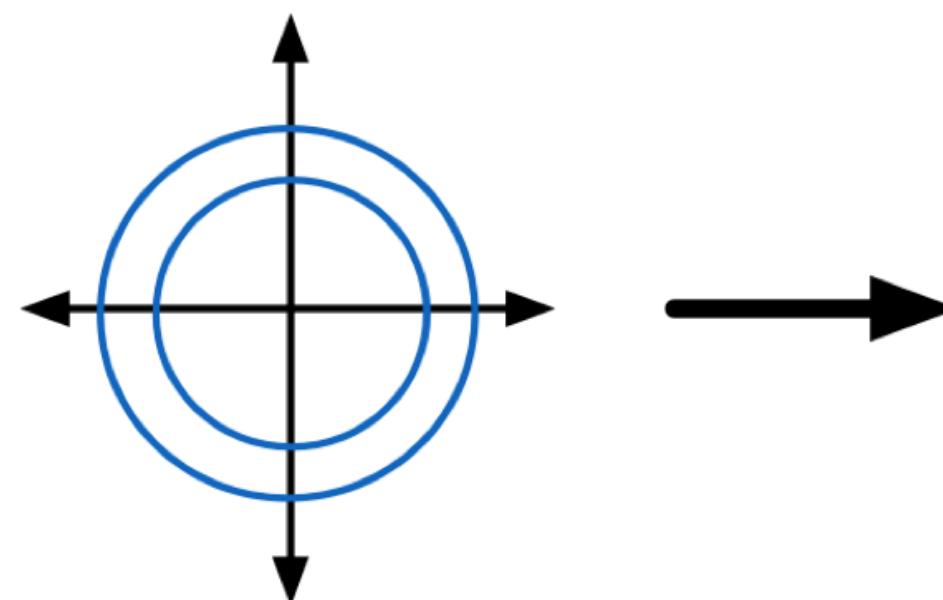


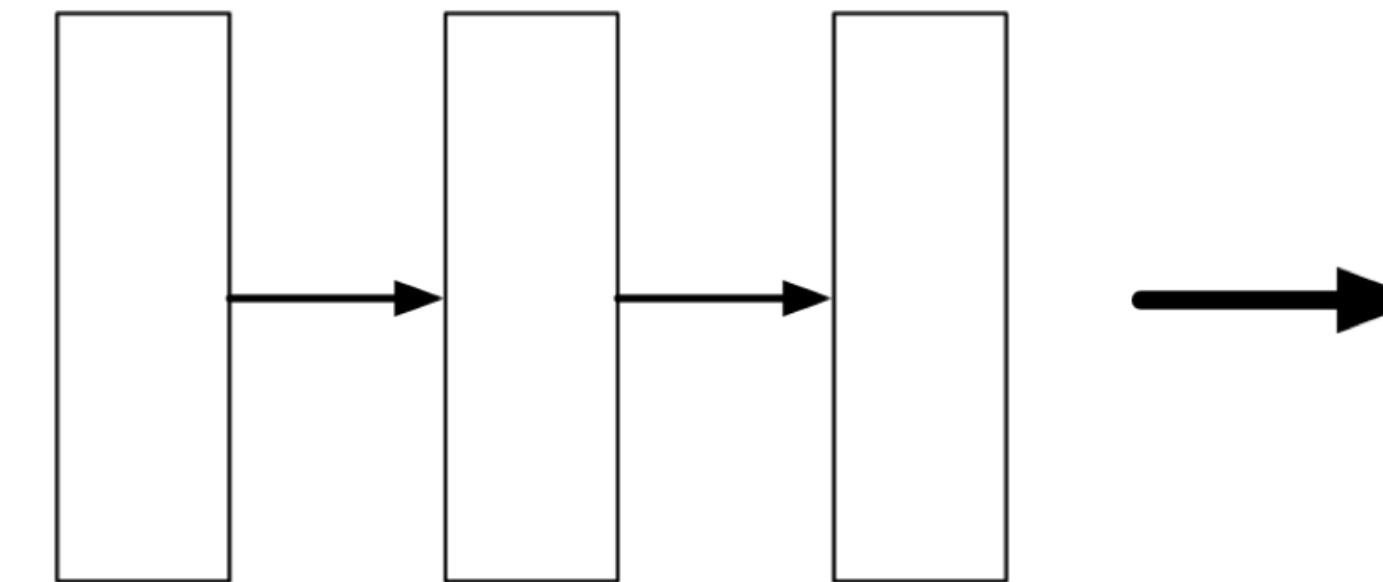
**GAN**

# Implicit Generative Models

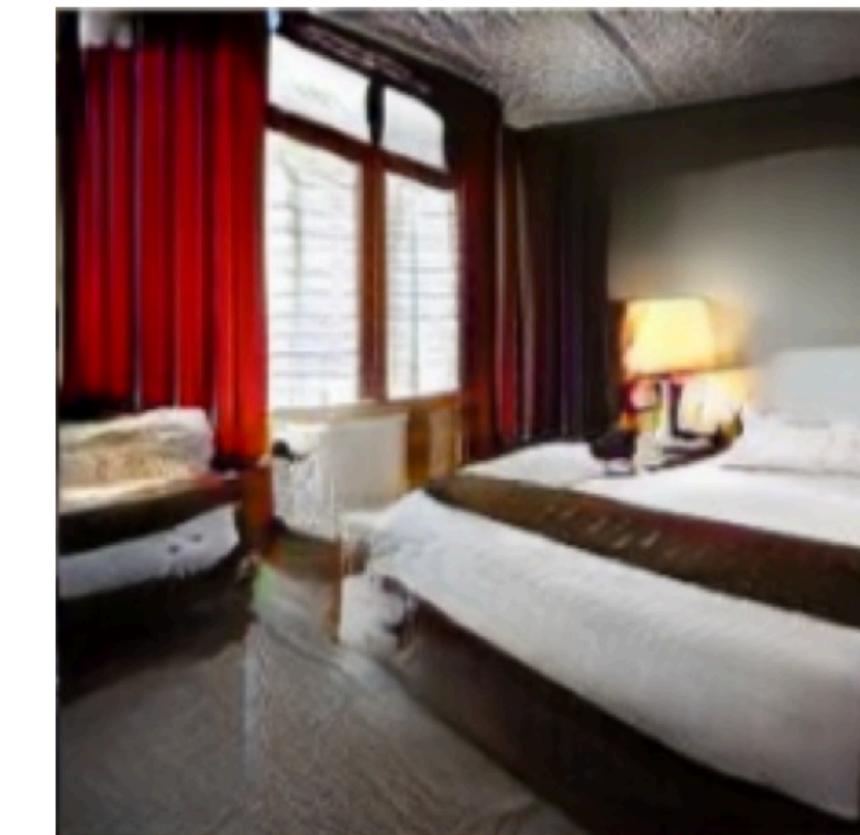
Recall: implicit generative models learn a mapping from random noise vectors to things that look like, e.g., images:



Each dimension of the code vector is sampled independently from a simple distribution, e.g. Gaussian or uniform.



This is fed to a (deterministic) generator network.

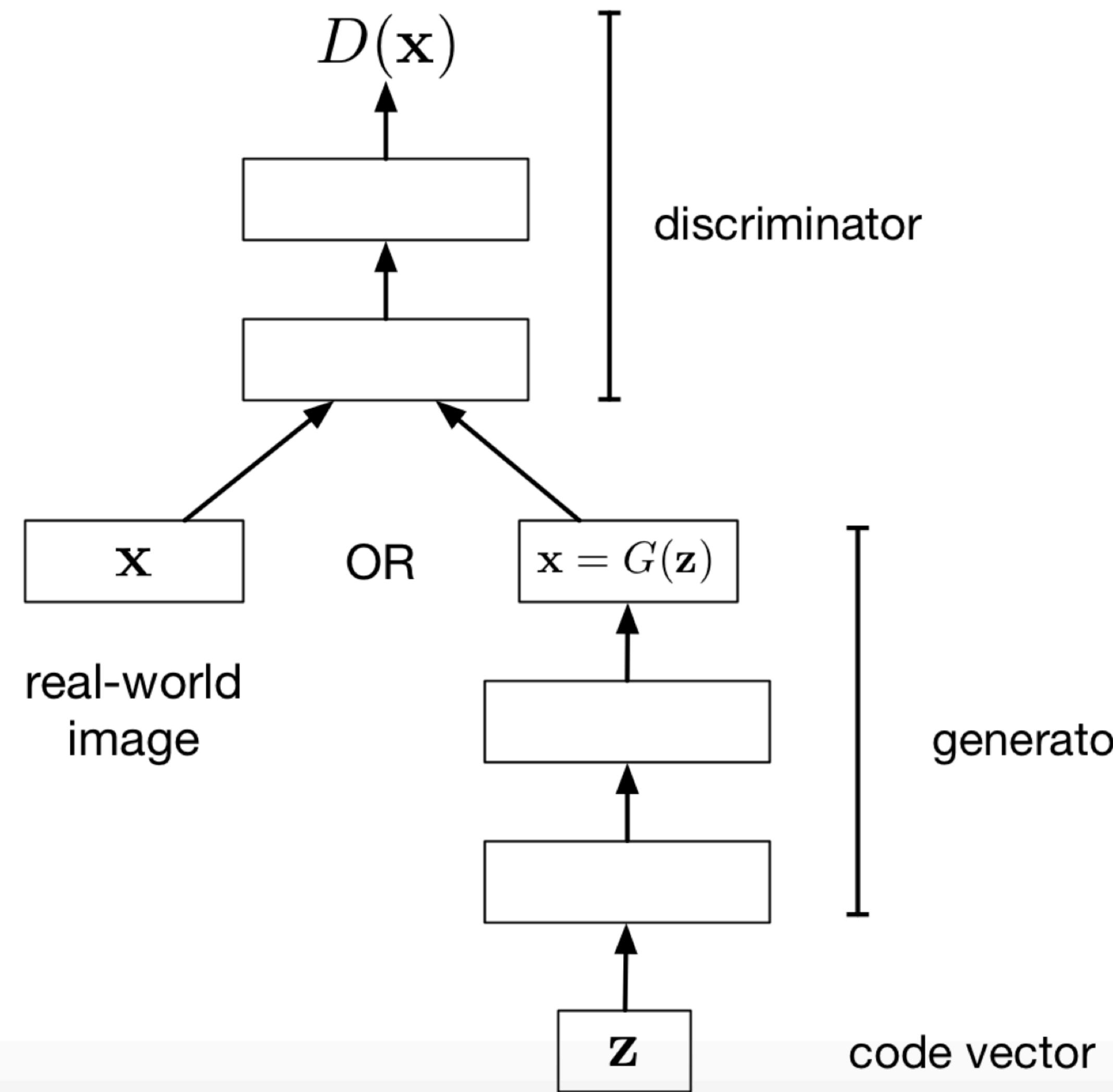


The network outputs an image.

