

COMP0169: Machine Learning for Visual Computing

Generative Modeling



Lectures will be Recorded

Recap

MLP

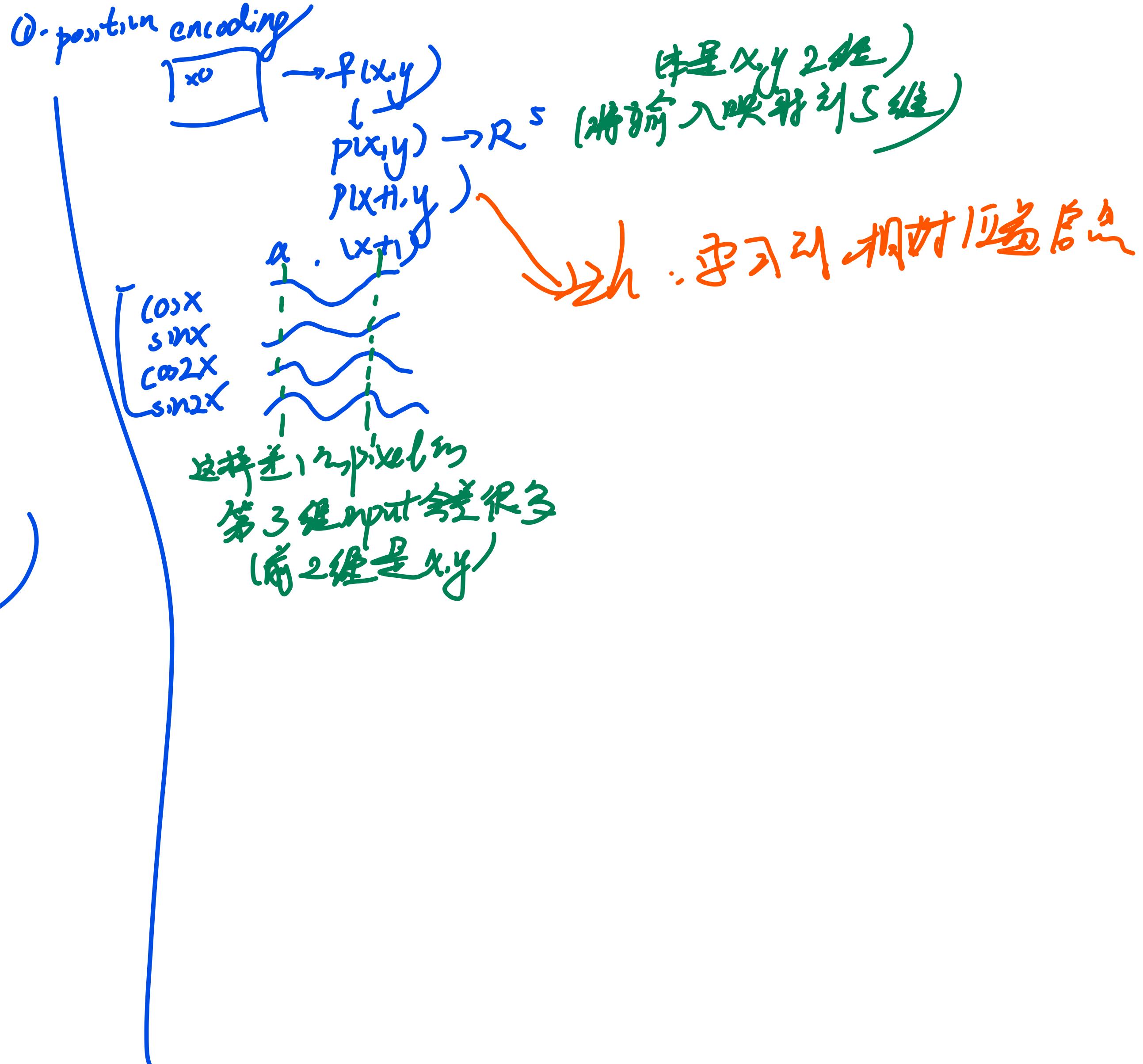
CNN

visualize

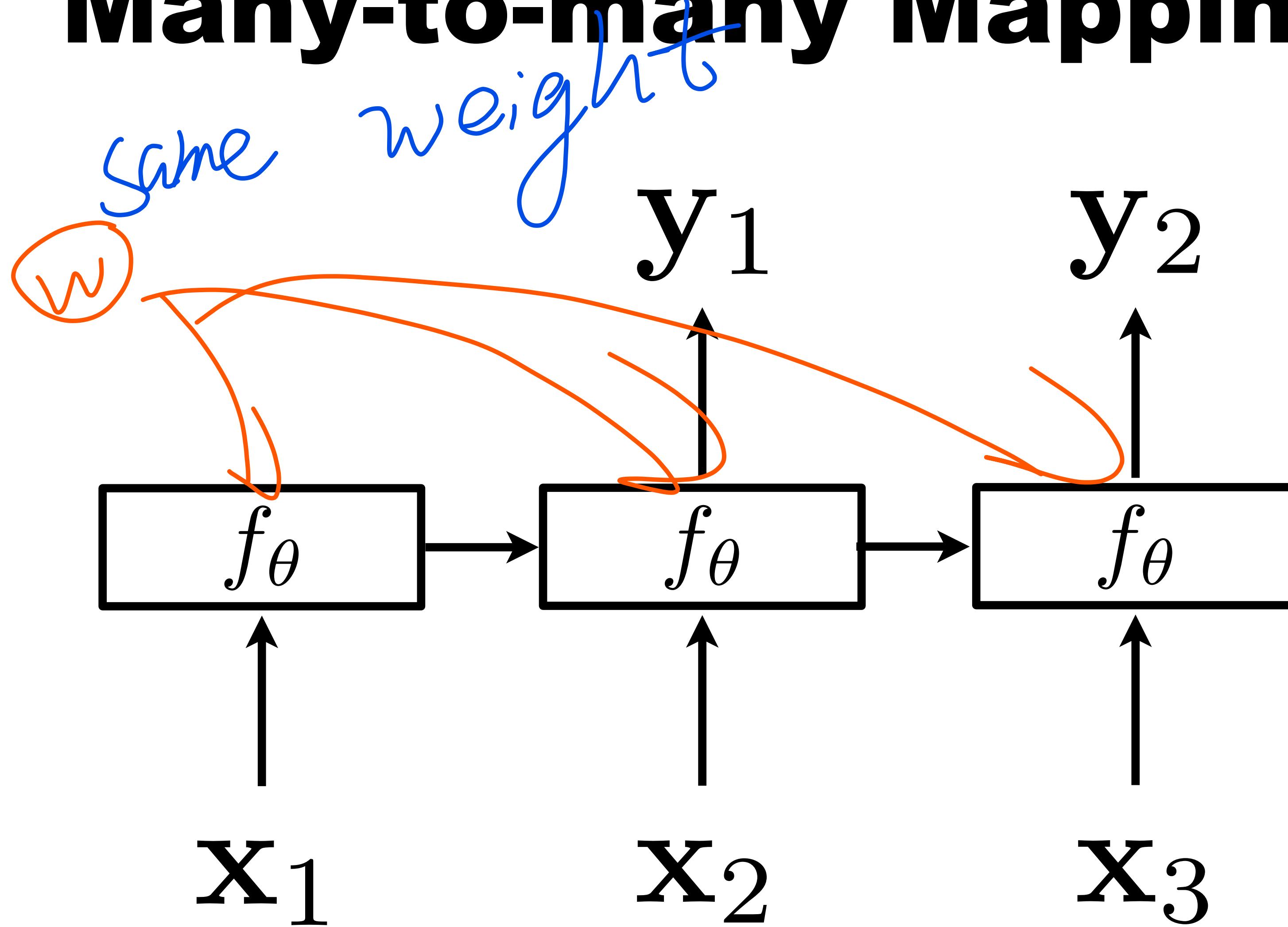
loss

sequence

CNN

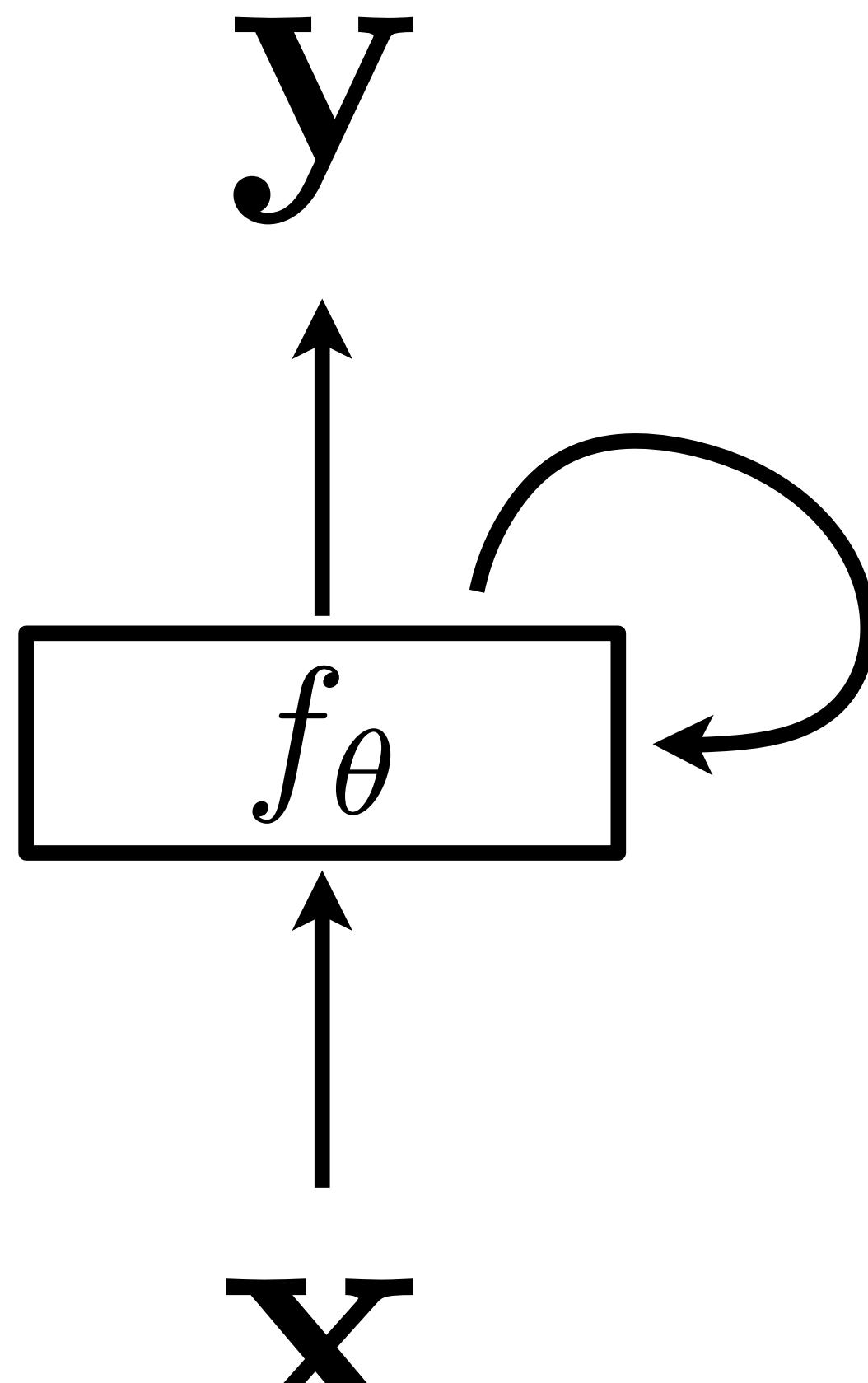


Many-to-many Mapping



- **Conditional Image Generation**
- **Machine Translation**
- **Animation Synthesis**
- **Skeleton to body animation**

Recurrent Neural Network (RNN)



$$h_t = f_W(h_{t-1}, x_t)$$
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$
$$y_t = W_{yh}h_t$$

Further reading: <https://pabloinsente.github.io/the-recurrent-net>



Cognitive Science

Volume 14, Issue 2, April–June 1990, Pages 179–211



Finding structure in time ★

Jeffrey L. Elman

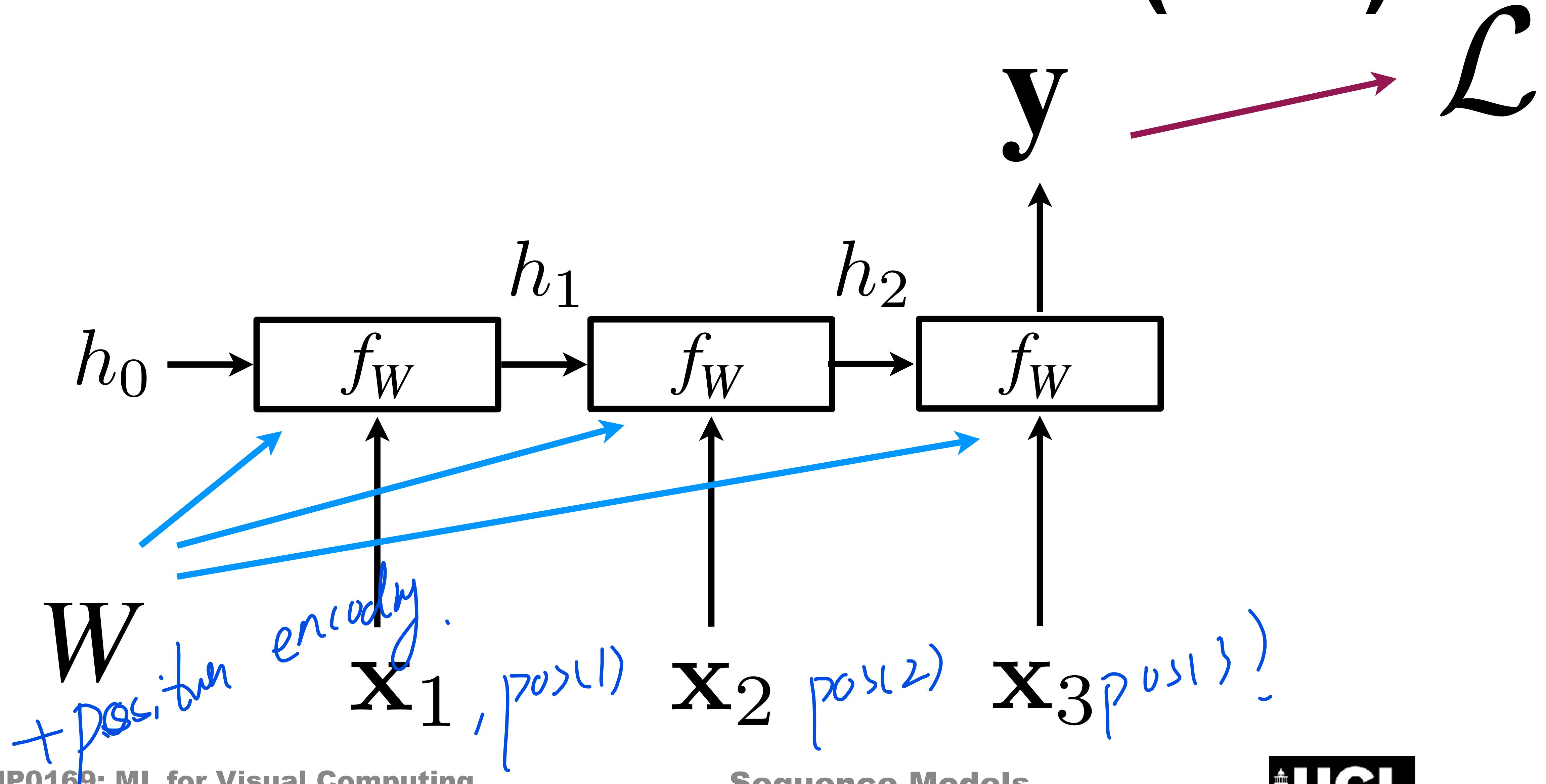
Show more

+ Add to Mendeley Share Cite

[https://doi.org/10.1016/0364-0213\(90\)90002-E](https://doi.org/10.1016/0364-0213(90)90002-E)

Get rights and content

Recurrent Neural Network (RNN)



Capturing Longterm Interactions

$h g / X_i \rightarrow [\dots]$

$X_j \rightarrow [\dots]$

$\|X_i - X_j\|^2$ - 直是 2
在其中一个字

solve

$Ex \rightarrow IR$

cat dog --- meitten

$dis(E_{cat}, E_{dog}) < dis(E_{cat}, E_{meitten})$

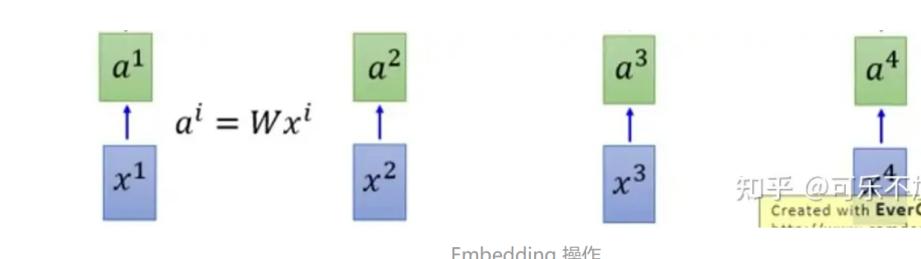
Embedding

Sequence Models

Embedding
(-d维 representation)

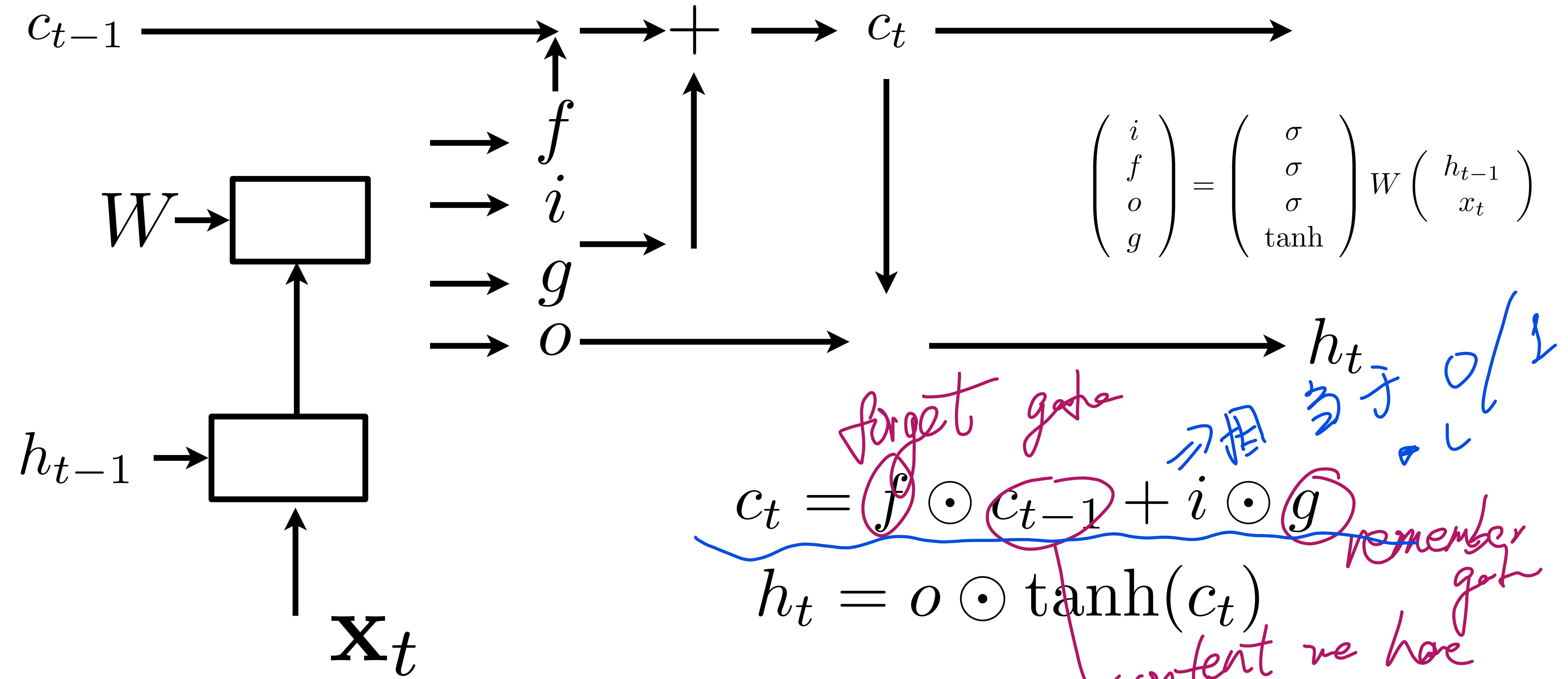
link cos, sim (x1, x2)

2.1 Embedding 操作
首先, 先对“你好机车”这段话进行Embedding操作(可以使用Word2Vec等方法), 得到新的向量, a^1, a^2, a^3, a^4 ; 如下图所示。

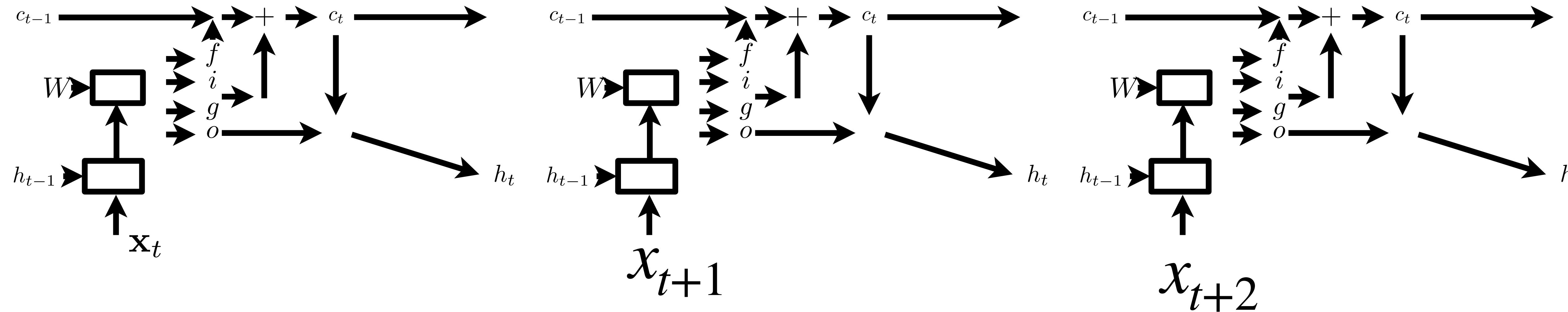


(E: word to point in
high dimension)
但这里一语双关
还是没有 context relevant
past for end. 1

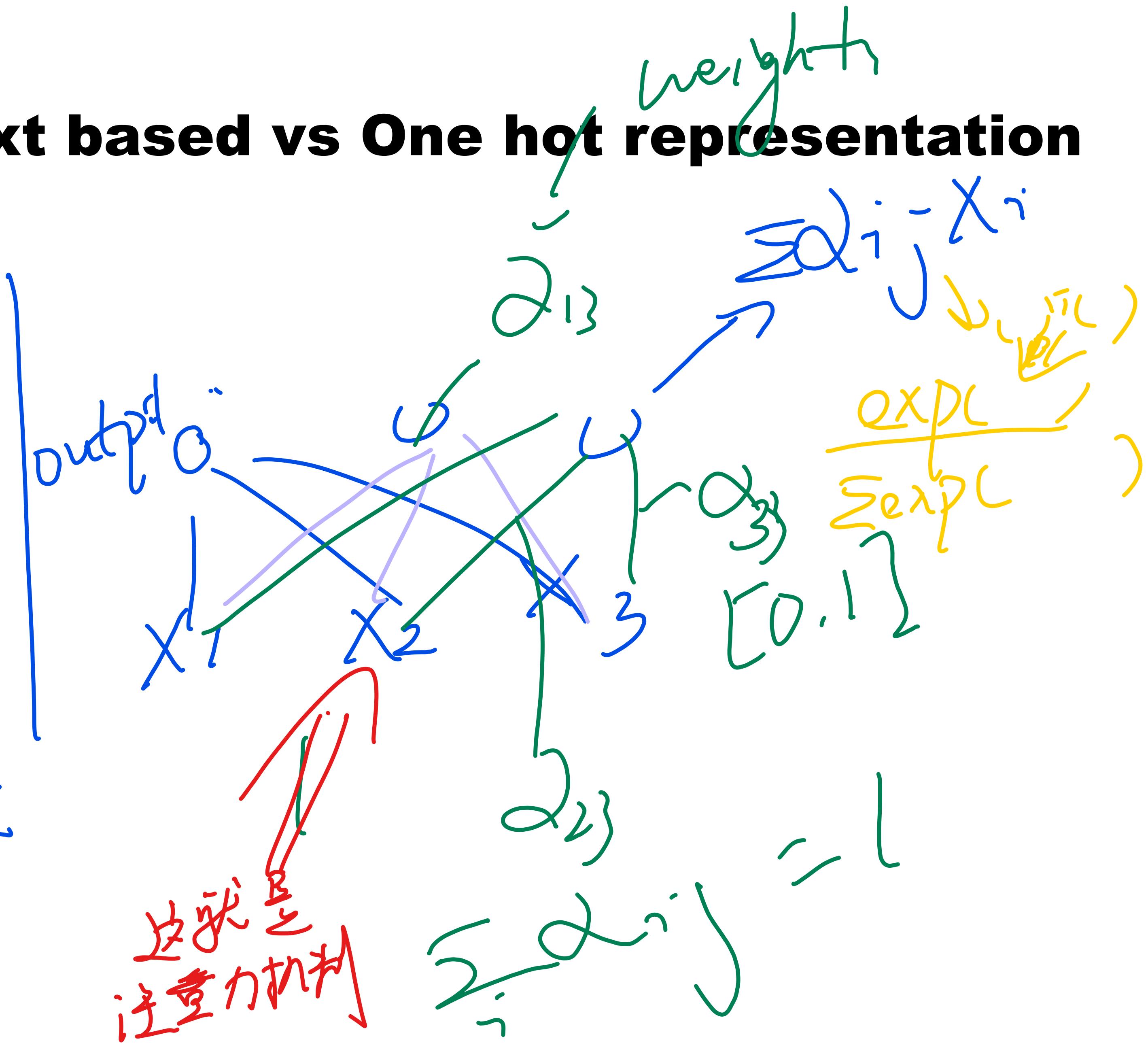
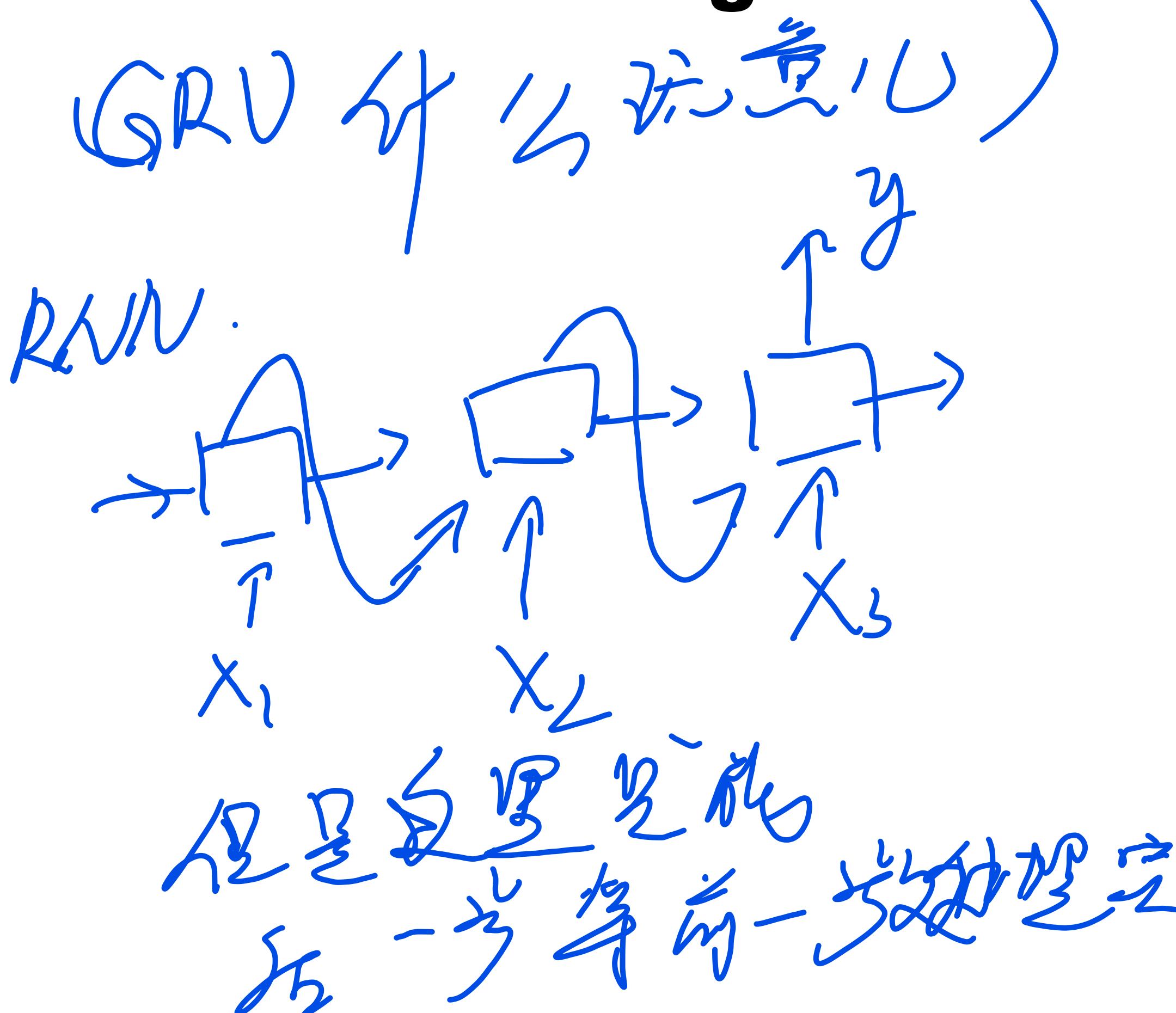
LSTM



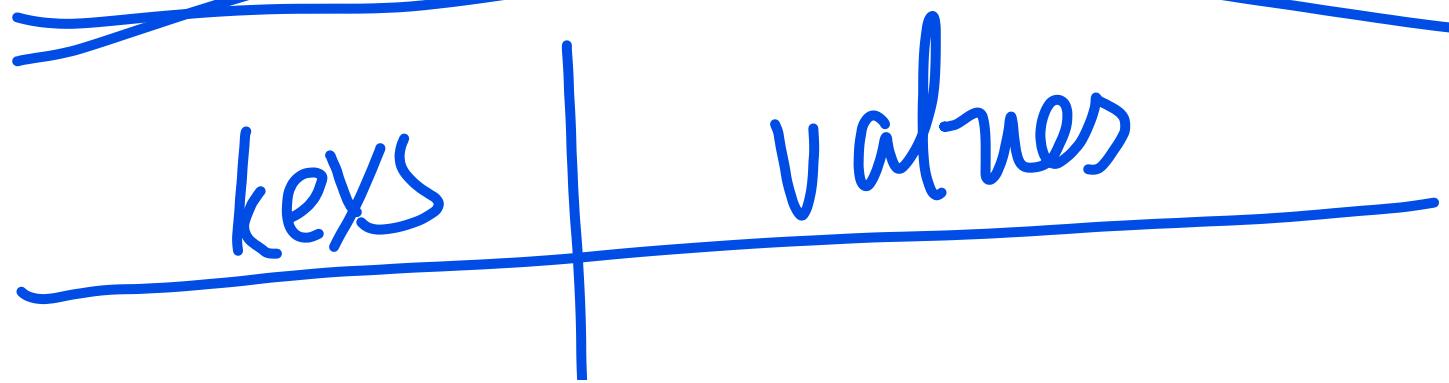
Chaining LSTMs



Word Embedding vs Context based vs One hot representation



① 2用字 ② state 1也可以基于 X₁ Attention Transformer

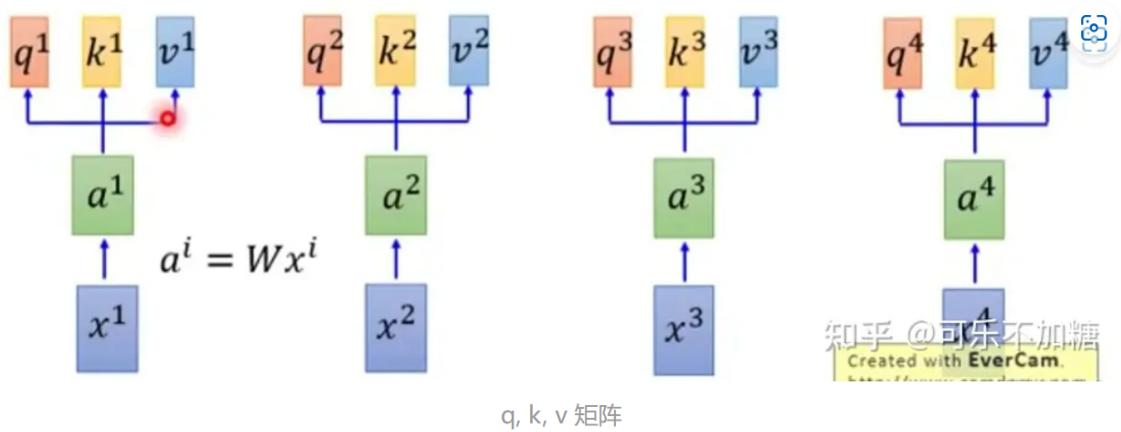


第一步，每个 a^1, a^2, a^3, a^4 都会分别乘以三个矩阵，分别是 q, k, v ；需要注意的是，矩阵 q, k, v 在整个过程中是共享的；公式如下：

$$q^i = W^q a^i$$

$$k^i = W^k a^i$$

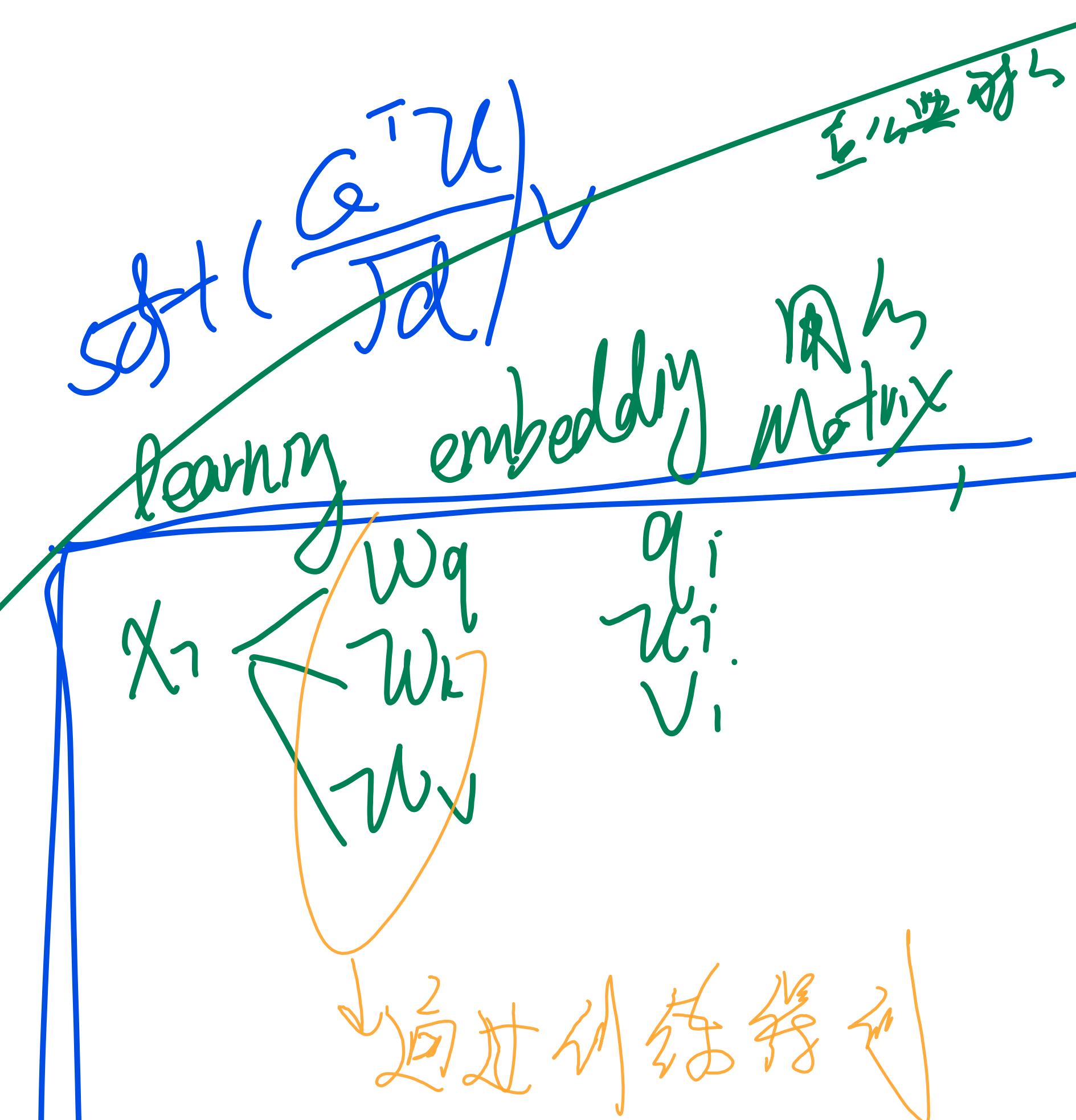
$$v^i = W^v a^i$$



其中， q (Query) 的含义一般的解释是用来和其他单词进行匹配，更准确地说是用来计算当前单词或字与其他的单词或字之间的关联或者关系； k (Key) 的含义则是被用来和 q 进行匹配，也可理解为单词或者字的关键信息。

