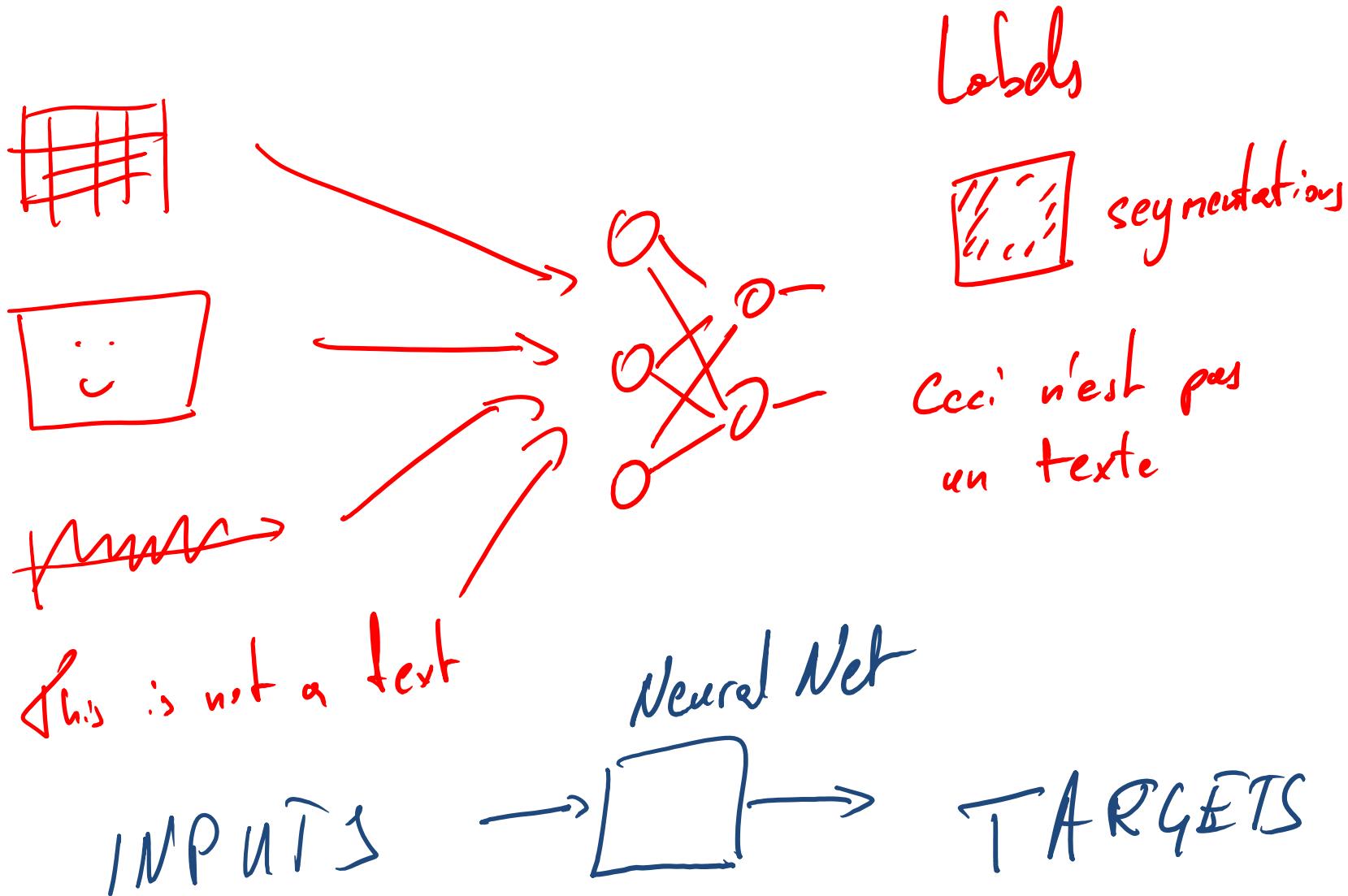


# Deep Unsupervised Learning

Deep Learning @ UWr 2021

Jan Chorowski

# Motivation: Deep Nets excel at supervised tasks



ollede

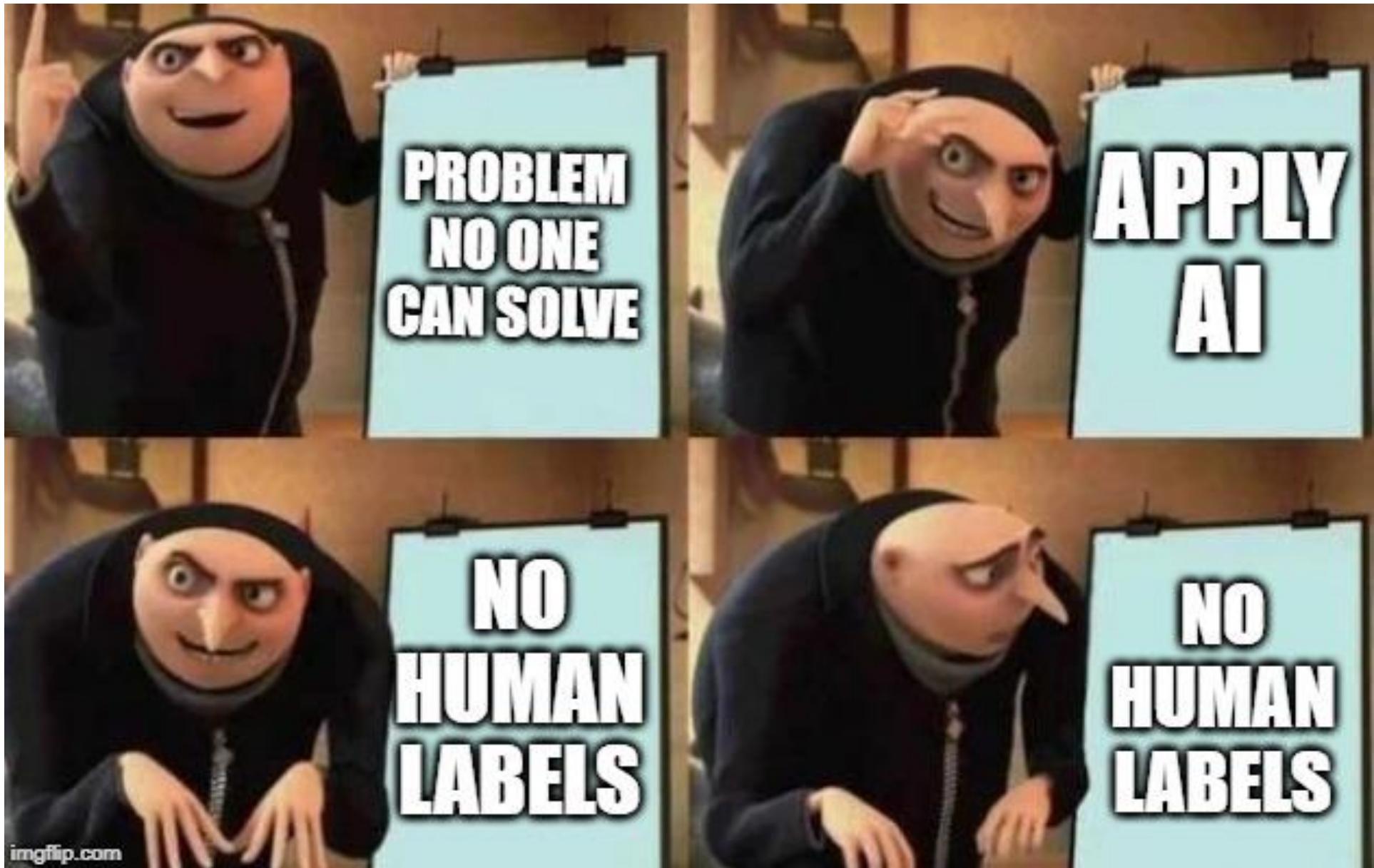
follede qollan 89 qollan oll  
 eed 48 oll de ollas oll lede das  
 das crede qollan qollan oll lede  
 & eed qollan er de das ollas oll  
 lede qollan das qollan qollan & das  
 das qollan eed qollan eed qollan  
 lede er de das qollan eed qollan &  
 lede oll lede ollas lede ollas oll  
 das qollan ollas er de das qollan &  
 das lede oll lede ollas er de das  
 qollan oll lede ollas & qollan  
 er de das das lede qollan qollan &  
 das qollan ollas das & ollas oll  
 & er de oll lede  
 lede oll lede oll  
 das & ollas  
 das qollan &  
 das qollan  
 eed qollan  
 er de das  
 & er de das  
 & er de das  
 oll lede  
 lede qollan &  
 das oll  
 oll lede &  
 das oll  
 & er de das

das oll

de lede



ollede qollan oll



# Unsupervised Learning Goals and Possibilities

## What we may want:

- Learn  $p(x)$

Then use it to:

- Generate new data
- Enhance data (denoise, transform an input sample into a more plausible one)

- Learn a useful representation

Then use it to:

- Train models with little labels
- Further analysis, e.g. clustering



## What losses we can use:

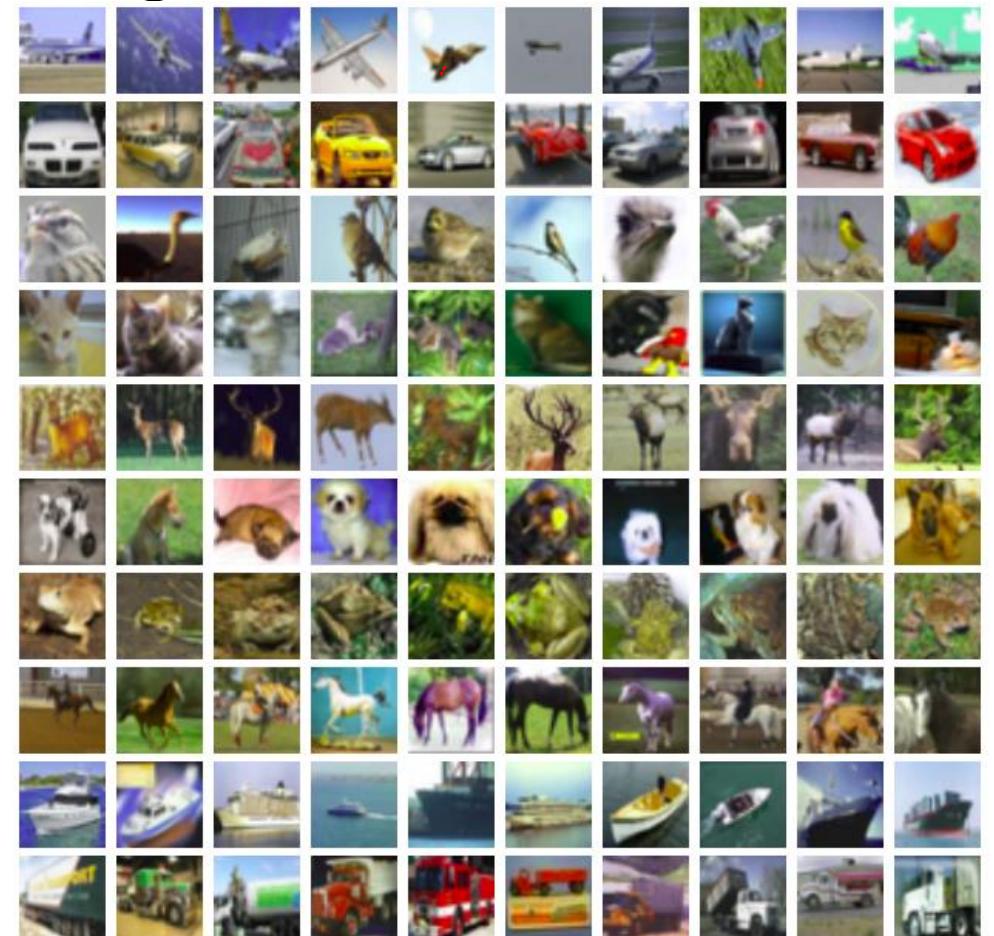
- Log-likelihood  $\log p(x; \theta)$
- Learn to produce plausibly looking samples (as judged by an auxiliary network)  $GAN$
- Learn to encode  $x \rightarrow z$  and decode  $z \rightarrow x'$  keeping  $x'$  close to  $x$
- Learn to solve a proxy task, for which labels are easy to get

# Data generation

Learning high dimensional prob. distributions is hard

- Assume we work with small (32x32) images
- Each data point is a real vector of size  $32 \times 32 \times 3$
- Data occupies only a tiny fraction of  $\mathbb{R}^{32 \times 32 \times 3}$
- Difficult to learn!

~~10<sup>300,000</sup>~~



QQ: what generative models we have  
already seen?

Autoregressive Models - Divide & Conquer

$$P(x_1, y_2, \dots, x_n) = \underbrace{p(x_1)}_{\text{f}} \underbrace{p(x_2 | x_1)}_{\text{g}} \dots \underbrace{p(x_n | x_1, \dots, x_{n-1})}_{\text{f}}$$

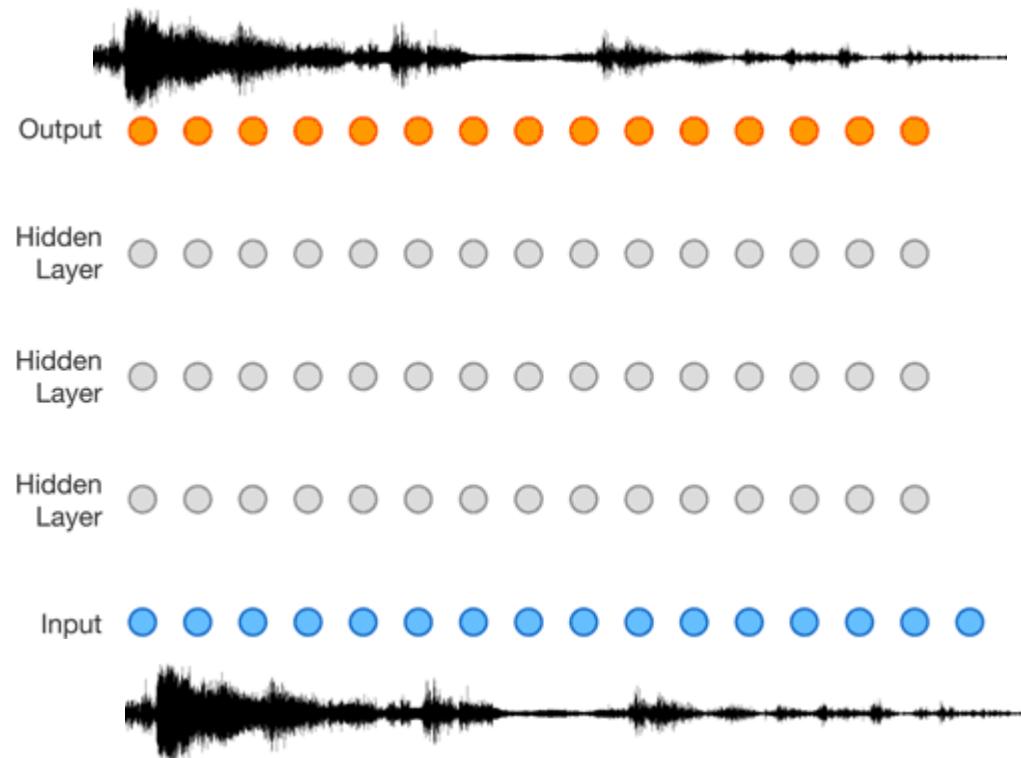
# **AUTOREGRESSIVE MODELS**

# Density Estimation Using Autoregressive Models: Language modeling

- Let  $x$  be a sequence of word ids.
- $p(x) = p(x_1, x_2, \dots, x_n) = \prod_i p(x_i | x_{<i})$ 
$$\approx \prod_i p(x_i | x_{i-k}, x_{i-k+1}, \dots, x_{i-1})$$
- $p(\text{It's a nice day}) = p(\text{It}) * p(\text{'s} | \text{it}) * p(\text{a} | \text{'s}) \dots$
- Classical n-gram models: cond. probs. estimated using counting
- Neural models: cond. probs. estimated using neural nets

# WaveNet: Autoregressive modeling of speech

- Treat speech as a sequence of samples!
- Predict each sample base on previous ones.



# PixelCNN:

## A “language model for images”

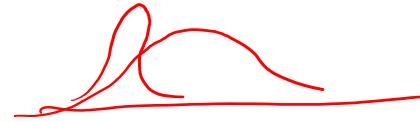


Pixels generated left-to-right,  
top-to-bottom.

$p(x_i|x_{<i})$  computed using  
convolutional neural nets

**Model supports auxiliary  
conditioning:  $p(x_i|x_{<i}, c)$**

# Modeling pixels



How to model pixel values:

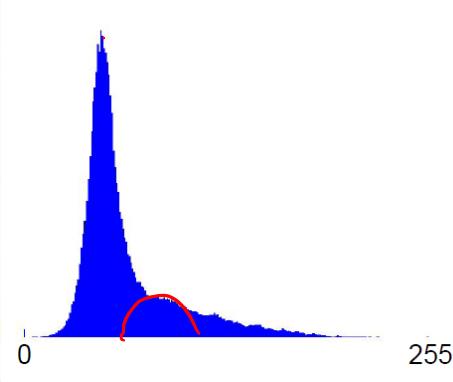
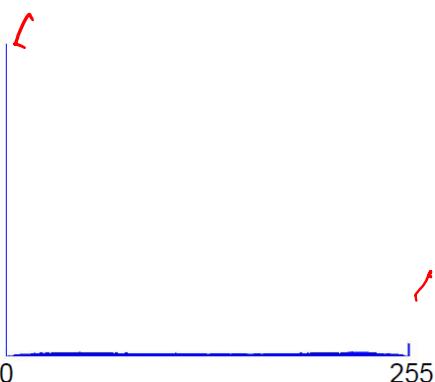
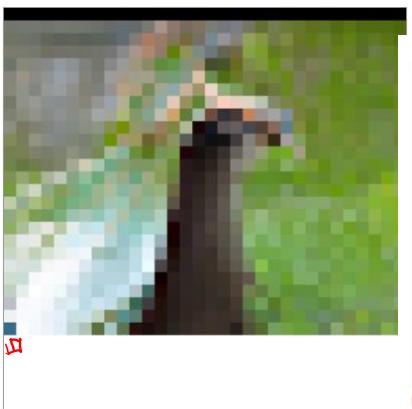
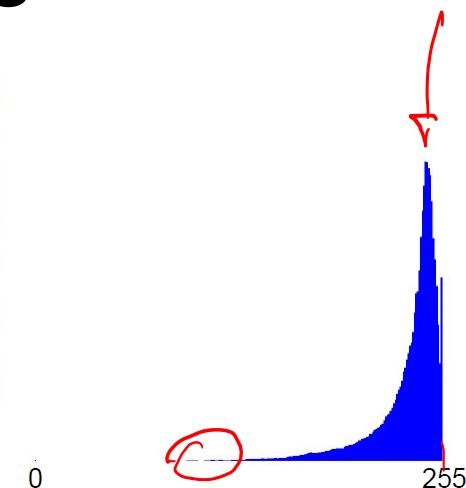
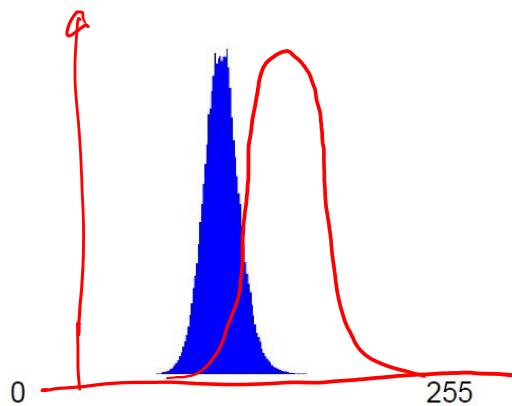
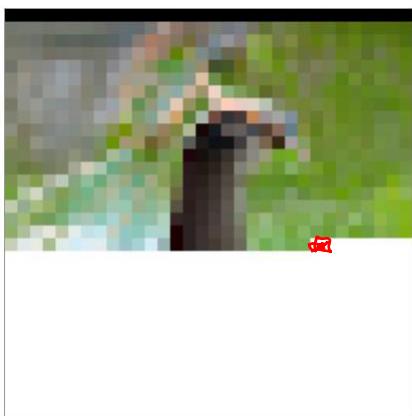
- A Gaussian with fixed st. dev?
- A Gaussian with tunable st. dev?
- A distribution over discrete levels [0,1,2,...255]?

$$p_{\text{pix}}(x_i | x_{c_i})$$

pix :   
 all pixels to  
 the left end  
 b doc

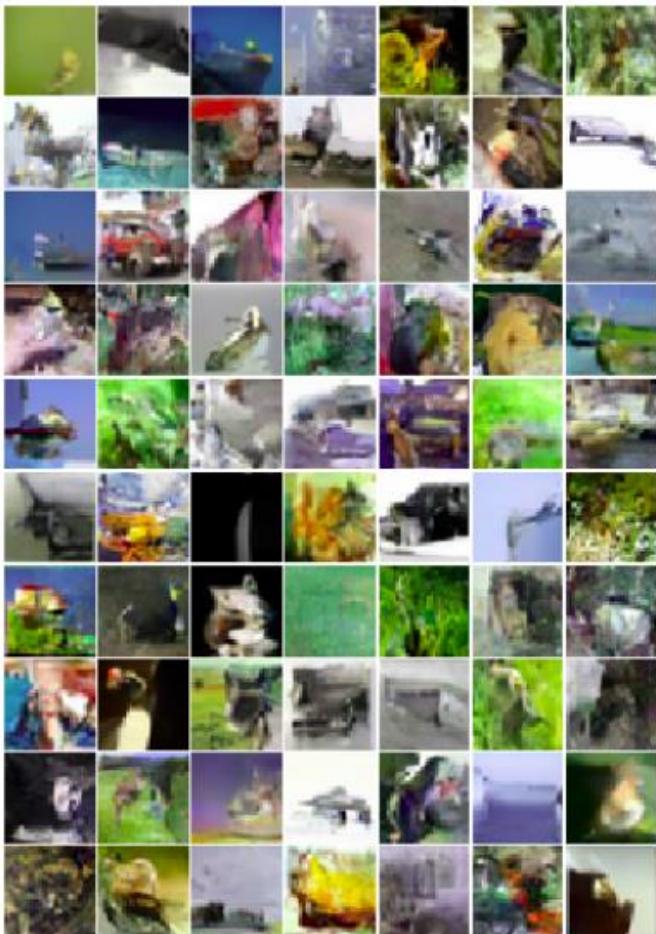
What are the implications?

# Modeling pixel values



Model works best with a flexible distribution: better to use a SoftMax over pixel values! 

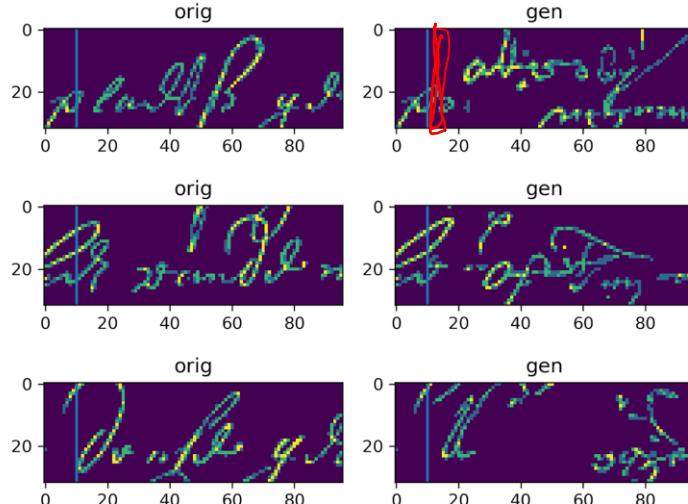
# PixelCNN samples & completions



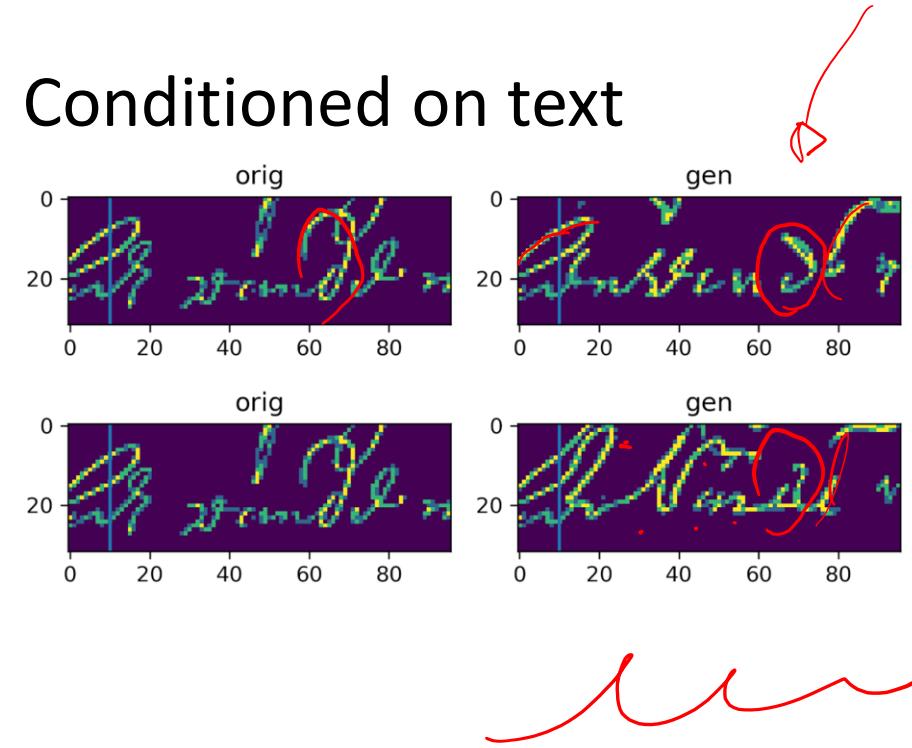
Salimans et al, “A PixelCNN Implementation with Discretized Logistic Mixture Likelihood and Other Modifications”  
van den Oord, A., et al. “Pixel Recurrent Neural Networks.” ICML (2016).

# PixelCNN samples

Unconditional



Conditioned on text



# Autoregressive Models Summary

The good:

- Easy to exploit correlations in data. ↗
- Reduce data generation to many small decisions
  - Simple define (just pick an ordering)
  - Trains like fully supervised
  - Model operations are deterministic, randomness needed during generation
- Often State-of-the-Art log-likelihood ↗

# Autoregressive Model Summary

The bad:

- Train/test mismatch (teacher forcing):  
trained on ground truth sequences  
but applied to own predictions
- Generation requires  $O(n)$  steps  
(Training can be sometimes parallelized)
- No compact intermediate data representation,  
not obvious how to use for downstream tasks.