# Generative Adversarial Networks

- The advantage of implicit generative models: if you have some criterion for evaluating the quality of samples, then you can compute its gradient with respect to the network parameters, and update the network's parameters to make the sample a little better
- The idea behind **Generative Adversarial Networks (GANs)**: train two different networks
  - The **generator network** tries to produce realistic-looking samples
  - The **discriminator network** tries to figure out whether an image came from the training set or the generator network
- The generator network tries to fool the discriminator network

# Generative Adversarial Networks



## Generative Adversarial Networks

- Let  $D$  denote the discriminator's predicted probability of being data
- Discriminator's cost function: cross-entropy loss for task of classifying real vs. fake images

$$\mathcal{J}_D = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[-\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z}}[-\log(1 - D(G(\mathbf{z})))]$$

- One possible cost function for the generator: the opposite of the discriminator's

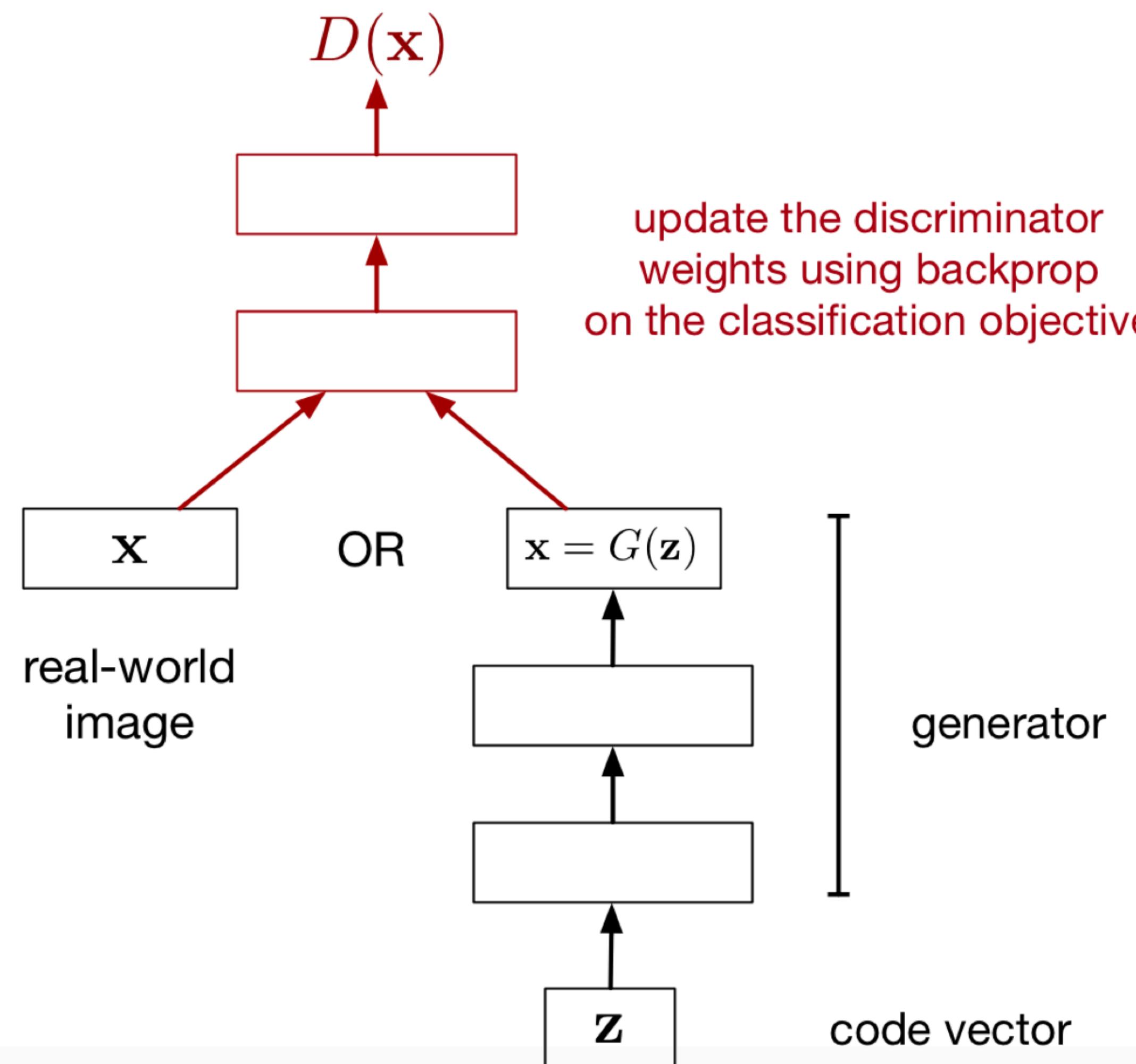
$$\begin{aligned}\mathcal{J}_G &= -\mathcal{J}_D \\ &= \text{const} + \mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]\end{aligned}$$

- This is called the **minimax formulation**, since the generator and discriminator are playing a **zero-sum game** against each other:

$$\max_G \min_D \mathcal{J}_D$$

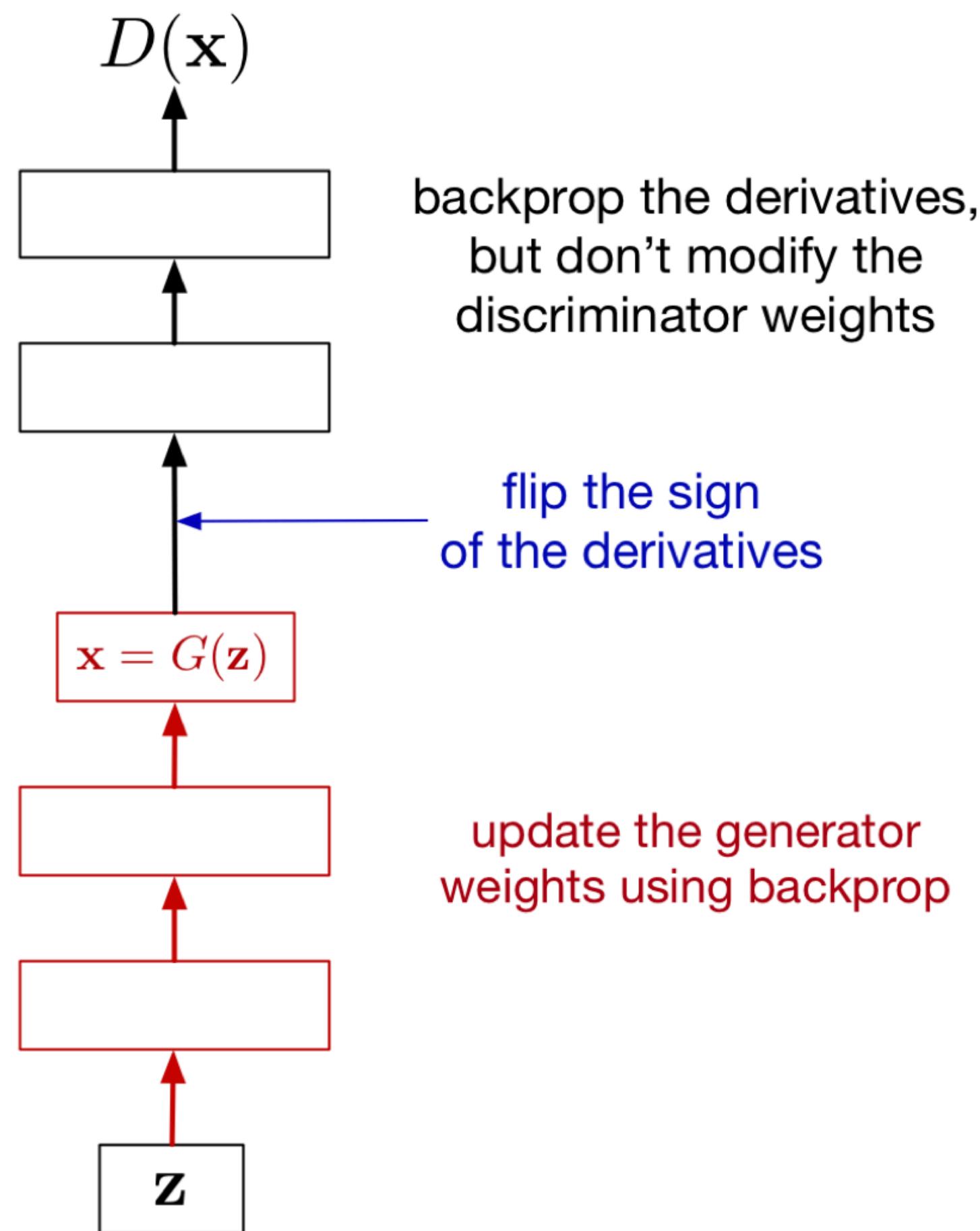
# Generative Adversarial Networks

# Updating the discriminator:



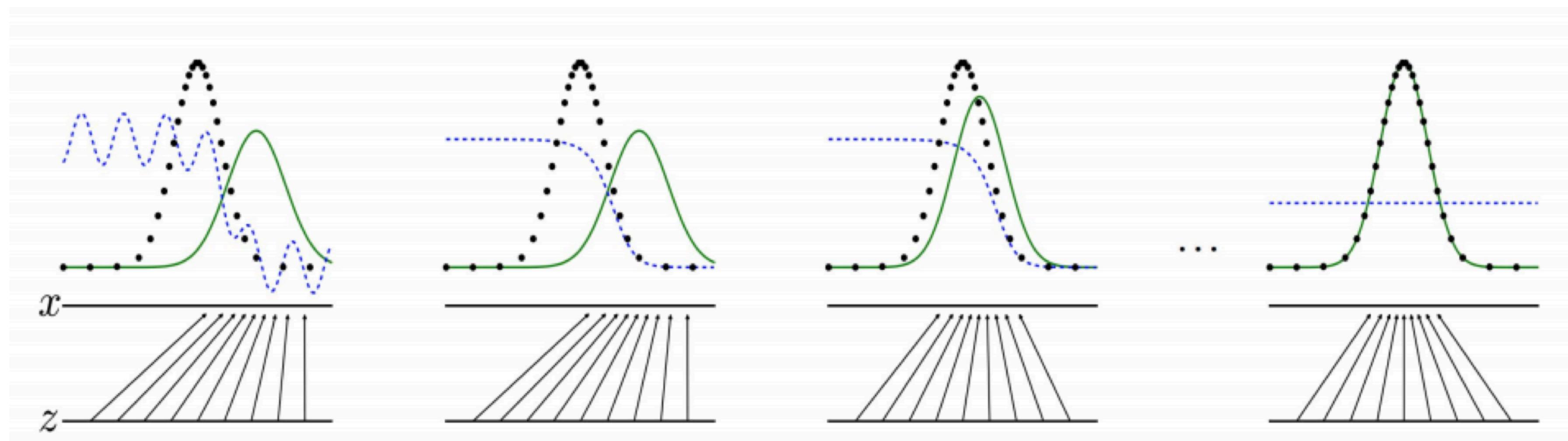
# Generative Adversarial Networks

Updating the generator:



# Generative Adversarial Networks

Alternating training of the generator and discriminator:



## A Better Cost Function

- We introduced the minimax cost function for the generator:

$$\mathcal{J}_G = \mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$$

- One problem with this is **saturation**.
- Recall from our lecture on classification: when the prediction is really wrong,
  - “Logistic + squared error” gets a weak gradient signal
  - “Logistic + cross-entropy” gets a strong gradient signal
- Here, if the generated sample is really bad, the discriminator’s prediction is close to 0, and the generator’s cost is flat.

# A Better Cost Function

- Original minimax cost:

$$\mathcal{J}_G = \mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$$

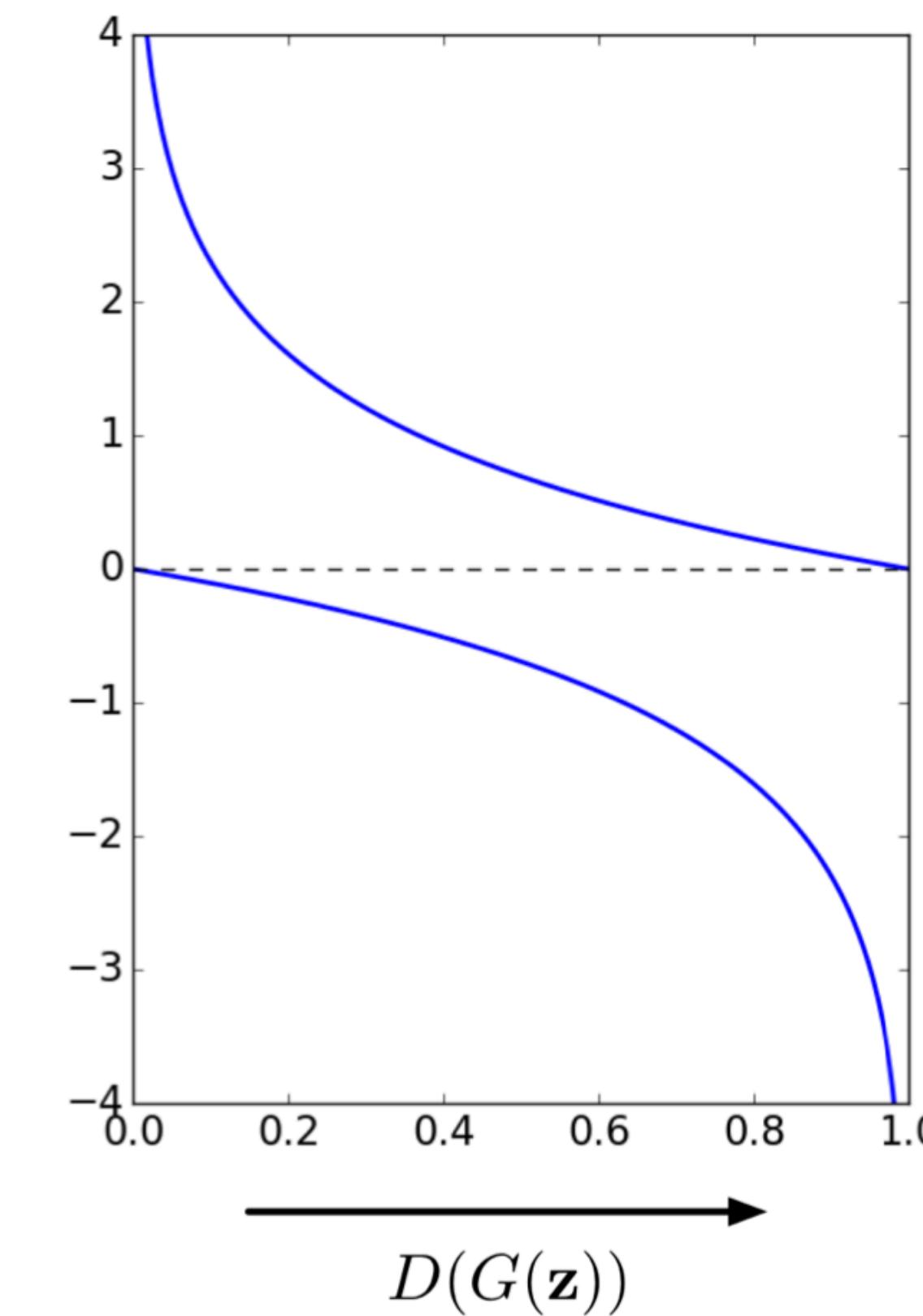
- Modified generator cost:

$$\mathcal{J}_G = \mathbb{E}_{\mathbf{z}}[-\log D(G(\mathbf{z}))]$$

- This fixes the saturation problem.

modified cost

minimax cost



(how well it fooled  
the discriminator)

# Generative Adversarial Networks

- Since GANs were introduced in 2014, there have been hundreds of papers introducing various architectures and training methods.
- Most modern architectures are based on the Deep Convolutional GAN (DC-GAN), where the generator and discriminator are both conv nets.
- GAN Zoo: <https://github.com/hindupuravinash/the-gan-zoo>
  - Good source of horrible puns (VEEGAN, Checkhov GAN, etc.)

# GAN Samples

Celebrities:



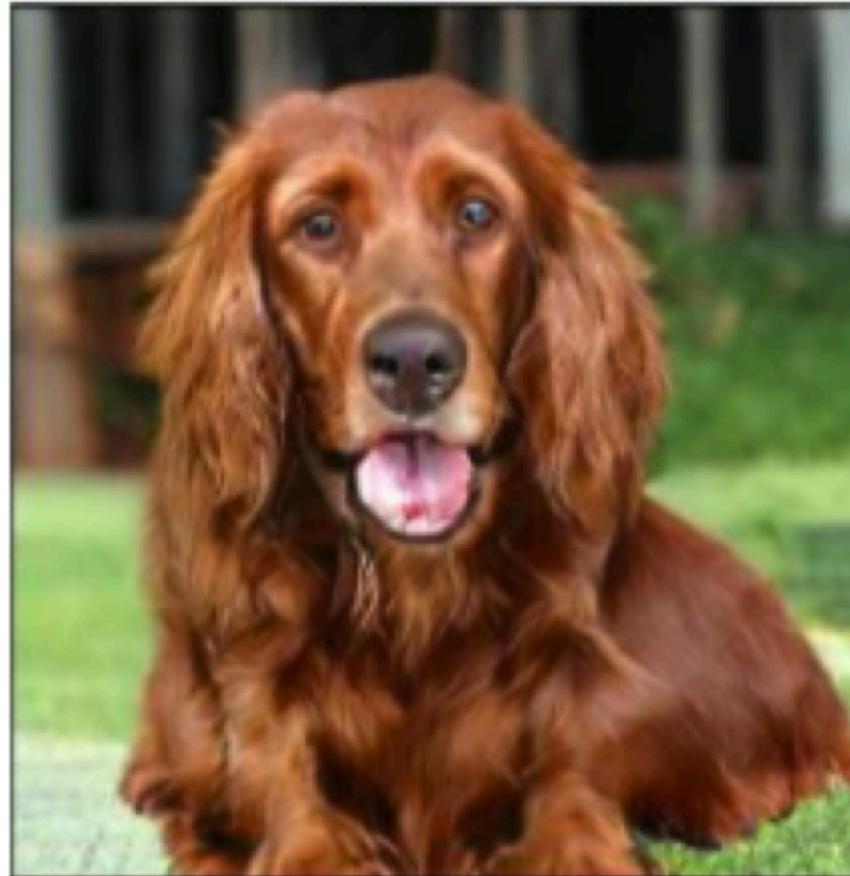
# GAN Samples

Bedrooms:



## GAN Samples

ImageNet object categories (by BigGAN, a much larger model with a bunch more engineering tricks):



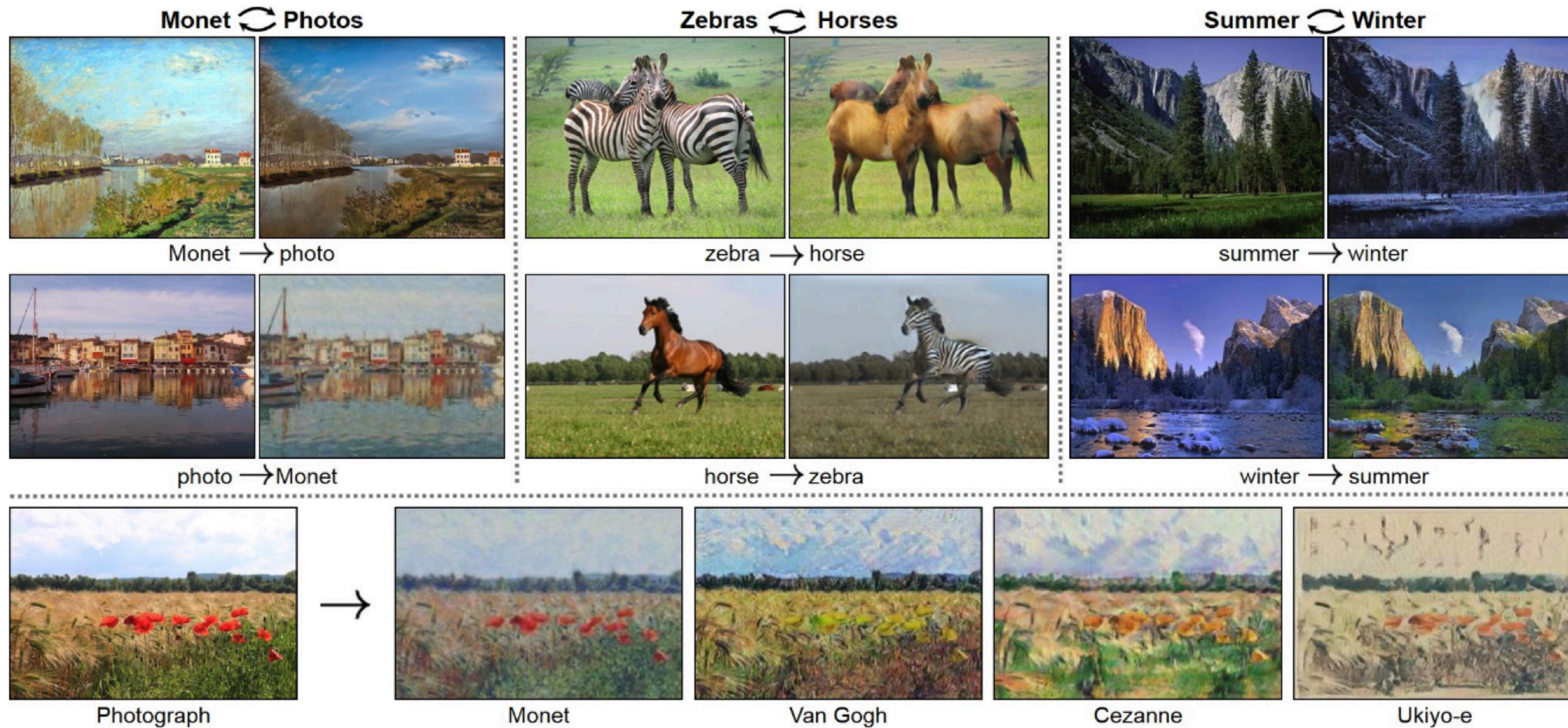
Brock et al., 2019. Large scale GAN training for high fidelity natural image synthesis.

## GAN Samples

- GANs revolutionized generative modeling by producing crisp, high-resolution images.
- The catch: we don't know how well they're modeling the distribution.
  - Can't measure the log-likelihood they assign to held-out data.
  - Could they be memorizing training examples? (E.g., maybe they sometimes produce photos of real celebrities?)
  - We have no way to tell if they are dropping important modes from the distribution.
  - See Wu et al., "On the quantitative analysis of decoder-based generative models" for partial answers to these questions.

# CycleGAN

Style transfer problem: change the style of an image while preserving the content.