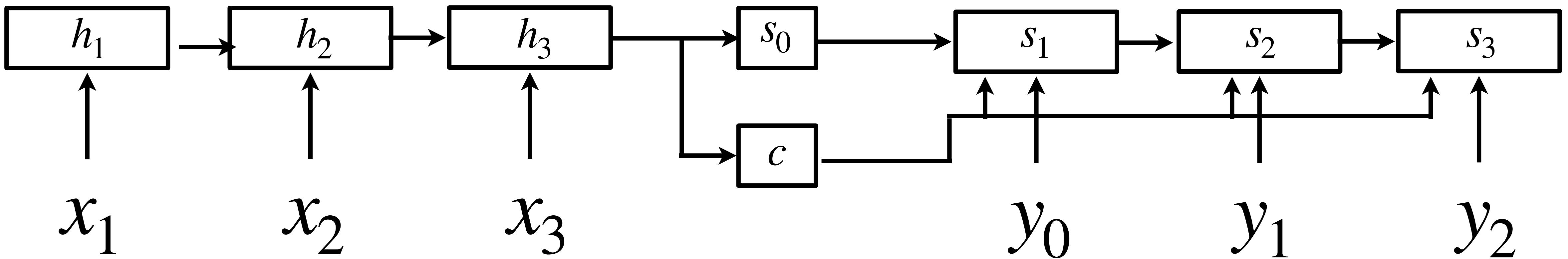
Query

$$q = \sum_i q^i u_i) v_i$$

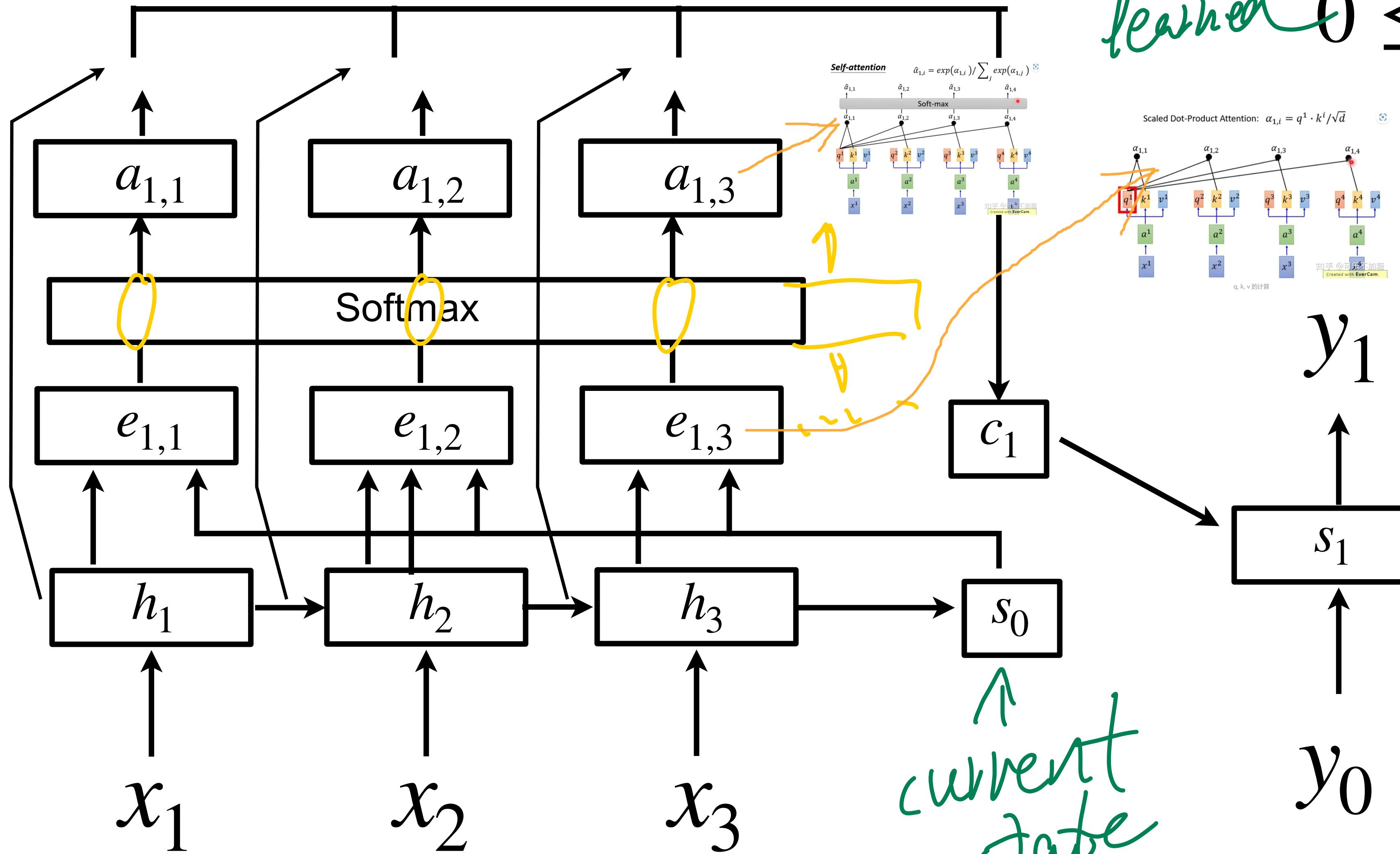


RNN Revisited

$$s_t \leftarrow g_\phi(y_{t-1}, h_{t-1}, c)$$

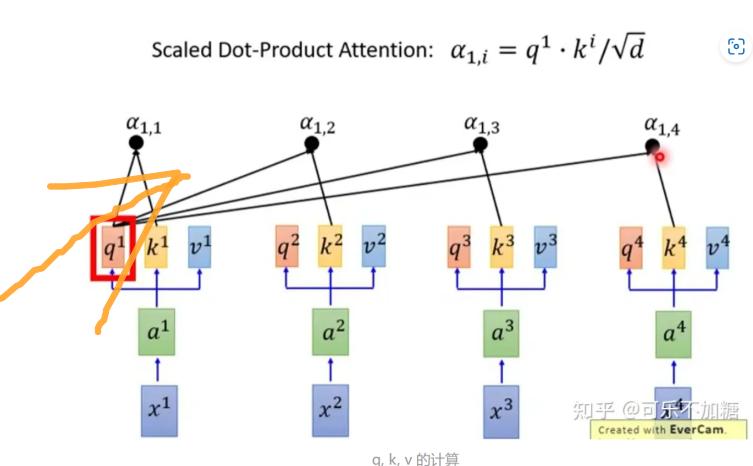


RNN with Attention



$$e_{t,i} \leftarrow g_{att}(s_{t-1}, h_i)$$

learned $0 \leq a_{t,i} \leq 1 \sum_i a_{t,i} = 1$



Handling Speech

Transformer

Unsupervised Learning

There is no direct ground truth for the quantity of interest

Focus on **generative** models:

- Variational Autoencoders (VAEs)
- Normalizing Flows
- Autoregressive Models
- Generative Adversarial Networks (GANs)

Generative Models

- **Assumption: the dataset are samples from an unknown distribution**
- **Goal: create a new sample from $p_{\text{data}}(x)$ that is not in the dataset**



Image credit: *A Style-Based Generator Architecture for Generative Adversarial Networks*, Karras et al.

Generative Models

- **Assumption: the dataset are samples from an unknown distribution**
- **Goal: create a new sample from $p_{\text{data}}(x)$ that is not in the dataset**



Image credit: *A Style-Based Generator Architecture for Generative Adversarial Networks*, Karras et al.

Generative Models

...



$$p_{\text{data}}(x) \approx p_{\theta}(x)$$

Generative model
with parameters θ



...

How do we measure the similarity of $p_{\theta}(\mathbf{x})$ and $p_{\text{data}}(\mathbf{x})$?

Which model?

Image credit: *A Style-Based Generator Architecture for Generative Adversarial Networks*, Karras et al.

Generative Models

How do we measure the similarity of $p_\theta(x)$ and $p_{\text{data}}(x)$?

1) Likelihood of data samples in $p_\theta(x)$

2) Adversarial game

Variational Autoencoders (VAEs)

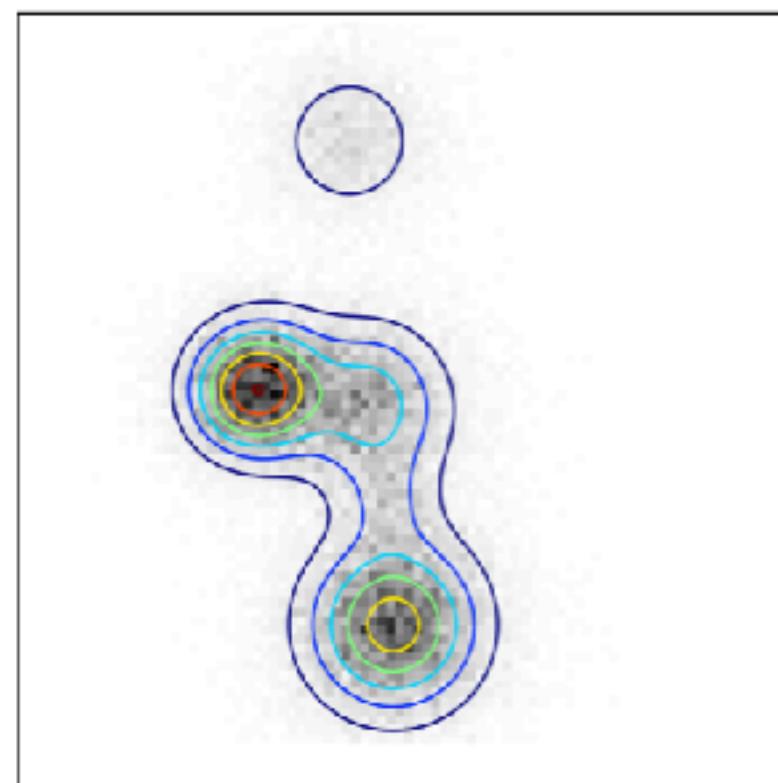
Normalizing Flows

Autoregressive Models

Generative Adversarial Networks (GANs)

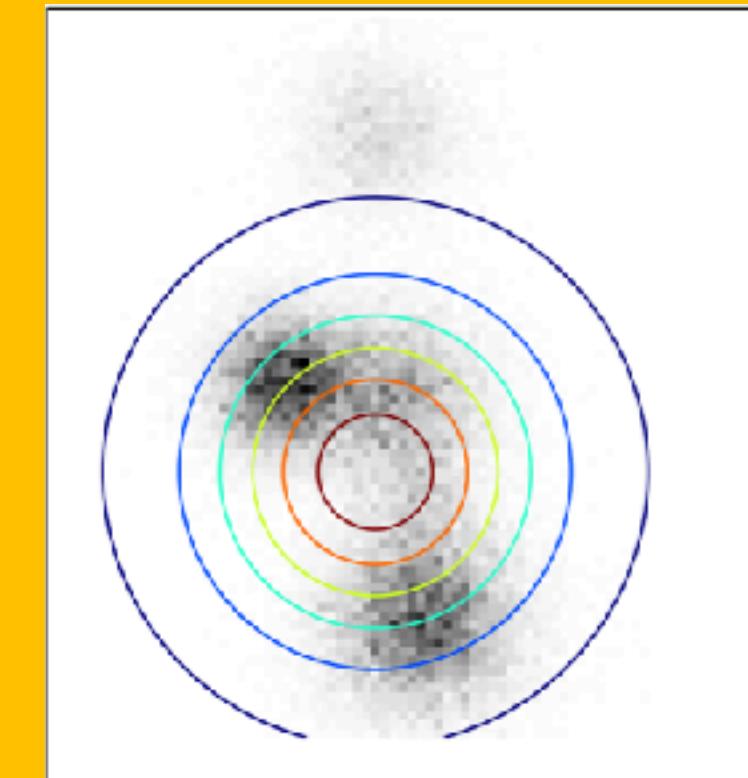
Generative Models

How do we measure the similarity of $p_\theta(x)$ and $p_{\text{data}}(x)$?



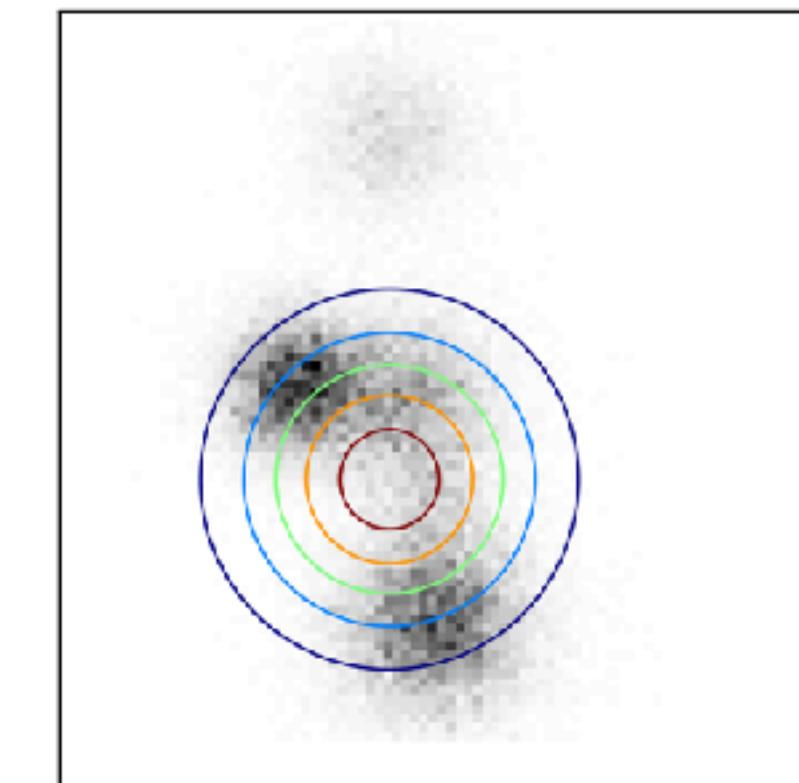
$$p_{\text{data}}(x)$$

1) Likelihood of data samples in $p_\theta(x)$



$$\approx KL(p_{\text{data}} \parallel p_\theta)$$

2) Adversarial game



$$\approx JS(p_{\text{data}} \parallel p_\theta)$$

Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

Likelihood-Based Models: Two Goals

1) Sample from the generative model



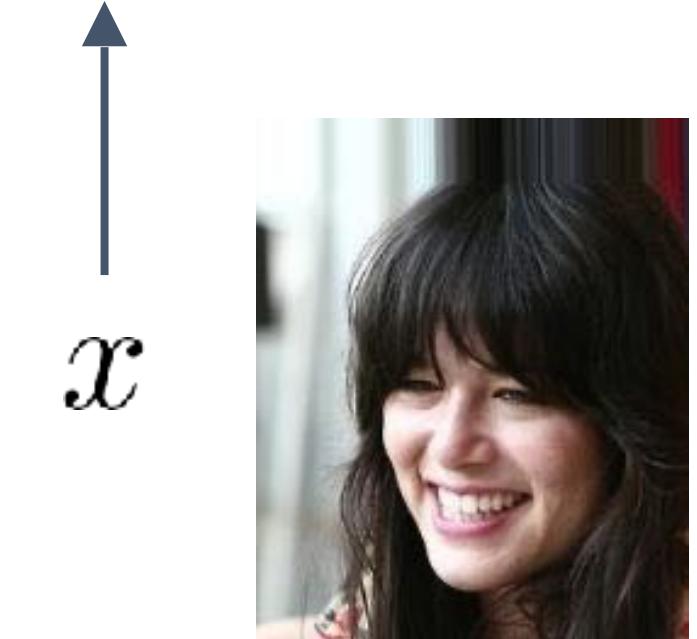
$$x \sim p_{\theta}$$

Generative model
with parameters θ

2) Evaluate the likelihood of a given sample in the model

$$p_{\theta}(x)$$

Generative model
with parameters θ



$$x$$

Feature Space

Data distribution in 2D feature space (colors are class labels)

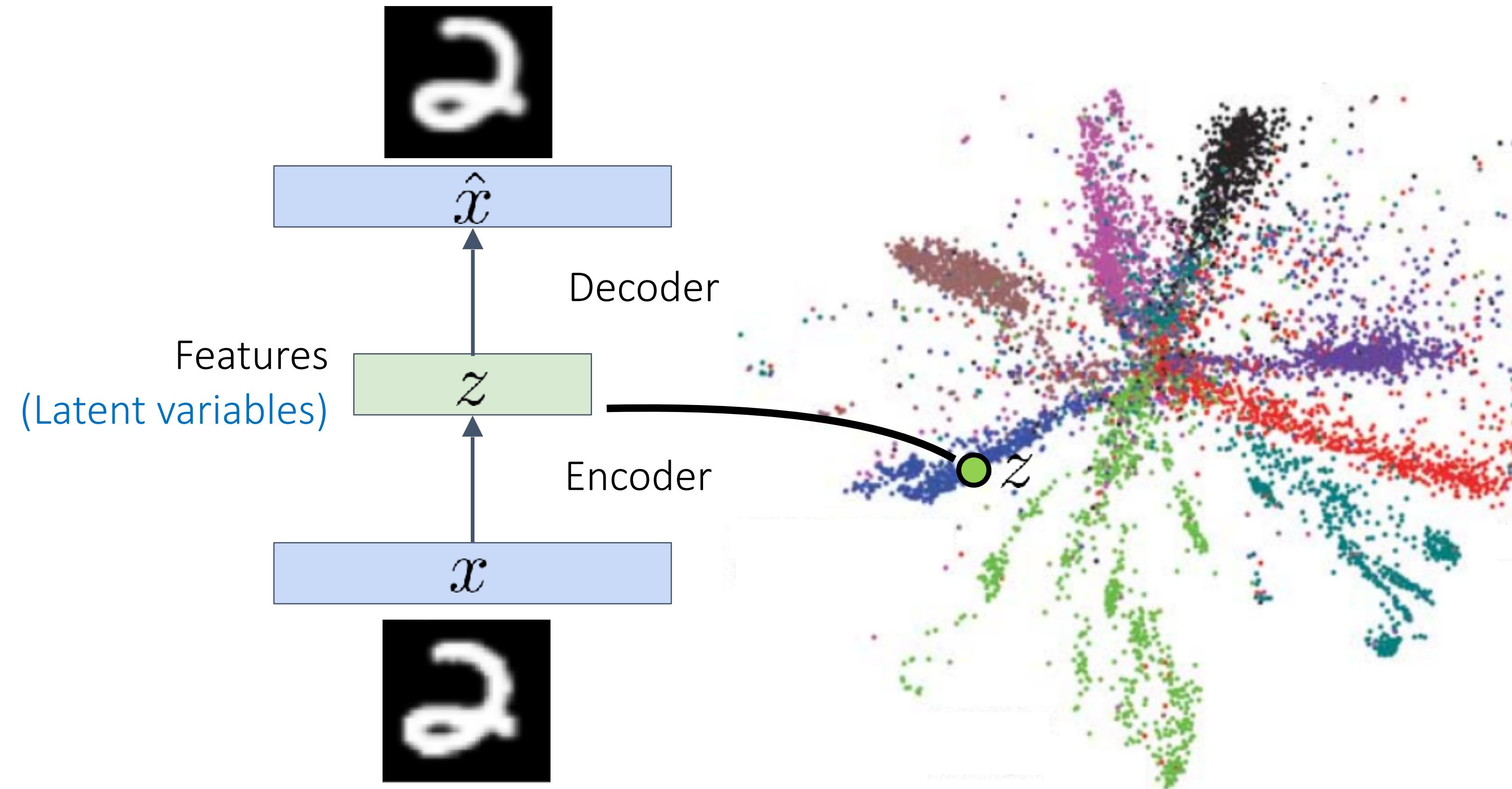


Image Credit: Reducing the Dimensionality of Data with Neural Networks, Hinton and Salakhutdinov

Autoencoders as Generative Models?

- Is a trained decoder a generative model?
- Can we generate a new sample?

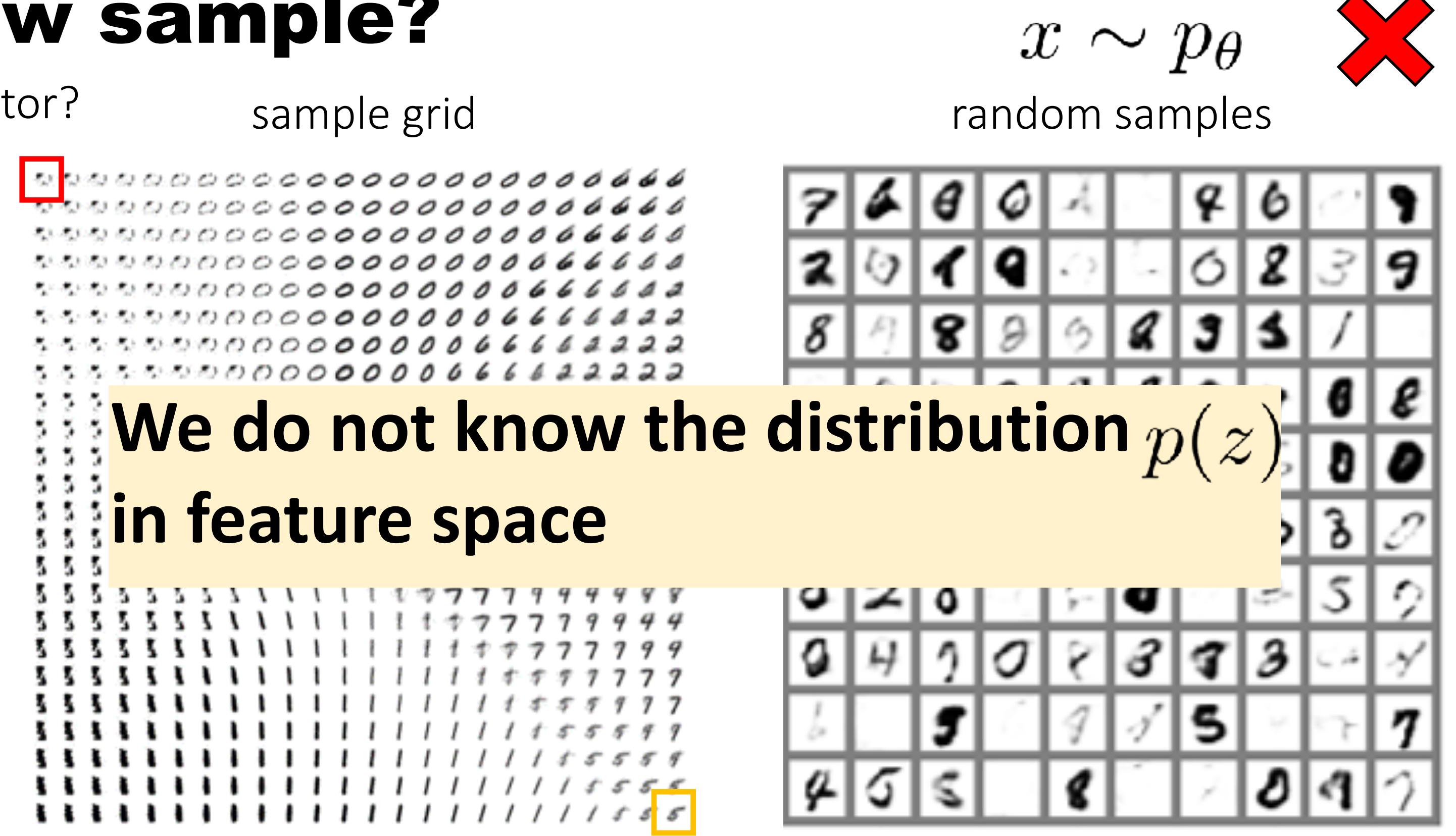
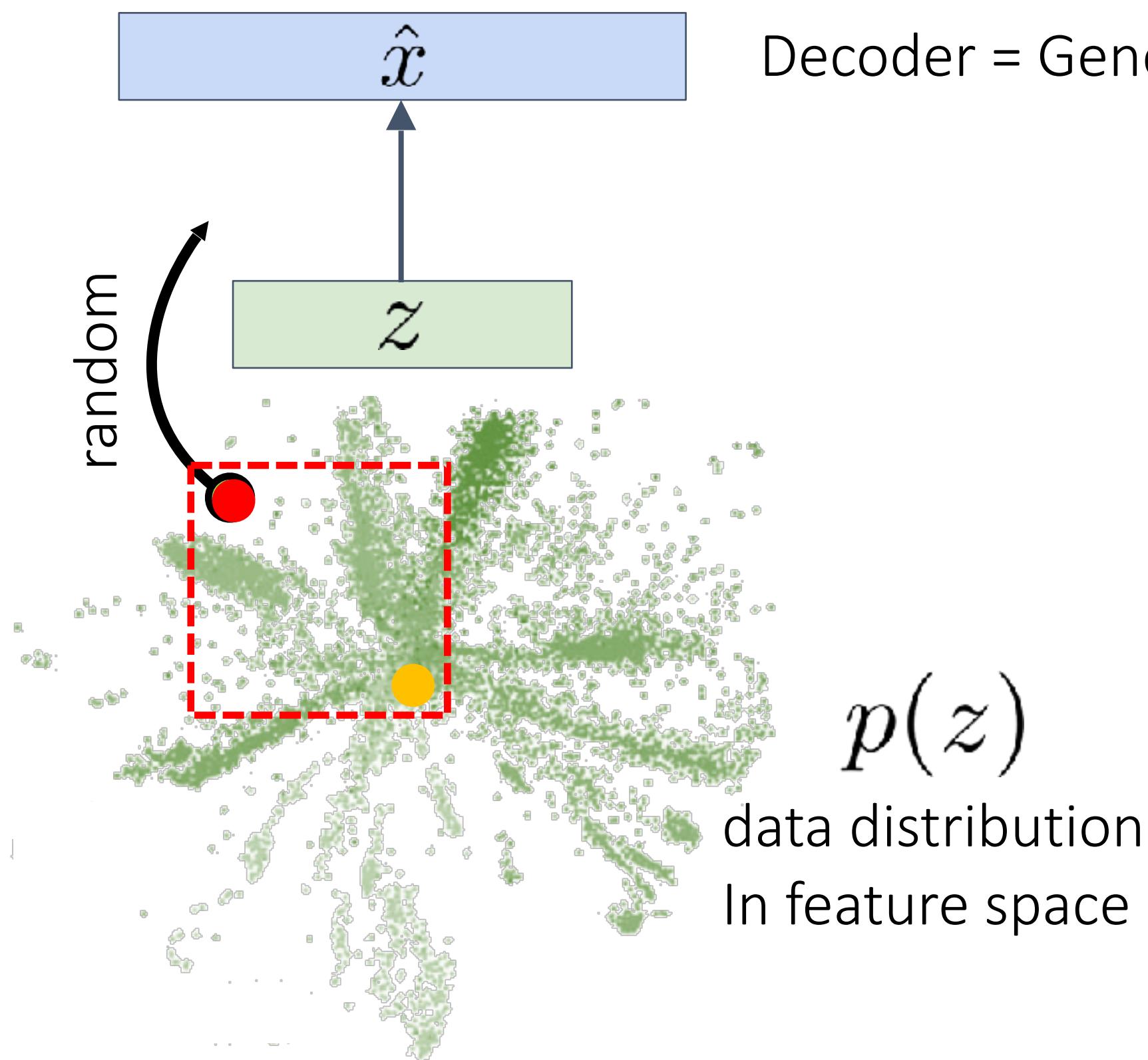
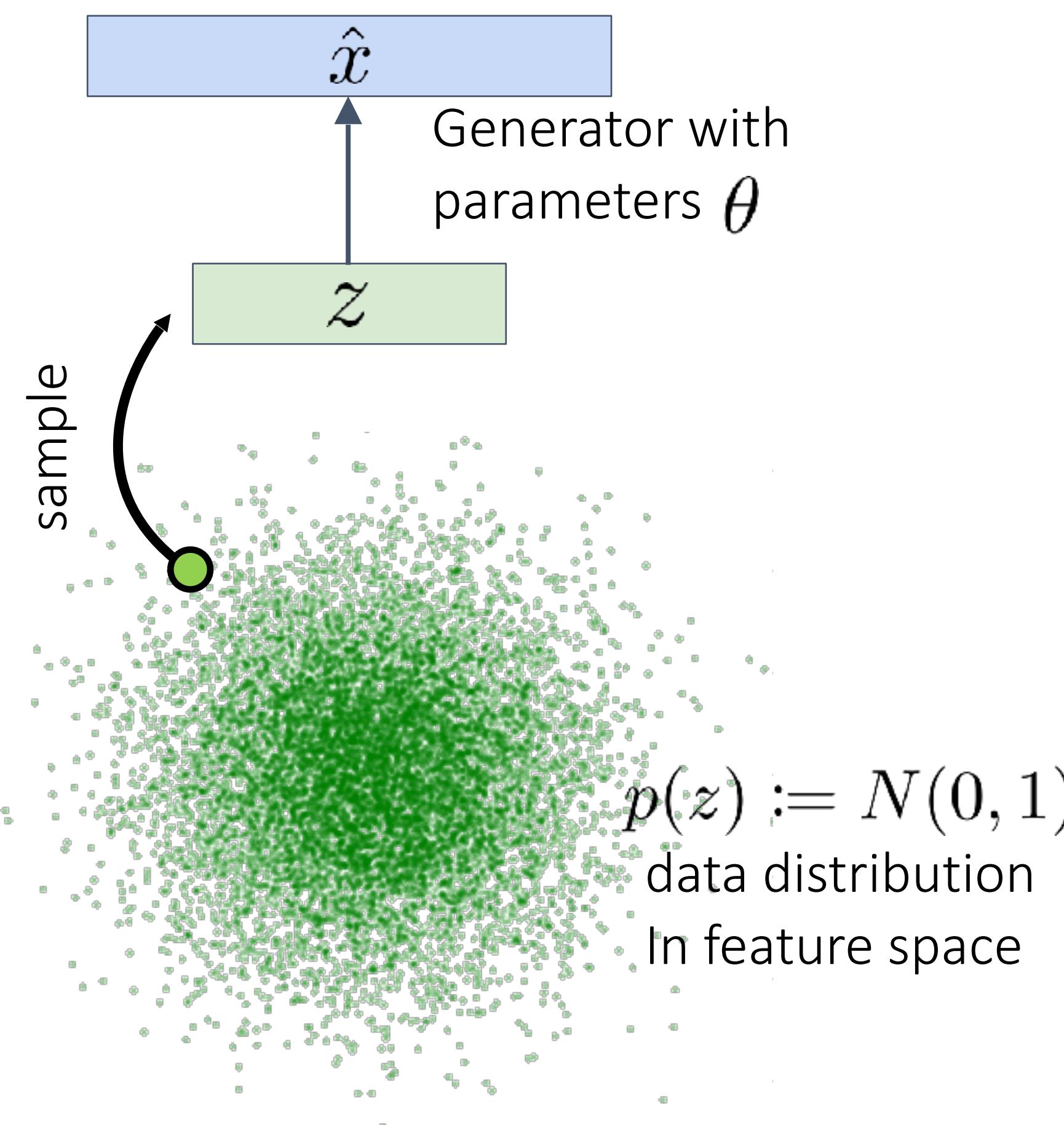


