

Image and Video Generations

Lecture 1: Introduction to Generative Models: GAN & VAE

劉育綸

Yu-Lun (Alex) Liu

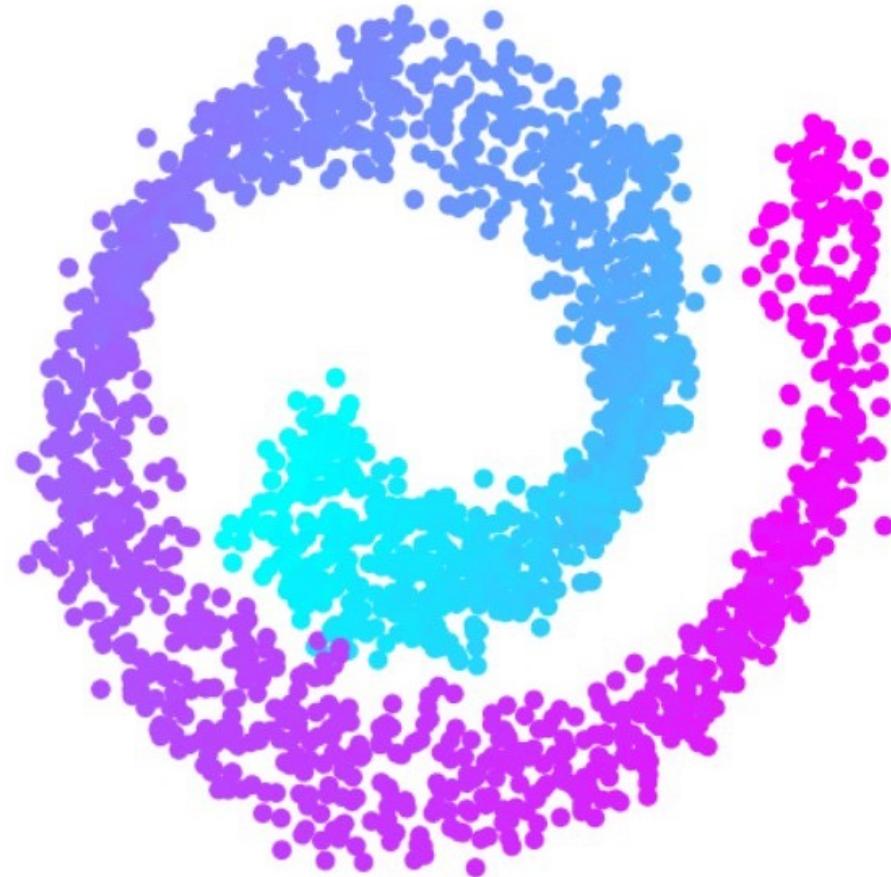
Let's consider a collection of real photos.



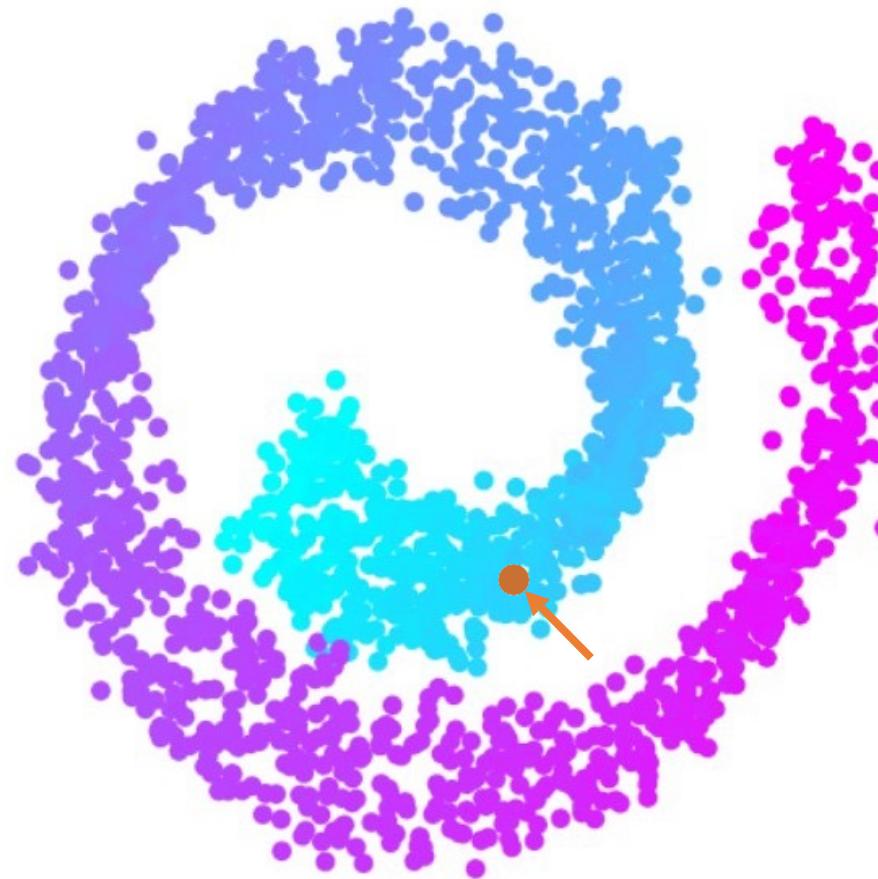
How can we generate a **new** realistic photo?



Let's consider a **simpler** case:
a collection of **2D points**.



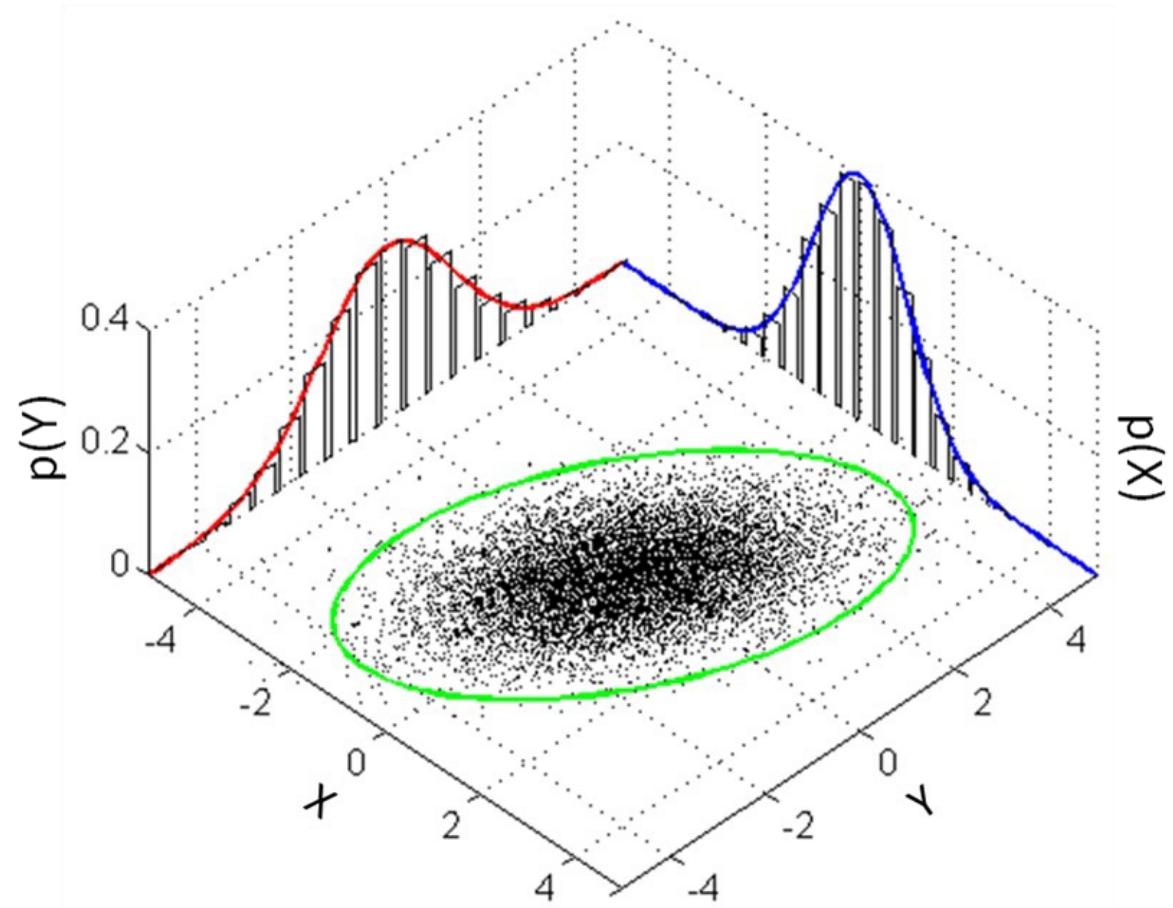
How can we sample
a new plausible 2D point?



Statistical Perspective

From a **statistical perspective**,
we will view this as

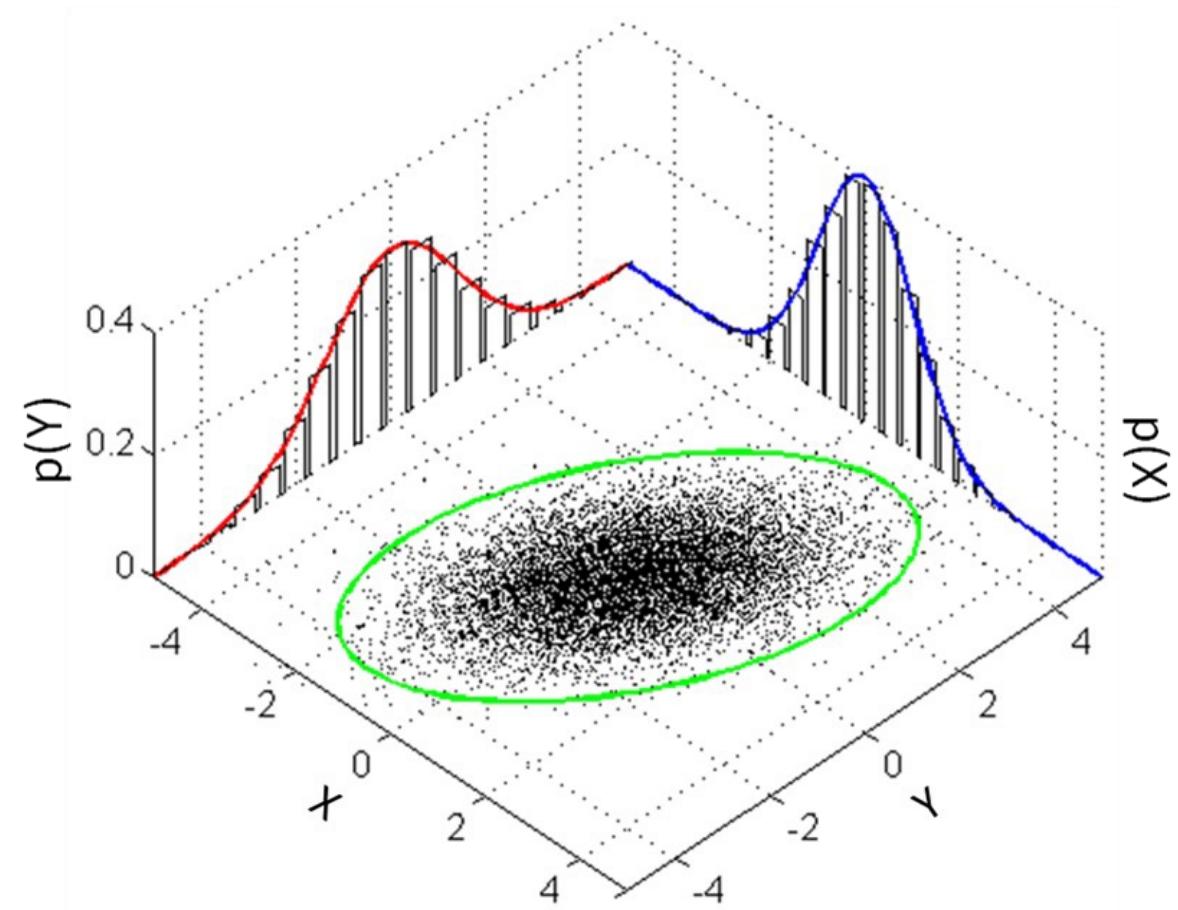
- there being a **probability distribution** of the data, and
- the given points are **samples** from the probability distribution.



Statistical Perspective

If the probability density function (PDF) of the distribution is known, we can sample from it directly.

How?

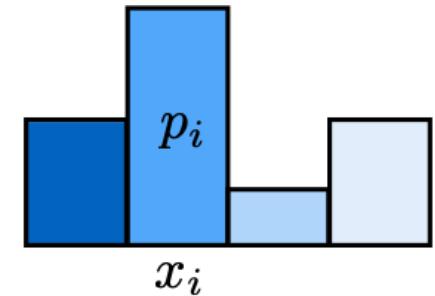


Discrete Probability Distributions

- n discrete values x_i with probability p_i .
- Requirements of a PDF:

$$p_i \geq 0$$

$$\sum_{i=1}^n p_i = 1$$



Q. How can we sample based on this PDF?

Cumulative Distribution Function (CDF)

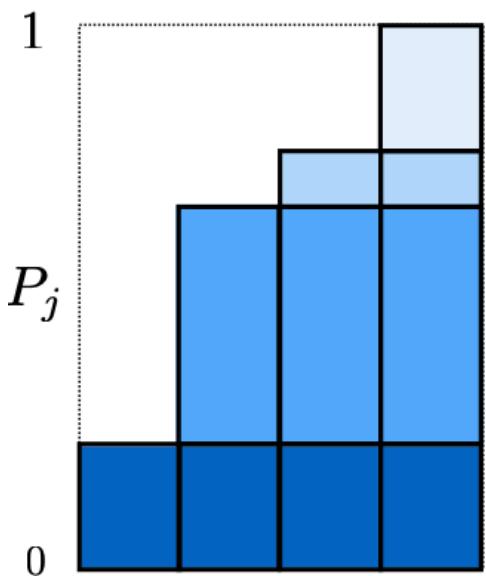
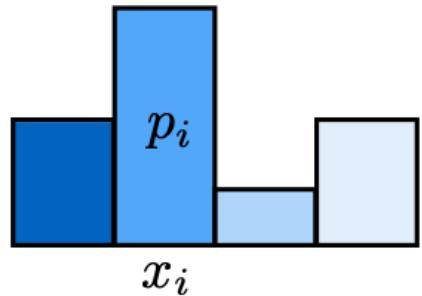
Cumulative distribution function:

$$P_i = \sum_{j=1}^i p_j$$

where

$$0 \leq P_i \leq 1$$

$$P_n = 1$$



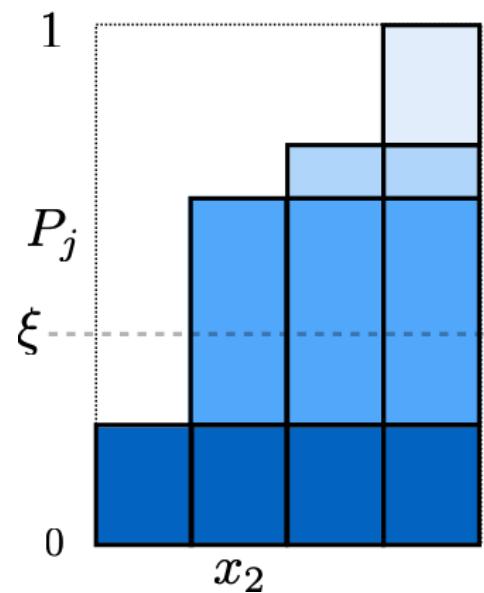
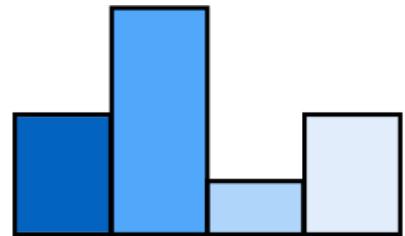
Inverse Transform Sampling

To randomly select an event,

- Draw $u \sim \mathcal{U}(0, 1)$.

- Take x_i if

$$P_{i-1} \leq u \leq P_i$$



Continuous Probability Distributions

Given a PDF $p(x)$,

- CDF $F_X(x)$

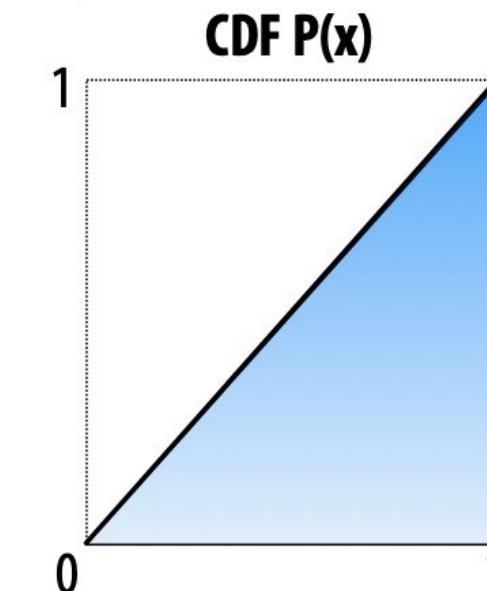
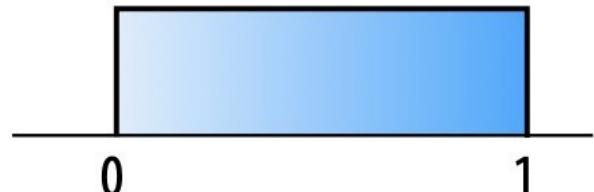
$$= \Pr(X \leq x) = \int_0^x p(t)dt$$

- $\Pr(a < X \leq b)$

$$= \int_a^b p(t)dt = F_X(b) - F_X(a)$$

- $F_X(1) = 1$

Uniform distribution: $p(x) = c$
(for random variable X defined on $[0,1]$ domain)

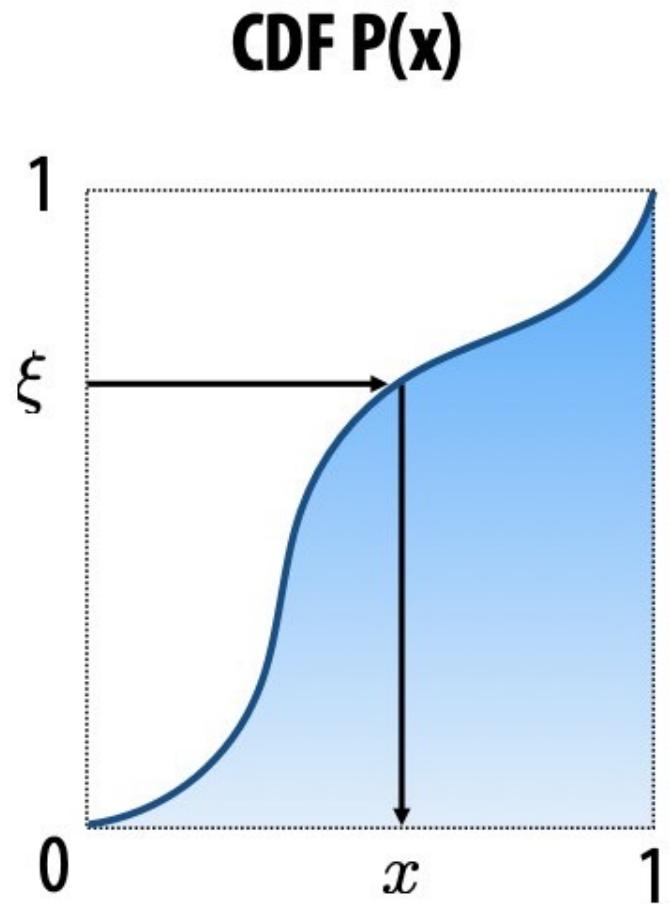


Inverse Transform Sampling

To randomly sample based on the given PDF,

- Compute **CDF** $F_X(x)$
- Draw $u \sim \mathcal{U}(0, 1)$.
- Take $x = F_X^{-1}(u)$.

Need to know the inverse function $F_X^{-1}(x)$.

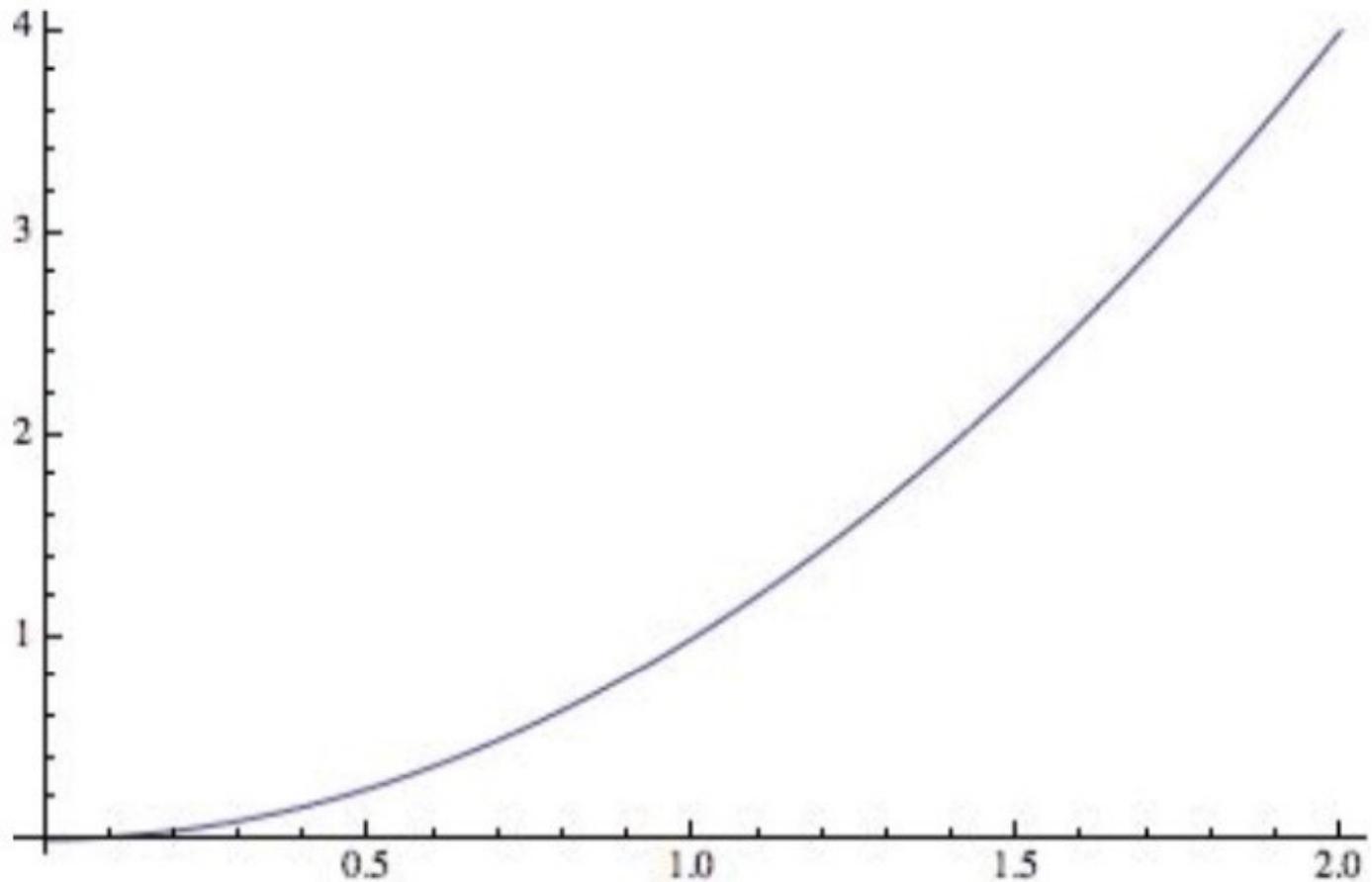


Inverse Transform Sampling – Example

PDF

$$p(x) = \frac{3}{8}x^2 \quad x \in [0,2]$$

Q. What is the inverse of CDF?



Inverse Transform Sampling – Example

PDF

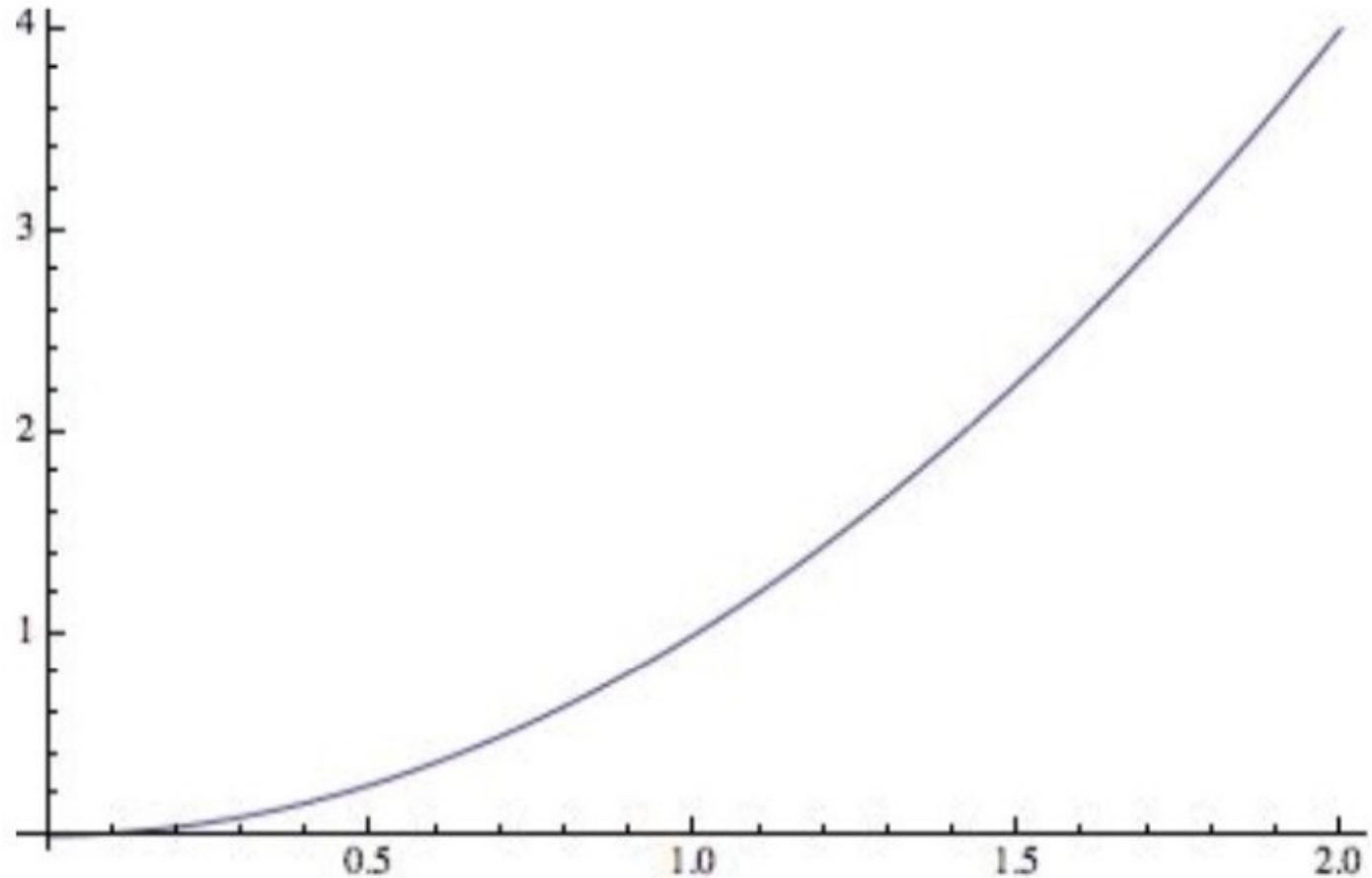
$$p(x) = \frac{3}{8}x^2 \quad x \in [0,2]$$