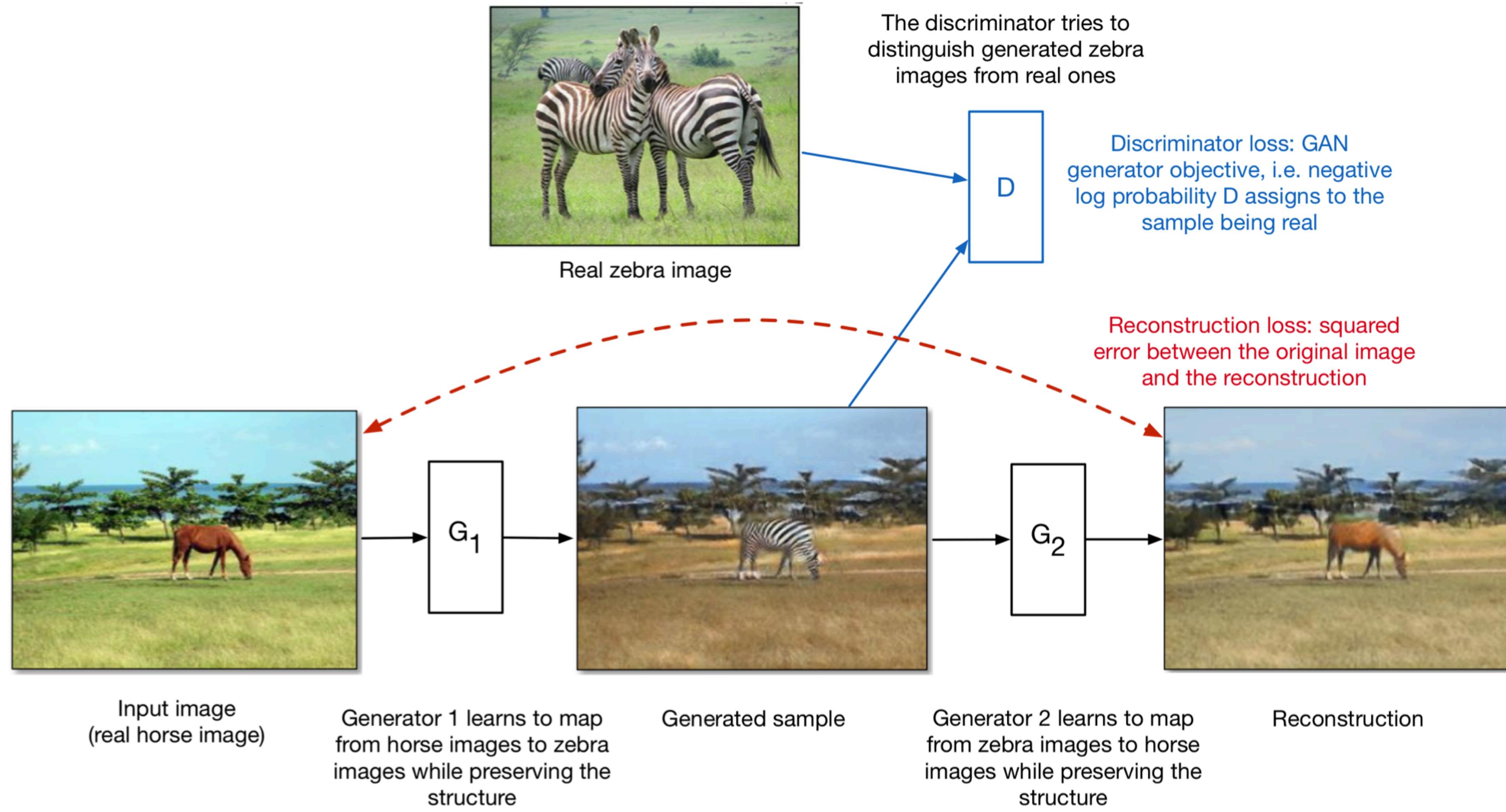


Data: Two unrelated collections of images, one for each style

# CycleGAN

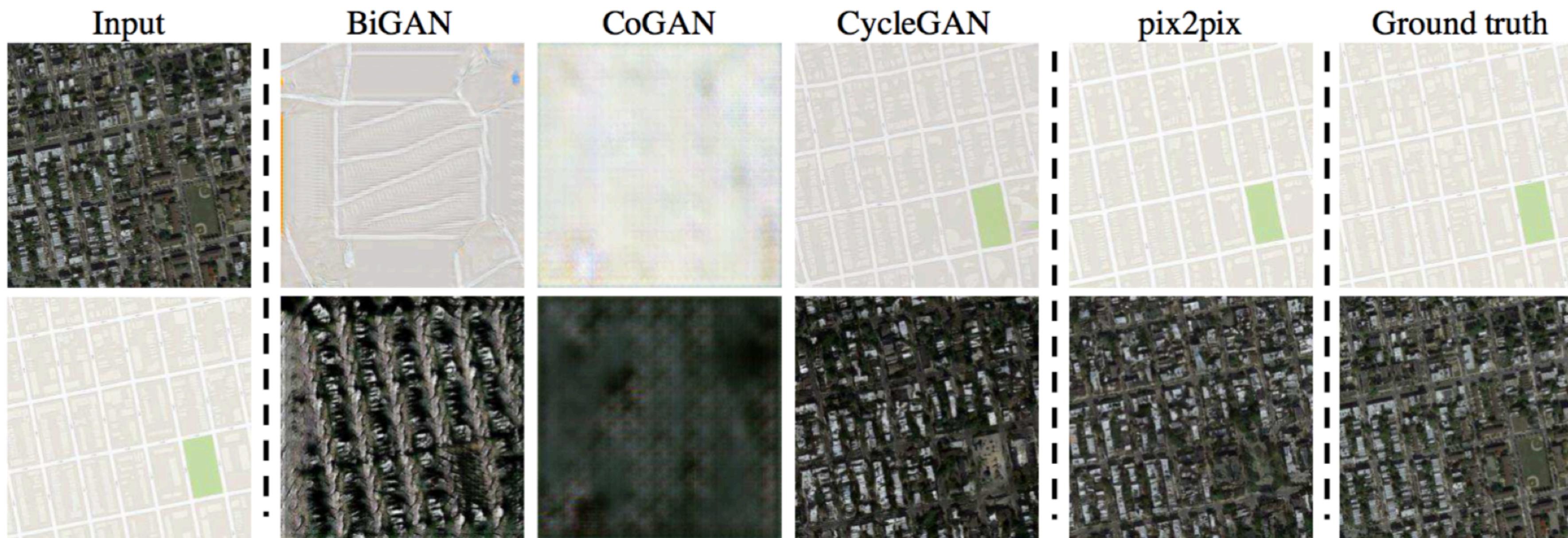
- If we had paired data (same content in both styles), this would be a supervised learning problem. But this is hard to find.
- The CycleGAN architecture learns to do it from unpaired data.
  - Train two different generator nets to go from style 1 to style 2, and vice versa.
  - Make sure the generated samples of style 2 are indistinguishable from real images by a discriminator net.
  - Make sure the generators are **cycle-consistent**: mapping from style 1 to style 2 and back again should give you almost the original image.

# CycleGAN



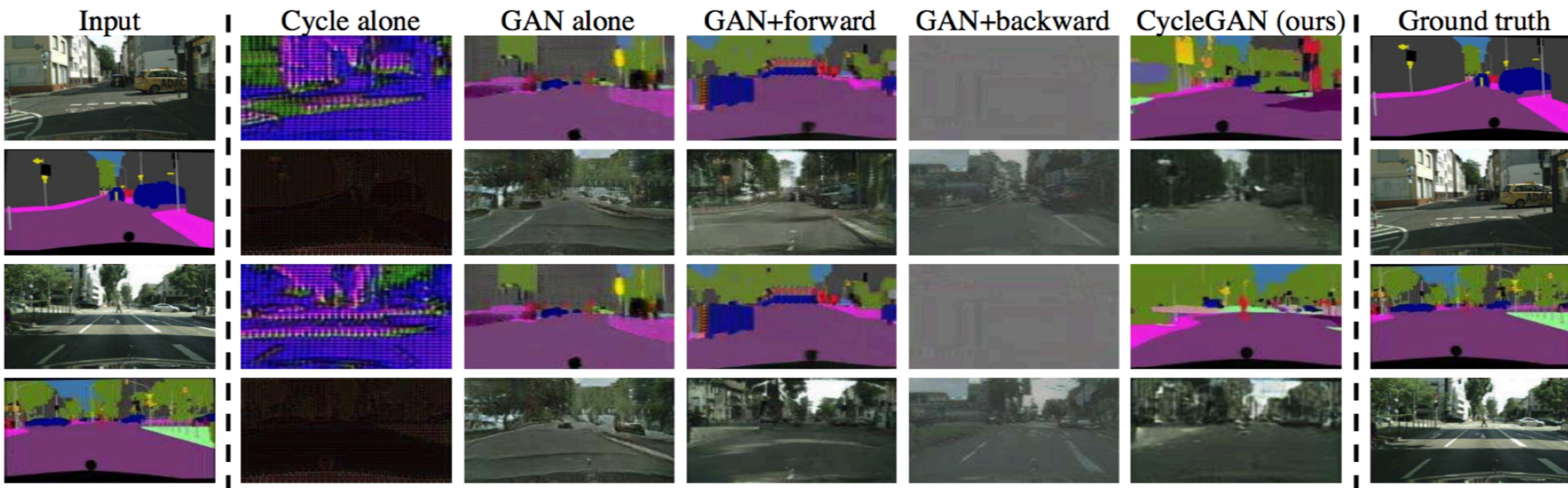
# CycleGAN

Style transfer between aerial photos and maps:



# CycleGAN

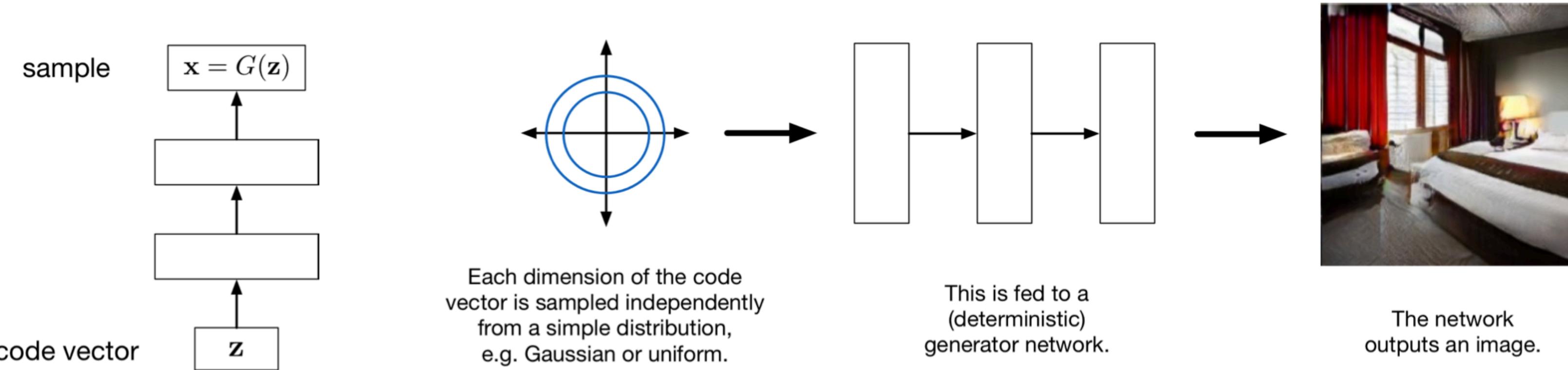
Style transfer between road scenes and semantic segmentations (labels of every pixel in an image by object category):



# Autoencoder

# Overview

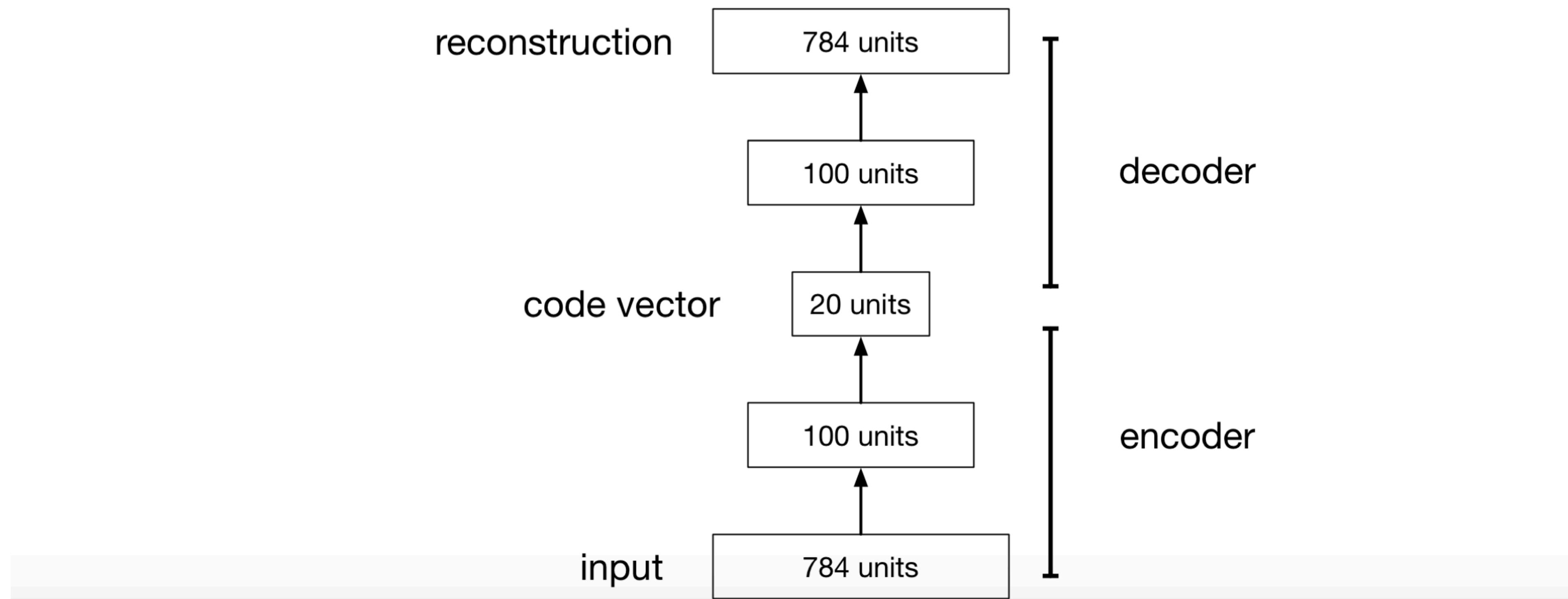
- Recall the generator network:



- One of the goals of unsupervised learning is to learn representations of images, sentences, etc.
- With reversible models,  $z$  and  $x$  must be the same size. Therefore, we can't reduce the dimensionality.
- Today, we'll cover the **variational autoencoder (VAE)**, a generative model that explicitly learns a low-dimensional representation.