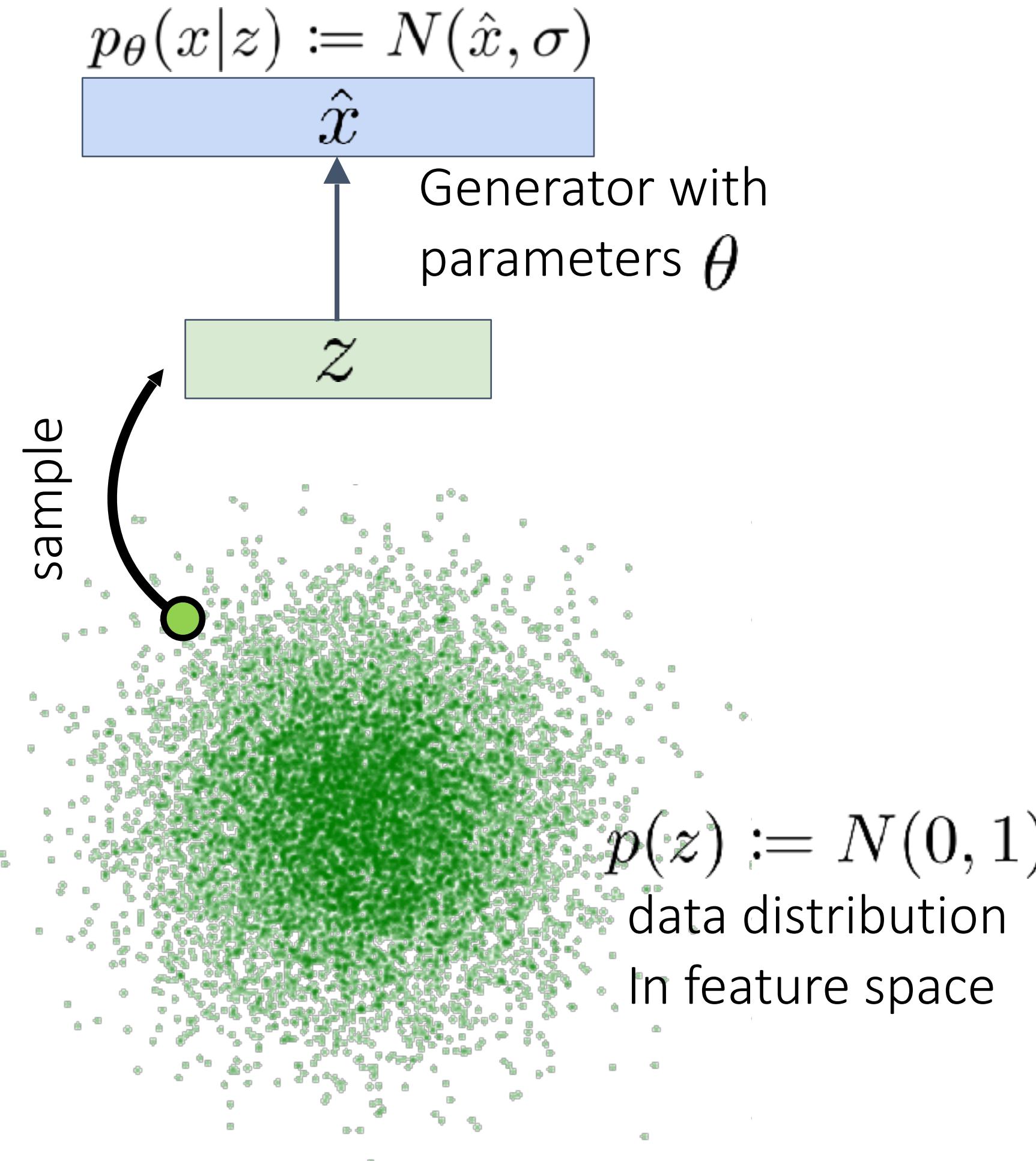
Image Credit: Reducing the Dimensionality of Data with Neural Networks, Hinton and Salakhutdinov

Latent Variable Model



- Define $p(z)$ as a known distribution

Latent Variable Model

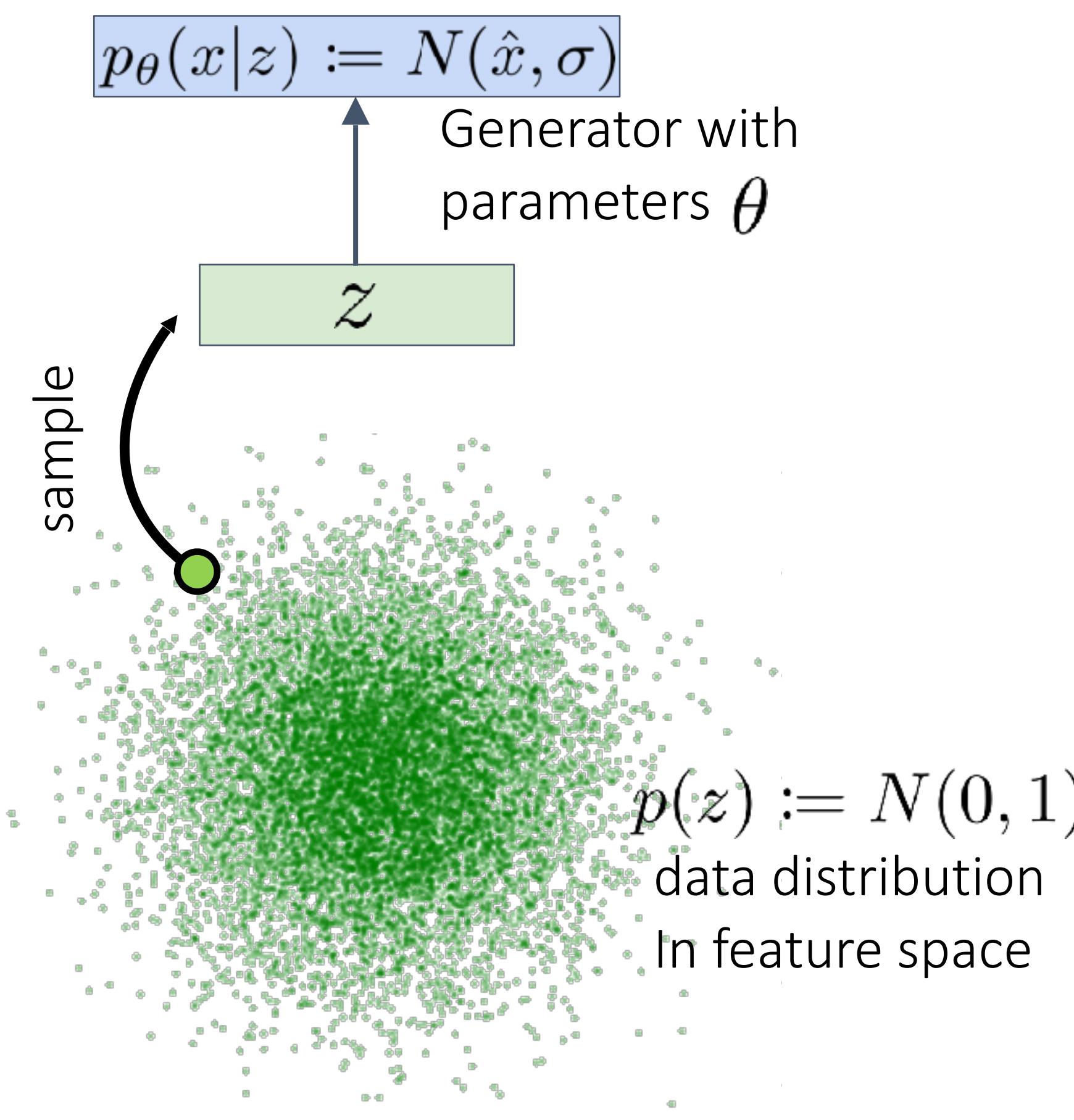


- Train generator with NLL of data as loss

$$-\log \sum_{x_i \in \text{data}} p_\theta(x_i)$$
- Can we compute the likelihood $p_\theta(x)$?

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

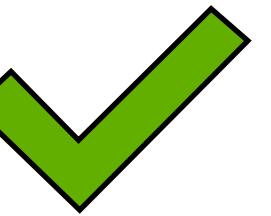
Latent Variable Model



- Train generator with NLL of data as loss

$$-\log \sum_{x_i \in \text{data}} p_\theta(x_i)$$
- Can we compute the likelihood $p_\theta(x)$?

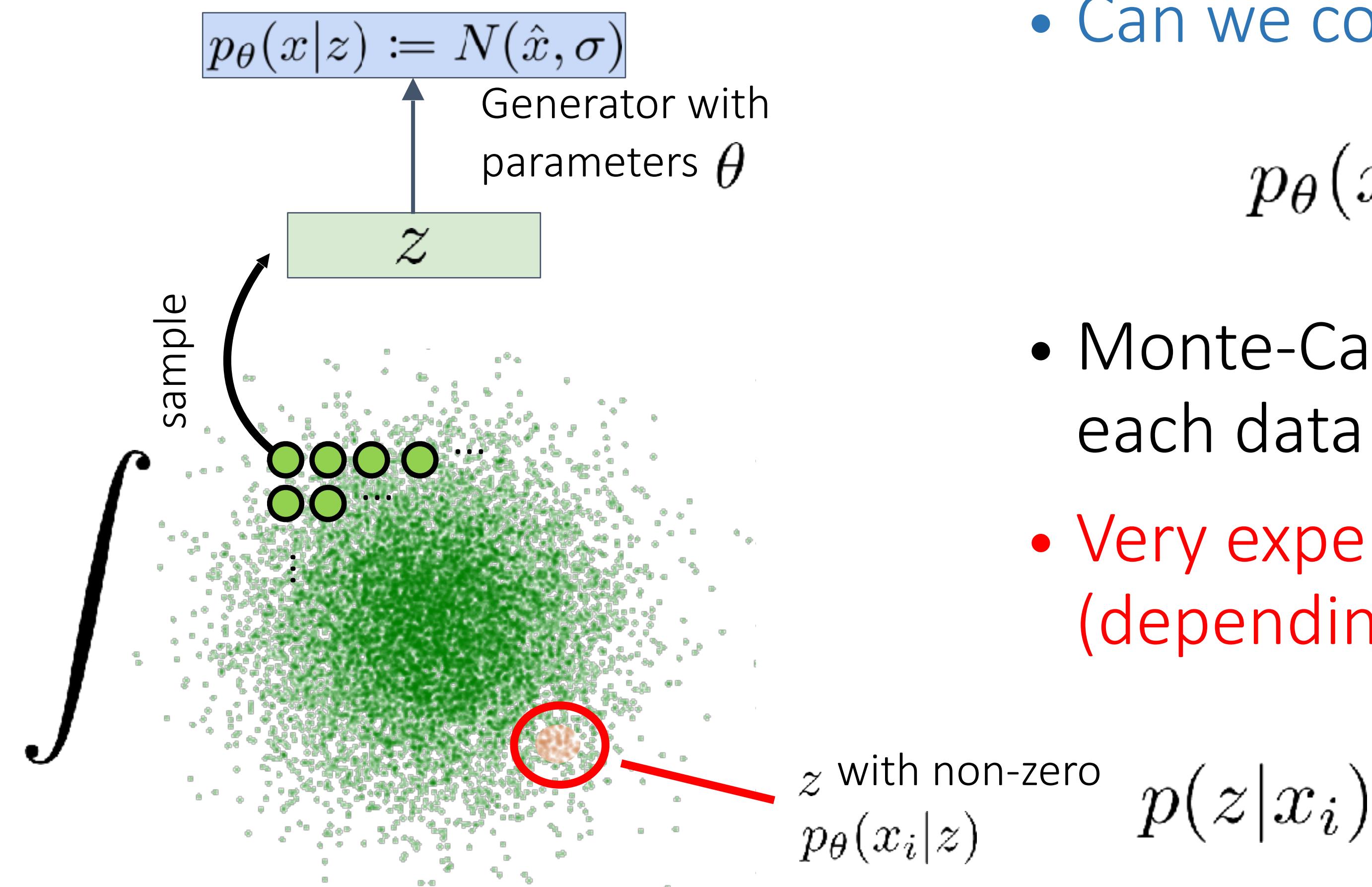
$$p_\theta(x) = \int p_\theta(x|z) p(z) dz$$



Variational Autoencoders (VAEs): The Encoder

- During training, another network can learn to $p(z | x_i)$ approximate the distribution
- $q_\phi(z | x_i)$ should be much smaller than $p(z)$
- Makes the computing the integral tractable

Latent Variable Model: Monte-Carlo



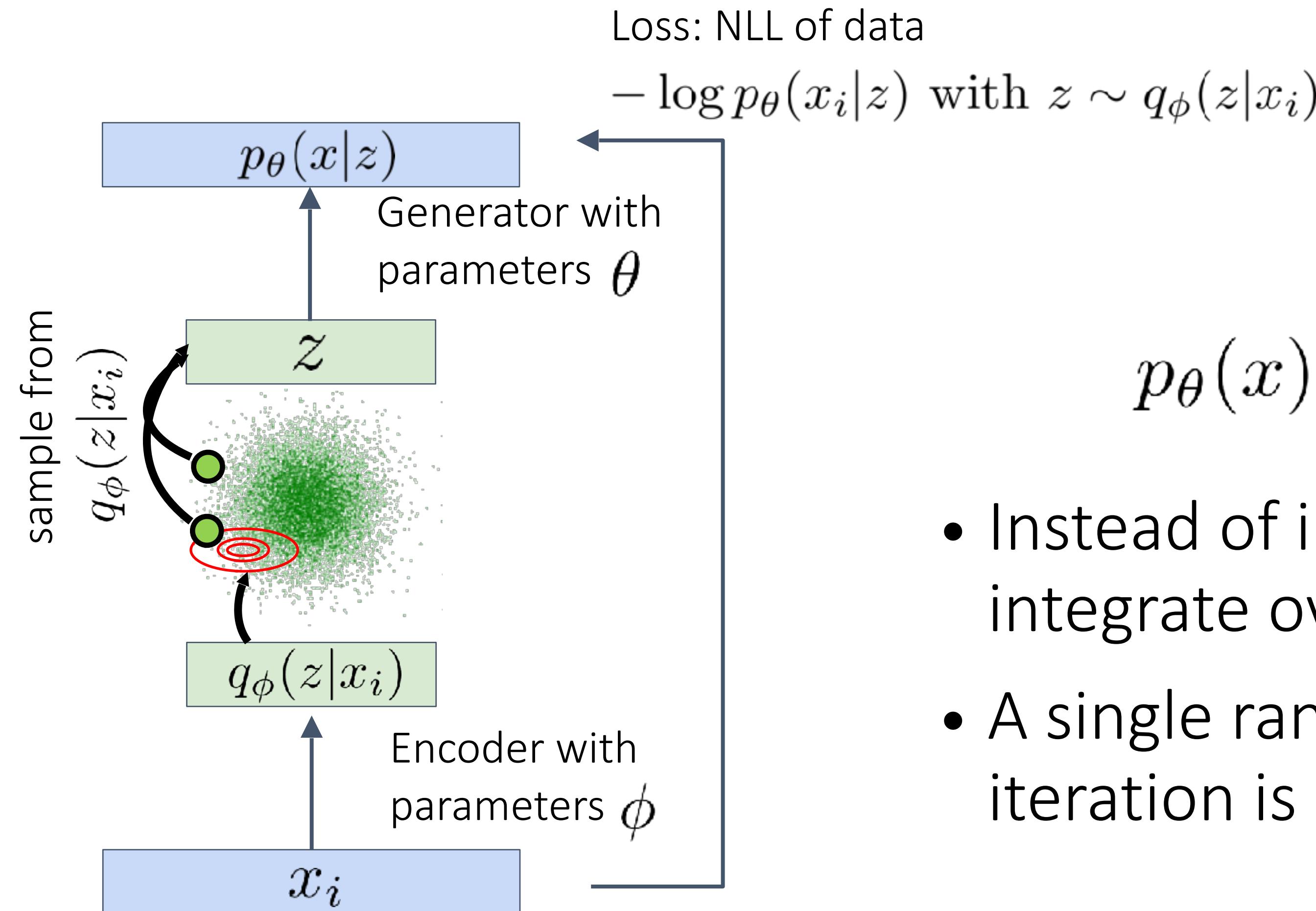
(Negative log-likelihood)

- Can we compute the likelihood $p_\theta(x)$?

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

- Monte-Carlo integration to solve integral for each data sample
- Very expensive, or very inaccurate (depending on sample count)

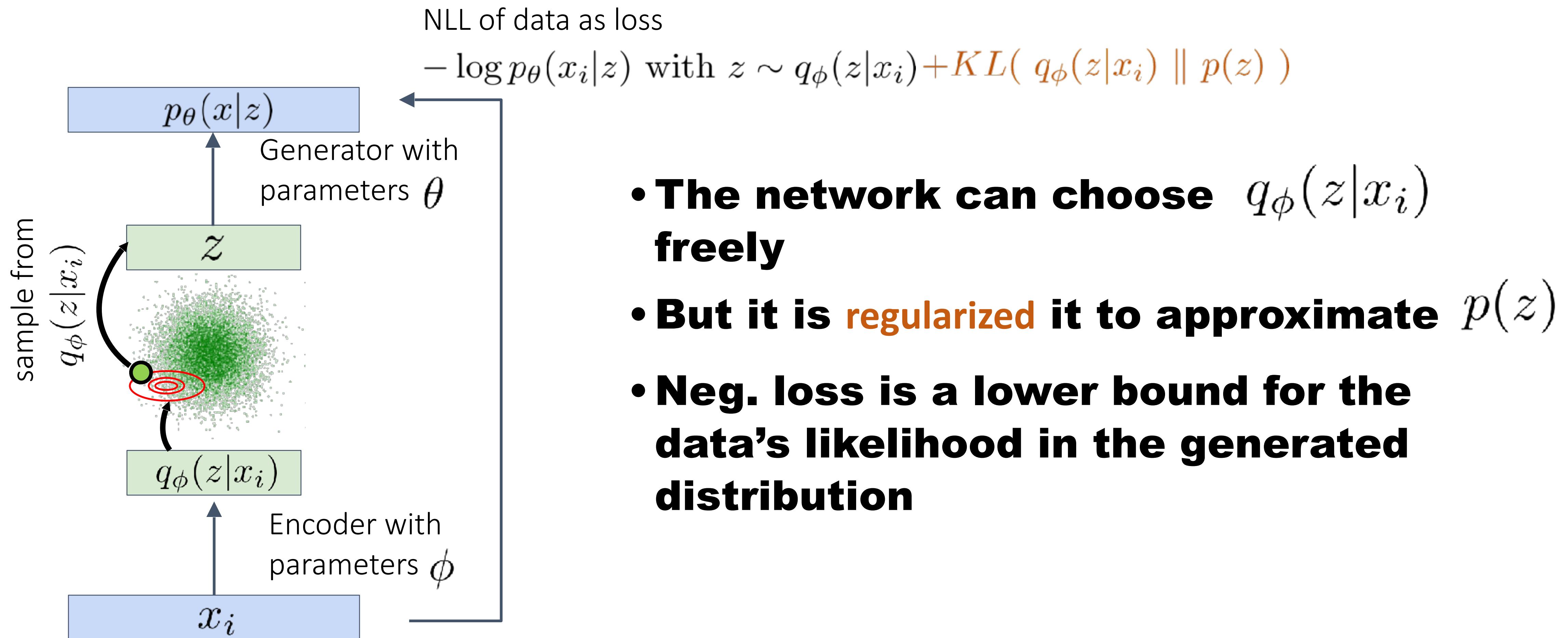
Variational Autoencoders (VAEs): Encoder



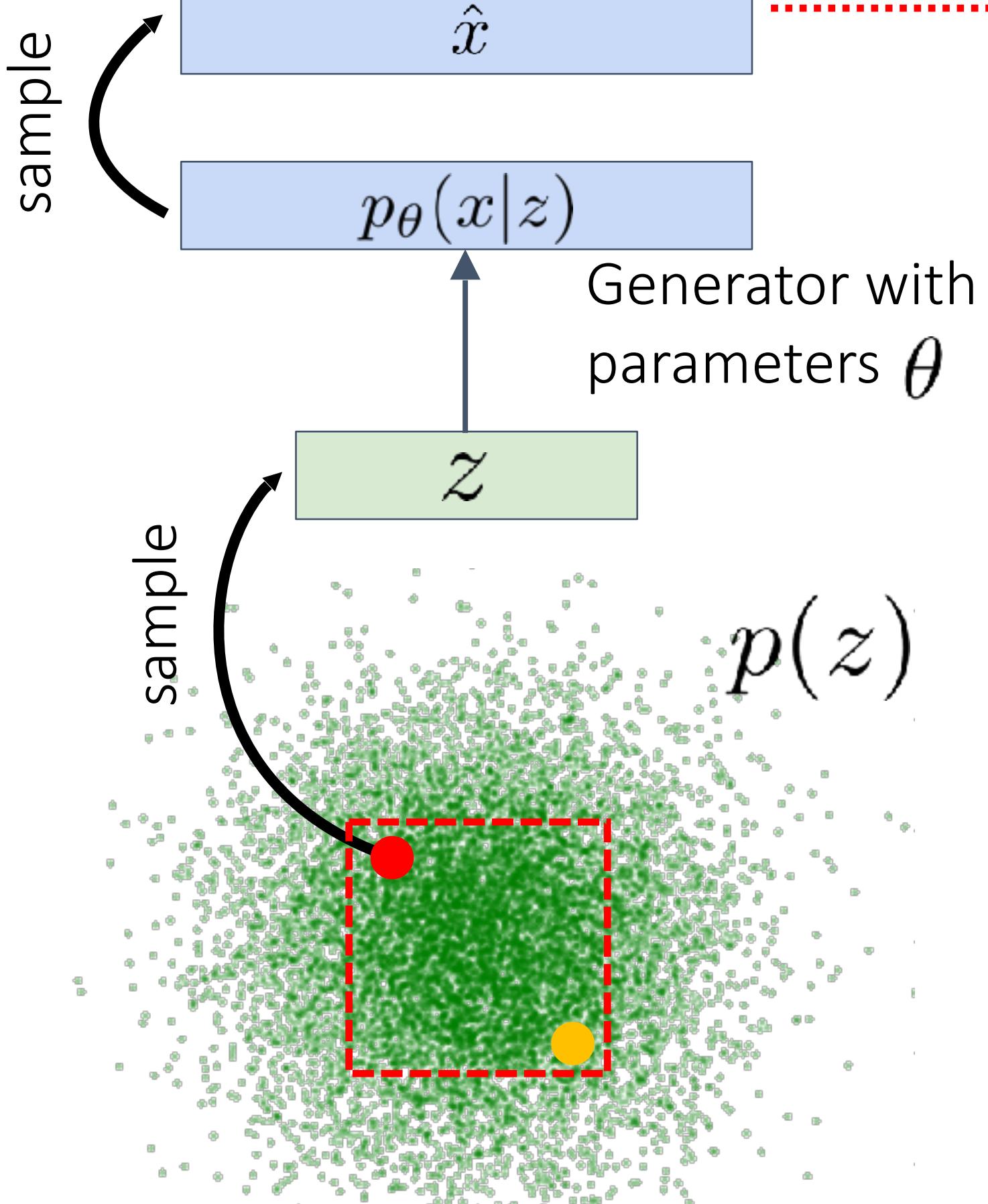
$$p_\theta(x) = \int p_\theta(x|z) p(z) dz$$

- Instead of integrating over all $p(z)$, integrate over $q_\phi(z|x_i)$ only
- A single random sample from $q_\phi(z|x_i)$ per iteration is usually enough

Variational Autoencoders (VAEs): Regularizer



Generating Data



The diagram shows two types of reconstructions: VAE and Autoencoder. The VAE reconstruction is a 10x10 grid of digits, where each digit is a 3x3 matrix. The first digit (top-left) is highlighted with a red box. The last digit (bottom-right) is highlighted with a yellow box. The Autoencoder reconstruction is a 10x10 grid of binary values (0s and 1s).

VAE Reconstruction (10x10 grid of 3x3 digits):

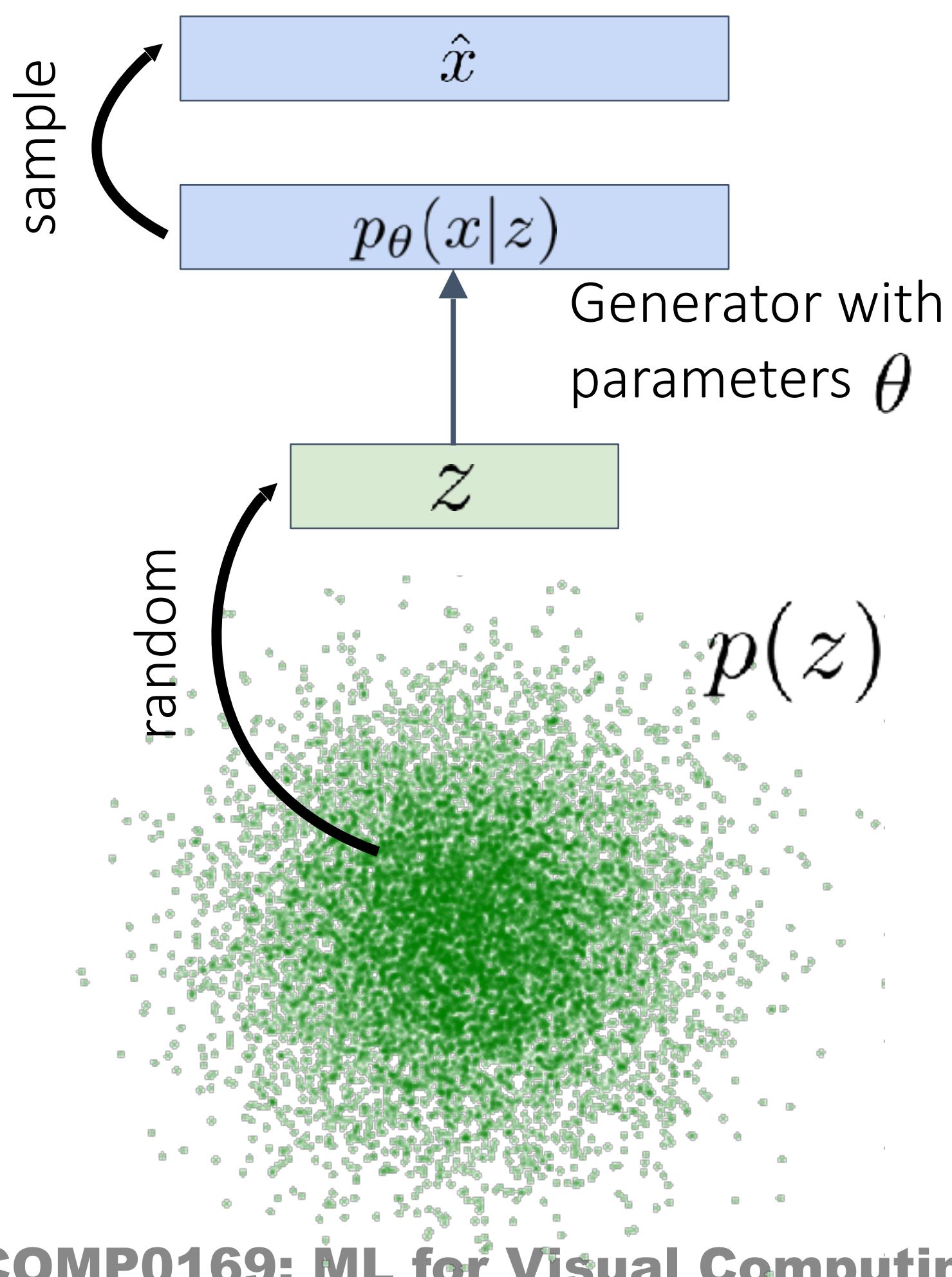
6	6	6	6	0	0	0	0	0	0
9	4	4	8	2	2	2	0	0	0
9	2	2	2	2	2	3	8	5	6
9	9	2	2	2	2	2	3	3	5
9	9	2	2	2	2	2	3	3	5
9	9	9	4	2	2	2	3	3	3
9	9	9	9	9	2	2	3	3	3
9	9	9	9	9	9	3	3	3	3
9	9	9	9	9	9	8	3	3	3
7	9	9	9	9	9	8	8	3	3

Autoencoder Reconstruction (10x10 grid of binary values):

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

Image Credit: Auto-Encoding Variational Bayes, Kingma and Welling

Generating Data



VAE

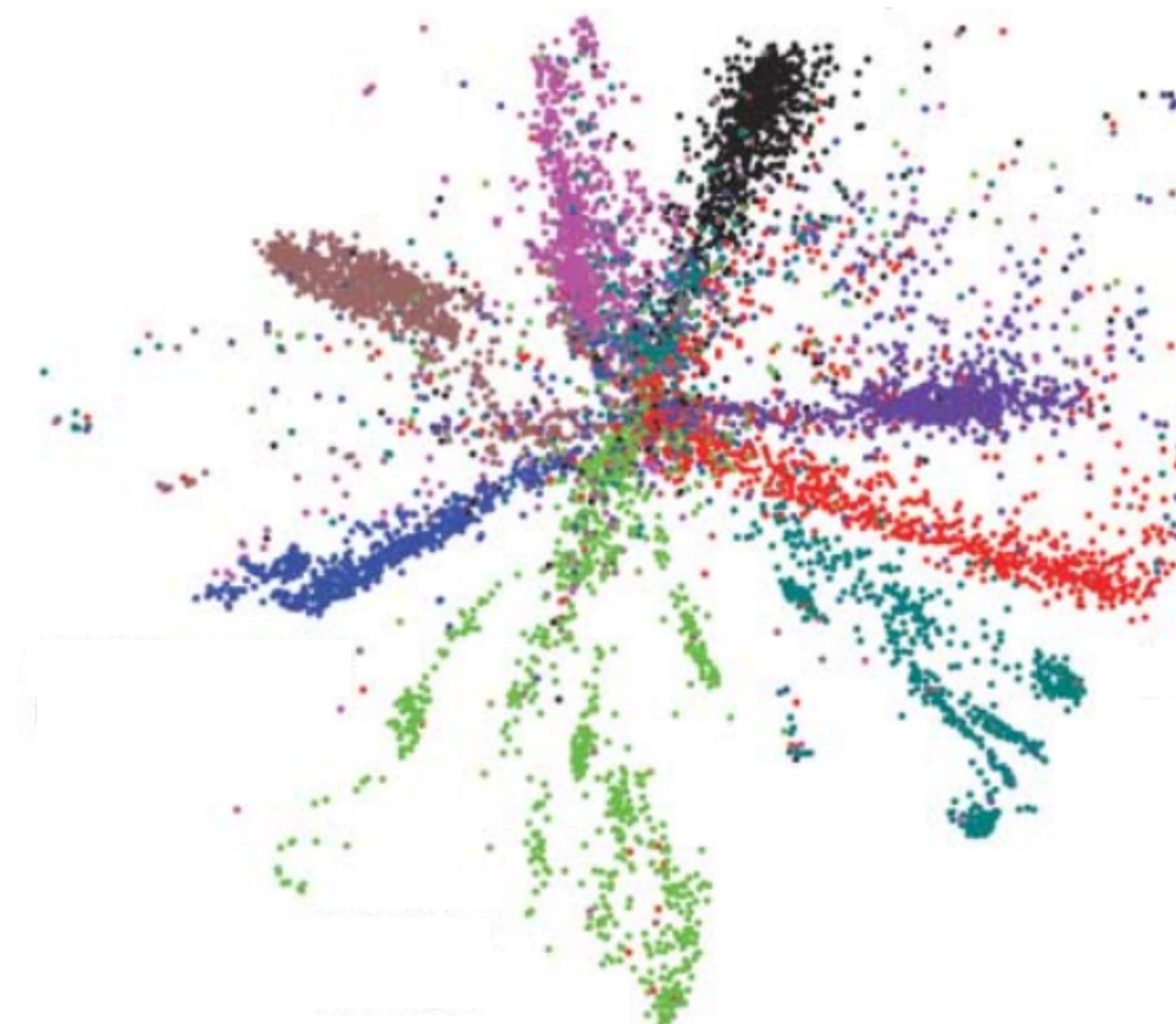
6 6 3 3 3 3 9 2 2 1
3 3 6 9 1 8 9 1 9 9
2 4 0 7 6 9 9 1 9 5
9 4 1 3 9 6 9 6 3 9
5 6 7 6 0 6 9 3 0 1
4 7 2 5 9 6 6 2 8 7
6 1 2 6 7 1 9 2 7 9
7 1 9 0 5 3 9 1 6 1
5 8 7 1 7 5 3 6 5
0 3 8 5 2 3 6 3 3 9