CDF

$$\begin{aligned} F_X(x) &= \int_0^x p(x)dx \\ &= \frac{1}{8}x^3 \end{aligned}$$

Inverse of CDF

$$F_X^{-1}(x) = 2\sqrt[3]{x}$$

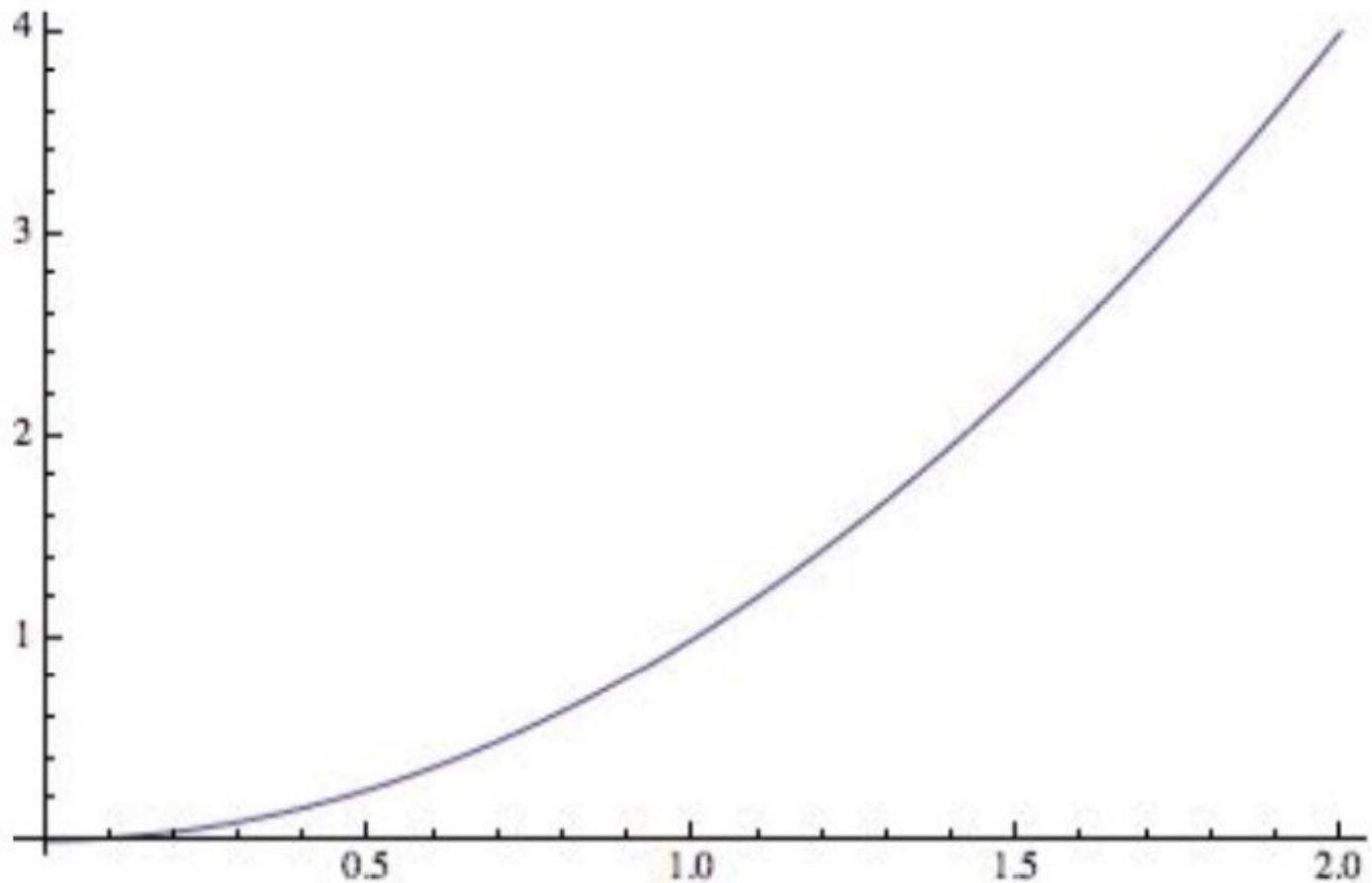


Inverse Transform Sampling – Example

Sampling from

$$p(x) = \frac{3}{8}x^2$$

1. Draw a sample $u \sim U(0, 1)$.
2. Take $2\sqrt[3]{x}$.

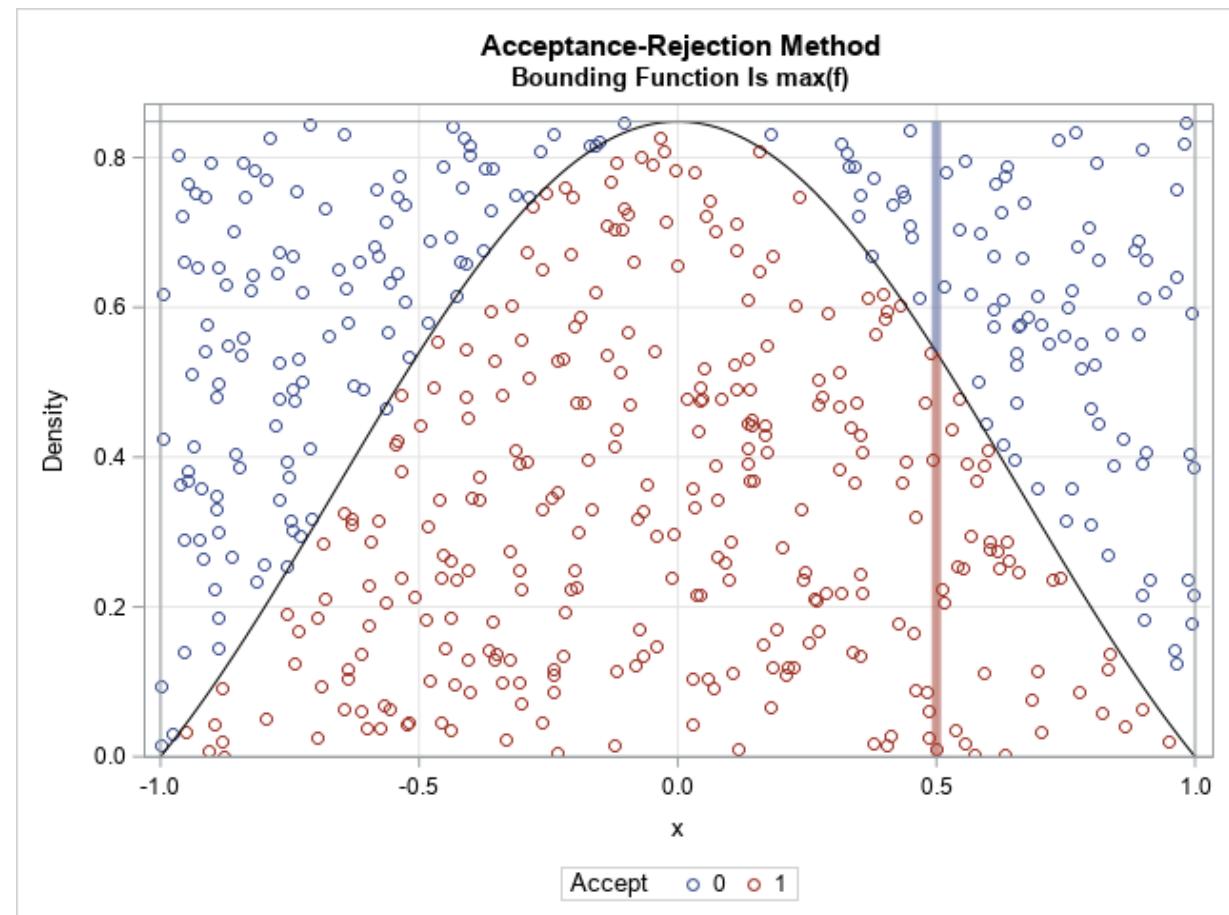


Rejection Sampling

If the inverse of the CDF cannot be computed (CDF 有時候很難積分、inverse function 沒有解析解):

1. Let $q(x)$ be an upper bound distribution: $\forall x \ q(x) \geq p(x)$.
2. Draw $x \sim q(x)$.
3. Draw a $h \sim \mathcal{U}(0, q(x))$.
4. Accept the sample x if $h \leq p(x)$; otherwise reject it.

試了再決定要不要



Reparameterization Trick

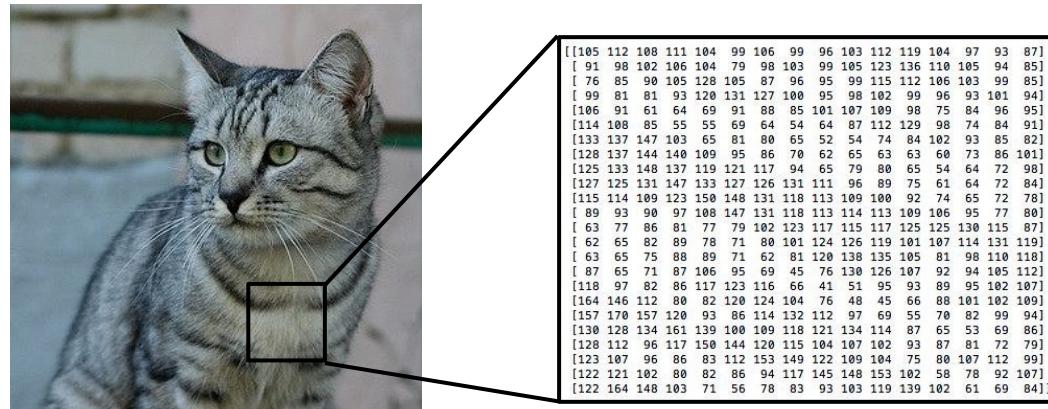
- A sample from a normal distribution $z \sim N(\mu, \Sigma)$ can be rewritten as follows:

$$z = \mu + \Sigma^{\frac{1}{2}}\epsilon \text{ where } \epsilon \sim N(0, I).$$

- We just need a standard normal sampler to sample from an arbitrary normal distribution.

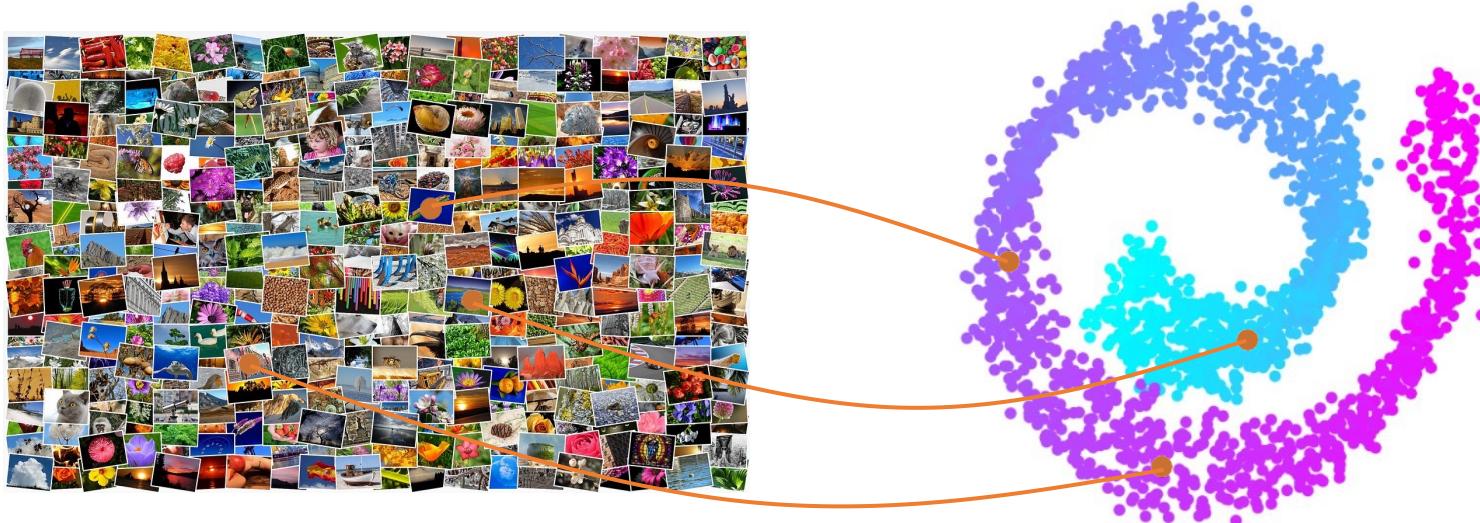
Statistical Perspective for Real Images

- Let's consider RGB images with a resolution of 256×256 .
- An image can be represented by a $256 \times 256 \times 3$ vector.
- This means that **an image is a point** in a $256 \times 256 \times 3$ -dimensional space.



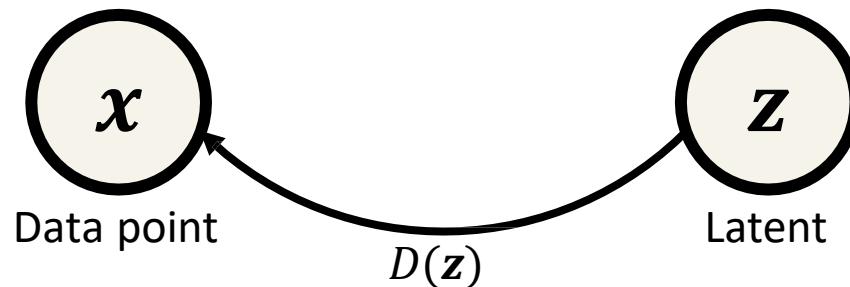
Statistical Perspective for Real Images

- Let us 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...?



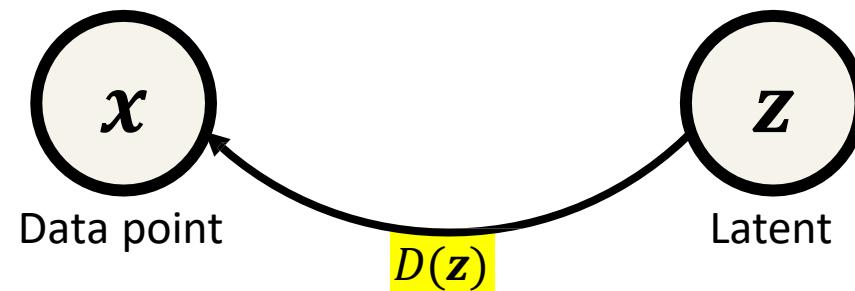
The Basic Idea

- Map a simple distribution $p(\mathbf{z})$ (e.g., a standard normal distribution $\mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$) 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.



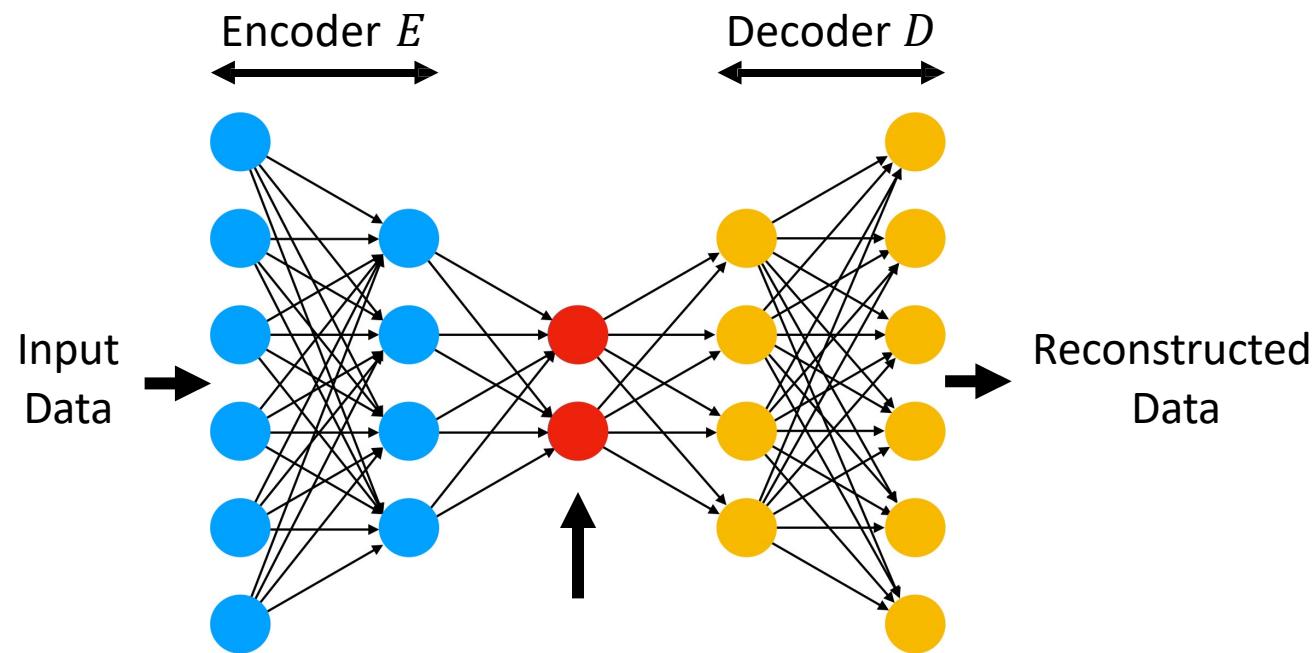
The Basic Idea

How to map a latent distribution $p(z)$ to the data distribution $p(x)$?
Let's use a **neural network**!



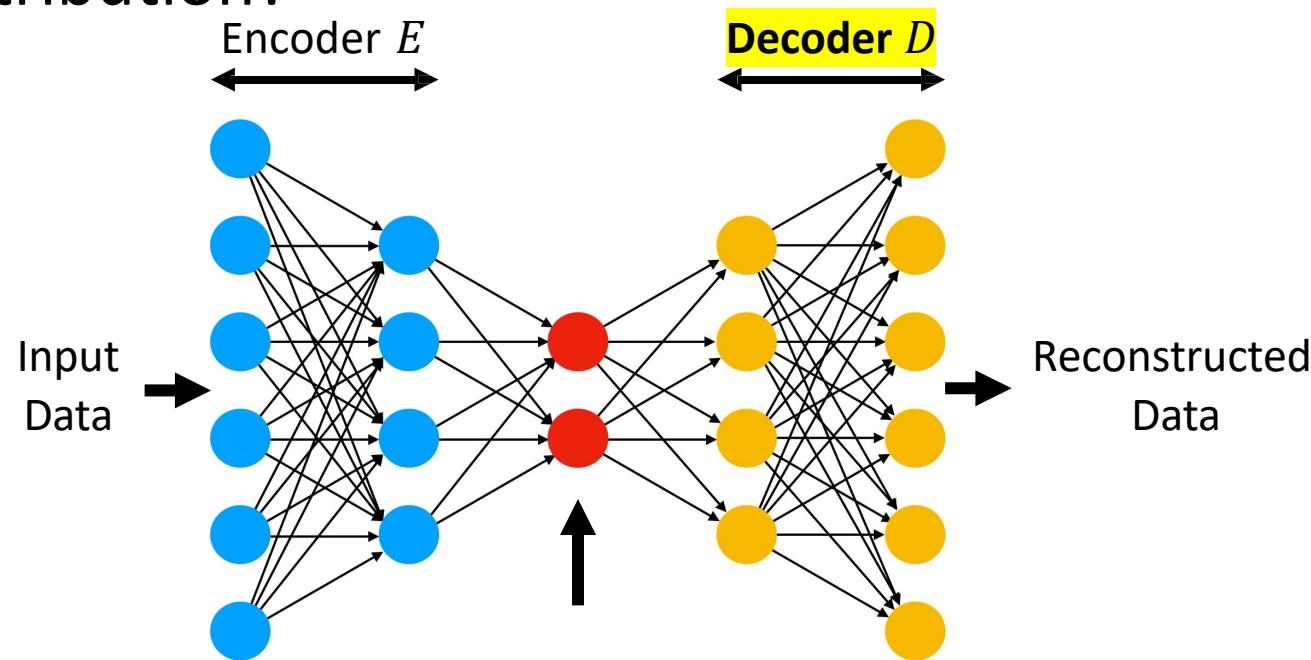
Autoencoder

An autoencoder is a neural network designed to **reconstruct the input data while encoding it into a lower-dimensional latent vector.**



Autoencoder

- What we need is the **decoder** ($\text{latent} \rightarrow \text{input data}$).
- But how do we **guarantee** that a latent is mapped to a data point of the data distribution?



Two Methods

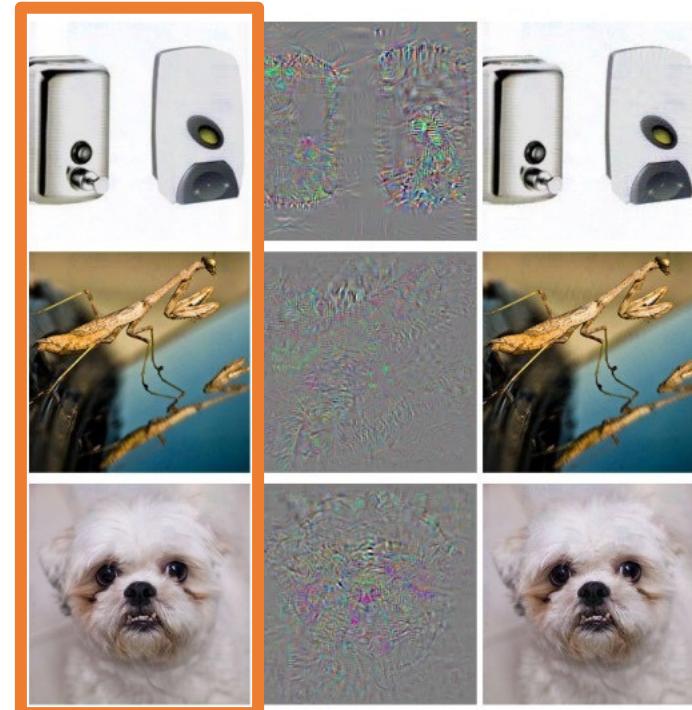
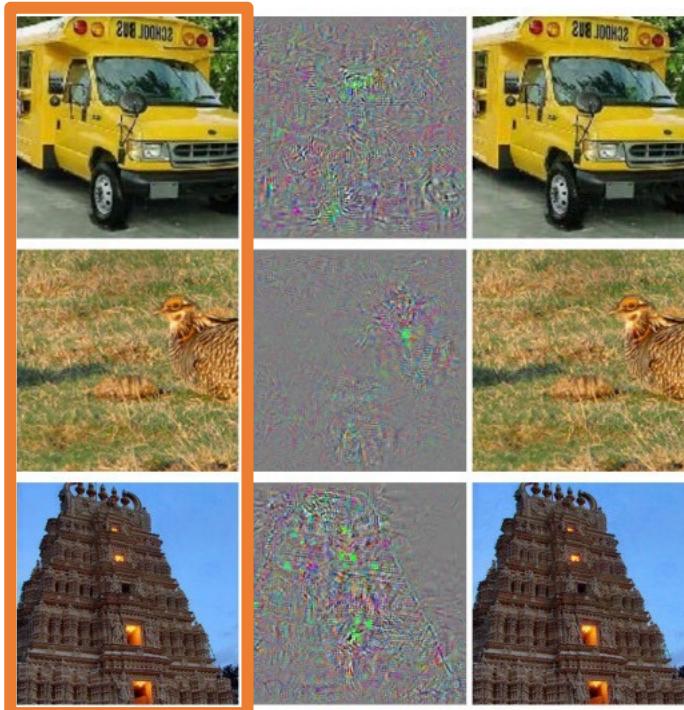
1. Generative Adversarial Network (GAN)
2. Variational Autoencoder (VAE)

Generative Adversarial Network (GAN)

Goodfellow et al., Generative Adversarial Networks, NeurIPS 2014.