# LATENT VARIABLE MODELS

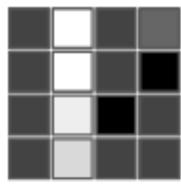
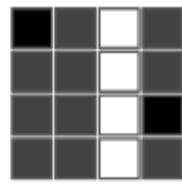
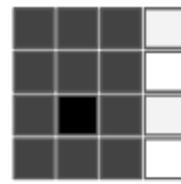
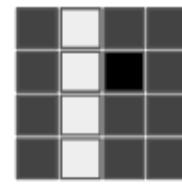
# Ad break

I got the many of following slides from Ulrich Paquet  
(DeepMind)

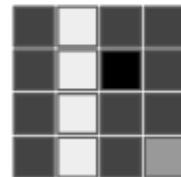
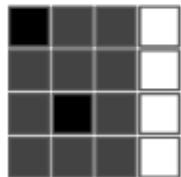
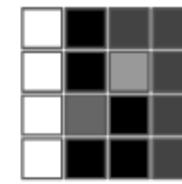
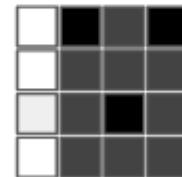
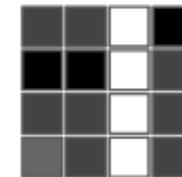
See <https://www.youtube.com/watch?v=xTsnNcctvmU> for a recording of his VAE explanation!

# 2 minute exercise

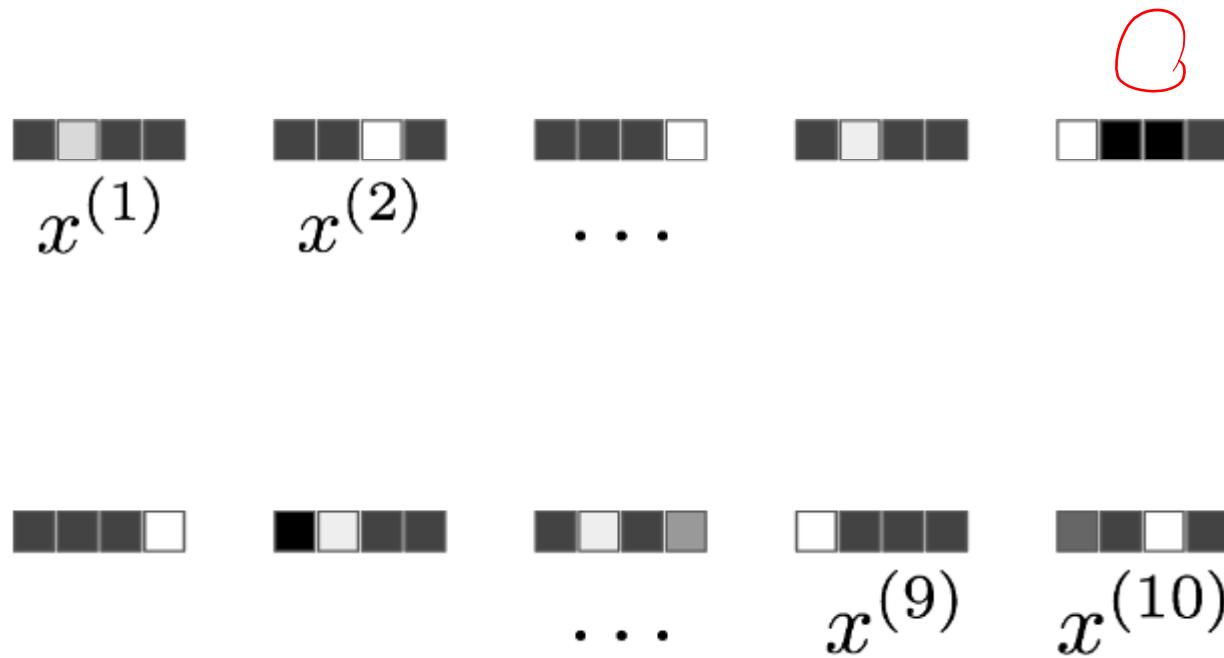
Think about all the properties of this dataset:

 $x^{(1)}$  $x^{(2)}$  $\dots$ 

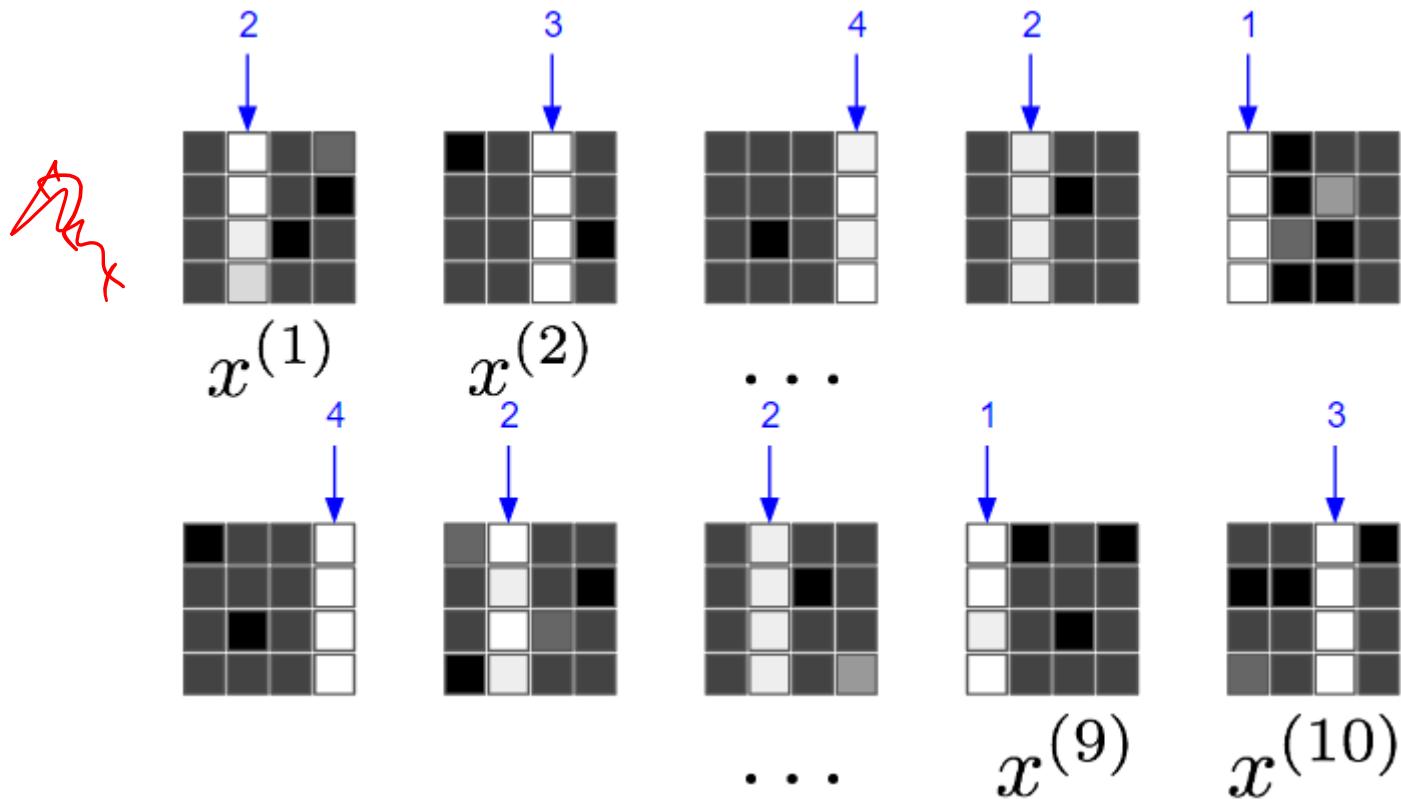
$\mathbb{R}^{6 \times 4}$

 $\dots$  $x^{(9)}$  $x^{(10)}$

# The rows are correlated



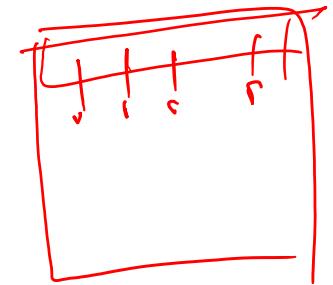
There's a simple encoding  
(the LATENT representation)



# Data structure

We can capture most of the variability in the data through **one** number

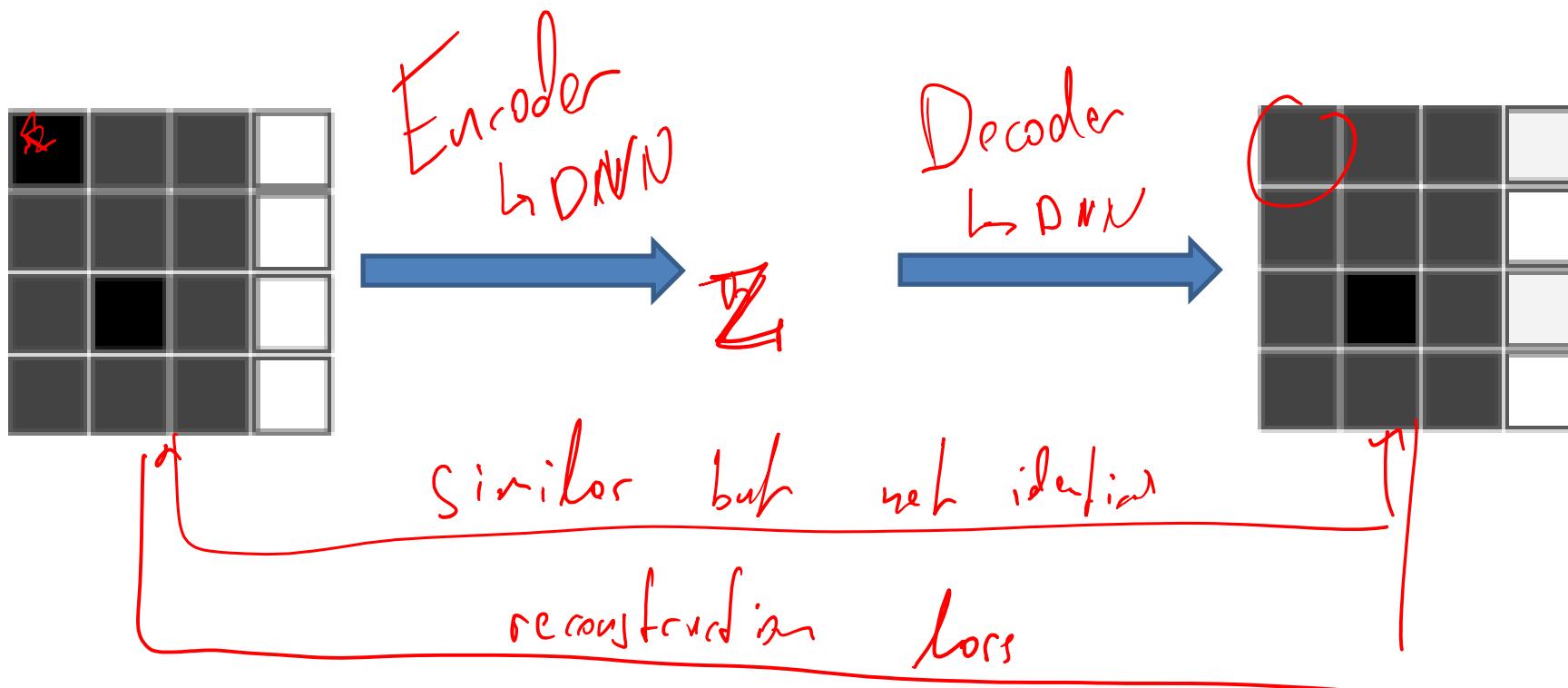
$$z^{(n)} = 1 \text{ or } 2, 3, 4$$



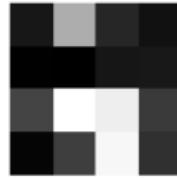
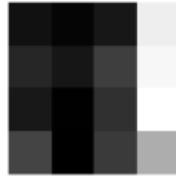
for each image  $n$ , even though each image is 16 dimensional

How to train a model to uncover it?

# Autoencoders

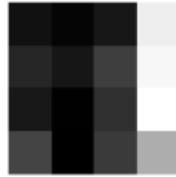


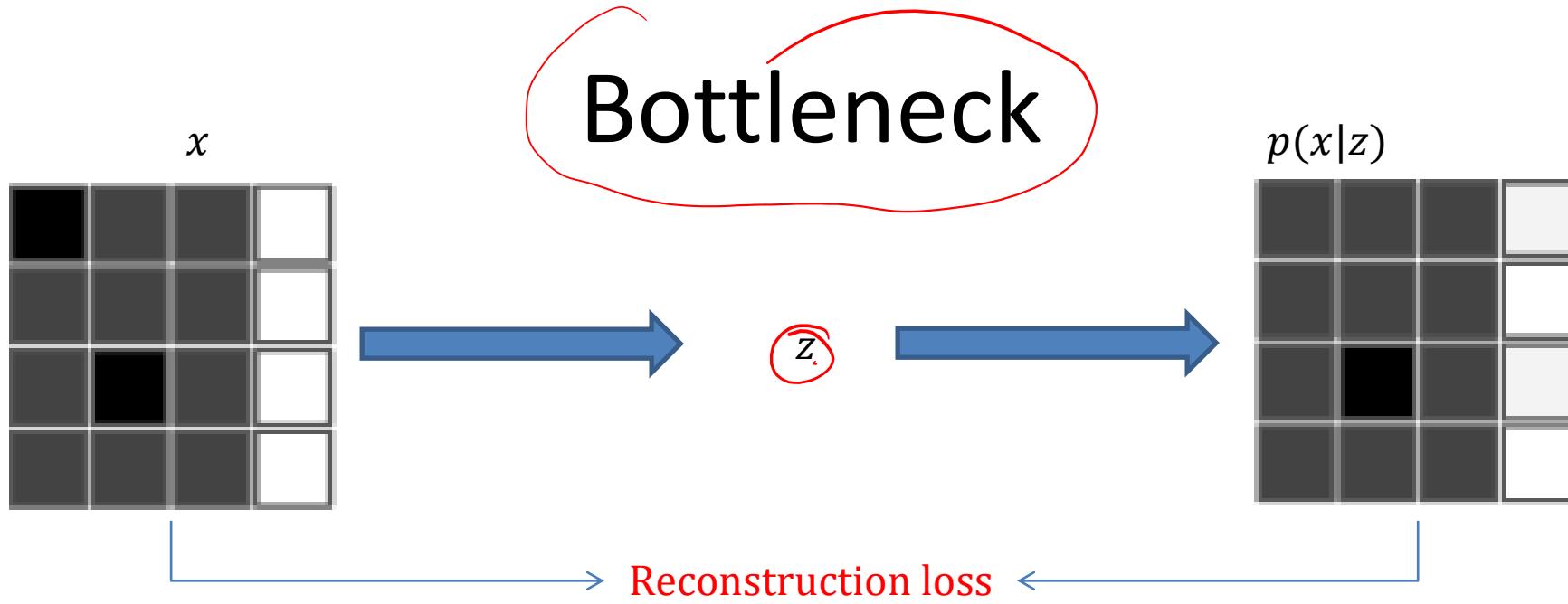
# Perfect Reconstruction === Useful Representations?



all info is  
there

but more  
difficult to  
work with





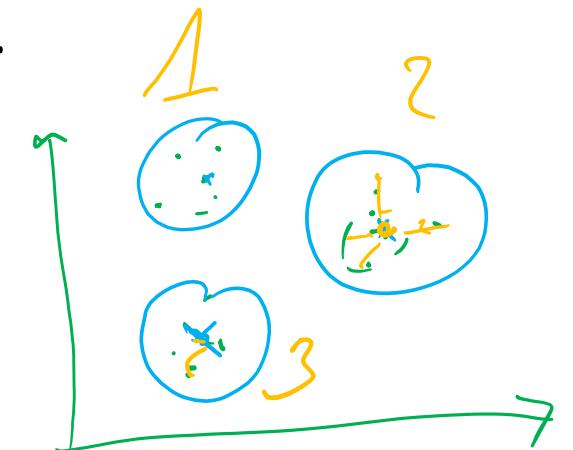
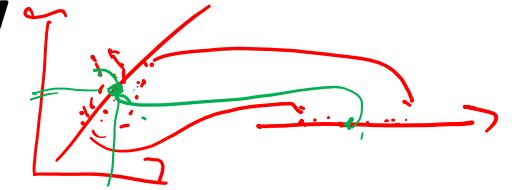
$p(x|z)$  shaped by  $z$ :

- dim reduction

-

# Auto-encoders that we know

- PCA and SVD find low-dimensional projections of data. They can be seen as a lossy autoencoder trained to minimize L2 reconstruction error under a dimensionality reduction bottleneck.
- K-Means assigns each data point to a cluster prototype. This is a lossy encoding. We “decode” a cluster ID into the location of the center. We can interpret K-Means as a lossy encoding scheme under a sparsity constraint (sample is assigned to only 1 cluster)
- Many other variants, e.g. the sparse autoencoder (examples in 13-autoencoders.ipynb)



61 6000 20

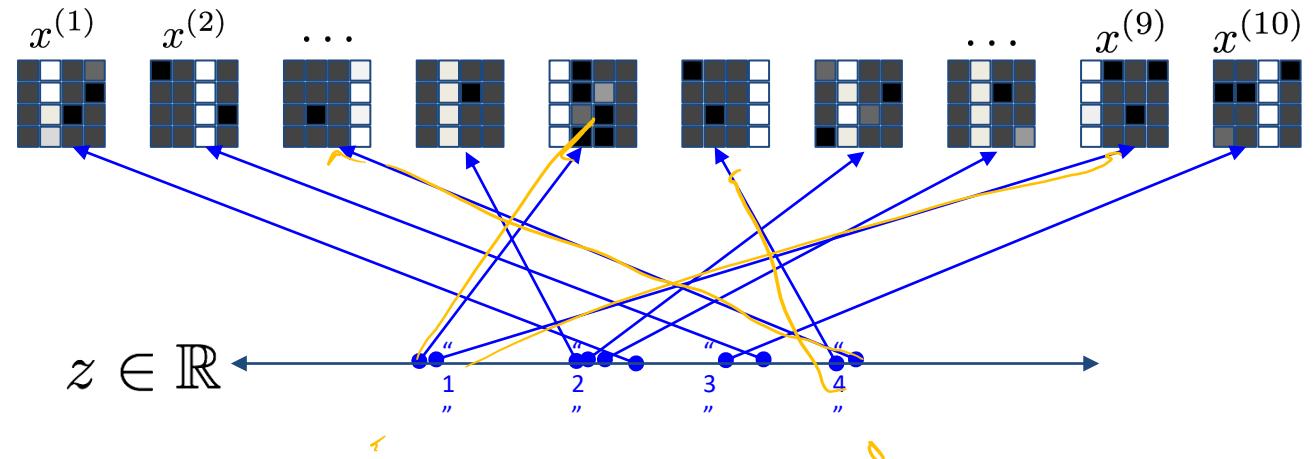
# Simple Auto-encoders summary

- + compact data representations
- +/- careful bottleneck design required to get desired information in the latent vector
- no generation capability

# **VARIATIONAL AUTOENCODER: A PRINCIPLED APPROACH TO AUTOENCODING AND DATA GENERATION**

# Intuition

- Complex data
- May have a simple latent structure



- Build a 2-stage generation process:

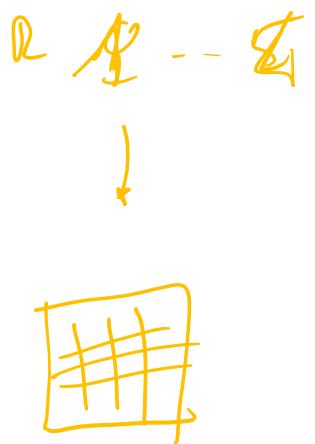
$z \sim \mathcal{N}(0,1)$

prior  $p(z)$  assumed to be simple

$x \sim p(x|z)$

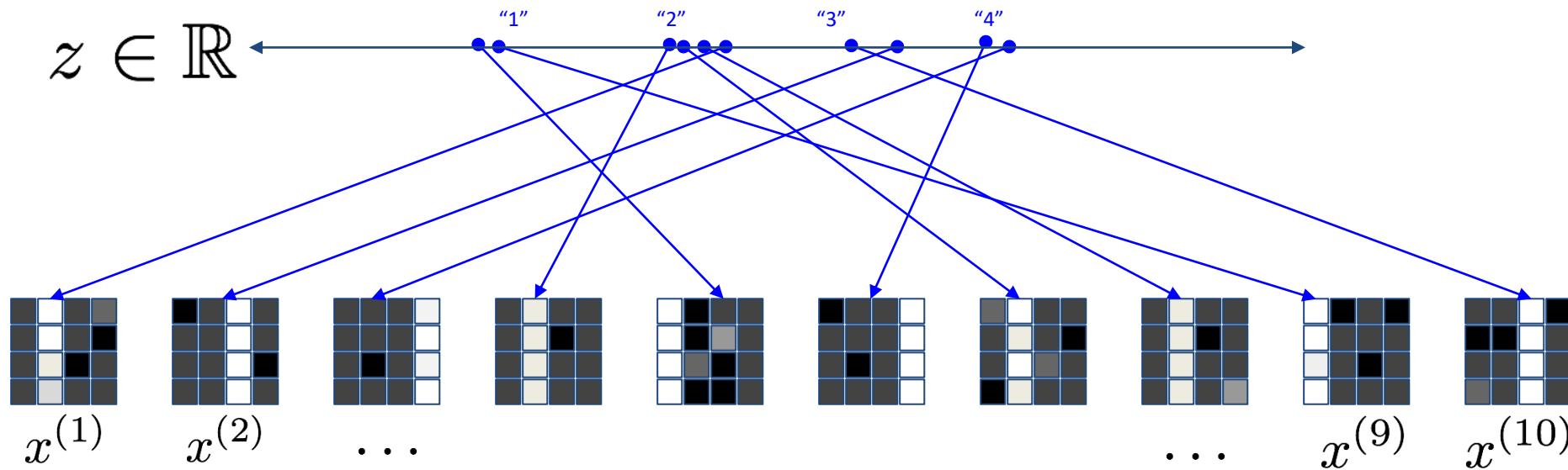
complicated transformation

implemented with a neural network



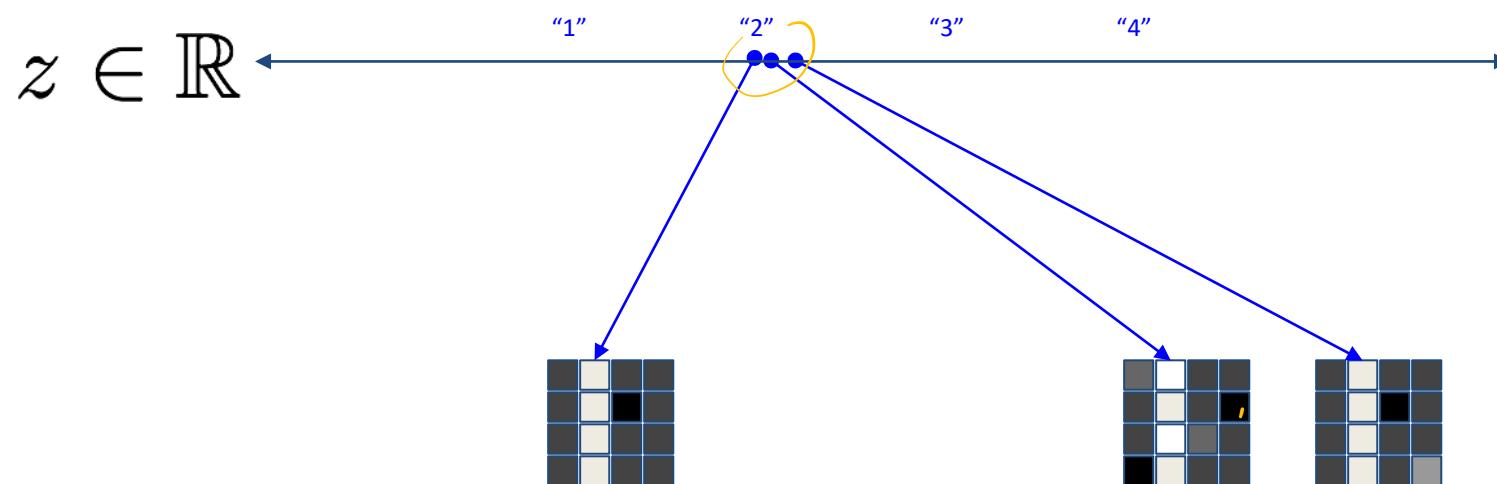
# Data manifold

The 16-dimensional images live on a 1-dimensional manifold, plus some “noise”



# and noise

The 16-dimensional images live on a 1-dimensional manifold, plus some “noise”



# Exercise

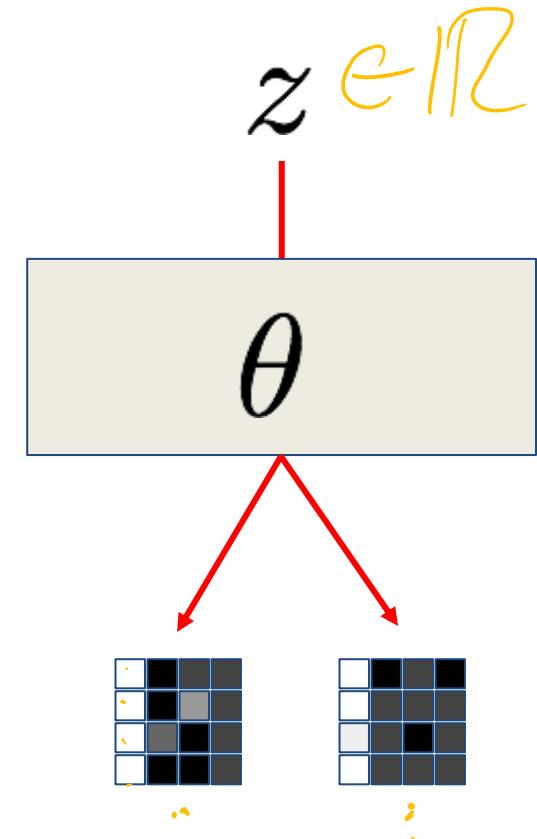
Think how to implement a neural network that takes  $z$  and produces a distribution over  $x$ :

$$p_{\Theta}(x|z)$$

Is your input one-dimensional?

Is your output 16-dimensional?

