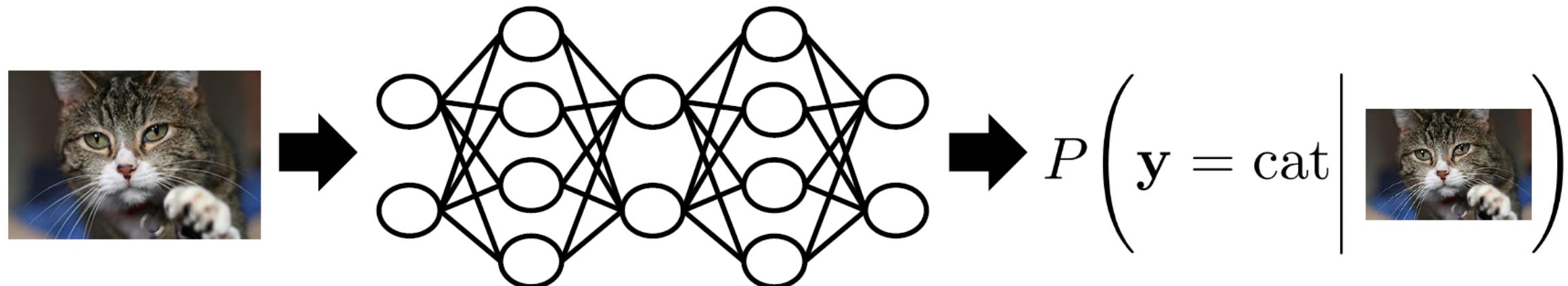


Deep Generative Modeling

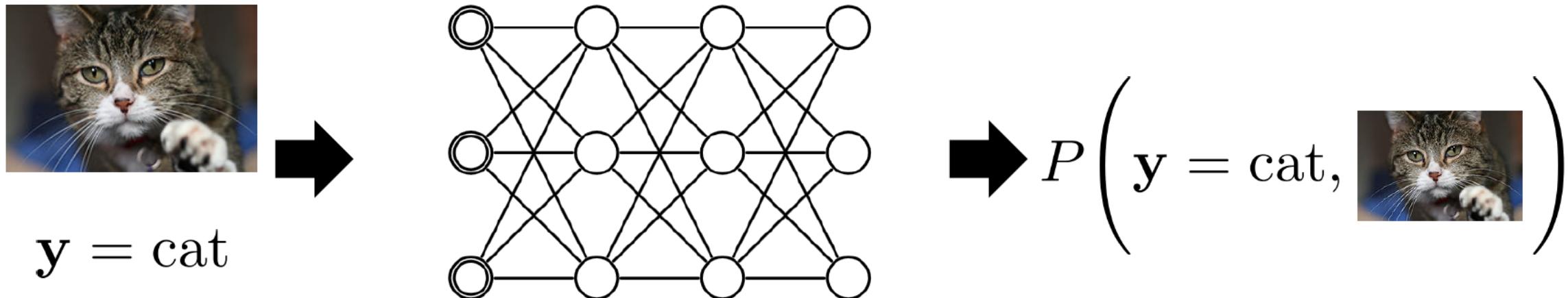
Discriminative vs. Generative Models

- Given an observed variable x and a target variable y ,
- Discriminative modeling estimates the conditional prob $p(y|x)$
 - E.g., logistic regression, ImageNet classifiers, etc.



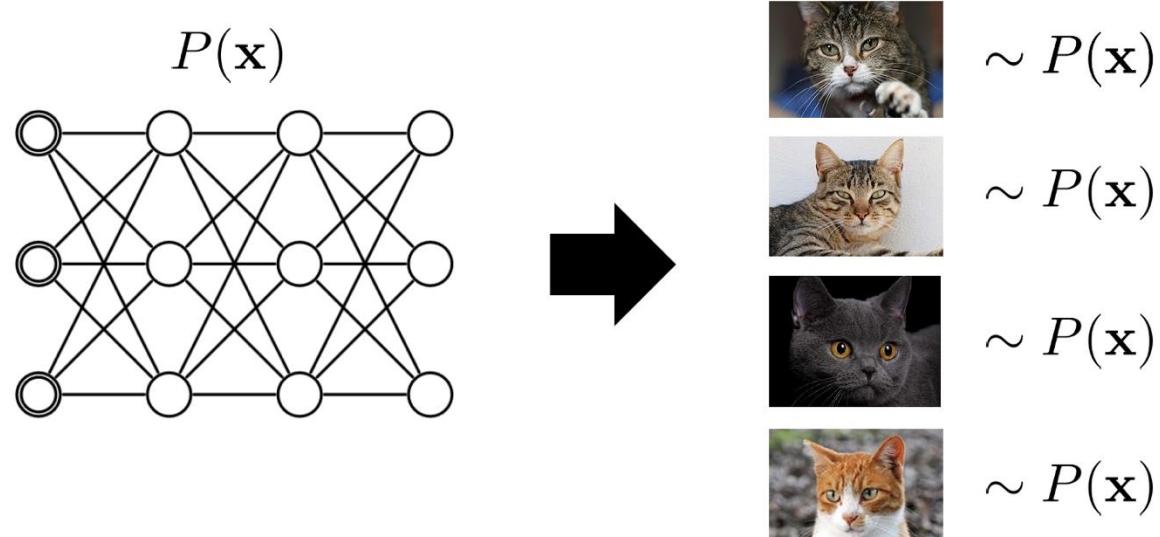
Discriminative vs. Generative Models

- Given an observed variable x and a target variable y ,
- Generative modeling estimates the joint distribution $p(x, y)$
 - E.g., Boltzmann machines, Gaussian mixture models, etc.



Why Generative Models?

- Without assuming y , generative models learn $p(\mathbf{x})$ from given data.
- $p(\mathbf{x})$ enables use to generate new data similar to the training dataset.
- We can use various sampling methods for generation based on $p(\mathbf{x})$.



Why Generative Models?

- Many real-world problems can be formulated as generation tasks.

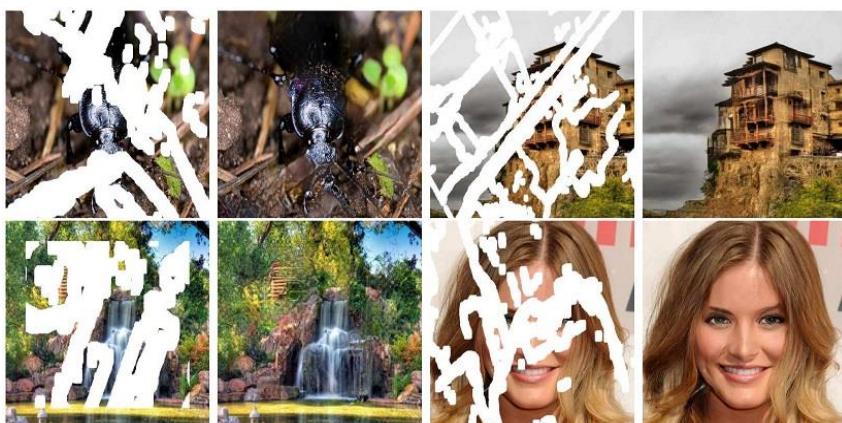
Image generation



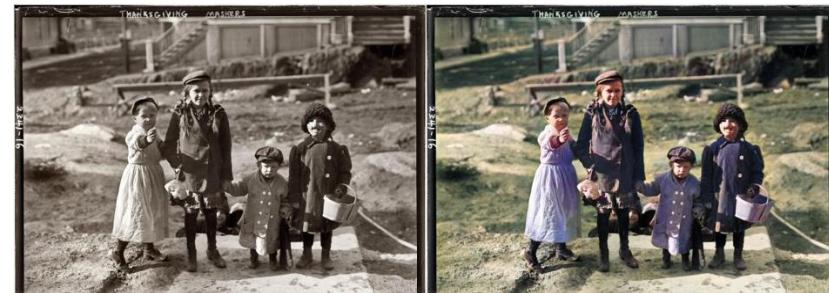
$p(\text{high res} / \text{low res})$



$p(\text{full image} / \text{masked image})$

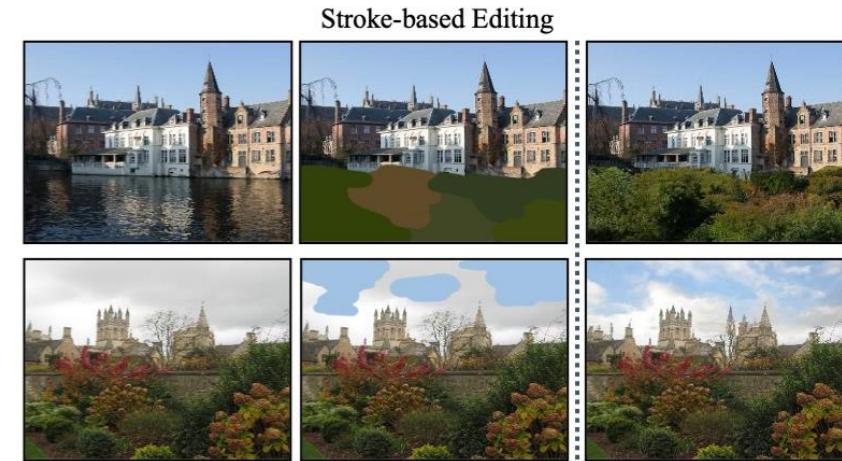
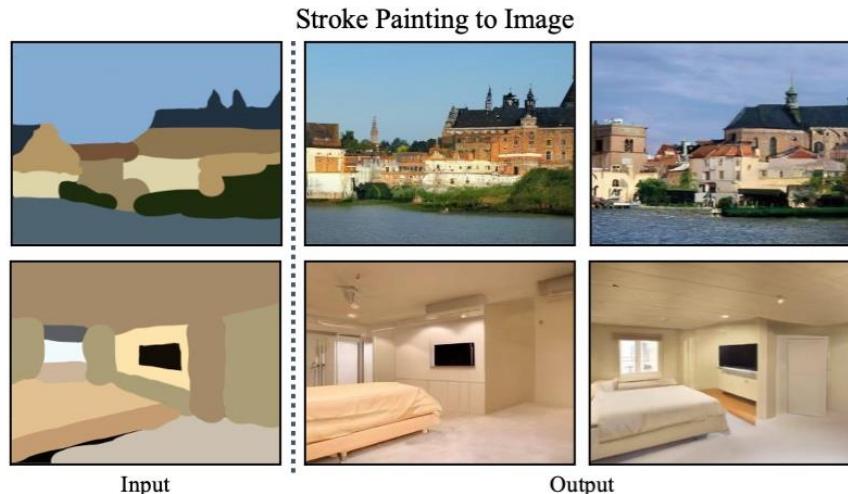


$p(\text{color image} / \text{grayscale})$



Why Generative Models?

- Many real-world problems can be formulated as generation tasks.



Why Generative Models?

- Many real-world problems can be formulated as generation tasks.