Off-Topic: Adversarial Examples

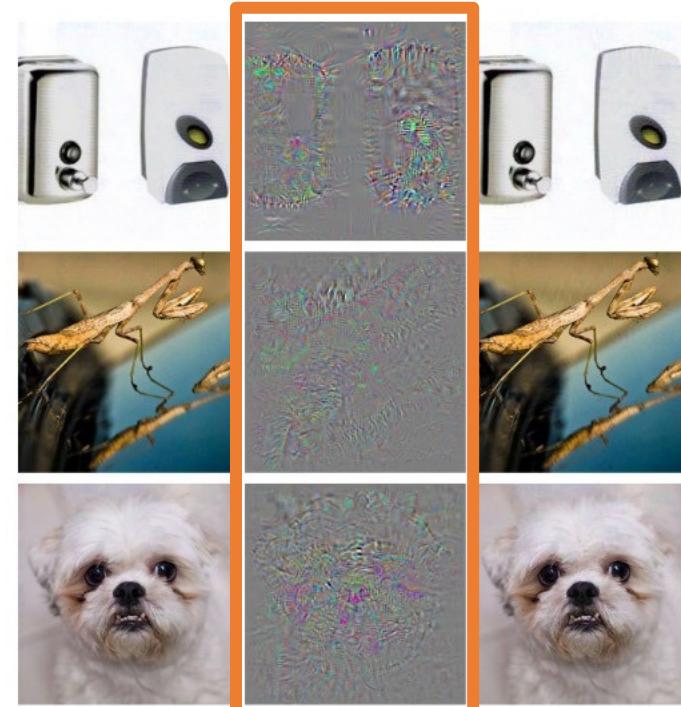
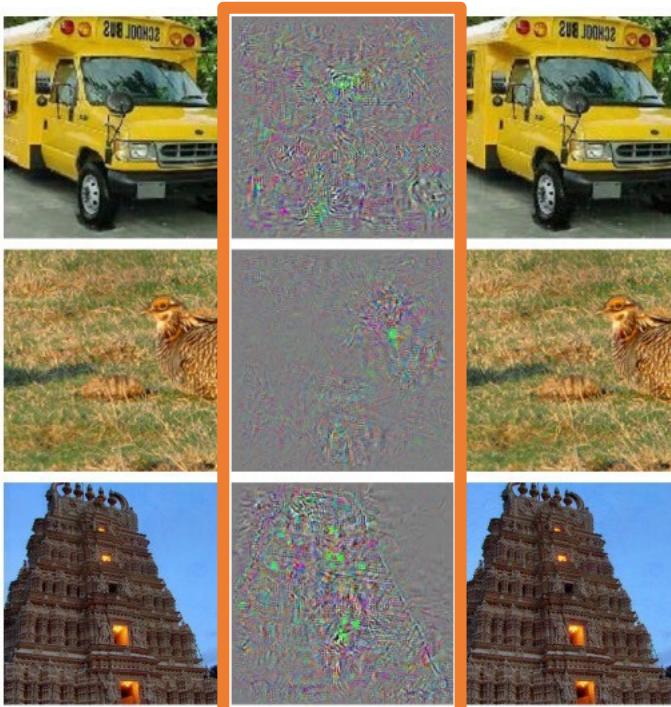
Neural networks work well, but how **vulnerable** is a neural network?



Images that are *correctly* classified

Off-Topic: Adversarial Examples

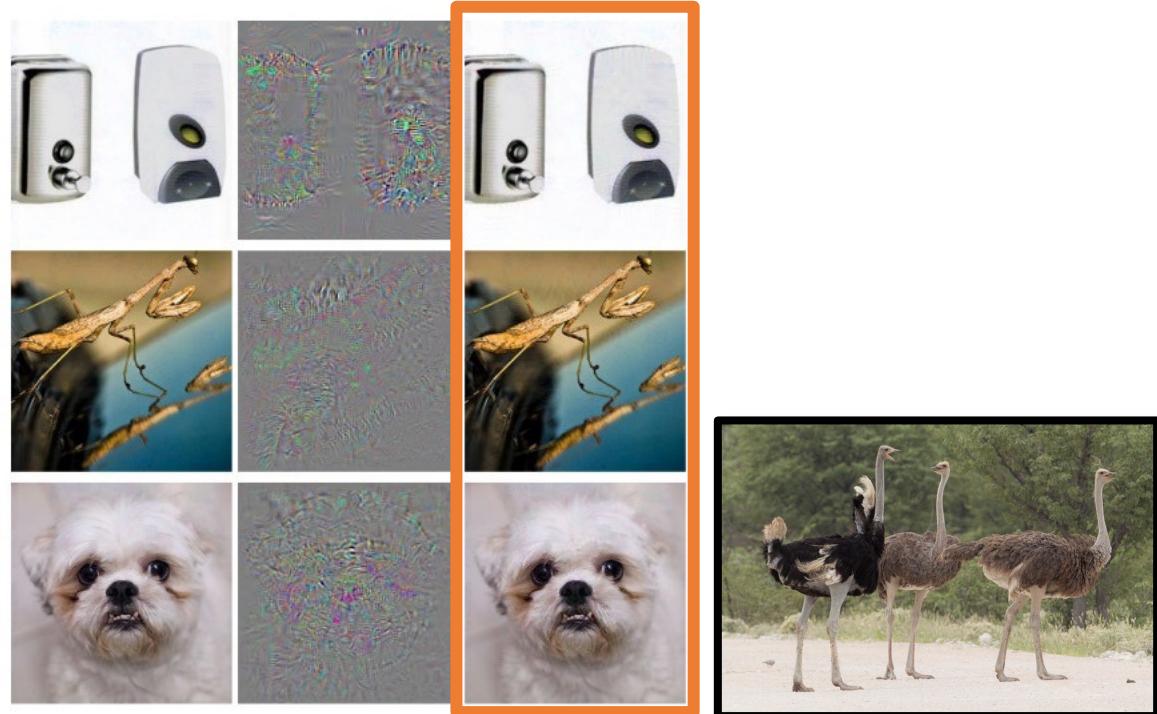
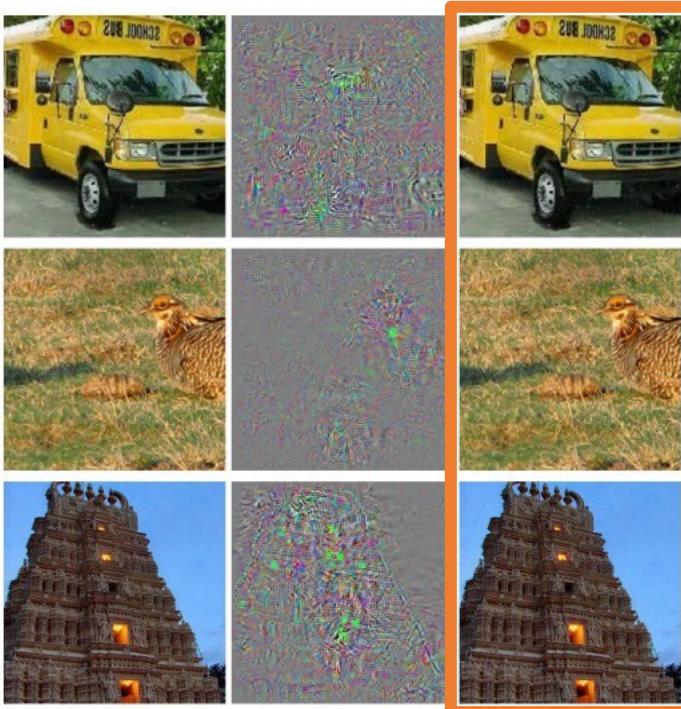
Neural networks work well, but how **vulnerable** is a neural network?



Difference between the two multiplied by 10 \times

Off-Topic: Adversarial Examples

Neural networks work well, but how **vulnerable** is a neural network?



Images that are *incorrectly* classified as ostrich

Wikipedia

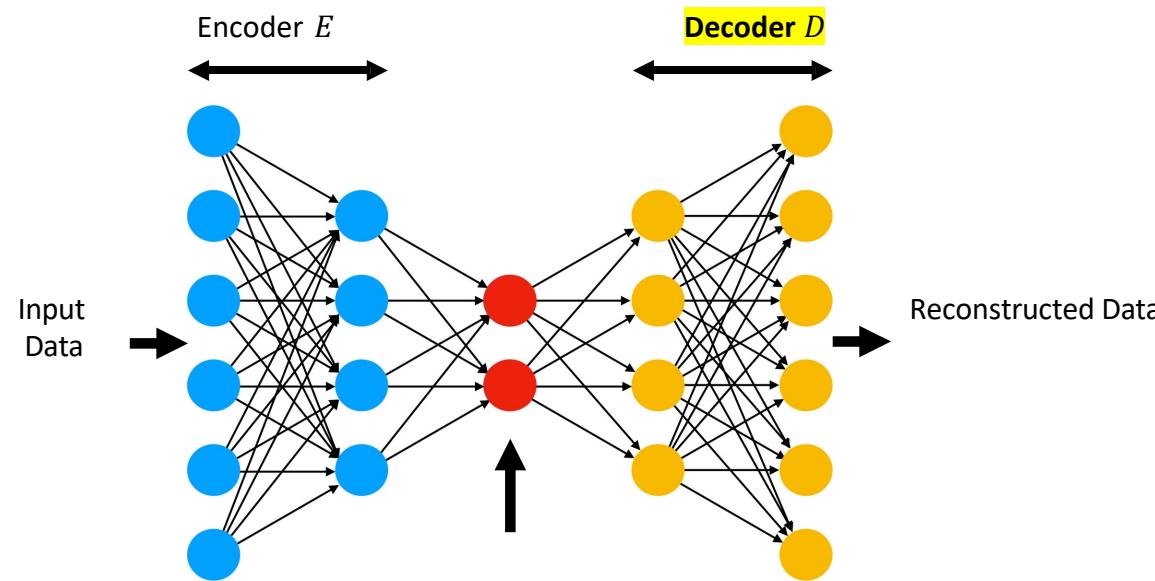
Off-Topic: Adversarial Examples

- Let's train a network to predict an image that causes the classifier to fail. **Adversarial attack!**
- Can we also **finetune the classifier** to prevent it from failing due to adversarial attacks?
- What happens if we make them **compete** against each other?

Generative Adversarial Network (GAN)

Think about

- a real/fake image **classifier** → **Discriminator**
- an **adversarial attack network** → **Generator (Decoder)**
that tries to make a real-like generated image.



Generative Adversarial Network (GAN)

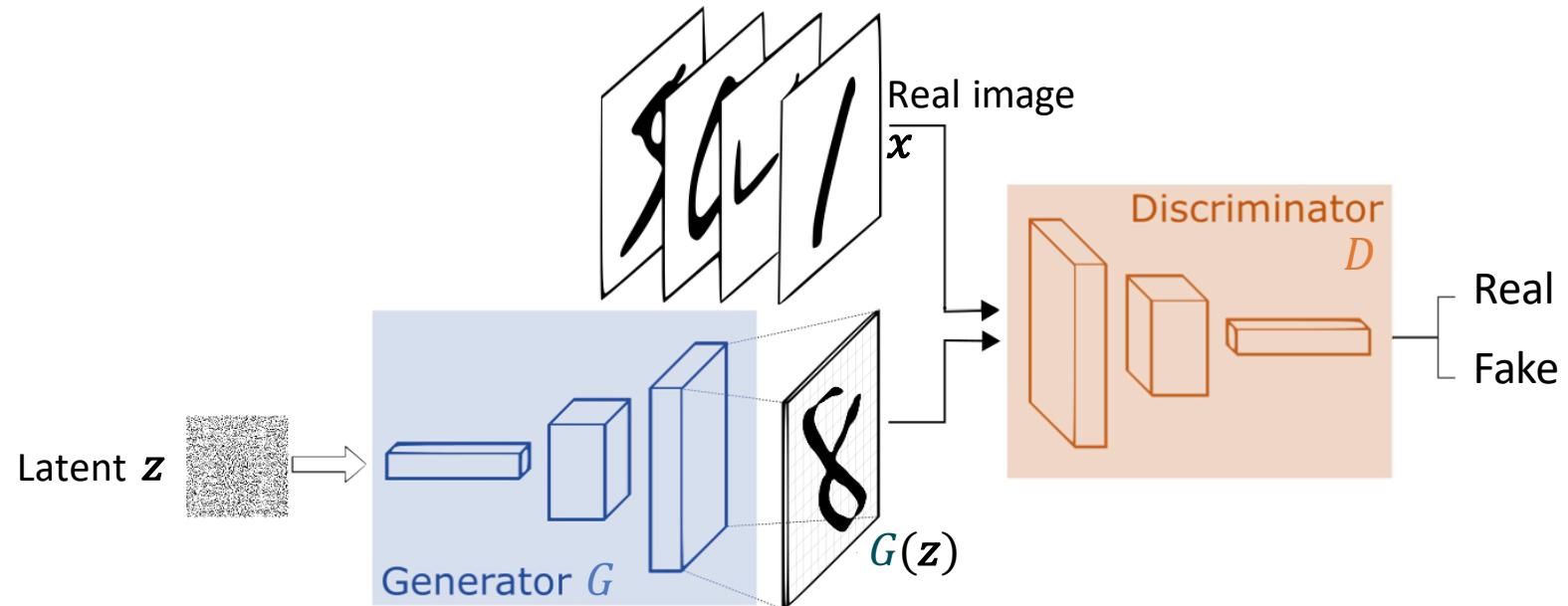
- **Generator** (or decoder) G takes a **latent** sampled from a unit Gaussian as input and generates a **synthetic** image.
- **Discriminator** D takes an image as input and **classifies** it as either **real** or **fake** (generated).

Make them compete against each other!

Generative Adversarial Network (GAN)

Loss function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p(x)}[\log D(x)] + \mathbb{E}_{z \sim p(z)}[\log(1 - D(G(z)))]$$



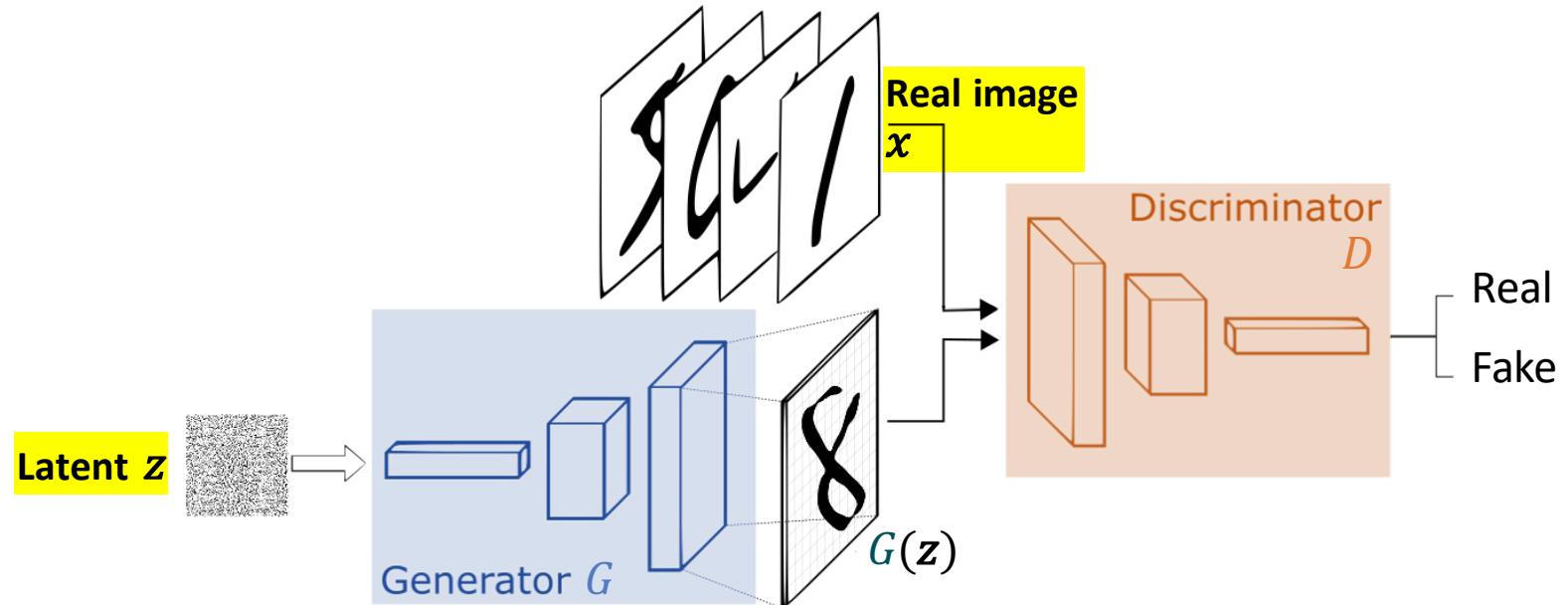
Generative Adversarial Network (GAN)

Loss function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

Real image sampled from
the data distribution

Latent sampled from
the latent distribution

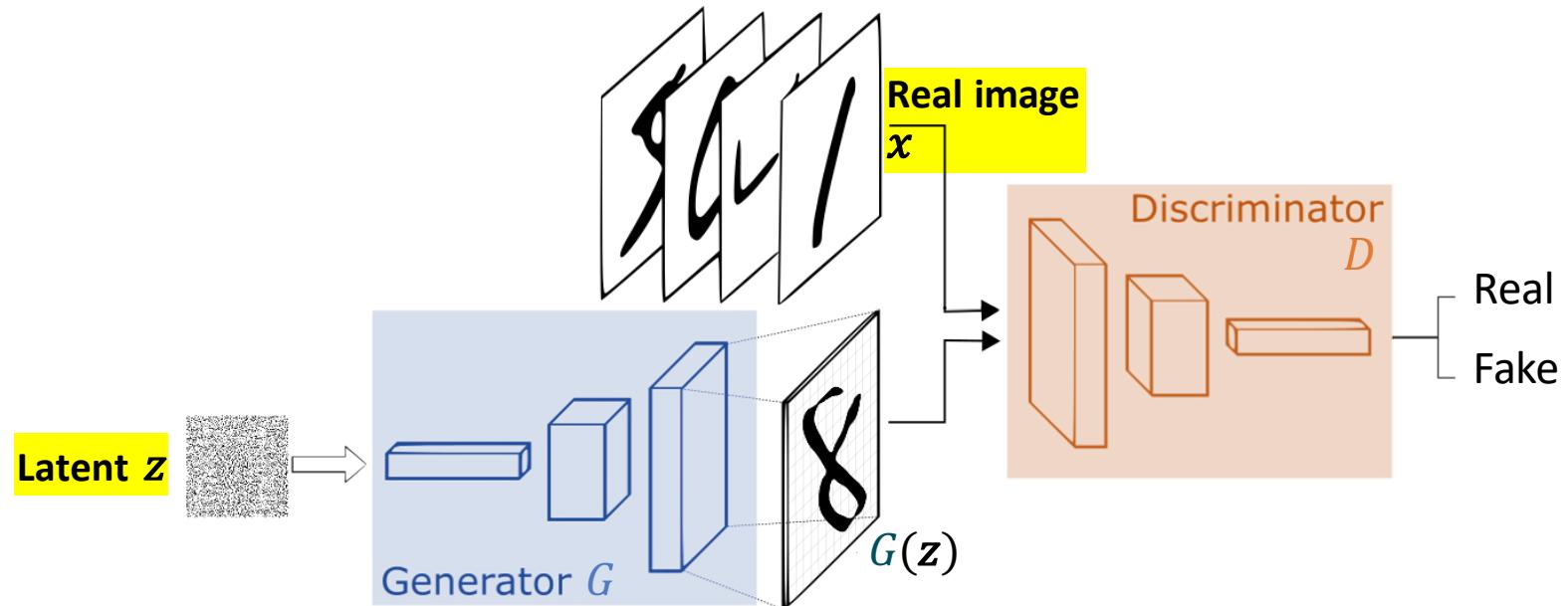


Generative Adversarial Network (GAN)

Loss function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

The *fake* image synthesized from the latent z .



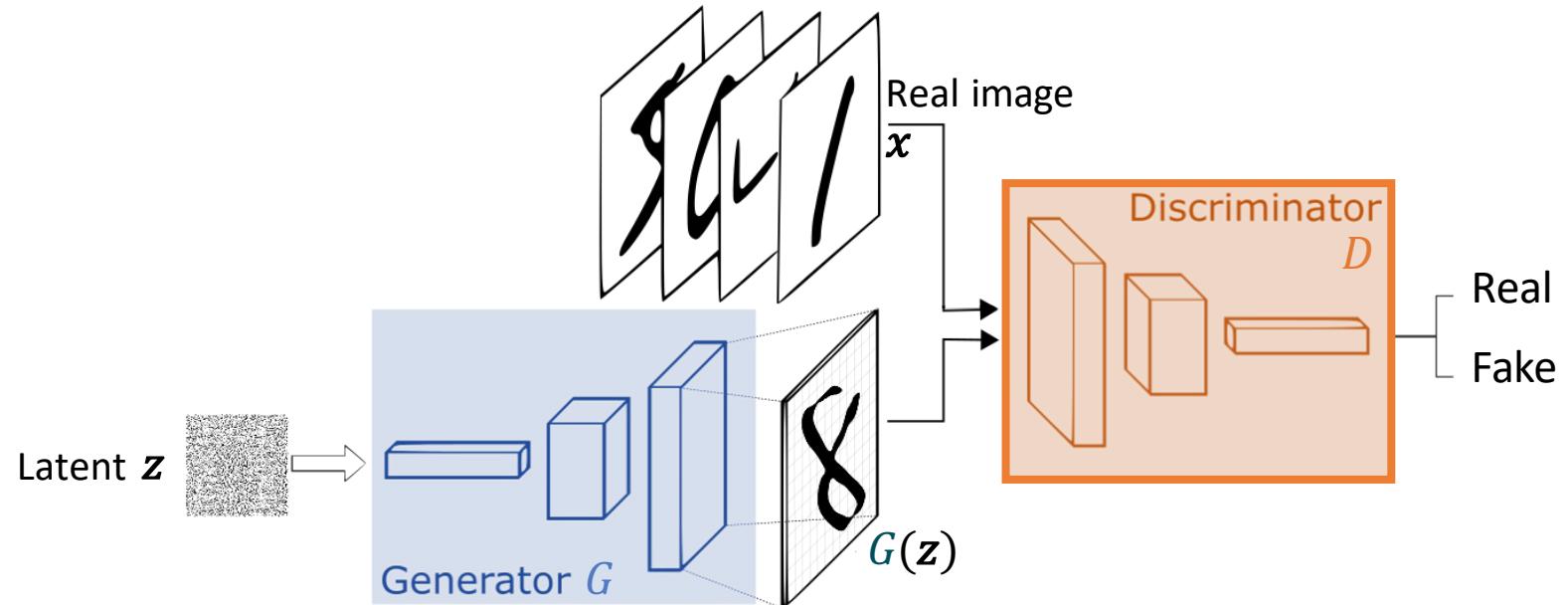
Generative Adversarial Network (GAN)

Loss function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

The probability of the *real* image x being a real image.

The probability of the *fake* image $G(z)$ being a real image.

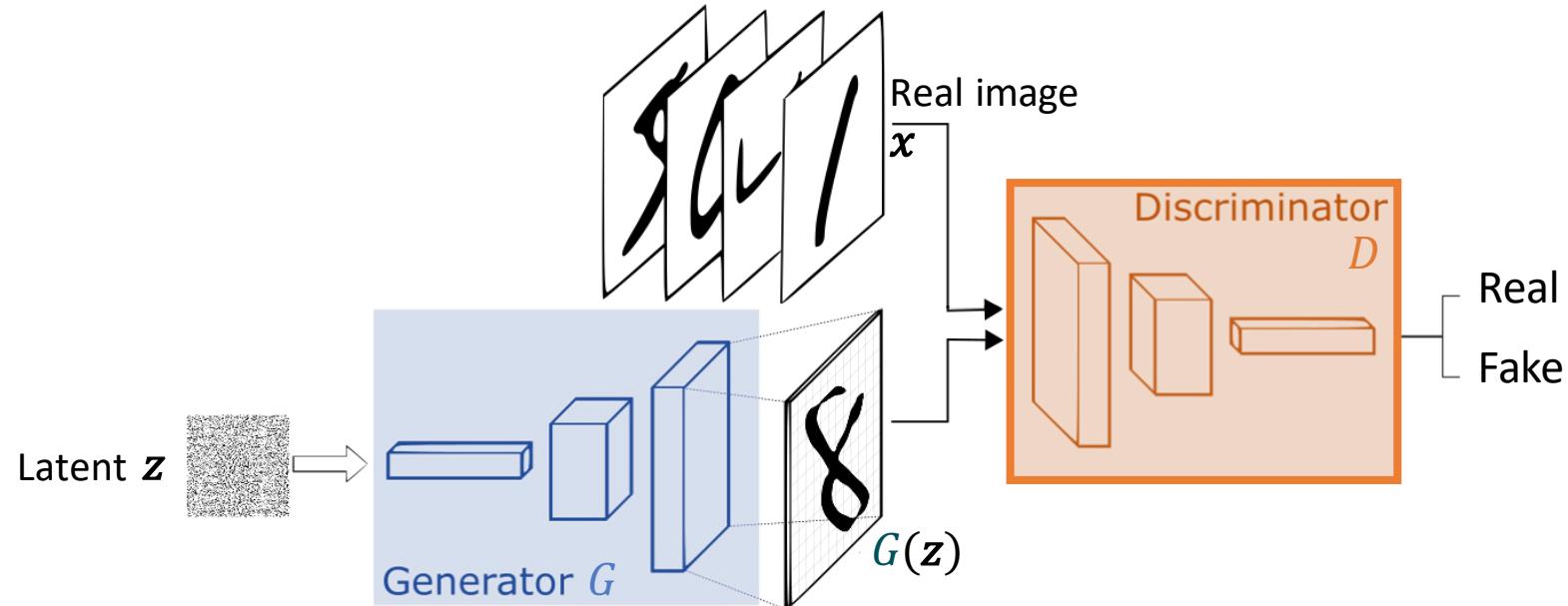


Generative Adversarial Network (GAN)

Loss function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

Discriminator *increases* the probability for the real image x ...
and *decreases* the probability for the synthetic image $G(z)$.

