(z))^2}{2\sigma^2(z)} - \frac{1}{2} \log \sigma(z) - \log \sqrt{\pi}$$



- PDF often depends on parameters  $\theta$  we are interested in
  - Write the density as  $f(x|\theta)$  or  $f(x; \theta)$
- Choose a family we know: Gaussian
- Model the parameters  $\theta$  as output of Neural Net



$$\max_{\theta, \psi} L(\theta, \psi)$$

$$L(\theta, \psi) = \sum_{z \sim q_{\psi}(Z|x)} \log p_{\theta}(x|z) - \sum_{z \sim q_{\psi}(Z|x)} \log \frac{q_{\psi}(z|x)}{\mathcal{N}(z; 0, 1)}$$

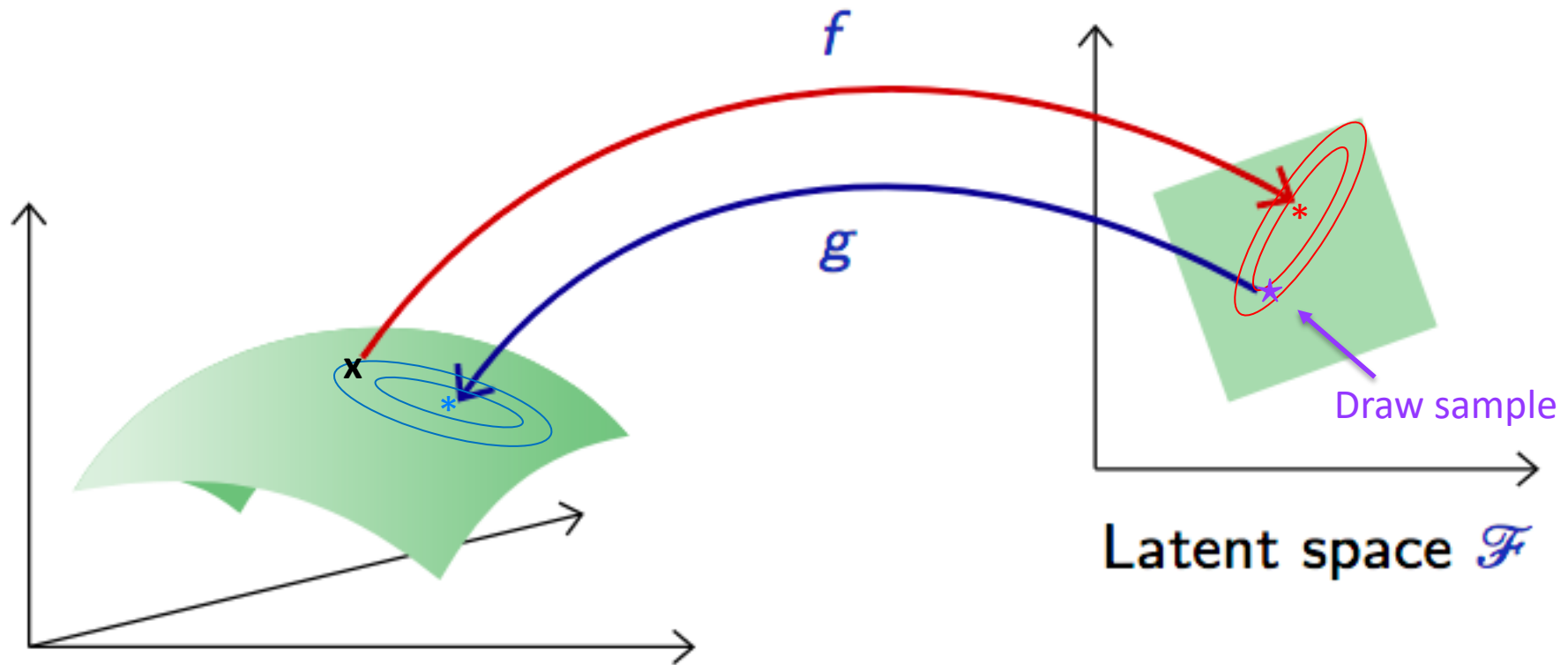
- First Term
  - Check compatibility with original data after encoding and decoding
- Second Term
  - Check compatibility of encoded data with prior
  - Constraint on latent distribution