Custom prompt

To get an A+ in deep generative models, students have to

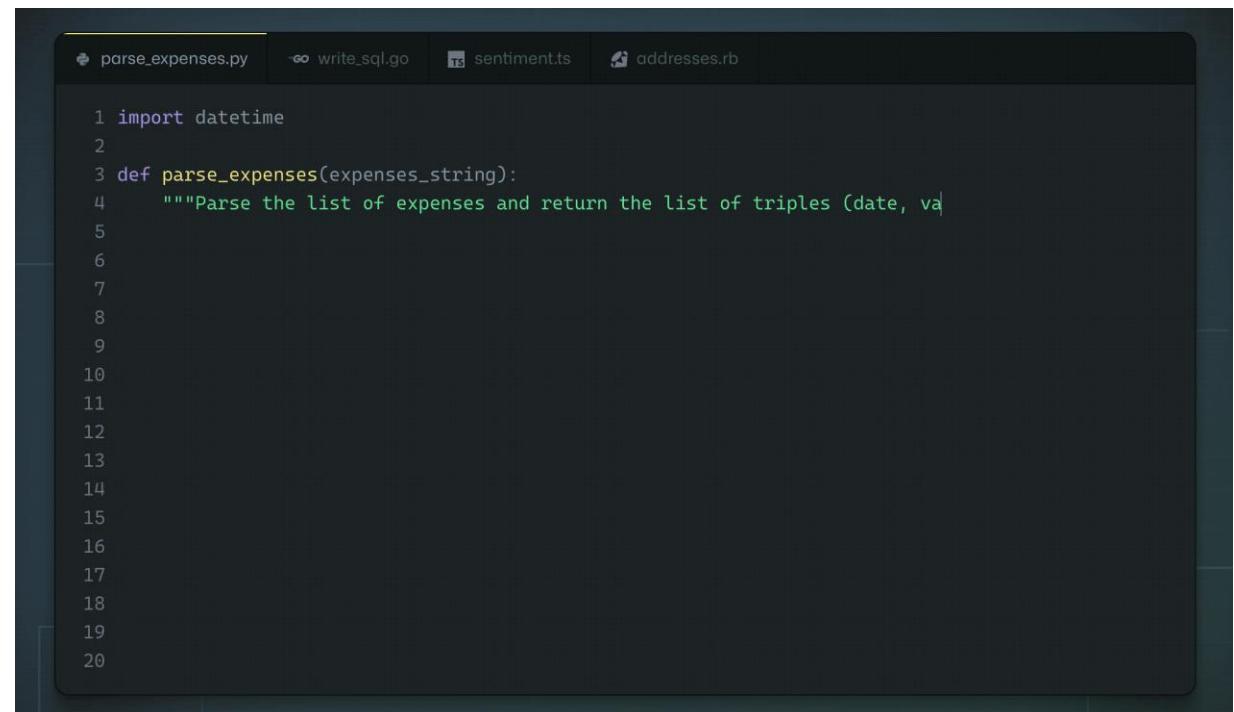
Completion

To get an A+ in deep generative models, students have to be willing to work with problems that are a whole lot more interesting than, say, the ones that most students work on in class. If you're a great student, the question above can be avoided and you'll be able to do great work, but if you're not, you will need to go beyond the basics before getting good.

Now to be clear, this advice is not just for the deep-learning crowd; it is good advice for any student who is taking his or her first course in machine learning.

The key point is that if you have a deep, deep brain of a computer scientist, that's just as important to you.

$$p(\text{next word} \mid \text{prev words})$$



```
1 import datetime
2
3 def parse_expenses(expenses_string):
4     """Parse the list of expenses and return the list of triples (date, va
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
```

Why Generative Models?

- Many real-world problems can be formulated as generation tasks.

TEXT PROMPT an armchair in the shape of an avocado....

AI-GENERATED
IMAGES



Edit prompt or view more images↓

$p(\text{image} \mid \text{caption})$

TEXT PROMPT a store front that has the word 'openai' written on it....

AI-GENERATED
IMAGES



Why Generative Models?

- Many real-world problems can be formulated as generation tasks.



February 2022



April 2022



July 2022



November 2022



March 2023



March 2023 (new version)



June 2023



December 2023

Midjourney: Prompt: “a hyper-realistic image of Harry Potter”

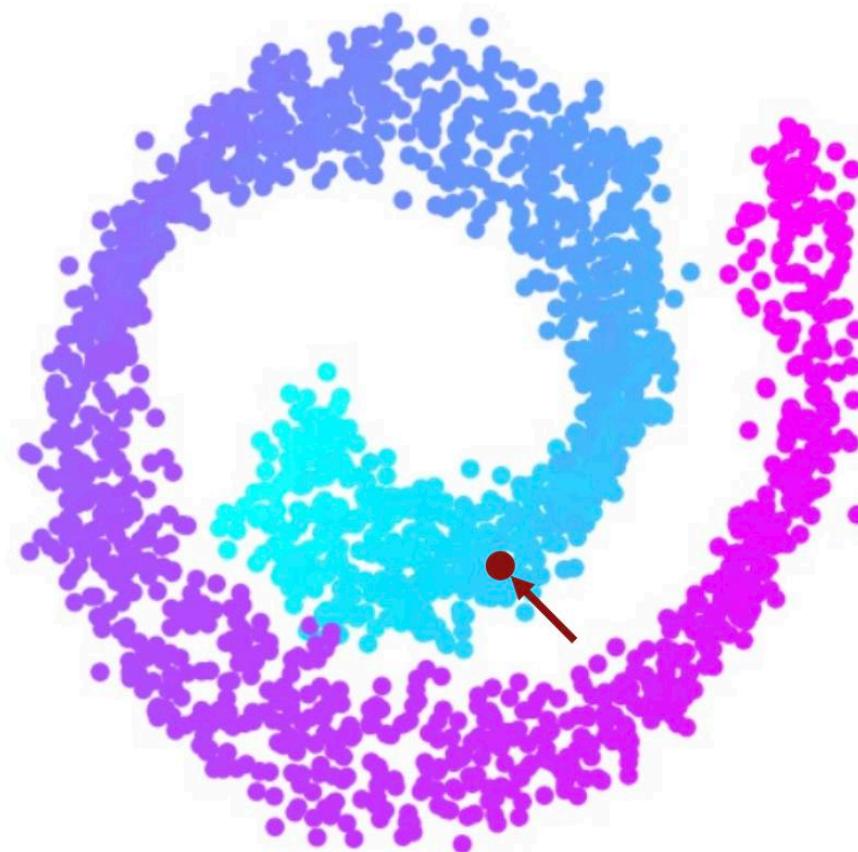
Generation and Sampling

- How can we generate a new realistic photo?



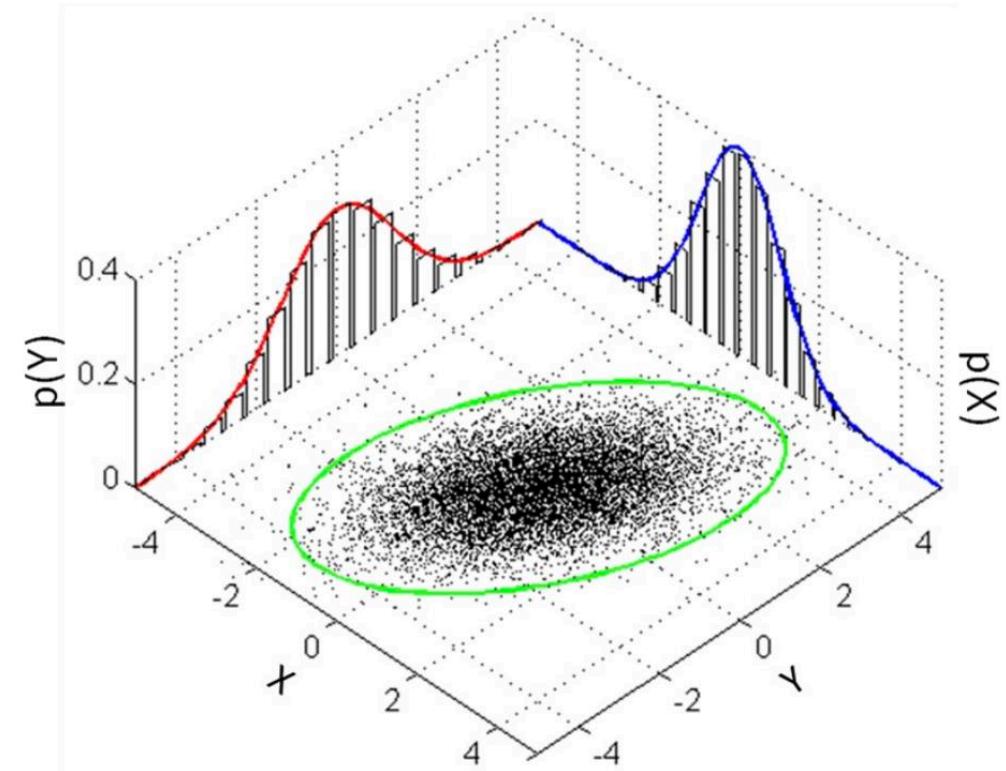
Generation and Sampling

- How can we sample a new plausible 2D point?



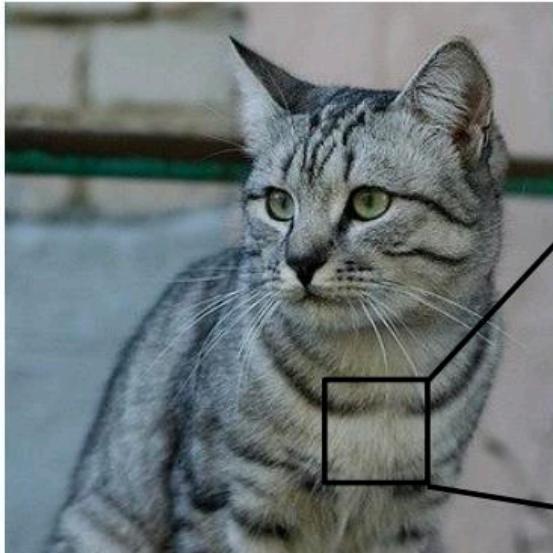
Statistical Perspective

- We will view generation/sampling as
 - there being a probability distribution of the data, and
 - the given points are samples from the probability distribution.
- If the probability distribution is known, then for some specific distributions, we can sample from them directly.