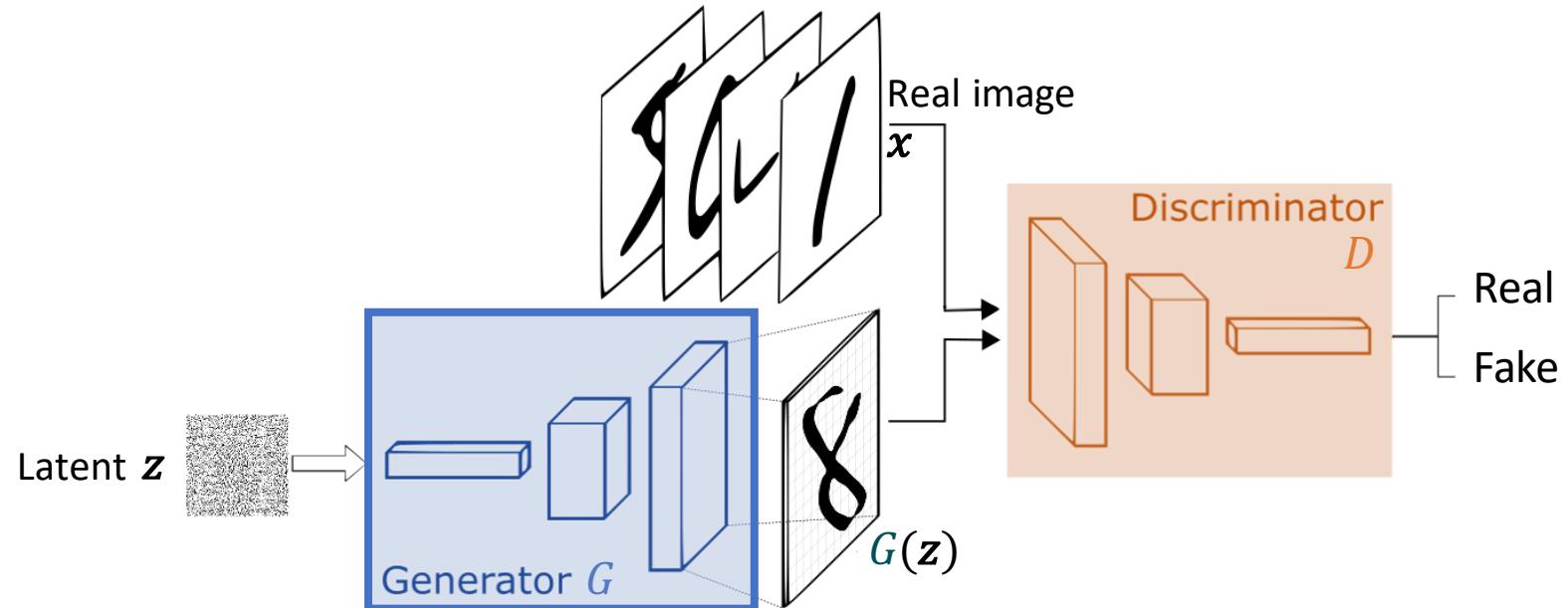


Generative Adversarial Network (GAN)

Loss function:

$$\min_{\underline{G}} \max_{\underline{D}} V(\underline{D}, \underline{G}) = \mathbb{E}_{x \sim p(x)} [\log \underline{D}(x)] + \underline{\mathbb{E}_{z \sim p(z)} [\log(1 - \underline{D}(\underline{G}(z)))]}$$

Generator *increases* the probability
for the synthetic image $G(z)$.



Generative Adversarial Network (GAN)

Loss function:

$$\min_{\underline{G}} \max_{\underline{D}} V(\underline{D}, \underline{G}) = \mathbb{E}_{x \sim p(x)} [\log \underline{D}(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - \underline{D}(\underline{G}(z)))]$$

A **min-max (minimax)** optimization problem, which is known as **very difficult to solve!**

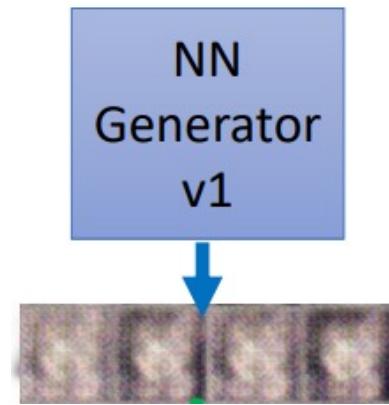
Generative Adversarial Network (GAN)

NN
Generator
v1

Real images:



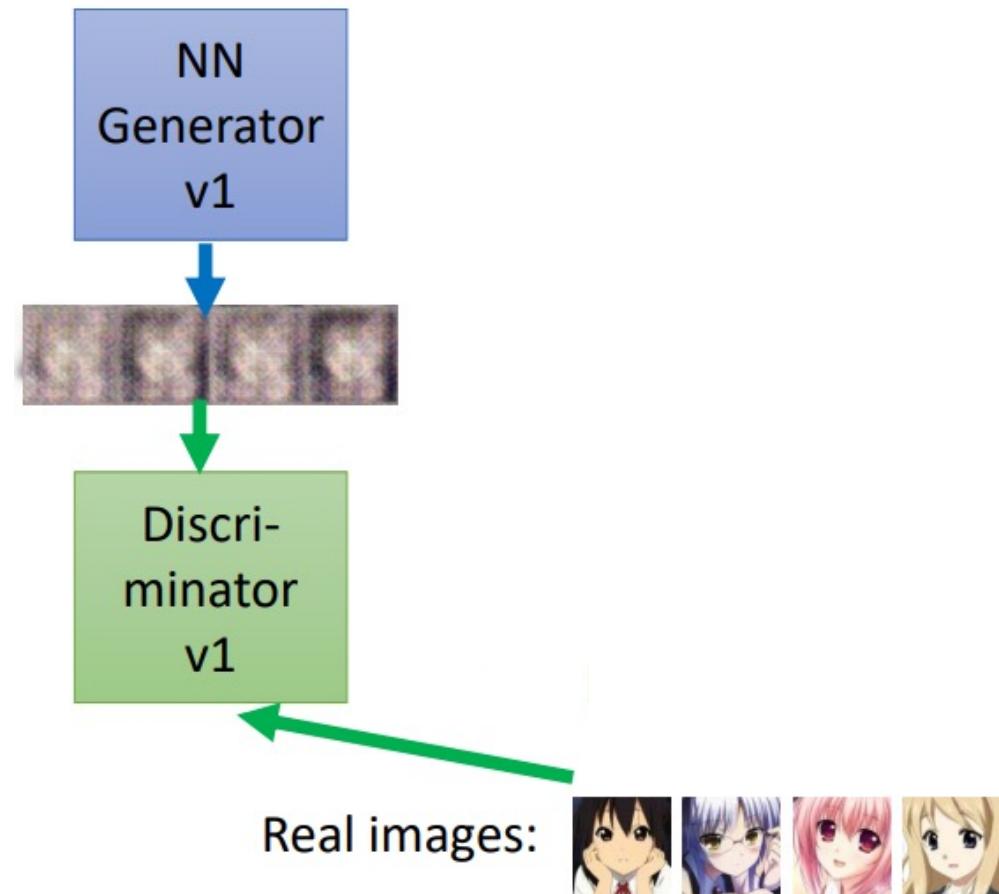
Generative Adversarial Network (GAN)



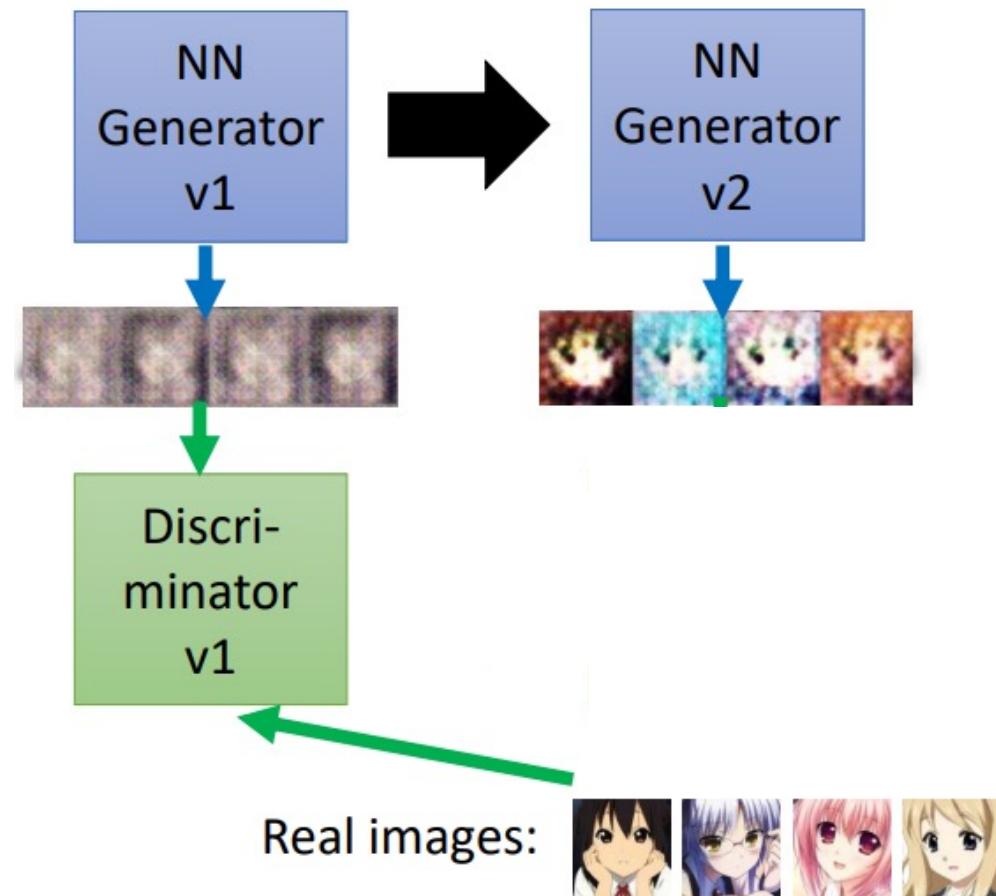
Real images:



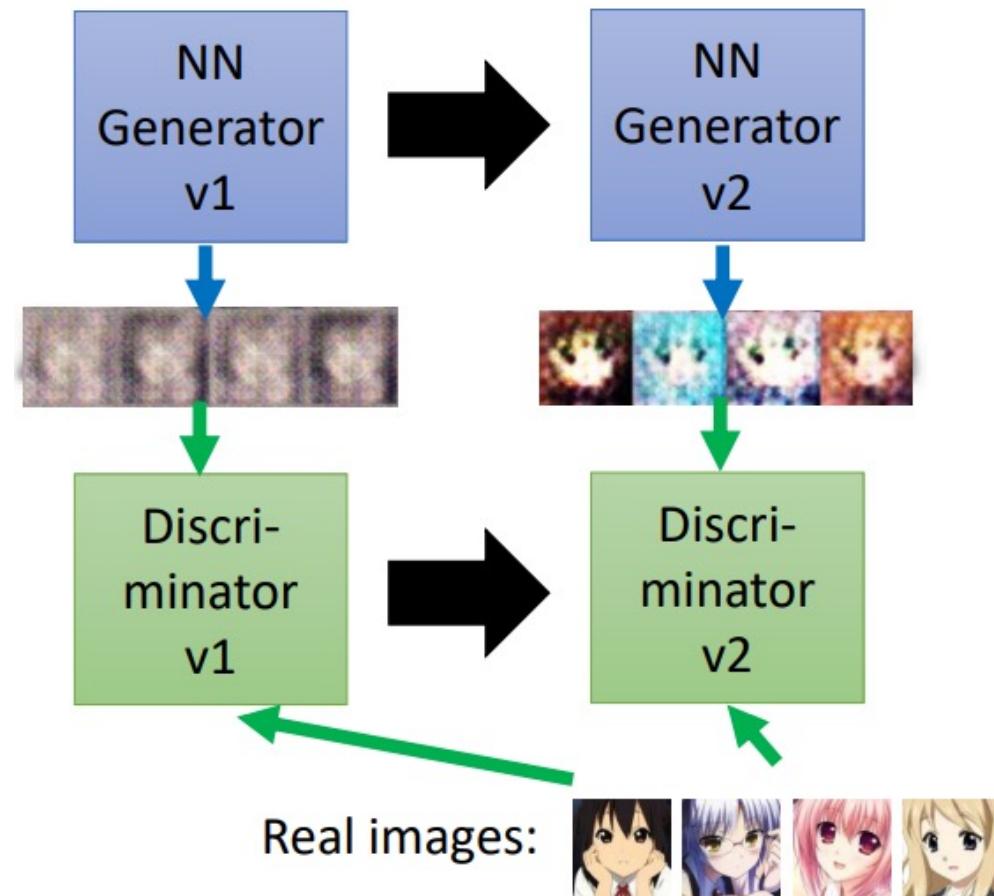
Generative Adversarial Network (GAN)



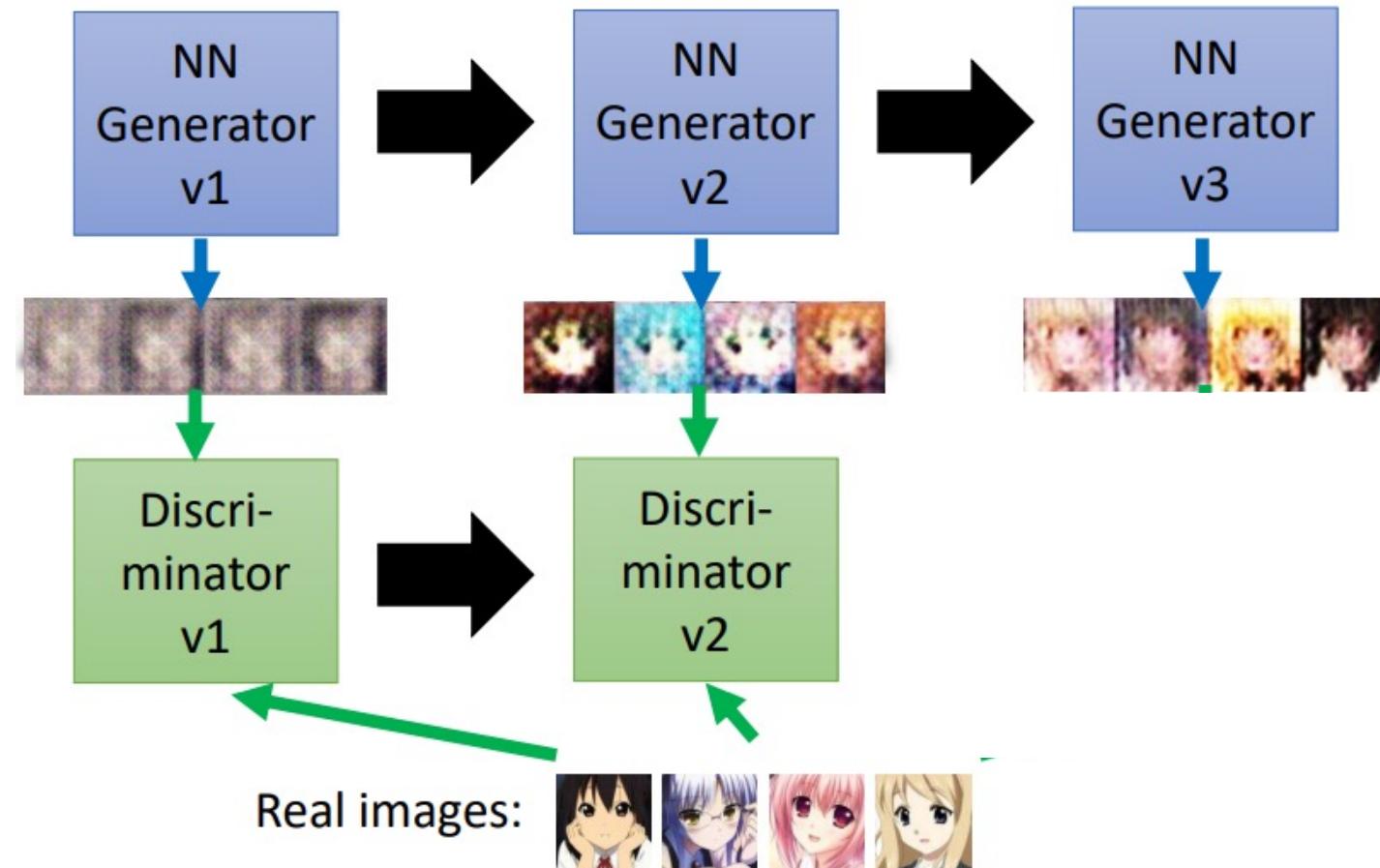
Generative Adversarial Network (GAN)



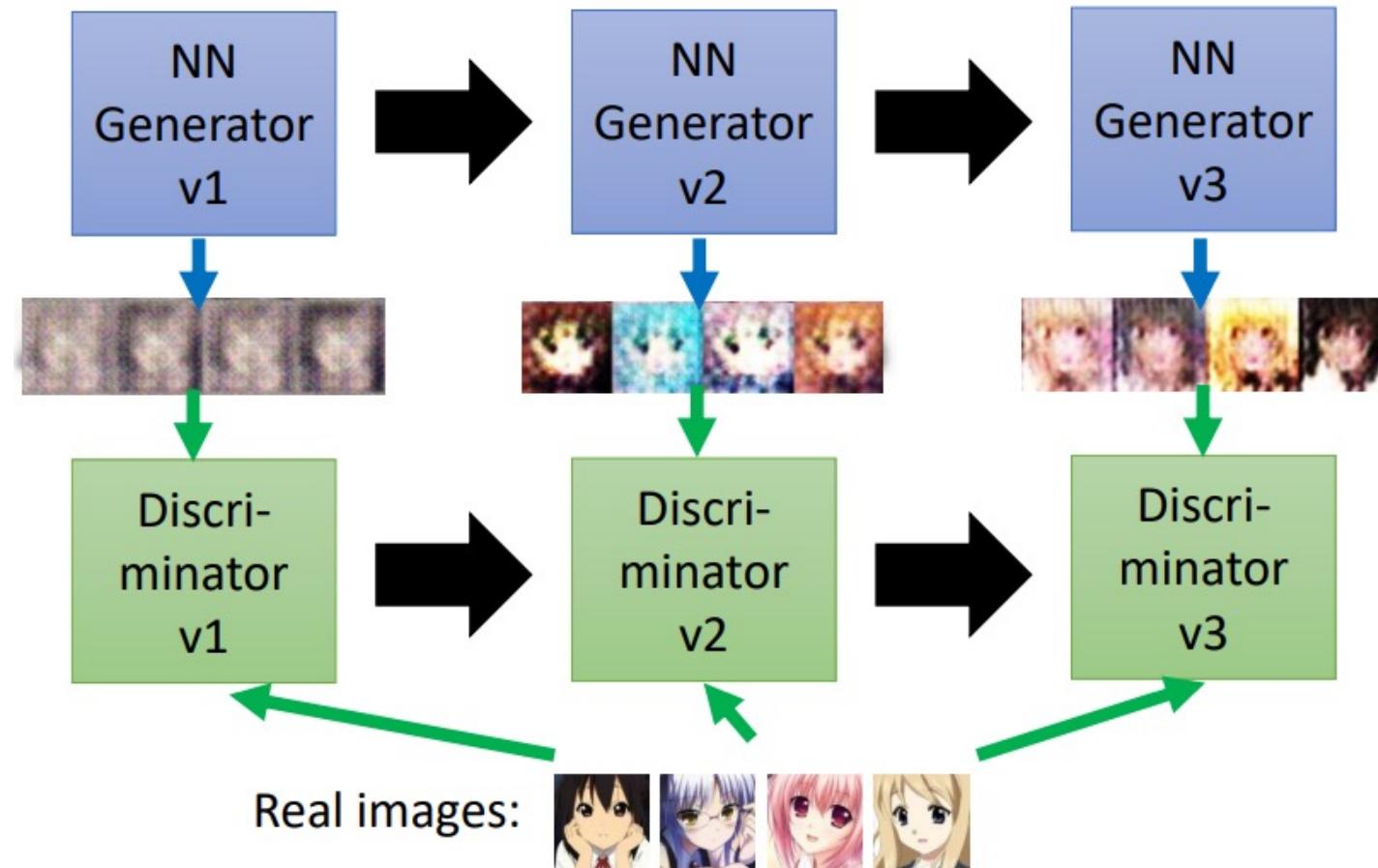
Generative Adversarial Network (GAN)



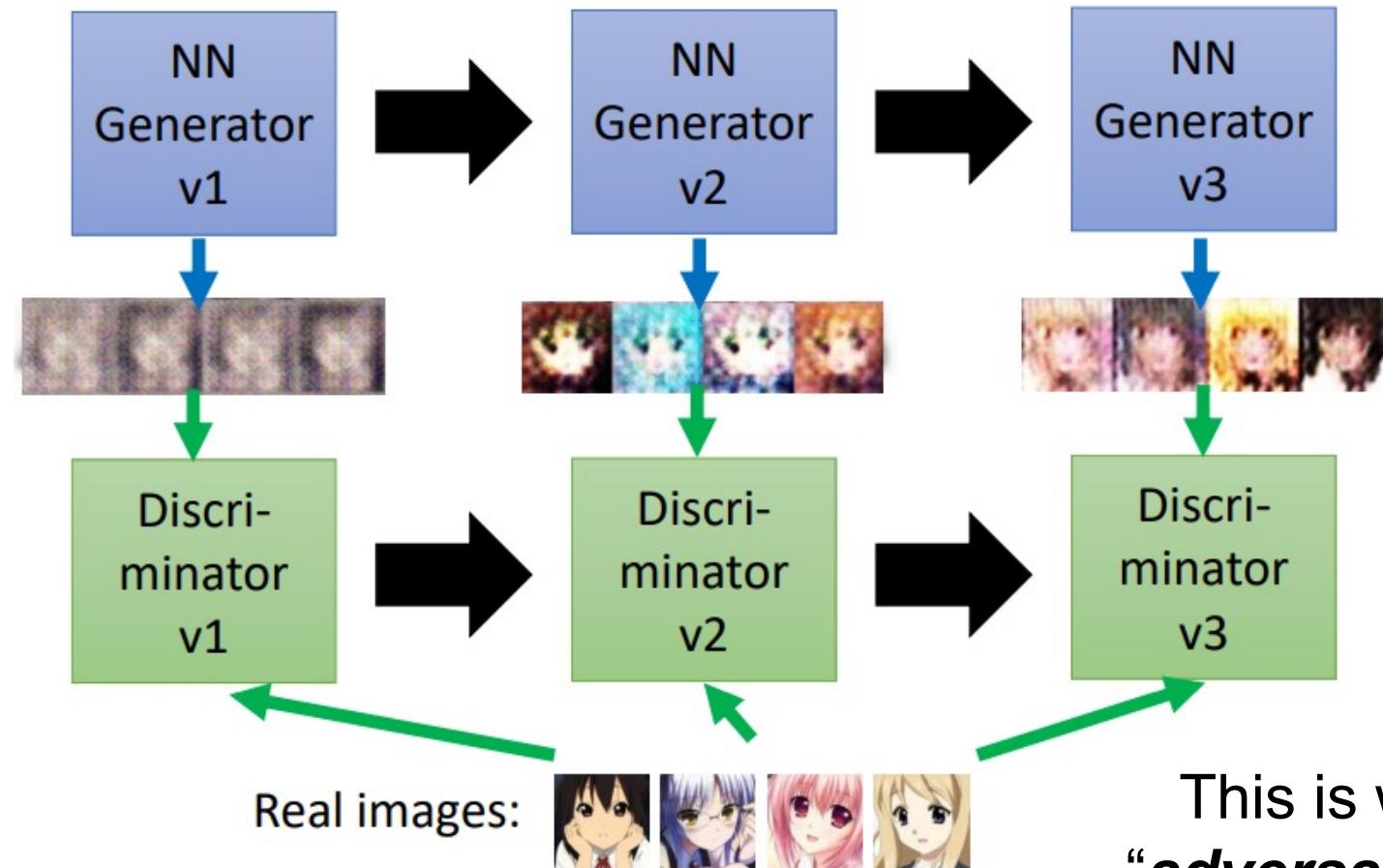
Generative Adversarial Network (GAN)



Generative Adversarial Network (GAN)



Generative Adversarial Network (GAN)

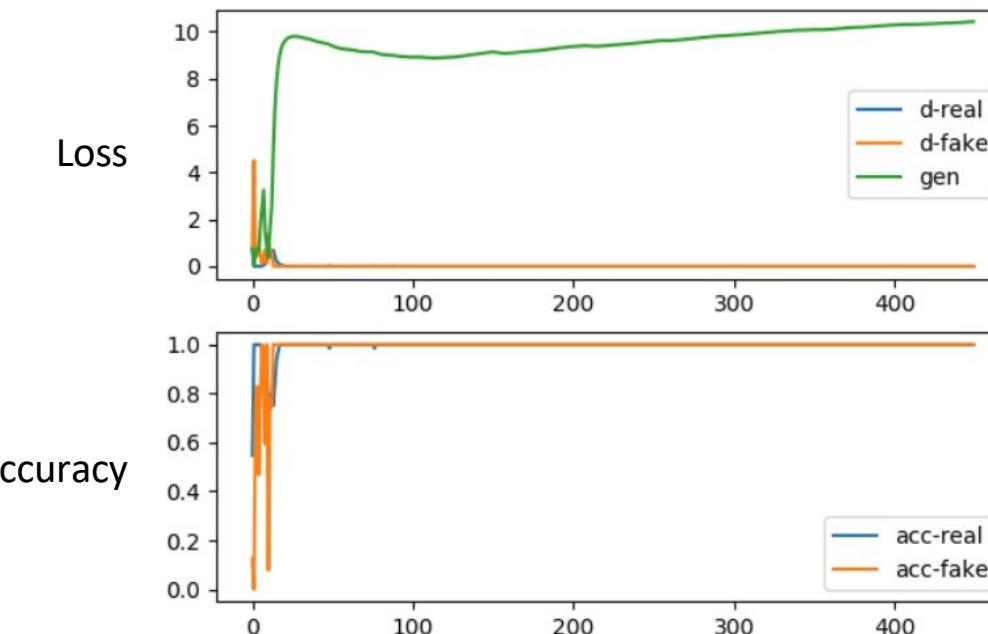


Challenges in GAN Training

1. Non-convergence & Instability
2. Mode collapse

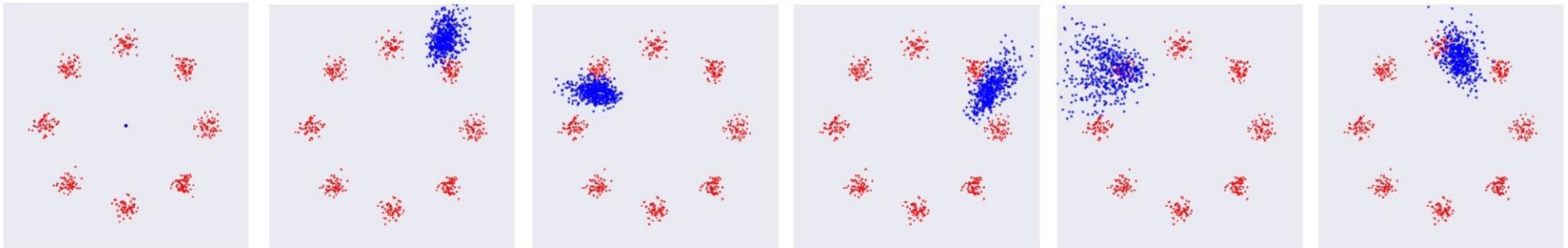
Non-Convergence & Instability

The **discriminator becomes too strong** and can easily distinguish between real and fake samples, leading to near-zero gradients for the generator and halting its learning.