# Autoencoders

- An **autoencoder** is a feed-forward neural net whose job it is to take an input  $x$  and predict  $x$ .
- To make this non-trivial, we need to add a **bottleneck layer** whose dimension is much smaller than the input.



# Autoencoders

## Why autoencoders?

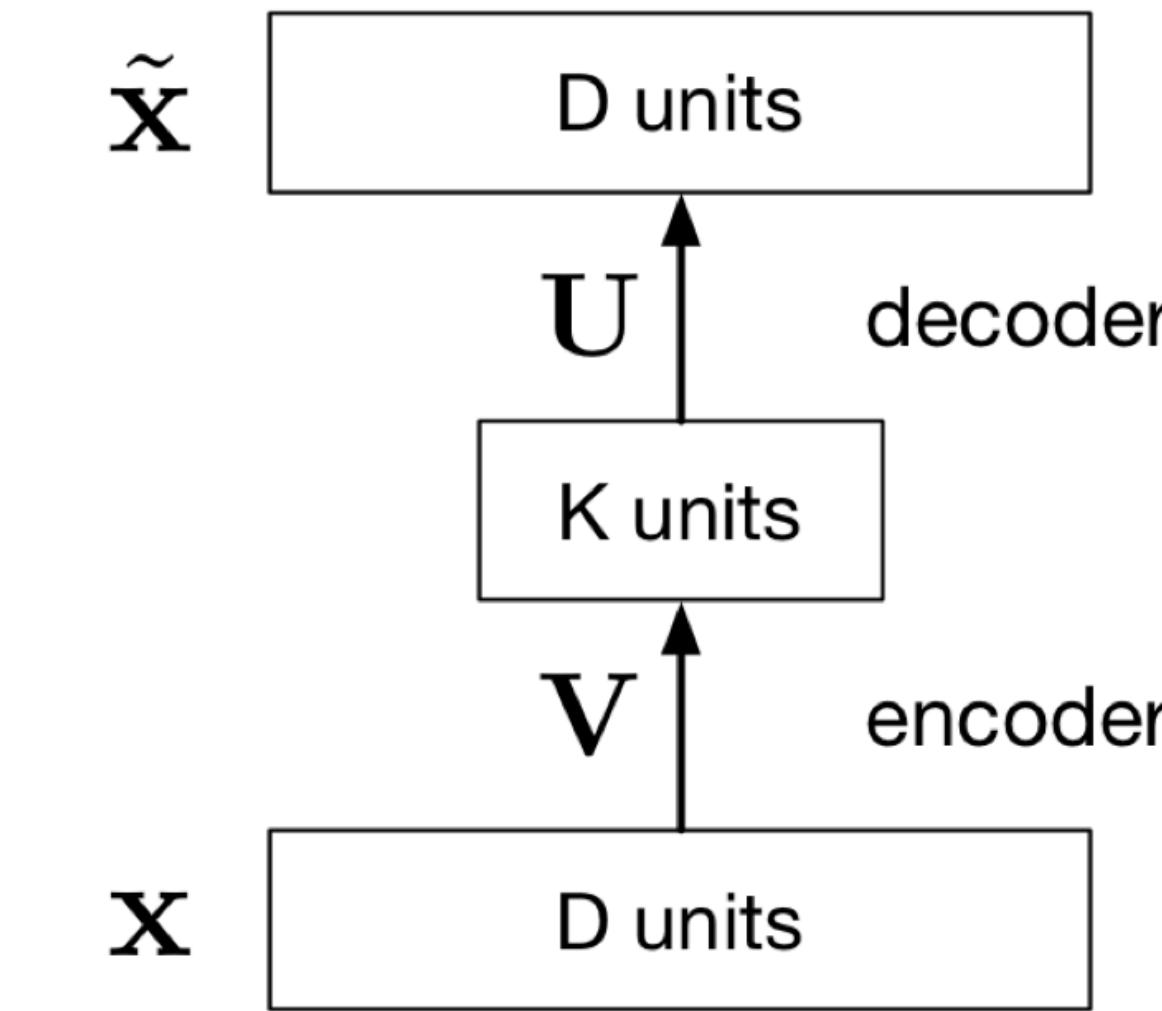
- Map high-dimensional data to two dimensions for visualization
- Compression (i.e. reducing the file size)
  - Note: this requires a VAE, not just an ordinary autoencoder.
- Learn abstract features in an unsupervised way so you can apply them to a supervised task
  - Unlabeled data can be much more plentiful than labeled data
- Learn a semantically meaningful representation where you can, e.g., interpolate between different images.

# Principal Component Analysis (optional)

- The simplest kind of autoencoder has one hidden layer, linear activations, and squared error loss.

$$\mathcal{L}(\mathbf{x}, \tilde{\mathbf{x}}) = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2$$

- This network computes  $\tilde{\mathbf{x}} = \mathbf{U}\mathbf{V}\mathbf{x}$ , which is a linear function.
- If  $K \geq D$ , we can choose  $\mathbf{U}$  and  $\mathbf{V}$  such that  $\mathbf{UV}$  is the identity. This isn't very interesting.
- But suppose  $K < D$ :
  - $\mathbf{V}$  maps  $\mathbf{x}$  to a  $K$ -dimensional space, so it's doing dimensionality reduction.
  - The output must lie in a  $K$ -dimensional subspace, namely the column space of  $\mathbf{U}$ .



# Principal Component Analysis (optional)

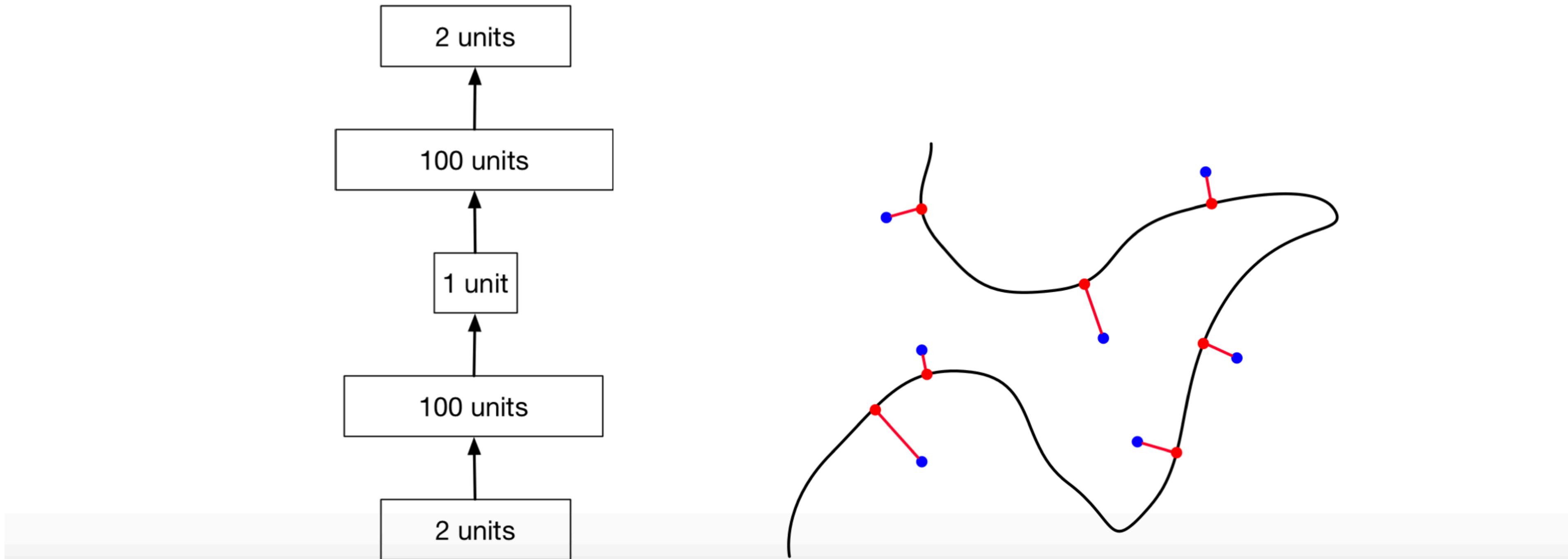
- autoencoders with squared error loss are equivalent to Principal Component Analysis (PCA).
- Two equivalent formulations:
  - Find the subspace that minimizes the reconstruction error.
  - Find the subspace that maximizes the projected variance.



“Eigenfaces”

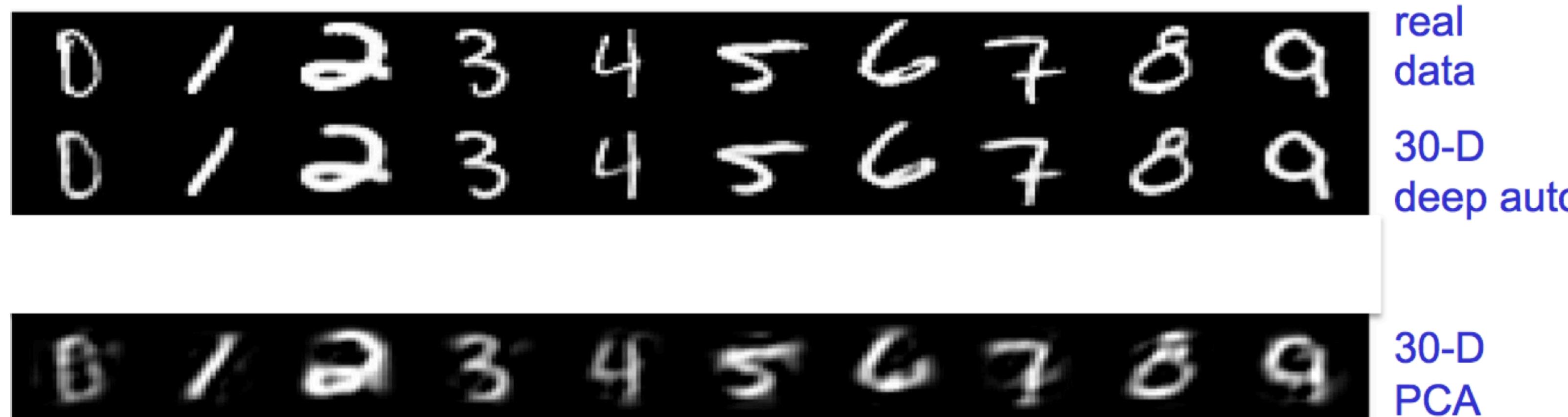
# Deep Autoencoders

- Deep nonlinear autoencoders learn to project the data, not onto a subspace, but onto a nonlinear **manifold**
- This manifold is the image of the decoder.
- This is a kind of **nonlinear dimensionality reduction**.



# Deep Autoencoders

- Nonlinear autoencoders can learn more powerful codes for a given dimensionality, compared with linear autoencoders (PCA)



# Deep Autoencoders

- Some limitations of autoencoders
  - They're not generative models, so they don't define a distribution
  - How to choose the latent dimension?