$$\max_{\theta, \psi} L(\theta, \psi)$$

$$L(\theta, \psi) = \sum_{z \sim q_{\psi}(Z|x)} \log p_{\theta}(x|z) - \sum_{z \sim q_{\psi}(Z|x)} \log \frac{q_{\psi}(z|x)}{\mathcal{N}(z; 0, 1)}$$

$$= \frac{1}{2} \sum_{z \sim q_{\psi}(Z|x)} \frac{(x - \mu(z))^2}{\sigma^2(z)} - \log \sigma(z)$$

$$- \frac{1}{2} \sum_{z \sim q_{\psi}(Z|x)} \sigma(x) - \mu^2(x) - 1 - \log \sigma(x)$$

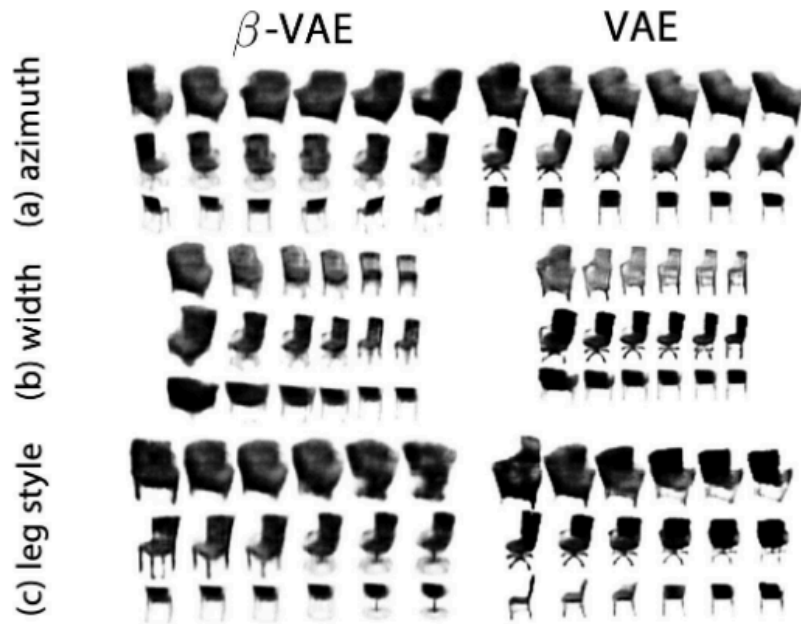


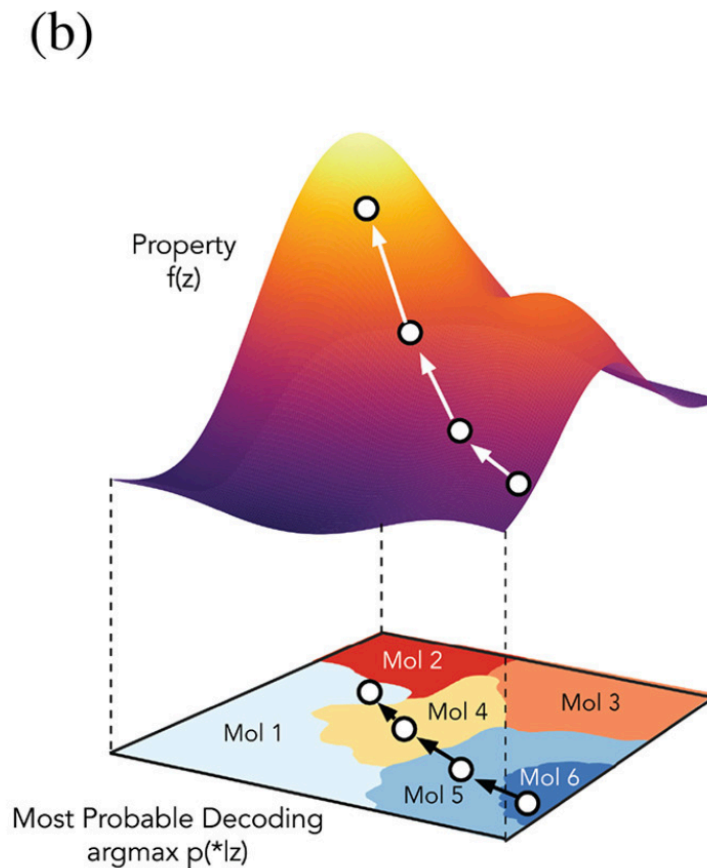
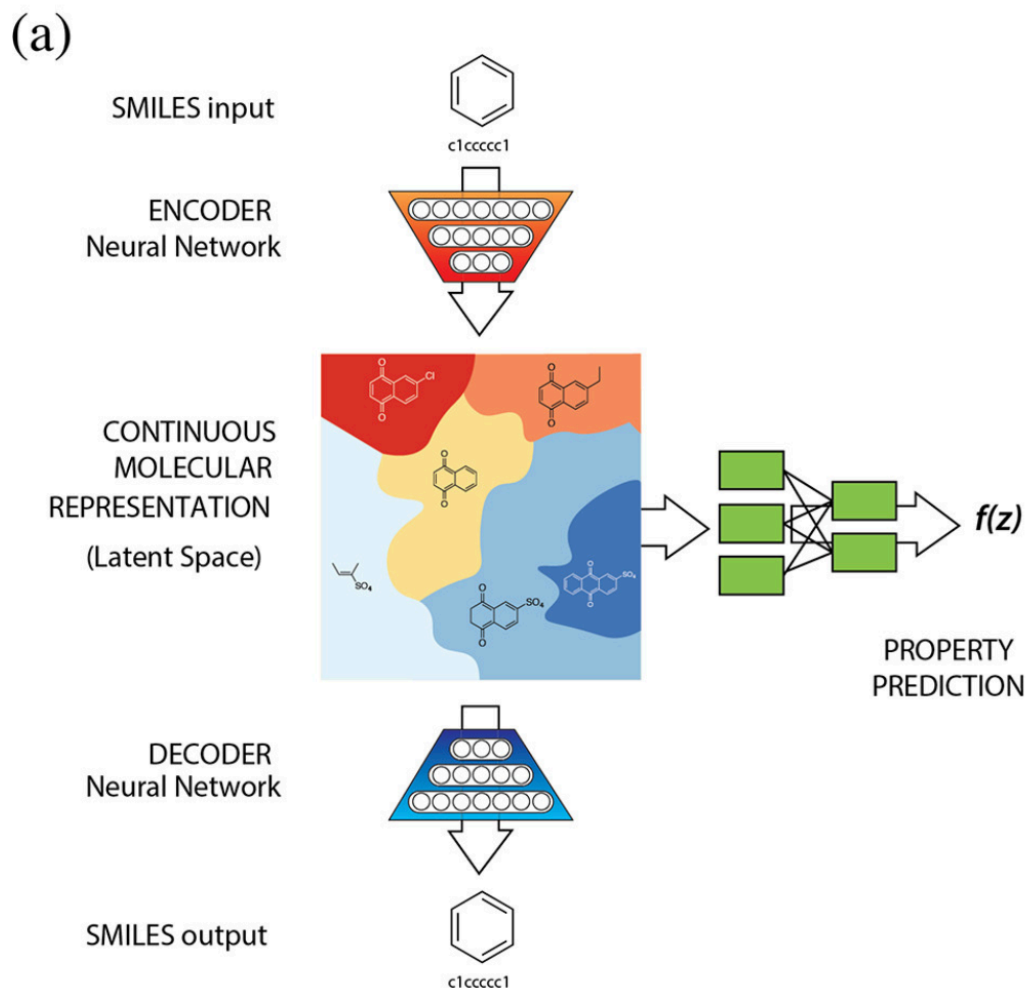
Original space  $\mathcal{X}$