Mode Collapse

The generator tends to **collapses** to a few modes of the data distribution instead of capturing the full diversity of the true data distribution.



Advances in GAN

Index	Software name	Language	Backend	Link
1	CGAN	Python	PyTorch	https://github.com/Lornatang/CGAN-PyTorch
2	DCGAN	Python	PyTorch	https://github.com/Natsu6767/DCGAN-PyTorch
3	AAEs	Python	TensorFlow	https://github.com/conan7882/adversarial-autoencoders
4	InfoGAN	Python	TensorFlow	https://github.com/openai/InfoGAN
5	SAD-GAN	—	—	—
6	LSGAN	Python	PyTorch	https://github.com/xudonmao/LSGAN
7	SRGAN	Python	TensorFlow	https://github.com/tensorlayer/SRGAN
8	WGAN	Python	PyTorch	https://github.com/Zeleni9/pytorch-wgan
9	CycleGAN	Python	TensorFlow	https://github.com/junyanz/CycleGAN
10	ProGAN	Python	PyTorch	https://github.com/tkarras/progressive_growing_of_gans
11	MidiNet	Python	TensorFlow	https://github.com/RichardYang40148/MidiNet
12	SN-GAN	Python	PyTorch	https://github.com/hanyoseob/pytorch-SNGAN
13	RGAN	Python	TensorFlow	https://github.com/ratschlab/RGAN
14	StarGAN	Python	PyTorch	https://github.com/yunjey/stargan
15	BigGAN	Python	PyTorch	https://github.com/ajbrock/BigGAN-PyTorch
16	MI-GAN	Python	TensorFlow	https://github.com/hazratali/MI-GAN
17	AttGAN	Python	TensorFlow	https://github.com/LynnHo/AttGAN-Tensorflow
18	PATE-GAN	Python	TensorFlow	https://github.com/vanderschaarlab/mlforhealthlabpub/tree/main/alg/pategan
19	DM-GAN	Python	PyTorch	https://github.com/MinfengZhu/DM-GAN
20	SinGAN	Python	PyTorch	https://github.com/tamarott/SinGAN
21	POLY-GAN	Python	PyTorch	https://github.com/nile649/POLY-GAN
22	MIEGAN	—	—	—
23	VQGAN	Python	PyTorch	https://github.com/dome272/VQGAN-pytorch
24	DALL-E	Python	PyTorch	https://github.com/lucidrains/DALLE-pytorch
25	CEGAN	—	—	—
26	Seismogen	Python	PyTorch	https://github.com/Miffka/seismogen
27	MetroGAN	Python	PyTorch	https://github.com/zwy-Giser/MetroGAN
28	M3GAN	Python	PyTorch	https://github.com/SLZWVICTOR/M3GAN
29	CNTS	Python	PyTorch	https://github.com/BomBooooo/CNTS/tree/main
30	RidgeGAN	Python	PyTorch	https://github.com/rahisha-thottolil/ridgegan

StyleGAN2





Variational Autoencoder (VAE)

Kingma and Welling , Auto-Encoding Variational Bayes, ICLR 2014.

Variational Autoencoder (VAE)

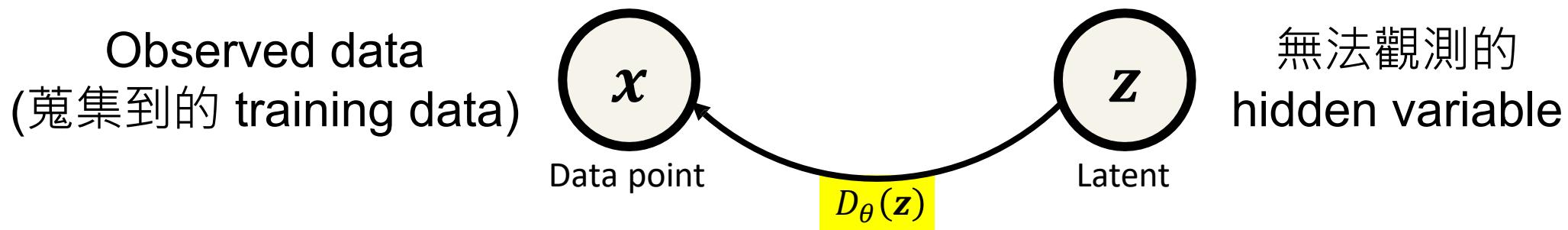
- How to make a generative model *without* solving a **minimax** problem?
- Let's represent the mapping from the latent distribution $p(\mathbf{z})$ to the data distribution $p(\mathbf{x})$ as a **conditional distribution** $p(\mathbf{x}|\mathbf{z})$.

Variational Autoencoder (VAE)

Let's consider the **decoder (generator)** as predicting the **mean** of the **conditional distribution** $p(x|z)$:

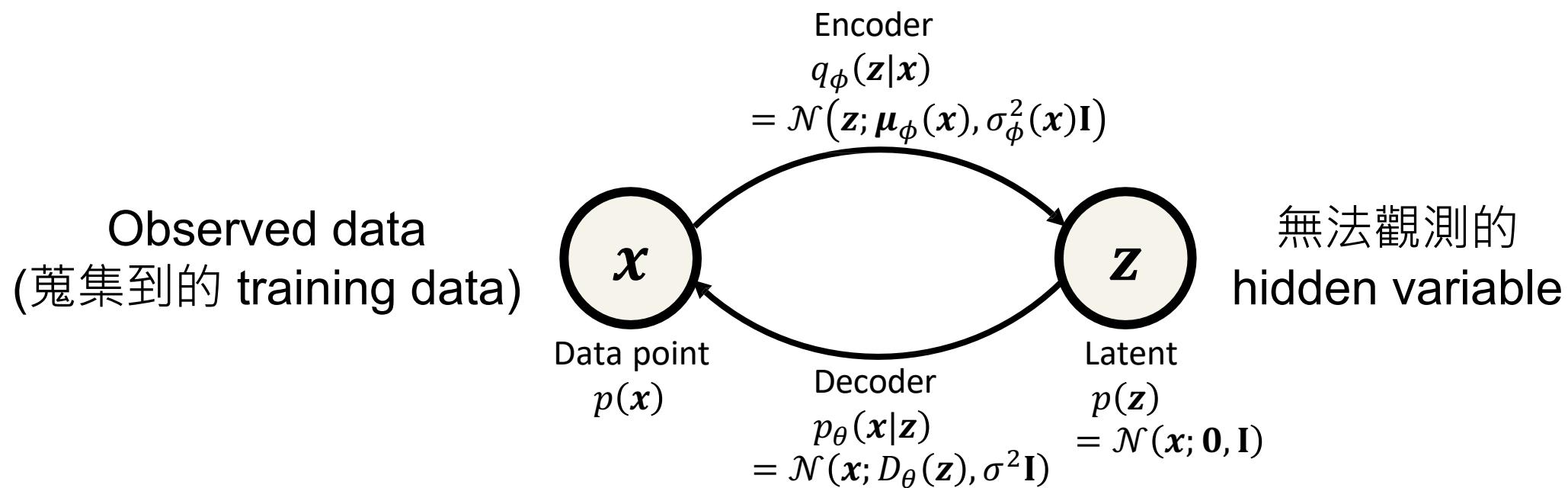
$$p(x|z) = \mathcal{N}(x; D_\theta(z), \underline{\sigma^2 I})$$

Fixed variance



Variational Autoencoder (VAE)

Introduce an **encoder** to learn a proxy of the posterior distribution.



Basics

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

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.

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


Take the integral over z to marginalize x .

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

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.

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

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

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

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

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

Kullback-Leibler (KL) Divergence

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

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

- $D_{KL}(p\|q) \geq 0$ with equality iff $p = q$ almost everywhere
- Not a true metric (asymmetric, doesn't satisfy triangle inequality)

Q. What is $D_{KL}(p\|p)$?

Kullback-Leibler (KL) Divergence

Q.

When $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2))$ and $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$,
What is $D_{KL}(p\|q)$?

Hint.

Closed-form solution for d -dimensional Gaussian distributions with

- $p \sim \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$
- $q \sim \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$

$$D_{KL}(p\|q) = \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_q|}{|\boldsymbol{\Sigma}_p|} - d + \text{tr}(\boldsymbol{\Sigma}_q^{-1} \boldsymbol{\Sigma}_p) + (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_q^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) \right]$$

Kullback-Leibler (KL) Divergence

A.

- $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2))$: General Gaussian
- $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$: Standard normal distribution

$$D_{KL}(p\|q) = \frac{1}{2} \sum_{i=1}^d (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2)$$

Jensen's Inequality

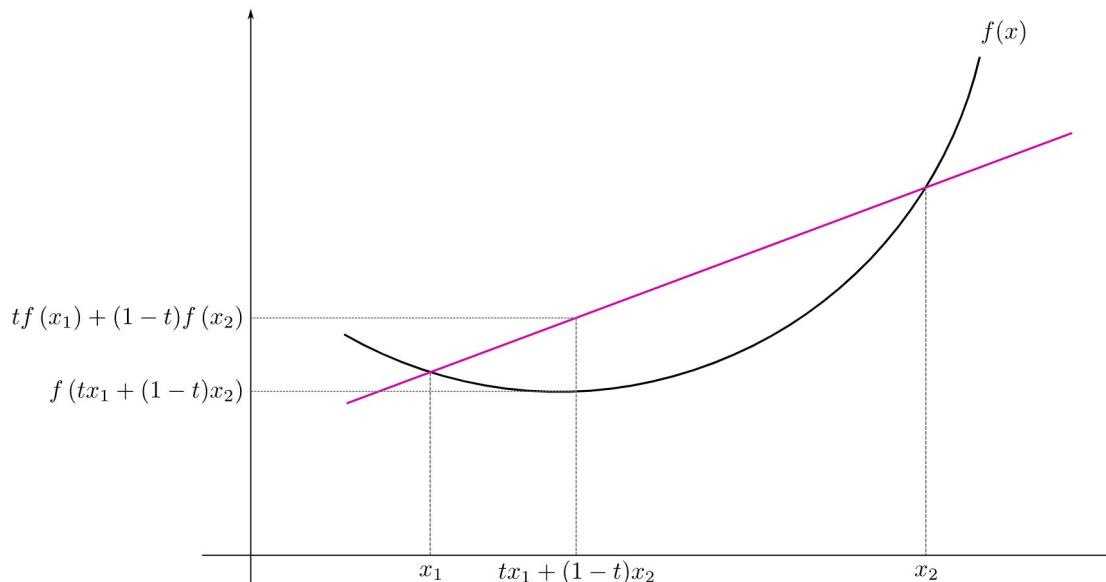


Johan Jensen

f is a **convex** function if

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

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



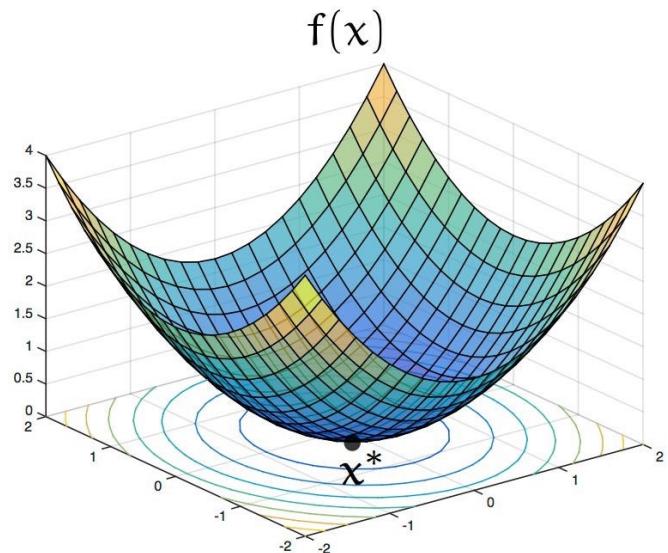
Jensen's Inequality

f is a **convex** function if

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

for all $x_i, t_i \in [0,1]$.

Affine combination 



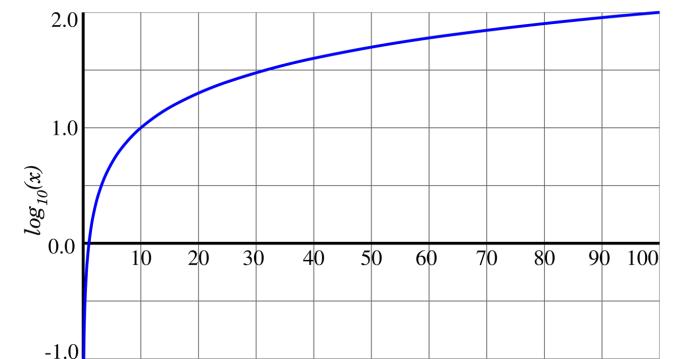
Jensen's Inequality

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 an affine combination.

Q. What if f is a **concave** function? E.g., \log ?



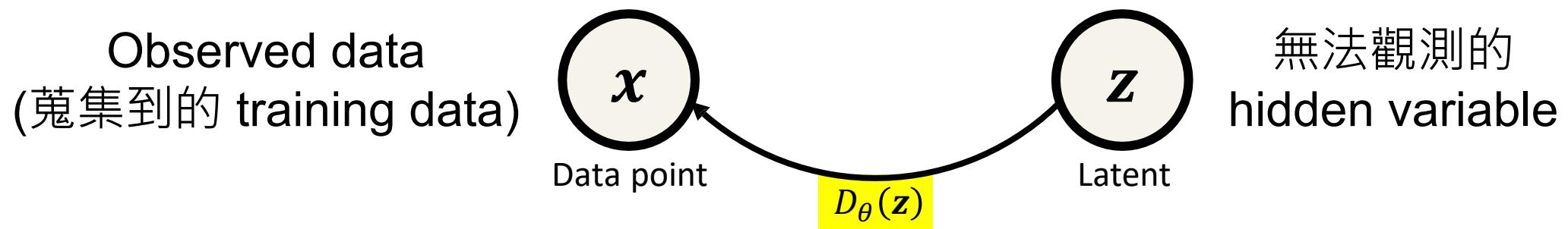
Back to VAE...

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

Variational Autoencoder (VAE)

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

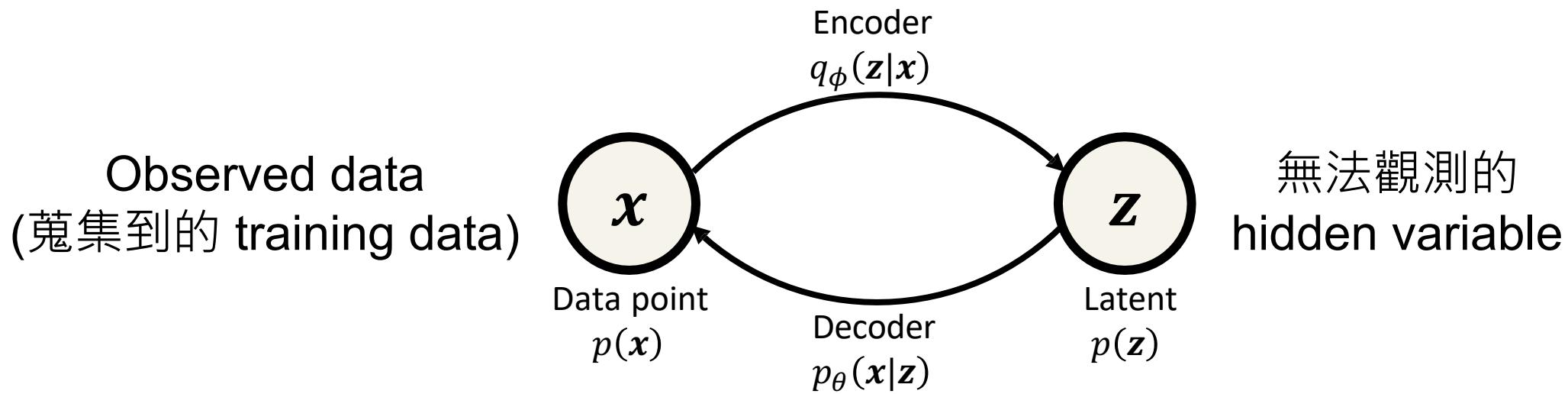
$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D_\theta(\mathbf{z}), \sigma^2 \mathbf{I})$$



Bayes' Rule

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

Likelihood Decoder Prior
Posterior Encoder Latent
Marginal Data



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

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

But this conditional distribution
is unknown.

Evidence Lower Bound (ELBO)

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

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

Auxiliary function
沒有特別指明是什麼樣的分佈

Evidence Lower Bound (ELBO)

- We cannot directly maximize $p(x) = \frac{p(x,z)}{p(z|x)}$ since $p(z|x)$ is unknown.
- Let's think about maximizing the **lower bound** of $\log p(x)$:

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

Auxiliary function
沒有特別指明是什麼樣的分佈

Evidence Lower Bound (ELBO)

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

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

Auxiliary function
沒有特別指明是什麼樣的分佈

$$\mathcal{N}(\mathbf{z}; \mathcal{M}(\mathbf{x}), \sigma^2(\mathbf{x}))$$

Evidence Lower Bound (ELBO)

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

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

Auxiliary function
沒有特別指明是什麼樣的分佈

$$\mathcal{N}(\mathbf{z}; \mathcal{M}(\mathbf{x}), \sigma^2(\mathbf{x}))$$

Evidence Lower Bound (ELBO)

- $q_\phi(\mathbf{z}|\mathbf{x})$ 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 don't know $p(\mathbf{z}|\mathbf{x})$.

Evidence Lower Bound (ELBO)

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

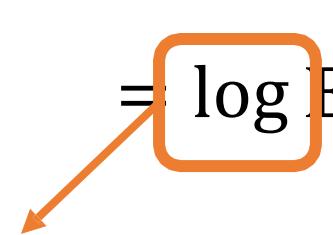
Evidence Lower Bound (ELBO)

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

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

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

Concave
function



Evidence Lower Bound (ELBO)

Q.

What is the **Jensen gap** $\left(\log p(x) - \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(x,z)}{q_\phi(z|x)} \right] \right)$?

Evidence Lower Bound (ELBO)

A.

The Jensen gap is the KL divergence between the variational posterior and the true posterior

$$\text{Jensen gap} = D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}|\mathbf{x}))$$

Key Insights:

- The Jensen gap is always ≥ 0 (KL divergence property)
- The gap equals 0 only when $q_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$ (perfect approximation)
- This explains why ELBO is a lower bound: $\log p(\mathbf{x}) = \text{ELBO} + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}|\mathbf{x})) \geq \text{ELBO}$
- In practice: We can't minimize this gap directly (requires knowing $p(\mathbf{z}|\mathbf{x})$, so we maximize ELBO instead)

Evidence Lower Bound (ELBO)

Let's decompose ELBO:

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

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

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

Reconstruction term
to be *maximized*.

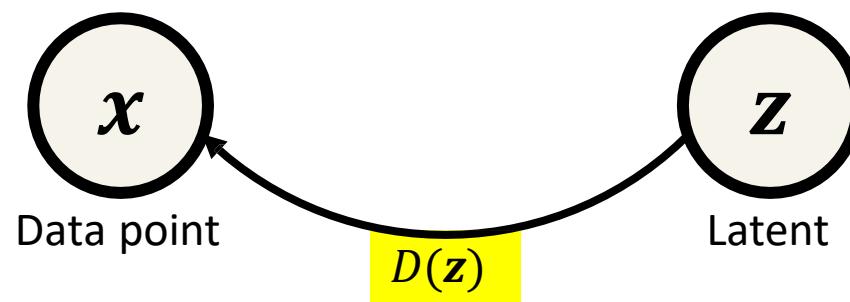
Prior matching term
to be *minimized*.

Back to VAE...