# Observation Model

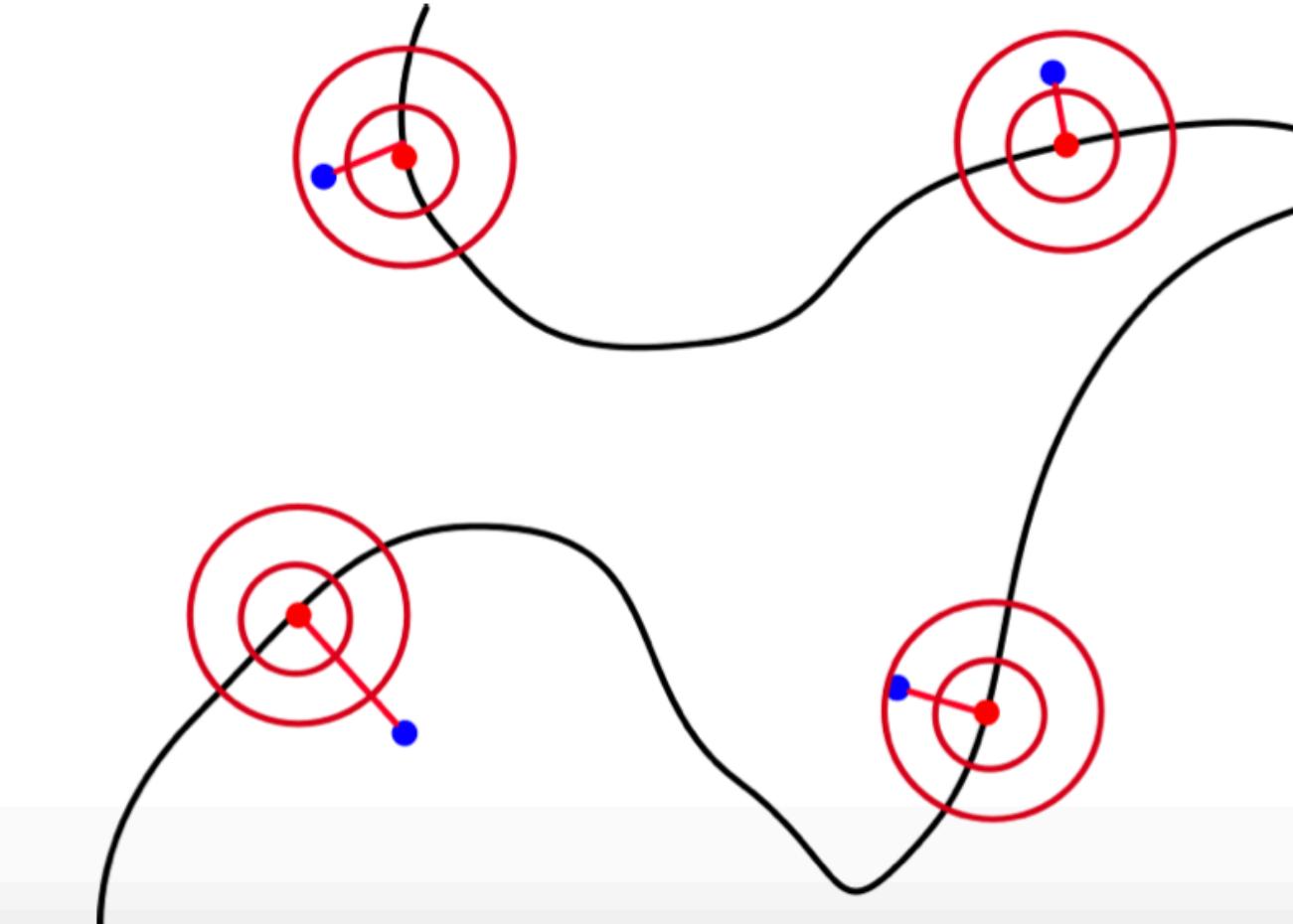
- Consider training a generator network with maximum likelihood.

$$p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) d\mathbf{z}$$

- One problem: if  $\mathbf{z}$  is low-dimensional and the decoder is deterministic, then  $p(\mathbf{x}) = 0$  almost everywhere!
  - The model only generates samples over a low-dimensional sub-manifold of  $\mathcal{X}$ .
- Solution: define a noisy observation model, e.g.

$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x}; G_\theta(\mathbf{z}), \eta \mathbf{I}),$$

where  $G_\theta$  is the function computed by the decoder with parameters  $\theta$ .



## Observation Model

- At least  $p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) d\mathbf{z}$  is well-defined, but how can we compute it?
  - The decoder function  $G_\theta(\mathbf{z})$  is very complicated, so there's no hope of finding a closed form.
- Instead, we will try to maximize a lower bound on  $\log p(\mathbf{x})$ .
  - The math is essentially the same as in the EM algorithm

# Variational Inference

- We obtain the lower bound using **Jensen's Inequality**: for a convex function  $h$  of a random variable  $X$ ,

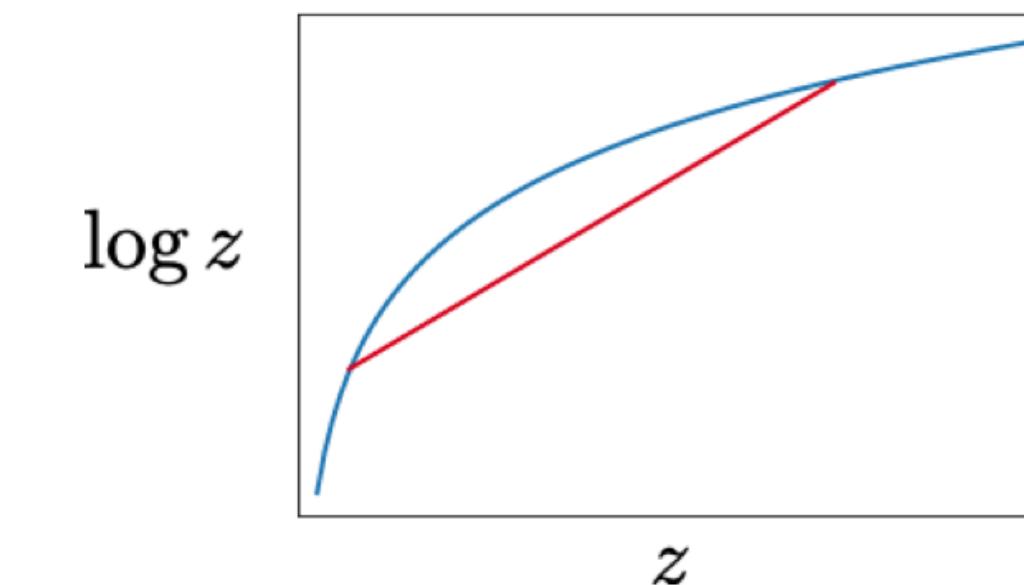
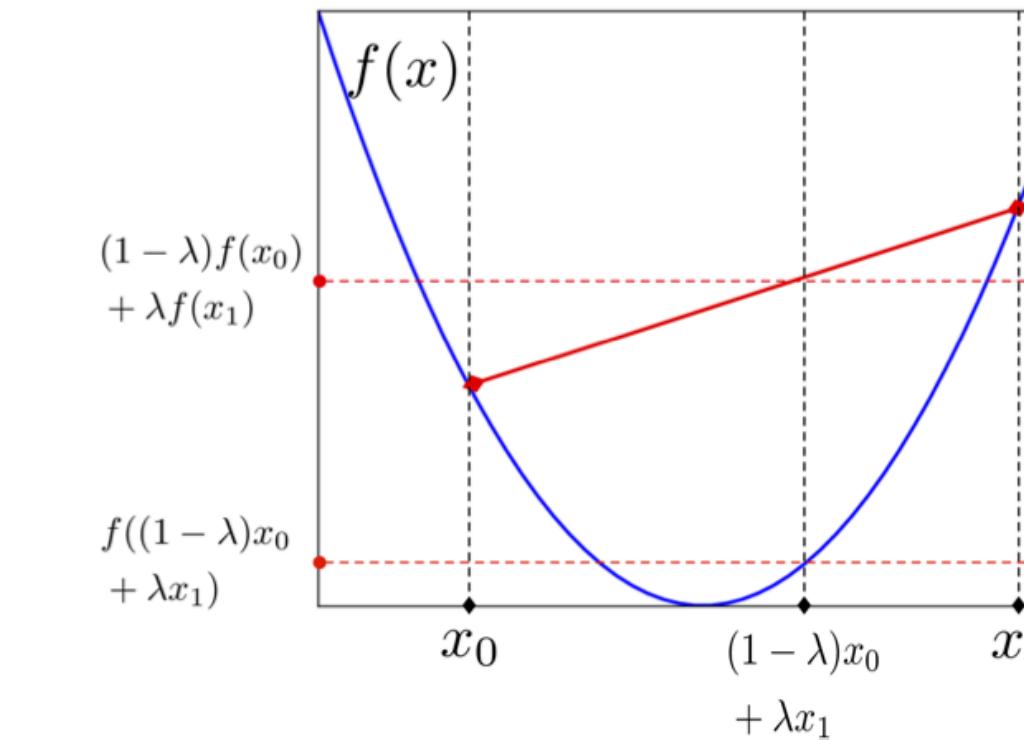
$$\mathbb{E}[h(X)] \geq h(\mathbb{E}[X])$$

Therefore, if  $h$  is **concave** (i.e.  $-h$  is convex),

$$\mathbb{E}[h(X)] \leq h(\mathbb{E}[X])$$

- The function  $\log z$  is concave. Therefore,

$$\mathbb{E}[\log X] \leq \log \mathbb{E}[X]$$



# Variational Inference

- Suppose we have some distribution  $q(\mathbf{z})$ . (We'll see later where this comes from.)
- We use Jensen's Inequality to obtain the lower bound.

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{z}) p(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &= \log \int q(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &\geq \int q(\mathbf{z}) \log \left[ \frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}|\mathbf{z}) \right] d\mathbf{z} \quad (\text{Jensen's Inequality}) \\ &= \mathbb{E}_q \left[ \log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right] + \mathbb{E}_q [\log p(\mathbf{x}|\mathbf{z})]\end{aligned}$$

- We'll look at these two terms in turn.

## Variational Inference

- The first term we'll look at is  $\mathbb{E}_q [\log p(\mathbf{x}|\mathbf{z})]$
- Since we assumed a Gaussian observation model,

$$\begin{aligned}\log p(\mathbf{x}|\mathbf{z}) &= \log \mathcal{N}(\mathbf{x}; G_{\theta}(\mathbf{z}), \eta \mathbf{I}) \\ &= \log \left[ \frac{1}{(2\pi\eta)^{D/2}} \exp \left( -\frac{1}{2\eta} \|\mathbf{x} - G_{\theta}(\mathbf{z})\|^2 \right) \right] \\ &= -\frac{1}{2\eta} \|\mathbf{x} - G_{\theta}(\mathbf{z})\|^2 + \text{const}\end{aligned}$$

- So this term is the expected squared error in reconstructing  $\mathbf{x}$  from  $\mathbf{z}$ . We call it the **reconstruction term**.

# Variational Inference

- The second term is  $\mathbb{E}_q \left[ \log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right]$ .
- This is just  $-D_{KL}(q(\mathbf{z}) \| p(\mathbf{z}))$ , where  $D_{KL}$  is the **Kullback-Leibler (KL) divergence**

$$D_{KL}(q(\mathbf{z}) \| p(\mathbf{z})) \triangleq \mathbb{E}_q \left[ \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \right]$$

- KL divergence is a widely used measure of distance between probability distributions, though it doesn't satisfy the axioms to be a distance metric.
- Typically,  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Hence, the KL term encourages  $q$  to be close to  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- We'll give the KL term a much more interesting interpretation when we discuss Bayesian neural nets.

## Variational Inference

- Hence, we're trying to maximize the **variational lower bound**, or **variational free energy**:

$$\log p(\mathbf{x}) \geq \mathcal{F}(\boldsymbol{\theta}, q) = \mathbb{E}_q [\log p(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q||p).$$

- The term “variational” is a historical accident: “variational inference” used to be done using variational calculus, but this isn’t how we train VAEs.
- We’d like to choose  $q$  to make the bound as tight as possible.
- It’s possible to show that the gap is given by:

$$\log p(\mathbf{x}) - \mathcal{F}(\boldsymbol{\theta}, q) = D_{\text{KL}}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})).$$

Therefore, we’d like  $q$  to be as close as possible to the posterior distribution  $p(\mathbf{z}|\mathbf{x})$ .

- Let's think about the role of each of the two terms.
- The reconstruction term

$$\mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] = -\frac{1}{2\sigma^2} \mathbb{E}_q[\|\mathbf{x} - G_\theta(\mathbf{z})\|^2] + \text{const}$$

is minimized when  $q$  is a **point mass** on

$$\mathbf{z}_* = \arg \min_{\mathbf{z}} \|\mathbf{x} - G_\theta(\mathbf{z})\|^2.$$

- But a point mass would have infinite KL divergence. (Exercise: check this.) So the KL term forces  $q$  to be more spread out.

# Reparameterization Trick

- To fit  $q$ , let's assign it a parametric form, in particular a Gaussian distribution:  $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$  and  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_K^2)$ .
- In general, it's hard to differentiate through an expectation. But for Gaussian  $q$ , we can apply the **reparameterization trick**:

$$z_i = \mu_i + \sigma_i \epsilon_i,$$

where  $\epsilon_i \sim \mathcal{N}(0, 1)$ .