Statistical Perspective for Real Images

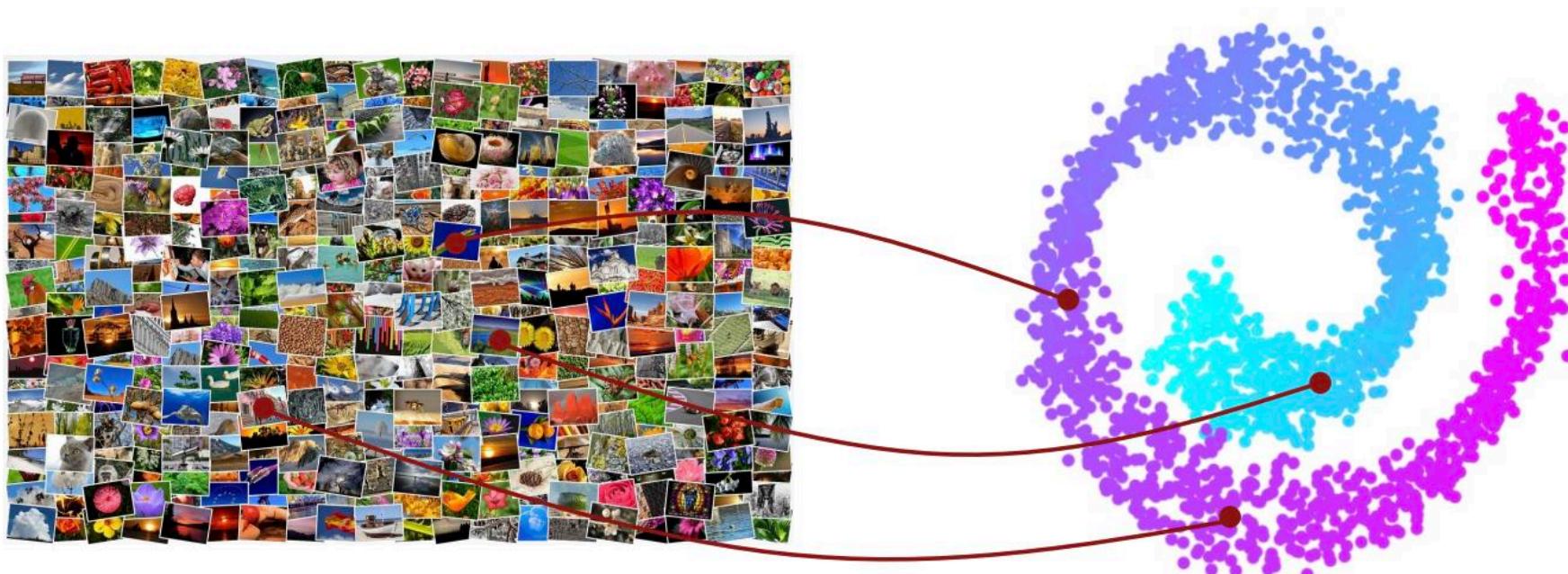
- Let's consider RGB images with a resolution of 256x256.
- An image can be represented by a 256x256x3 vector.
- This means that an image is a point in a 256x256x3-dimensional space.



```
[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 99 85]
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]
[128 137 144 140 109 95 86 70 62 65 63 63 60 73 86 101]
[125 133 148 137 119 121 117 94 65 79 80 65 54 64 72 98]
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]
[ 89 93 90 97 108 147 131 118 113 114 113 109 106 95 77 88]
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]
[128 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 84]]
```

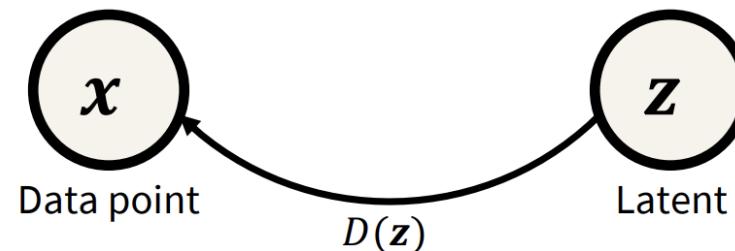
Statistical Perspective for Real Images

- Let's consider real images $\{x_1, x_2, \dots, x_n\}$ as samples from a data distribution $p(x)$ that measures how likely it is for an image to be a real photo.
- Can we derive the PDF of the data distribution?



The Basic Idea

- Map a simple distribution $p(\mathbf{z})$ (e.g., a standard normal distribution) to the data distribution $p(\mathbf{x})$
 - \mathbf{z} : Latent variable, $p(\mathbf{z})$: Latent distribution
- Sample from $p(\mathbf{z})$ and map it to a data point.



- How to map a latent distribution $p(\mathbf{z})$ to the data distribution $p(\mathbf{x})$?

Let's use a neural network for $D(\mathbf{z})$!



Variational Autoencoder

Autoencoder (AE)

- A starting point is the autoencoder (AE).
 - AE consists of two networks: a deterministic encoder and a deterministic decoder.

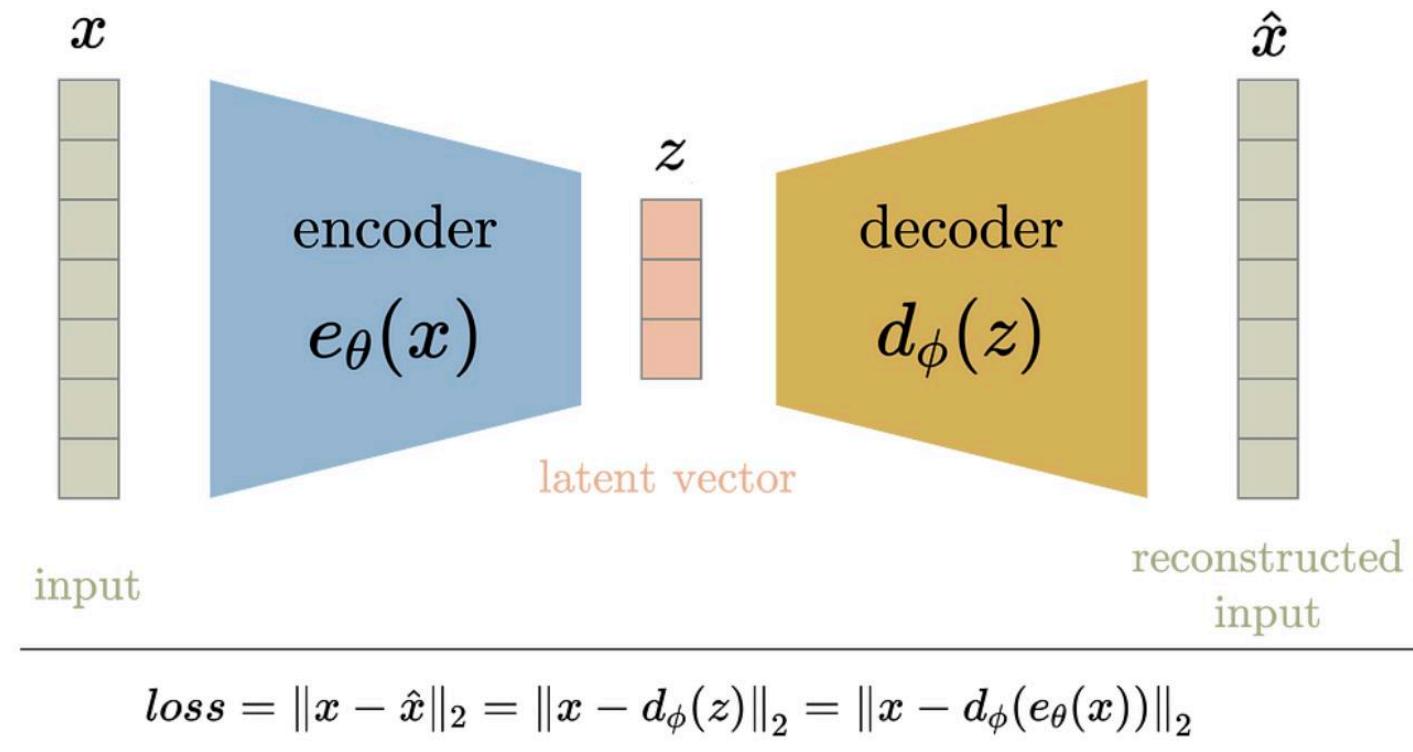


Figure: <https://towardsdatascience.com/difference-between-autoencoder-ae-and-variational-autoencoder-vae-ed7be1c038f2/>

Autoencoder (AE)

- What we need is the decoder (latent \rightarrow input data).
 - AE enables accurate reconstruction, but the latent space is unstructured: randomly sampling latent codes usually produces meaningless outputs.
 - How do we guarantee that a latent is mapped to a data point?

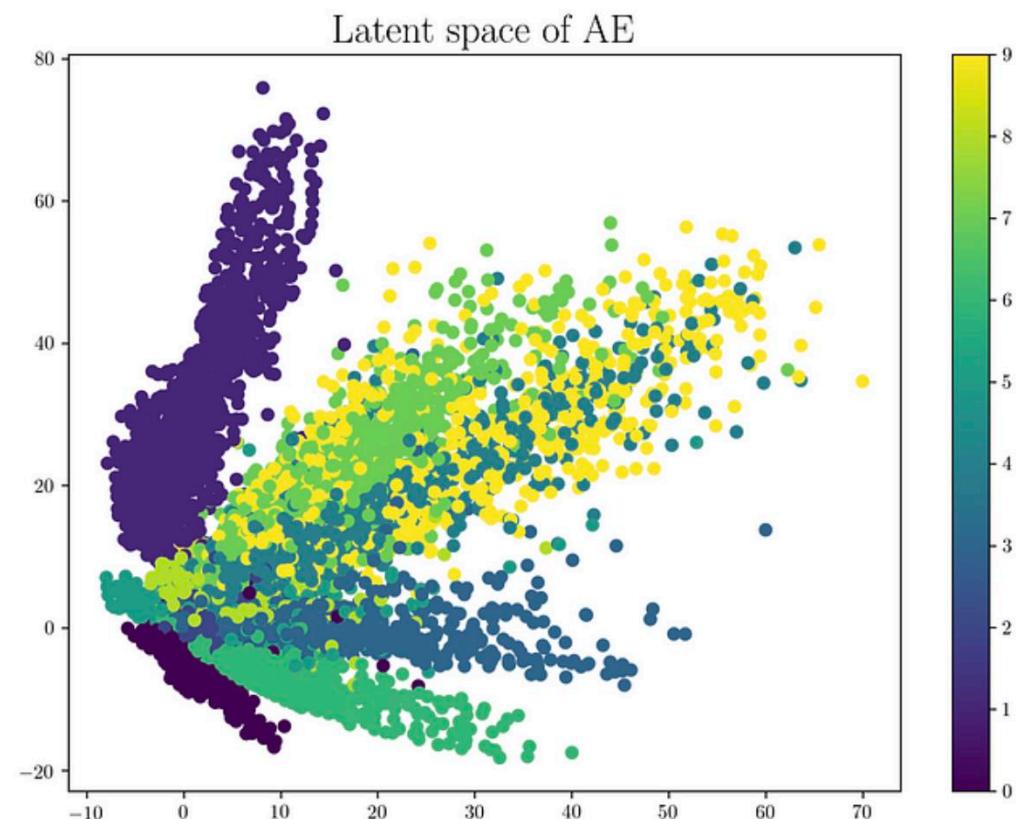


Figure: <https://towardsdatascience.com/difference-between-autoencoder-ae-and-variational-autoencoder-vae-ed7be1c038f2/>

Variational Autoencoder (VAE)

- VAE imposes a probabilistic structure on the latent space.
 - AE (simple reconstruction) → a generative model