In VAE, we want model $p(x)$ as

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

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

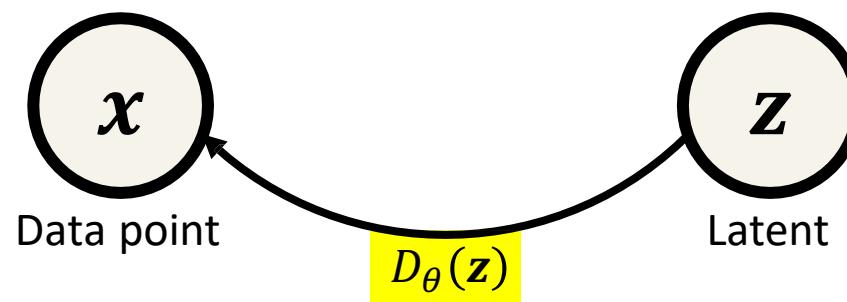


Variational Autoencoder (VAE)

In VAE, we want model $p(x)$ as

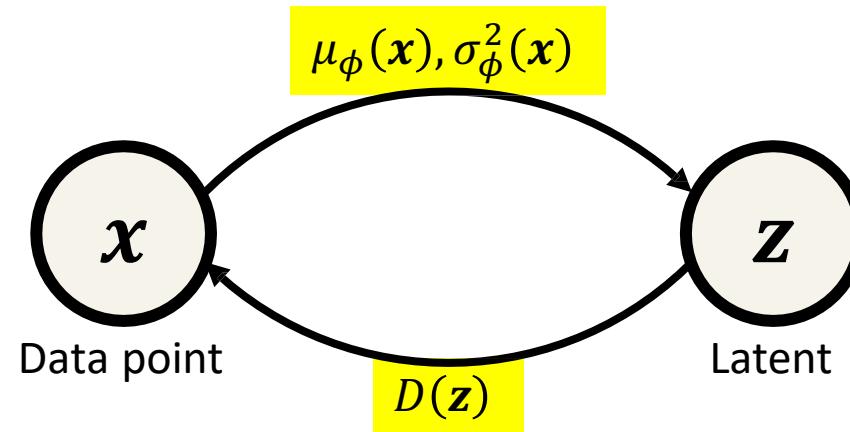
$$p(x) = \int p(x, z) dz = \int p_{\theta}(x|z)p(z) dz$$

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



Variational Autoencoder (VAE)

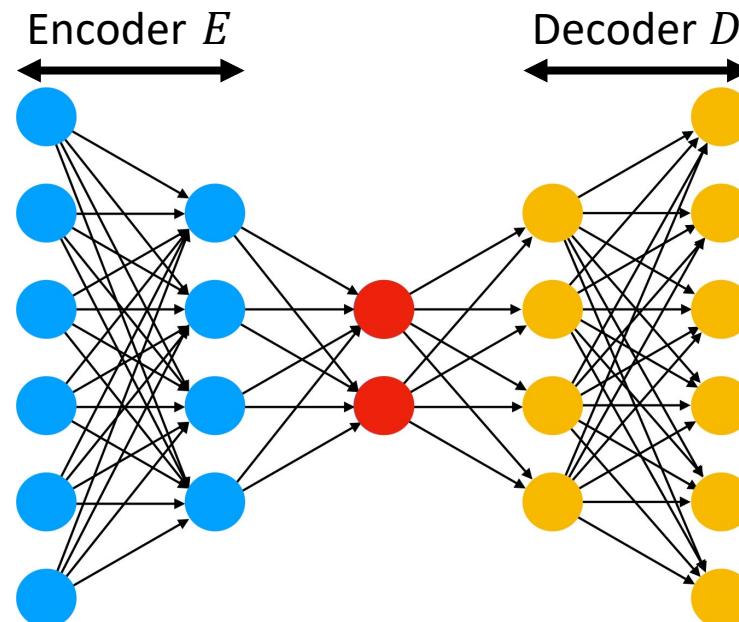
- How to maximize $p(x)$? Maximize **ELBO**!
- We need the **proxy** distribution $q_\phi(z|x) \rightarrow \text{Encoder}$
- Let $q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x)\mathbf{I})$



Variational Autoencoder (VAE)

Same with autoencoder but,

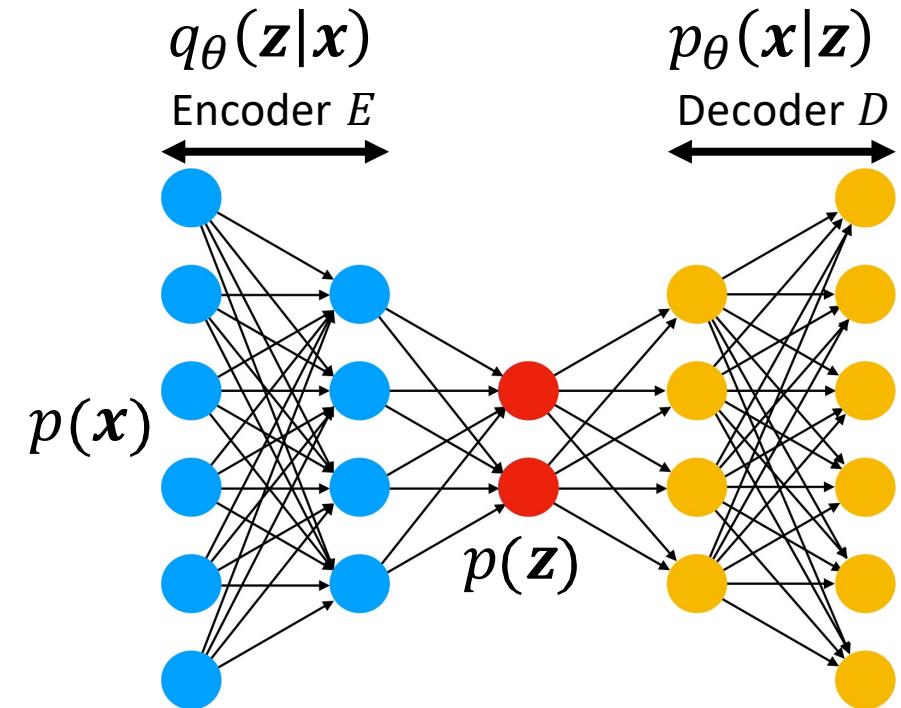
- From input x , the **encoder** predicts $\mu_\phi(x), \sigma_\phi^2(x)$.
- The **decoder** takes a sample $z \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x)\mathbf{I})$ as input.



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



Training

How to maximize **ELBO**?

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

Approximates using Monte Carlo estimate:

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

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

Training

How to maximize ELBO?

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

Approximates using Monte Carlo estimate:

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

Reparameterization Trick

where $\mathbf{z}^{(i)} \sim \mu_{\phi}(\mathbf{x}) + \sigma_{\phi}(\mathbf{x})\boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Training

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

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

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

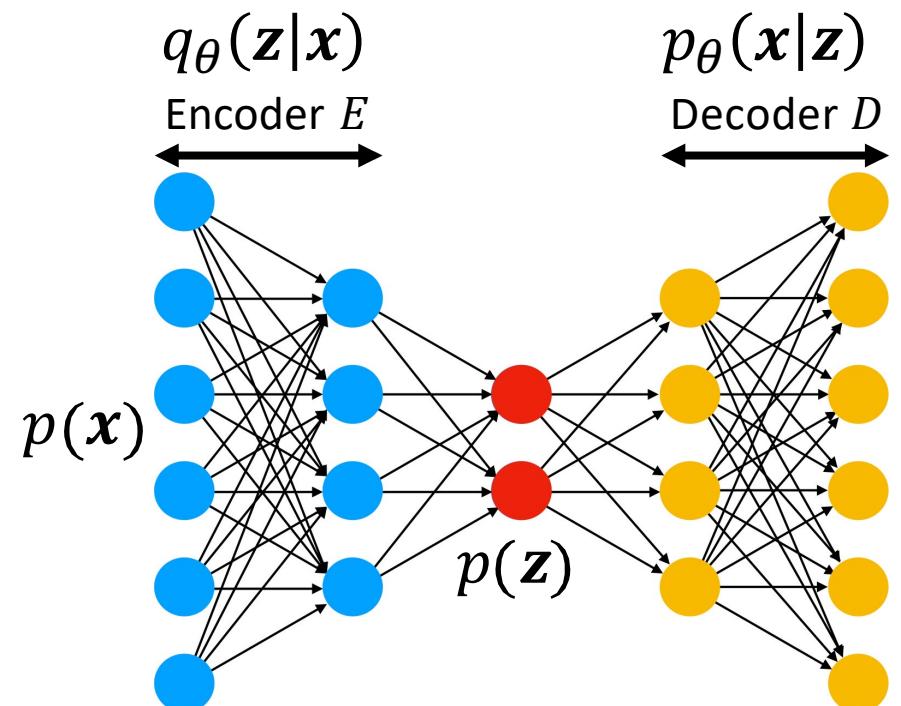
This is why it is called
the reconstruction term.

Constant

Training

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

Q. Why is the **sampling** differentiable?



Training

How to maximize **ELBO**?

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

Approximates using Monte Carlo estimate:

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

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

Training

How to maximize ELBO?

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

Approximates using Monte Carlo estimate:

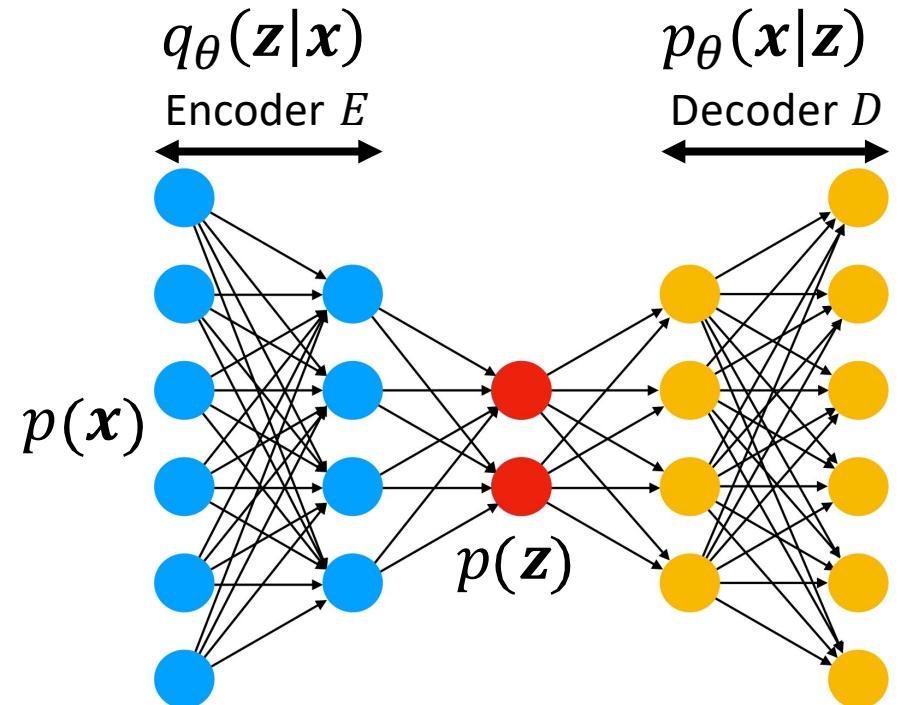
$$\operatorname{argmax}_{\theta, \phi} \frac{1}{N} \sum_i \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Reparameterization Trick

where $\mathbf{z}^{(i)} \sim \mu_{\phi}(\mathbf{x}) + \sigma_{\phi}(\mathbf{x})\boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Generation

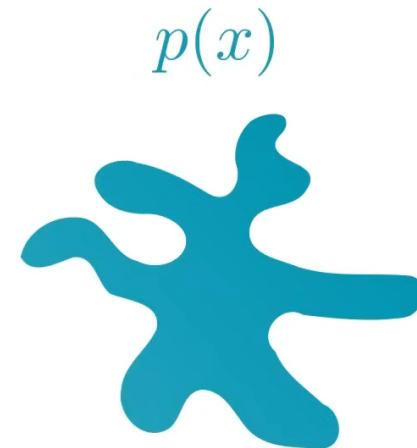
1. Sample a latent variable \mathbf{z} from $p(\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$.
2. Feed \mathbf{z} to the decoder to predict $\hat{\mathbf{x}} = D_\theta(\mathbf{z})$.



Dimensions

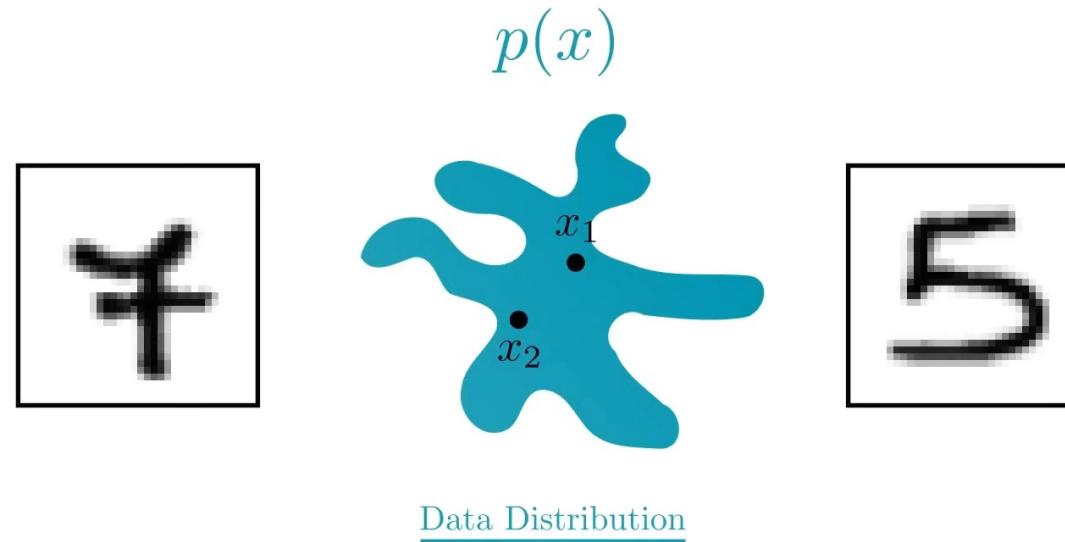
Q. Should the **dimensions** of the input data and the latent variables be the same?

Variational Autoencoders (VAE)

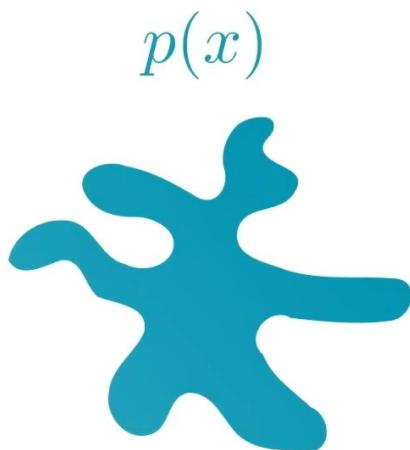


Data Distribution

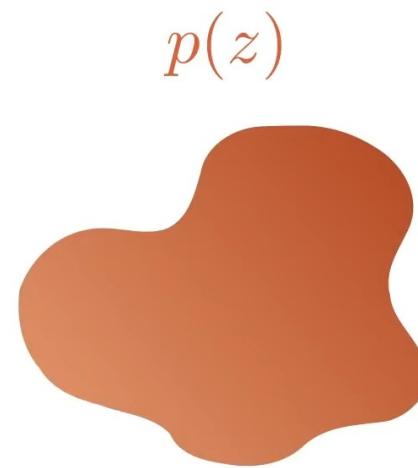
Variational Autoencoders (VAE)



Variational Autoencoders (VAE)

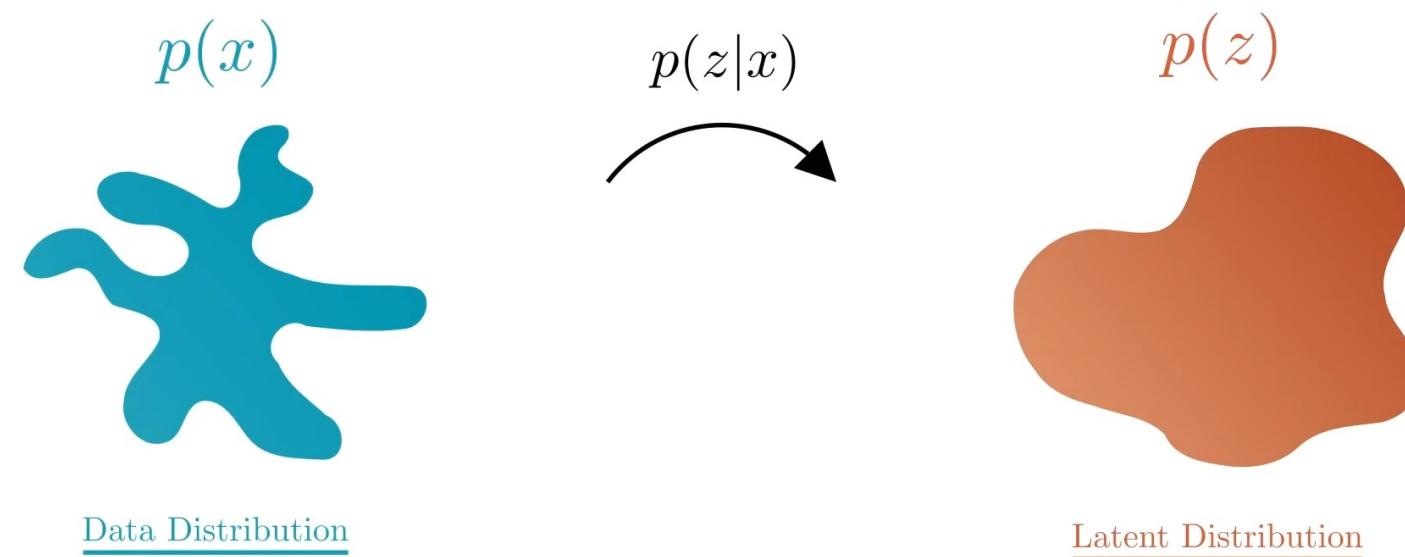


Data Distribution

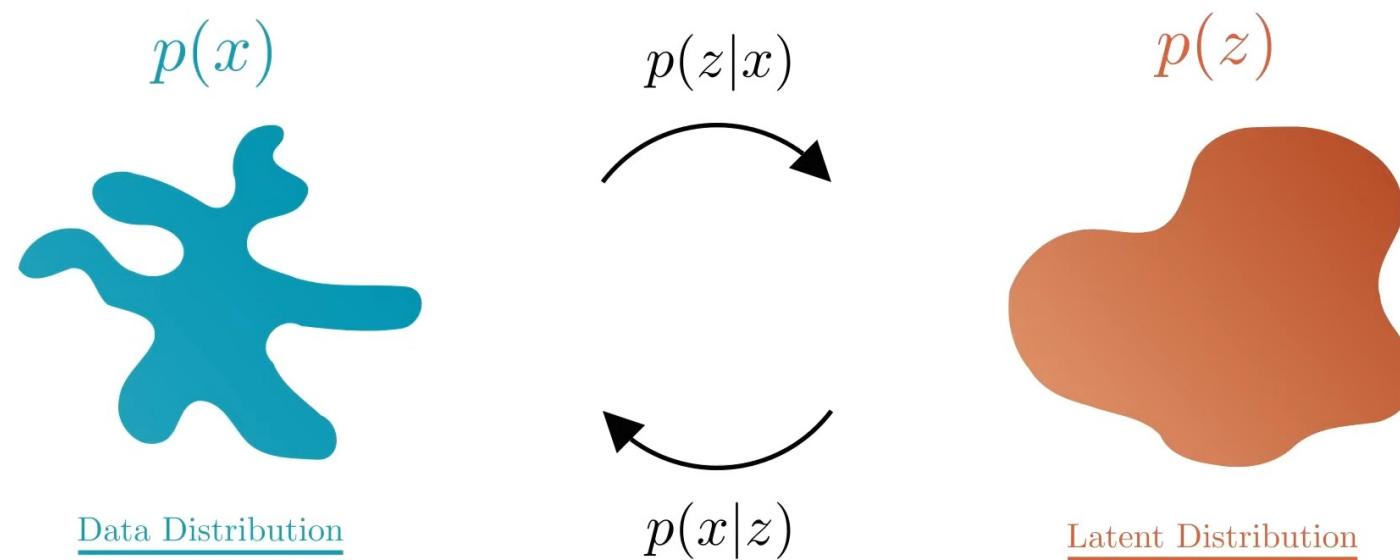


Latent Distribution

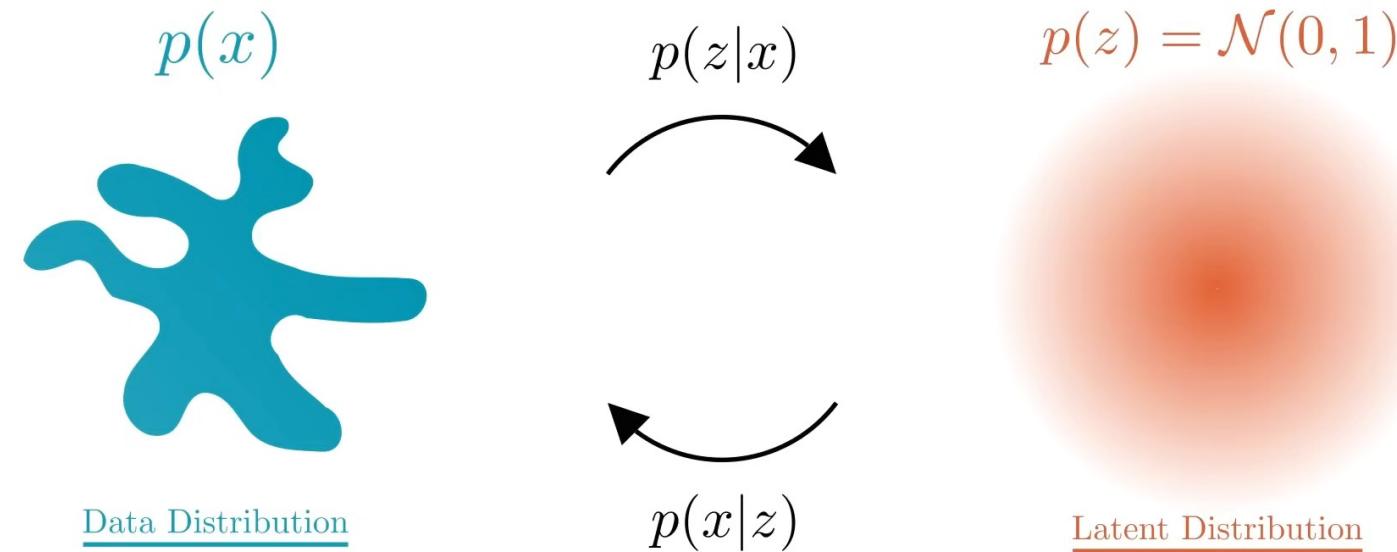
Variational Autoencoders (VAE)



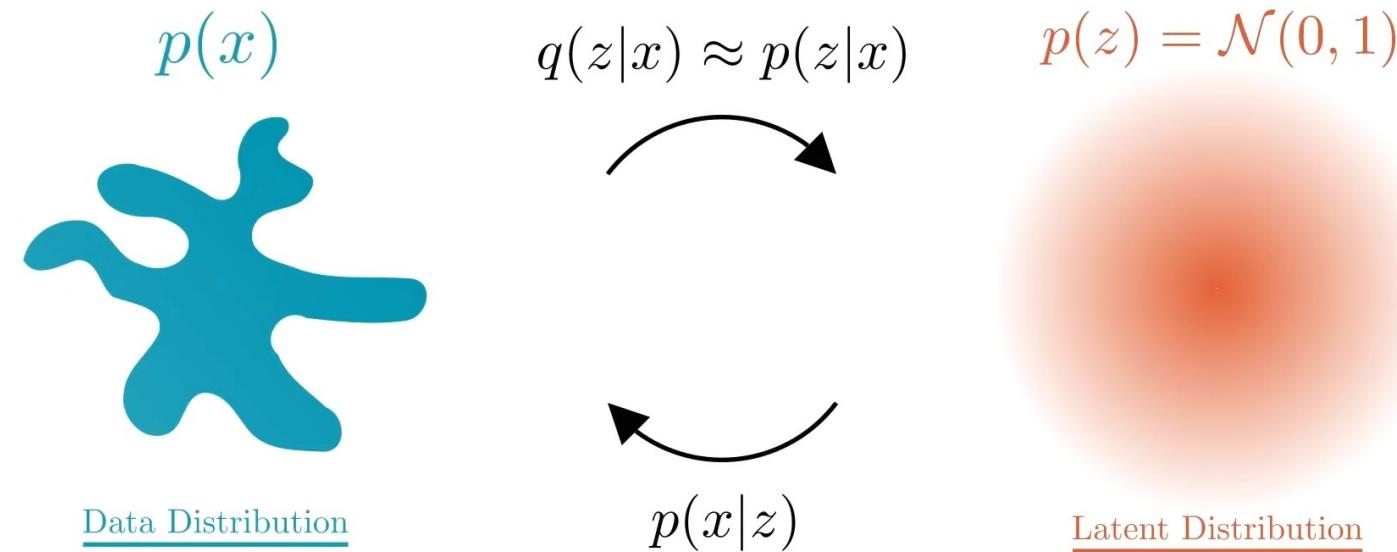
Variational Autoencoders (VAE)



Variational Autoencoders (VAE)

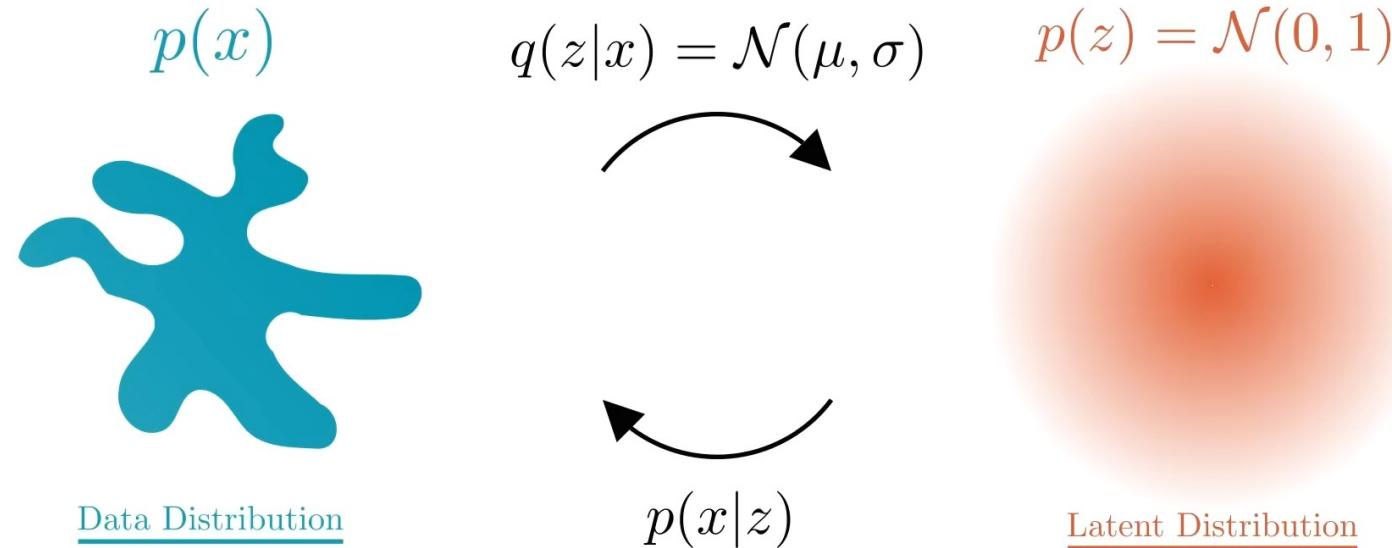


Variational Autoencoders (VAE)