- Hence,

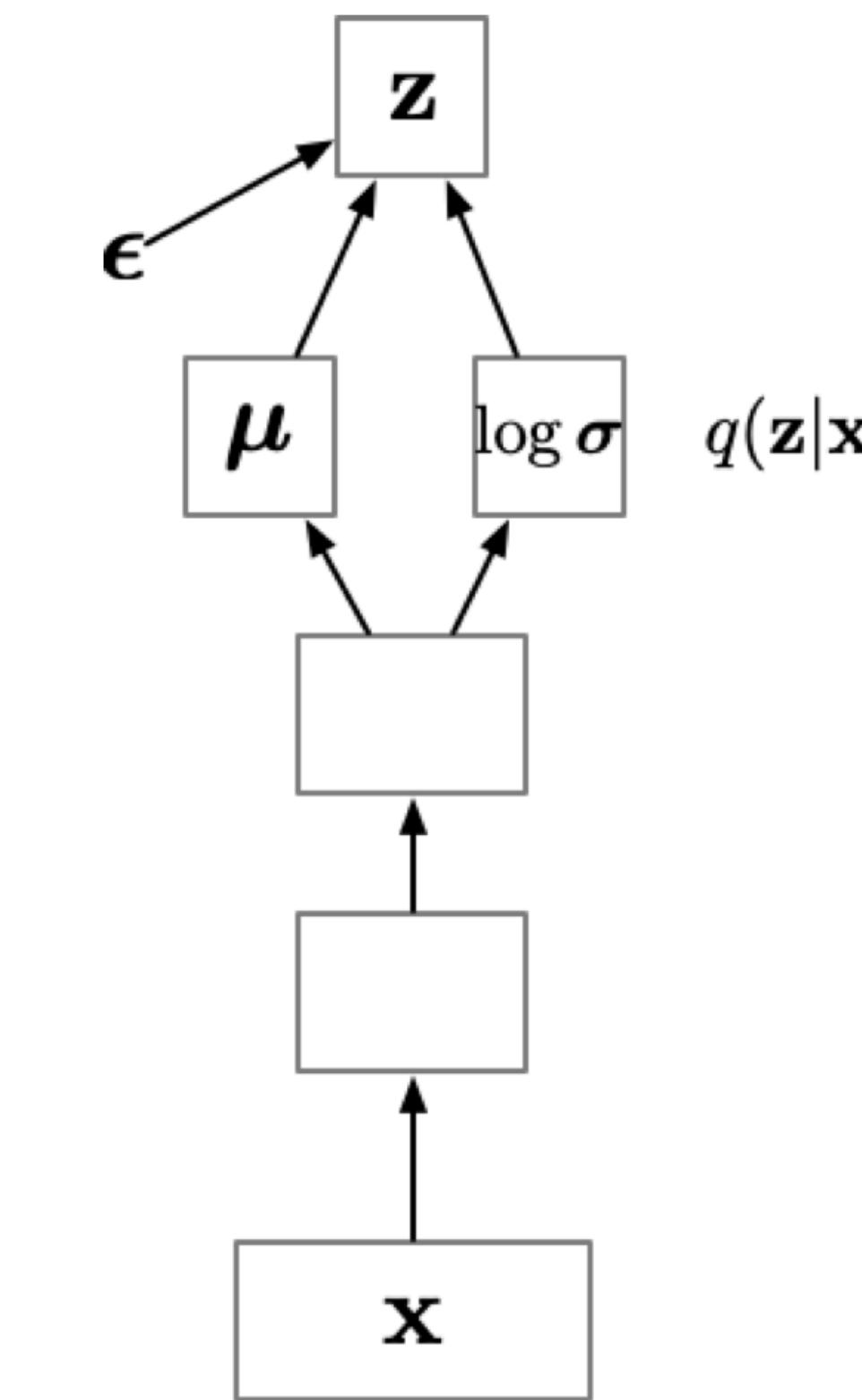
$$\overline{\mu_i} = \overline{z_i} \quad \overline{\sigma_i} = \overline{z_i} \epsilon_i.$$

# Amortization

- This suggests one strategy for learning the decoder. For each training example,
  - ➊ Fit  $q$  to approximate the posterior for the current  $\mathbf{x}$  by doing many steps of gradient ascent on  $\mathcal{F}$ .
  - ➋ Update the decoder parameters  $\theta$  with gradient ascent on  $\mathcal{F}$ .
- **Problem:** this requires an expensive iterative procedure for every training example, so it will take a long time to process the whole training set.

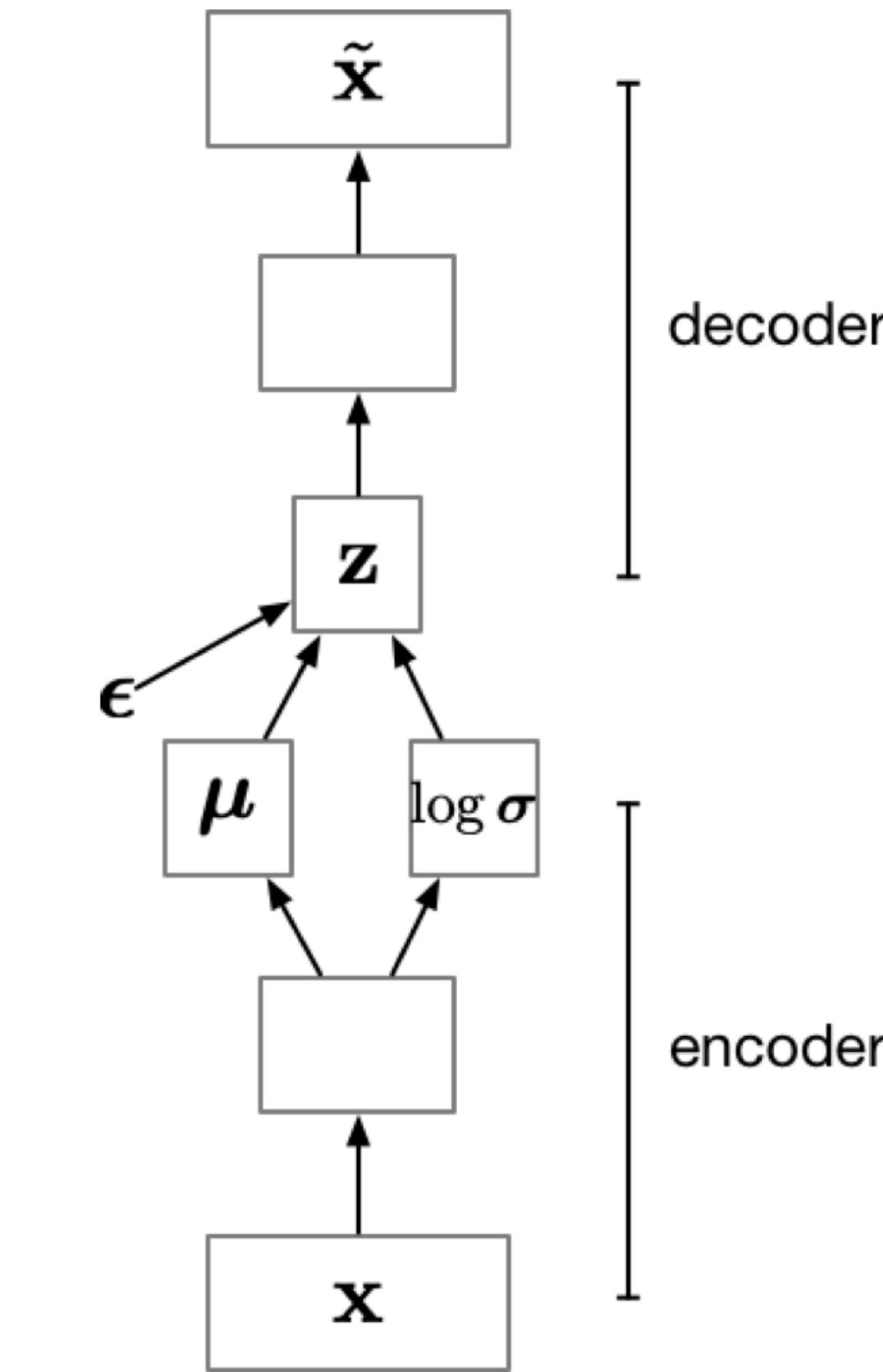
# Amortization

- Idea: amortize the cost of inference by learning an inference network which predicts  $(\mu, \Sigma)$  as a function of  $\mathbf{x}$ .
- The outputs of the inference net are  $\mu$  and  $\log \sigma$ . (The log representation ensures  $\sigma > 0$ .)
- If  $\sigma \approx 0$ , then this network essentially computes  $\mathbf{z}$  deterministically, by way of  $\mu$ .
  - But the KL term encourages  $\sigma > 0$ , so in general  $\mathbf{z}$  will be noisy.
  - The notation  $q(\mathbf{z}|\mathbf{x})$  emphasizes that  $q$  depends on  $\mathbf{x}$ , even though it's not actually a conditional distribution.



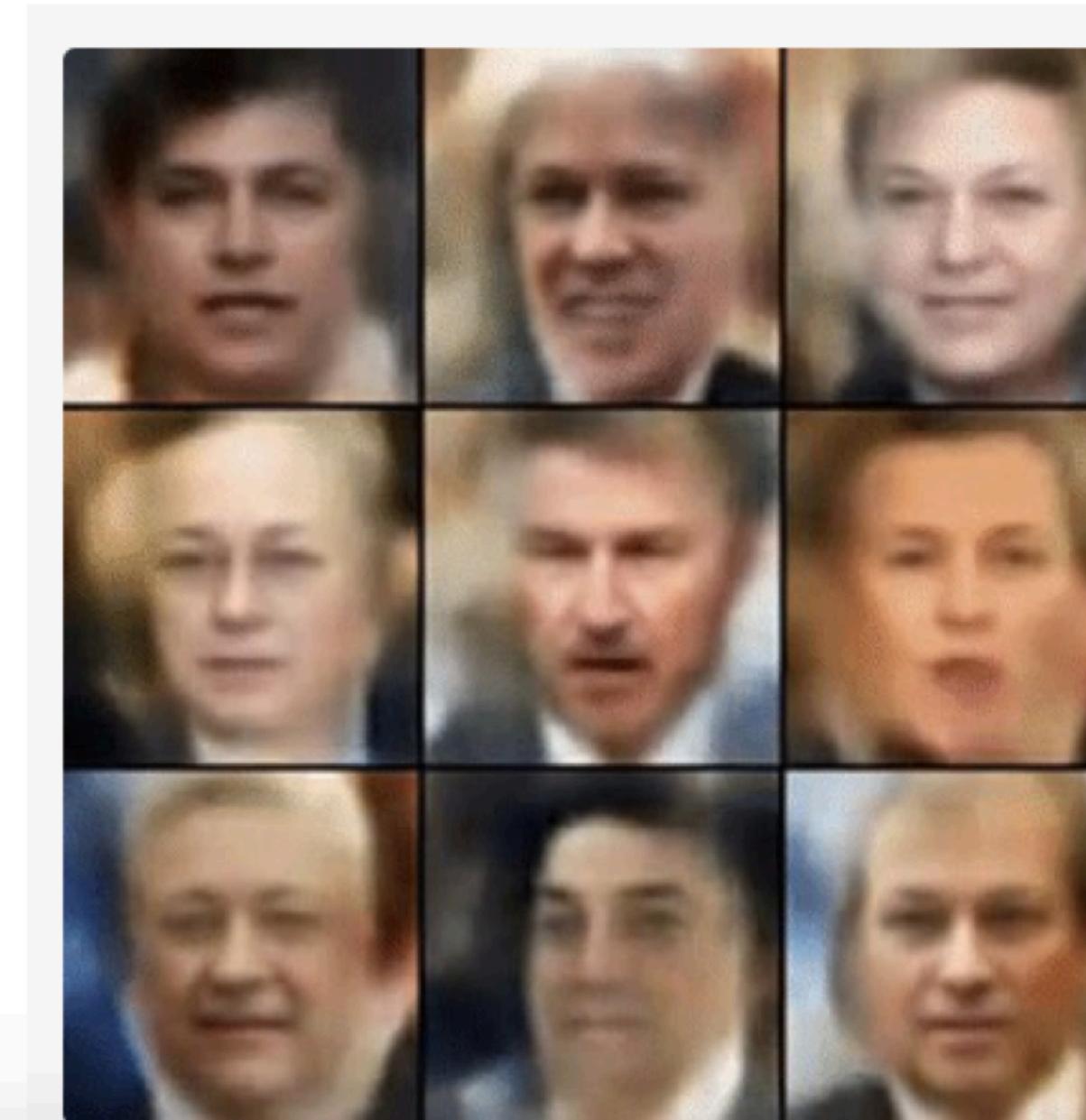
# Amortization

- Combining this with the decoder network, we see the structure closely resembles an ordinary autoencoder. The inference net is like an encoder.
- Hence, this architecture is known as a **variational autoencoder (VAE)**.
- The parameters of both the encoder and decoder networks are updated using a single pass of ordinary backprop.
  - The reconstruction term corresponds to squared error  $\|\mathbf{x} - \tilde{\mathbf{x}}\|^2$ , like in an ordinary VAE.
  - The KL term regularizes the representation by encouraging  $\mathbf{z}$  to be more stochastic.