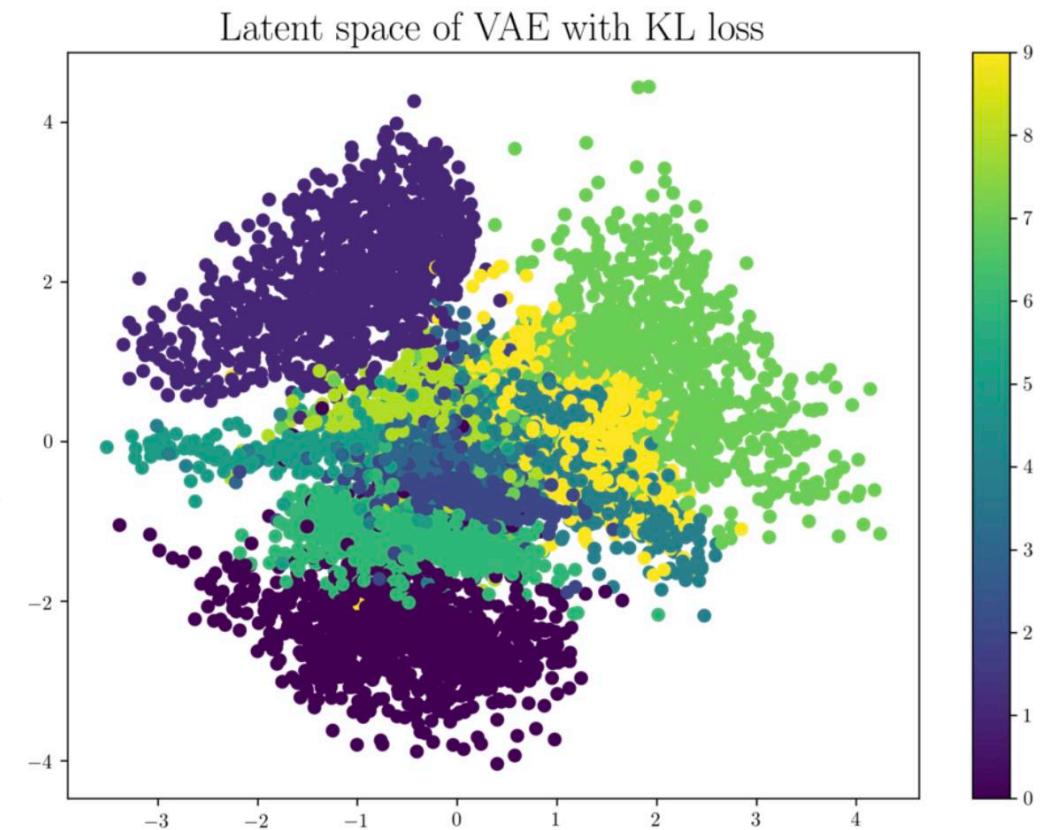
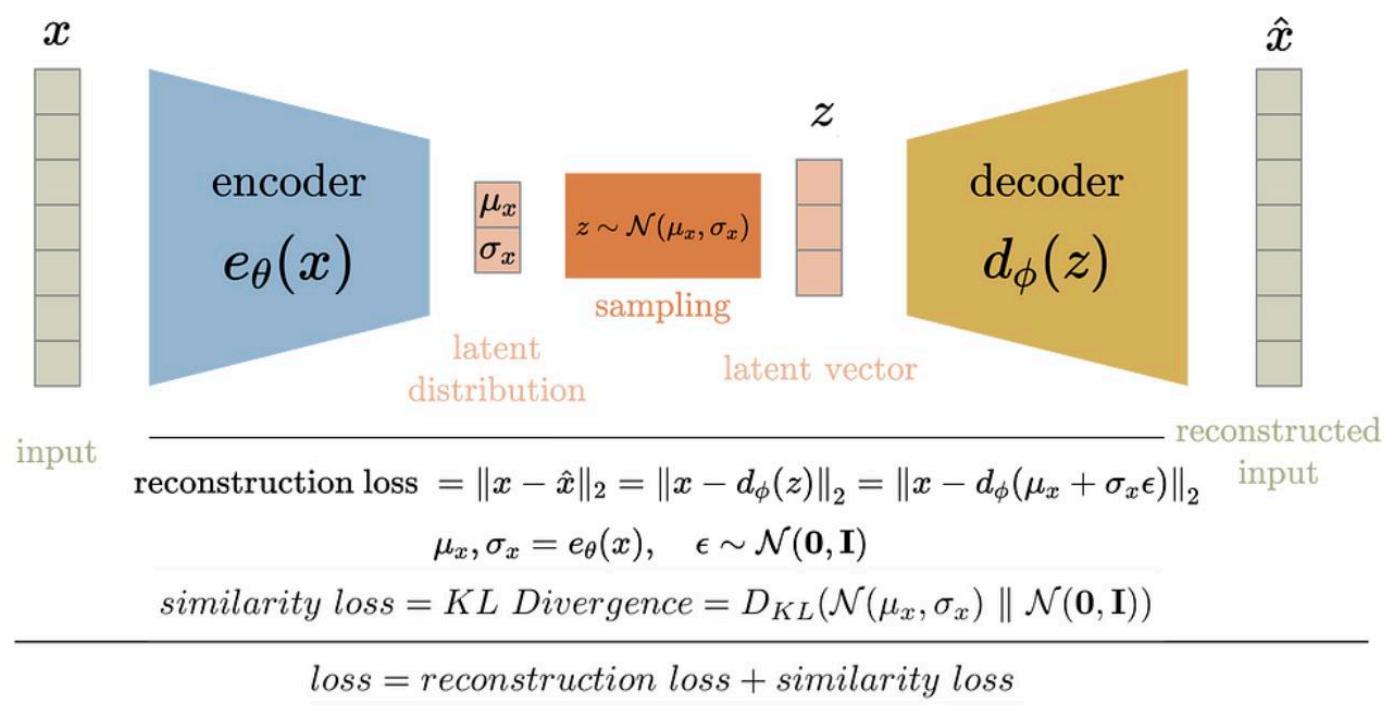
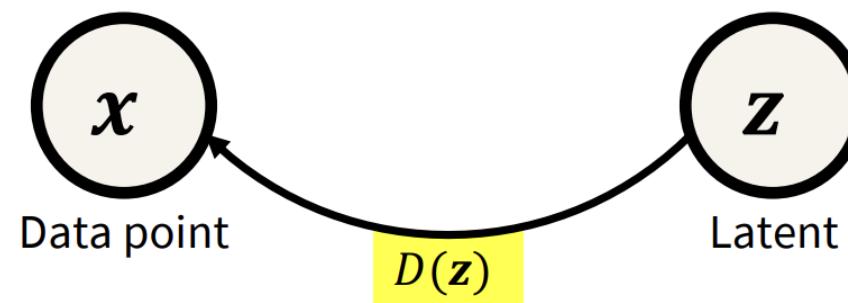


Figure: <https://towardsdatascience.com/difference-between-autoencoder-ae-and-variational-autoencoder-vae-ed7be1c038f2/>

Variational Autoencoder (VAE)

- Let's represent the mapping from the latent distribution $p(\mathbf{z})$ to the data distribution $p(\mathbf{x})$ as a conditional distribution $p(\mathbf{x}|\mathbf{z})$.
- Consider the decoder (generator) as predicting the mean of the conditional distribution $p(\mathbf{x}|\mathbf{z})$:

$$p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}; D(\mathbf{z}), \sigma^2 \mathbf{I})$$



(Appendix) Basics

- Marginal distribution
- Expected value
- Bayes' rule
- Kullback-Leibler (KL) Divergence
- Jensen's inequality

(Appendix) Marginal Distribution

- The marginal distribution of a subset of a set of random variables is the probability distribution of the variables contained in the subset.

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


Take the integral over \mathbf{z} to marginalize \mathbf{x} .

(Appendix) Marginal Distribution

- Q) Given the joint probability table below, what is $p(x = 1)$?

		y		
		1	2	3
x	1	0.32	0.03	0.01
	2	0.06	0.24	0.02
	3	0.02	0.03	0.27

(Appendix) Expected Value

- The expected value is the arithmetic mean of the possible values a random variable can take, weighted by the probability of those outcomes.

$$\mathbb{E}_{p(x)}[x] = \int x \cdot p(x) dx$$

(Appendix) Expected Value

- Q) What is $\mathbb{E}_{p(x)}[x]$?

		y		
		1	2	3
x	1	0.32	0.03	0.01
	2	0.06	0.24	0.02
	3	0.02	0.03	0.27

(Appendix) Bayes' Rule

- Bayes' rule is a mathematical formula used to determine the conditional probability of events.

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Posterior

Likelihood

Prior

Marginal

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

(Appendix) Bayes' Rule

- Q) What is $p(y = 2|x = 1)$?

		y		
		1	2	3
x	1	0.32	0.03	0.01
	2	0.06	0.24	0.02
	3	0.02	0.03	0.27

(Appendix) Kullback-Leibler Divergence

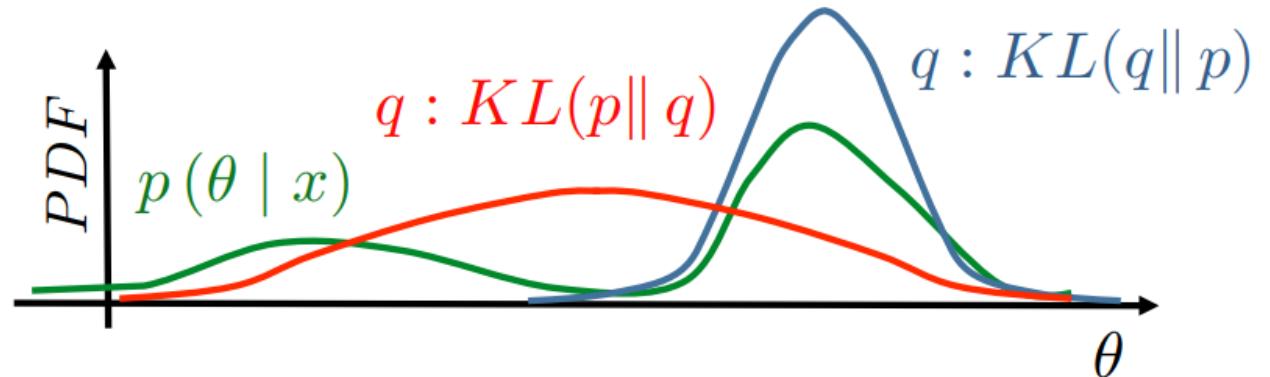
- Kullback-Leibler (KL) divergence is a measure of how one probability distribution p is different from a reference probability distribution q :

$$D_{\text{KL}}(p \parallel q) = \int p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) d\mathbf{x} = \mathbb{E}_{p(\mathbf{x})} \left[\log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right]$$

(Appendix) Kullback-Leibler Divergence

- Properties

- $D_{\text{KL}}(p \parallel q) \geq 0$
- $D_{\text{KL}}(p \parallel q) = 0 \Leftrightarrow p(x) = q(x)$
- $D_{\text{KL}}(p \parallel q) \neq D_{\text{KL}}(q \parallel p)$



- In general, $D_{\text{KL}}(p \parallel q)$ is called forward KL, whereas $D_{\text{KL}}(q \parallel p)$ is called reverse KL.

- Forward KL is known as zero avoiding while reverse KL is known as zero forcing.

$\overbrace{\min_{\hat{p}} \text{KL}}$ $\overbrace{\log\left(\frac{p}{q}\right)}$
 $g \rightarrow 0 \quad \log \rightarrow \infty.$

$\overbrace{\min_{\hat{p}} \text{KL}}$ $\overbrace{\log\left(\frac{q}{p}\right)}$
 $g \rightarrow 0 \quad \text{Ignore } p.$

(Appendix) Kullback-Leibler Divergence

- Q)

When $p(\mathbf{x}) = N(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ and $q(\mathbf{x}) = N(\mathbf{x}; \mathbf{0}, \mathbf{I})$, what is $D_{\text{KL}}(p || q)$?