Autoencoder



Feature Space of Autoencoders vs. VAEs

Autoencoder



VAE

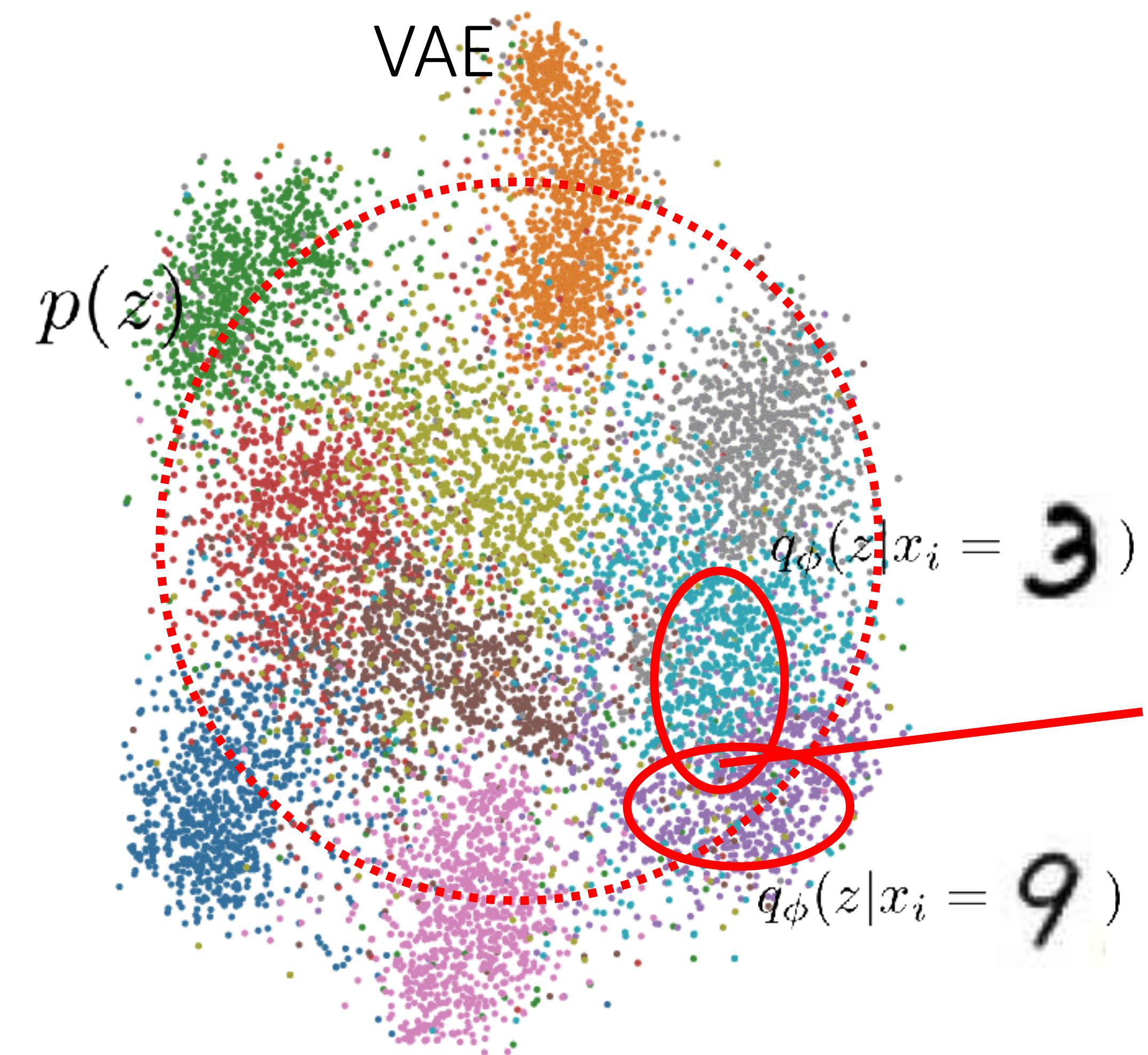


Image Credit: Reducing the Dimensionality of Data with Neural Networks, Hinton and Salakhutdinov

Summary: Variational Autoencoders (VAEs)

Positives

- **Creates a feature space**
- **Relatively stable to train**

Negatives

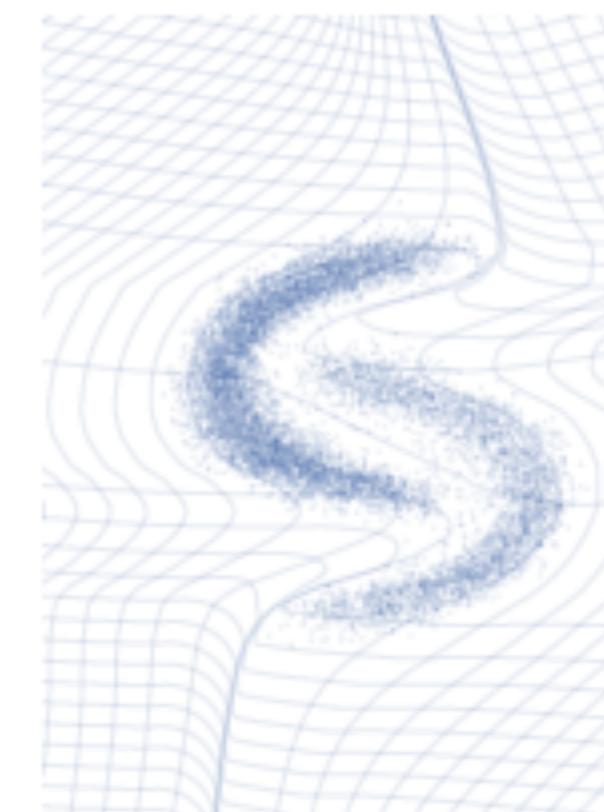
- **Likelihood evaluation can only be approximated**
- **Projection of sample into feature space can only be approximated**
- **Regularization makes the results a bit blurry**

Normalizing Flows

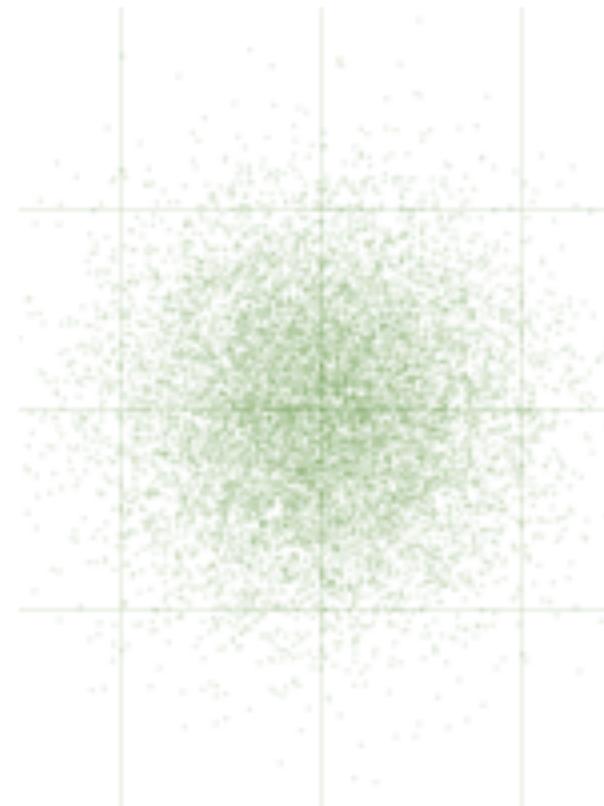
$$x \sim p_\theta(x)$$

f_θ Generator with parameters θ

$$z \sim p_z(z)$$



$$x = f_\theta(z)$$



$$p_z(z) \text{ (known)}$$

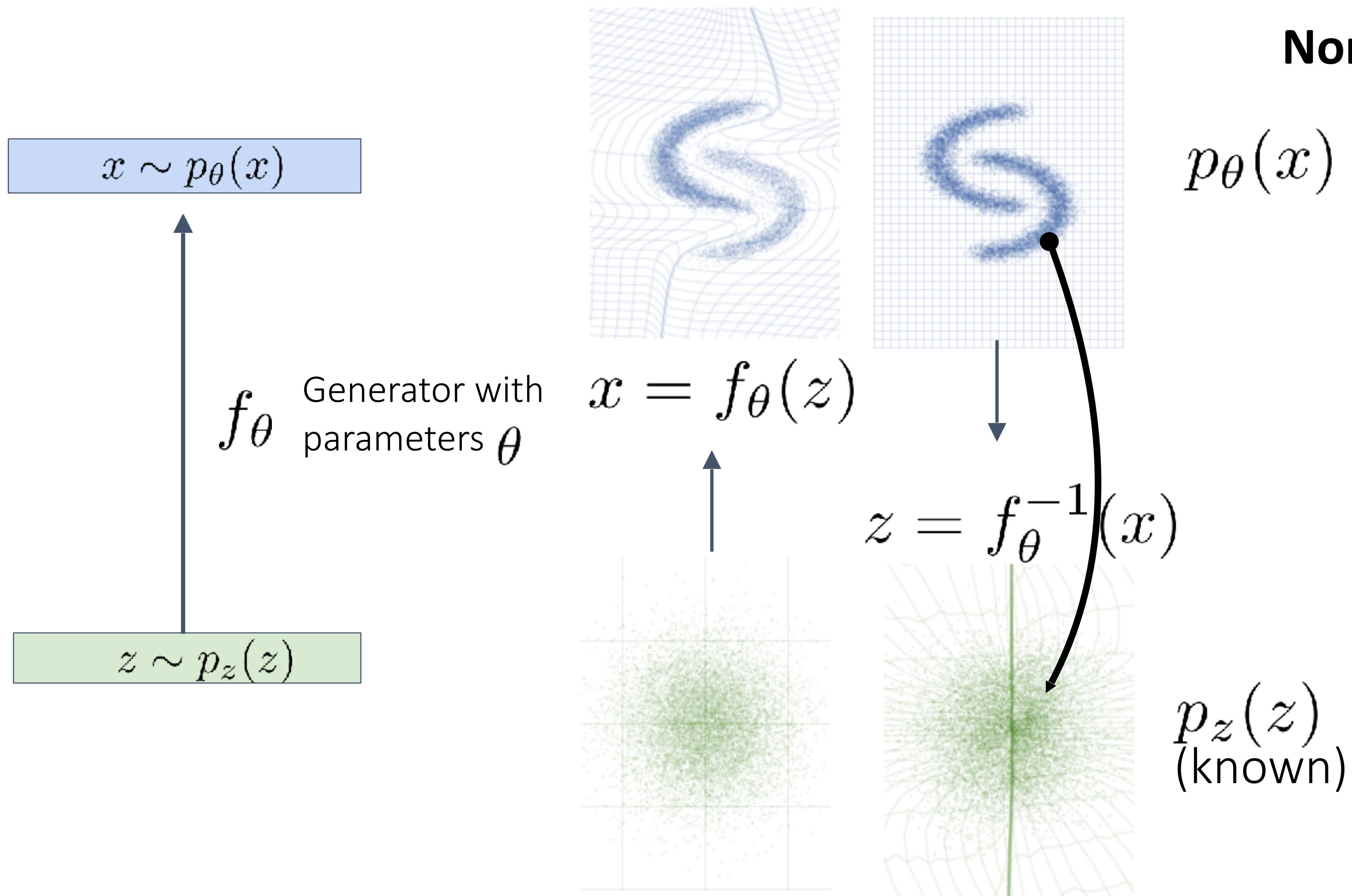
Variational Autoencoders (VAEs):

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

Image Credit: DENSITY ESTIMATION USING REAL NVP:, Dinh et al.

Normalizing Flows

Sample? 
Evaluate? 



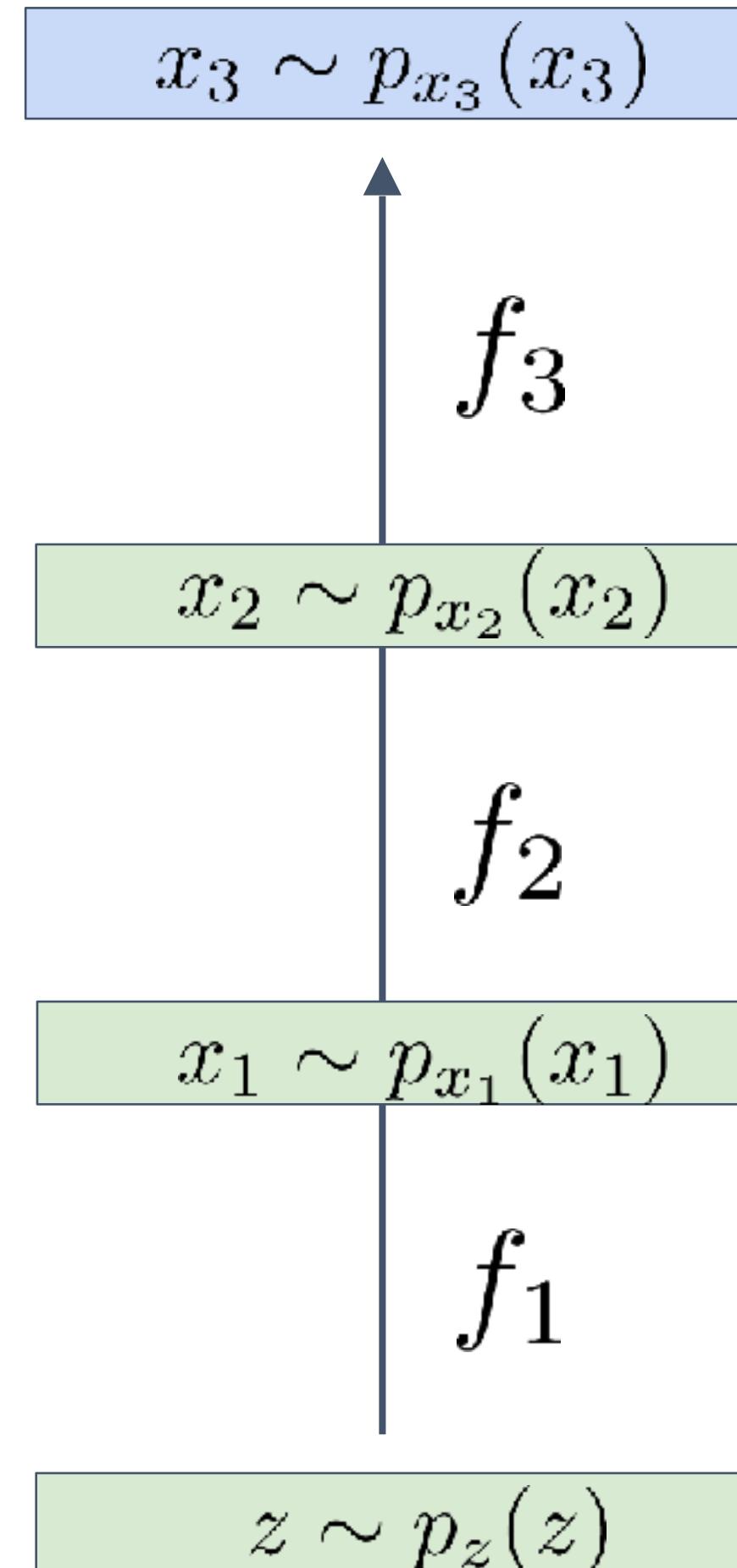
Normalizing Flows:

$$p_\theta(x) = p_z(f_\theta^{-1}(x))$$

... times the local density change caused by f_θ

$$\left| \det \left(\frac{\partial f_\theta^{-1}(x)}{\partial x} \right) \right|$$

Normalizing Flows: Chaining



$$p_3(x_3) = p_z(f_1^{-1} \circ f_2^{-1} \circ f_3^{-1}(x_3))$$

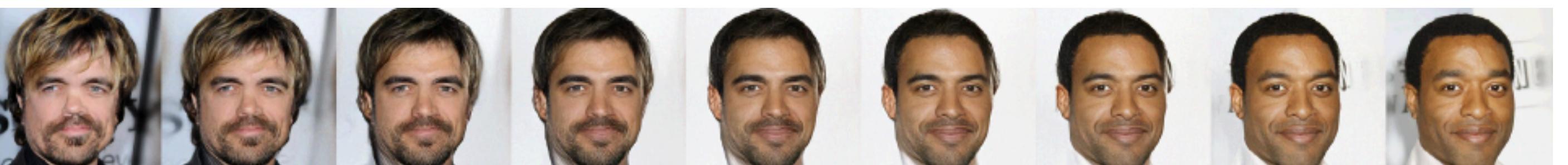
... times the local density change caused
by the chain of transformations

$$\prod_{i=1}^3 \left| \det \left(\frac{\partial f_i^{-1}}{\partial x_i} \right) \right|$$

Example: Glow

Invertible functions:

- **Linear (1x1 conv) layer with weight matrix parameterized by its LU decomposition**
- **Affine coupling layer to propagate information between pixels**



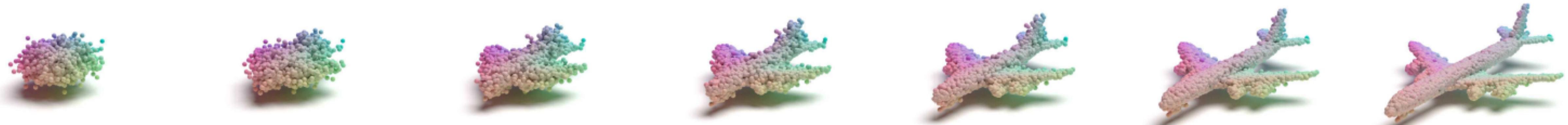
Training: 40 GPUs, 2 weeks

Image Credit: Glow: Generative Flow with Invertible 1x1 Convolutions, Kingma and Dhariwal

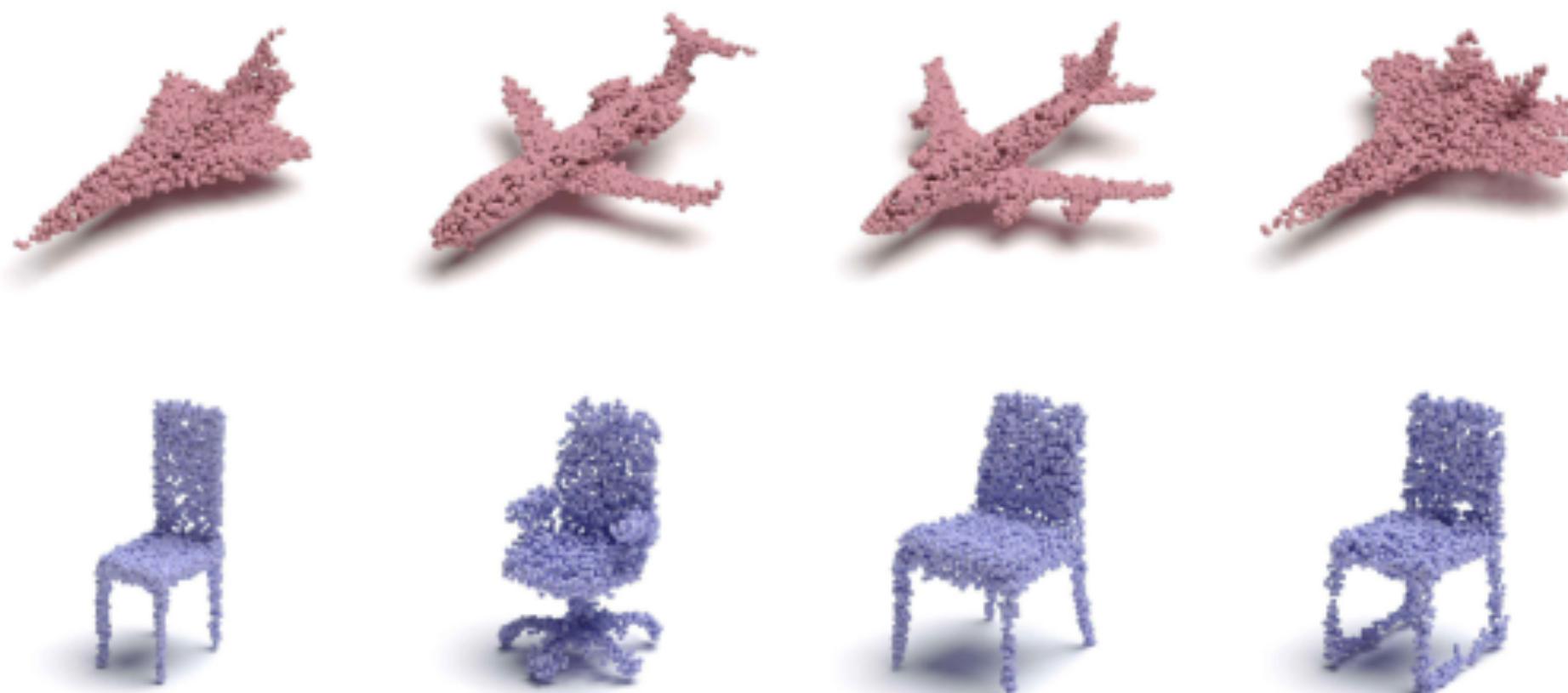
Example: PointFlow

Two flows: One to create the distribution of shape feature vectors, one for the distribution of points on a shape

Shape generation flow



Free shape generation



Shape interpolation

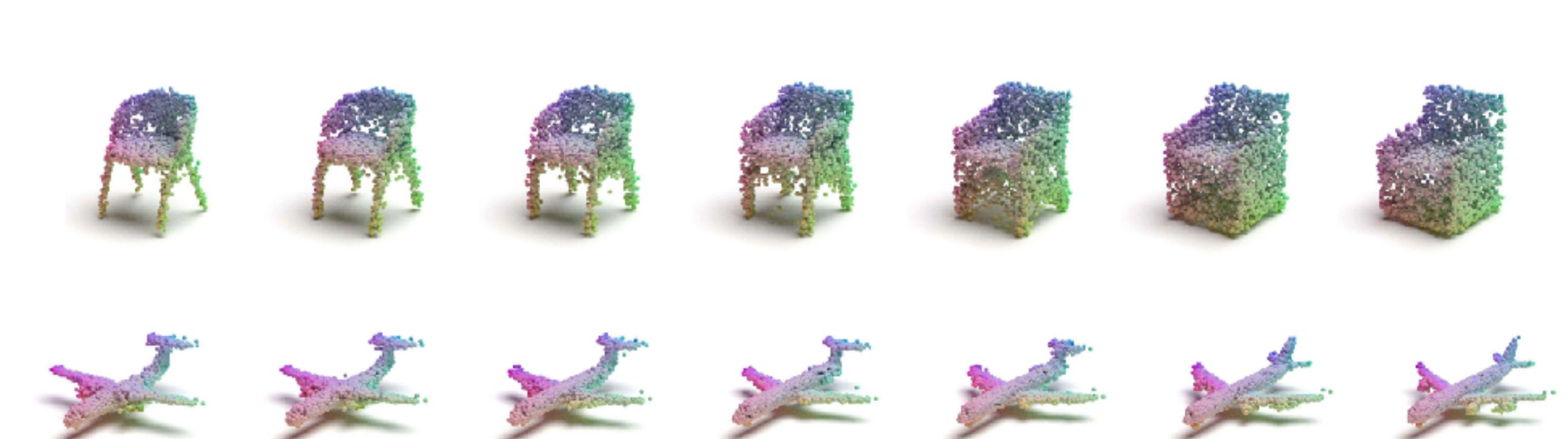


Image Credit: PointFlow: 3D Point Cloud Generation with Continuous Normalizing Flows, Yang et al.

Summary: Normalizing Flows

Positives

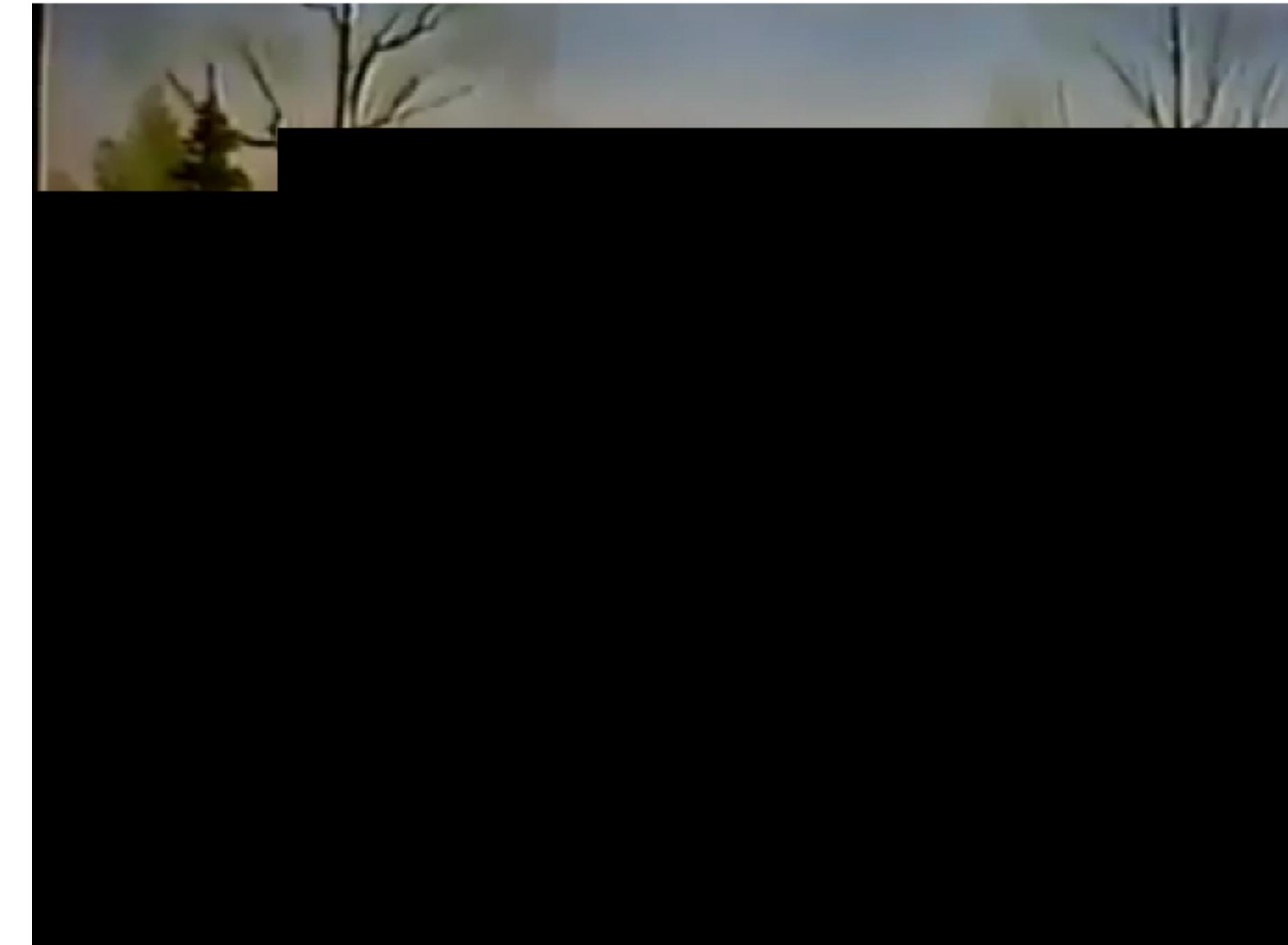
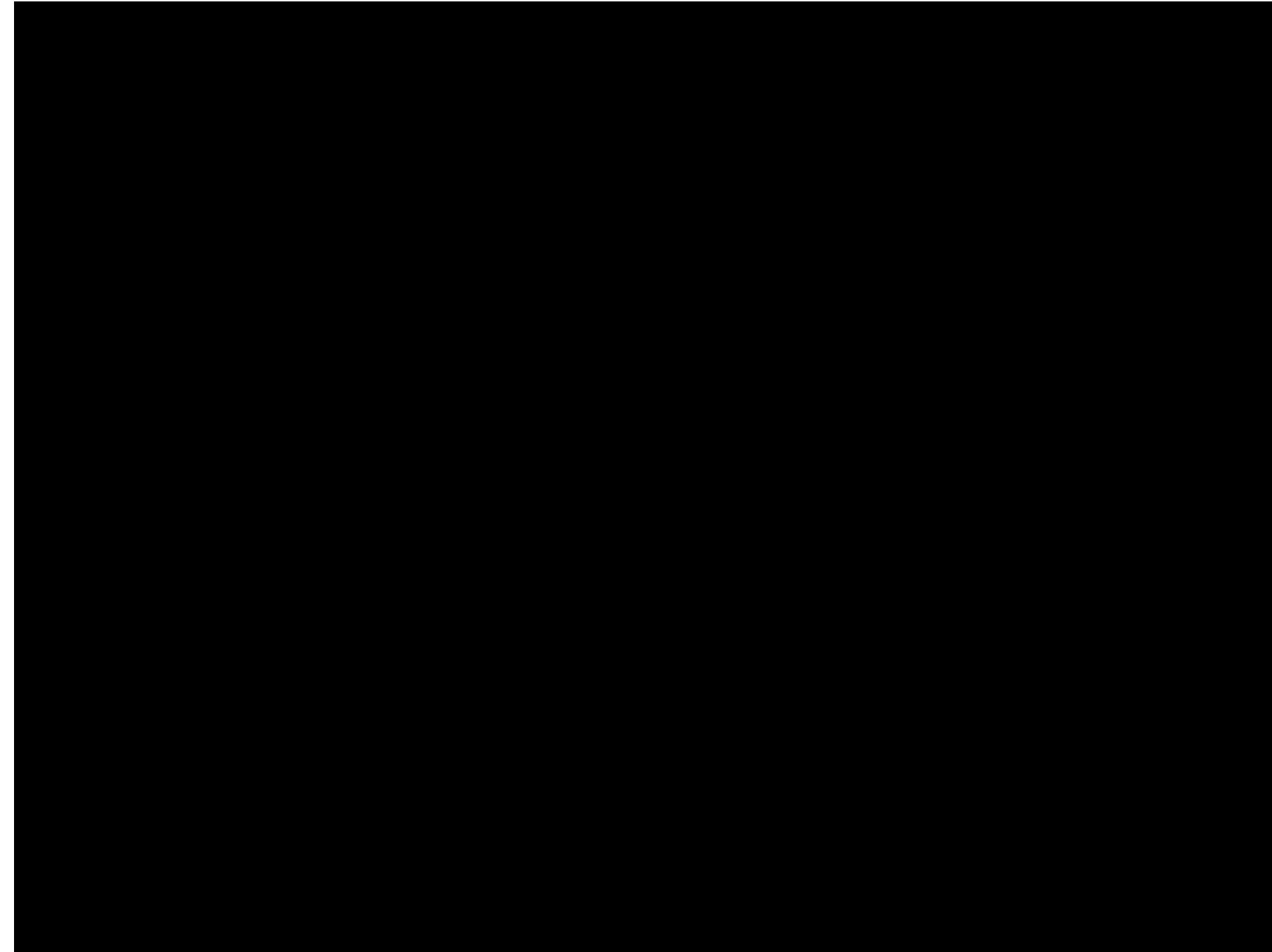
- **Creates a feature space**
- **Exact projection to feature space**
- **Exact likelihood evaluation**

Negatives

- **Only a limited set of functions (invertible and Jacobian easy to compute)**
- **Currently takes longer and needs more parameters than GANs for same quality**

Autoregressive Models

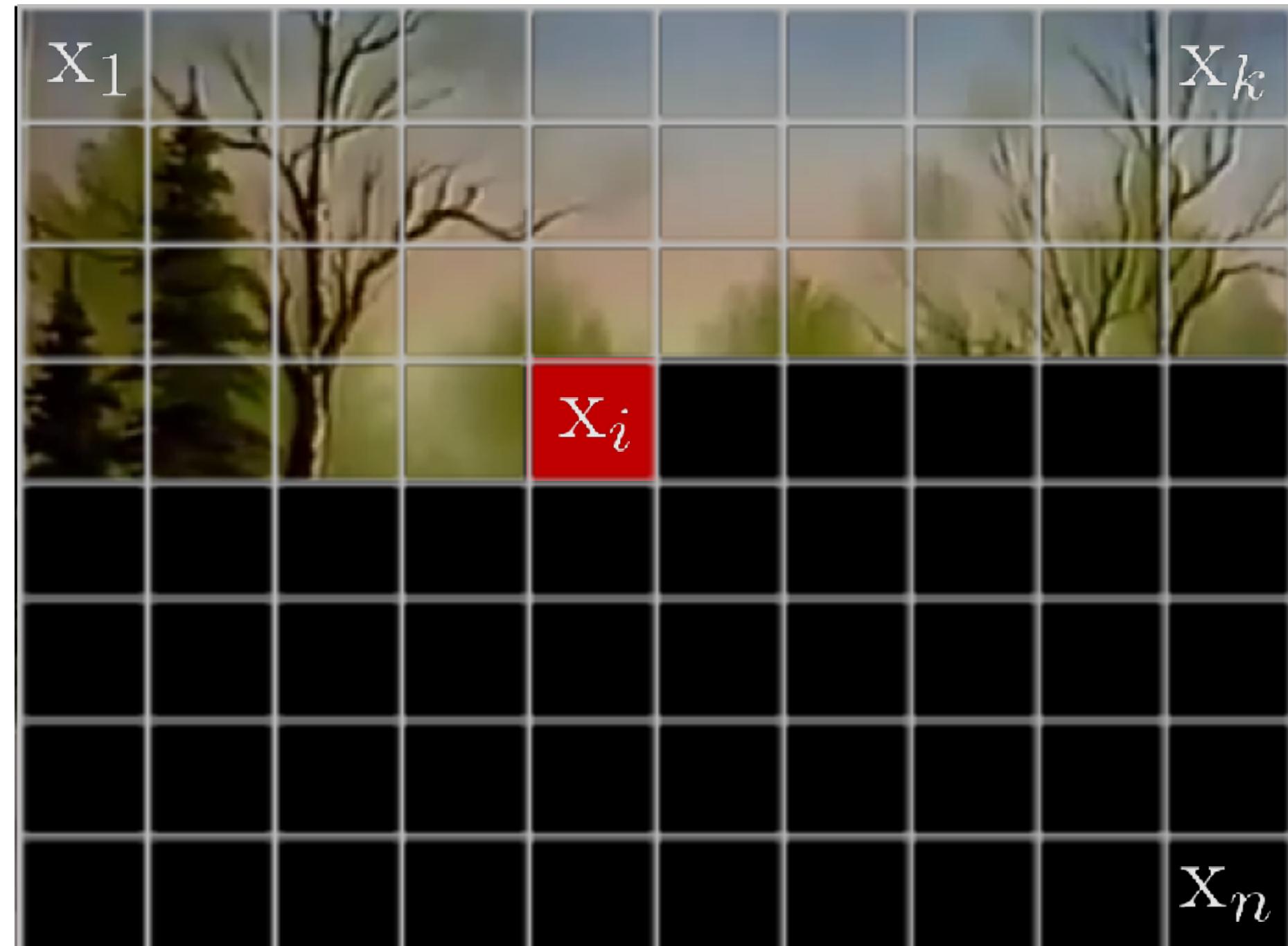
- **Create output step-by-step**
- **Each step depends on the output of all previous steps**