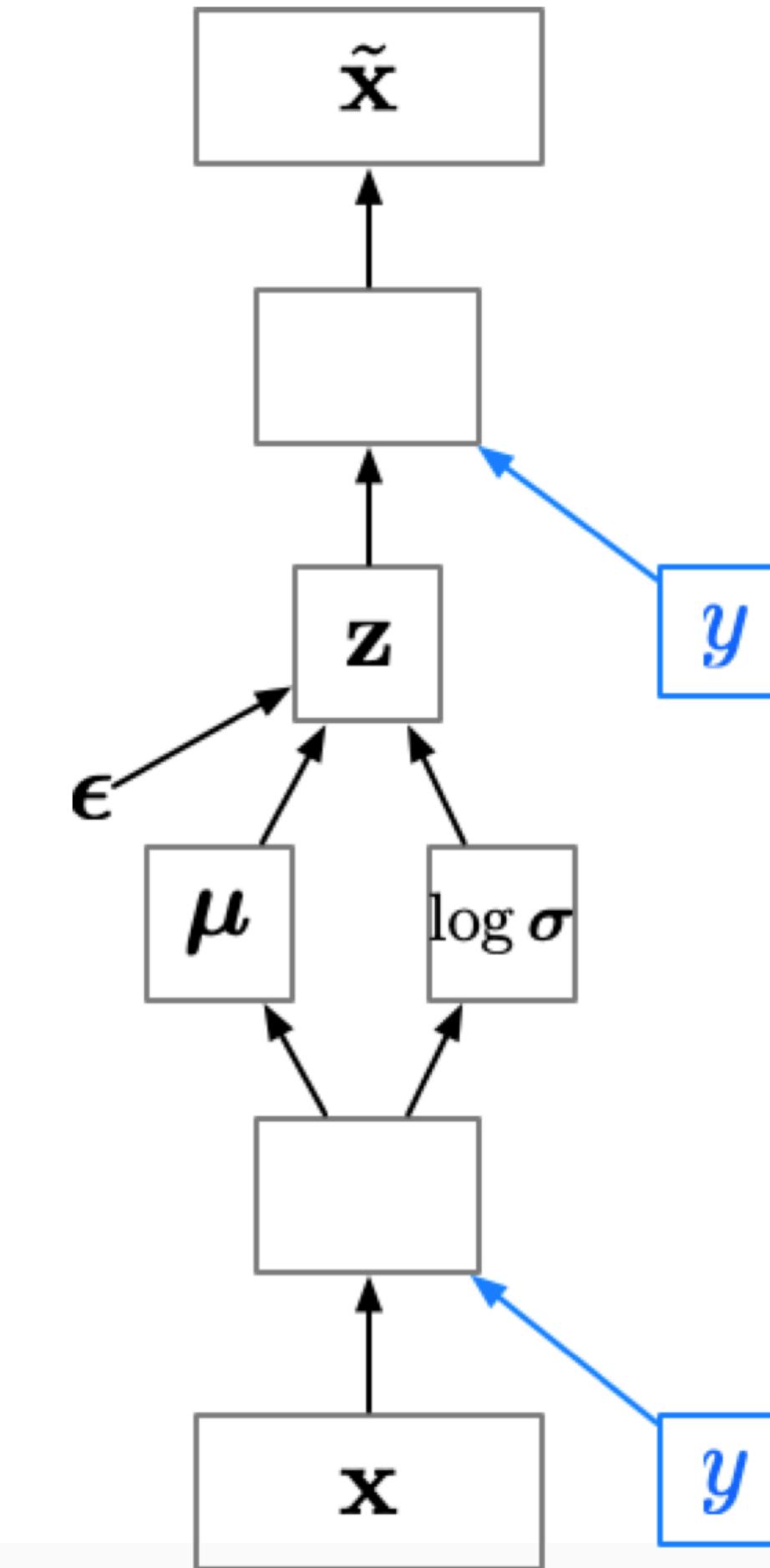
# VAEs vs. Other Generative Models

- In short, a VAE is like an autoencoder, except that it's also a generative model (defines a distribution  $p(\mathbf{x})$ ).
- Unlike autoregressive models, generation only requires one forward pass.
- Unlike reversible models, we can fit a low-dimensional latent representation. We'll see we can do interesting things with this...



# Class-Conditional VAE

- So far, we haven't used the labels  $y$ . A **class-conditional VAE** provides the labels to both the encoder and the decoder.
- Since the latent code  $z$  no longer has to model the image category, it can focus on modeling the stylistic features.
- If we're lucky, this lets us **disentangle** style and content. (Note: disentanglement is still a dark art.)
- See Kingma et al., "Semi-supervised learning with deep generative models."

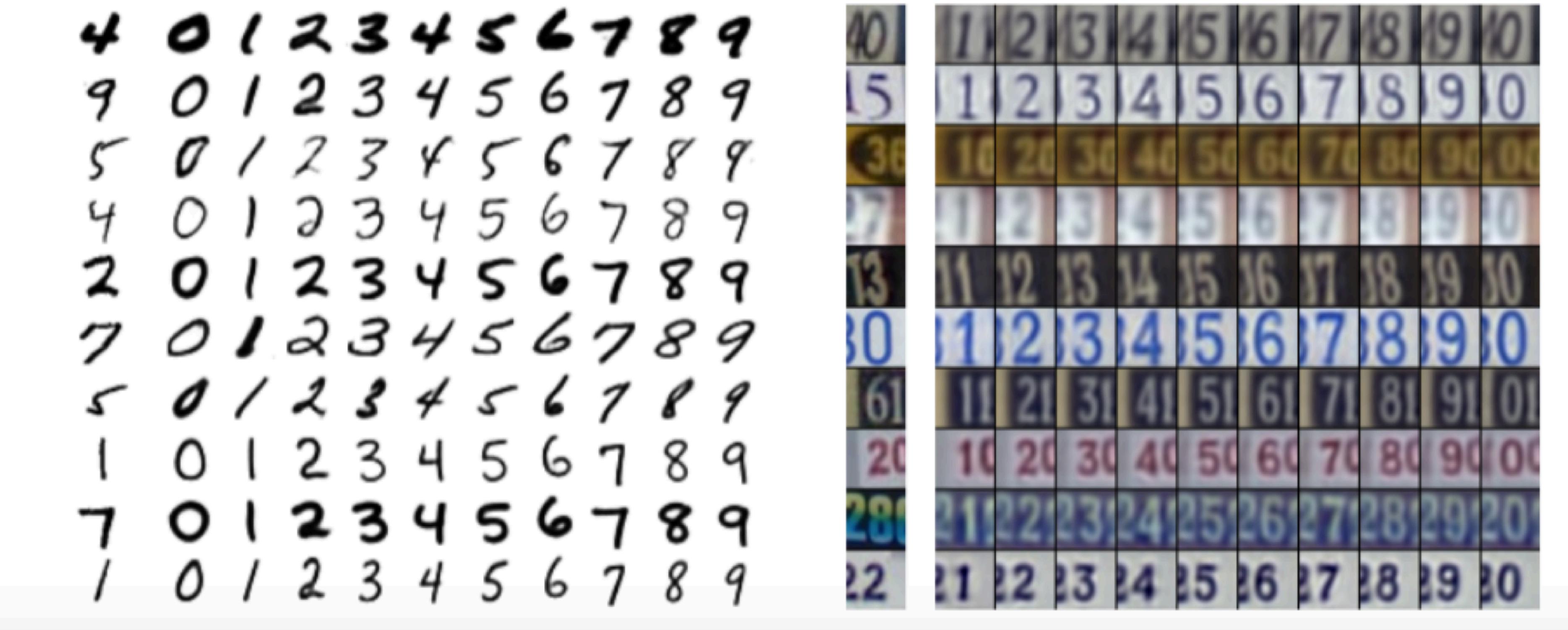


# Class-Conditional VAE

- By varying two **latent dimensions** (i.e. dimensions of  $z$ ) while holding  $y$  fixed, we can visualize the **latent space**.

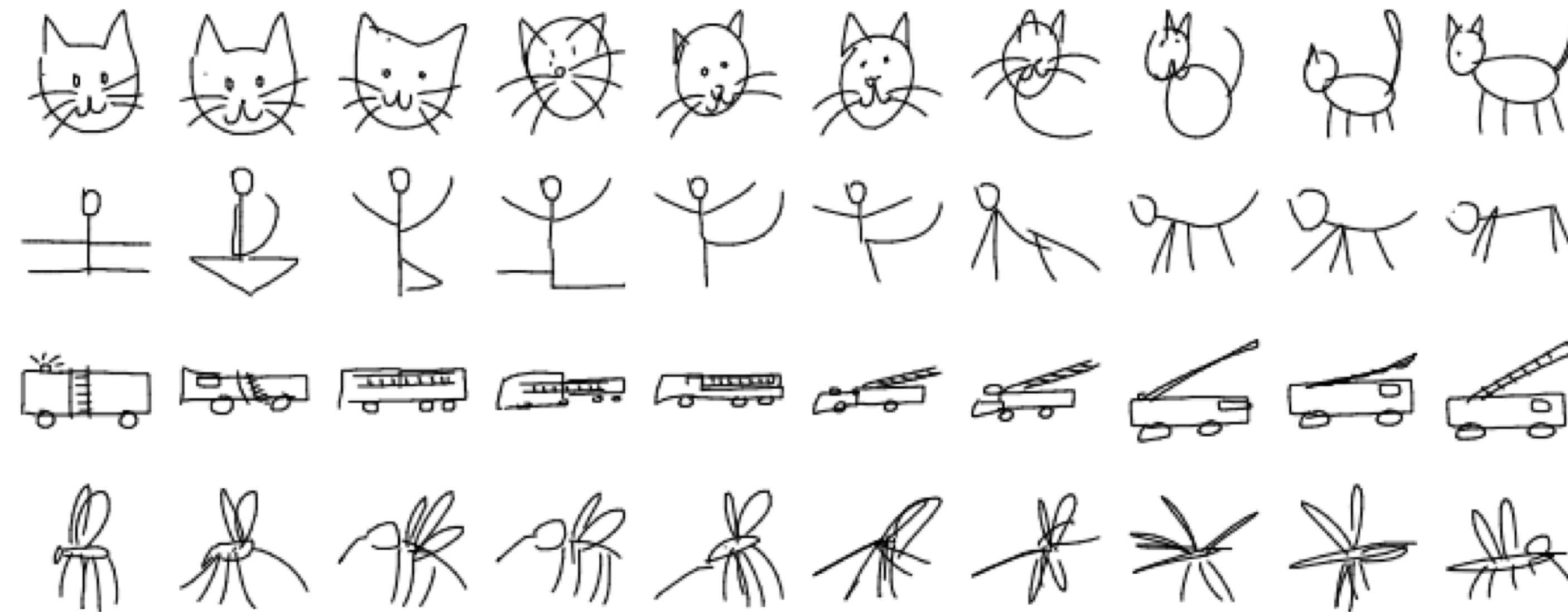
# Class-Conditional VAE

- By varying the label  $y$  while holding  $z$  fixed, we can solve image analogies.



# Latent Space Interpolations

- You can often get interesting results by interpolating between two vectors in the latent space:



Ha and Eck, "A neural representation of sketch drawings"

# Latent Space Interpolations

- Latent space interpolation of music:  
<https://magenta.tensorflow.org/music-vae>