Identify all the “tunable” parameters  $\Theta$  of your function

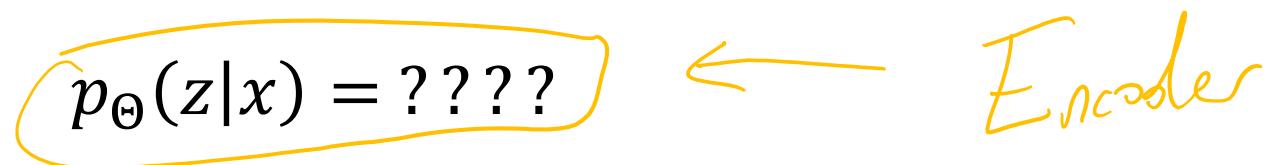


# Inference

Generation:



Inference:

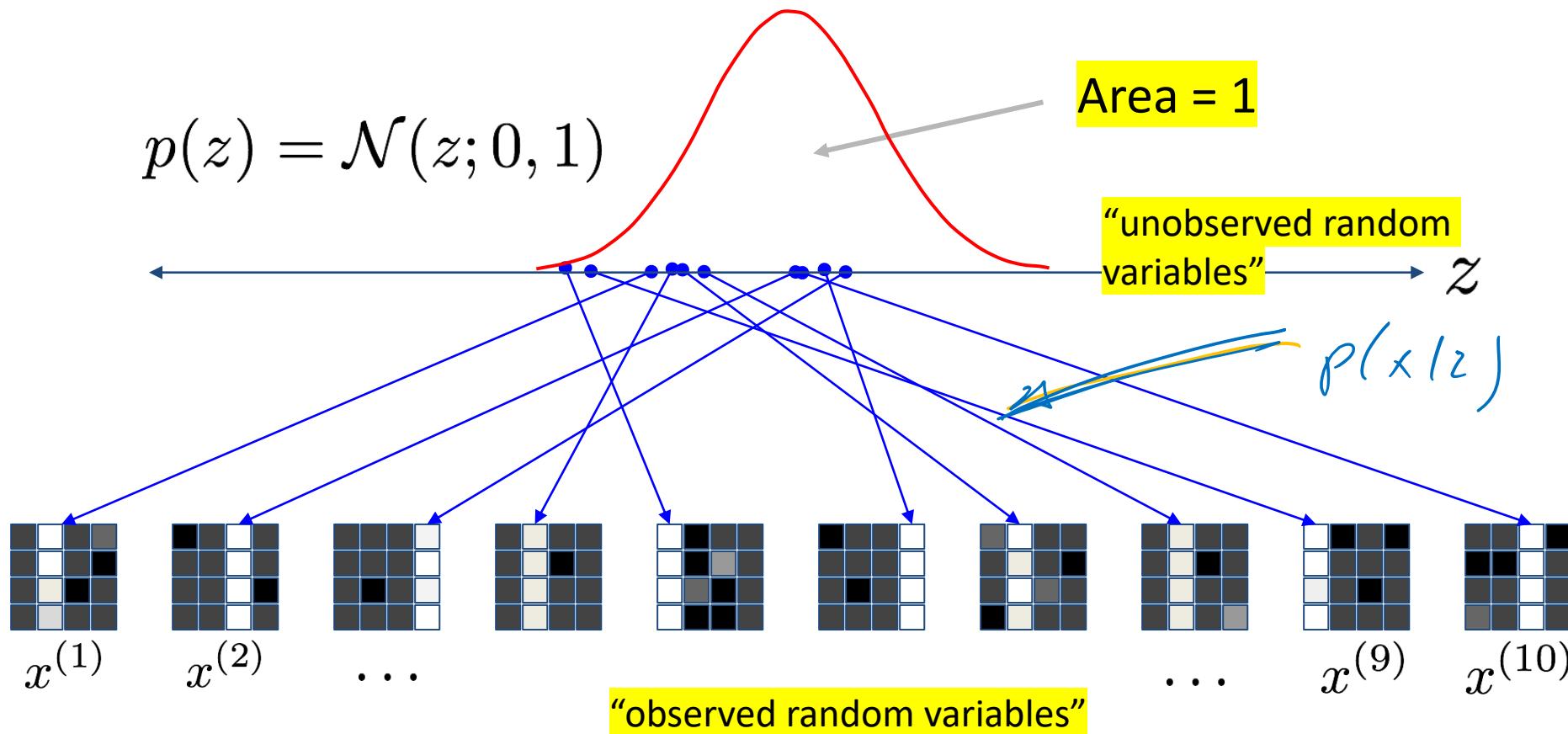


Bayes says:

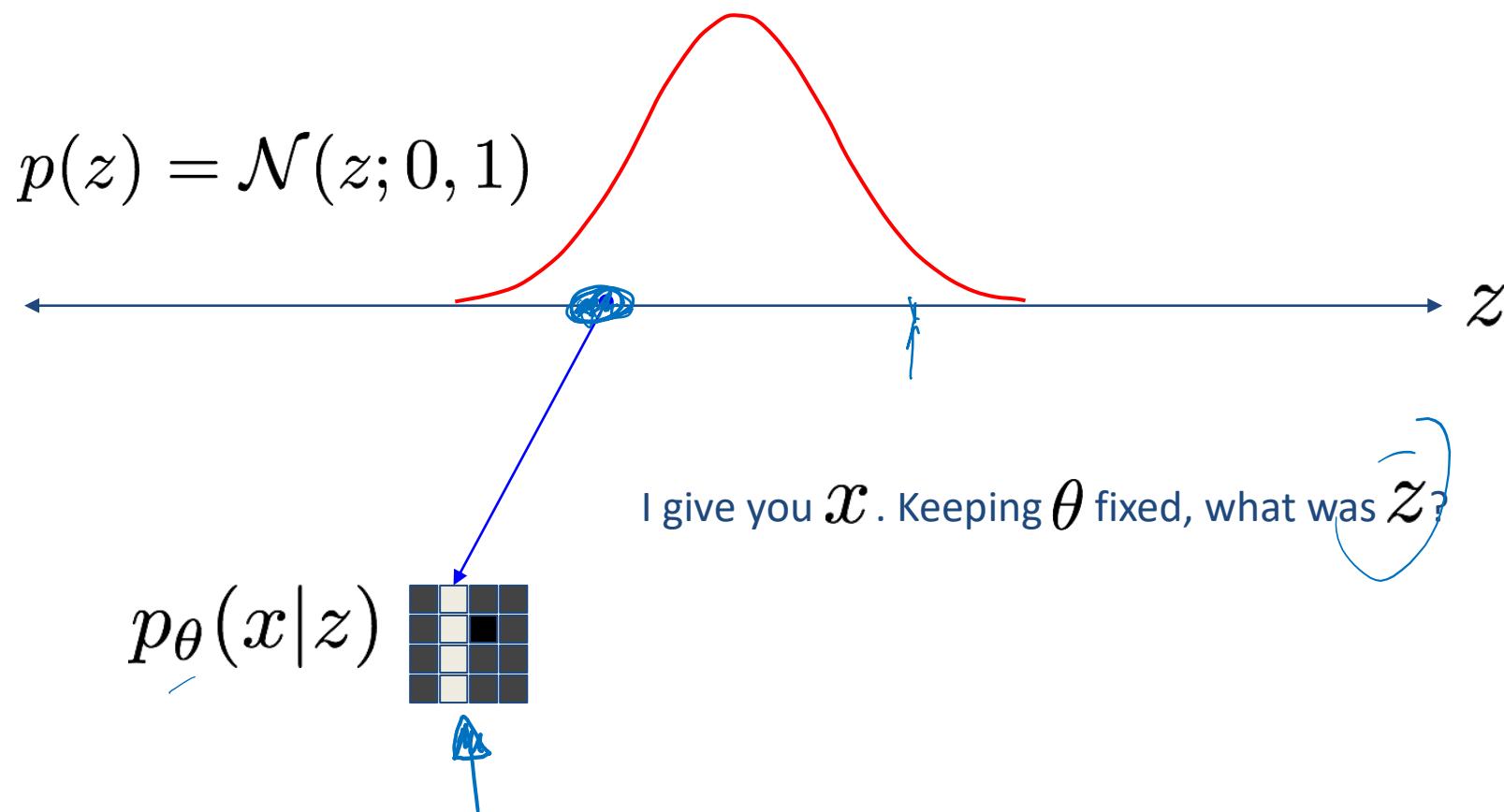
$$p_{\Theta}(z|x) = \frac{p_{\Theta}(x|z)p(z)}{\int dz' p_{\Theta}(x|z')p(z')}$$

Below the equation, a yellow arrow points from the term  $\int dz' p_{\Theta}(x|z')p(z')$  to the expression  $\int d\gamma' p_{\Theta}(x, z')p(z')$ . This expression is further simplified to  $\rho(x)$  below it.

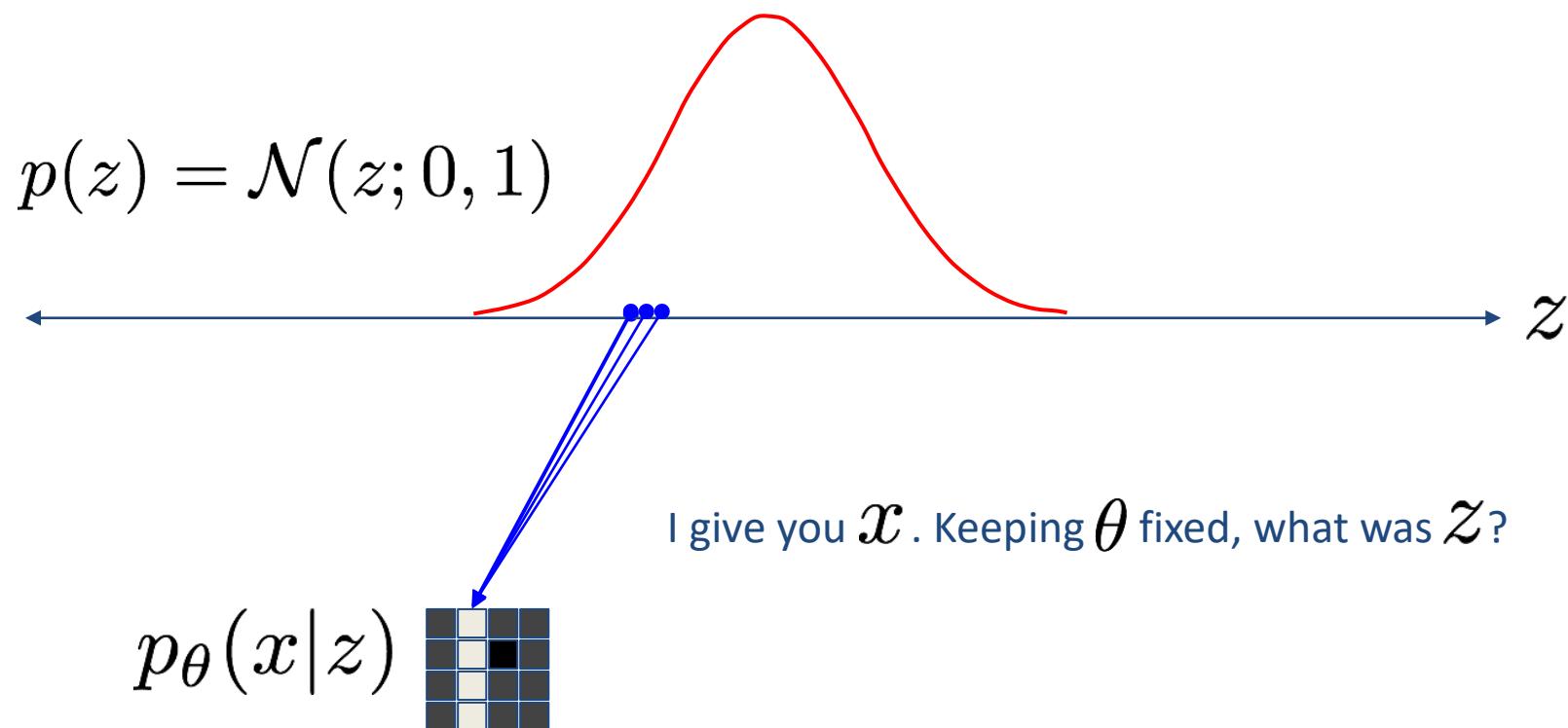
# Inference starts with priors



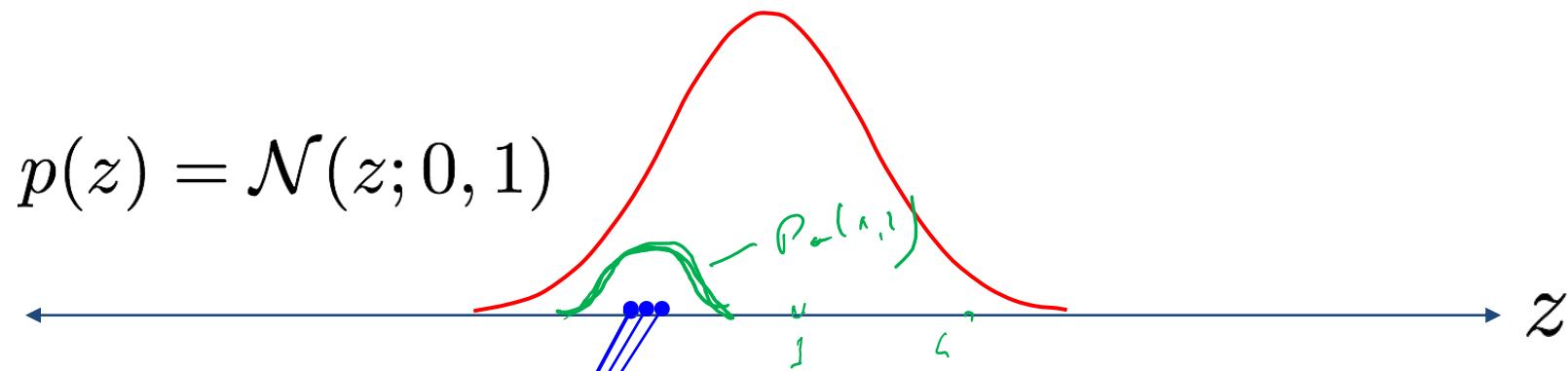
# Inference



# Inference



# Exercise



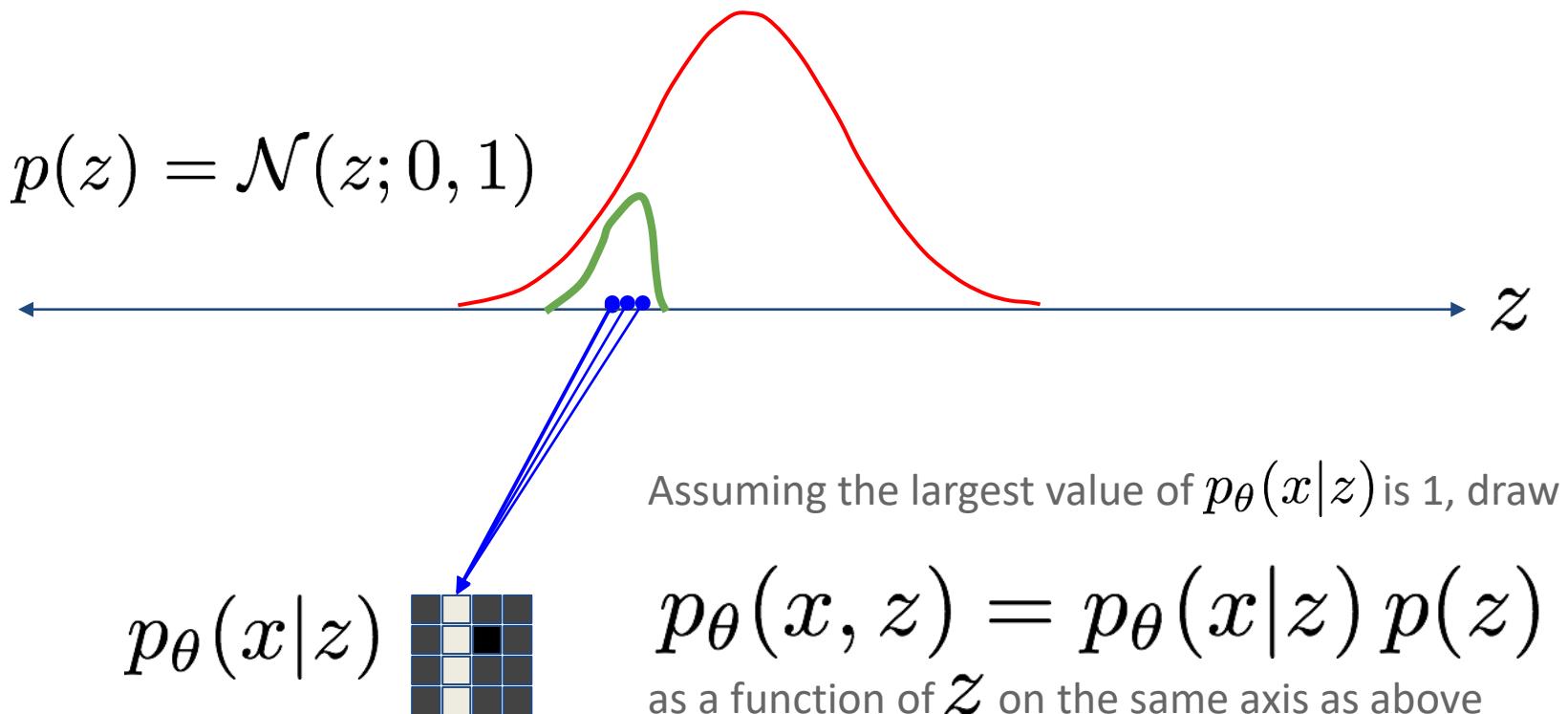
Assuming the largest value of  $p_{\theta}(x|z)$  is 1, draw

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

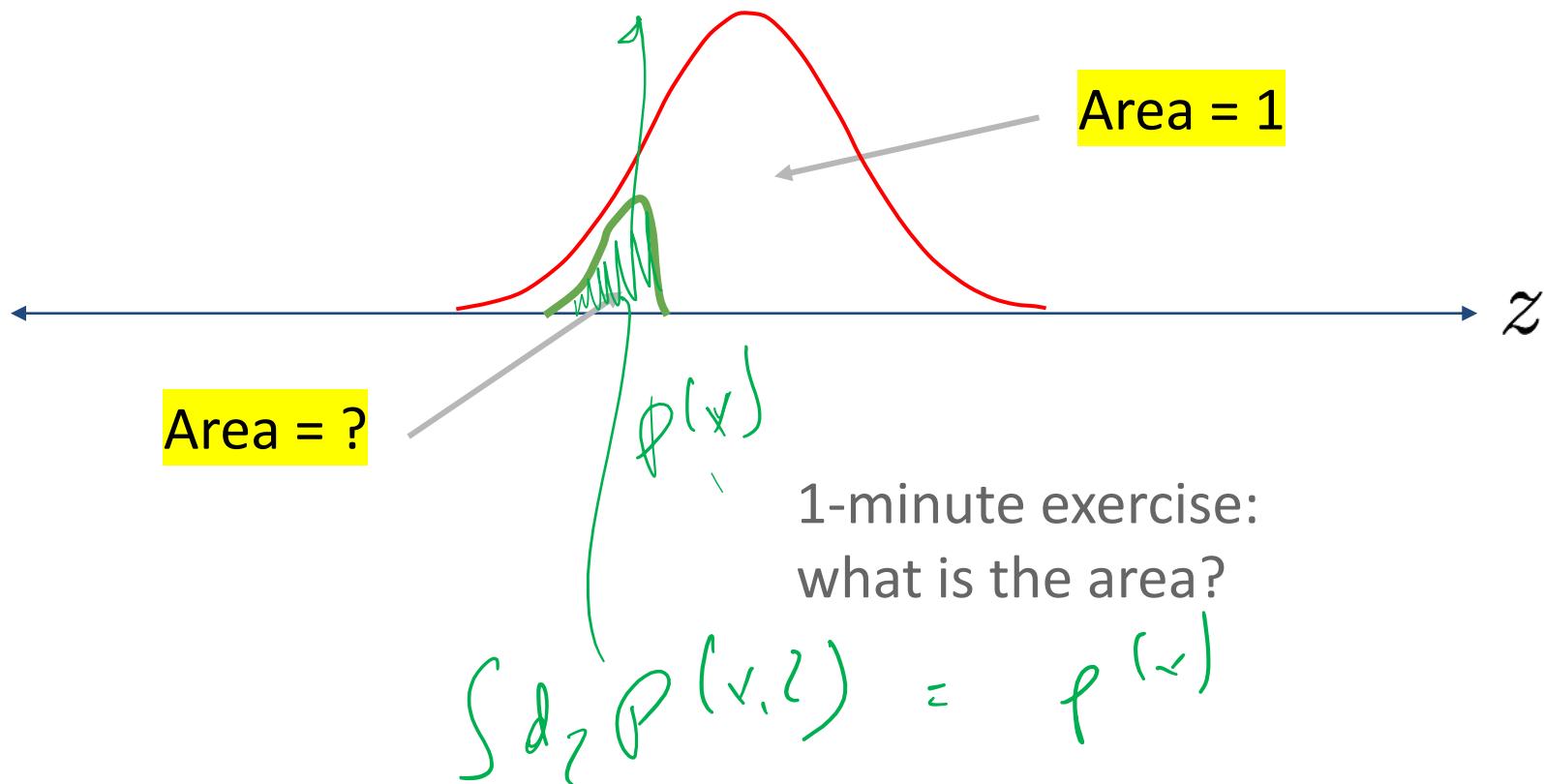
$$p_{\theta}(x, z) = p_{\theta}(x|z) p(z)$$

as a function of  $z$  on the same axis as above

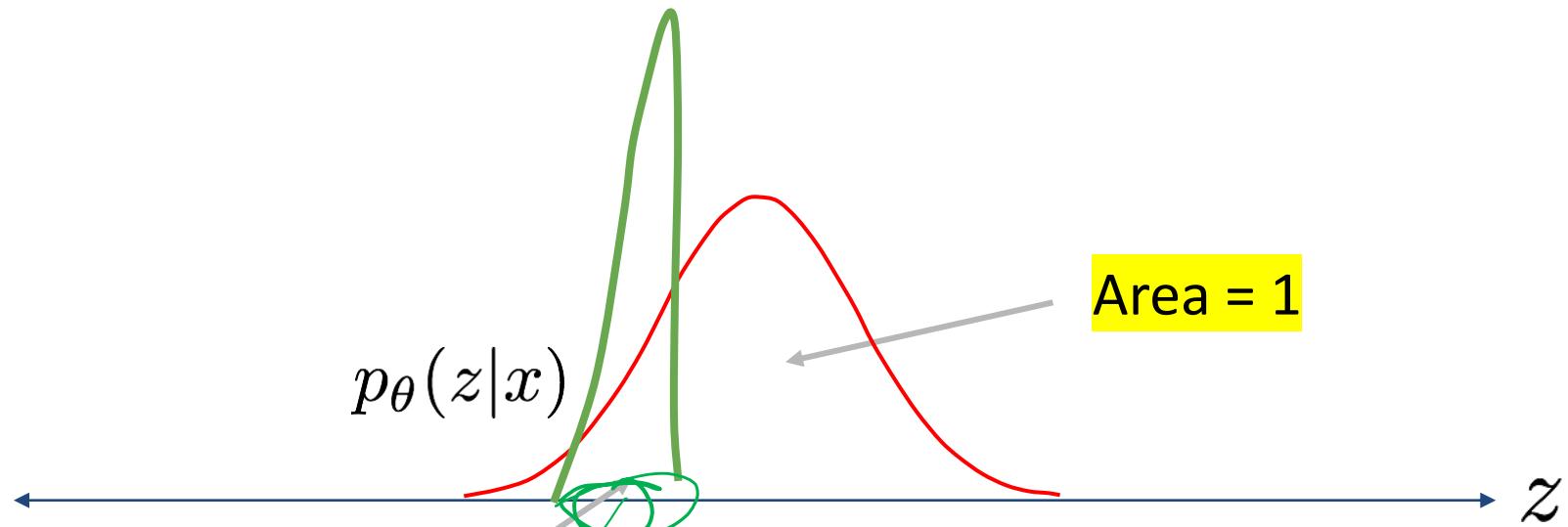
# Joint density (with $x$ observed)



# Joint density (with $x$ observed)



# Posterior



Area = 1

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

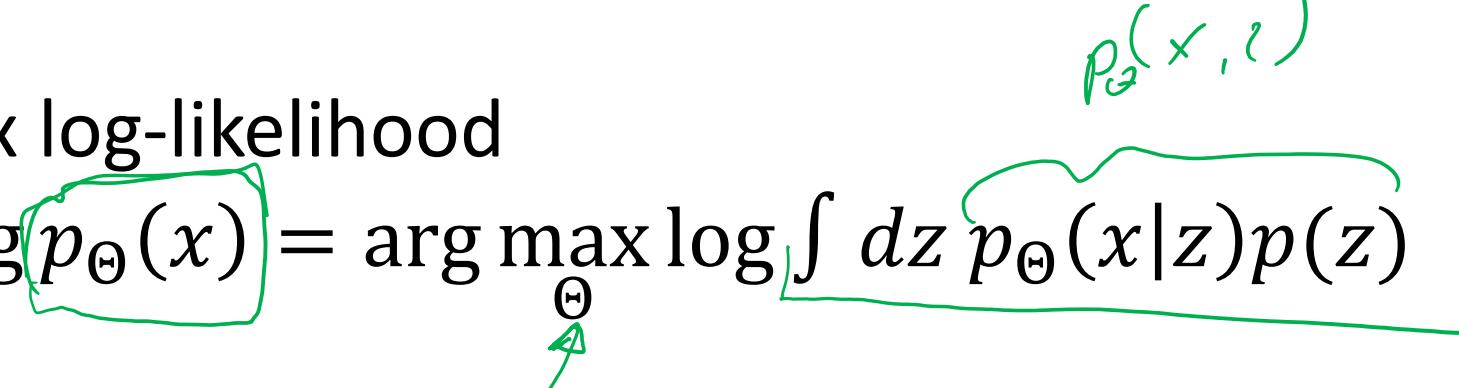
Dividing by the marginal likelihood (evidence)  
scales the area back to 1...