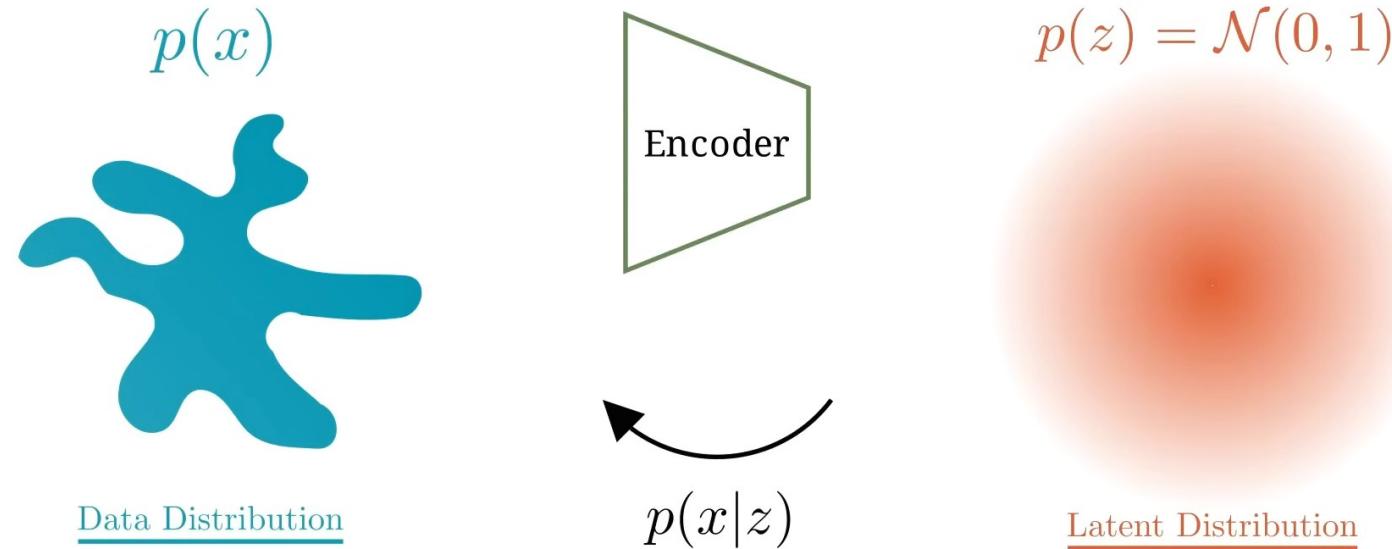
Variational Autoencoders (VAE)

Variational Bayes



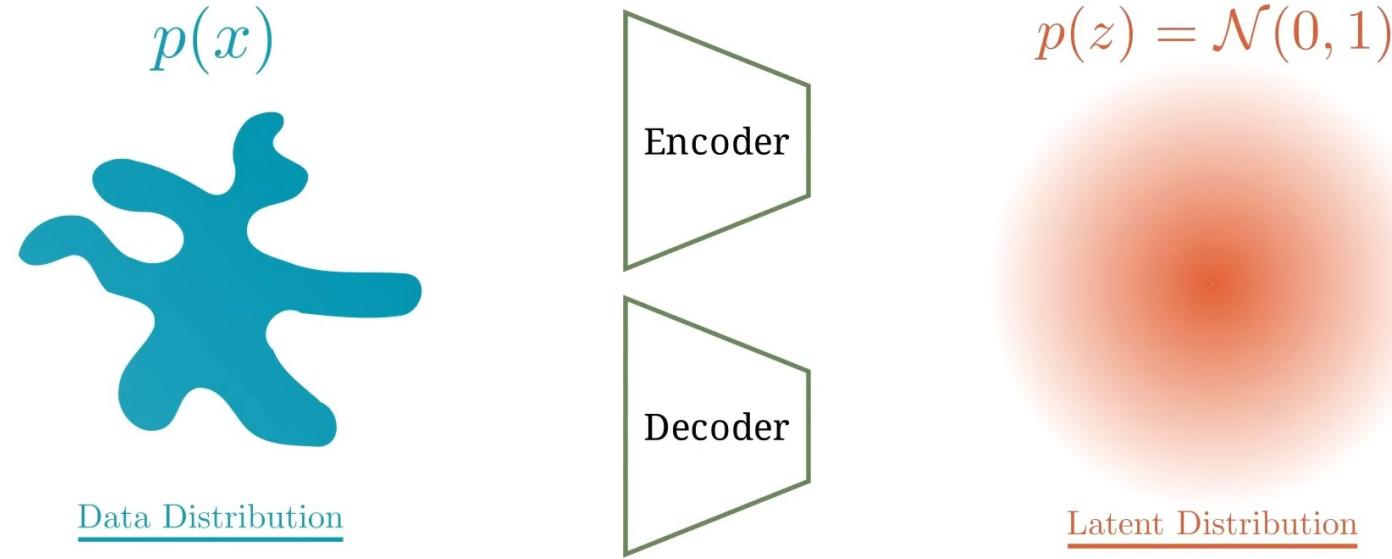
Variational Autoencoders (VAE)

Autoencoder Variational Bayes



Variational Autoencoders (VAE)

Autoencoder Variational Bayes



Variational Autoencoders (VAE)

How do we train this autoencoder ?

Variational Autoencoders (VAE)

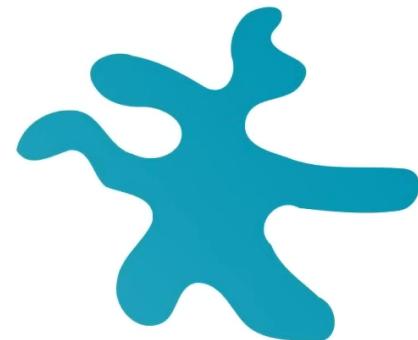
$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$

Variational Autoencoders (VAE)

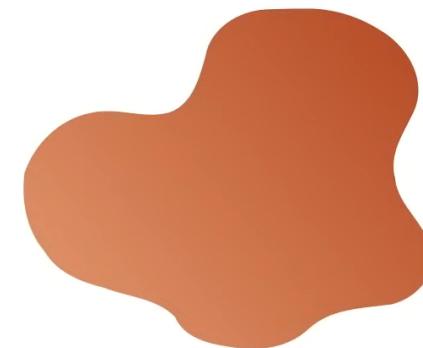
$$\mathcal{L}(x) = \underbrace{\mathbb{E}_{q(z|x)} [\log p(x|z)]}_{\text{Data consistency}} - \overbrace{\text{KL}(q(z|x) \mid p(z))}^{\text{Regularization}}$$

Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$



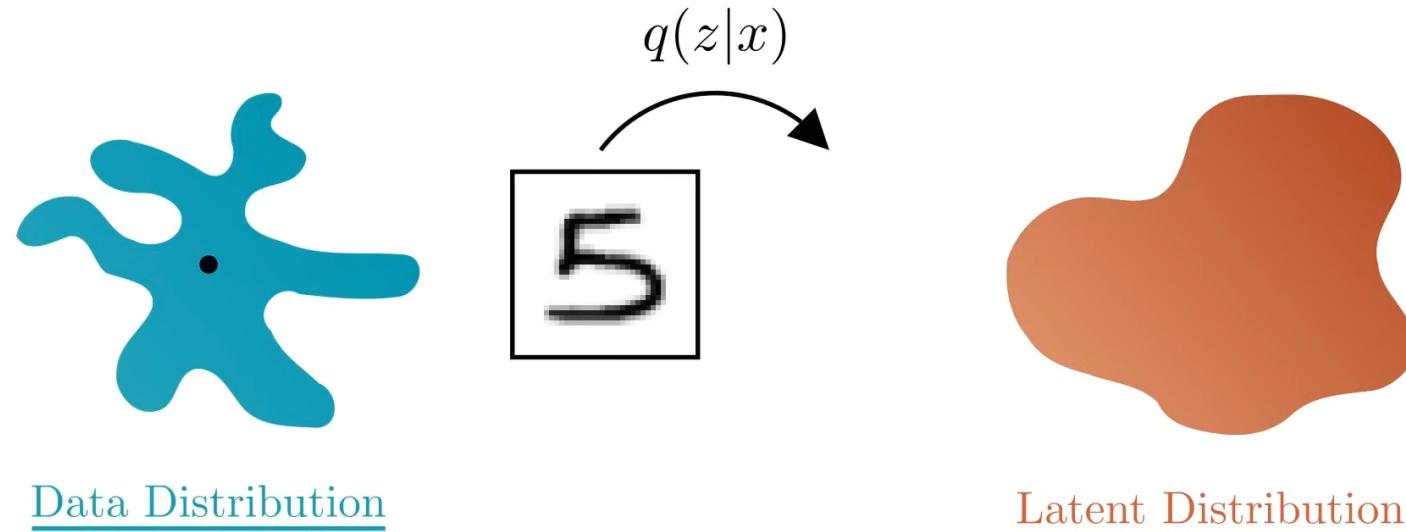
Data Distribution



Latent Distribution

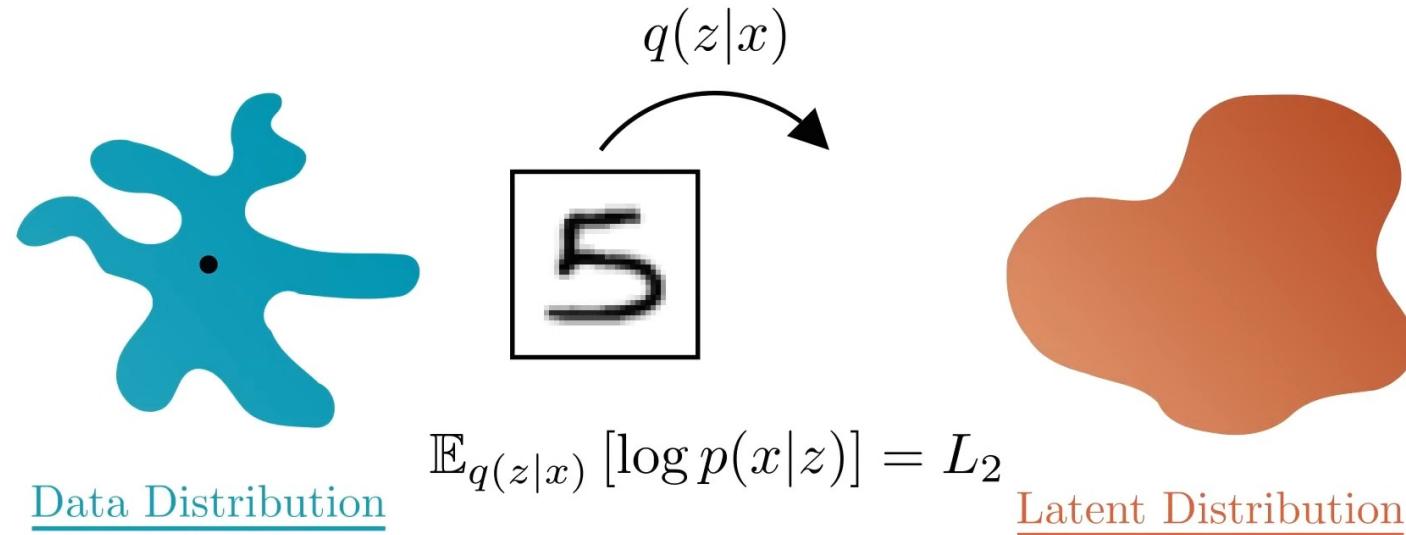
Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$



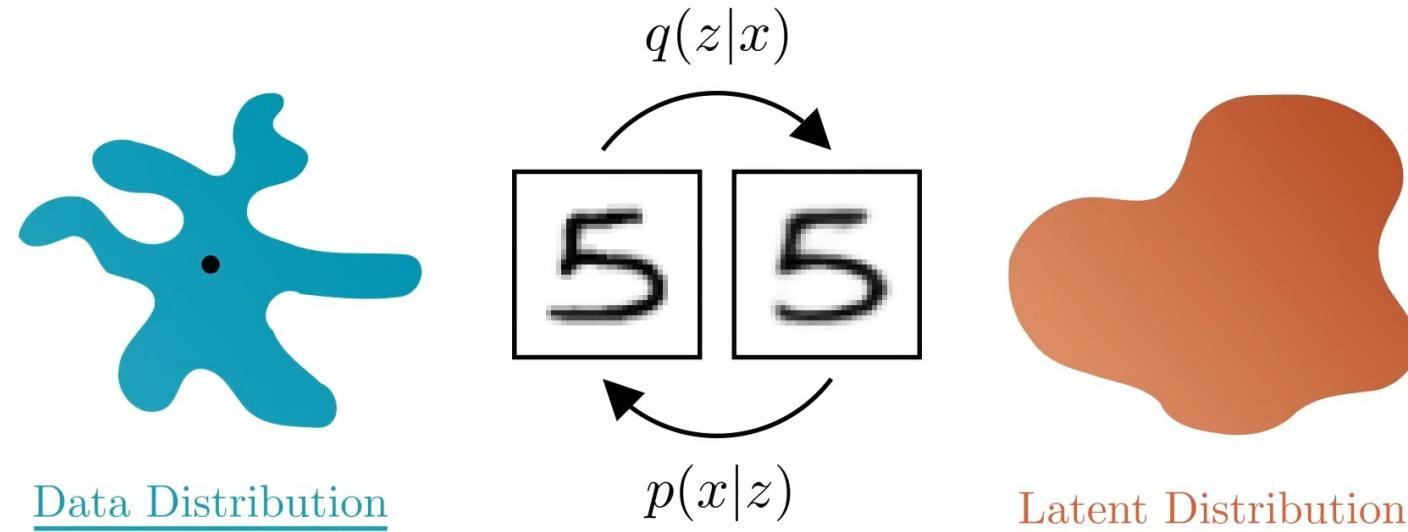
Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$



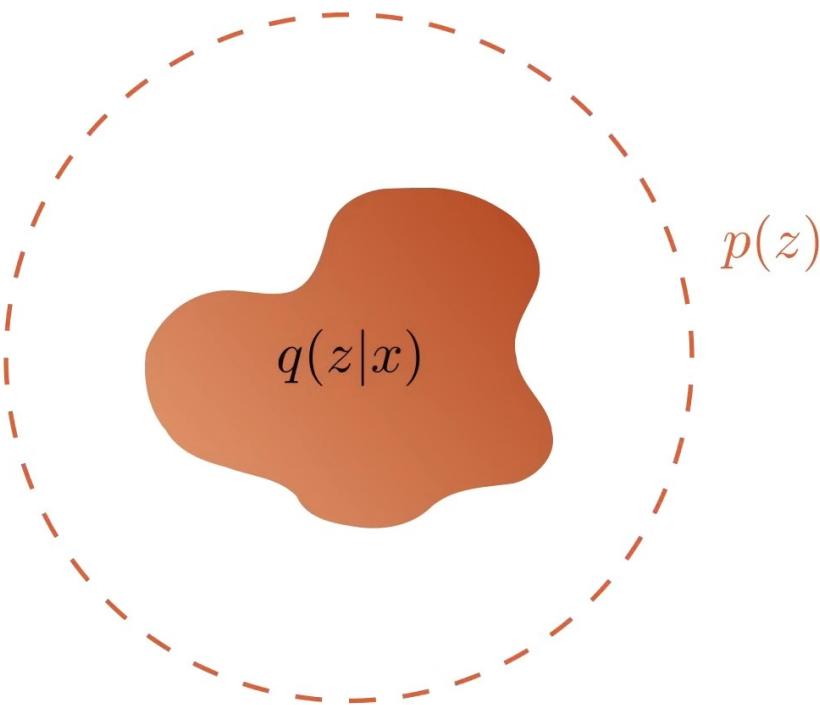
Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$



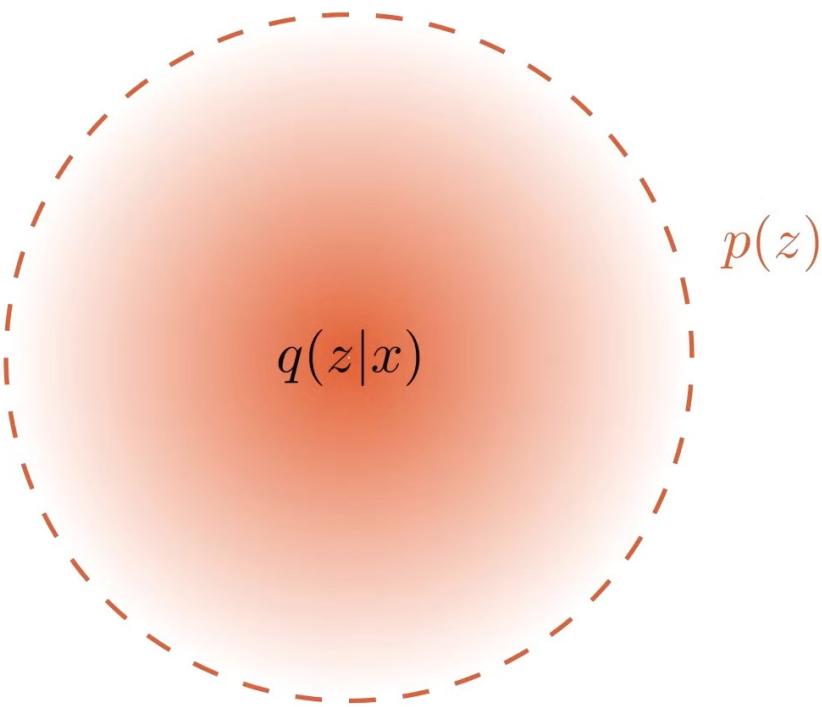
Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$



Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$



Variational Autoencoders (VAE)

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) \mid p(z))$$

Variational Autoencoders (VAE)

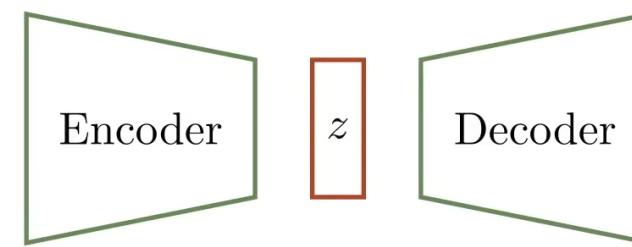
Latent space regularization

$$\mathcal{L}(x) = \underbrace{\mathbb{E}_{q(z|x)} [\log p(x|z)]}_{\text{L2}} - \overbrace{\text{KL}(q(z|x) \mid p(z))}^{\text{Latent space regularization}}$$

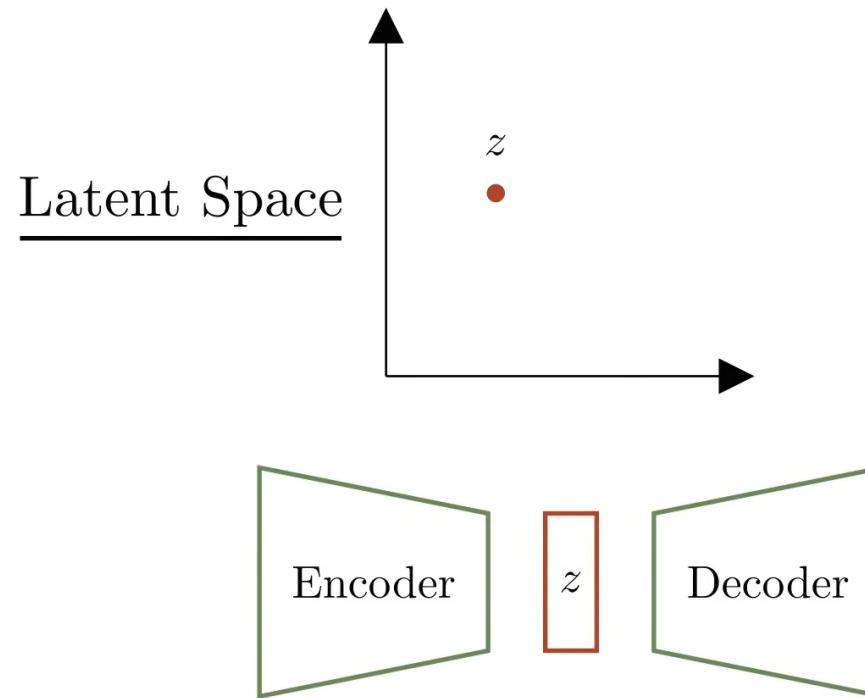
Variational Autoencoders (VAE)

How do we do in practice ?

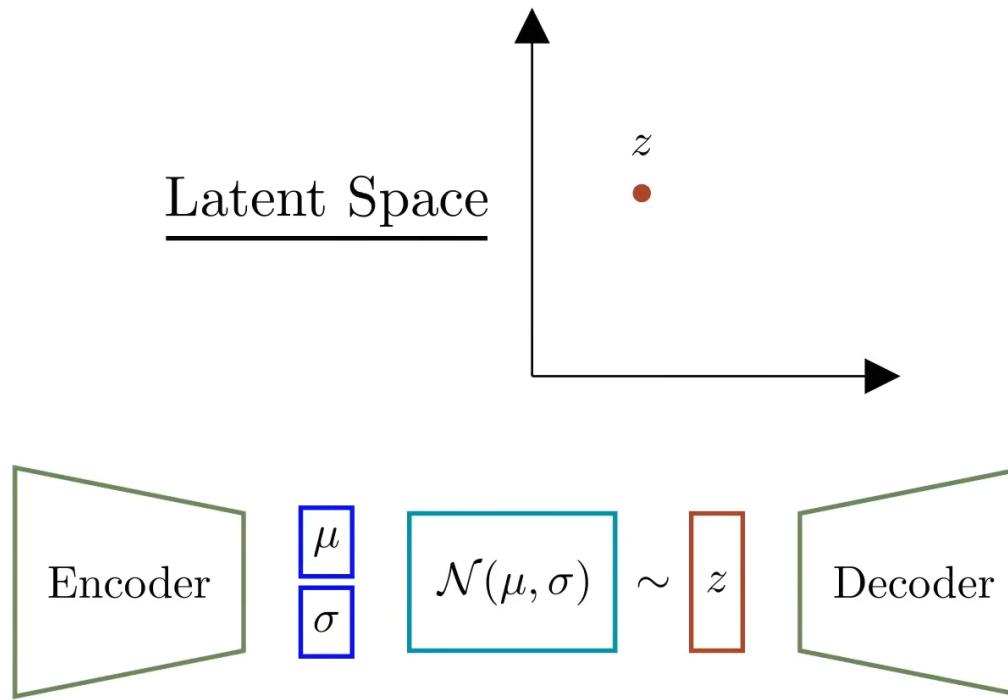
Variational Autoencoders (VAE)



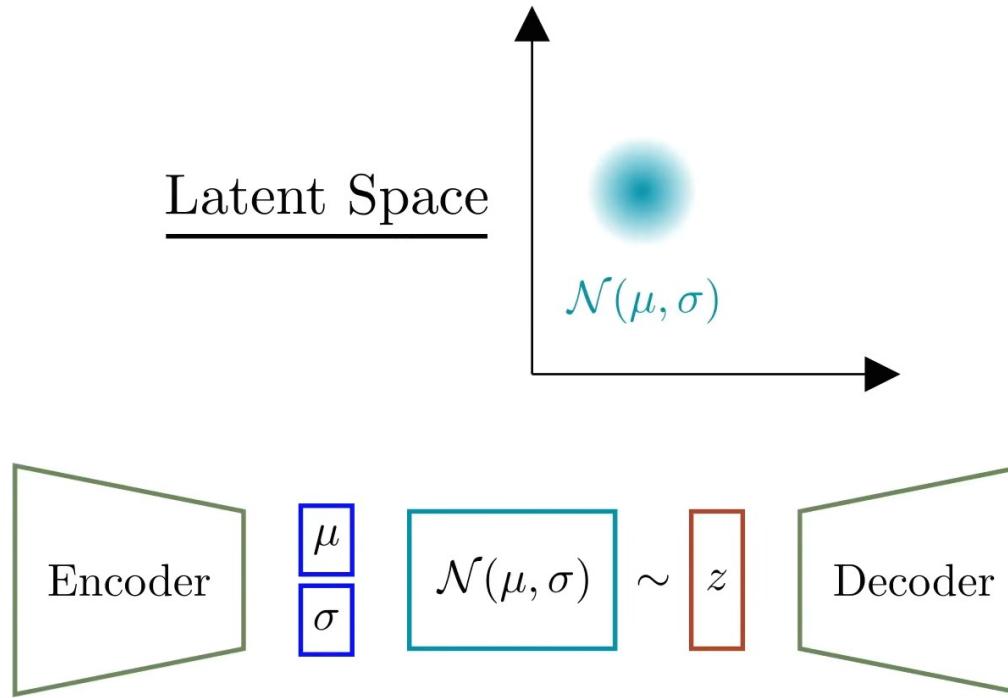
Variational Autoencoders (VAE)