If $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_p, \sigma_p^2 \mathbf{I})$, $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_q, \sigma_q^2 \mathbf{I})$, $\mathbf{x} \in \mathbb{R}^d$,

$$D_{\text{KL}}(p \parallel q) = \frac{1}{2} \left[d \left(\frac{\sigma_p^2}{\sigma_q^2} - 1 - \ln \frac{\sigma_p^2}{\sigma_q^2} \right) + \frac{\|\boldsymbol{\mu}_p - \boldsymbol{\mu}_q\|^2}{\sigma_q^2} \right].$$

When $\boldsymbol{\mu}_q = \mathbf{0}$, $\sigma_q = 1$,

$$D_{\text{KL}}(p \parallel q) = \frac{1}{2} \left[d(\sigma_p^2 - 1 - \ln \sigma_p^2) + \|\boldsymbol{\mu}_p\|^2 \right]$$

(Appendix) Jensen's Inequality

- f is a convex function if

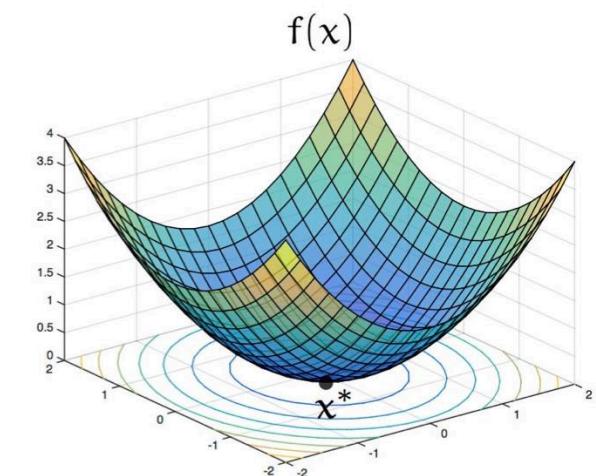
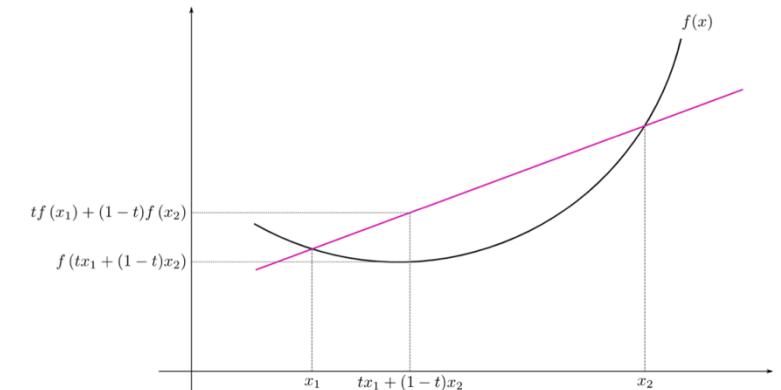
$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

convex combination

for all x_1, x_2 and $t \in [0,1]$

$$f\left(\sum_i t_i x_i\right) \leq \sum_i t_i f(x_i)$$

for all x_i such that $t_i \geq 0$ and $\sum_i t_i = 1$



(Appendix) Jensen's Inequality

- If x is a random variable and f is a convex function, then

$$f(\mathbb{E}_{p(x)}[x]) \leq \mathbb{E}_{p(x)}[f(x)]$$

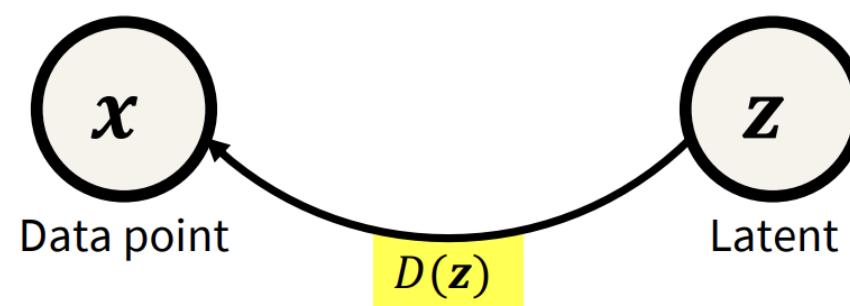
since the expected value is a convex combination.

- Q) What if f is a concave function, such as \log ?

Back to Variational Autoencoder (VAE)

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

$$p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}; D(\mathbf{z}), \sigma^2 \mathbf{I})$$



Variational Autoencoder (VAE)

- For all given real images x , we want to maximize the marginal probability:

$$p(x) = \int p(x, z) dz = \int p(x|z)p(z) dz$$

- How to compute the integral?
 - Monte-Carlo method for x and z takes too much time. → Intractable

Variational Autoencoder (VAE)

- Or, we can compute

$$p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\boxed{p(\mathbf{z}|\mathbf{x})}}$$

But this conditional distribution
is unknown.

- We cannot directly maximize $p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})}$ since $p(\mathbf{z}|\mathbf{x})$ is unknown.

Variational Autoencoder (VAE)

- Let's think about the lower bound of $\log p(\mathbf{x})$:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

- $q_\phi(\mathbf{z}|\mathbf{x})$ is a variational distribution with parameters ϕ .
- E.g., a Gaussian distribution with mean and variance as parameters
- Consider $q_\phi(\mathbf{z}|\mathbf{x})$ as an arbitrary conditional distribution that may not be the same with $p(\mathbf{z}|\mathbf{x})$.
- We use the proxy distribution $q_\phi(\mathbf{z}|\mathbf{x})$ since we do not know $p(\mathbf{z}|\mathbf{x})$.

Variational Autoencoder (VAE)

- $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$ is the expected value over \mathbf{z} sampled from the variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$.
- Why $\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$?
 - Because of the Jensen's inequality.

Evidence Lower Bound (ELBO)

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

$$= \log \int p(\mathbf{x}, \mathbf{z}) \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

Concave function

$$\begin{aligned}&= \boxed{\log} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\&\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$

Evidence Lower Bound (ELBO)

- Let's decompose ELBO:

$$\begin{aligned}\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} \right] \\ &= \boxed{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})]} - \boxed{D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}) \right)}\end{aligned}$$

Reconstruction term
to be *maximized*.

Prior matching term
to be *minimized*.

Evidence Lower Bound (ELBO)

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

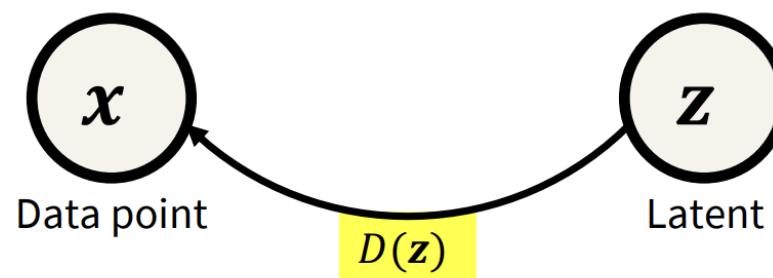
- Reconstruction
 - Accurate recovery of \mathbf{x} from its latent code \mathbf{z}
 - With Gaussian assumptions, it reduces the familiar reconstruction loss of AE.
- Latent KL
 - Encourages $q_{\phi}(\mathbf{z}|\mathbf{x})$ to stay close to (a simple Gaussian) prior $p(\mathbf{z})$.
 - Shapes the latent space into a smooth and continuous structure.

Variational Autoencoder (VAE)

- In VAE, we want model $p(\mathbf{x})$ as

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

where $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D(\mathbf{z}), \sigma^2 I)$ and $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.

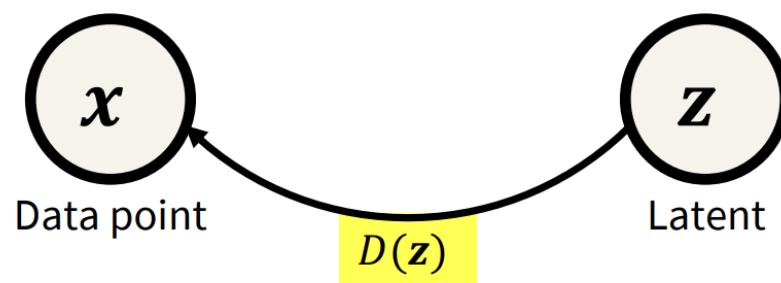


Variational Autoencoder (VAE)

- In VAE, we want model $p(\mathbf{x})$ as

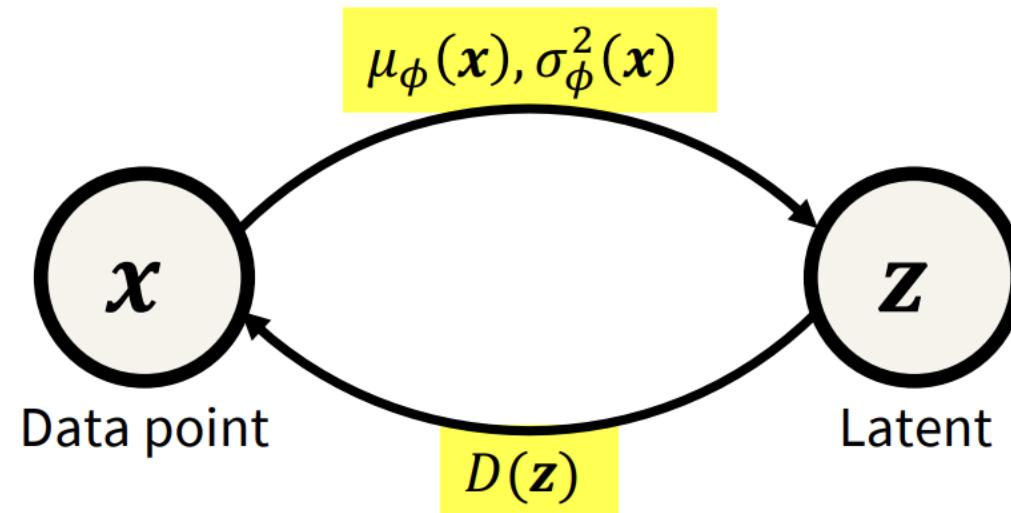
$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

where $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D_{\theta}(\mathbf{z}), \sigma^2 I)$ and $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.



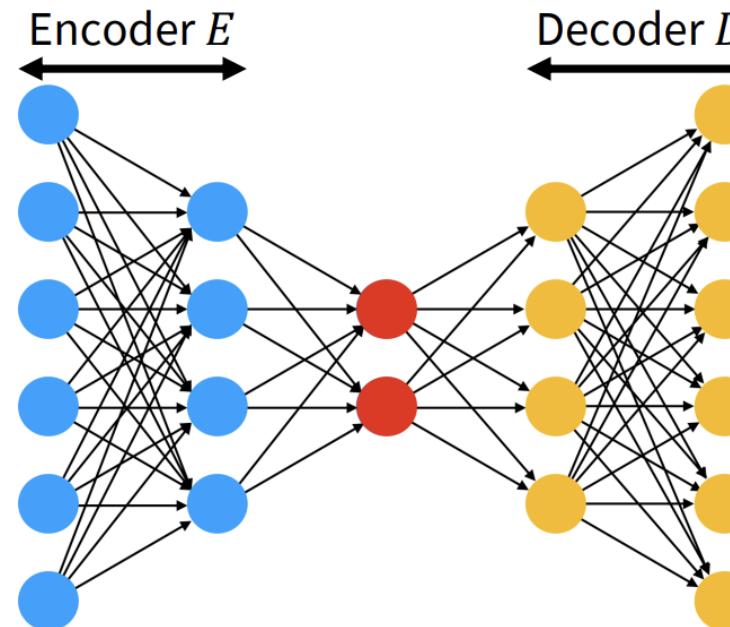
Evidence Lower Bound (ELBO)

- How to maximize $\log p(x)$? Maximize ELBO!
- We need the proxy distribution $q_\phi(z|x) \rightarrow$ Encoder
- Let $q_\phi(z|x) = N(z; \mu_\phi(x), \sigma_\phi^2(x)\mathbf{I})$



Variational Autoencoder (VAE)