# Model Training Intuitions

Generation:

$$\begin{aligned} p(z) \\ p_{\Theta}(x|z) \end{aligned}$$

Training by max log-likelihood

$$\arg \max_{\Theta} \log [p_{\Theta}(x)] = \arg \max_{\Theta} \log \int dz p_{\Theta}(x|z)p(z)$$


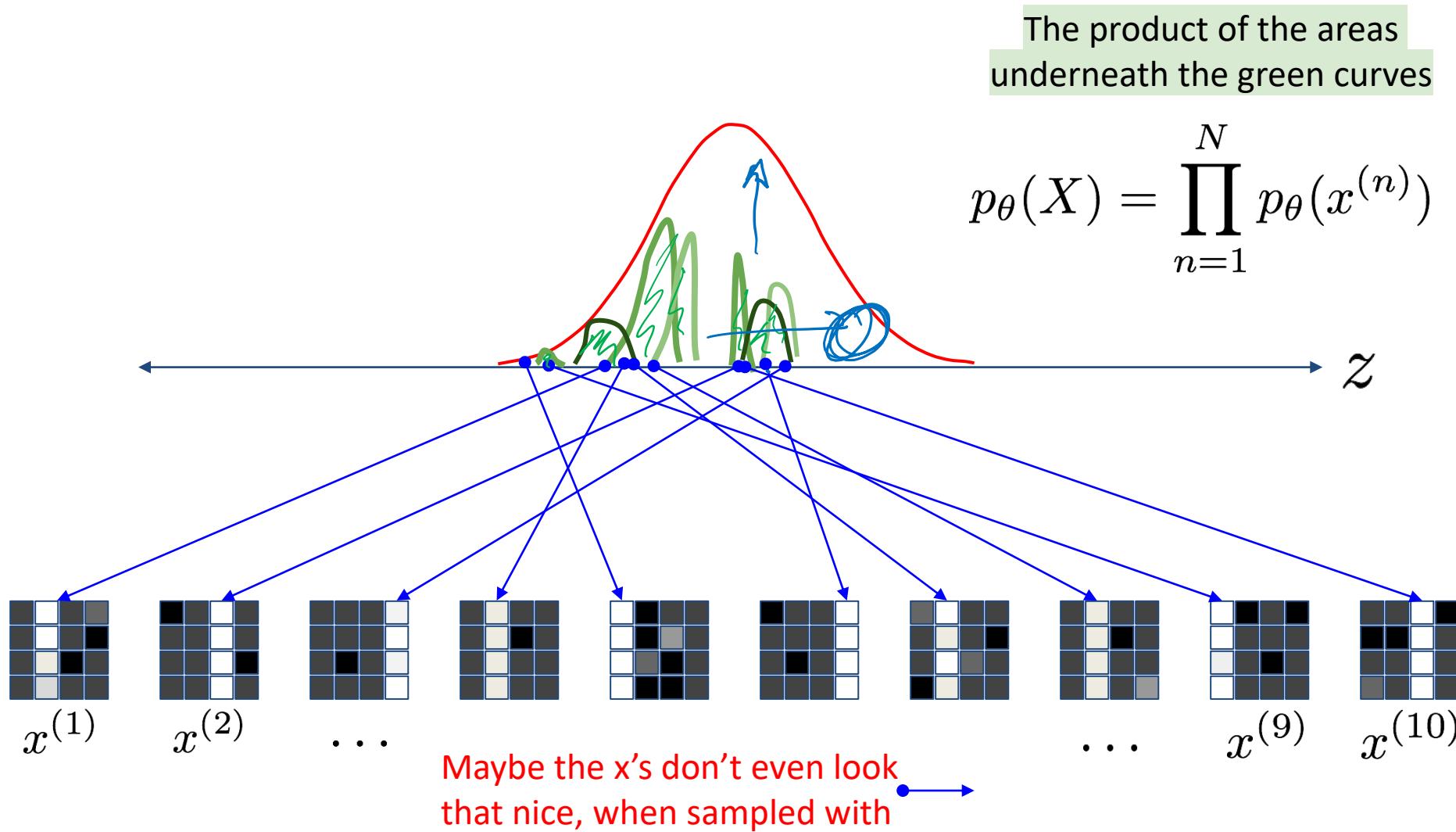
# Evidence of all data points

$$p_{\theta}(X) = \prod_{n=1}^N p_{\theta}(x^{(n)})$$

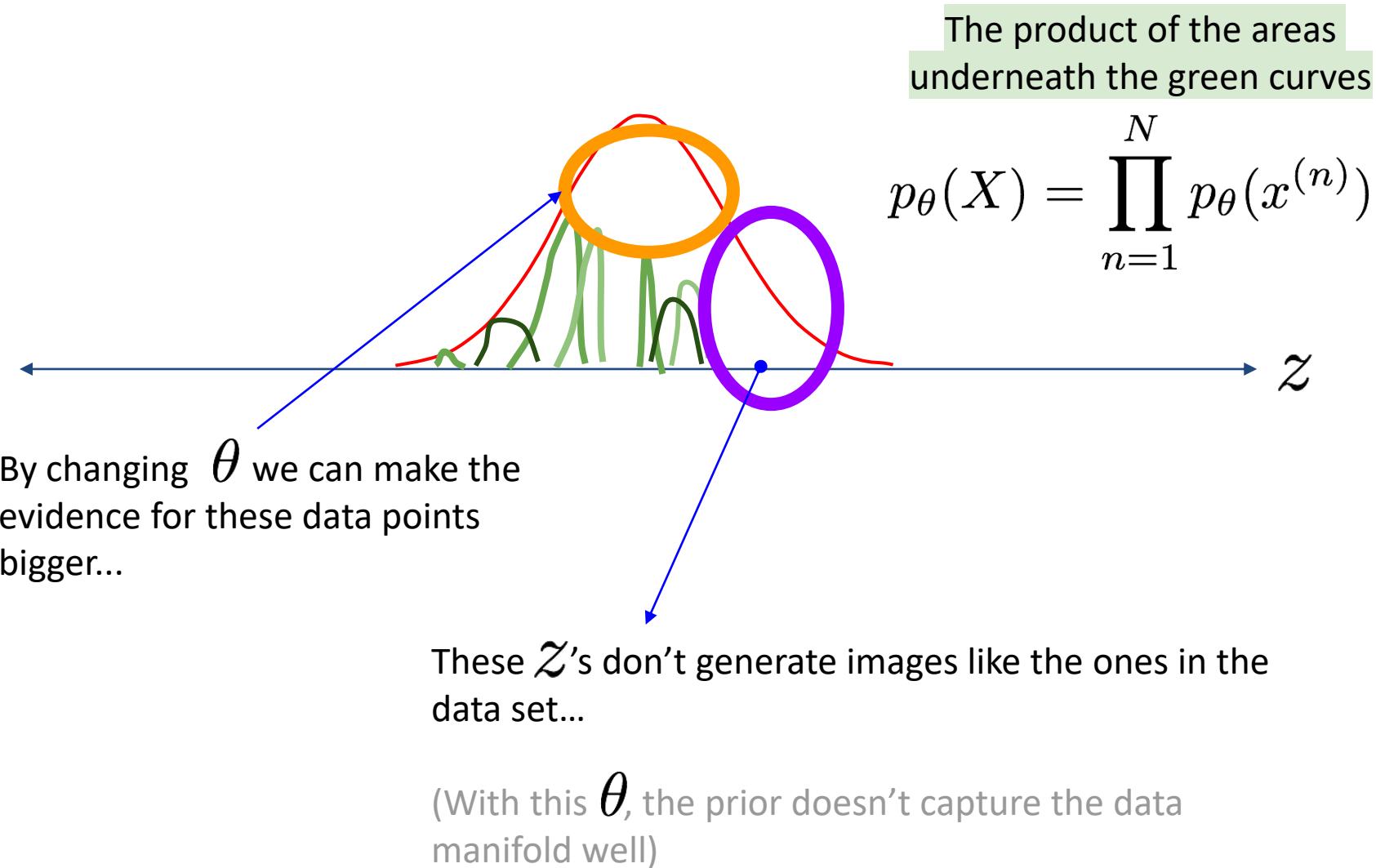
Area for data  
point  $n$

$$\log p_{\theta}(X) = \sum_{n=1}^N \log p_{\theta}(x^{(n)})$$

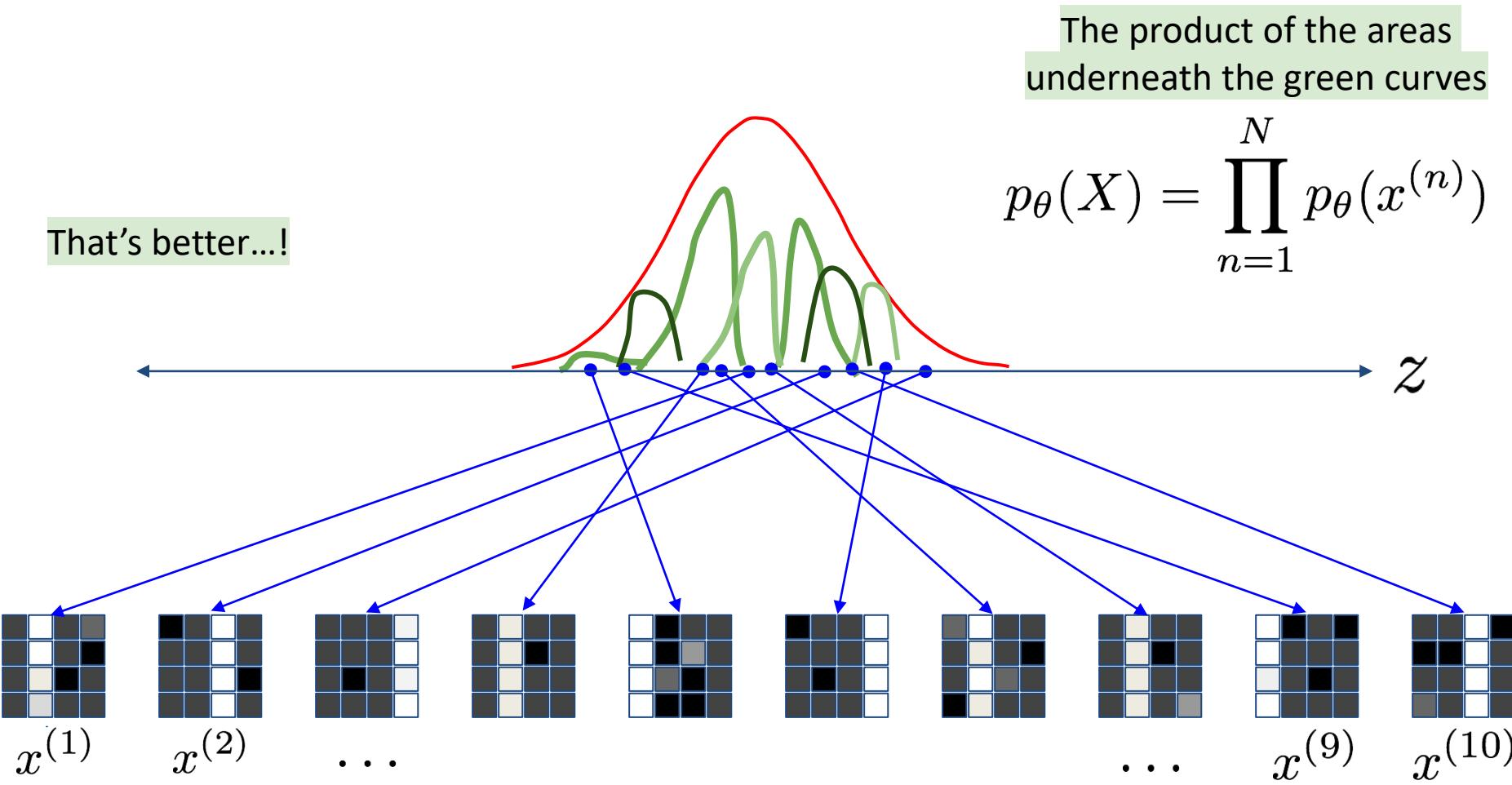
# Evidence for all data points



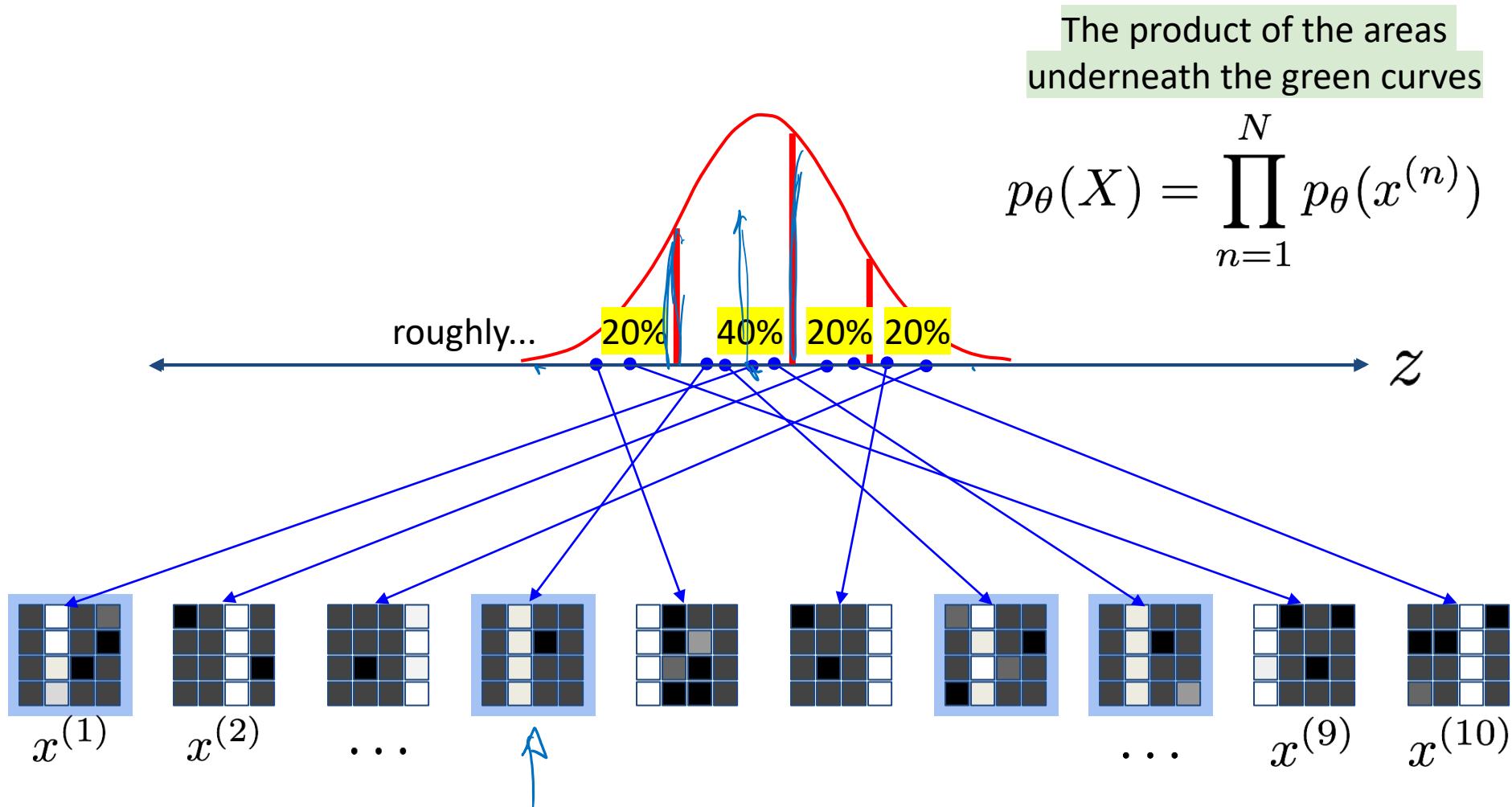
# Maximizing the evidence



# Maximizing the evidence



# For the sharp-sighted



# Training – practical aspects

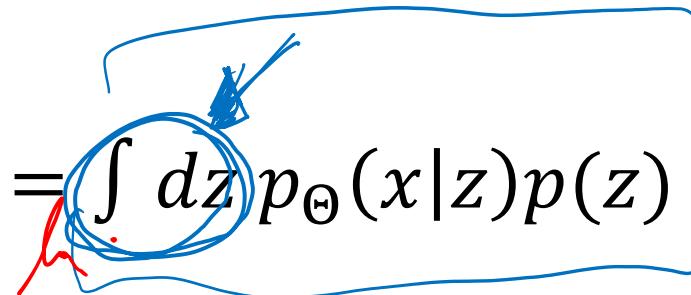
Generation:

$$\begin{aligned} p(z) \\ p_{\Theta}(x|z) \end{aligned}$$


Training by max log-likelihood

$$\arg \max_{\Theta} \log p_{\Theta}(x)$$

But

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


# Approximate likelihood optimization

Our approach:

- Lower bound  $\log p_\Theta(x)$
- Push the lower-bound up...  
... hoping to increase  $\log p_\Theta(x)$

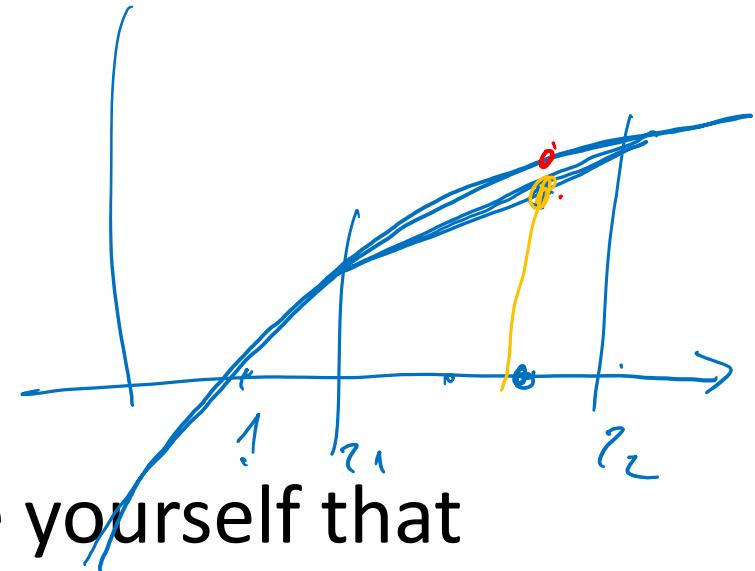
# Exercise

Jensen's inequality

Draw  $\log(\dots)$  as a function, convince yourself that

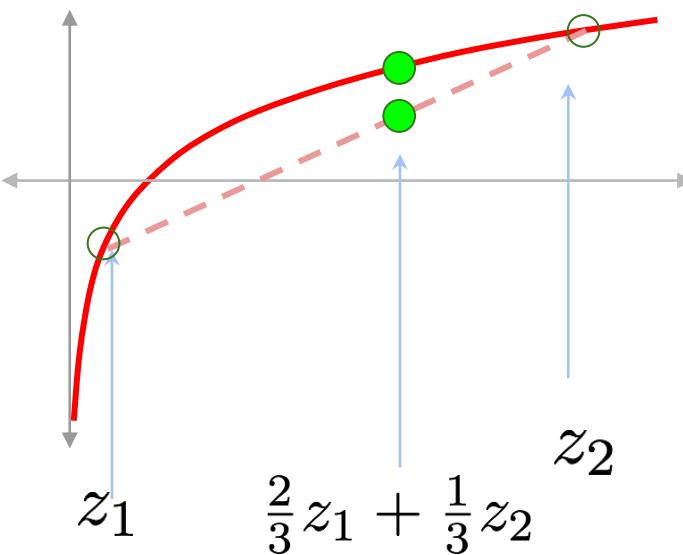
$$\log\left(\frac{2}{3}z_1 + \frac{1}{3}z_2\right) \geq \frac{2}{3}\log z_1 + \frac{1}{3}\log z_2$$

is true for any (nonnegative) setting of  $z_1$  and  $z_2$ .



# Jensen's inequality

$$\log\left(\frac{2}{3}z_1 + \frac{1}{3}z_2\right) \geq \frac{2}{3}\log(z_1) + \frac{1}{3}\log(z_2)$$



$$\log \int dz q(z) f(z) \geq \int dz q(z) \log f(z)$$

# ELBO: A likelihood bound

$$\log p_{\Theta}(x) = \log \int dz p_{\Theta}(x, z) =$$

$$= \log \int dz q_{\Phi}(z|x) \frac{p_{\Theta}(x, z)}{q_{\Phi}(z|x)}$$

← Works as long as  
 $q_{\Phi}(z|x) \neq 0$

$$\geq \int dz q_{\Phi}(z|x) \log \frac{p_{\Theta}(x, z)}{q_{\Phi}(z|x)} - \text{Jensen}$$

$\int dz q_{\Phi}(z|x) = 1$

$$= \mathbb{E}_{q_{\Phi}(z|x)} \left[ \log \frac{p_{\Theta}(x|z)p(z)}{q_{\Phi}(z|x)} \right]$$

$$= \mathbb{E}_{q_{\Phi}(z|x)} [\log p_{\Theta}(x|z)]$$

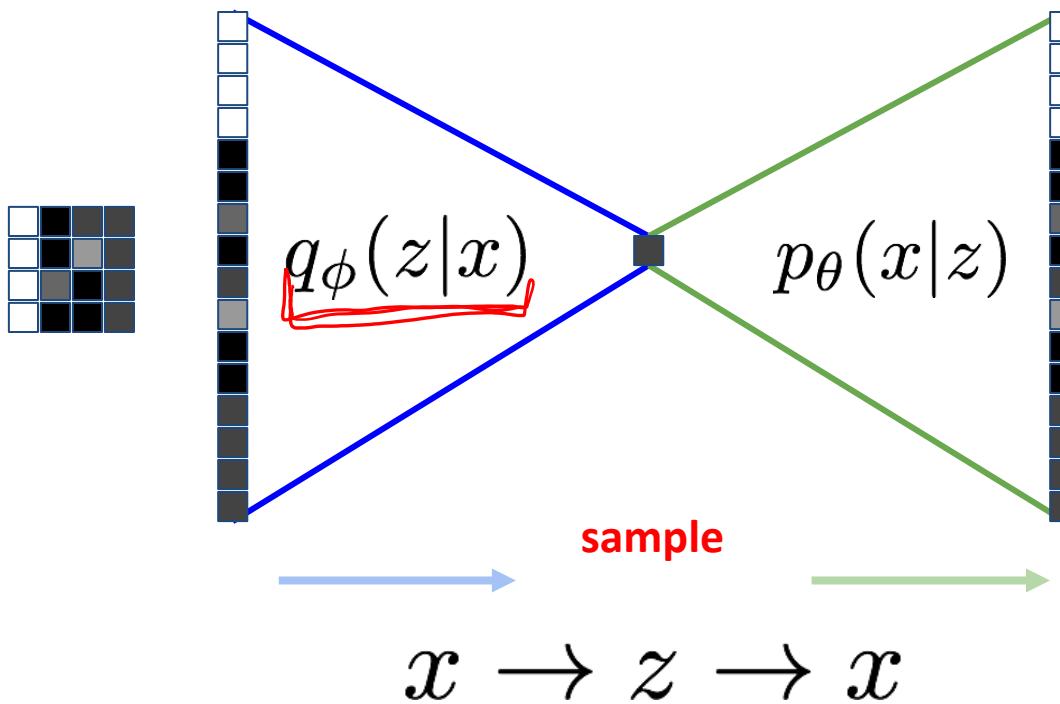
$$- \mathbb{E}_{q_{\Phi}(z|x)} \left[ \log \frac{q_{\Phi}(z|x)}{p(z)} \right]$$

$$= \mathbb{E}_{q_{\Phi}(z|x)} [\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) || p(z))$$

# ELBO interpretation

$$\log p_{\Theta}(x) \geq \underbrace{\mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)]}_{\text{auto-encoding term!}} - KL(q_{\Phi}(z|x) \parallel p(z))$$

$\mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)]$ : auto-encoding term!