Video Credit: YouTube user *karwan kalary*, *Bob Ross Time Lapse*

Autoregressive Models

- **Create output step-by-step**
- **Each step depends on the output of all previous steps**



$$x = [x_1, x_2, \dots, x_n]$$

$$p_\theta(x) = p_\theta(x_1, x_2, \dots, x_n)$$

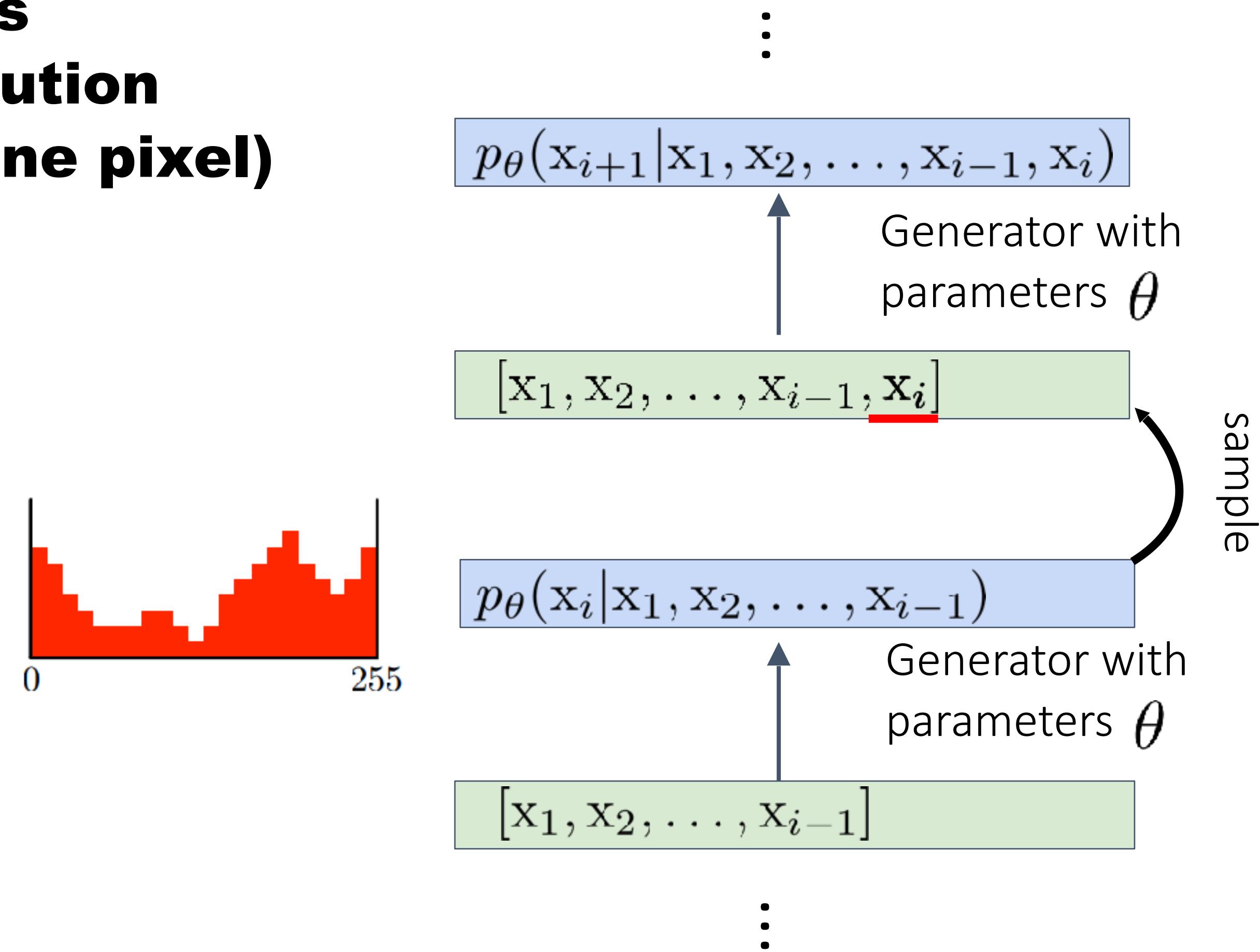
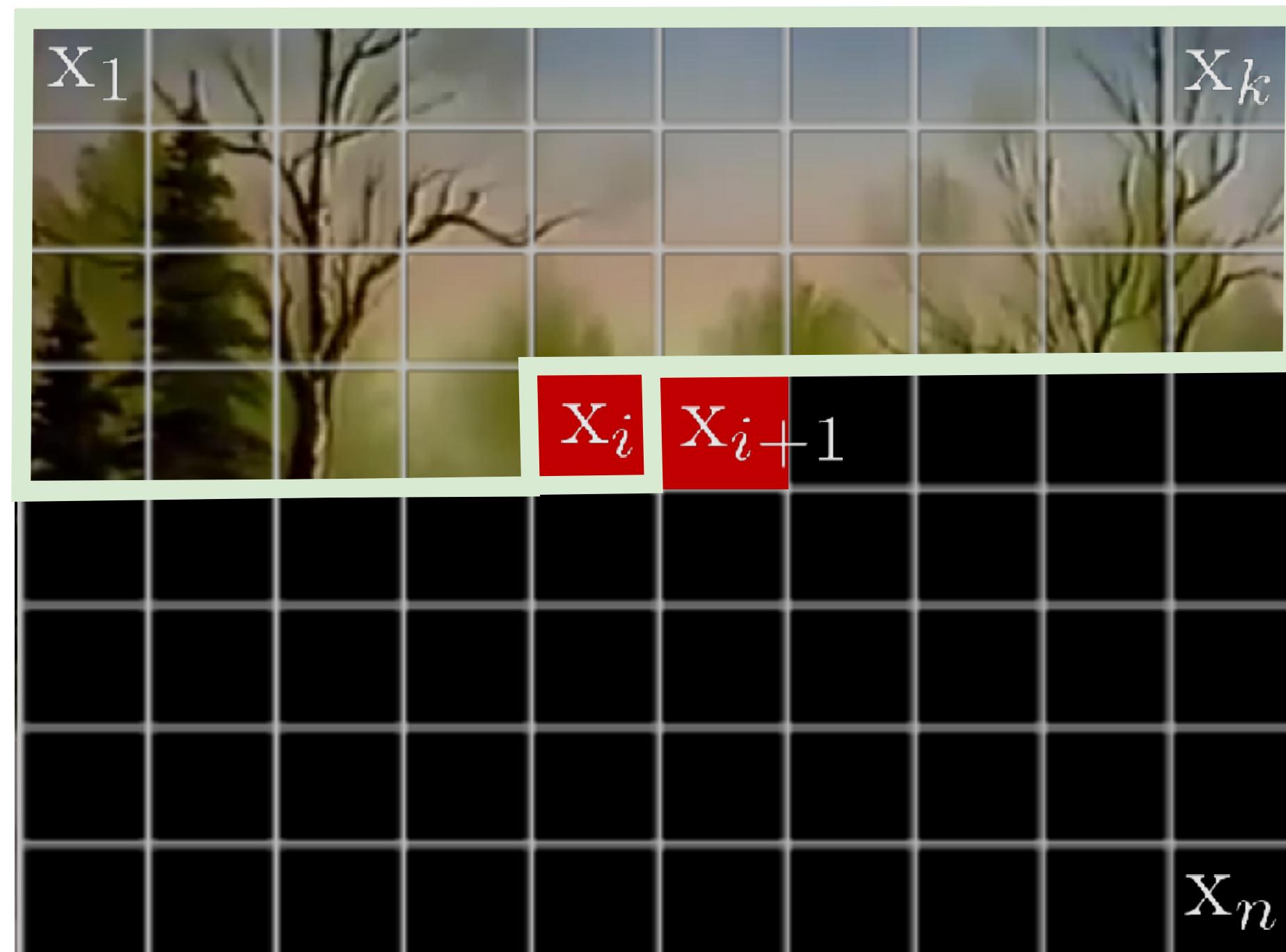
Chain rule of probability:

$$p_\theta(x) = \prod_{i=1}^n p_\theta(x_i | x_1, x_2, \dots, x_{i-1})$$

Video Credit: YouTube user *karwan kalary*, *Bob Ross Time Lapse*

Autoregressive Models

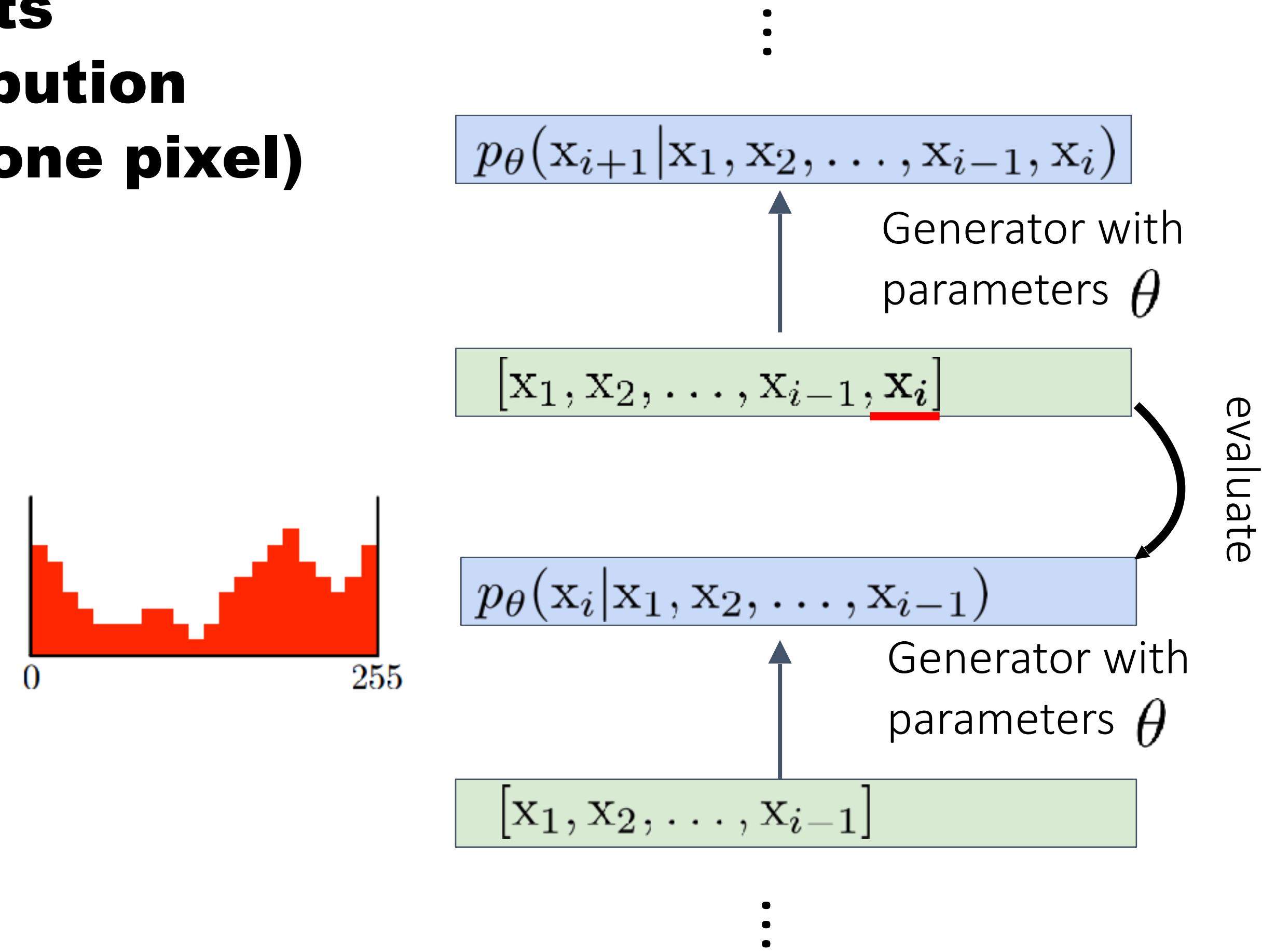
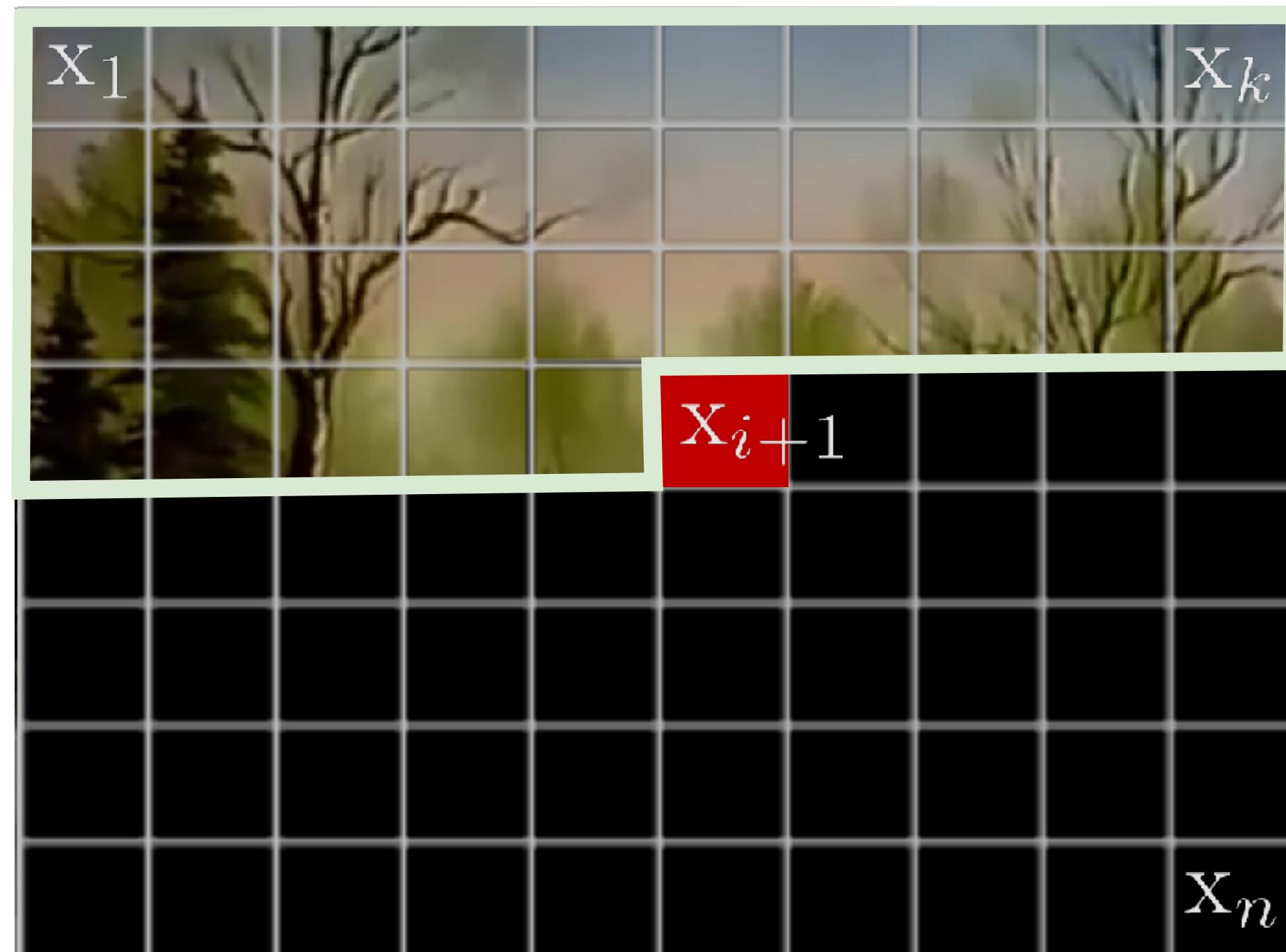
- In each step, the model outputs a low-dimensional prob. distribution (e.g. over intensity values for one pixel)**



Video Credit: YouTube user *karwan kalary*, *Bob Ross Time Lapse*

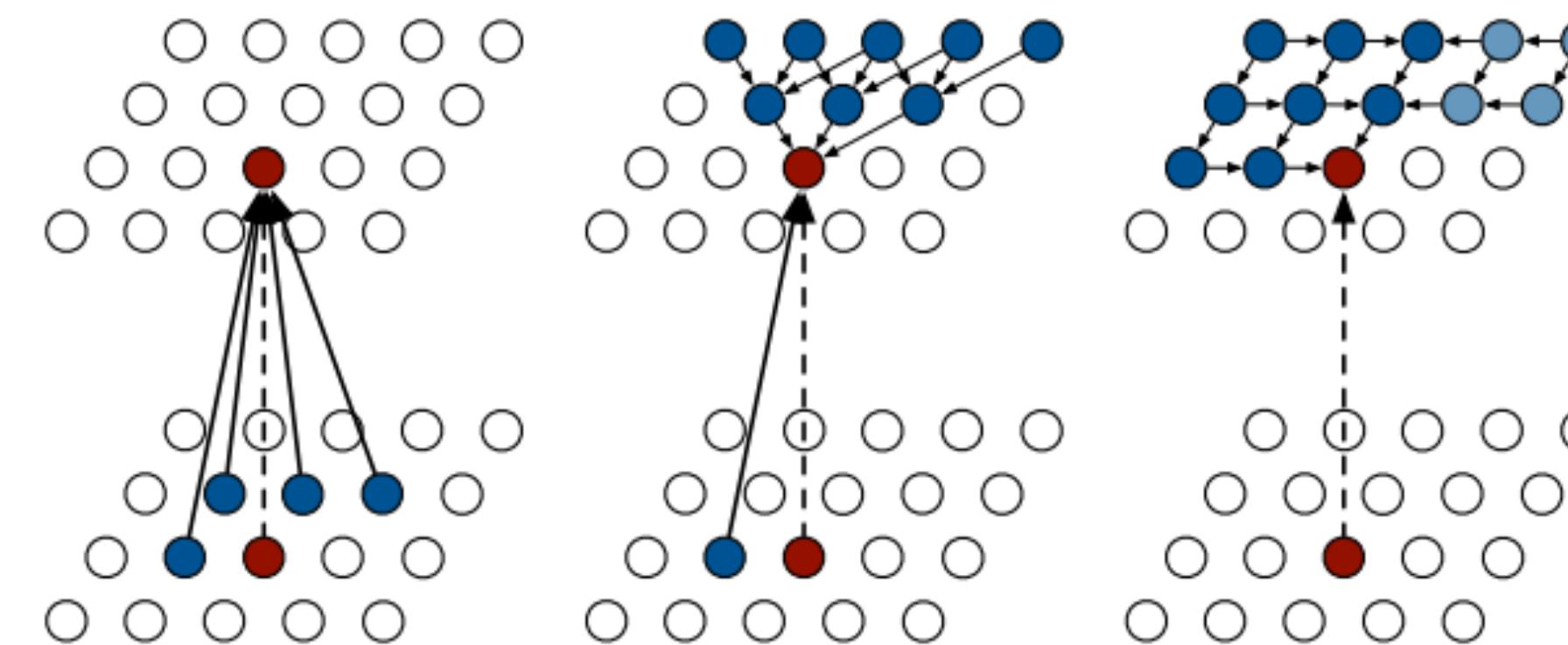
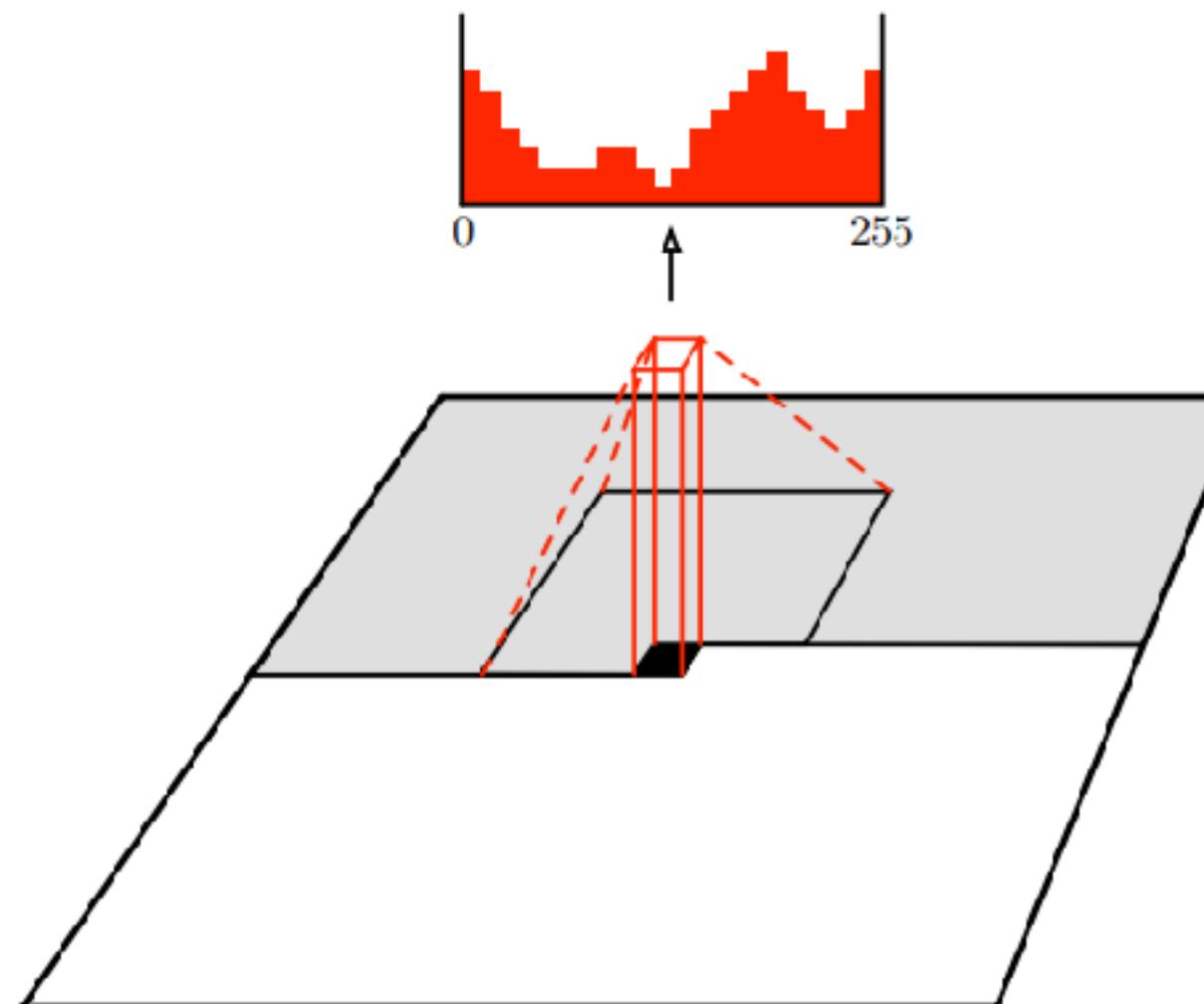
Autoregressive Models

- In each step, the model outputs a low-dimensional prob. distribution (e.g. over intensity values for one pixel)**

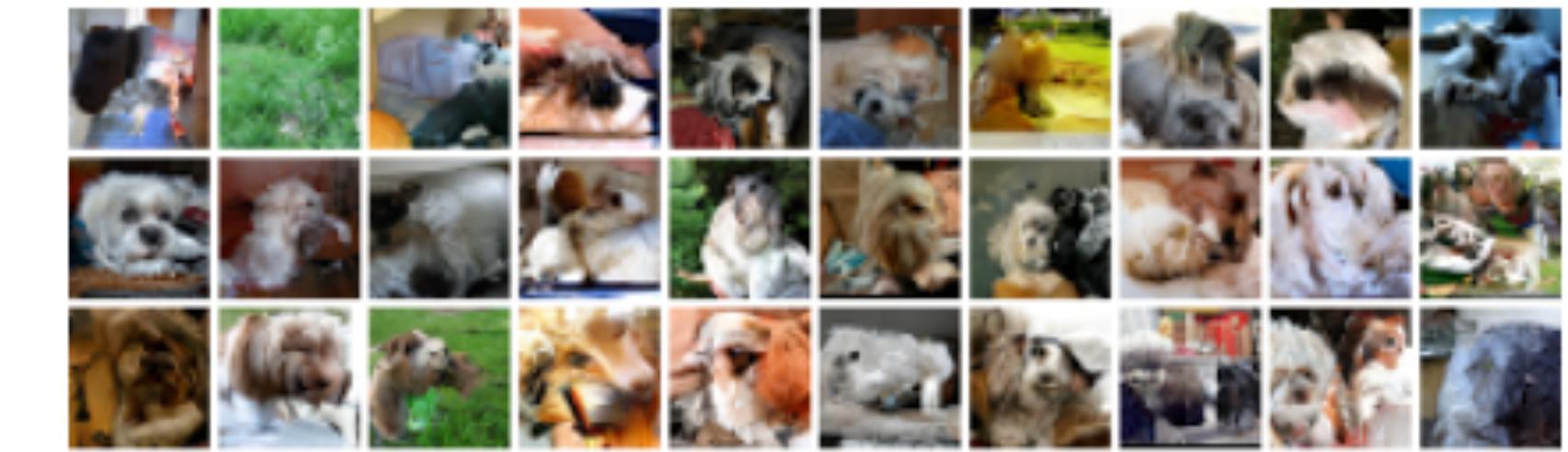


Example: PixelRNN and PixelCNN

- Recursive network that has an input and a state (LSTM)
- Only recent steps are used as input, the state summarizes older steps



Sandbar



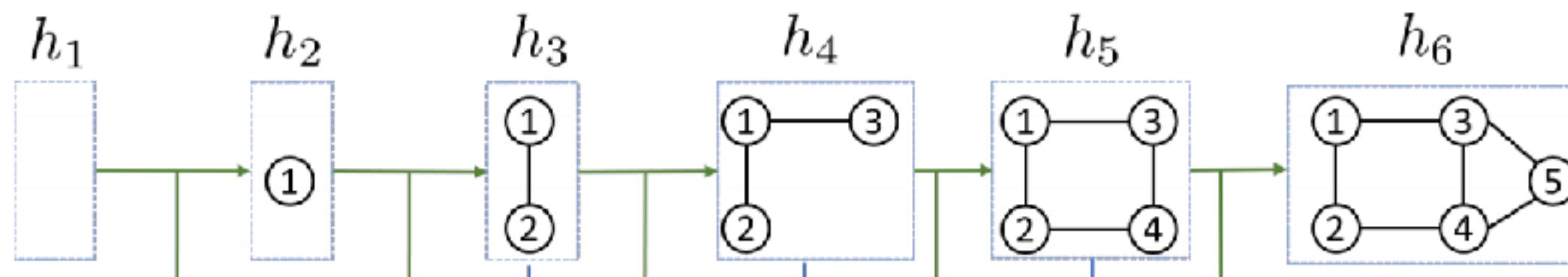
Lhasa Apso (dog)



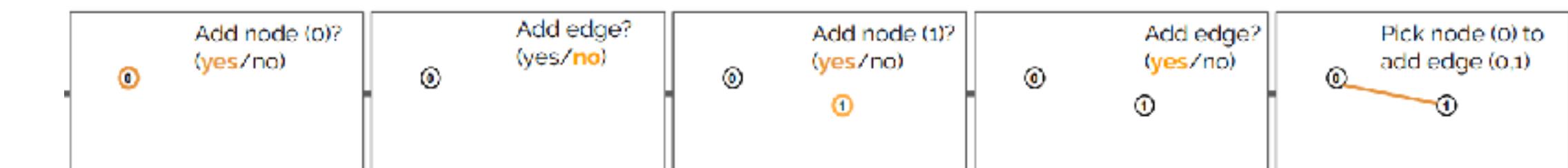
Brown bear

Example: Graph Generation

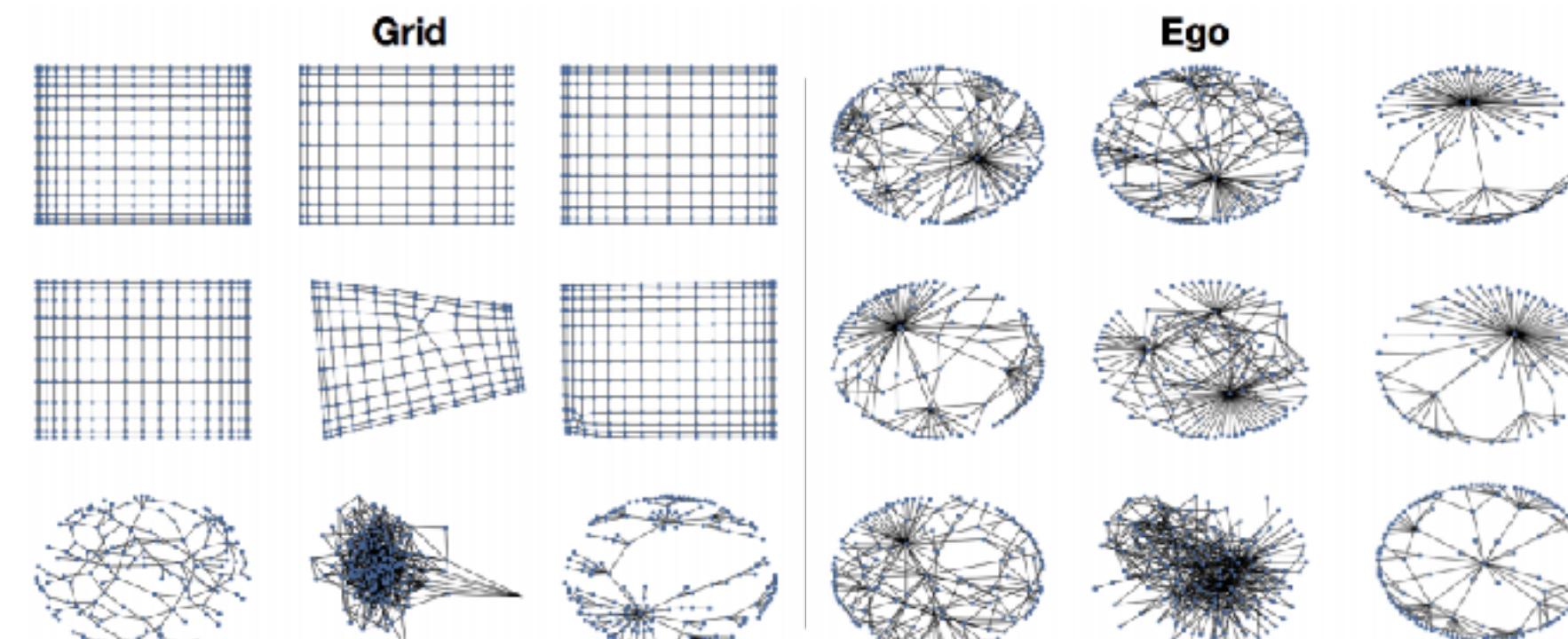
GraphRNN



Learning Deep Generative Models of Graphs



Training Set



Training Set

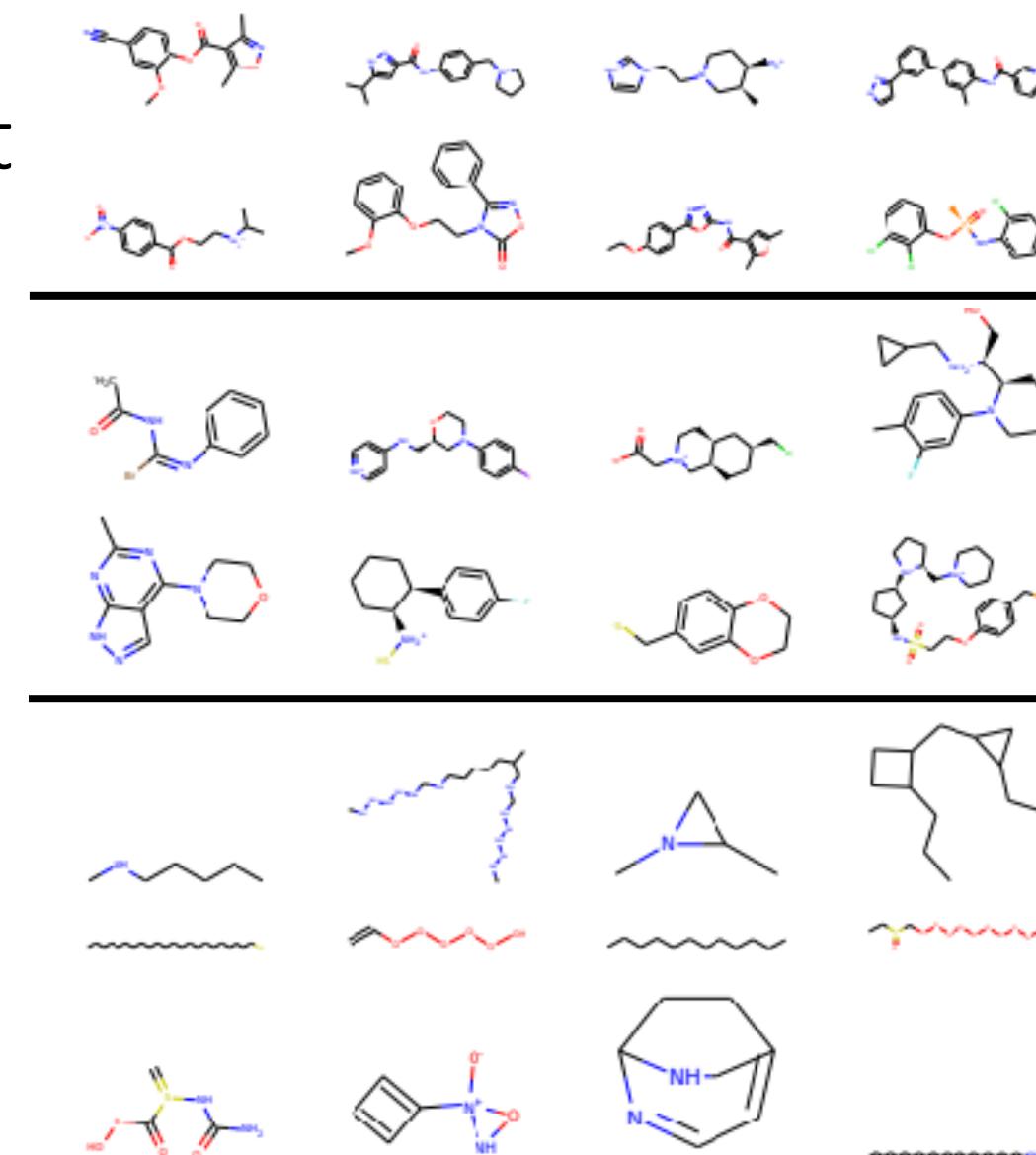


Image Credit:

Conditional Image Generation with PixelCNN Decoders, Oord et al.
GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models, You et al.
Learning Deep Generative Models of Graphs , Li et al.