Latent space  $\mathcal{F}$

$$\sum_{z \sim q_{\psi}(z|x)} \log p_{\theta}(x|z)$$

$$\sum_{z \sim q_{\psi}(z|x)} \log \frac{q_{\psi}(z|x)}{\mathcal{N}(z; 0, 1)}$$

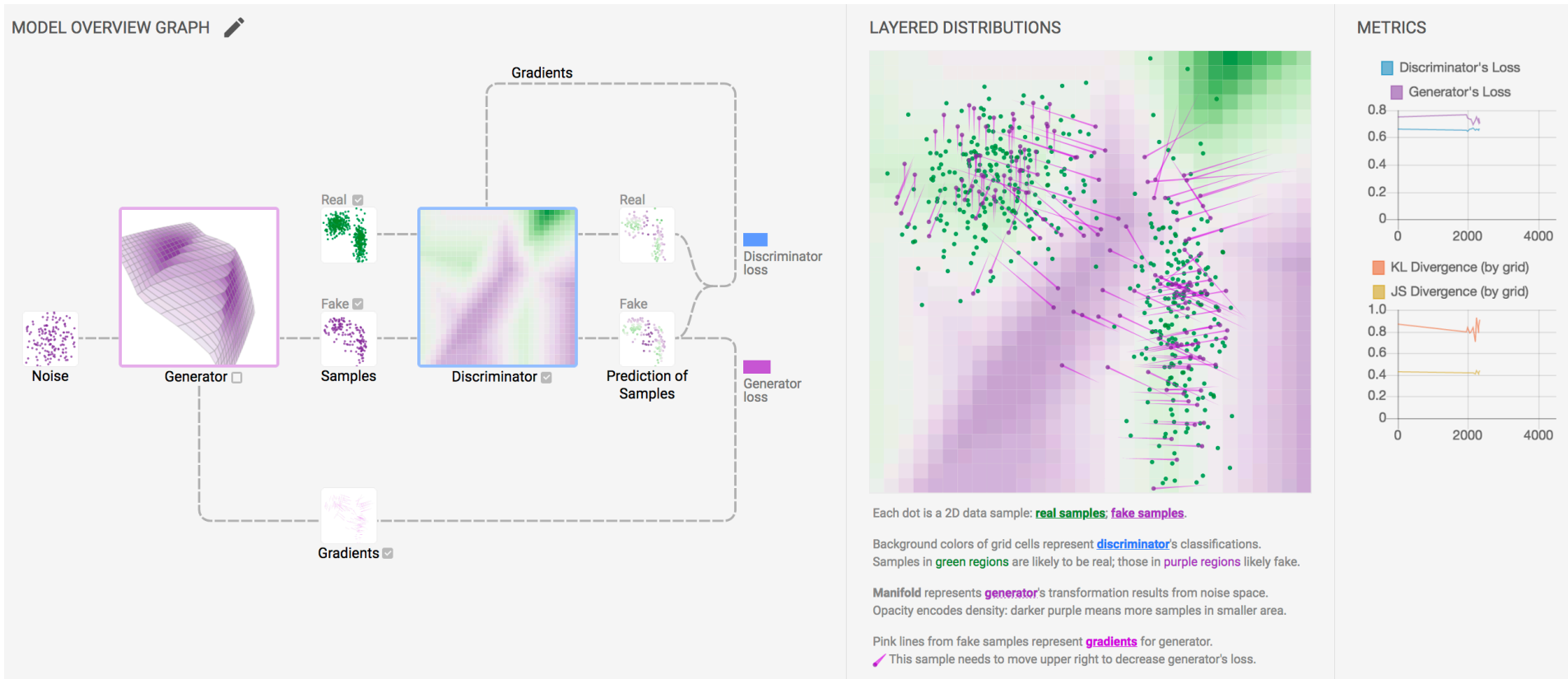




Design of new molecules with desired chemical properties.  
 (Gomez-Bombarelli et al, 2016)



# GAN Training Example



[GAN Lab Demo](#)



- 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!