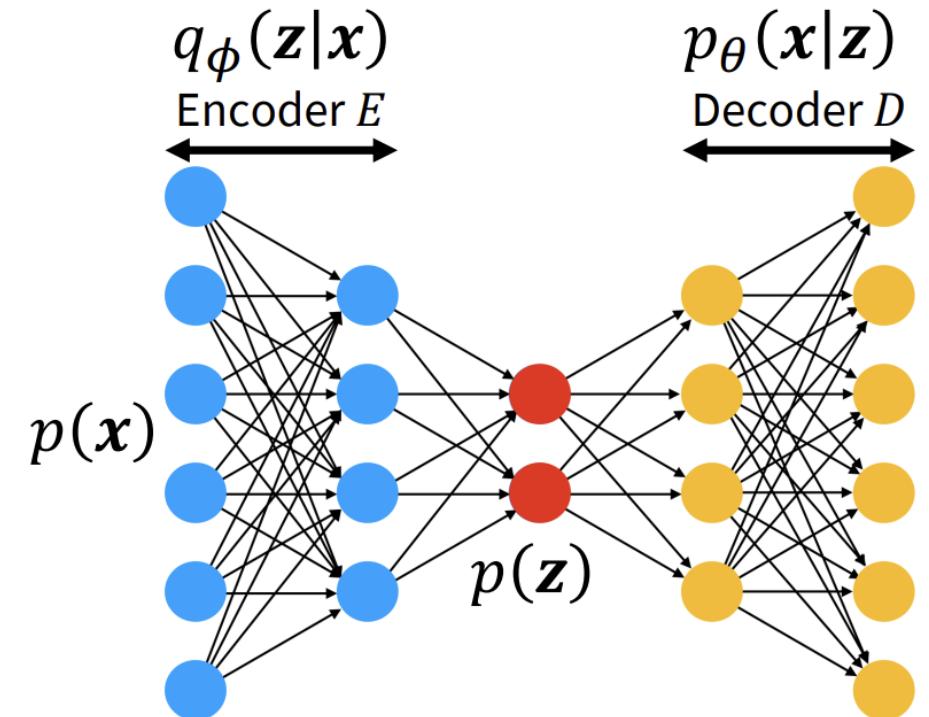
- Same with autoencoder but,
- From input x , the encoder predicts $\mu_\phi(x), \sigma_\phi^2(x)\mathbf{I}$.
- The decoder takes a sample $\mathbf{z} \sim N(\mathbf{z}; \mu_\phi(x), \sigma_\phi^2(x)\mathbf{I})$ as input.



Variational Autoencoder (VAE)

- Summary

Data distribution	$p(\mathbf{x})$
Encoder	$q_{\phi}(\mathbf{z} \mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x})\mathbf{I})$
Latent distribution	$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$
Decoder	$p_{\theta}(\mathbf{x} \mathbf{z}) = \mathcal{N}(\mathbf{x}; D_{\theta}(\mathbf{z}), \sigma^2 I)$



VAE Training

- How to maximize ELBO?

$$\operatorname{argmax}_{\theta, \phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

- Approximates using a Monte Carlo estimate:

$$\operatorname{argmax}_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

where $\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$ for the given \mathbf{x}

VAE Training

- How to maximize ELBO?

$$\operatorname{argmax}_{\theta, \phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

- Approximates using a Monte Carlo estimate:

$$\operatorname{argmax}_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Reparameterization trick

where $\mathbf{z}^{(i)} = \boldsymbol{\mu}_{\phi}(\mathbf{x}) + \sigma_{\phi}(\mathbf{x})\boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I})$

VAE Training

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

$$\log p_\theta(\mathbf{x}|\mathbf{z}^{(i)}) = \log \left(\frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp \left(-\frac{\|\mathbf{x} - D_\theta(\mathbf{z})\|^2}{2\sigma^2} \right) \right)$$

$$= -\frac{1}{2\sigma^2} \boxed{\|\mathbf{x} - D_\theta(\mathbf{z})\|^2} - \log \sqrt{(2\pi\sigma^2)^d}$$

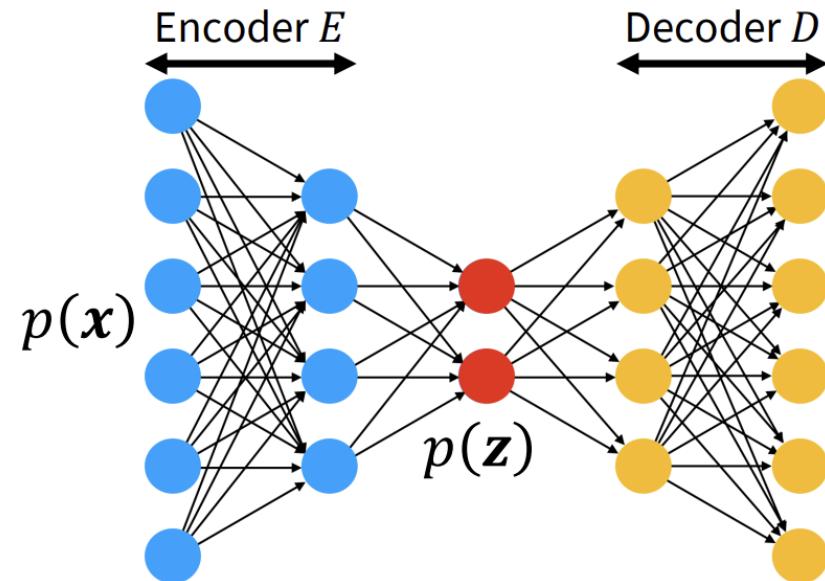
This is why it is called
the reconstruction term.

Constant

VAE Training

1. Feed a data point x to the encoder to predict and $\mu_\phi(x)$ and $\sigma_\phi^2(x)$.
2. Sample a latent variable z from $q_\phi(z|x) = N(z; \mu_\phi(x), \sigma_\phi^2(x)\mathbf{I})$.
3. Feed z to the decoder to predict $\hat{x} = D_\theta(z)$.
4. Perform the gradient descent through the negative ELBO.

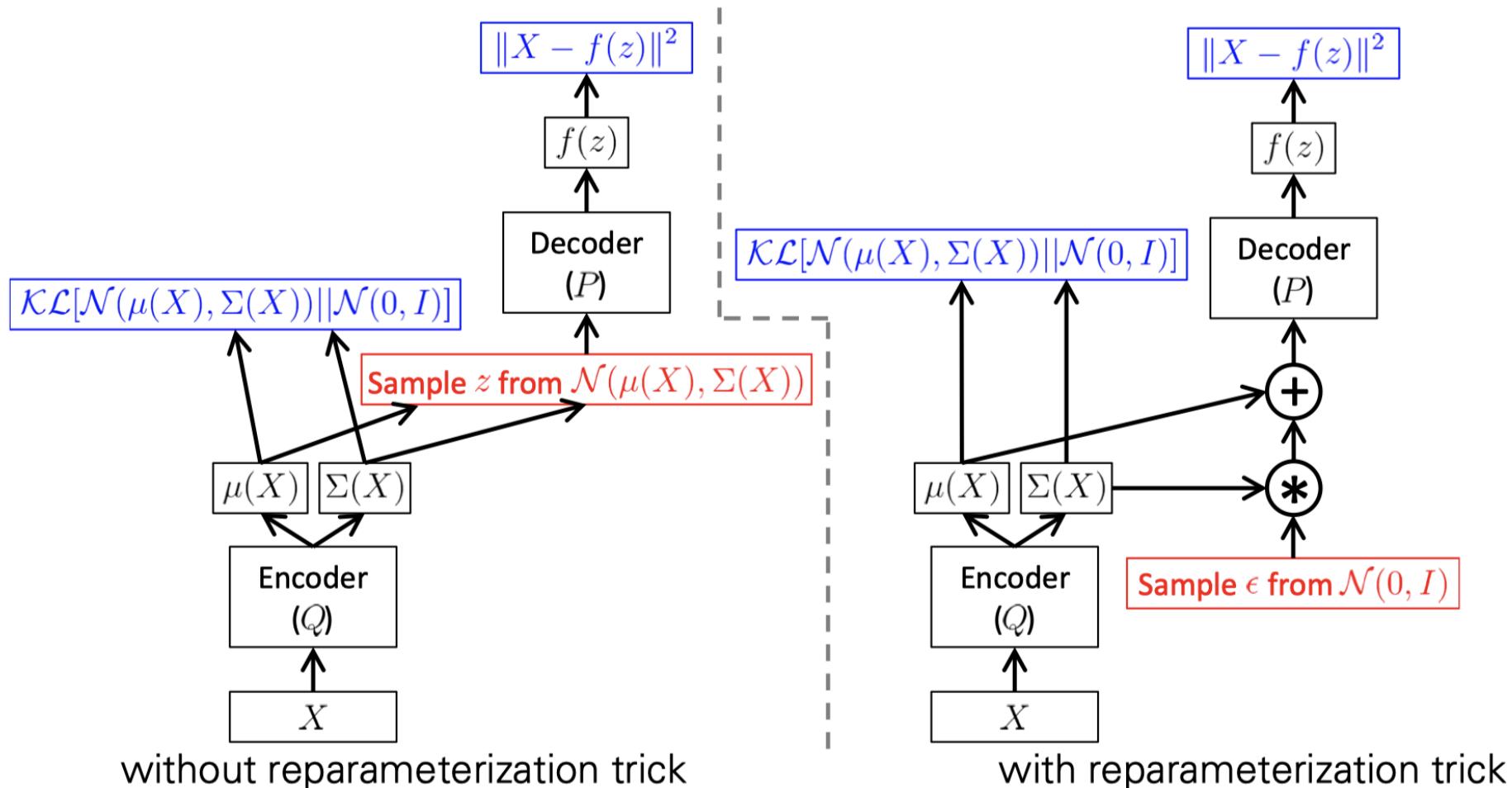
Q) Why is the sampling differentiable?



VAE Training

$$z = \mu + \sigma \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

Q) Why is the sampling differentiable?