Summary: Autoregressive Models

Positives

- **Flexible output length**
- **Exact likelihood evaluation**

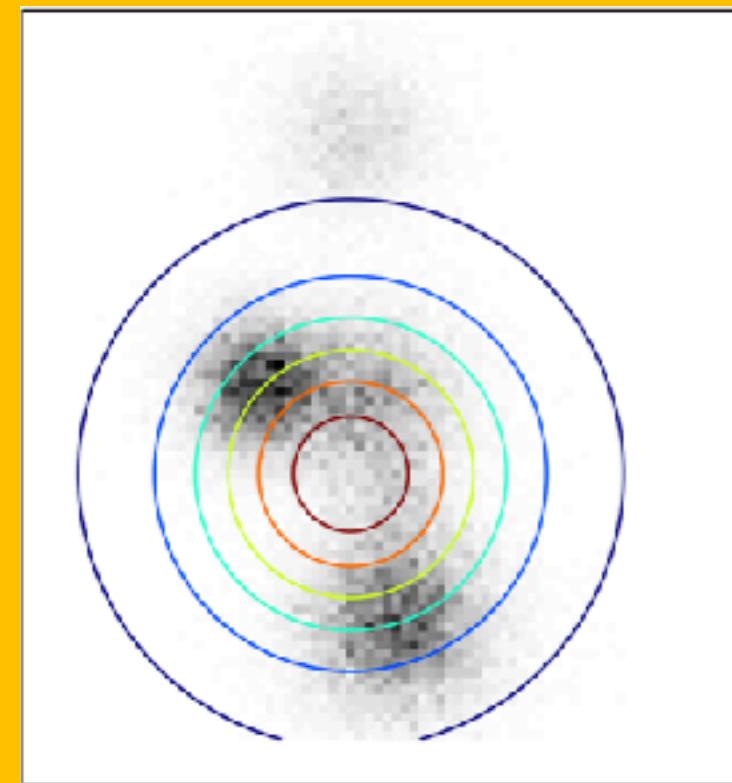
Negatives

- **No feature space**
- **Sequential generation (usually slow)**

Generative Models

How do we measure the similarity of $p_\theta(x)$ and $p_{\text{data}}(x)$?

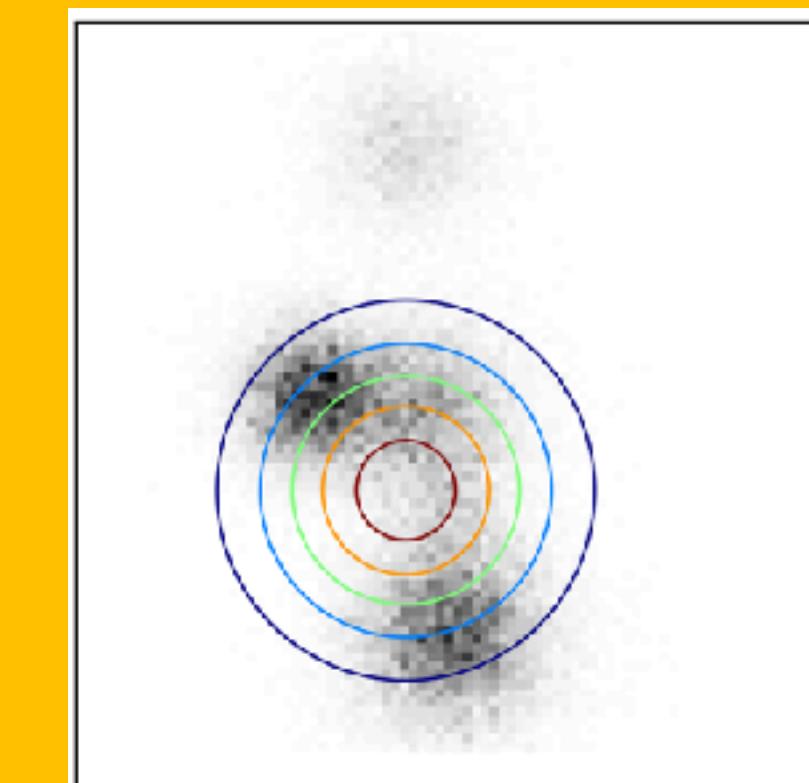
1) Likelihood of data samples in $p_\theta(x)$



$$p_{\text{data}}(x)$$

$$\approx KL(p_{\text{data}} \parallel p_\theta)$$

2) Adversarial game



$$\approx JS(p_{\text{data}} \parallel p_\theta)$$

Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

Generative Adversarial Networks (GANs)

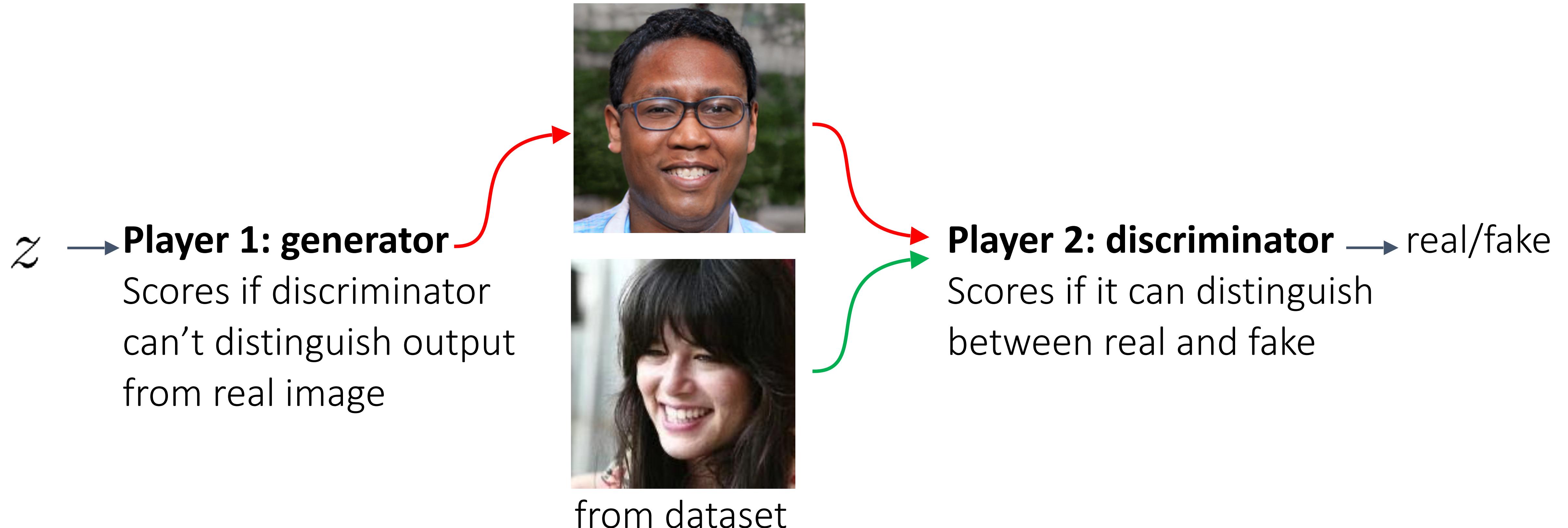


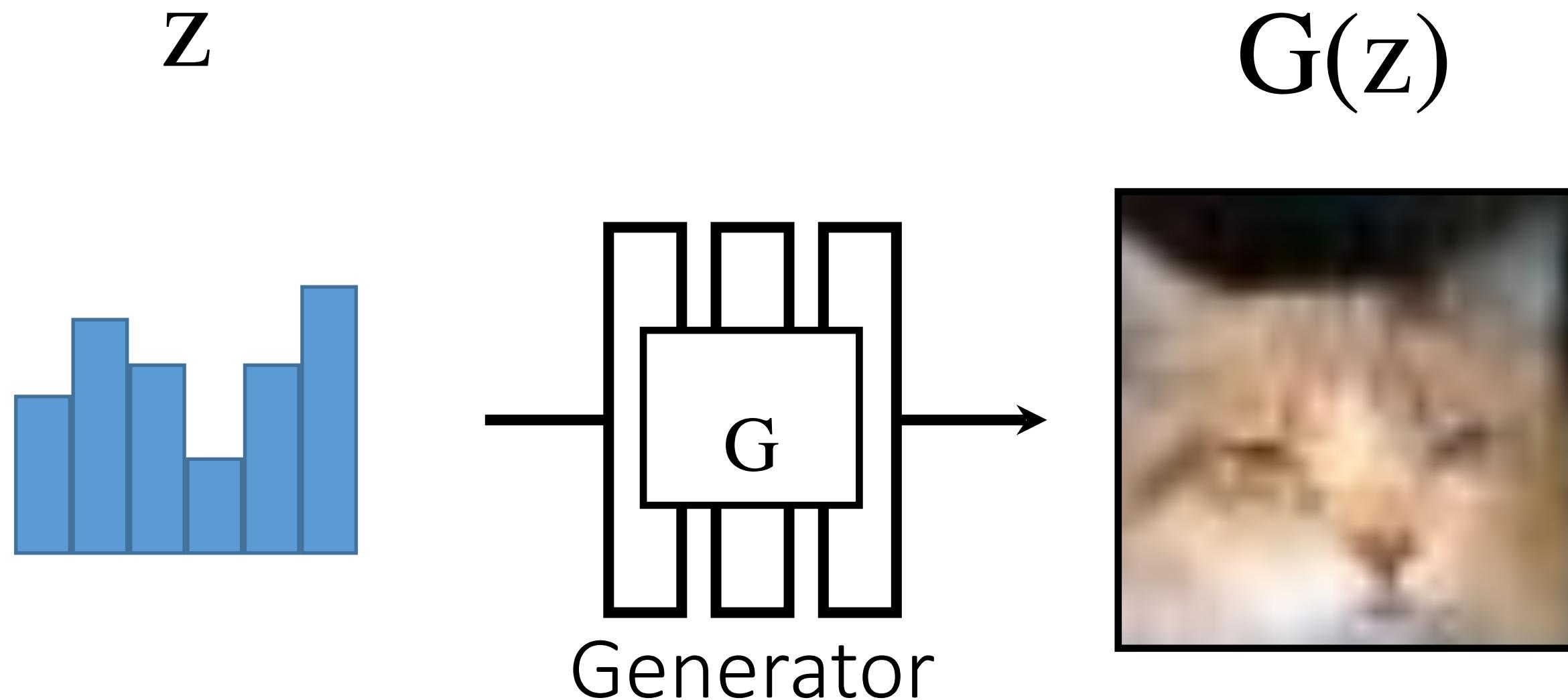
Image credit: *A Style-Based Generator Architecture for Generative Adversarial Networks*, Karras et al.

Generative Adversarial Networks (GANs)

Player 1: generator
Scores if discriminator
can't distinguish output
from real image

Player 2: discriminator
Scores if it can distinguish
between real and fake

Generator



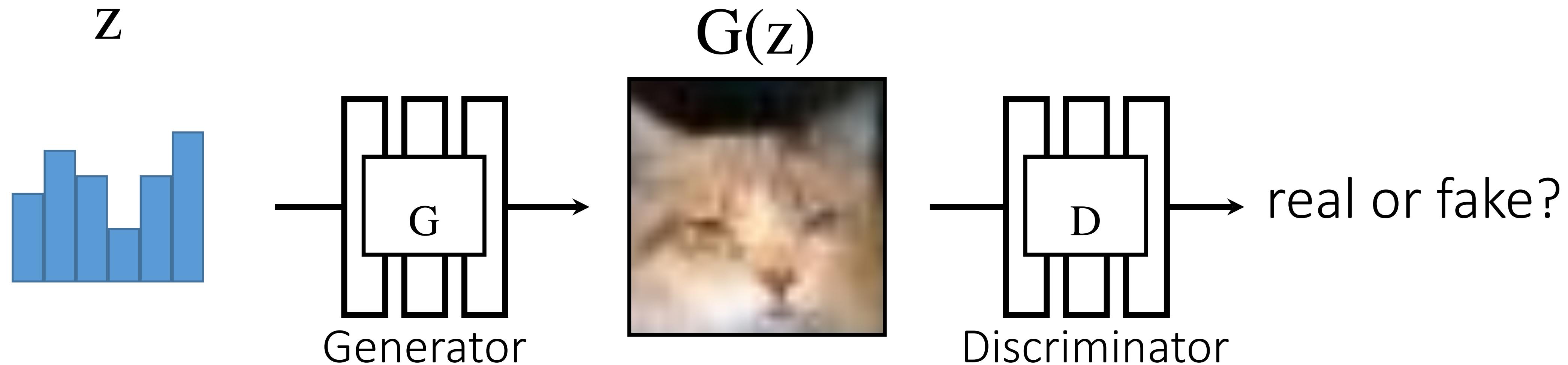
G : generate fake samples that can fool D

cat credit: aleju/cat-generator



slide credit: Phillip Isola & Jun-Yan Zhu

Generator + Discriminator



slide credit: Phillip Isola & Jun-Yan Zhu



$$\min_{G} \max_{D} \mathbb{E}_{z,x} [$$

Z

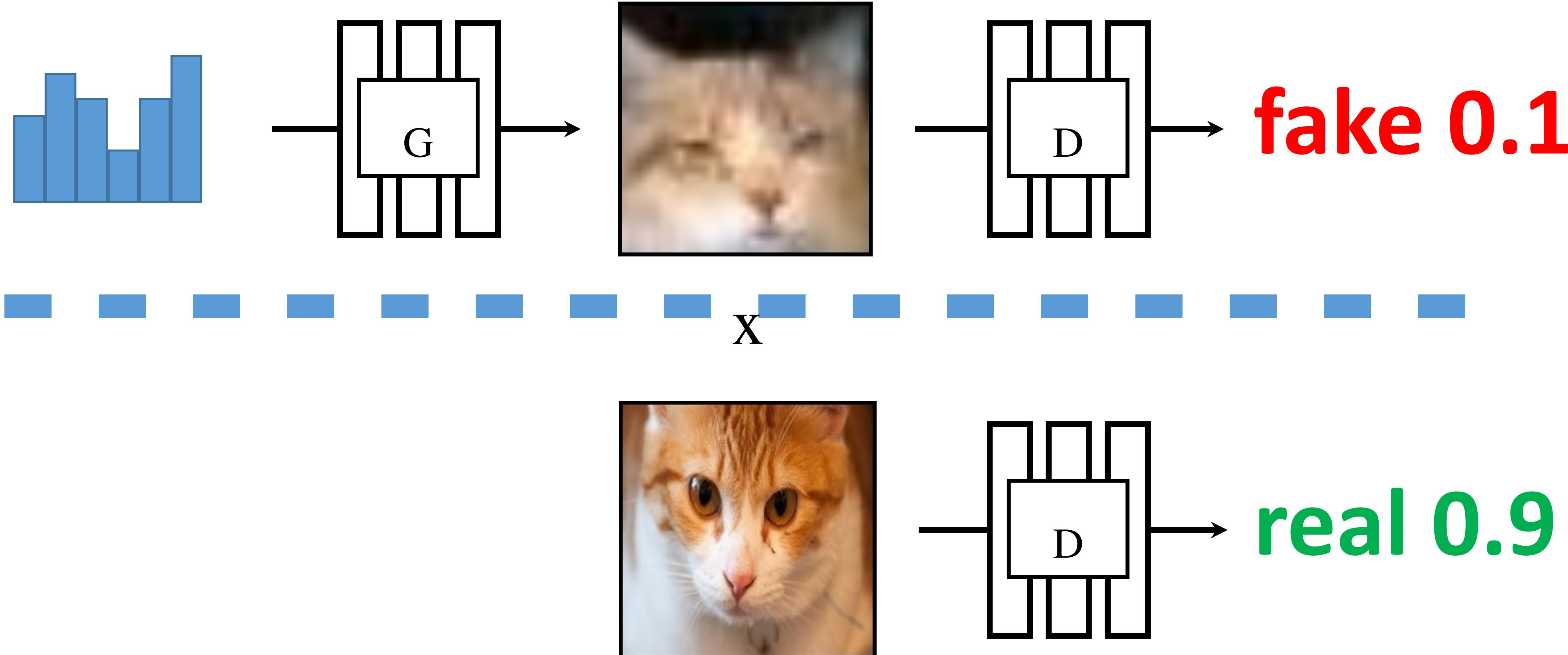
G(z)



$$\min_{G} \max_{D} \mathbb{E}_{z,x} \left[\log D(G(z)) - \text{fake} \right]$$

Z

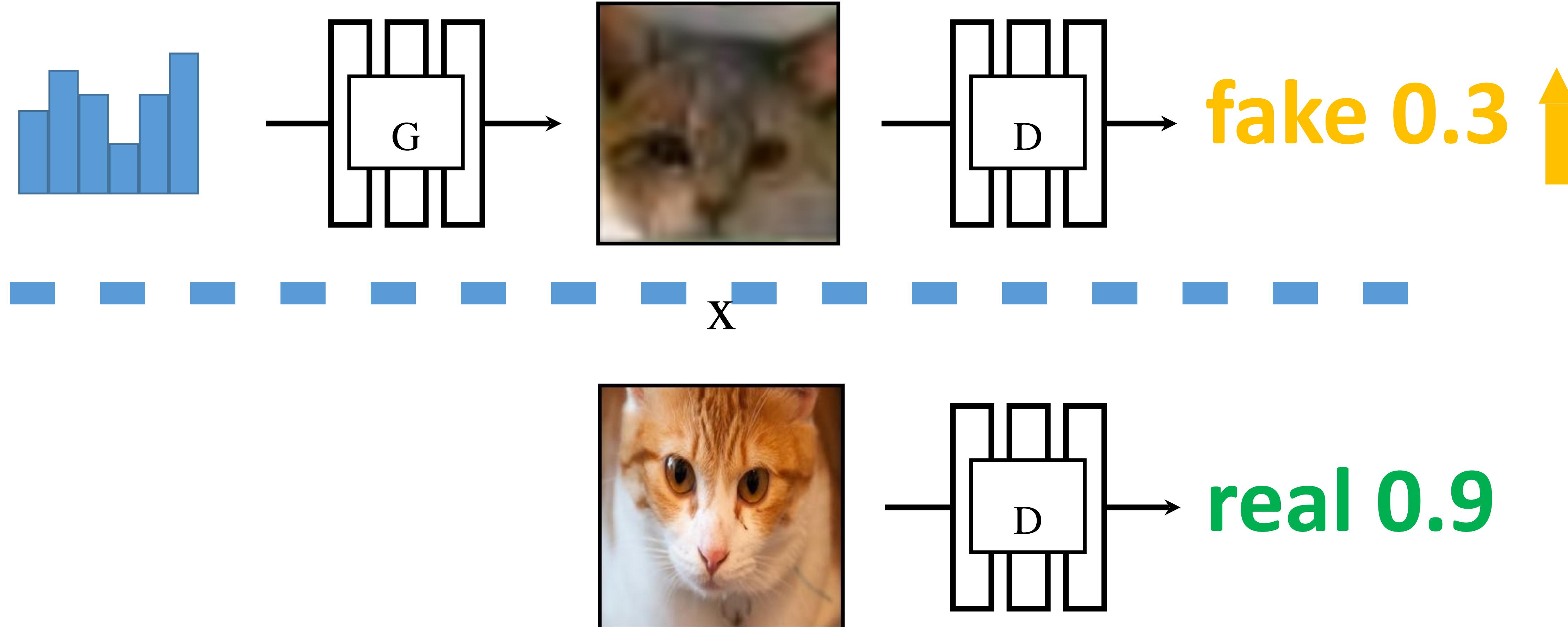
G(z)



$$\min_{G} \max_{D} \mathbb{E}_{z,x} \left[\underbrace{\log D(G(z))}_{\text{fake}} + \underbrace{\log(1 - D(x))}_{\text{real}} \right]$$

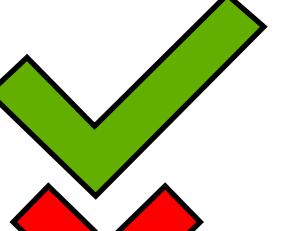
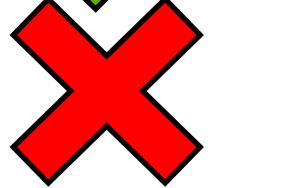
Z

G(z)



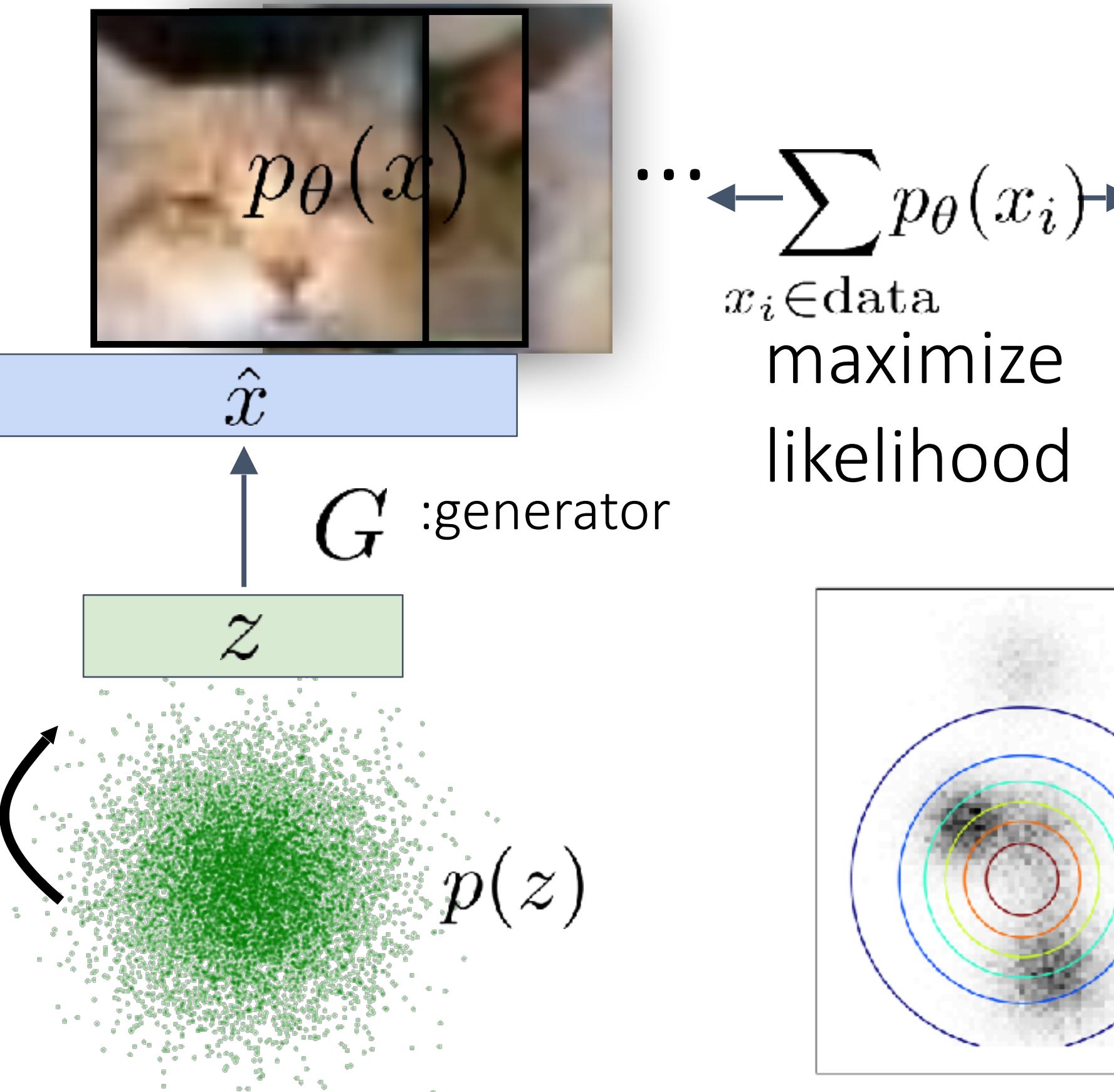
$$\min_{G} \max_{D} \mathbb{E}_{z,x} \left[\log D(G(z)) + \log(1 - D(x)) \right]$$

Update G

Sample? 
Evaluate? 

Generated Distributions of GANs vs ML

VAEs or Norm. Flows



$$\approx KL(p_{\text{data}} \parallel p_\theta)$$

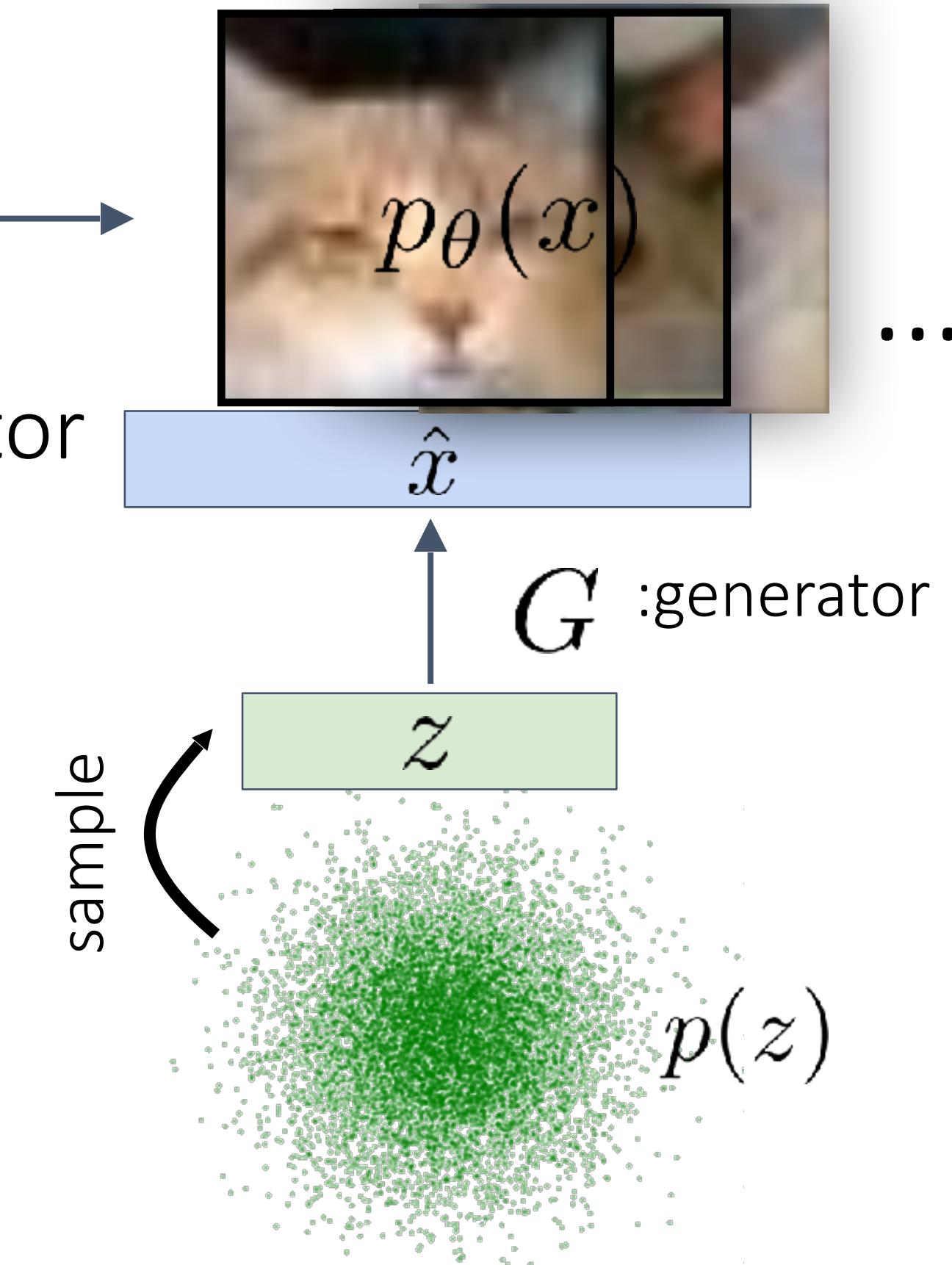


$\cdots \leftarrow D \rightarrow$

discriminator

Generative Modeling

GANs



$$\approx JS(p_{\text{data}} \parallel p_\theta)$$

StyleGAN

Additional Tricks:

- **Coarse-to-fine training**
- **Transformation of $p(z)$ to a more complex distr.**
- ...

style



Image credit: *A Style-Based Generator Architecture for Generative Adversarial Networks*, Karras et al.

Summary: GANs

Positives

- **Creates a feature space**
- **Currently highest-quality results**

Negatives

- **Can be unstable to train**
- **Not guaranteed to cover all of the data distribution**
- **Cannot evaluate likelihood**

Conditional GAN: Pix2Pix

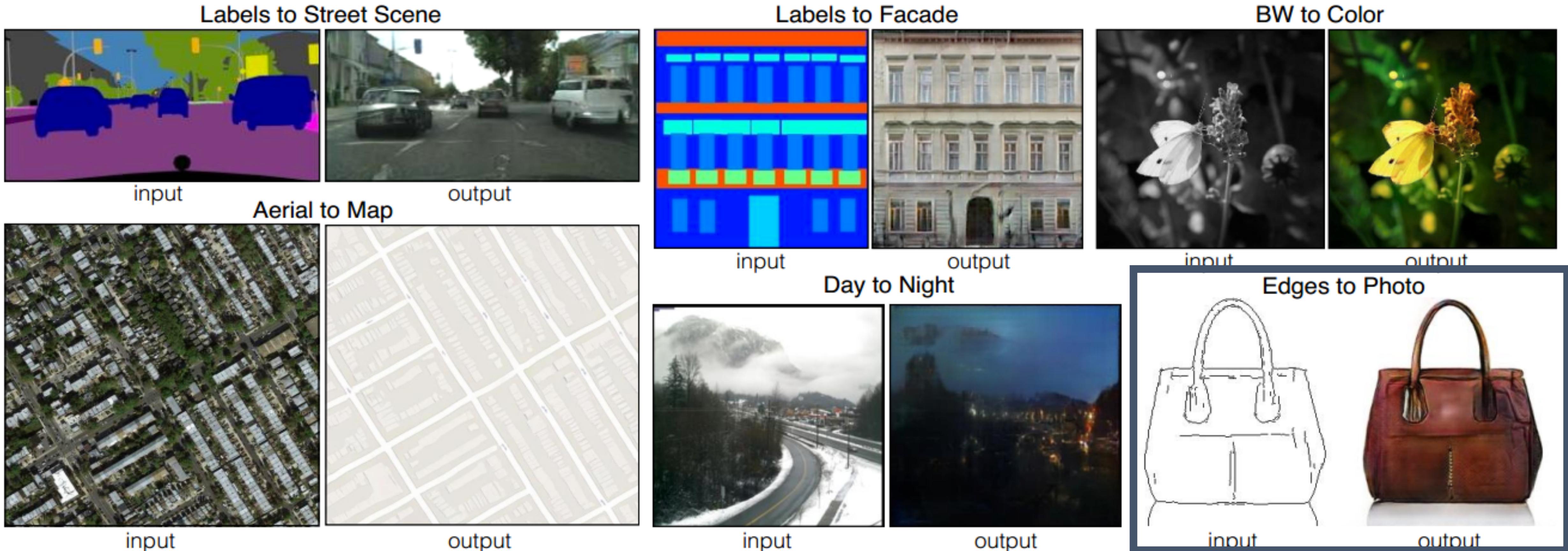
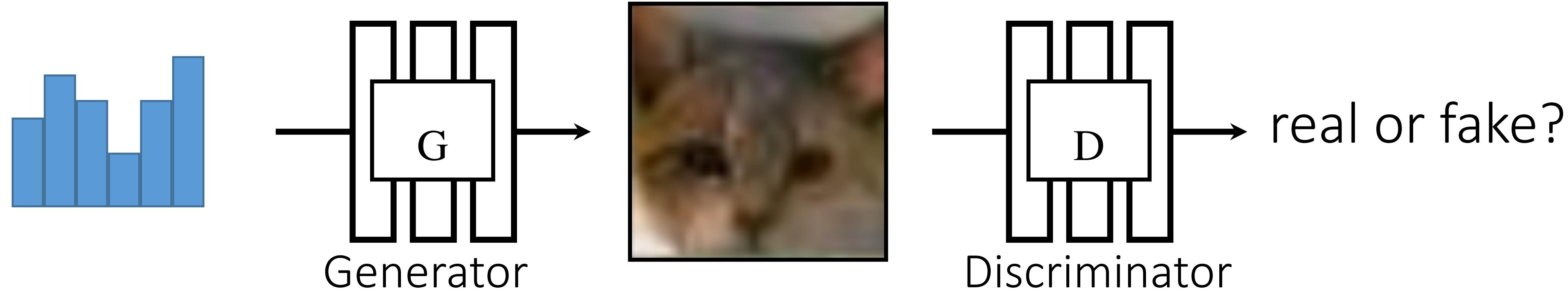


Image-to-image Translation with Conditional Adversarial Nets
Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, Alexei A. Efros. CVPR 2017

Z

G(z)



$$\min_{G} \max_{D} \mathbb{E}_{x,y} \left[\log D(G(x)) + \log(1 - D(y)) \right]$$

X

G(x)



$$\min_{G} \max_{D} \mathbb{E}_{x,y} \left[\log D(G(x)) + \log(1 - D(y)) \right]$$