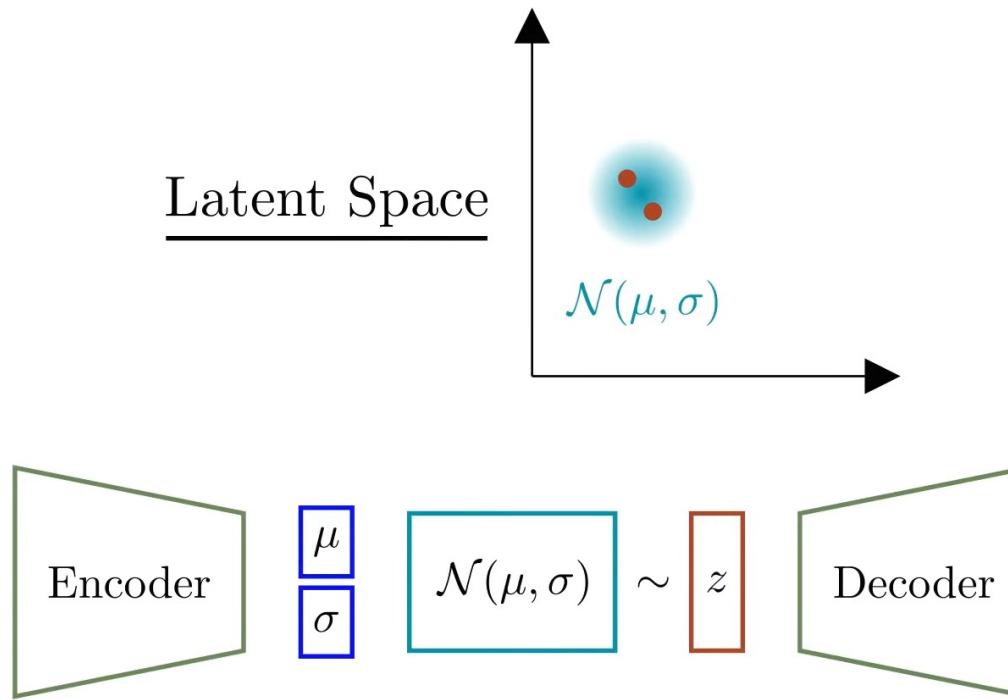
Variational Autoencoders (VAE)



Variational Autoencoders (VAE)

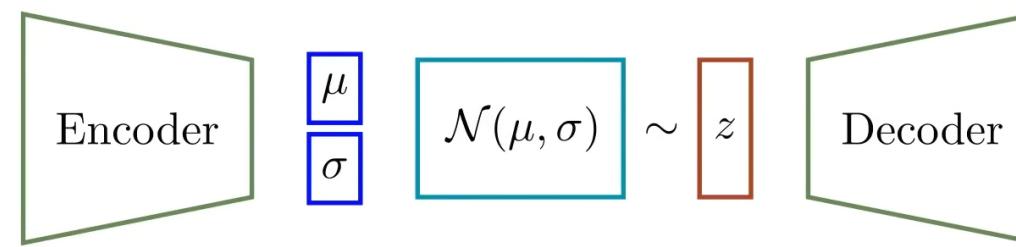


Variational Autoencoders (VAE)



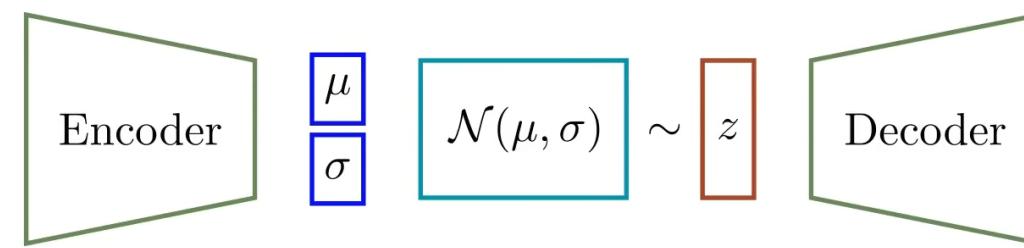
Variational Autoencoders (VAE)

$$\mathcal{L}(x, x') = L_2(x, x') + D_{KL}(\mathcal{N}(\mu, \sigma) \mid \mathcal{N}(0, 1))$$



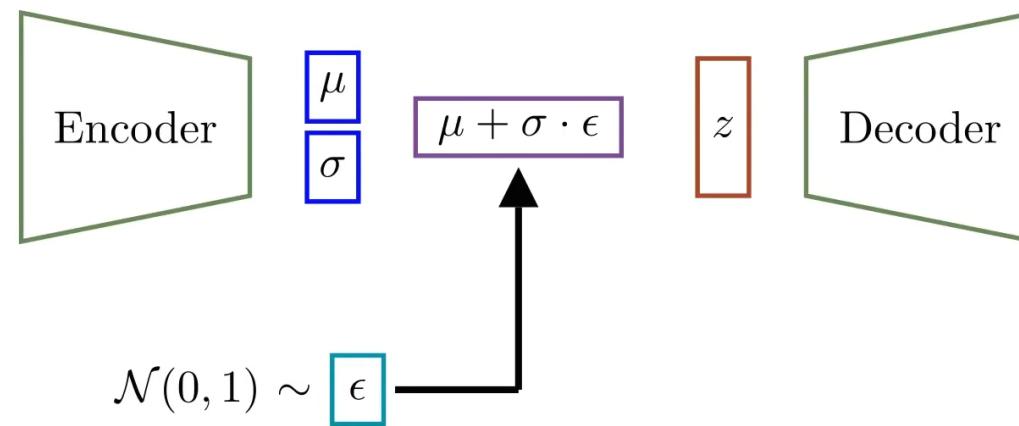
Variational Autoencoders (VAE)

How to backpropagate through the sampling process?

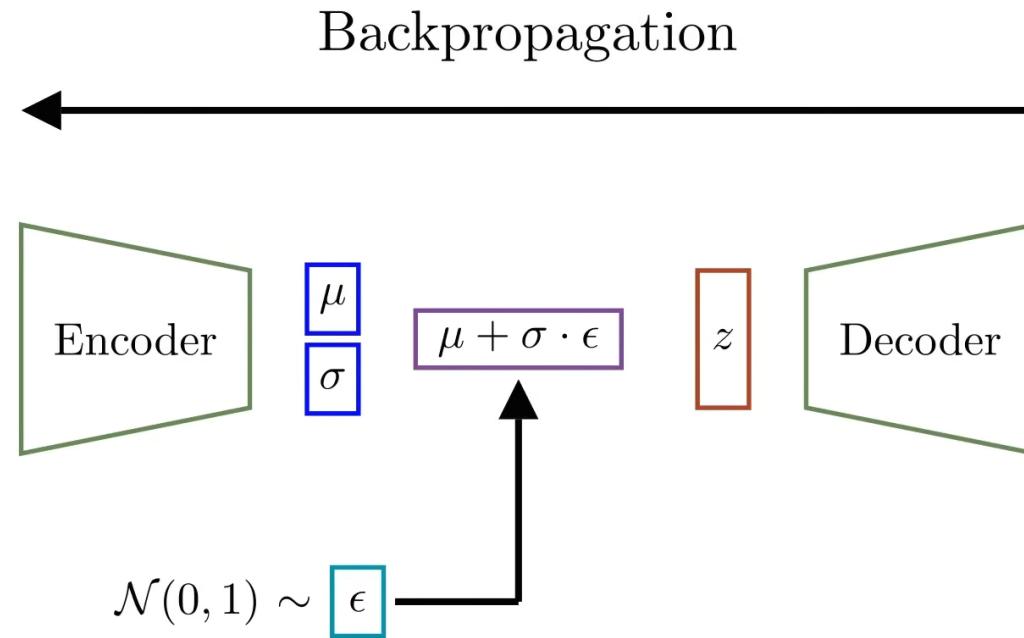


Variational Autoencoders (VAE)

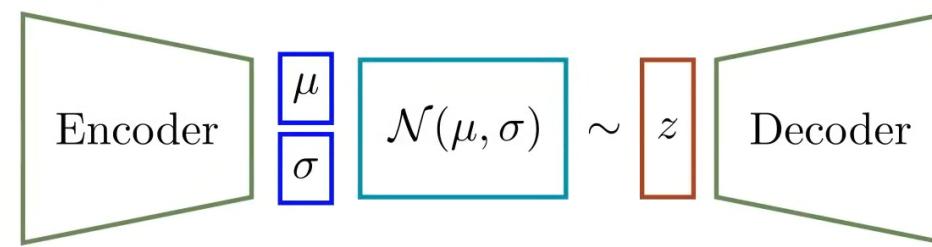
Reparameterization Trick



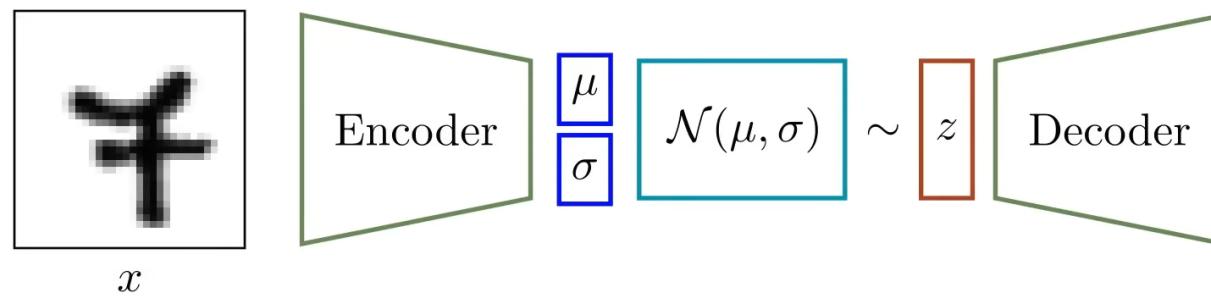
Variational Autoencoders (VAE)



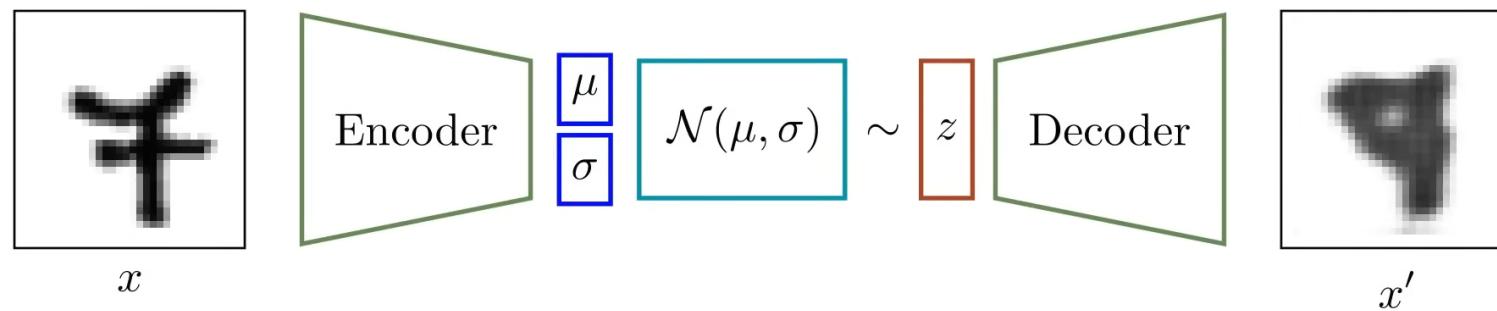
Variational Autoencoders (VAE)



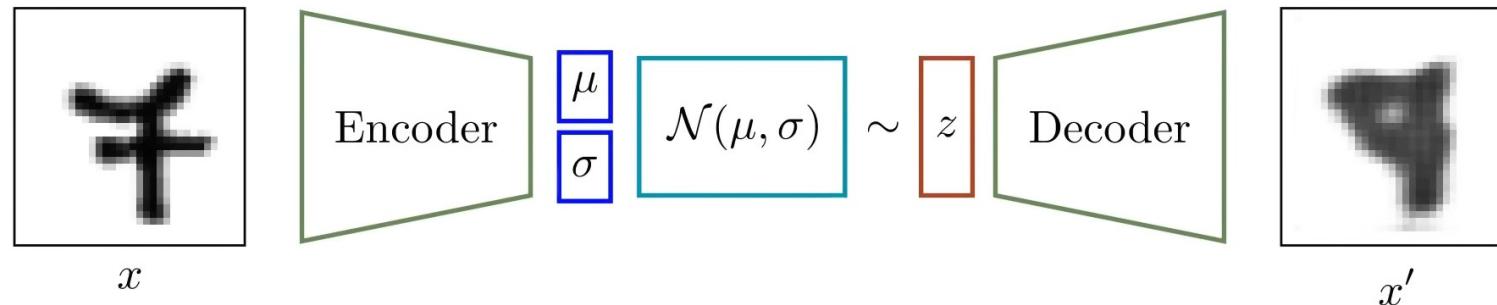
Variational Autoencoders (VAE)



Variational Autoencoders (VAE)

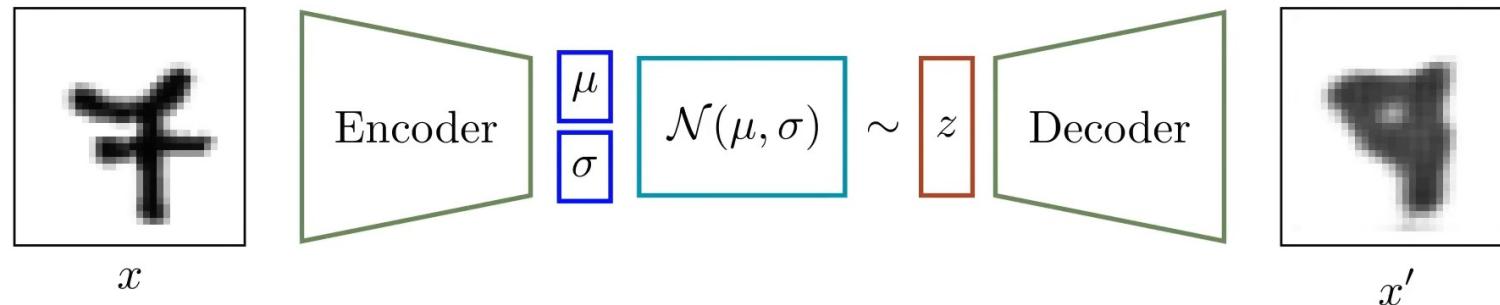


Variational Autoencoders (VAE)



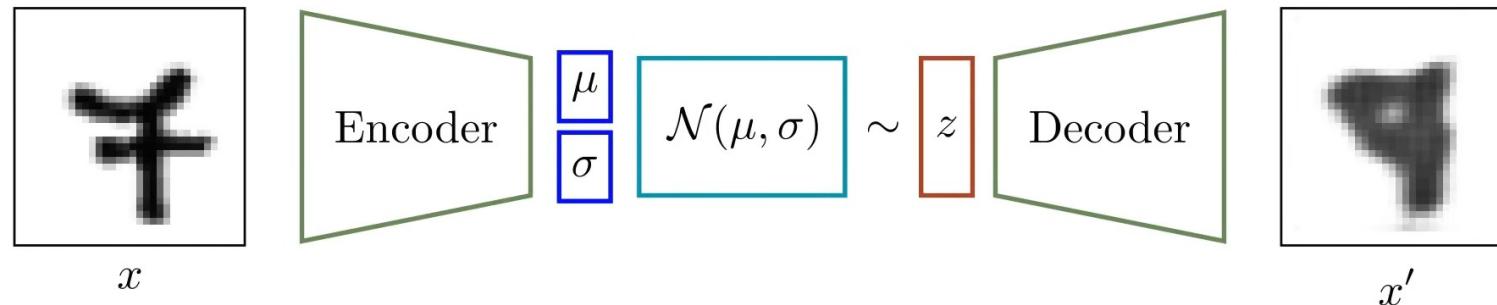
$$\mathcal{L} = \mathcal{L}_{KL}(\mathcal{N}(\mu, \sigma) \mid \mathcal{N}(0, 1)) + \boxed{\mathcal{L}_2(x, x')}$$

Variational Autoencoders (VAE)



$$\mathcal{L} = [\mathcal{L}_{KL}(\mathcal{N}(\mu, \sigma) \mid \mathcal{N}(0, 1)) + \mathcal{L}_2(x, x')]$$

Variational Autoencoders (VAE)



$$\mathcal{L} = [\mathcal{L}_{KL}(\mathcal{N}(\mu, \sigma) \mid \mathcal{N}(0, 1)) + \mathcal{L}_2(x, x')]$$

$$\mathcal{L}_{KL} = -\frac{1}{2}(1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$

Limitations of VAEs

Typical failure cases of VAEs

Real
images



Generated
Images

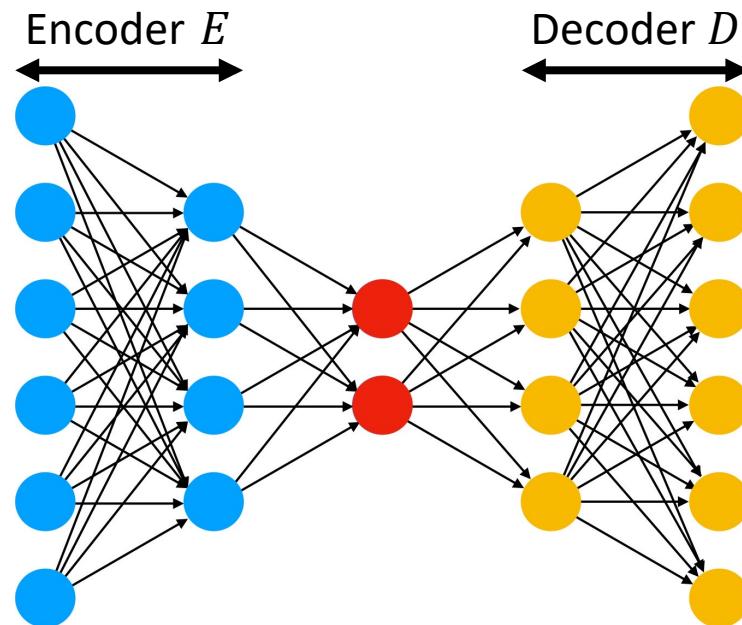


Limitations of VAEs

Is a **Gaussian** distribution sufficient as the variational approximation for the posterior distribution?

$$q_{\theta}(z|x)$$

$$= \mathcal{N}(z; \mu_{\theta}(x), \sigma_{\theta}^2(x)\mathbf{I})$$



Limitations of VAE

We maximize *not* $\log p(\mathbf{x})$ but its **lower bound** (ELBO).

Q. What is the difference between the two?

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

Limitations of VAE

$$\text{A. } \log p(x) - \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(x,z)}{q_\phi(z|x)} \right] = D_{KL}(q_\phi(z|x) || p(z|x)).$$

- The lower bound becomes **tight** when the variational distribution $q_\phi(z|x)$ is **identical** to the true posterior distribution $p(z|x)$.
- Will the true posterior distribution $p(z|x)$ be close to a Gaussian distribution...?

Limitations of VAE

What is a better method for approximating the posterior distribution in a variational way?

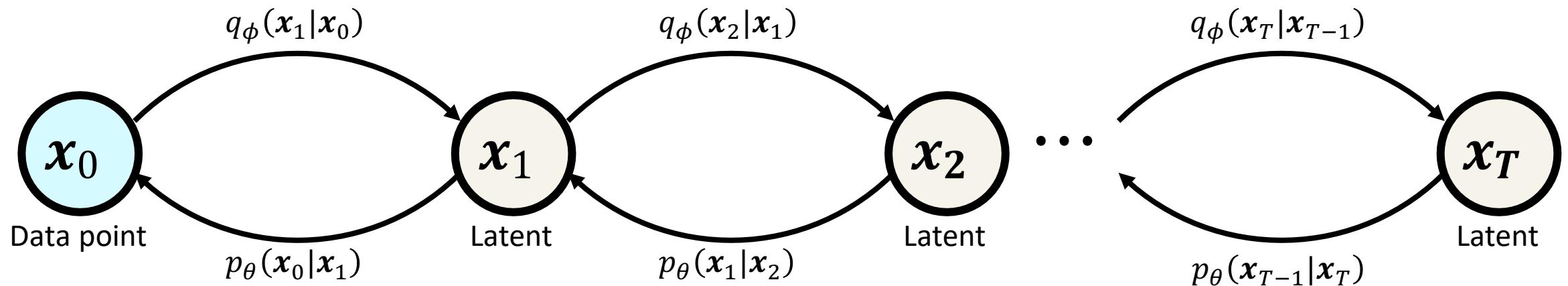
VAE Variants

- Vector-Quantized VAE (VQ-VAE) [Oord et al., NeurIPS 2017]
- Beta-VAE [Higgins et al., ICLR 2017]
- Wasserstein VAE (WAE) [Tolstikhin et al., ICLR 2018]
- VAE-GAN [Larsen et al., ICML 2016]
- Normalizing Flow VAE [Rezende and Mohamed, ICML 2015]

Hierarchical VAEs

Make a recursive (hierarchical) VAE.

- Data point: $x \rightarrow \textcolor{brown}{x}_0$
- Latent variable(s): $z \rightarrow \textcolor{brown}{x}_{1:T}$



Markovian Hierarchical VAEs

Let's consider a **Markovian** process.

A Markov process is a stochastic process where the probability of each event **depends only on the previous** state.

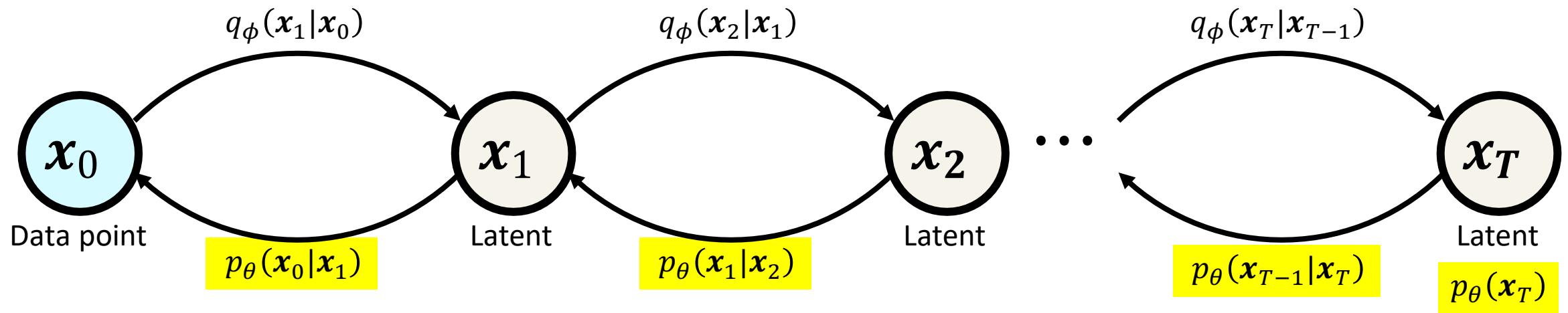
“What happens next depends only on the state of affairs now!”

Markovian Hierarchical VAEs

Let's consider a Markovian process.

Joint distribution:

$$p_{\theta}(\mathbf{x}_{0:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$



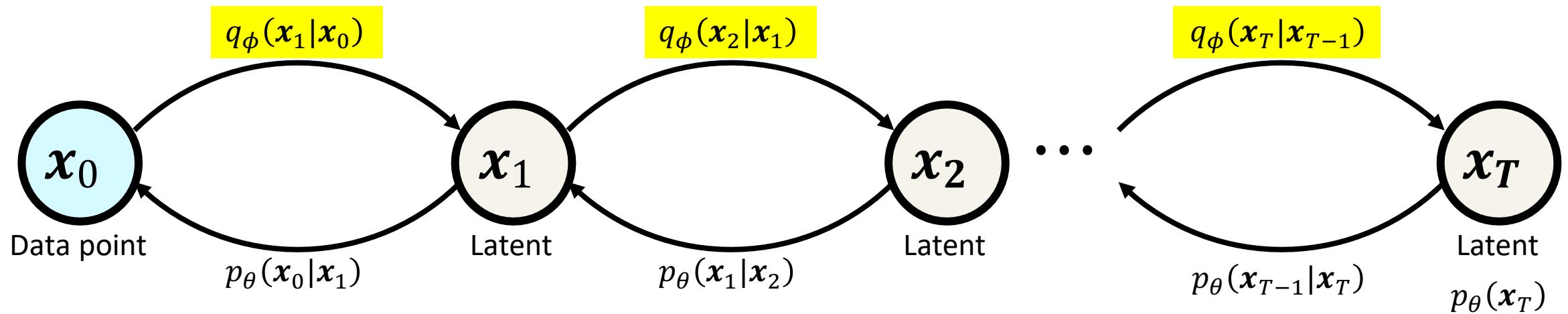
Markovian Hierarchical VAEs

Let's consider a Markovian process.

Variational posterior:

While each $q_\phi(x_{t+1}|x_t)$ is a normal distribution,
 $q_\phi(x_T|x_0)$ can be a more complex distribution.

$$q_\phi(x_{1:T}|x_0) = \prod_{t=1}^T q_\phi(x_t|x_{t-1})$$



Markovian Hierarchical VAEs

$$\begin{aligned}\log p(\mathbf{x}_0) &= \log \int p_{\theta}(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\ &= \log \int p_{\theta}(\mathbf{x}_{0:T}) \frac{q_{\phi}(\mathbf{x}_{1:T} | \mathbf{x}_0)}{q_{\phi}(\mathbf{x}_{1:T} | \mathbf{x}_0)} d\mathbf{x}_{1:T} \\ &= \log \mathbb{E}_{q_{\phi}(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\frac{p_{\theta}(\mathbf{x}_{0:T})}{q_{\phi}(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \\ &\geq \mathbb{E}_{q_{\phi}(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q_{\phi}(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right]\end{aligned}$$

VAEs → Diffusion Models

Flow-Based Models

- Normalizing flow
- Nonlinear Independent Components Estimation (NICE)
- Real Non-Volume Preserving (Real NVP)
- Generative Flow (Glow)
- Masked autoregressive flow (MAF)
- Continuous Normalizing Flow (CNF)

We'll revisit this later!