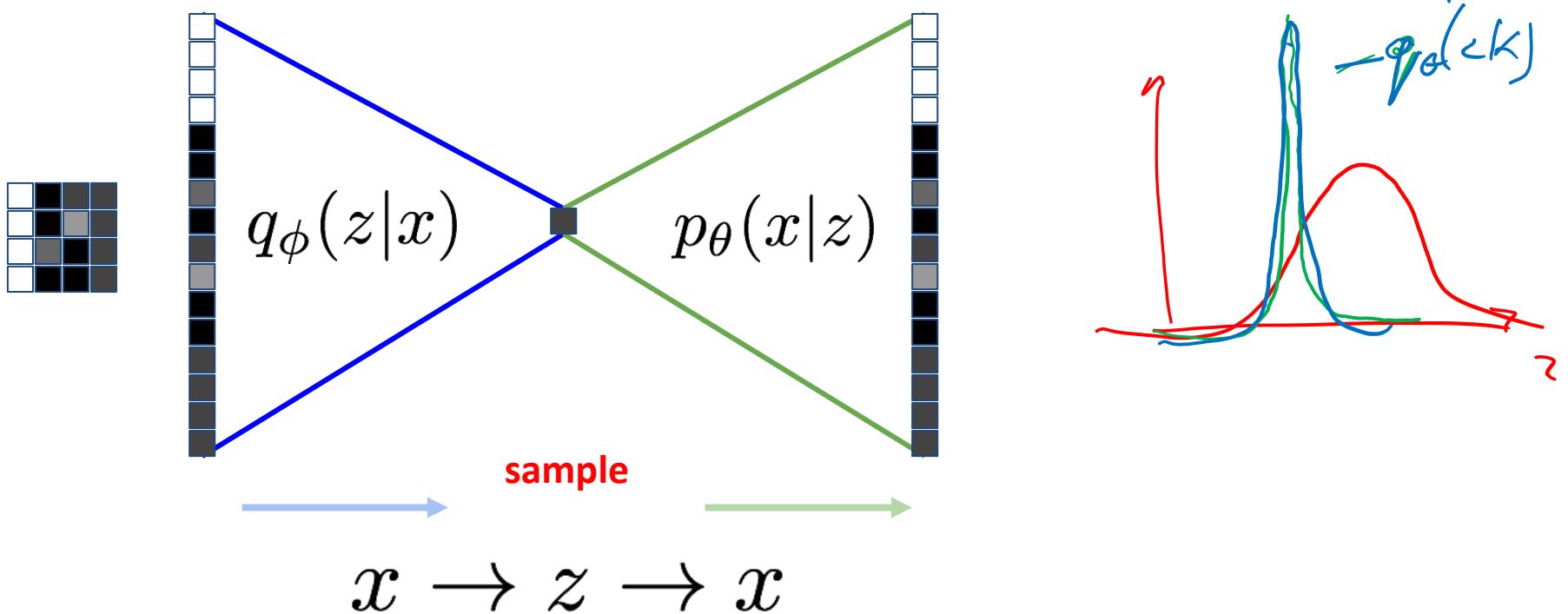


# ELBO interpretation

$$KL > 0$$

$$\log p_{\Theta}(x) \geq \mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z))$$

$KL(q_{\Phi}(z|x) \parallel p(z))$ : Amount of information transmitted in  $q_{\Phi}(z|x)$



# ELBO interpretation

Evidence lower bound

ELBO, or evidence lower bound:

$$\log p(x) \geq \mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\Phi(z|x) \parallel p(z))$$

where:

$\mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)]$  reconstruction quality:

how many nats we need to reconstruct  $x$ ,  
when someone gives us  $q(z|x)$

$KL(q_\Phi(z|x) \parallel p(z))$  code transmission cost:

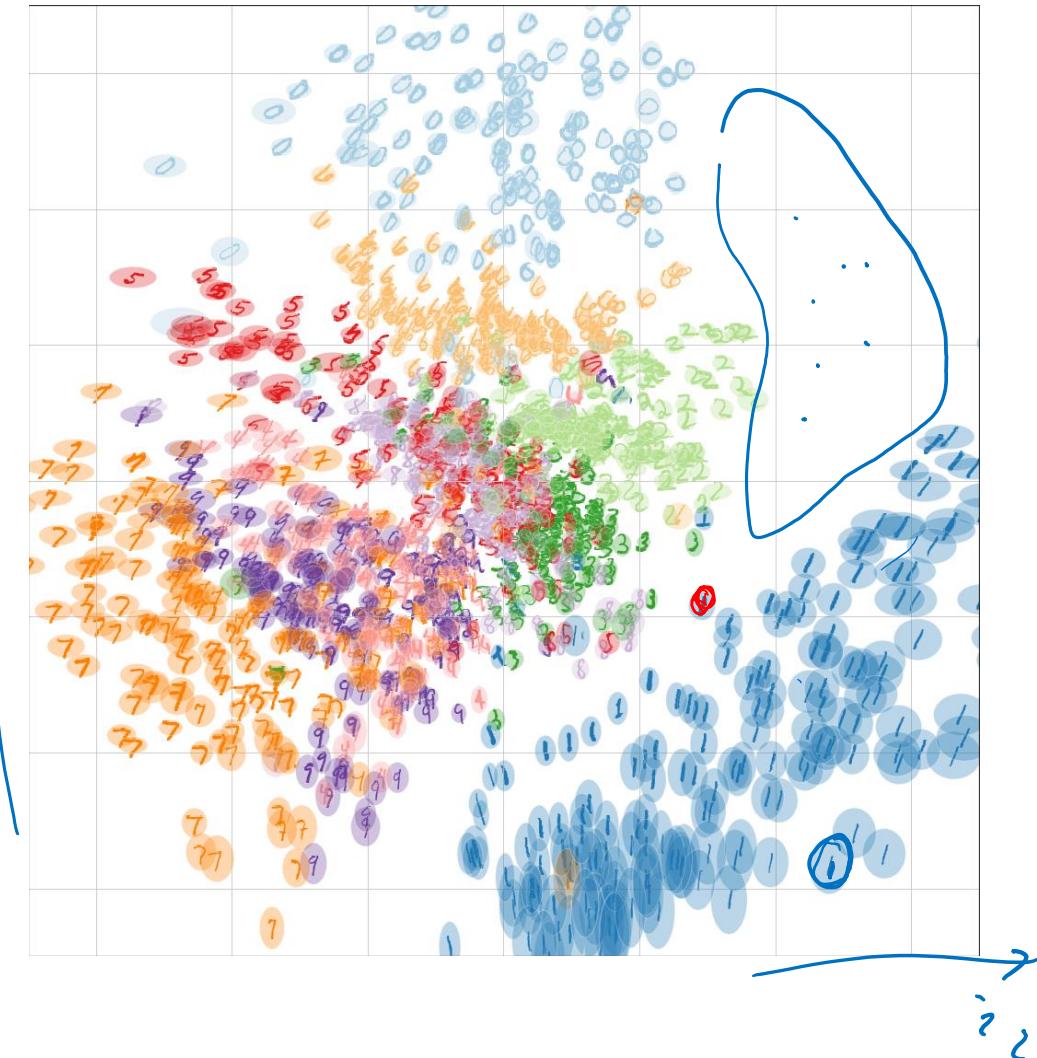
how many nats we transmit about  $x$  in  $q_\Phi(z|x)$  rather than  $p(z)$

Interpretation: do well at reconstructing  $x$ , limiting the amount of information about  $x$  encoded in  $z$ .

# VAE is an Information Bottleneck

Each sample is represented as a Gaussian

This discards information  
(latent representation has low precision)



# ELBO Evaluation

Compute:

$$\underbrace{\log p(x)}_{\text{?}} \geq \mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\Phi(z|x) \| p(z))$$

$KL(q_\Phi(z|x) \| p(z))$  has closed form for simple  $q_\Phi(z|x)$

$\mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)]$  can be approximated:

$$\mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)] \approx \sum_i^N \log p_\Theta(x|z_i)$$

Where  $z_i$  drawn from  $q_\Phi(z|x)$

$$\text{for } i \in 1 \dots N$$

# ELBO optimization

$$\log p_{\Theta}(x) \geq \mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z))$$

ELBO is a function of  $x$ ,  $\Theta$ , and  $\Phi$

What it means to maximize ELBO over  $\Phi$ ?

It can't change  $\log p_{\Theta}(x)$ ...

It tries to make the bound tight!

# Exercise

Recall Jensen's inequality:

$$\log \int dz q(z) f(z) \geq \int dz q(z) \log f(z)$$

When is it an **equality**?

When  $f(z) = \text{const}$

# When is ELBO tight?

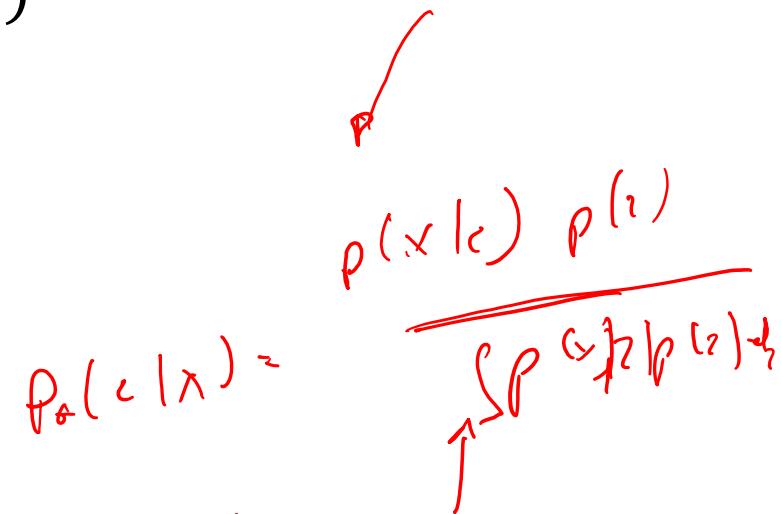
$$\begin{aligned}\log p_{\Theta}(x) &= \log \int dz q_{\Phi}(z|x) \frac{p_{\Theta}(x,z)}{q_{\Phi}(z|x)} \\ &\geq \int dz q_{\Phi}(z|x) \log \frac{p_{\Theta}(x,z)}{q_{\Phi}(z|x)} = \text{ELBO}\end{aligned}$$

When  $\underbrace{\frac{p_{\Theta}(x,z)}{q_{\Phi}(z|x)}}_{=} = \text{const!}$

What does it mean?

$$\frac{p_{\Theta}(x,z)}{q_{\Phi}(z|x)} = \frac{p_{\Theta}(z|x)p(x)}{q_{\Phi}(z|x)} = \text{const} \Rightarrow \boxed{p_{\Theta}(x|z) = q_{\Phi}(z|x)}$$

ELBO is tight when  $q_{\Phi}(z|x)$  does exact inference!



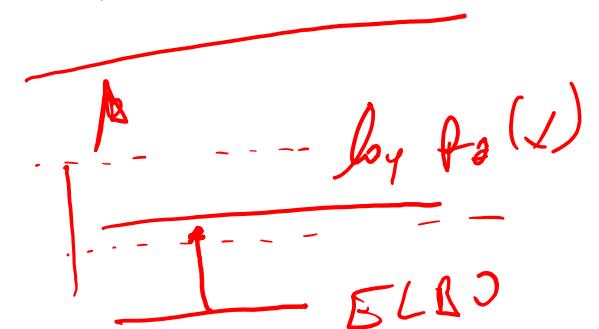
# ELBO optimization

$$\log p_{\Theta}(x) \geq \mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z))$$

ELBO is a function of  $x$ ,  $\Theta$ , and  $\Phi$

What it means to maximize ELBO over  $\Theta$ ?

Can only affect  $\mathbb{E}_{q_{\Phi}(z|x)}[p_{\Theta}(x|z)]$ !



Makes  $p_{\Theta}(x|z)$  generate back our  $x$ !

This affects  $\log p_{\Theta}(x)$ ...

...making room for improving  $q$ !

# ELBO optimization

$$\log p_{\Theta}(x) \geq \mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z))$$

Change  $\Phi$  to maximize the bound,  
making  $q_{\Phi}(z|x) \approx p_{\Theta}(z|x)$

Similar to E step 

Change  $\Theta$  to (if bound sufficiently tight)  
improve  $\log p_{\Theta}(x)$

Similar to M step

But we tune  $\Phi$  and  $\Theta$  at the same time!

# Alternative VAE Explanation: EM

Want:

$$\max \log p_{\Theta}(x) = \max \log \sum_z p_{\Theta}(x, z) = \max \log \sum_z p_{\Theta}(x|z)p_{\Theta}(z)$$

Iterate:

E-step: find  $q_{\Phi}(z|x) = p_{\Theta}(z|x)$

M-step:  $\max_{\Theta} \mathbb{E}_{z \sim q_{\Phi}(Z|x)} \left[ \log \left( \frac{p_{\Theta}(z,x)}{q_{\Phi}(z|x)} \right) \right]$  (keep  $\Phi$  fixed)

Why?

$$\log p(x) = KL(q_{\Phi}(z|x) \parallel p_{\Theta}(z|x)) + \mathbb{E}_{z \sim q_{\Phi}(Z|x)} \left[ \log \left( \frac{p_{\Theta}(z,x)}{q_{\Phi}(z|x)} \right) \right]$$

# Alternative VAE Explanation: ELBO Proof

$$\begin{aligned} & KL(q_\Phi(z|x) \parallel p_\Theta(z|x)) + \mathbb{E}_{z \sim q_\Phi(z|x)} \left[ \log \left( \frac{p_\Theta(z,x)}{q_\Phi(z|x)} \right) \right] = \\ &= \mathbb{E}_{z \sim q_\Phi(z|x)} \left[ \log \frac{q_\Phi(z|x)}{p_\Theta(z|x)} \right] + \mathbb{E}_{z \sim q_\Phi(z|x)} \left[ \log \left( \frac{p_\Theta(z,x)}{q_\Phi(z|x)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\Phi(z|x)} \left[ \log \frac{q_\Phi(z|x)}{p_\Theta(z|x)} \frac{p_\Theta(z|x)p_\Theta(x)}{q_\Phi(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_\Phi(z|x)} \log p_\Theta(x) = \log p_\Theta(x) \end{aligned}$$

# Alternative VAE Explanation: ELBO Proof

We have shown that for any  $q_\Phi(z|x)$

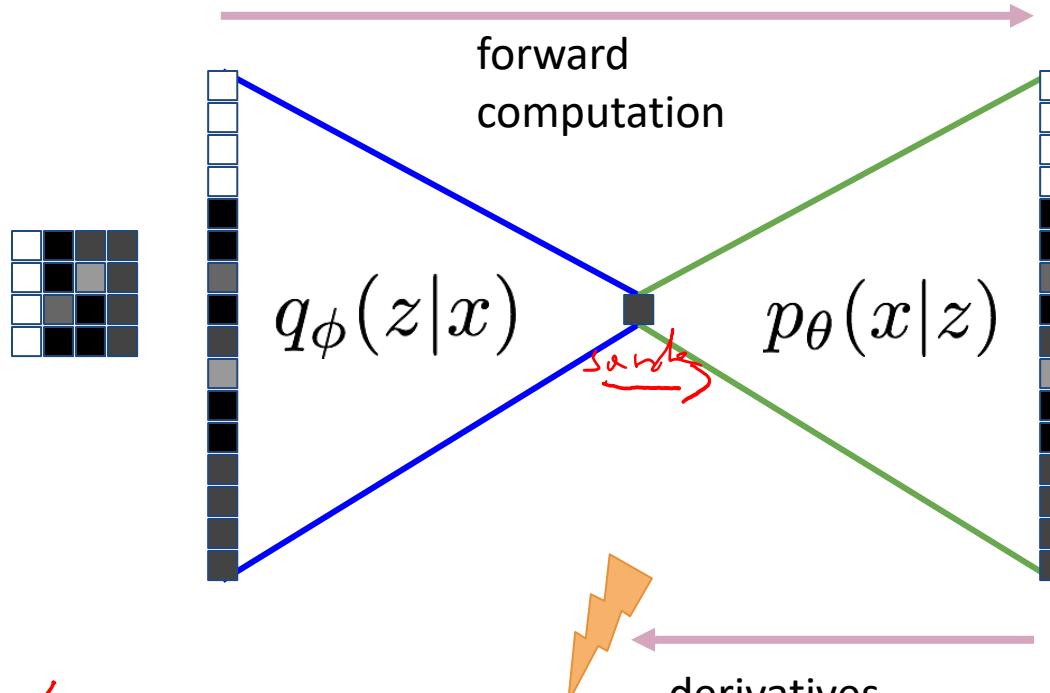
$$\begin{aligned}\log p_\Theta(x) &= KL(q_\Phi(z|x) \parallel p_\Theta(z|x)) + \mathbb{E}_{z \sim q_\Phi(z|x)} \left[ \log \left( \frac{p_\Theta(z,x)}{q_\Phi(z|x)} \right) \right] \\ &\geq \mathbb{E}_{z \sim q_\Phi(z|x)} \left[ \log \left( \frac{p_\Theta(z|x)p_\Theta(x)}{q_\Phi(z|x)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\Phi(z|x) \parallel p_\Theta(z))\end{aligned}$$

The bound is tight for  $p(z|x) = q(z|x)$  (Then  $KL(q_\Phi(z|x) \parallel p_\Theta(z|x)) = 0$  ).

Idea (VAE): Optimize ELBO using gradient techniques jointly over  $\Theta$  and  $\Phi$

# How to train a VAE?

We want to  
max ELBO



- $\log p(x) \geq \mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\Phi(z|x) \parallel p(z))$
- Forward computation involves drawing samples
- Can't backprop (get  $\frac{\partial \log p(x)}{\partial \Phi}$ ) ☹

# Reparameterization exercise

Assume that  $q_{\Phi}(z|x) = \mathcal{N}(\mu_z, \sigma_z)$ .

Exercise:

you can sample from  $\mathcal{N}(0,1)$

Q: how to draw samples from  $\mathcal{N}(\mu_z, \sigma_z)$

# Reparameterization exercise

Assume that  $q_{\Phi}(z|x) = \mathcal{N}(\mu_z, \sigma_z)$ .

Exercise:

you can sample from  $\mathcal{N}(0,1)$

Q: how to draw samples from  $\mathcal{N}(\mu_z, \sigma_z)$  —

A:

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

$$z_i = \mu_z + \sigma_z \epsilon$$

# Reparametrization to the rescue

Assume that  $q_\Phi(z|x) = \mathcal{N}(\mu_z, \sigma_z)$ .

Then:

$$\begin{aligned} & \mathbb{E}_{z \sim q_\Phi(z|x)} [\log p_\Theta(x|z)] \\ &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} [\log p_\Theta(x|\mu_z + \sigma_z \epsilon)] \end{aligned}$$

$\epsilon$  is drawn from a fixed distribution.

With  $\epsilon$  given, the computation graph is deterministic -> we can backprop!

$\mu_z, \sigma_z$  are computed  
by  $q_\Phi(z|x)$

# VAE in action

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

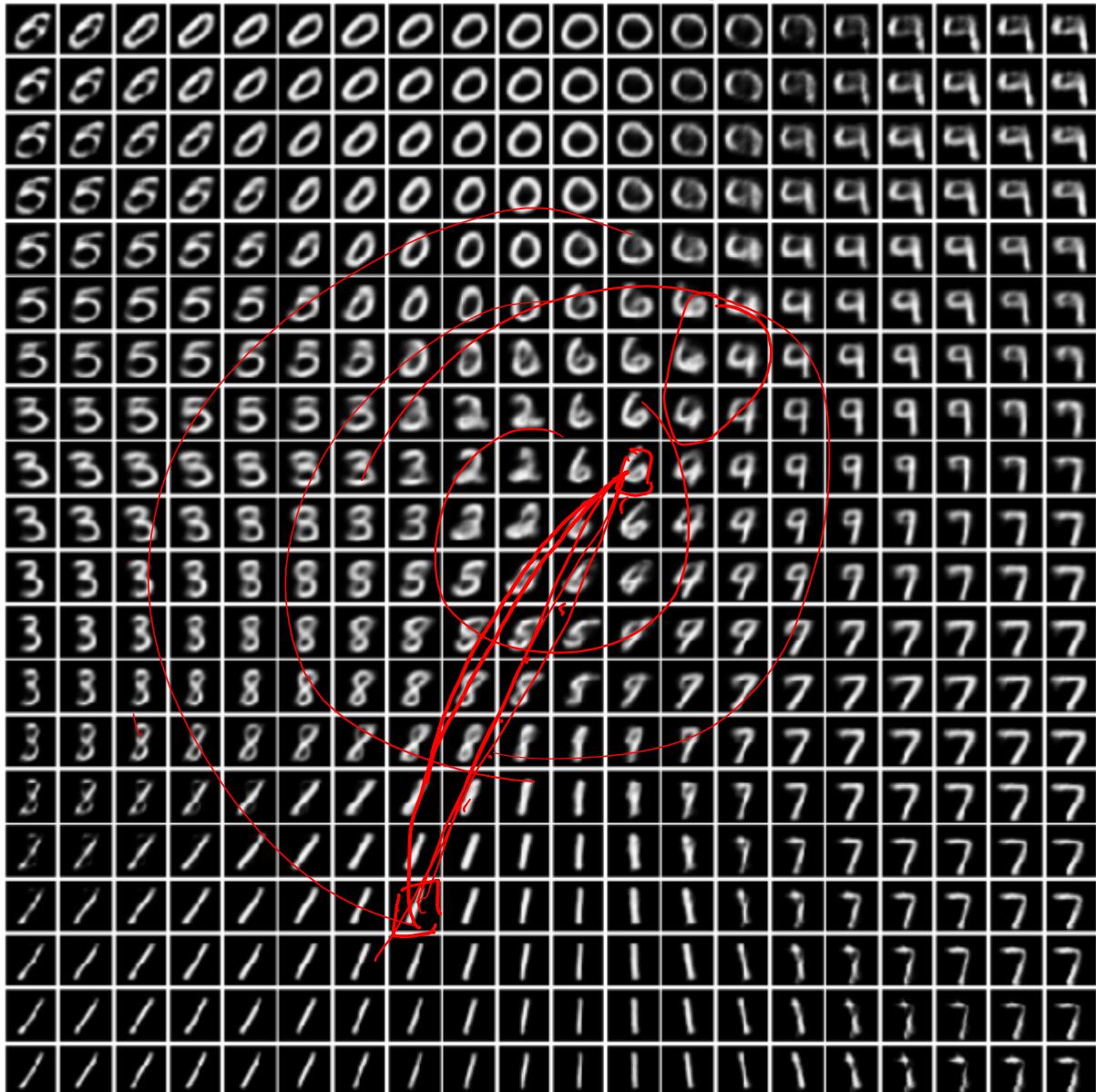
Reconstructions using the latent expected value									
9	0	2	7	0	2	0			
3	6	8	5	0	2	6			
7	9	8	3	0	3	2			
6	6	3	4						

Reconstructions from multiple latent samples

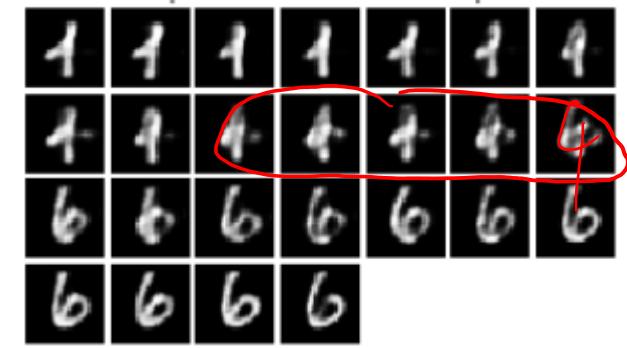
2	2	2	2	2
2	2	2	2	2
2	2	2	2	2
2	2	2	2	2

# VAE in action

Reconstructions from a 2D latent space



interpolation between two samples

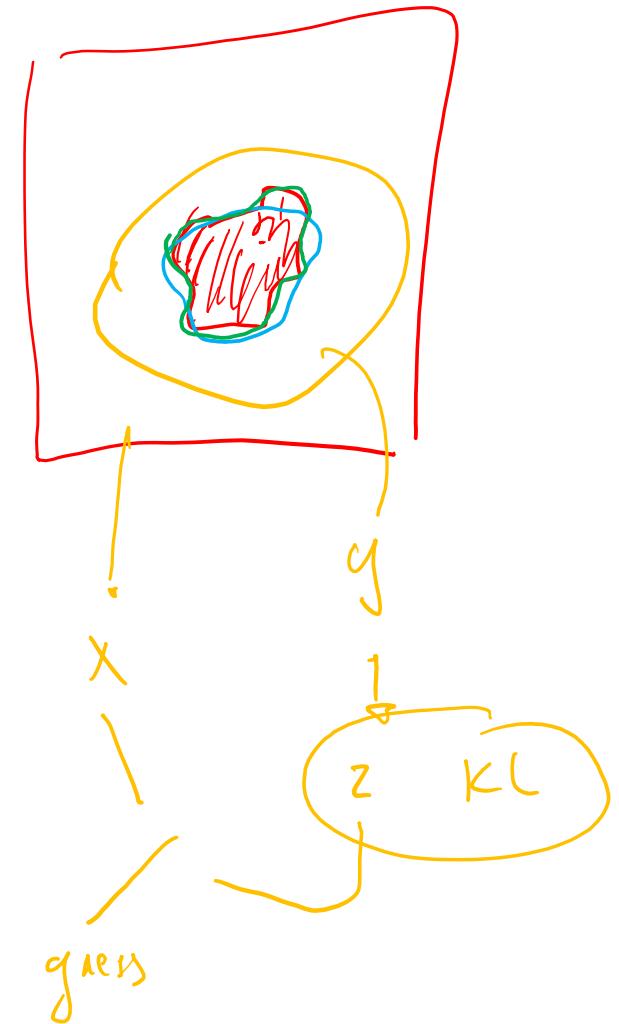
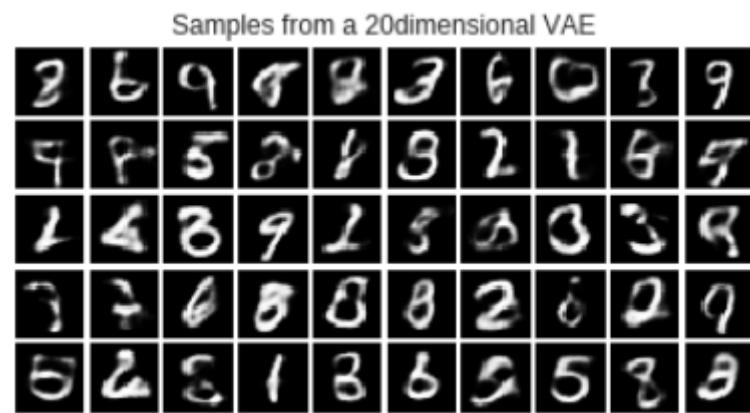


$$\mu_1, \sigma_1 = q(\cdot|v_1)$$

$$\mu_2, \sigma_2 = \phi(\cdot|x_2)$$

$$z_s = \mu_1 + 0.9\mu_2 + 0.1\mu_3 + \dots - \mu_n$$

# VAE in action



# **NORMALIZING FLOWS**

# Exercise

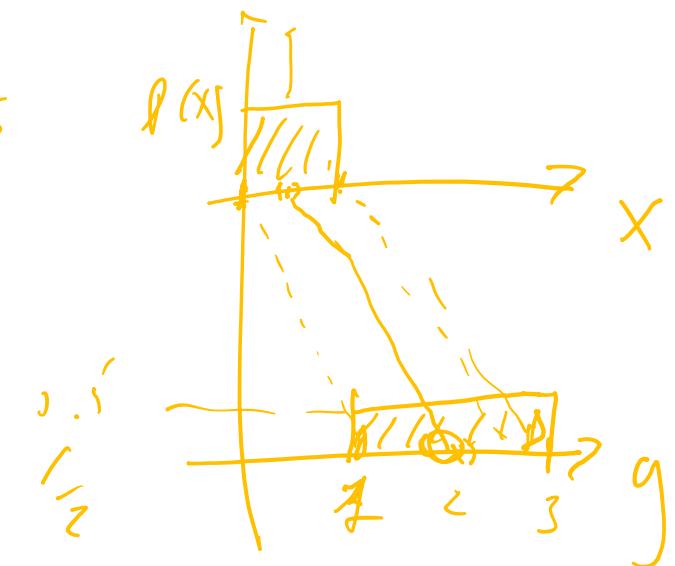
$$p(x) = \text{uniform}(0,1) = \begin{cases} 1 & \text{for } x \in (0,1) \\ 0 & \text{otherwise} \end{cases}$$

$$y = 2x + 1$$

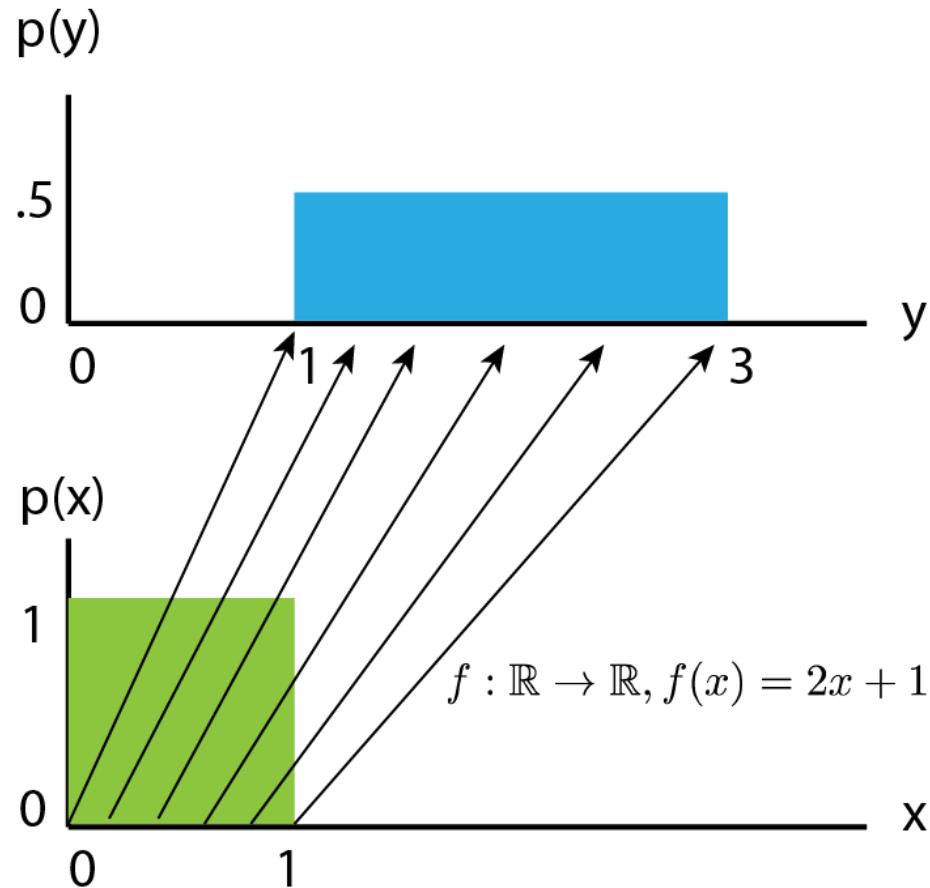
$p(y) = ?$

$$\begin{cases} 1 & \text{for } y \in (1,3) \\ 0 & \text{otherwise} \end{cases}$$

$$y = x^2$$



# Transforming distributions



$y$  also has a  
uniform  
distribution!

$$\begin{aligned} p(y) &= \text{uniform}(1, 3) \\ &= \begin{cases} 1/2 & \text{for } y \in (1, 3) \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

# The Jacobian: stretching space

Let  $y = f(x)$

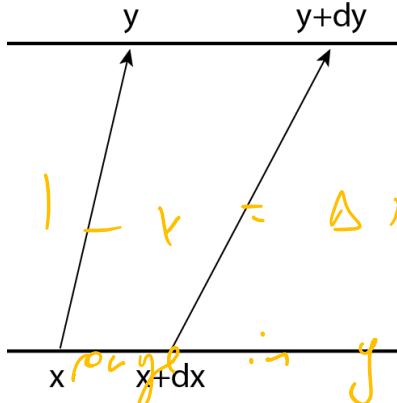
Take a range  $(x, x + \Delta x)$

Question:

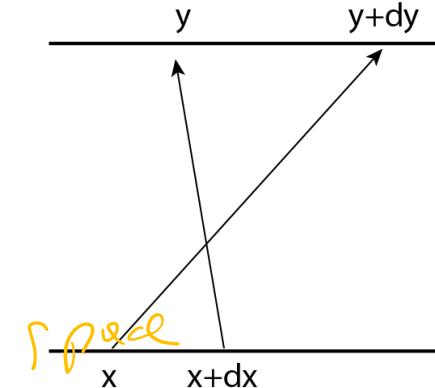
How long is this range?

$$\Delta x = (x + \Delta x) - x = \Delta x$$

*The same range in y*



$$\frac{dy}{dx} < 0$$



Question:

How long is  $(f(x), f(x + \Delta x))$ ?