X

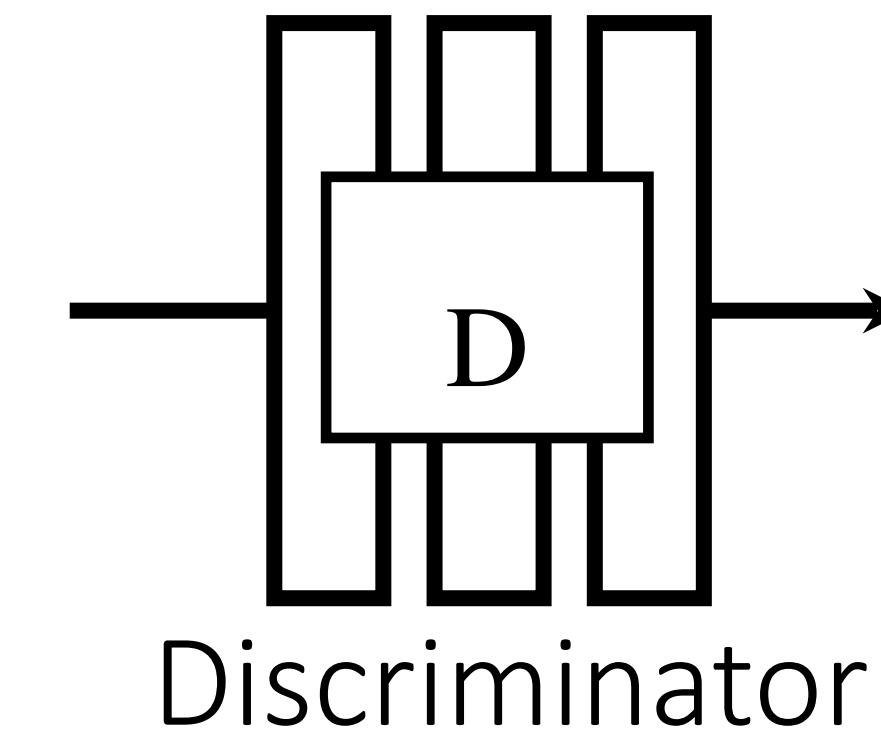
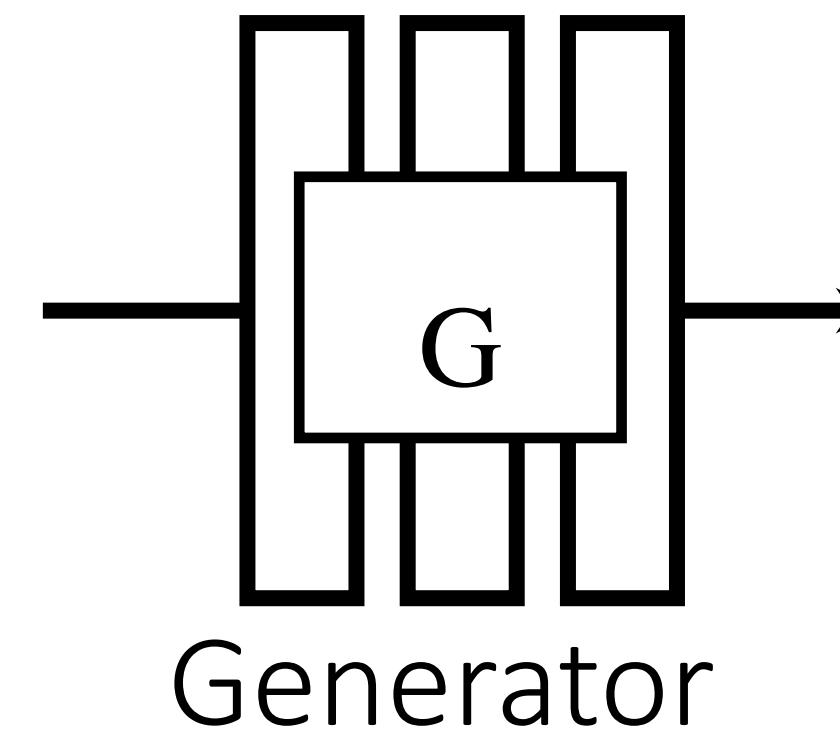
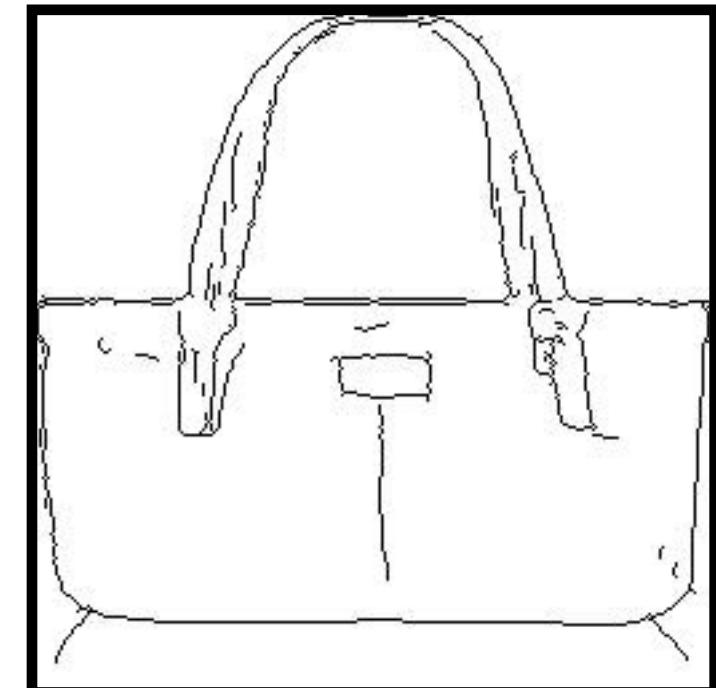
G(x)



$$\min_{G} \max_{D} \mathbb{E}_{x,y} \left[\log D(G(x)) + \log(1 - D(y)) \right]$$

X

G(x)

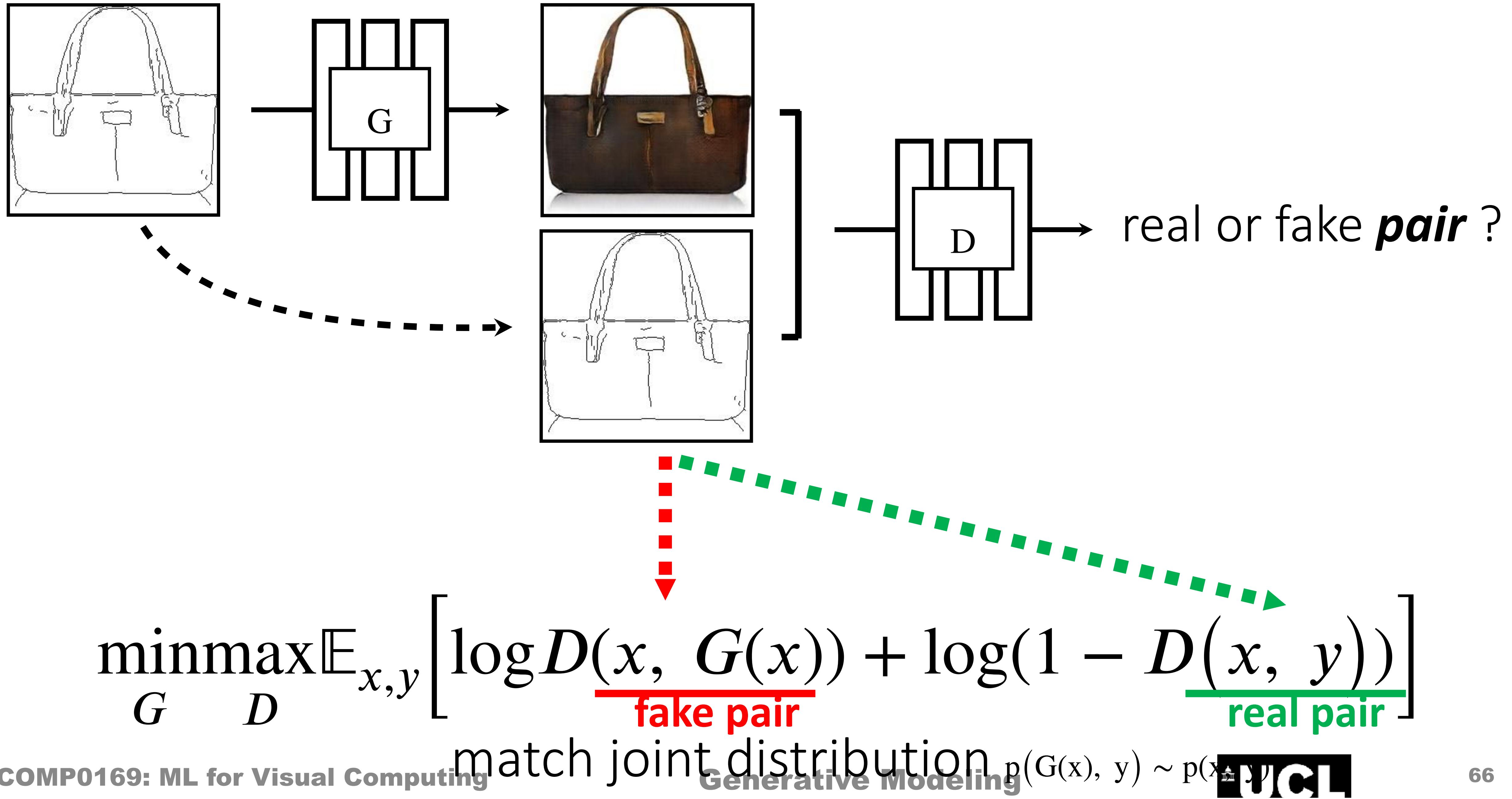


Real too!

$$\min_{G} \max_{D} \mathbb{E}_{x,y} \left[\log D(G(x)) + \log(1 - D(y)) \right]$$

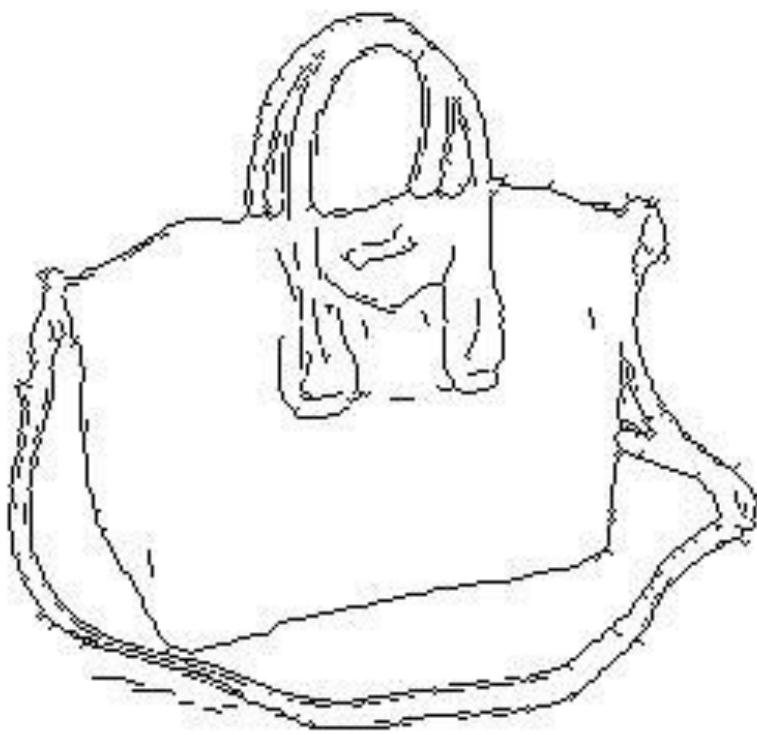
X

G(x)



Edges → Images

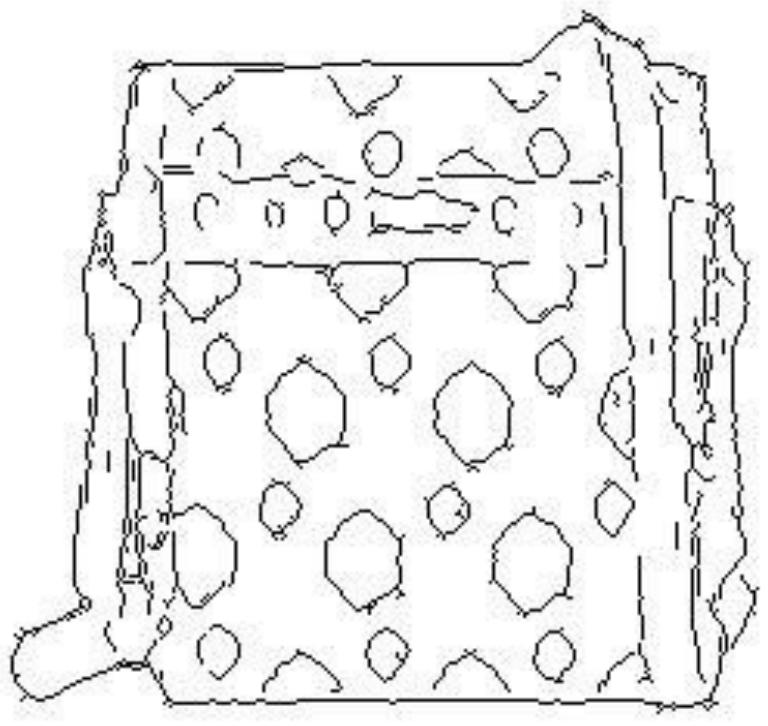
Input



Output



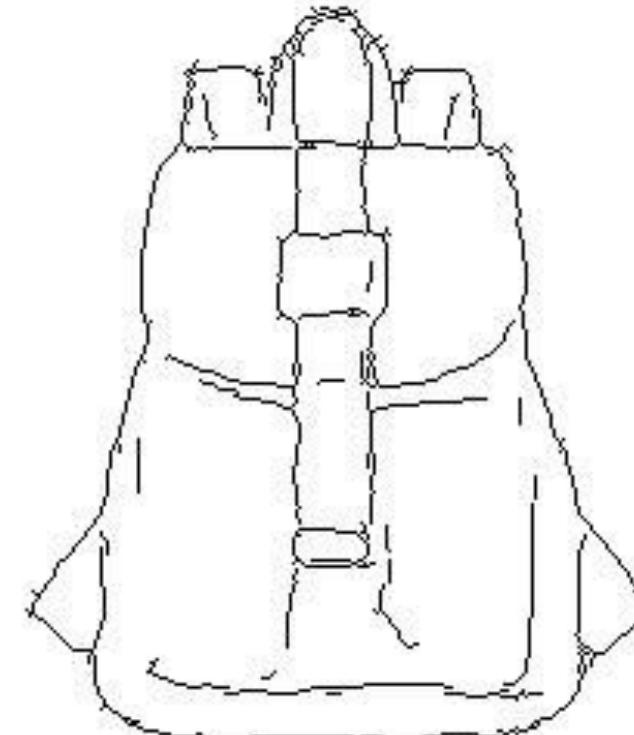
Input



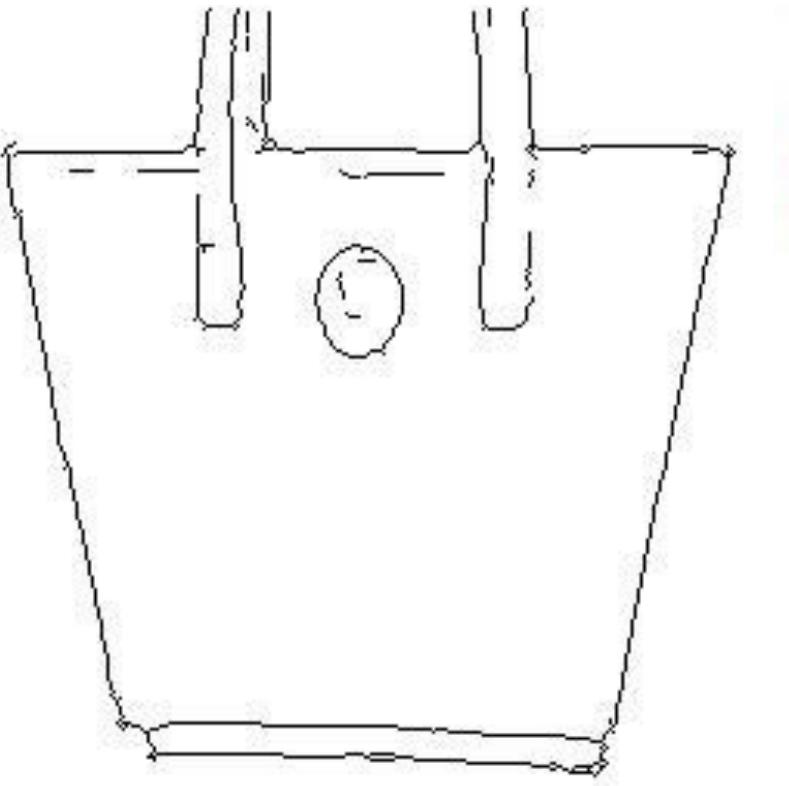
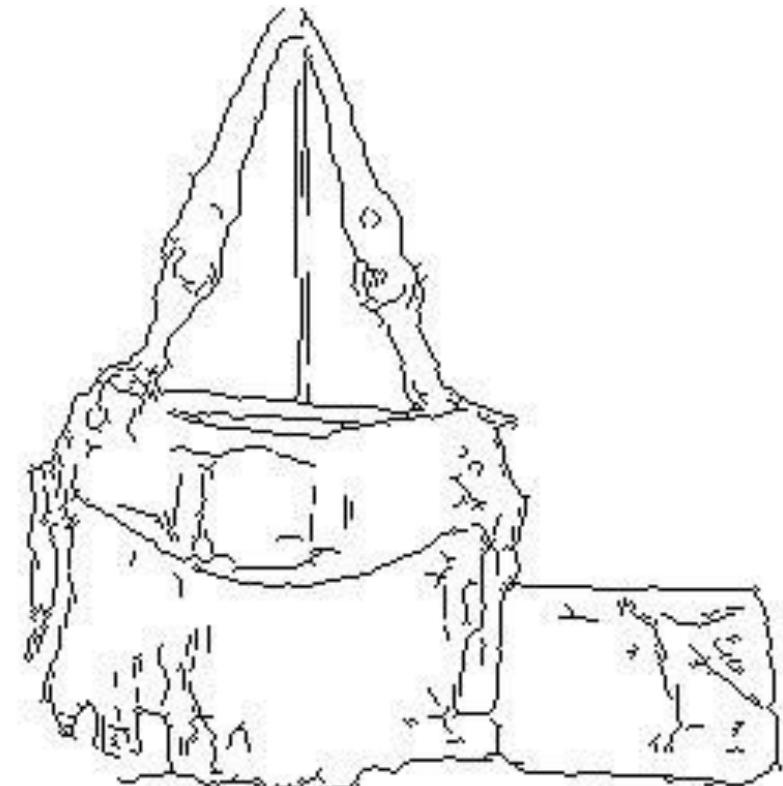
Output



Input



Output

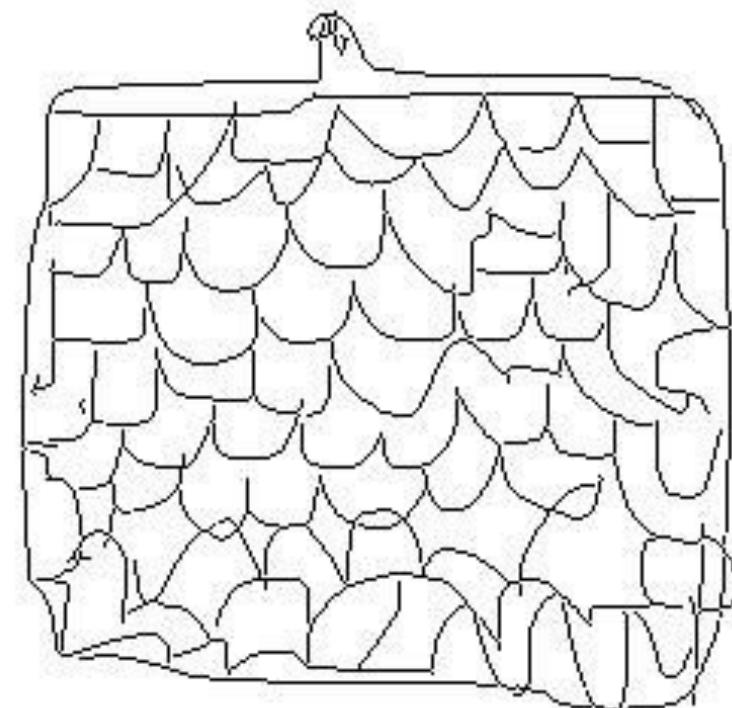


Edges from [Xie & Tu, 2015]

slide credit: Phillip Isola & Jun-Yan Zhu

Sketches → Images

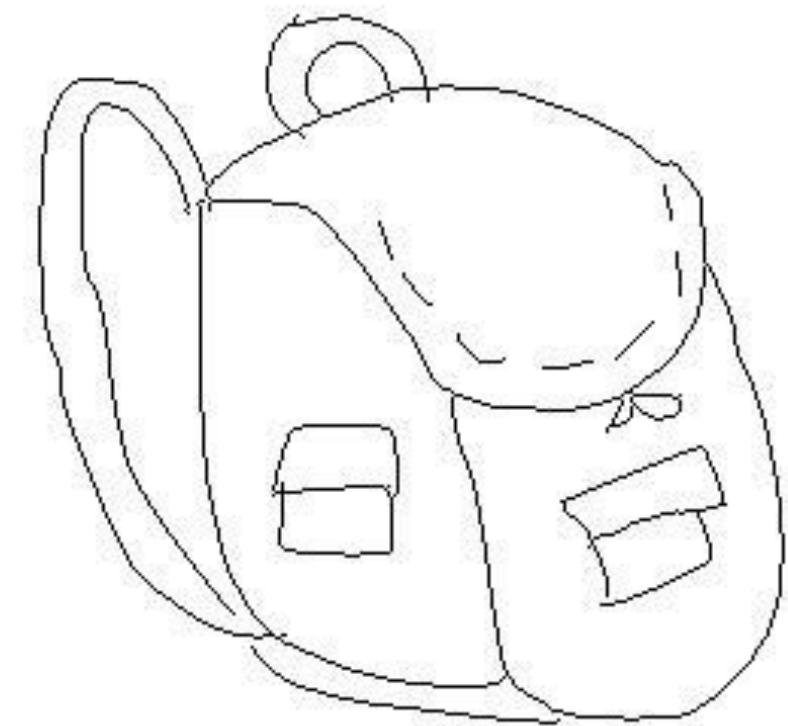
Input



Output



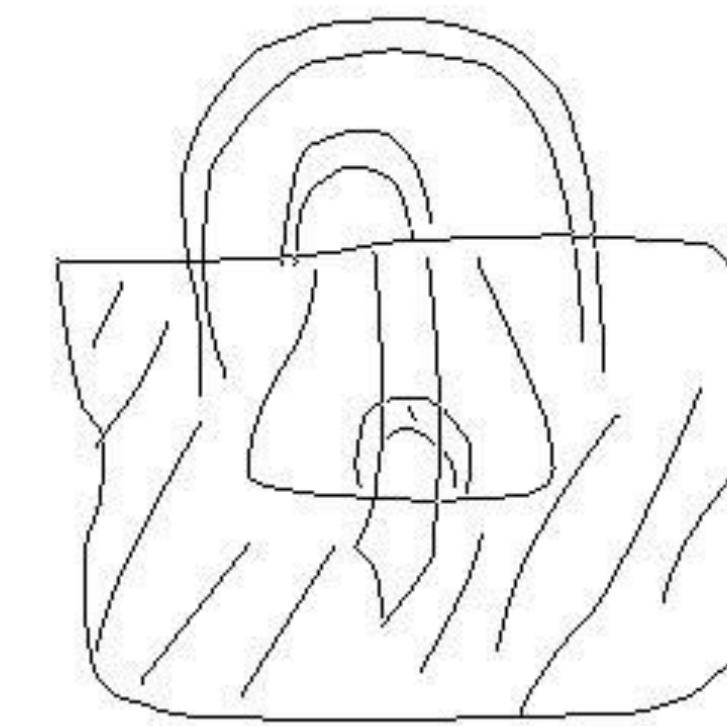
Input



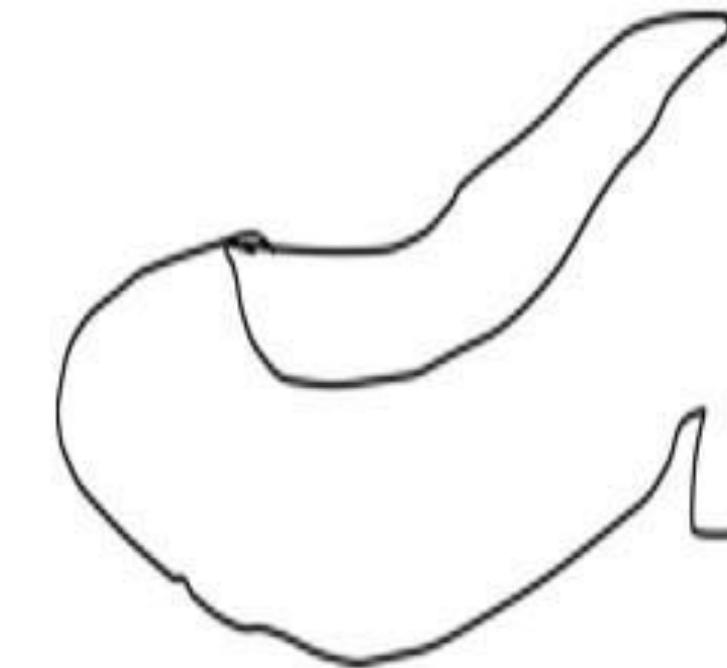
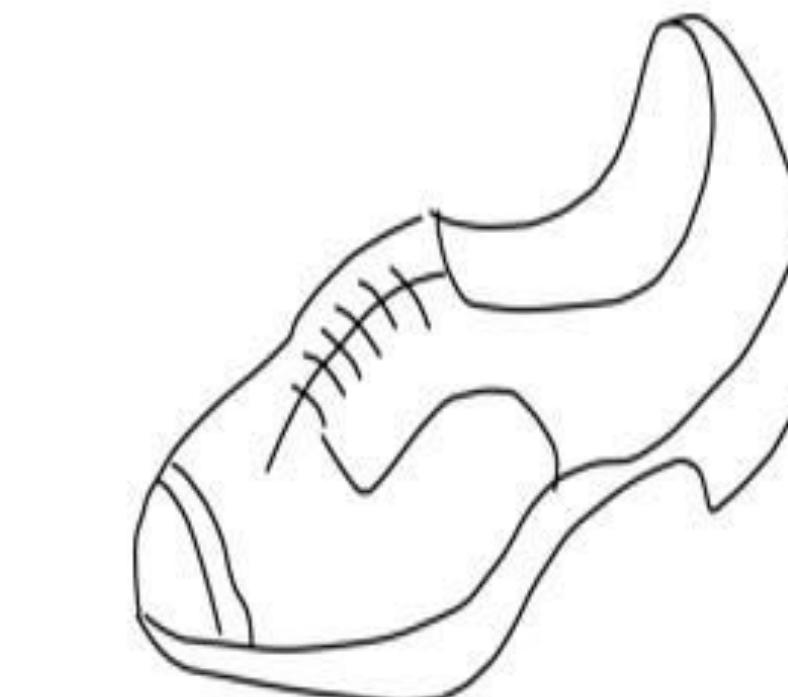
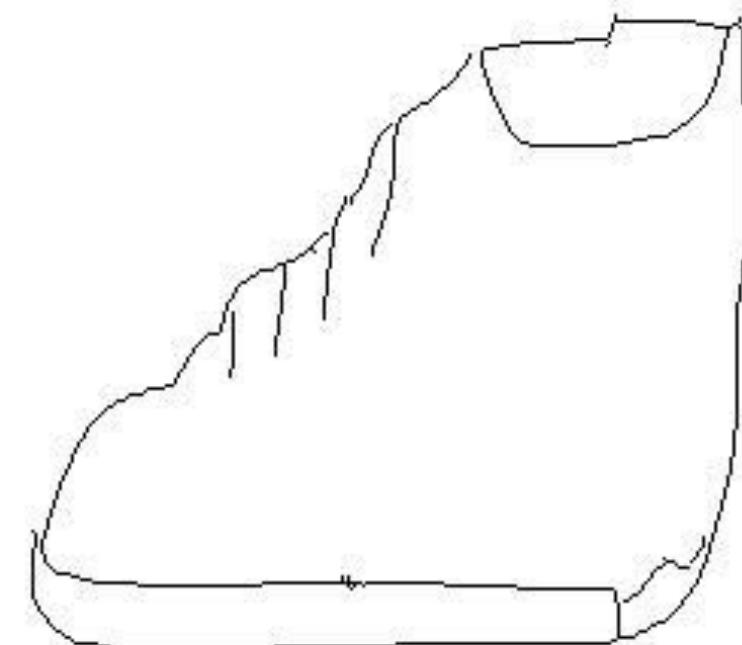
Output



Input



Output

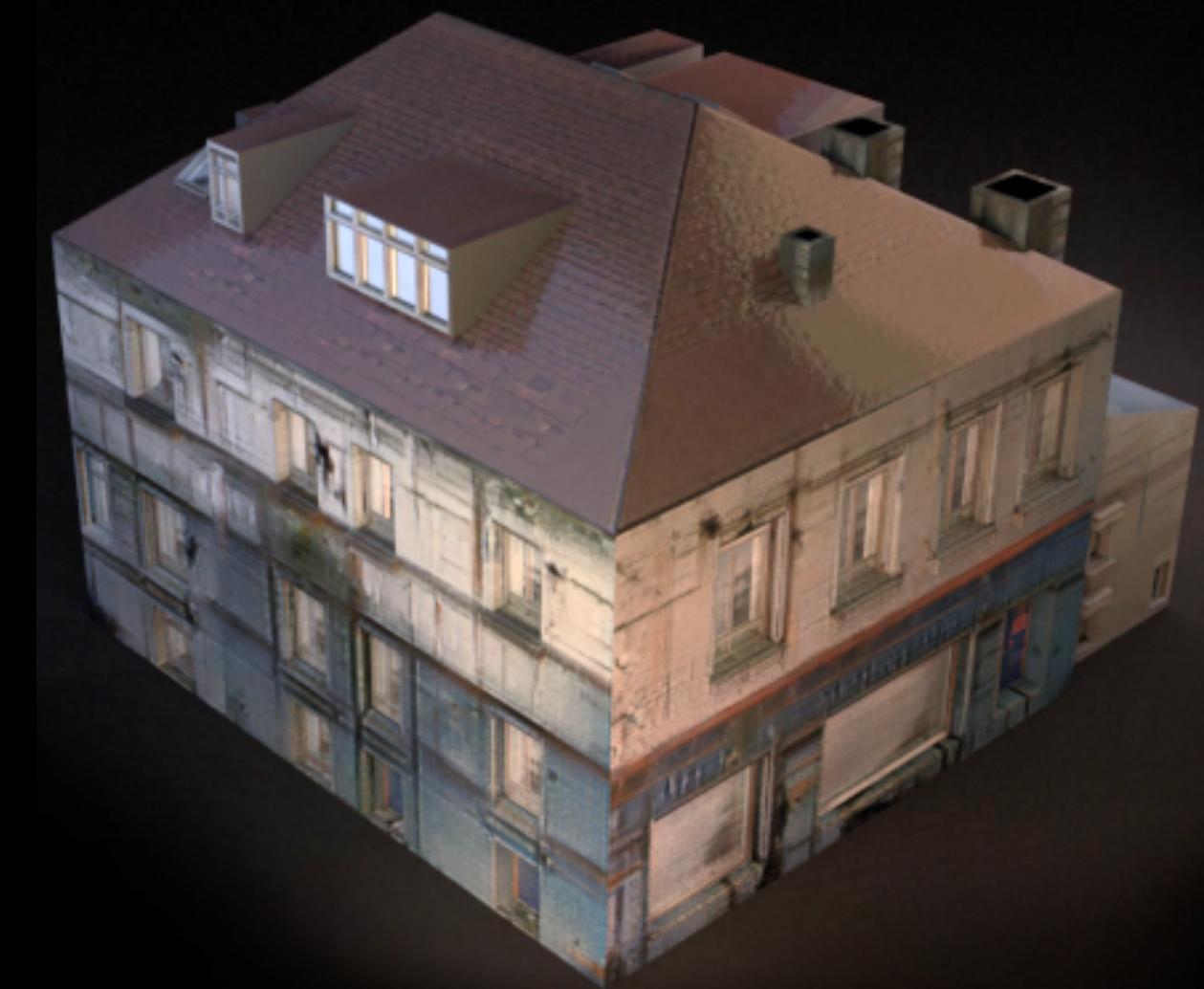
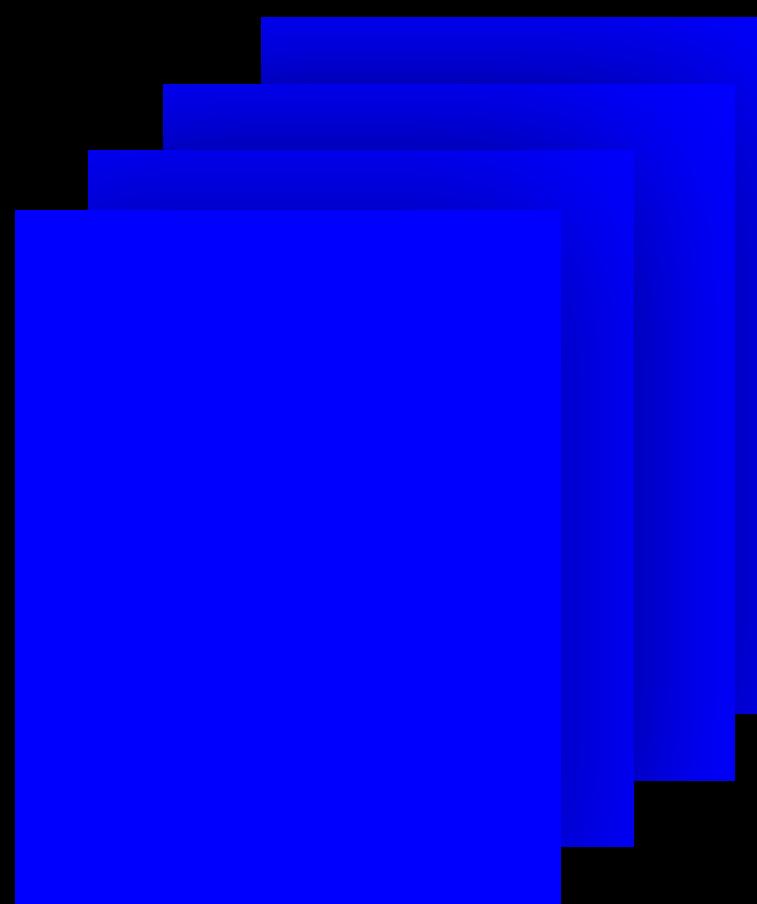