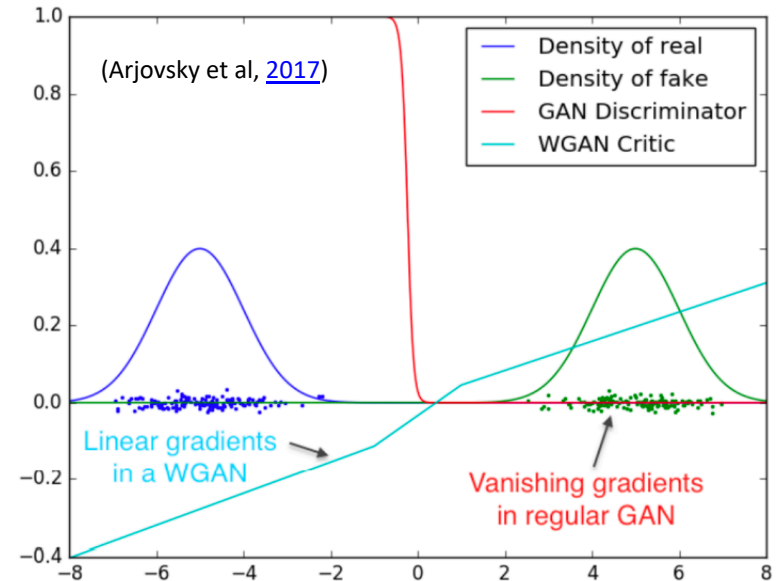
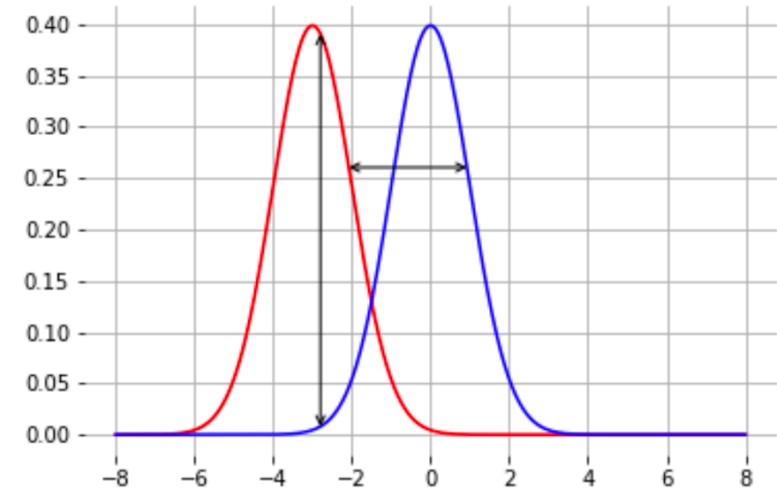
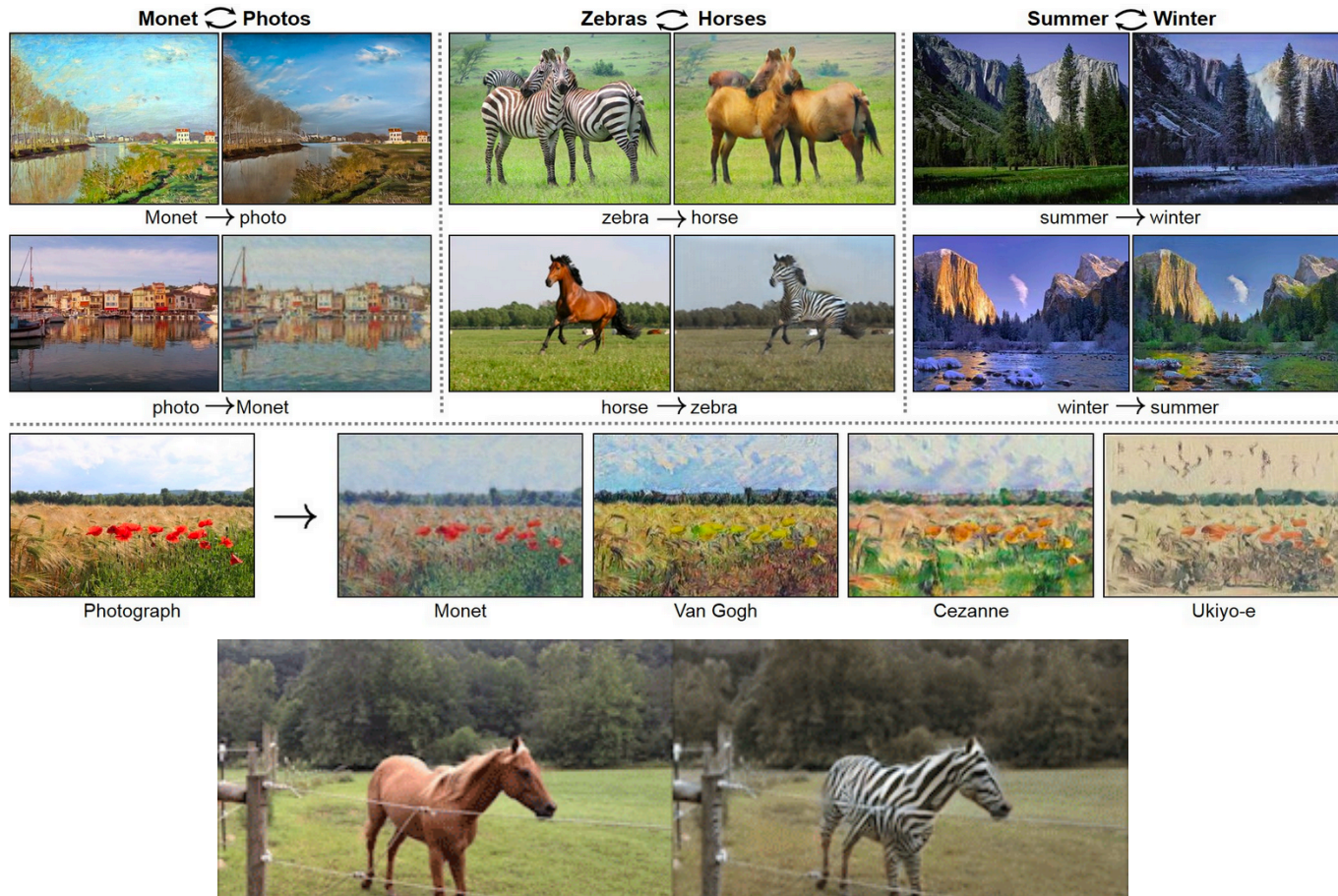


Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.



(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  $A$





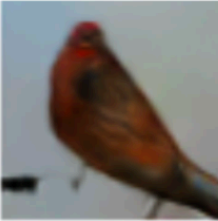

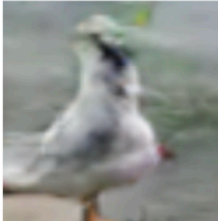
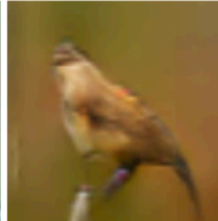
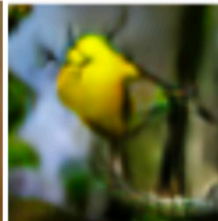
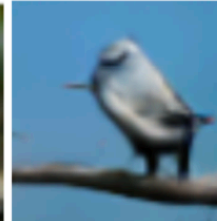


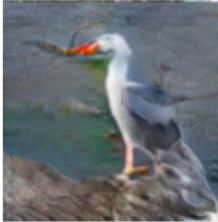


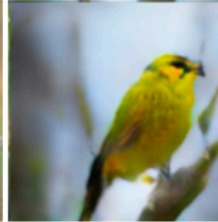



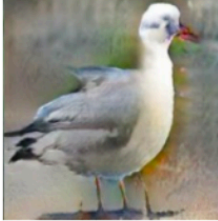
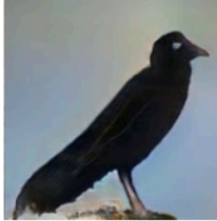
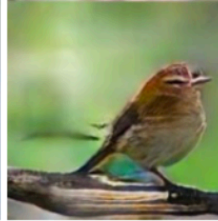

Text description	This bird is red and brown in color, with a stubby beak	The bird is short and stubby with yellow on its body	A bird with a medium orange bill white body gray wings and webbed feet	This small black bird has a short, slightly curved bill and long legs	A small bird with varying shades of 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)