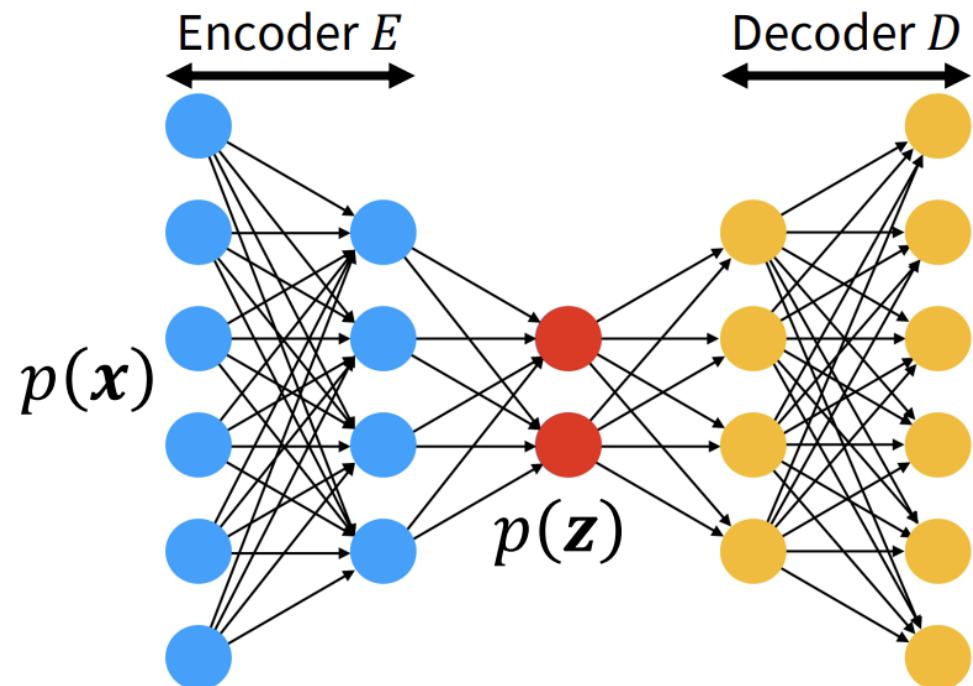


without reparameterization trick

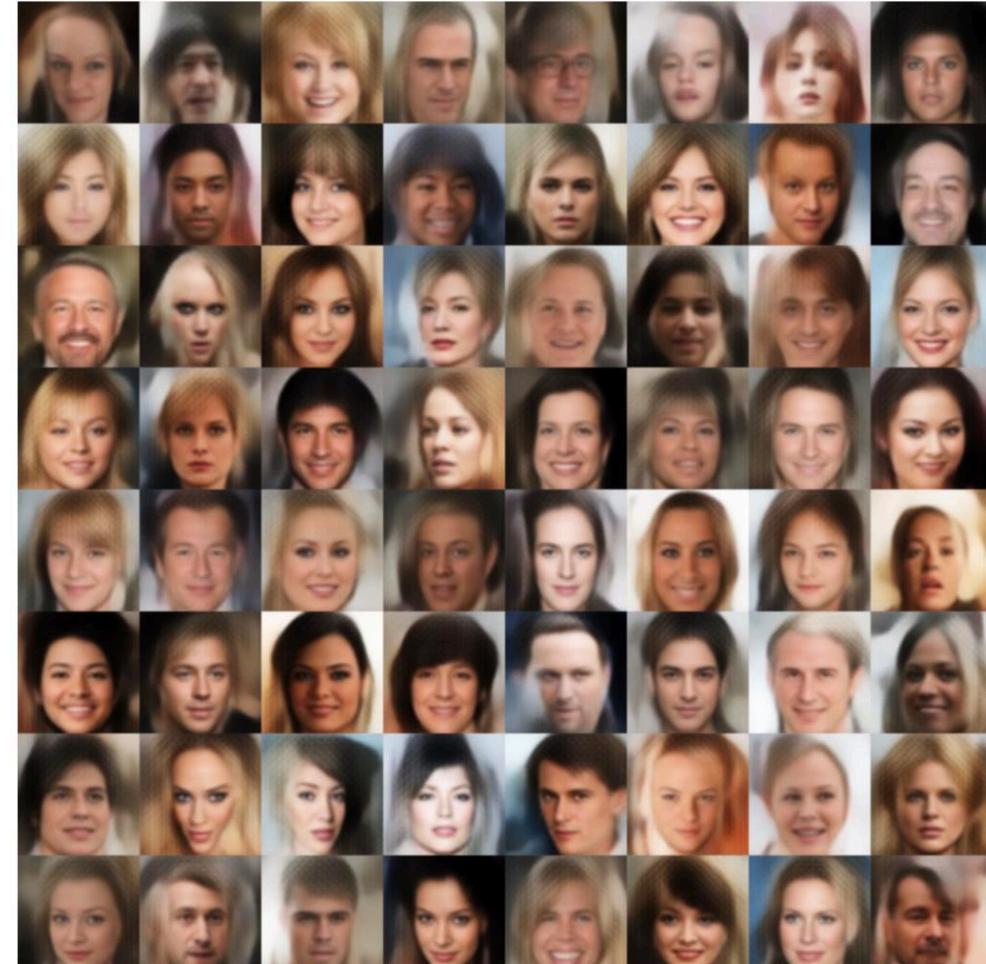
with reparameterization trick

VAE Generation

1. Sample a latent variable \mathbf{z} from $p(\mathbf{z}) = N(\mathbf{z}; \mathbf{0}, \mathbf{I})$.
2. Feed \mathbf{z} to the decoder to predict $\hat{\mathbf{x}} = D_\theta(\mathbf{z})$.



Samples from VAEs

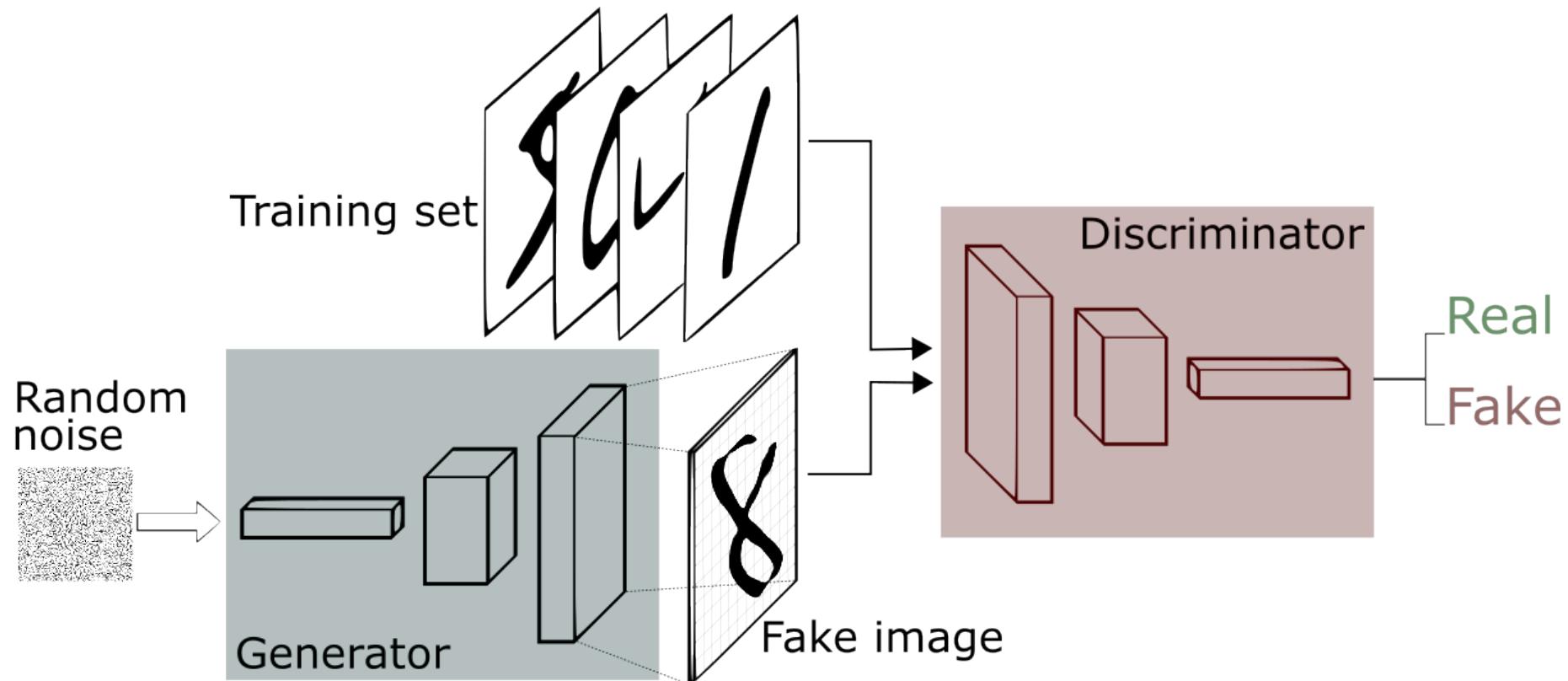


VAE vs AE diff Eq.

Generative Adversarial Networks

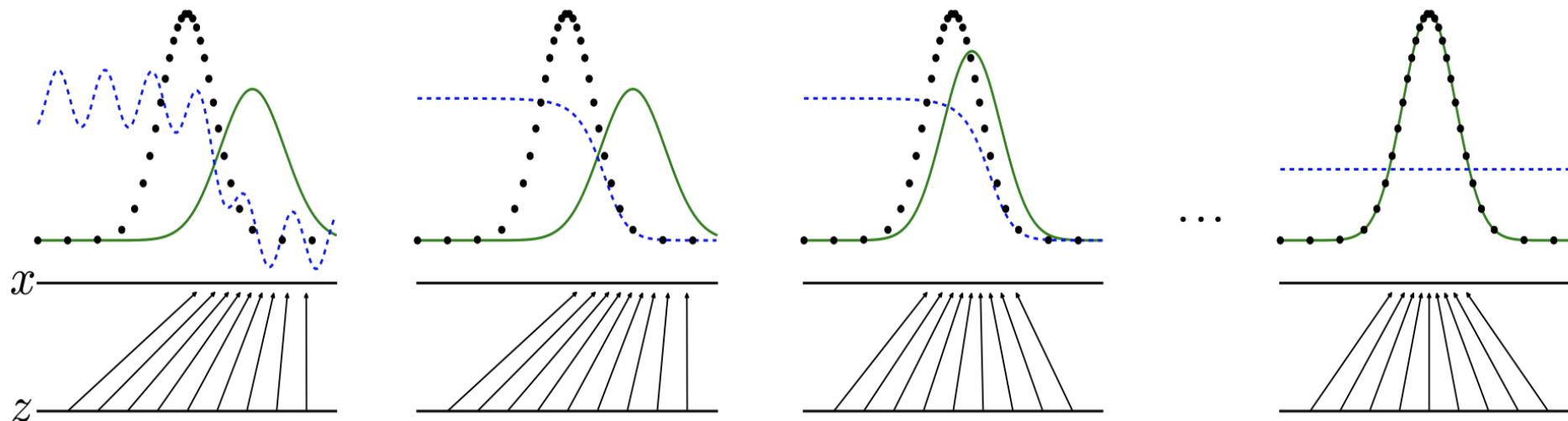
Generative Adversarial Networks

- Two player minimax game



Training Objective of GAN

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

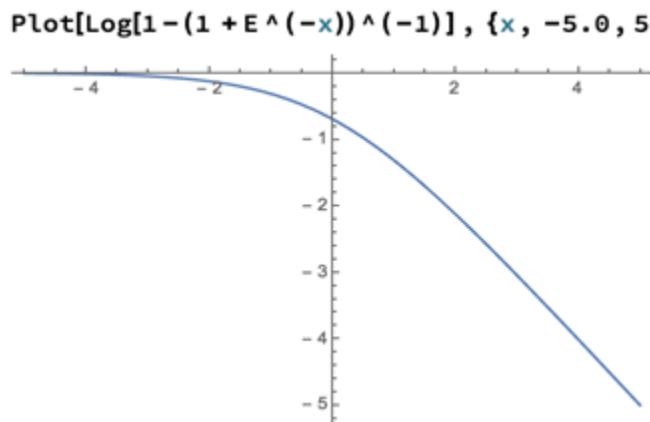


Training GAN

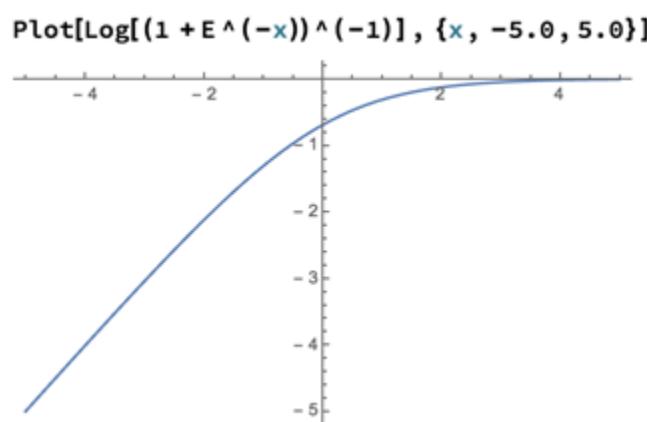
- Heuristics for training

In practice, equation 1 may not provide sufficient gradient for G to learn well. Early in learning, when G is poor, D can reject samples with high confidence because they are clearly different from the training data. In this case, $\log(1 - D(G(z)))$ saturates. Rather than training G to minimize $\log(1 - D(G(z)))$ we can train G to maximize $\log D(G(z))$. This objective function results in the same fixed point of the dynamics of G and D but provides much stronger gradients early in learning.