Trained on Edges → Images

Data from [Eitz, Hays, Alexa, 2012]
slide credit: Phillip Isola & Jun-Yan Zhu

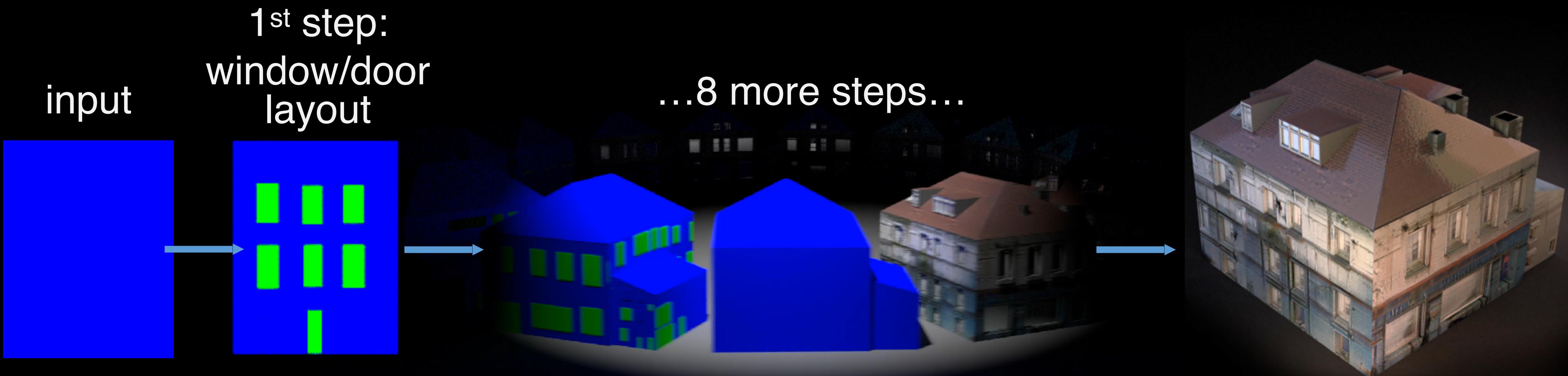
Decomposition into Steps: FrankenGAN

input

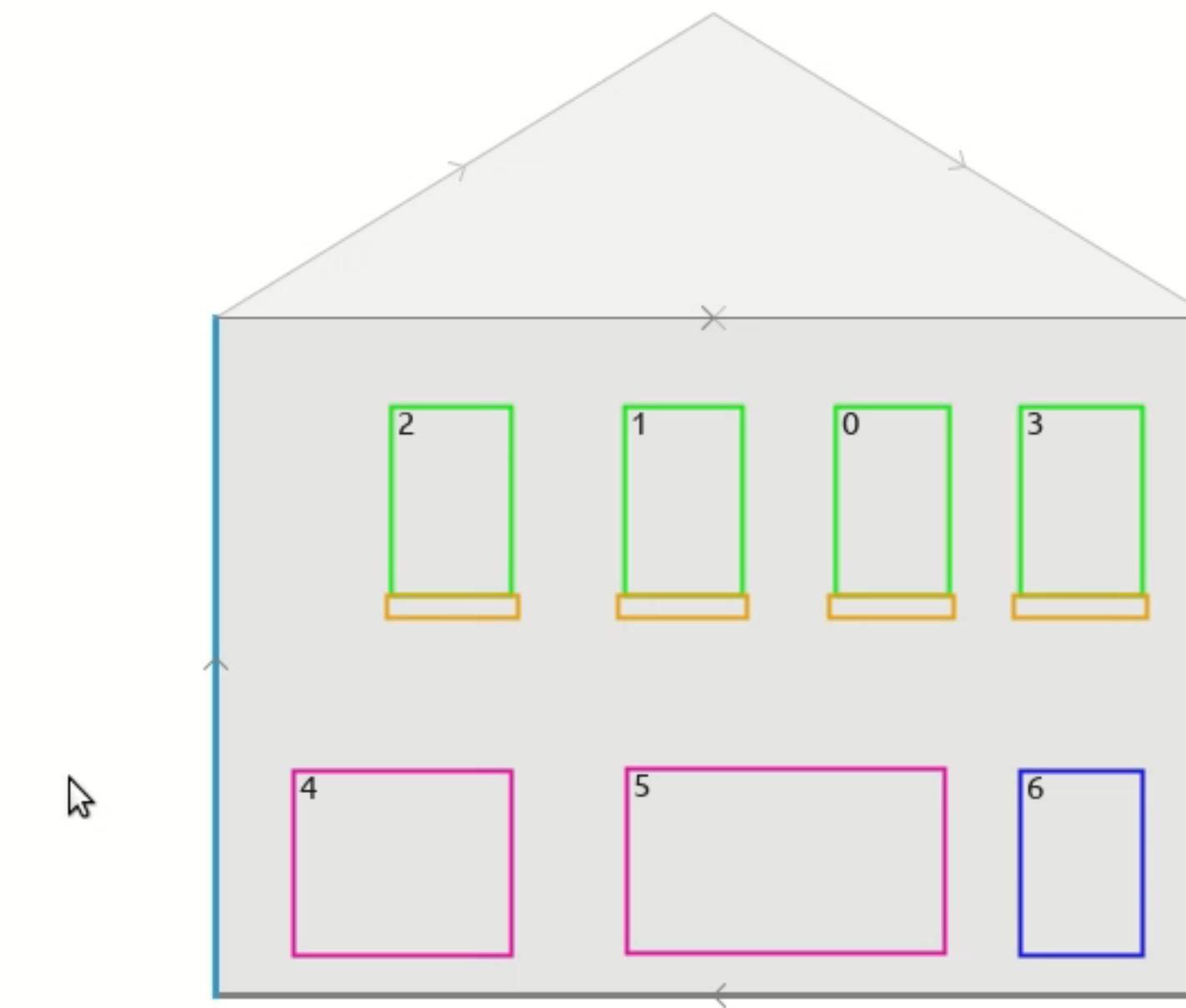
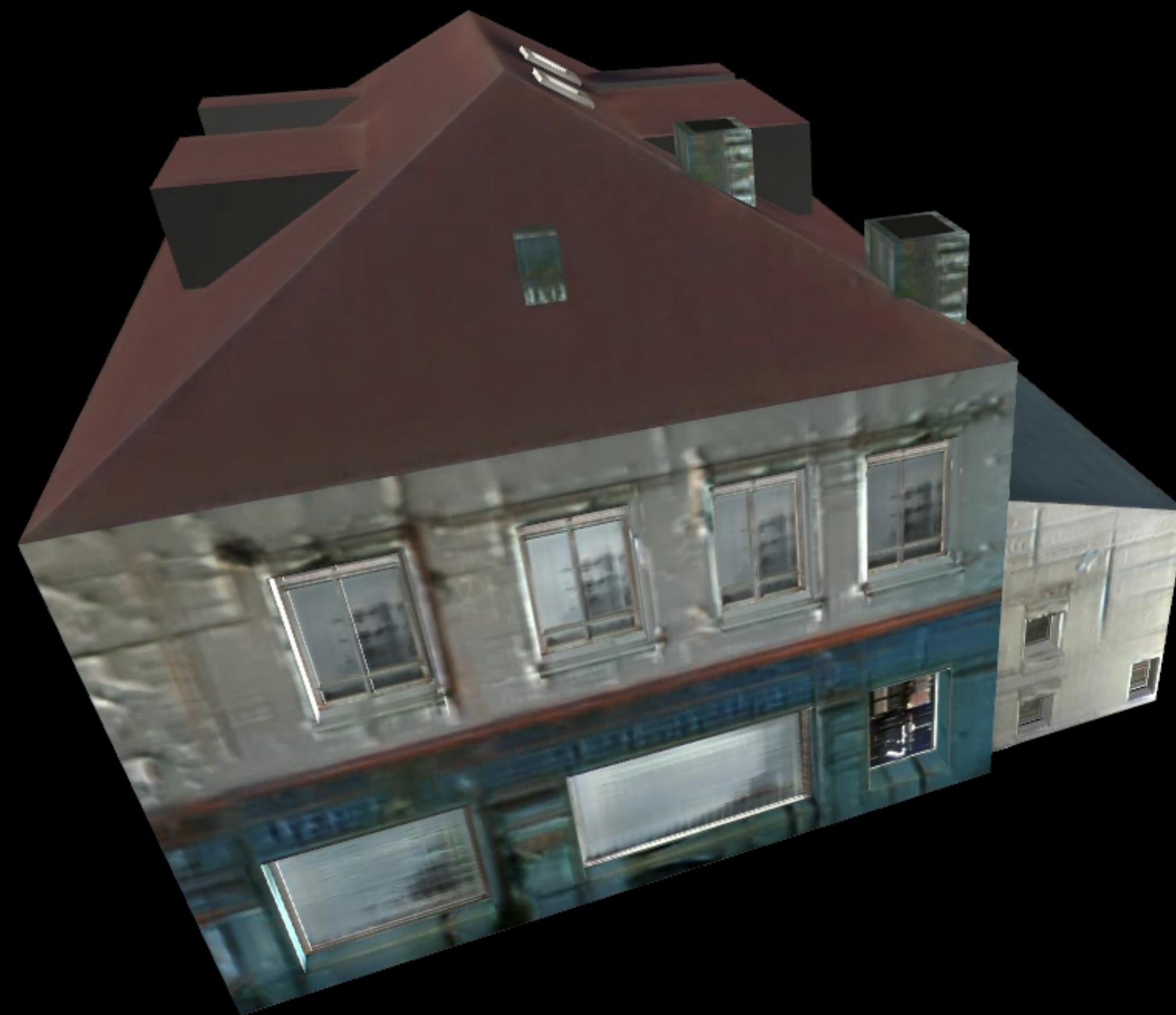


Slide Credit: FrankenGAN, Kelly and Guerrero et al.

Decomposition into Steps: FrankenGAN



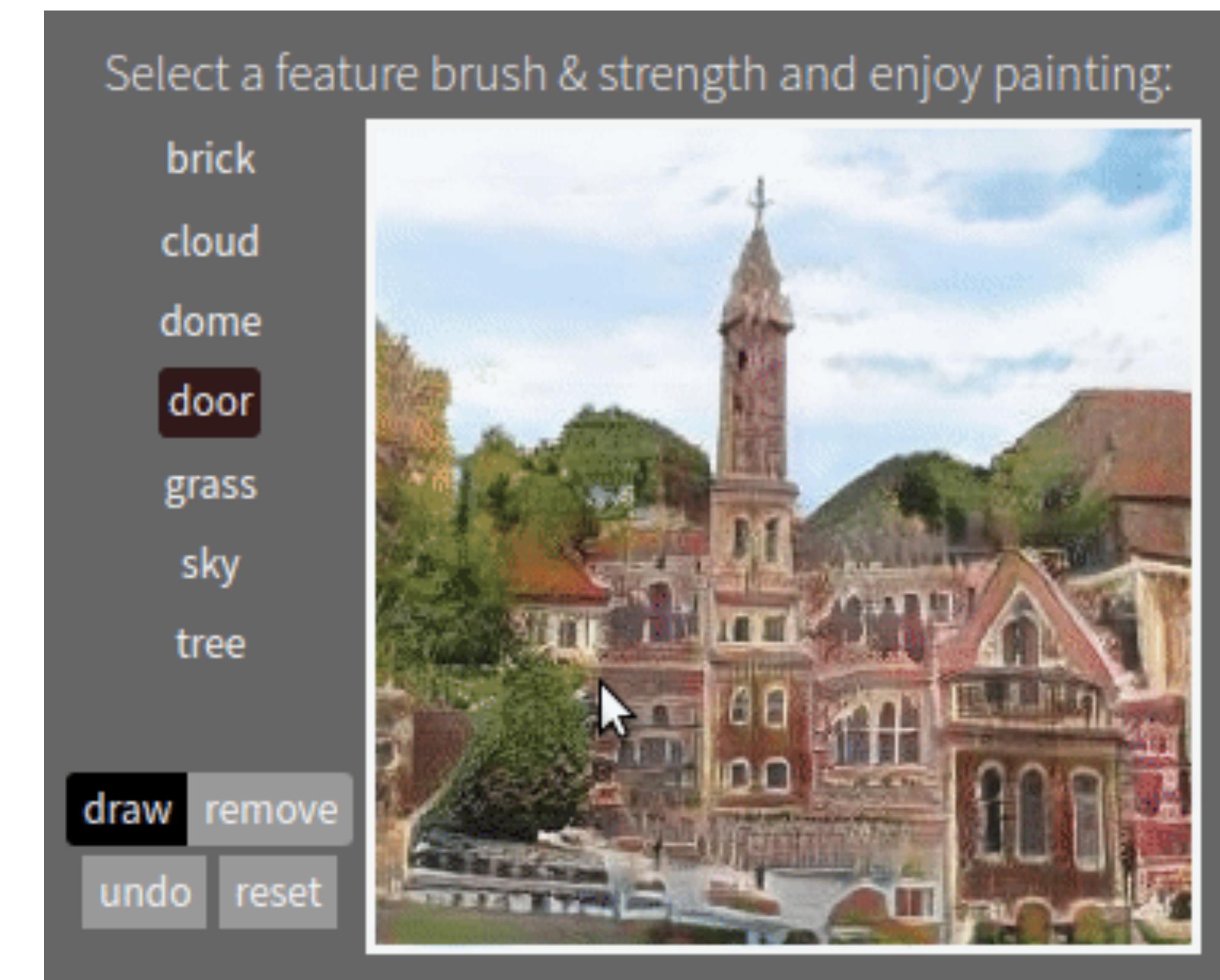
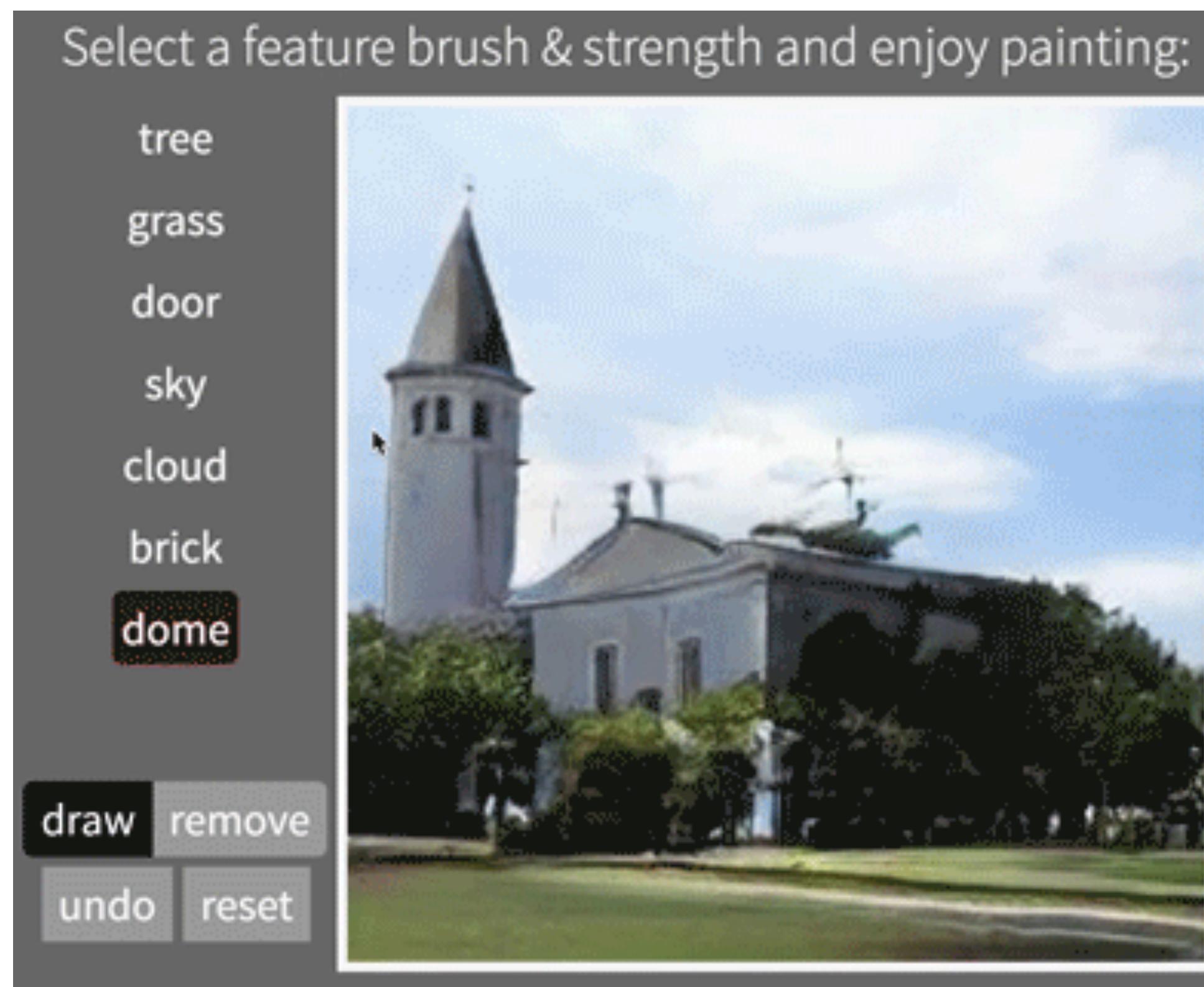
Manual Control



GAN Dissection

Question: How does a GAN create an image? What do individual neurons do?

Insight: Neurons are specialized to create objects of specific types



StructureNet

- **Consistent structure across the dataset**

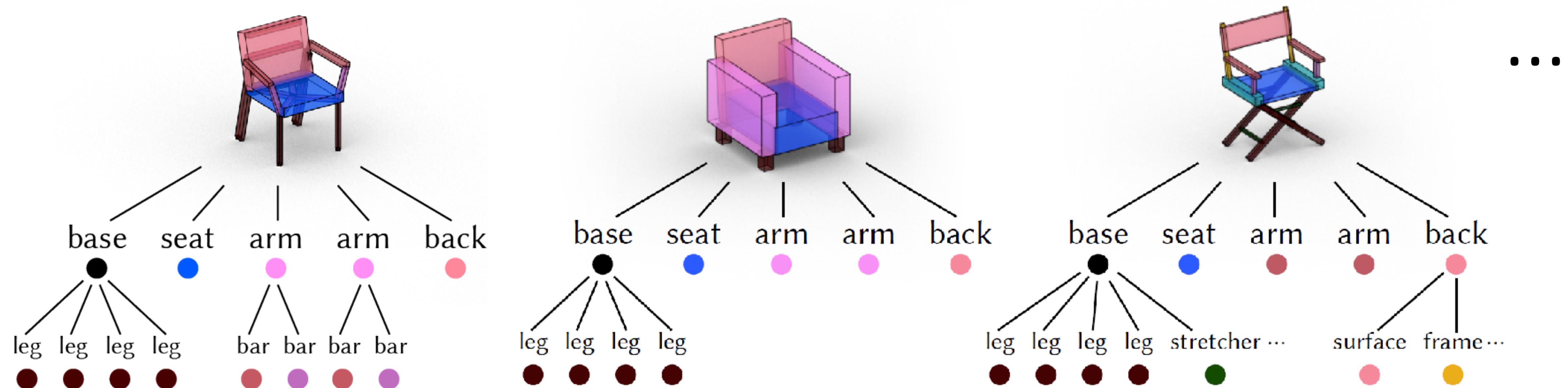
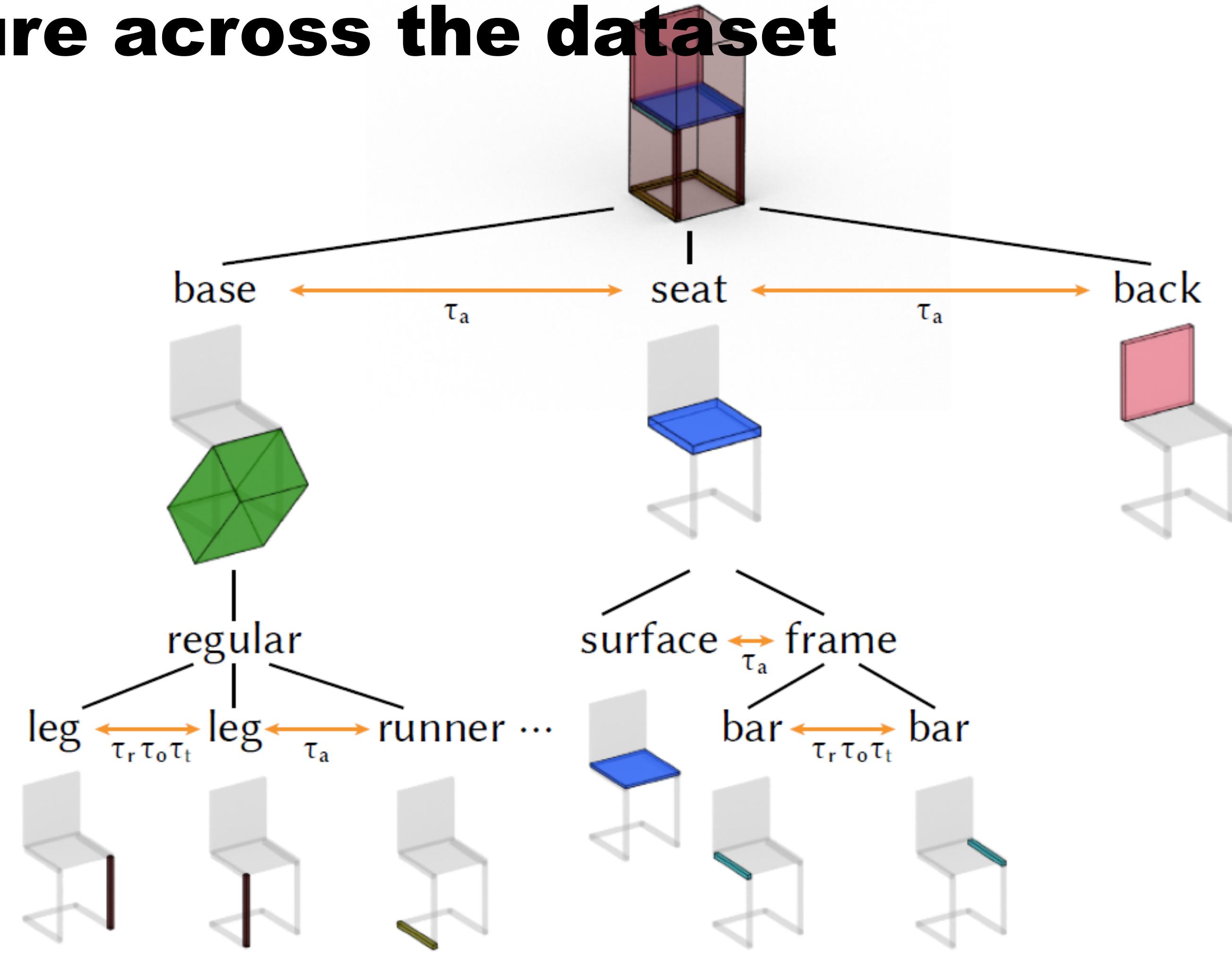


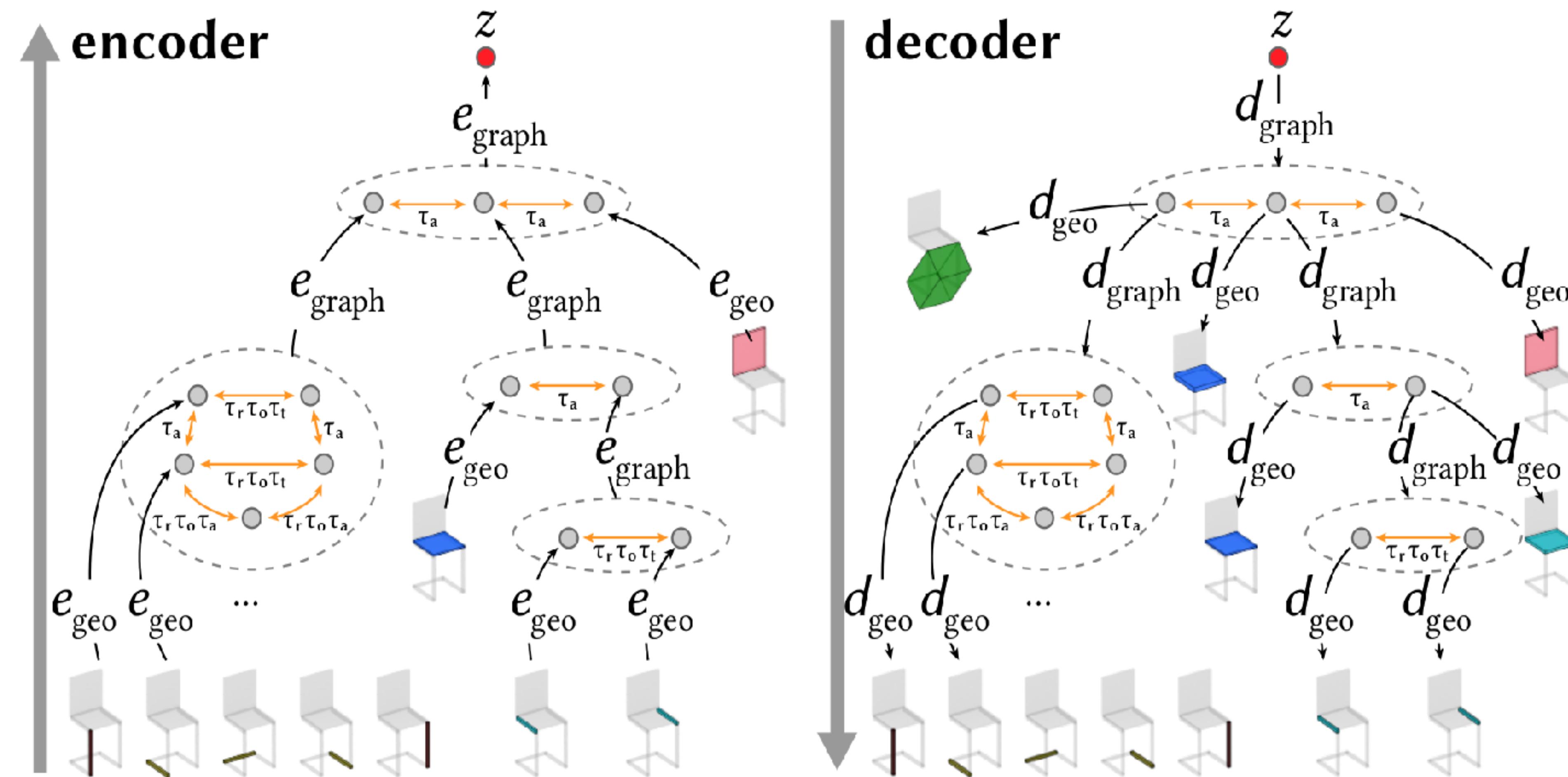
Image Credit: StructureNet: Hierarchical Graph Networks for 3D Shape Generation , Mo and Guerrero et al.

StructureNet

- **Consistent structure across the dataset**
- **Richer structure**

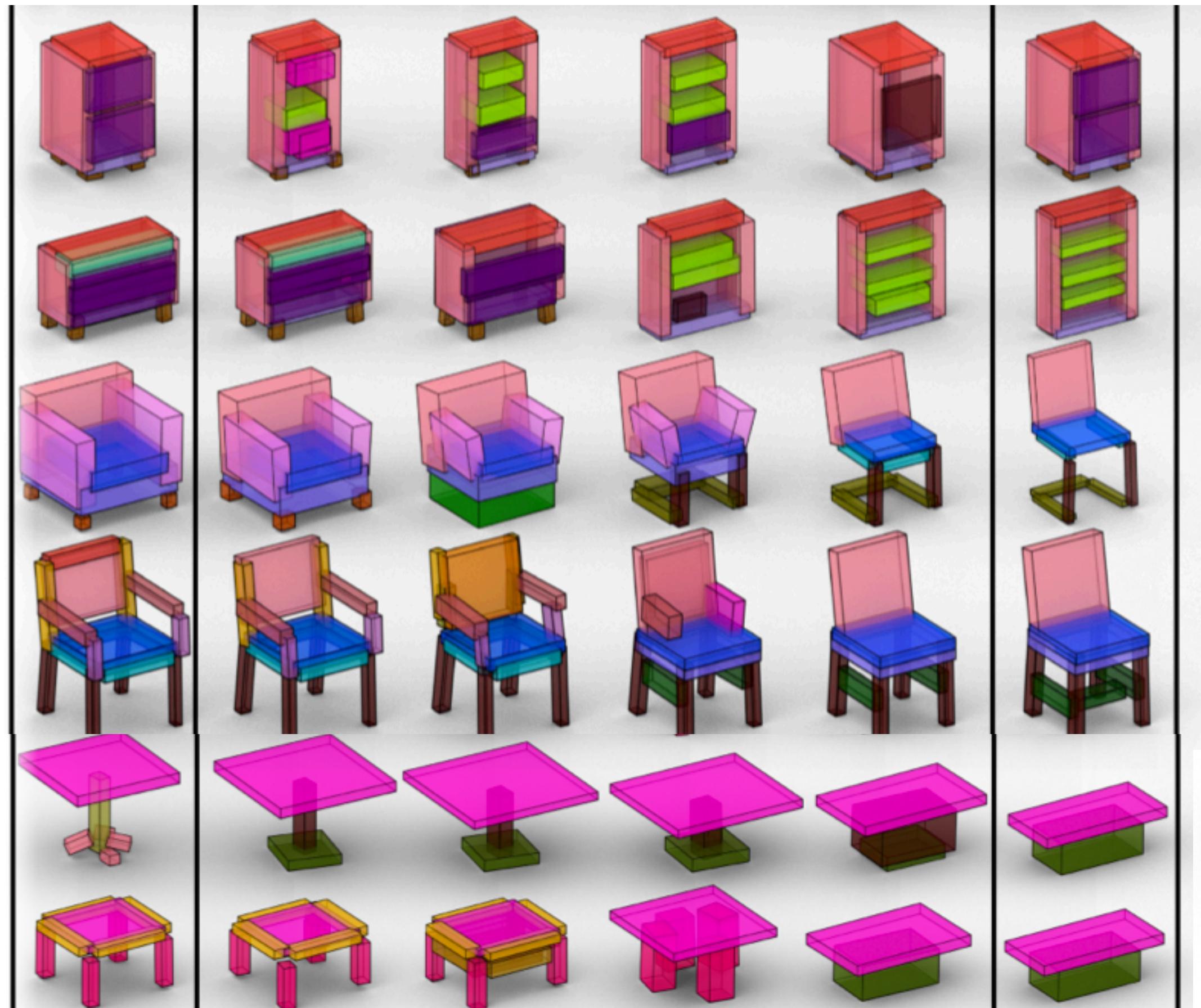


Hierarchical Graph Networks as Encoder and Decoder



StructureNet Results

interpolation



free generation



Image Credit: StructureNet: Hierarchical Graph Networks for 3D Shape Generation , Mo and Guerrero et al.

GAN Training Convergence

GAN training can be unstable

- **Generator and discriminator do now always converge (Nash equilibrium)**
- **Vanishing discriminator gradients**
- **Mode Collapse**

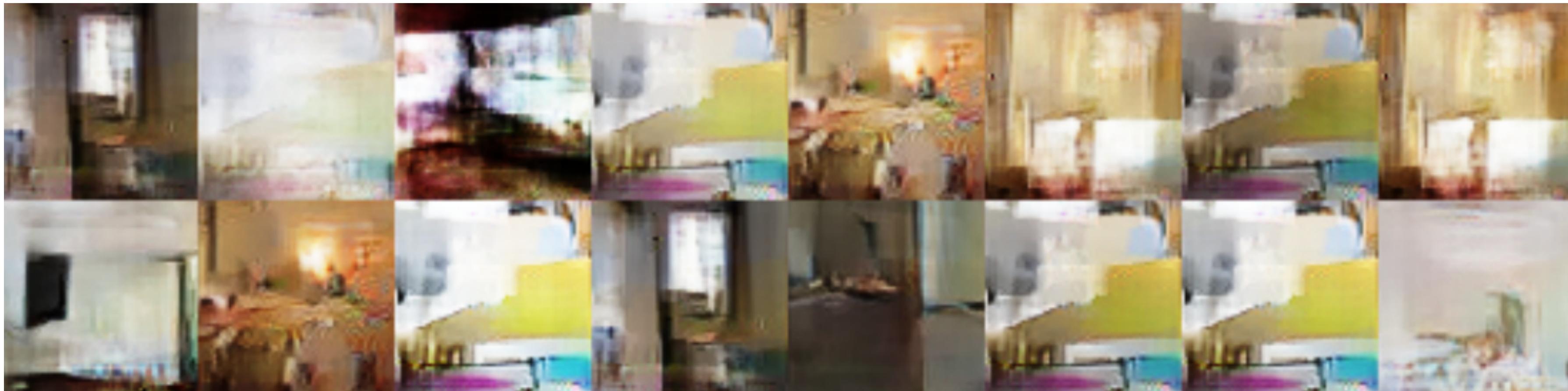
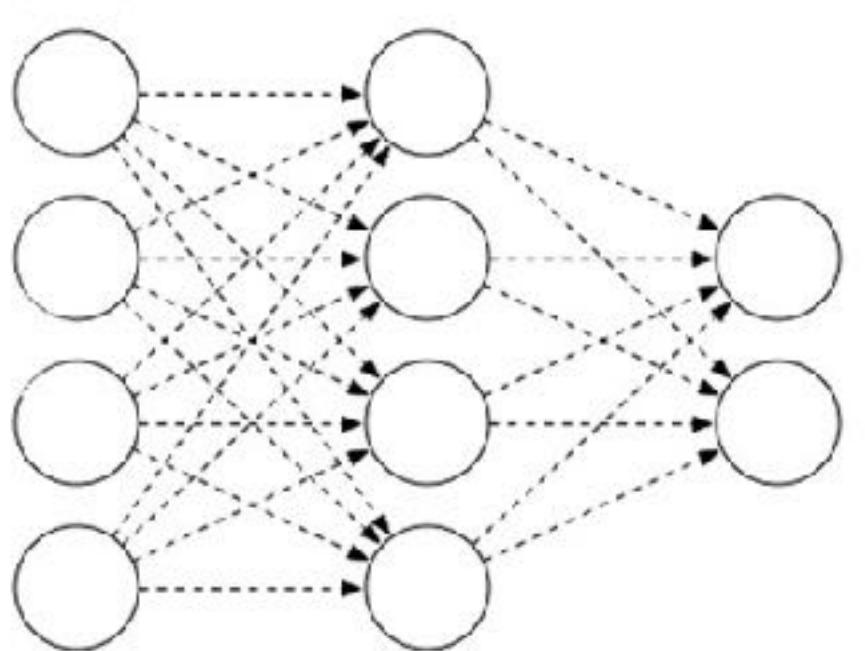


Image Credit: Wasserstein GAN, Arjovsky et al.



<http://geometry.cs.ucl.ac.uk/creativeai/>