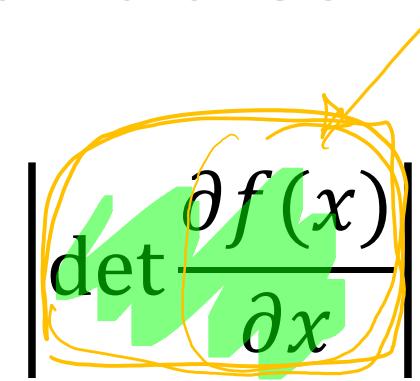
$$f(x + \Delta x) \approx \underline{f(x)} + \frac{\partial f}{\partial x} \Delta x$$

$$\left| \frac{\partial f}{\partial x} \Delta x \right|$$

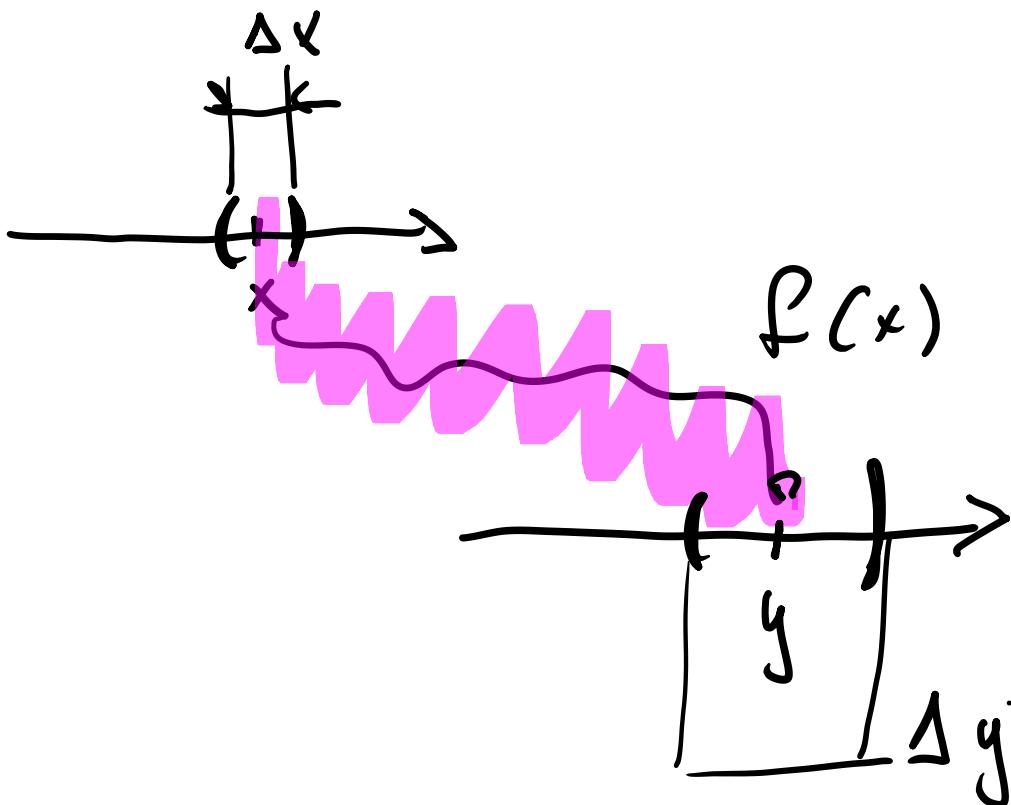
# Change of variables

$y = f(x)$ ,  $f$  is a bijection

$$p_x(x) = p_y(f(x))$$



$f$  transforms  $x \pm \frac{\Delta x}{2}$  to  $y \pm \frac{\Delta y}{2}$



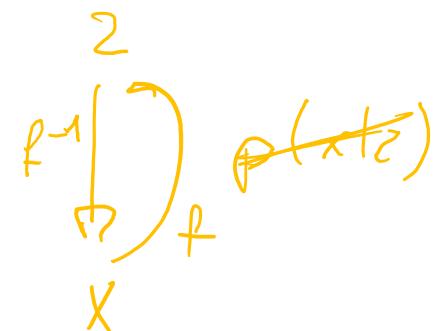
The space is stretched by

$$\frac{\partial f(x)}{\partial x} \approx \frac{\Delta y}{\Delta x}$$

# Idea

Start with  $z \sim \mathcal{N}(0,1)$

Then  $x = f^{-1}(z)$  (equivalently  $z = f(x)$ )



$$p_x(x) = p_z(f(x)) \left| \det \frac{\partial f(x)}{x} \right|$$

Tractable when:

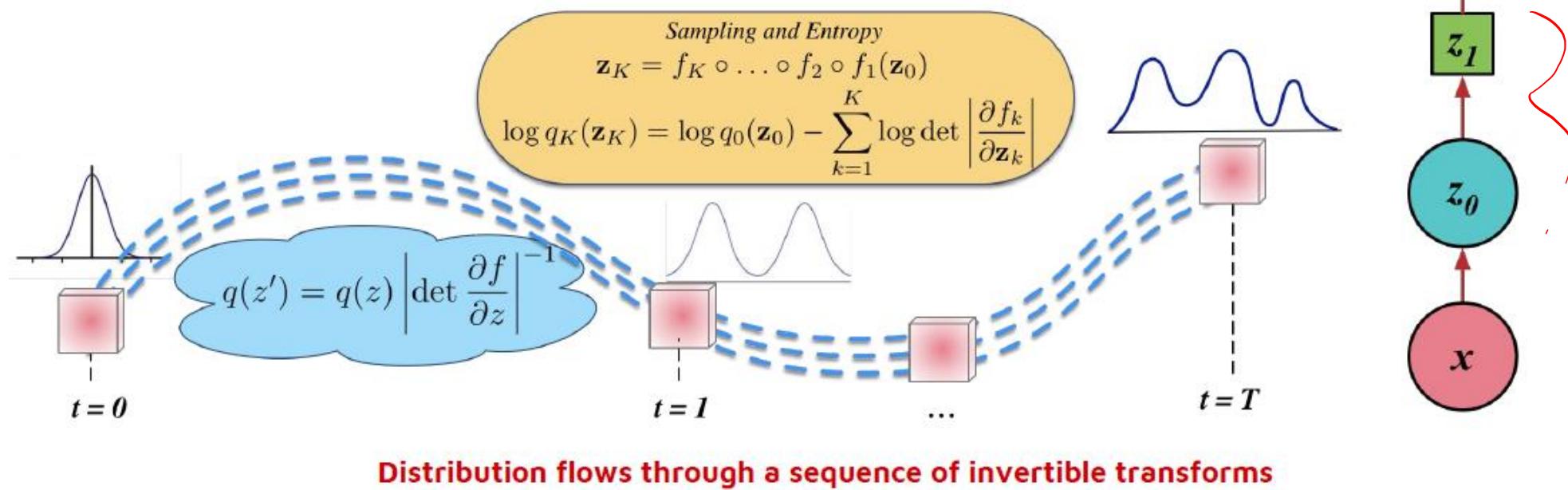
- $f$  is easy to exactly invert  $\cancel{P}$
- $f$  and  $f^{-1}$  form an exact auto-encoder!
- $\det \frac{\partial f(x)}{x}$  is easy to compute  $\cancel{A}$

=> We need special  $f$ !

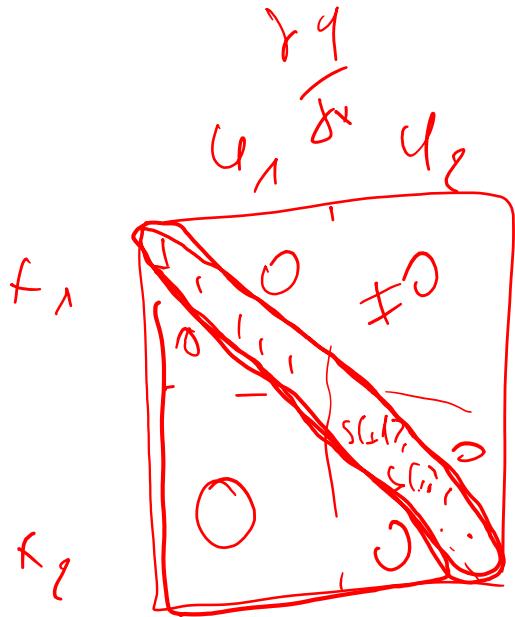
# Normalizing Flows

Exploit the rule for change of variables:

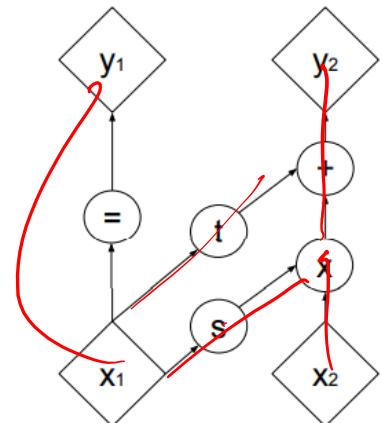
- Begin with an initial distribution
- Apply a sequence of K invertible transforms



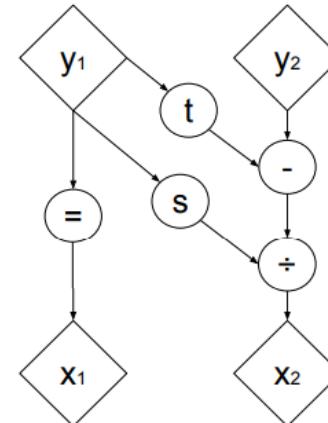
Rezende and Mohamed, 2015



# A special form of $f$



(a) Forward propagation



(b) Inverse propagation

$$x \rightarrow y$$

$$\begin{aligned} x_1, x_2 &= \text{split}(x) \\ y_1 &= x_1 \\ y_2 &= x_2 s(x_1) + t(x_1) \end{aligned}$$

Trivial to invert!

$\frac{\partial y}{\partial x}$  is diagonal, determinant is easy!

L. Dinh et al "NICE", L. Dinh et al "Real NVP"

VAE  $z \sim p(z)$   
 Flow  $N(x)$   
 $p_\theta(x|z)$   
 $x = f(z)$   
 Train  
 ELBO  
 try to find  
 of "inverting"  
 $p(x|z)$   
 train with  
 exact  $\det J_f$   
 word...  
 generic  
 neural net  
 - constain  
 nowal if  
 - inv  
 - easy gradient

# Normalizing flows in action

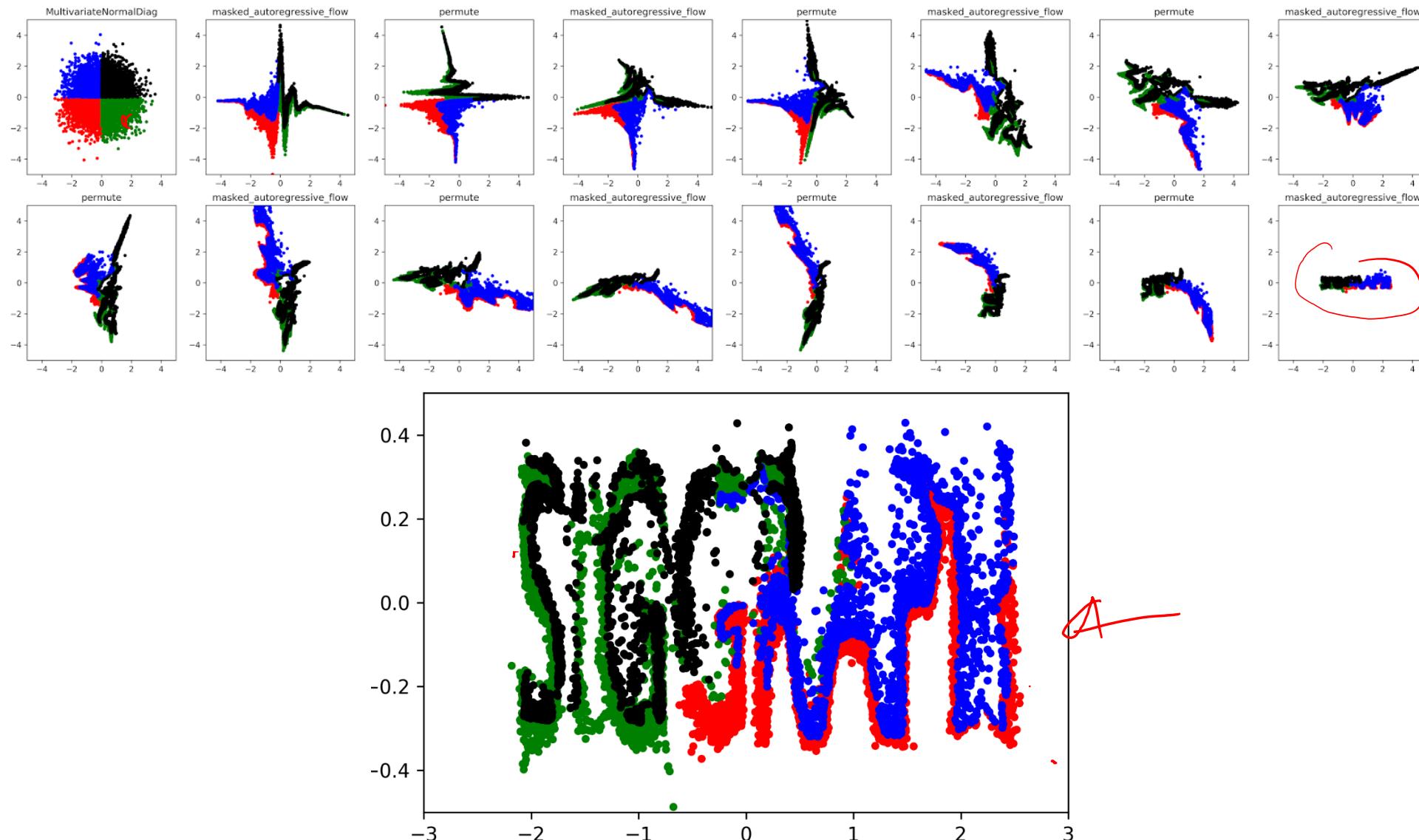
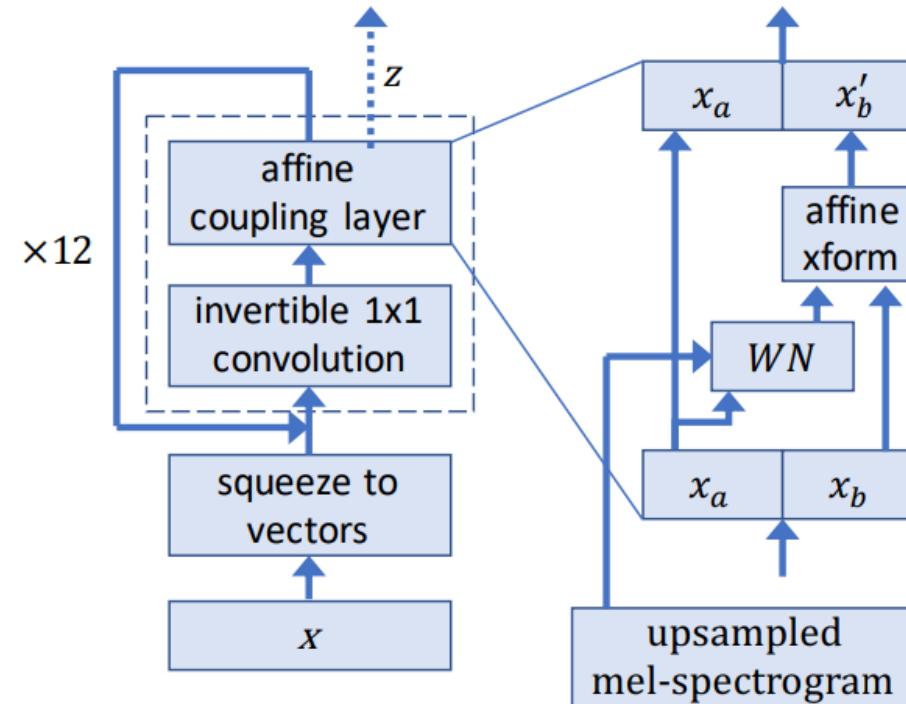


Image credit Eric Jang, <https://blog.evjang.com/2018/01/nf2.html>

# Normalizing flows in action



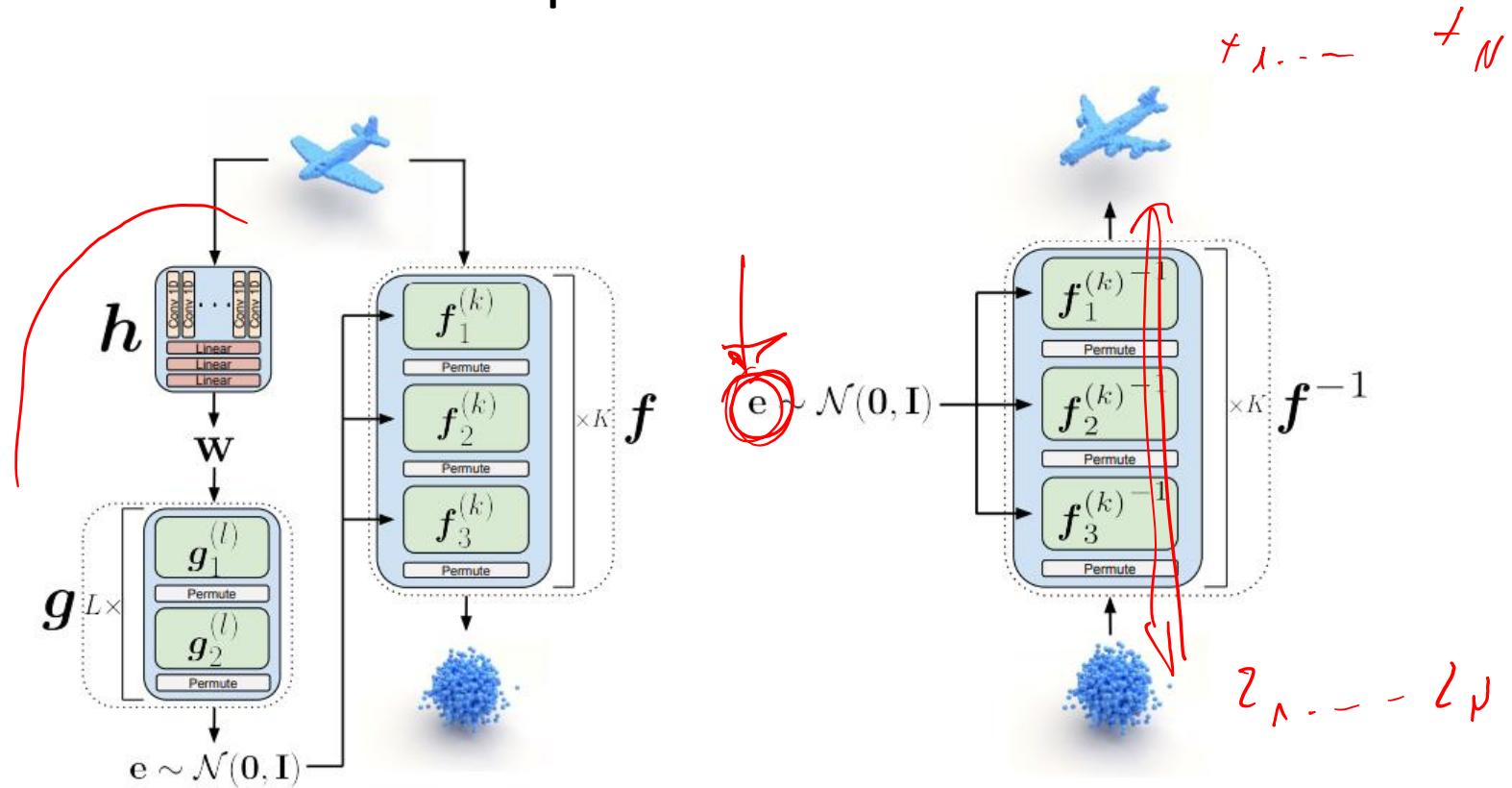
Kingma et al, “GLOW”



Prenger et al, “WAVEGLOW”

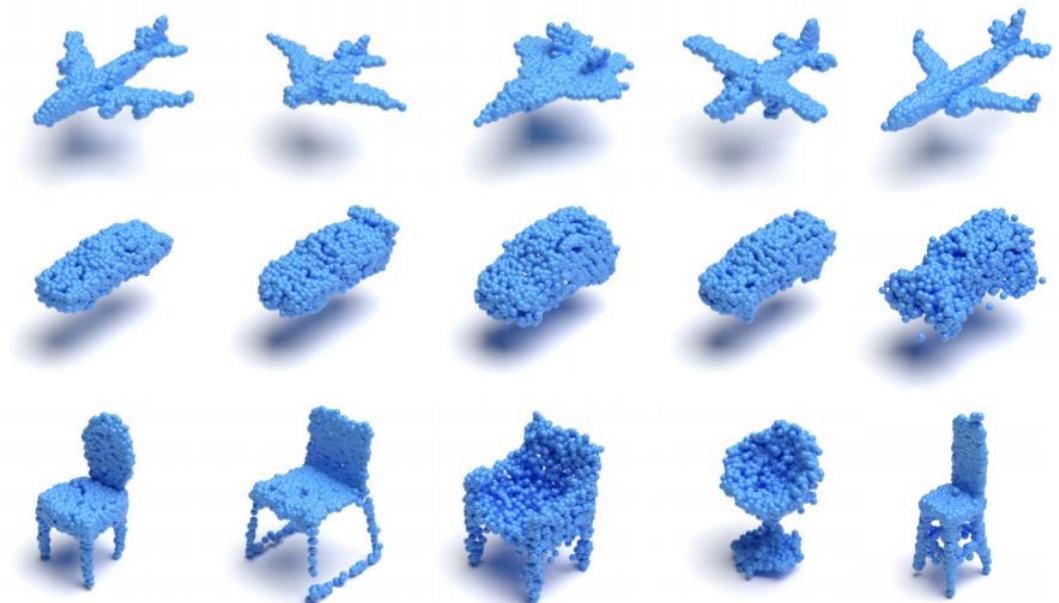
# Normalizing Flows in Action

A generative model for point clouds:

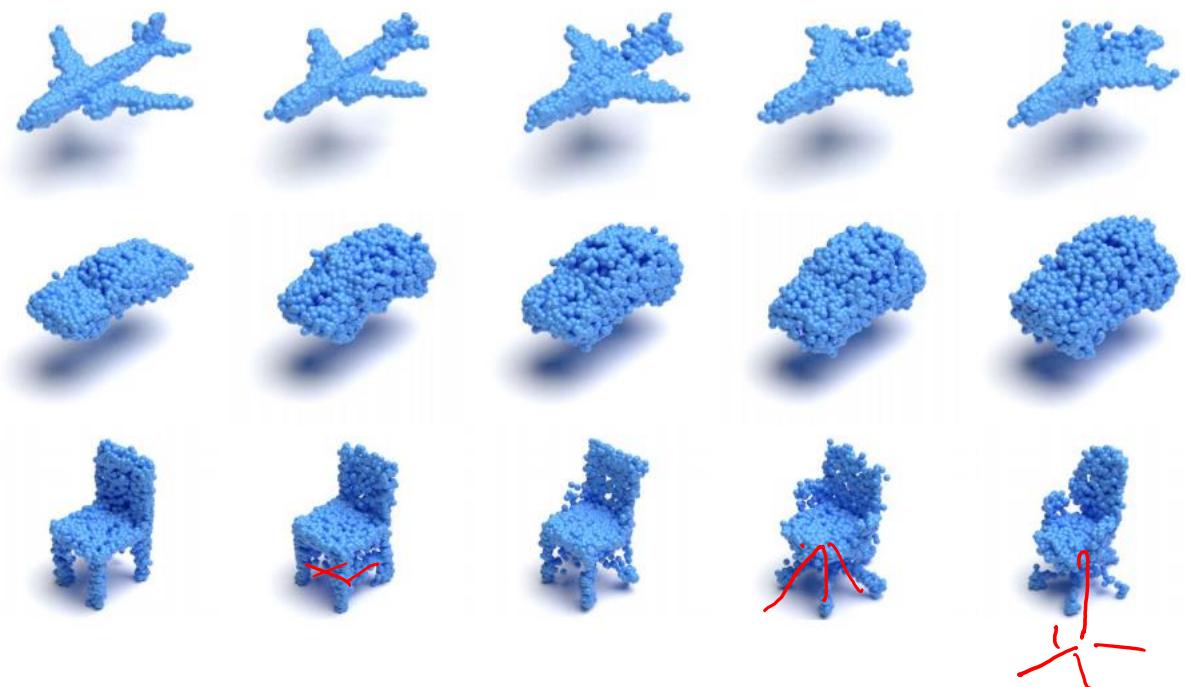


# Normalizing Flows in Action

Samples for different  $\epsilon$

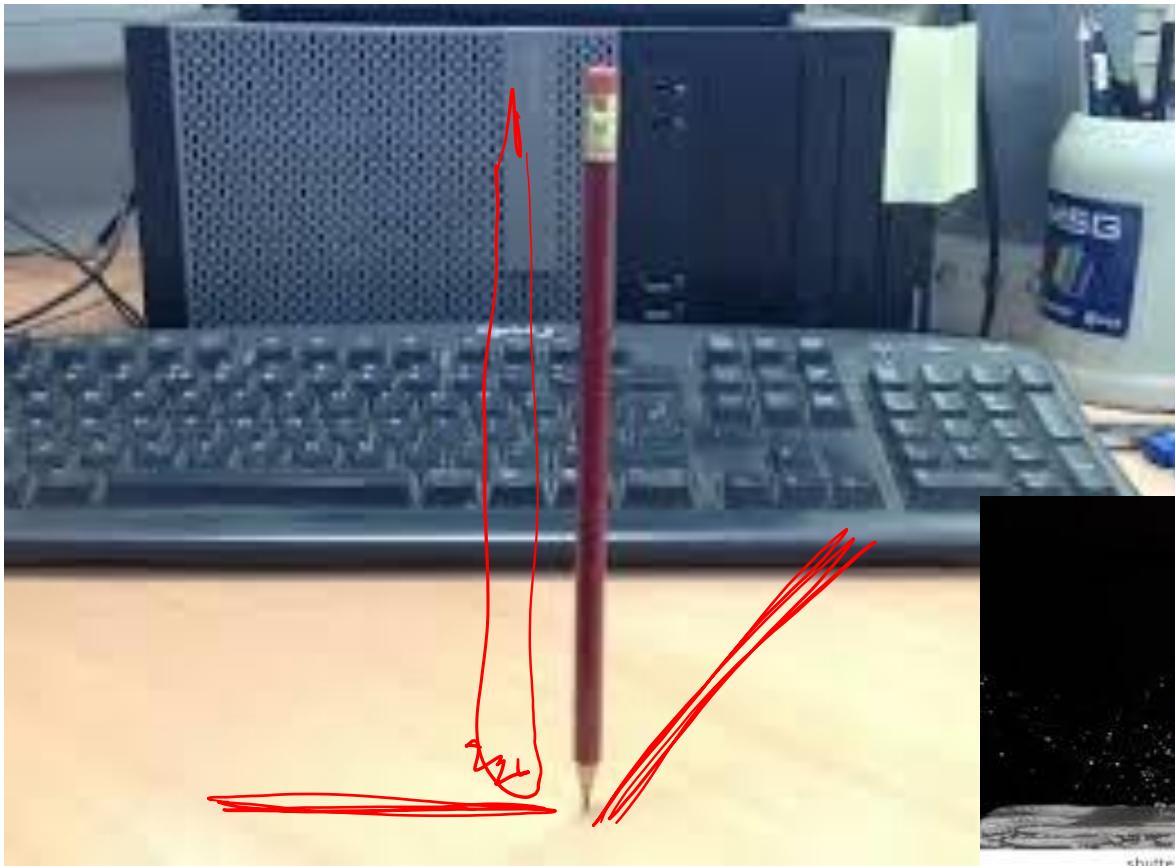


Interpolations  $\epsilon$



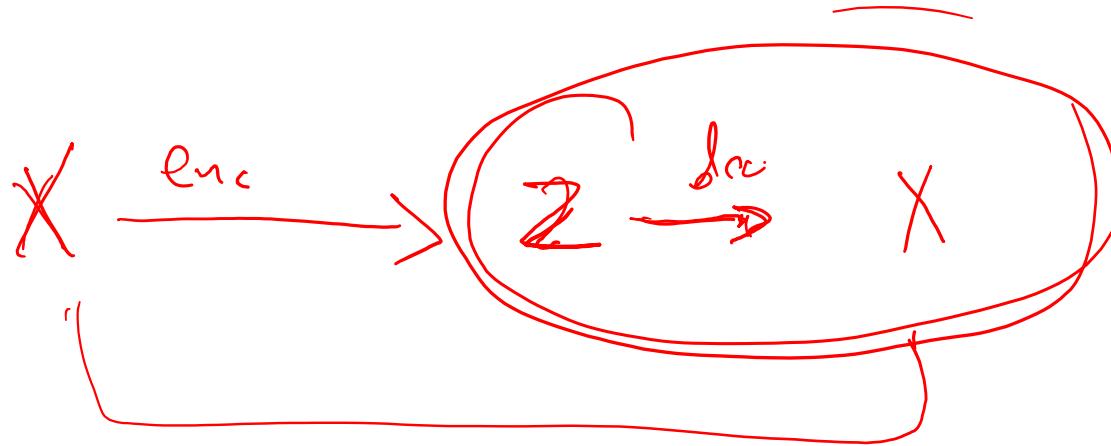
# **GAN**

# What will happen next?



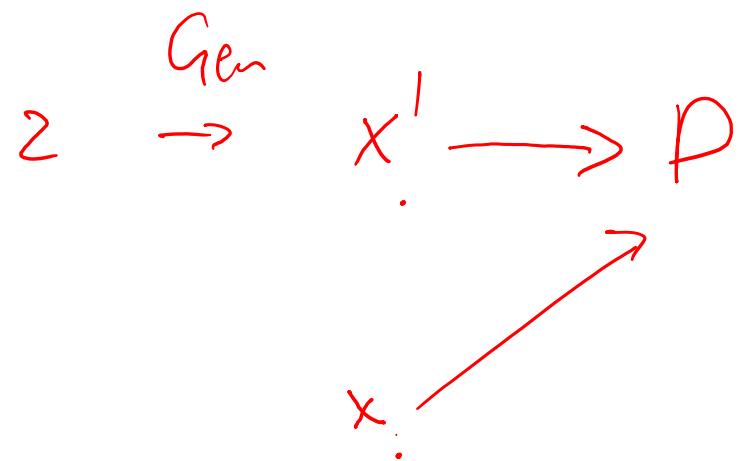
shutterstock.com • 767078788

# Trouble with autoencoders

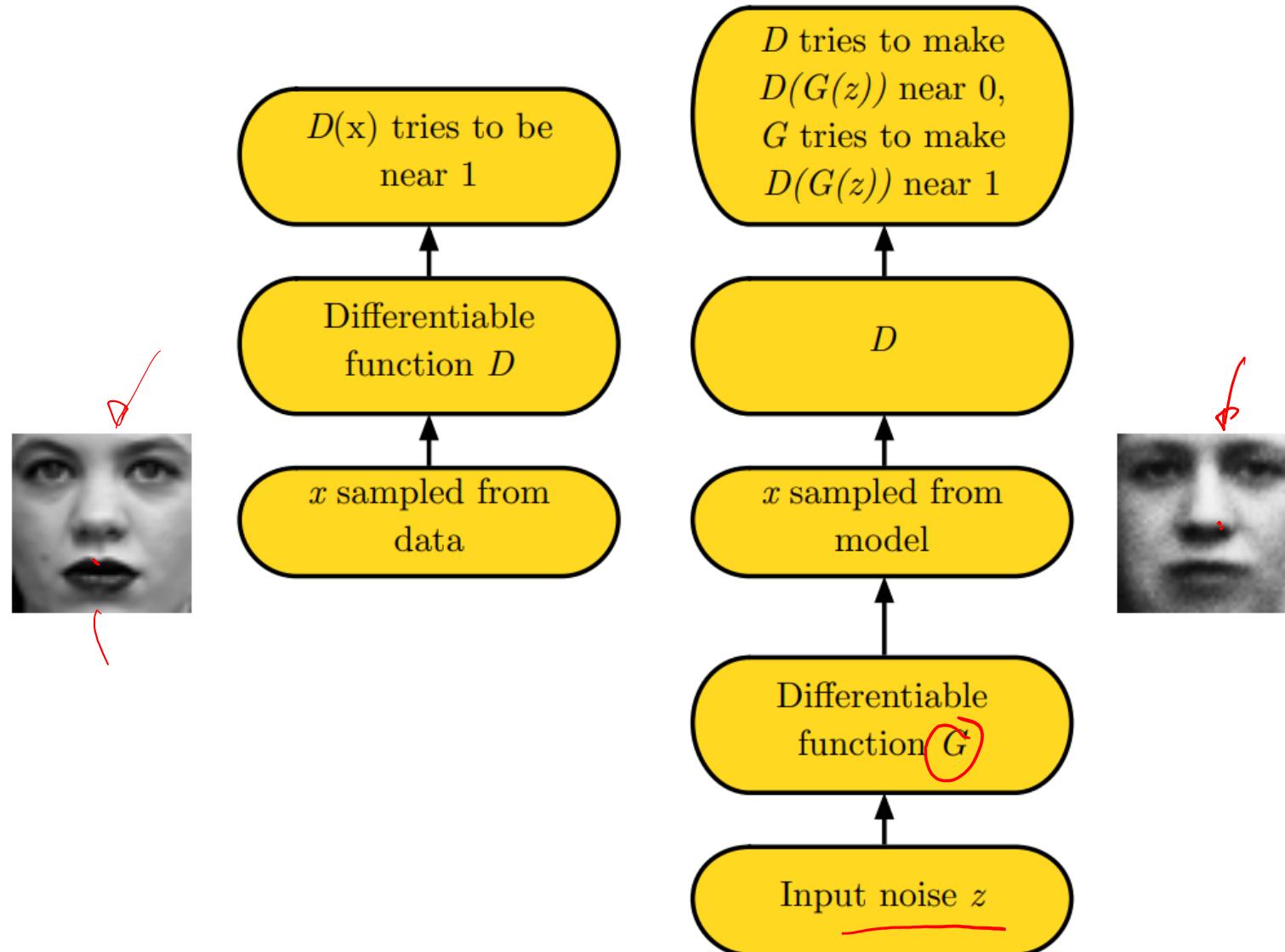


# GAN idea

- Learn to tell plausible (i.e. like data) samples from other ones.

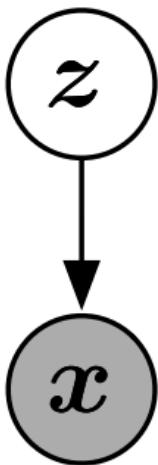


# Adversarial Nets Framework



(Goodfellow 2016)

# Generator Network



$$x = G(z; \theta^{(G)})$$

*q*

↑      ↑

- Must be differentiable
- No invertibility requirement
- Trainable for any size of  $z$
- Some guarantees require  $z$  to have higher dimension than  $x$
- Can make  $x$  conditionally Gaussian given  $z$  but need not do so

(Goodfellow 2016)

# Training Procedure

- Use SGD-like algorithm of choice (Adam) on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples
- Optional: run  $k$  steps of one player for every step of the other player.

# Minimax Game

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} \log D(x) - \frac{1}{2} \mathbb{E}_z \log (1 - D(G(z)))$$
$$J^{(G)} = -J^{(D)}$$

*x - Punkt von u  
1 : f real  
0 : f fikt*

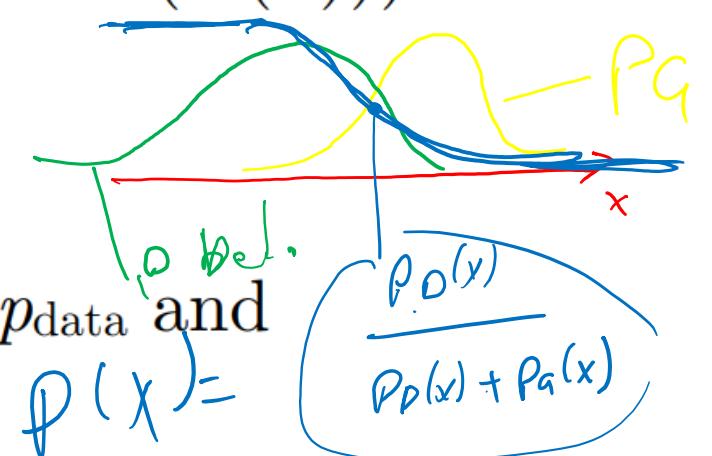
- Equilibrium is a saddle point of the discriminator loss
- Resembles Jensen-Shannon divergence
- Generator minimizes the log-probability of the discriminator being correct

# Exercise 1

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

$$J^{(G)} = -J^{(D)}$$

- What is the solution to  $D(x)$  in terms of  $p_{\text{data}}$  and  $p_{\text{generator}}$ ?
- What assumptions are needed to obtain this solution?



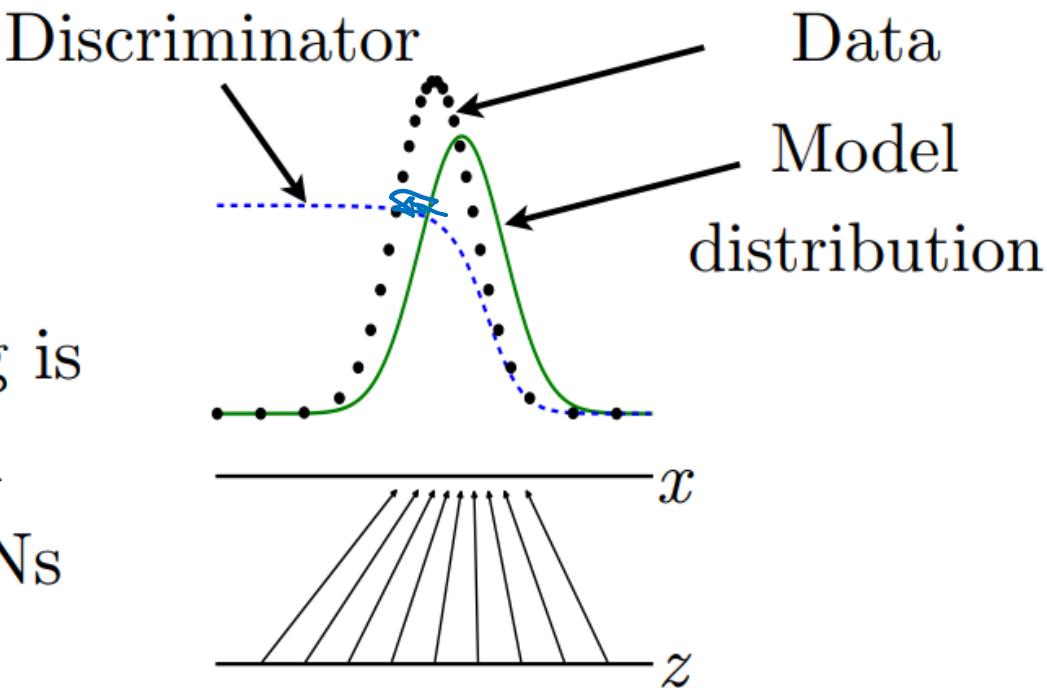
$$x : p_D(x) = p_G(x)$$

# Discriminator Strategy

Optimal  $D(\mathbf{x})$  for any  $p_{\text{data}}(\mathbf{x})$  and  $p_{\text{model}}(\mathbf{x})$  is always

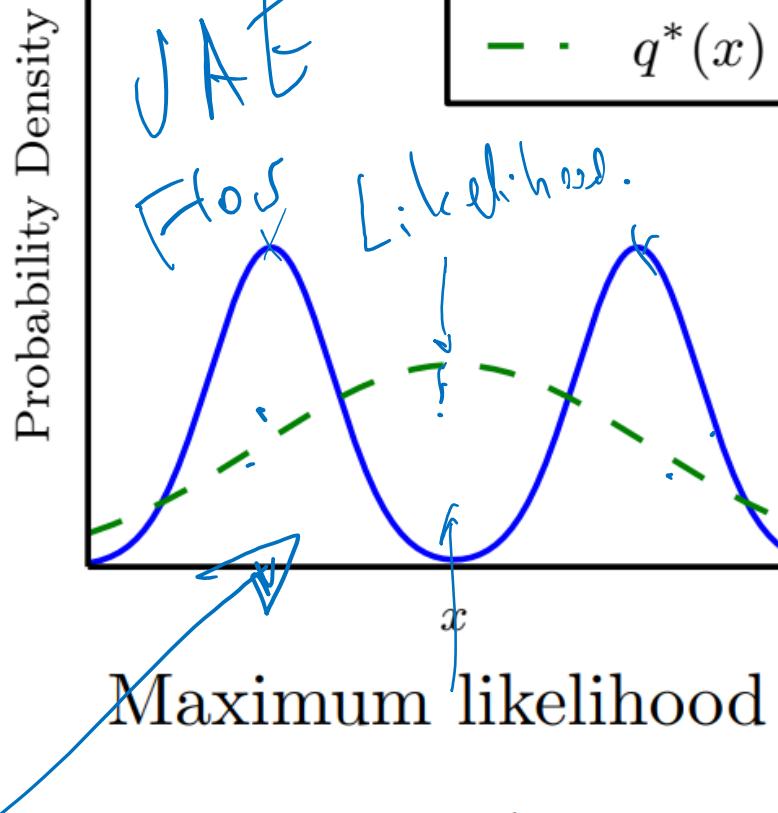
$$D(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}$$

Estimating this ratio  
using supervised learning is  
the key approximation  
mechanism used by GANs

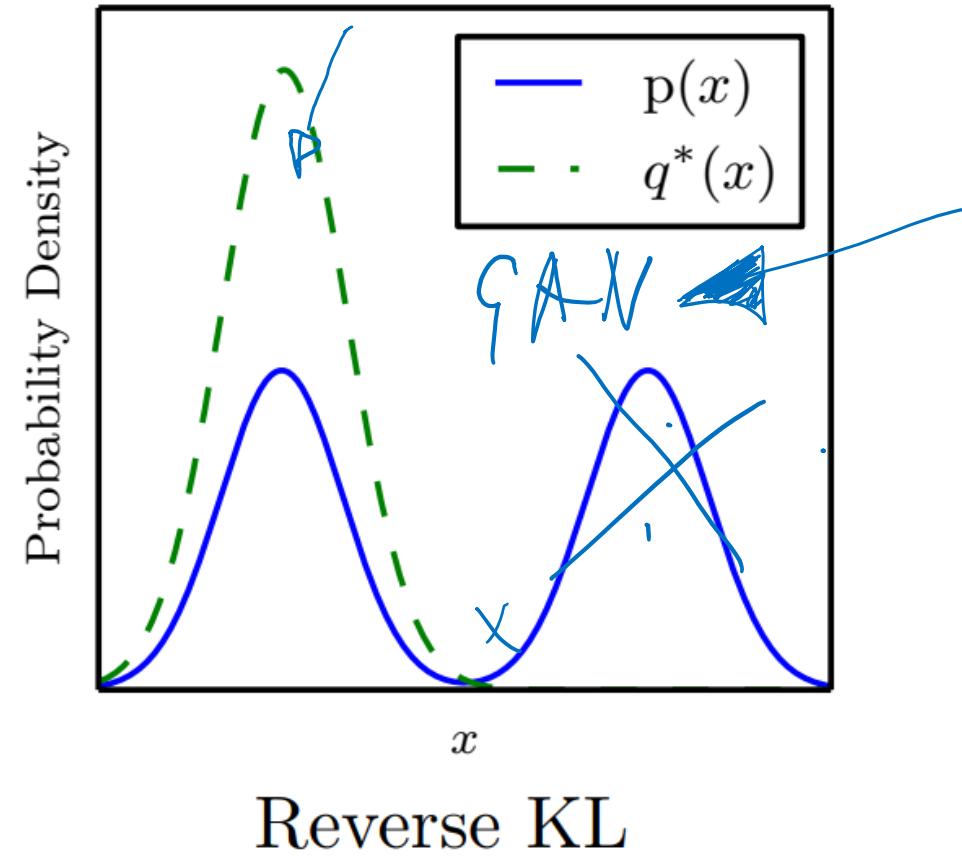


# Is the divergence important?

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$



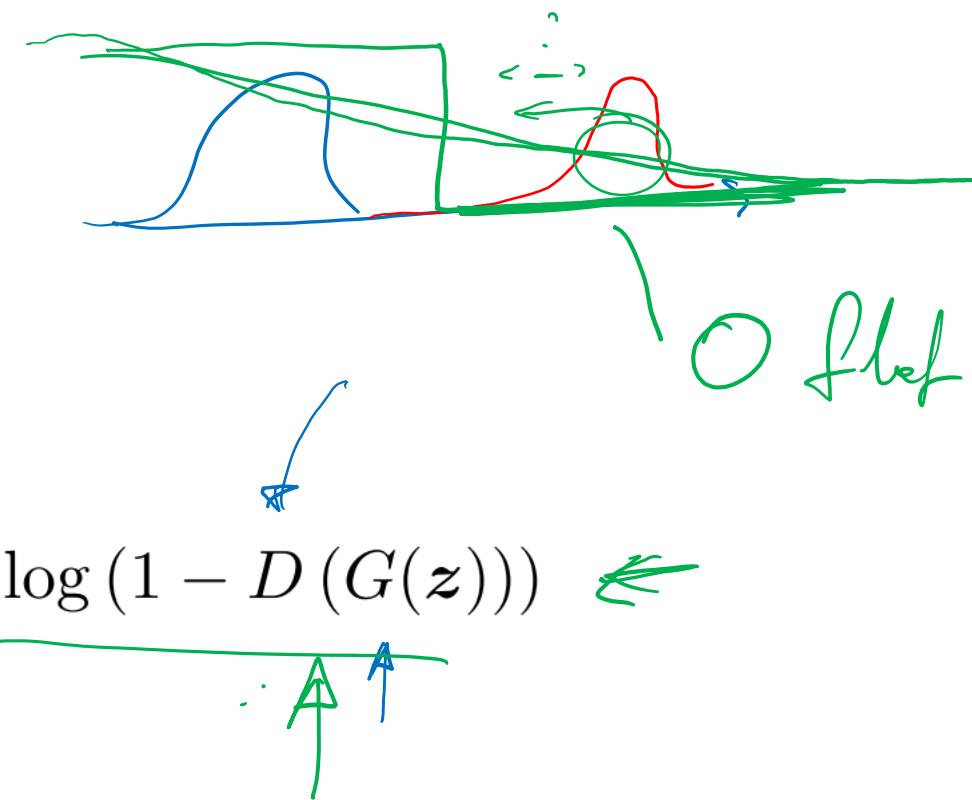
(Goodfellow et al 2016)

# GAN Trick

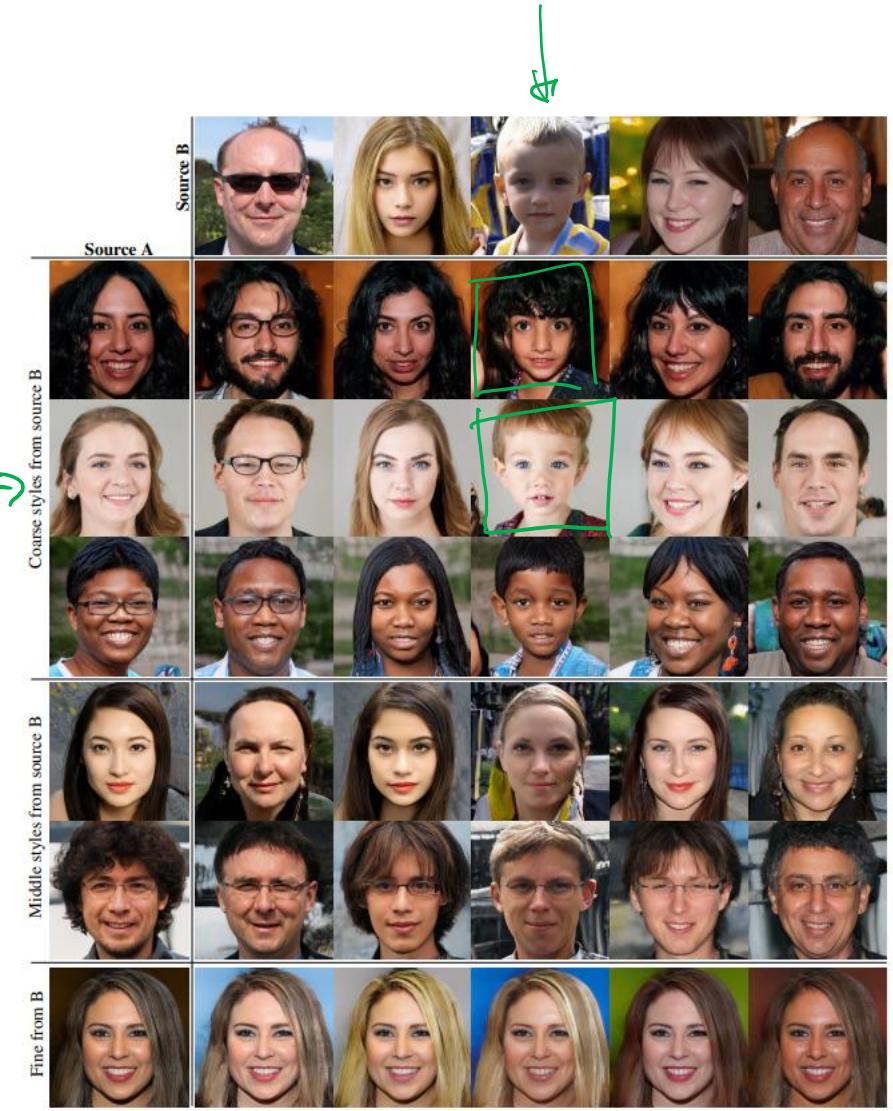
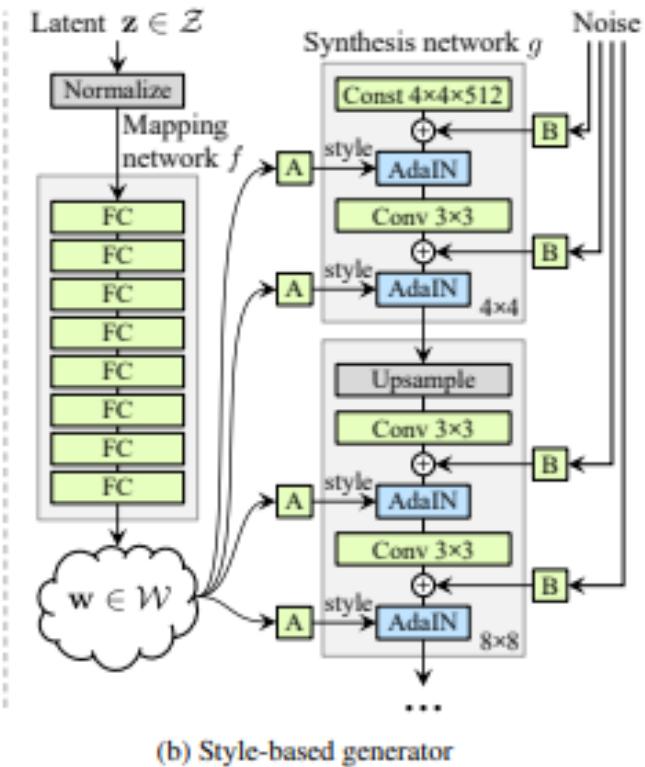
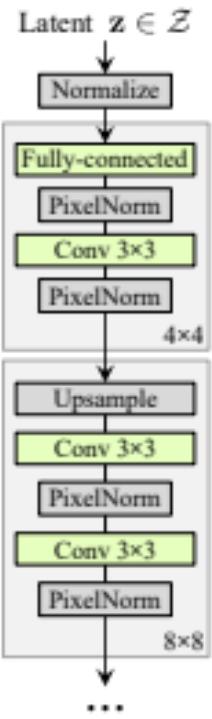
Non-saturating loss