$\log(1 - x)$



$\log(x)$



- Instead of flipping the sign on $V(D,G)$, we flip the target!
- Not minimax game
- Theoretical analysis does not hold

Training GAN

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

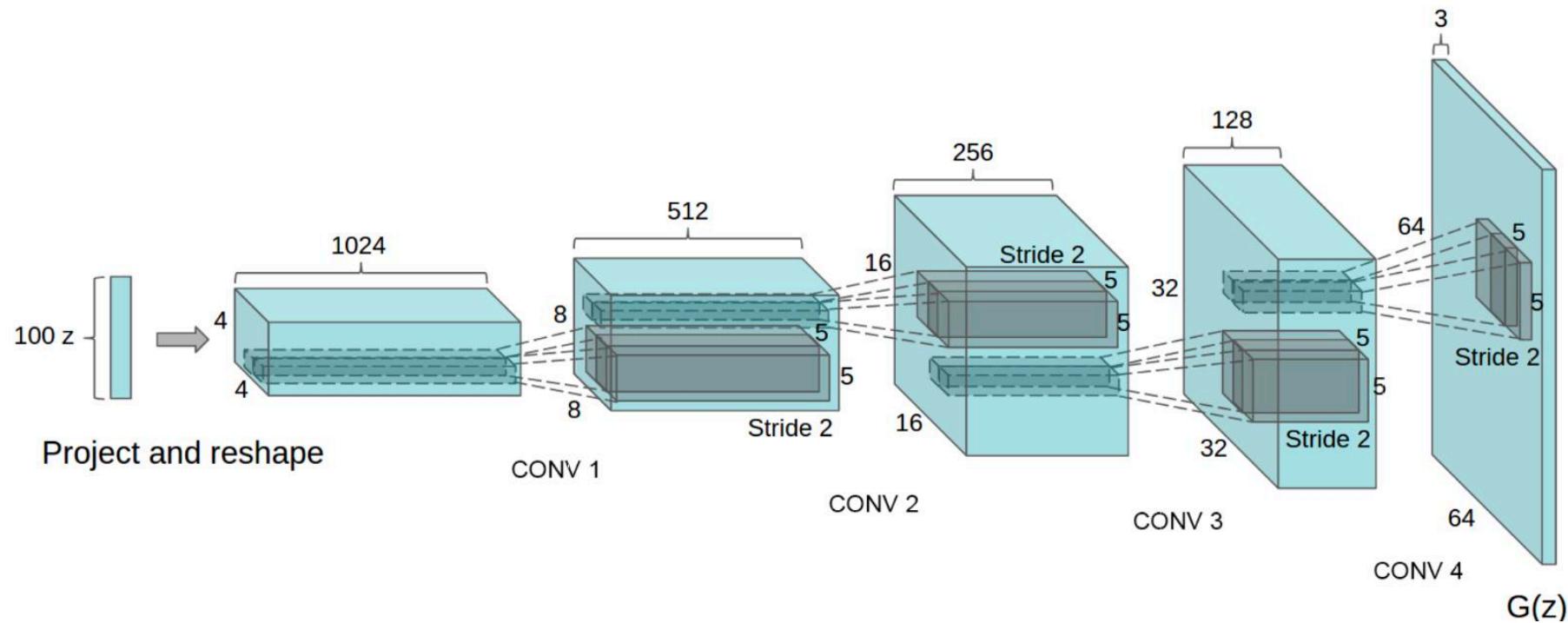
Note that this algorithm does not work in practice, but for the theoretical analysis.

end for

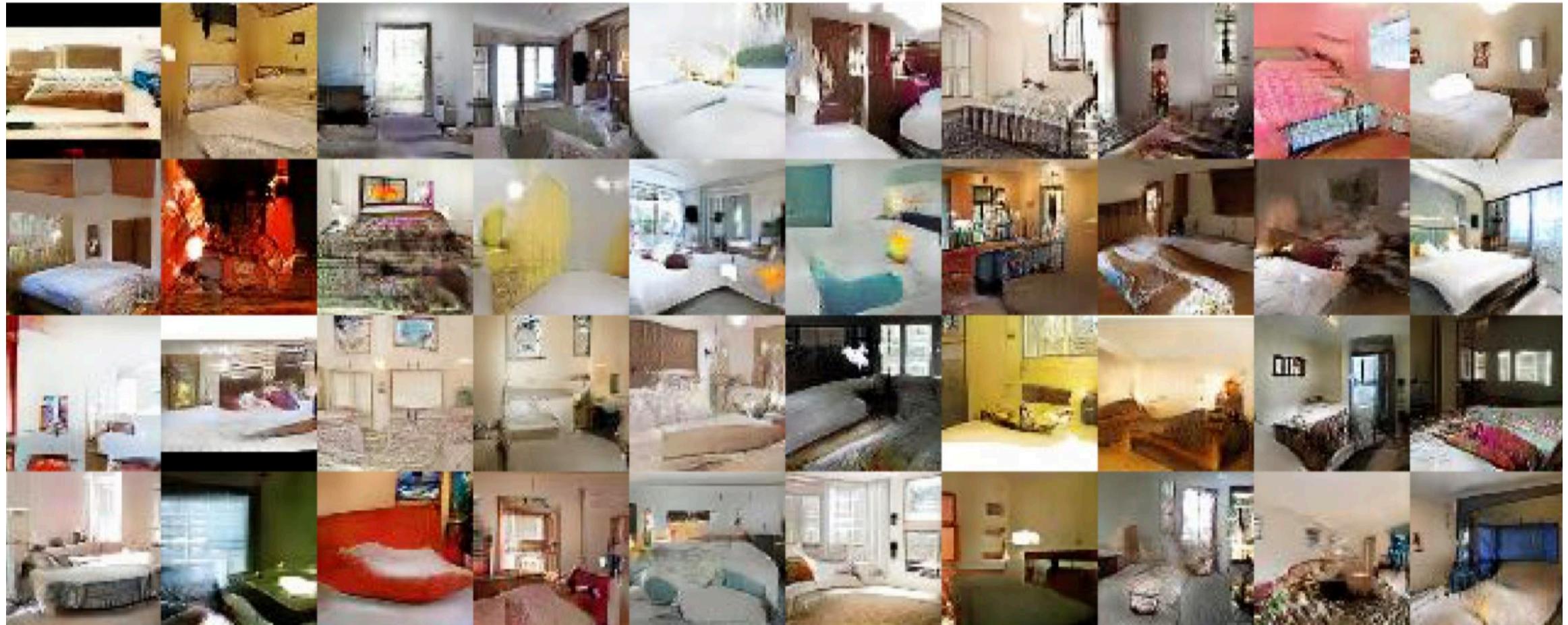
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Deep Convolutional GAN

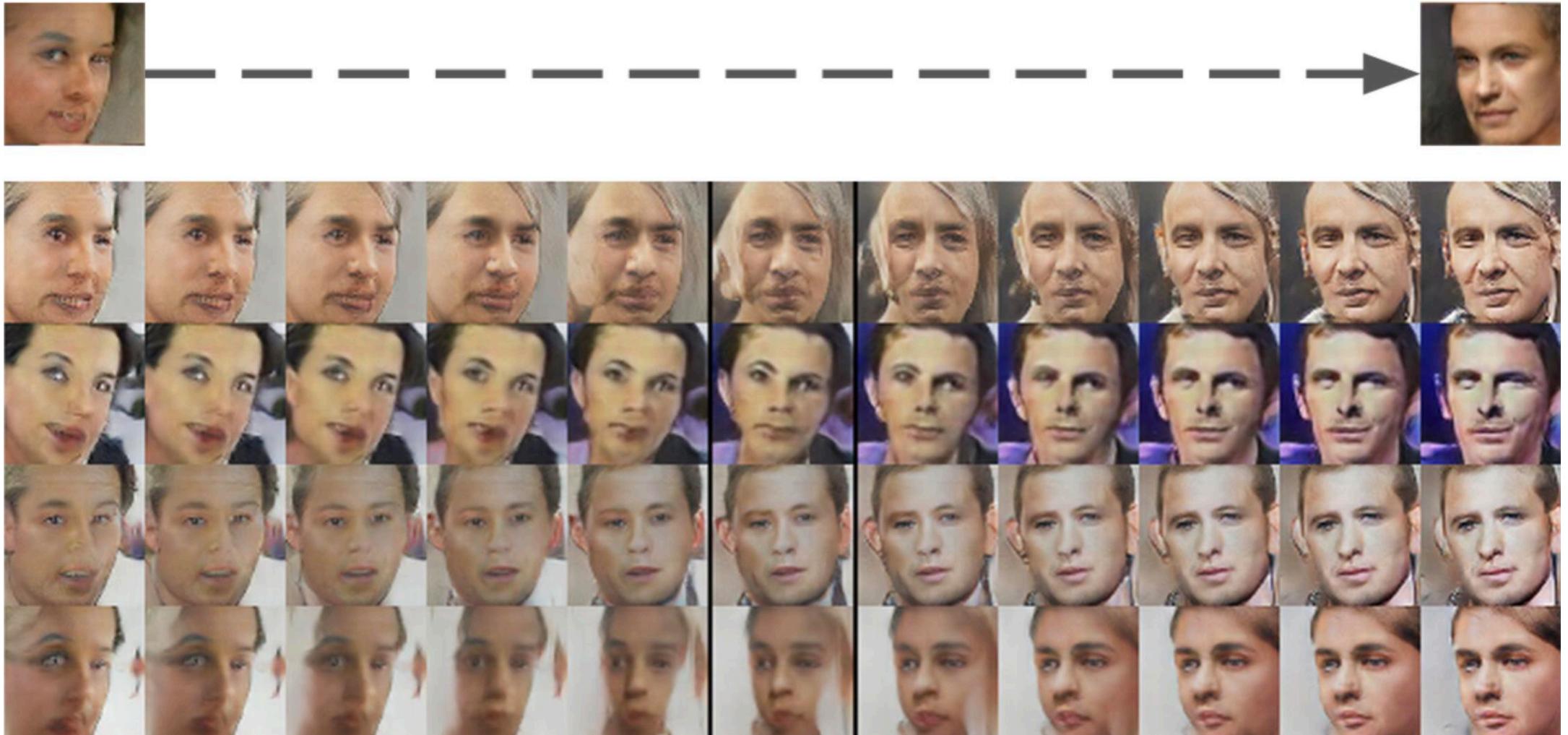
- Deep Convolutional GAN (Radford et al., 2015)
 - introduces some tips to make GAN training more stable
 - fully convolutional layers, leaky ReLU in D, Adam optimizer, batch norm, etc.



Deep Convolutional GAN (Radford et al., 2015)



Deep Convolutional GAN (Radford et al., 2015)



Deep Convolutional GAN (Radford et al., 2015)