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

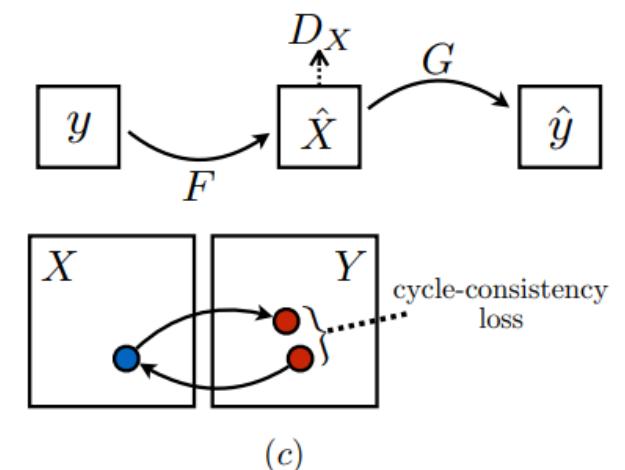
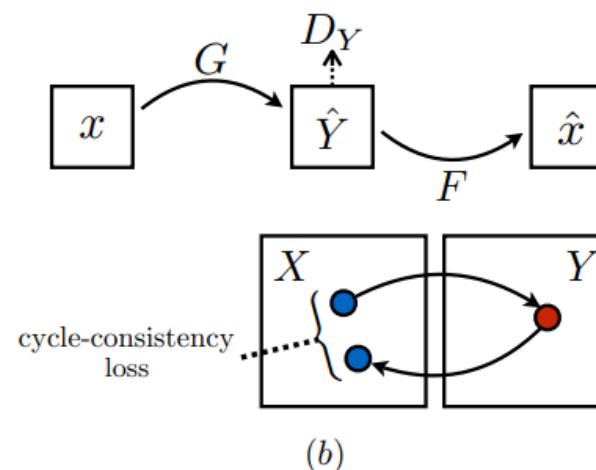
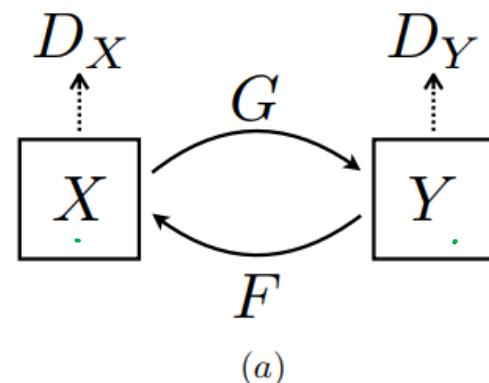
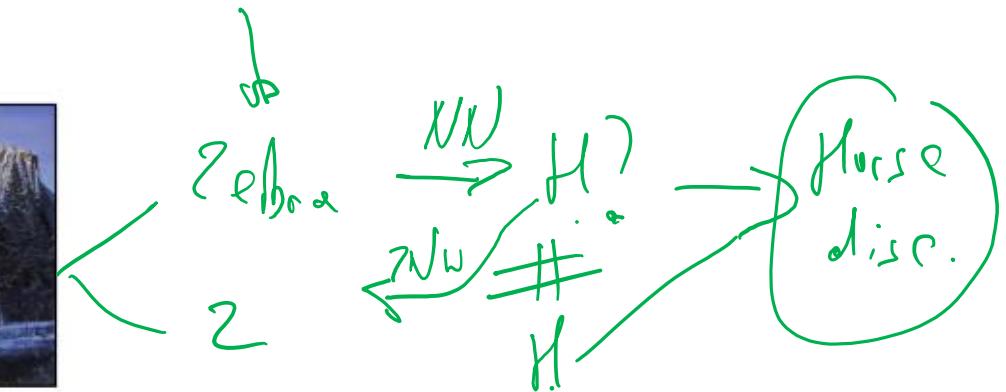
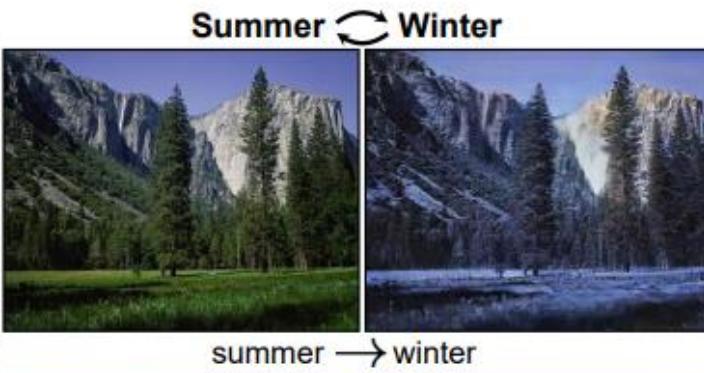
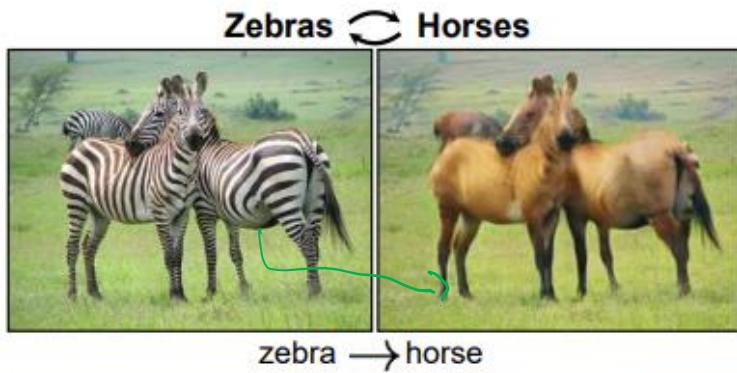
$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$



# GAN Uses: image generation

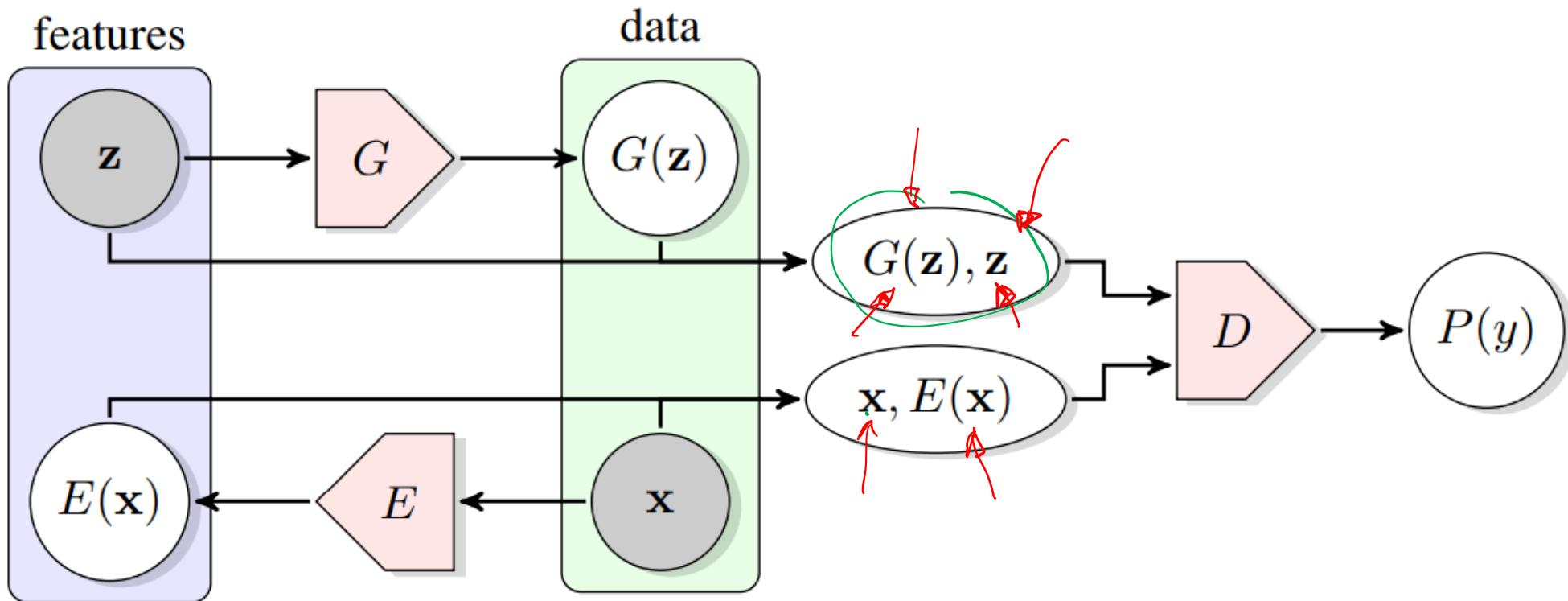


# GAN Uses: unpaired image-image translation



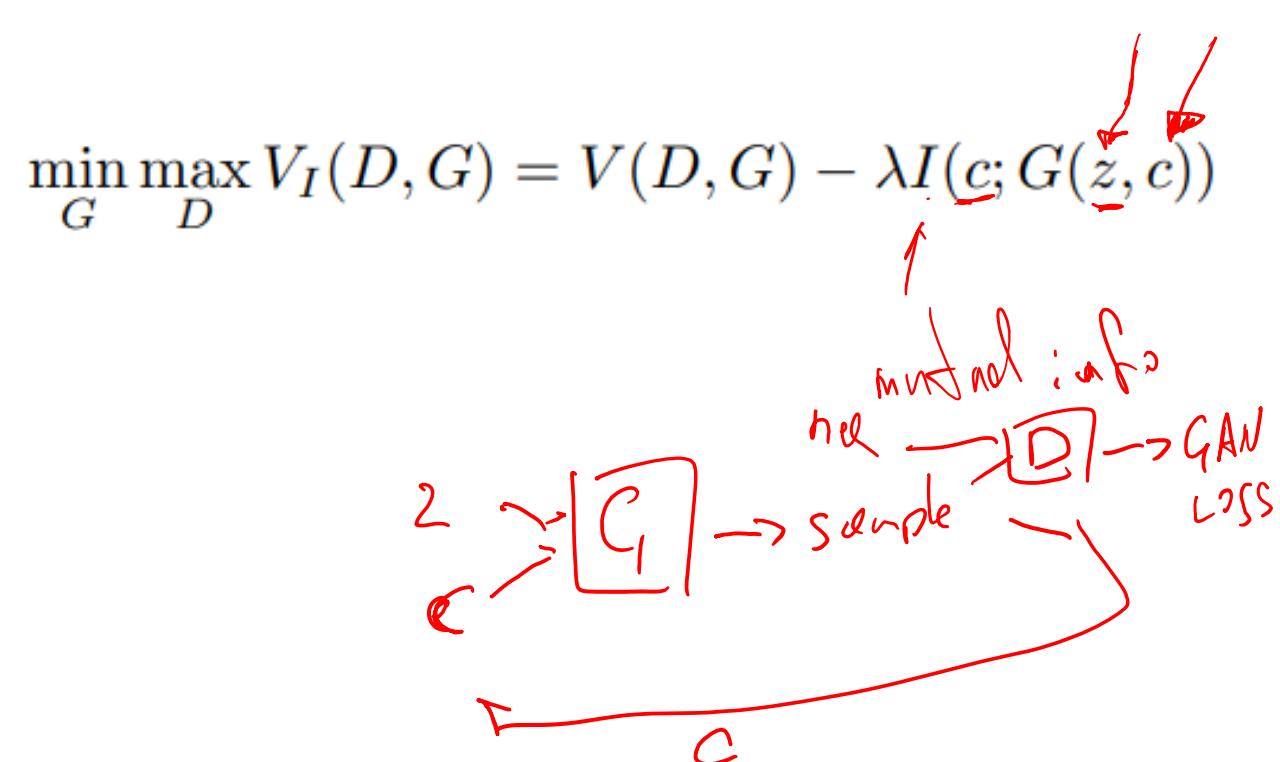
# GANs and Latent Variables 1/3

ADVERSARILY LEARNED INFERENCE (<https://arxiv.org/pdf/1606.00704.pdf>)  
aka BiGAN (<https://arxiv.org/abs/1605.09782>)



# GANs and Latent Variables 2/3

InfoGAN: <https://arxiv.org/pdf/1606.03657.pdf>



	(a) Varying $c_1$ on InfoGAN (Digit type)
	(b) Varying $c_1$ on regular GAN (No clear meaning)

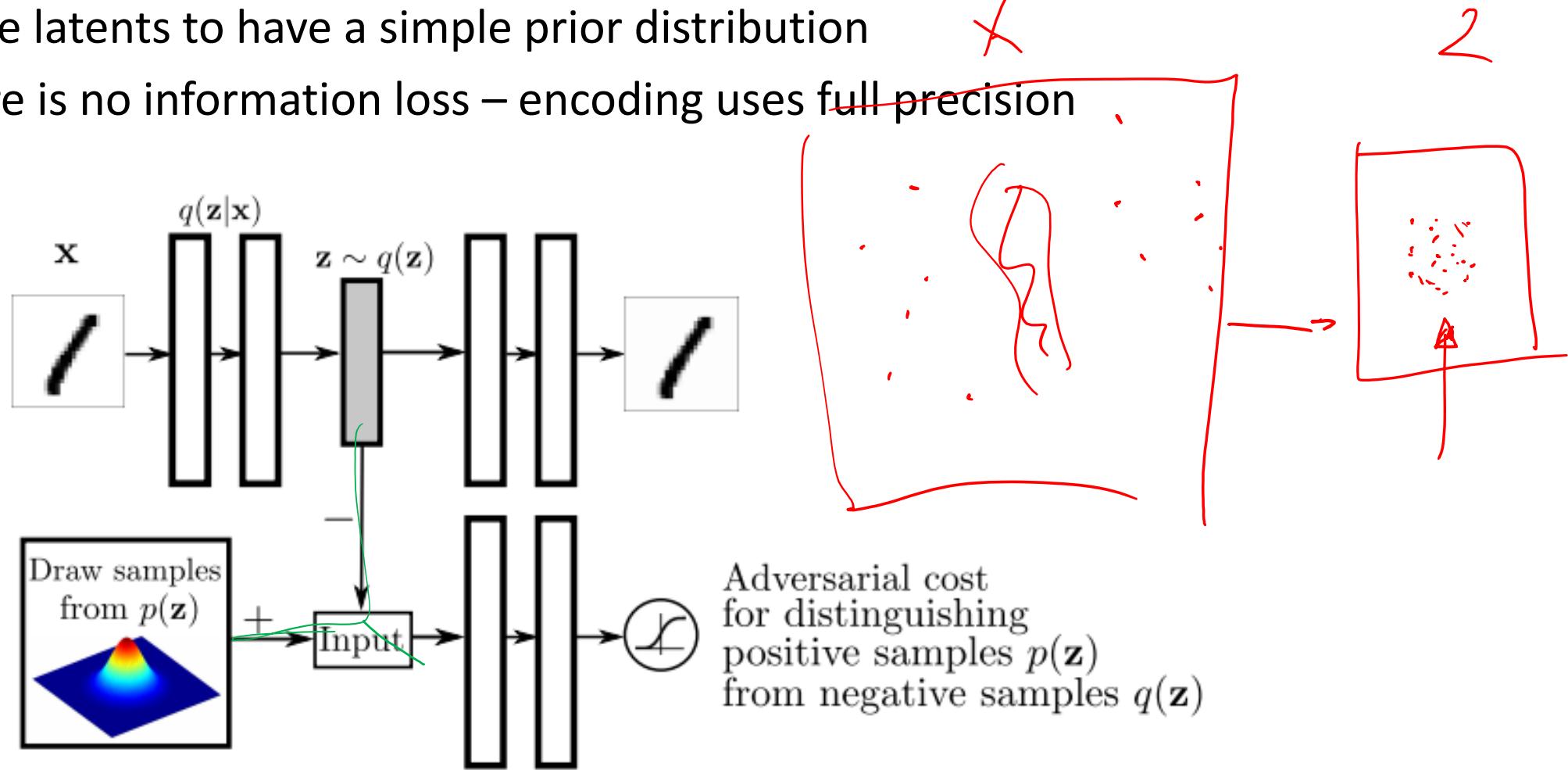
	(c) Varying $c_2$ from $-2$ to $2$ on InfoGAN (Rotation)
	(d) Varying $c_3$ from $-2$ to $2$ on InfoGAN (Width)

# GANs and Latent Variables 3/3

Adversarial Autoencoders (<https://arxiv.org/pdf/1511.05644.pdf>)

Like VAE, enforce latents to have a simple prior distribution X

Unlike VAE, there is no information loss – encoding uses ~~full precision~~ 2



$x$   $x$   $\rightarrow$  batch codes  $\rightarrow$  solve new tasks.  
 $x$   $y$

## **THE FUTURE OF UNSUPERVISED: SELF-SUPERVISED LEARNING**

for self  
learning

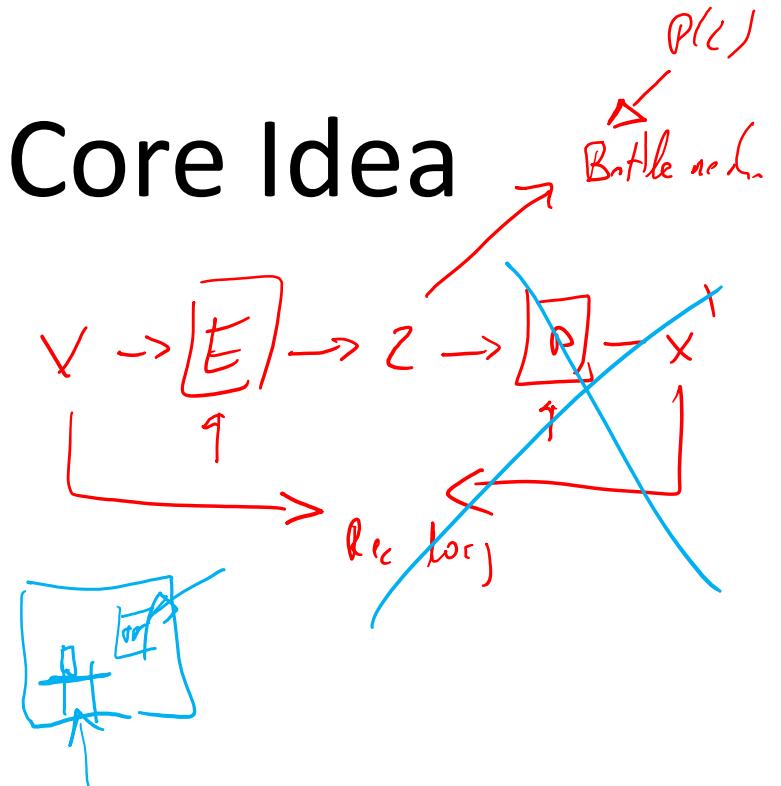
# Self-supervised learning: Core Idea

We usually train an Encoder-Decoder pair

But later, we only use/keep the encoder!!

This is wasteful – the decoder is often complex (and may even be more expensive to train than the encoder)

Q: Can we train an encoder-only system?



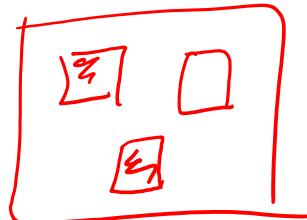
# Self-supervised learning: Core Idea

Q: Can we train an encoder-only system?

Yes, then use its output to solve a task requiring data understanding, for which labels are easy to get.

E.g.:

- cut image into parts, predict their relative location.
- cut video into parts, predict time direction ↪
- cut sentence into words, predict neighbors  
(quick, what was the name of this model??) → Word & Vec
- Take a sentence, mask some words, train the model to impute them  
(again – what was the name of the model?) ← BERT



# Self-supervised learning: good proxy tasks

Some domains have known invariances

E.g. images: rotation, scaling, recolorization, mild crop, etc.

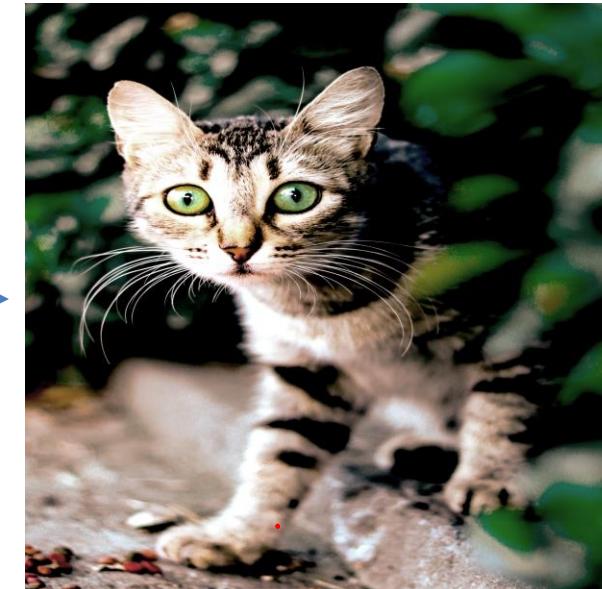
only slightly change the contents of an image



Random:

- Color transformation
- Affine transform
- Crop

Easy to implement,  
No labels needed!



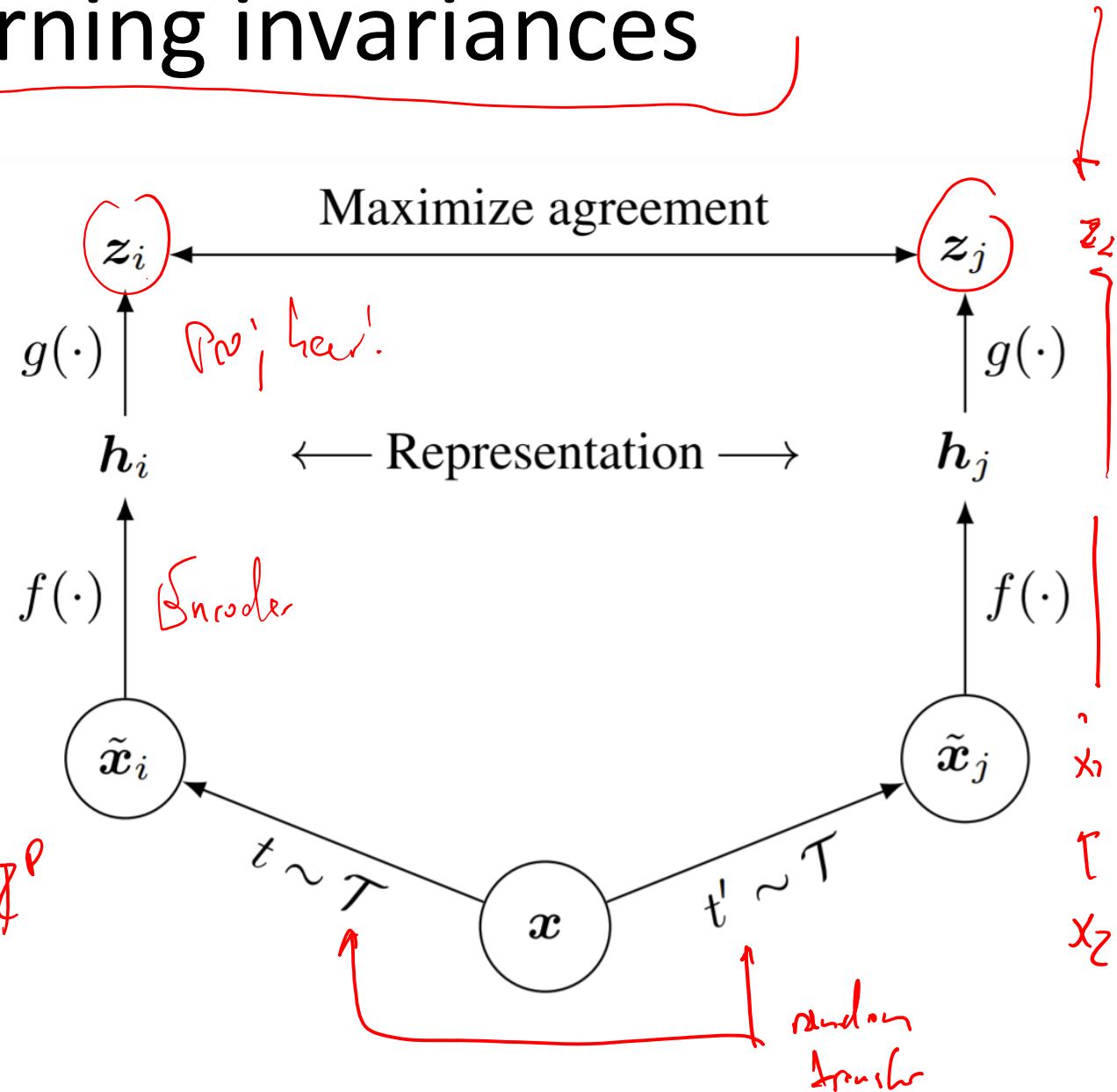
Q: how to teach the model to be invariant?

# SimCLR: Self-learning invariances

Idea:

- Maximize agreement between two augmentations of the same input
- Q: what prevents the model from learning  $f(x) = 0$ ?

$$f(x) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

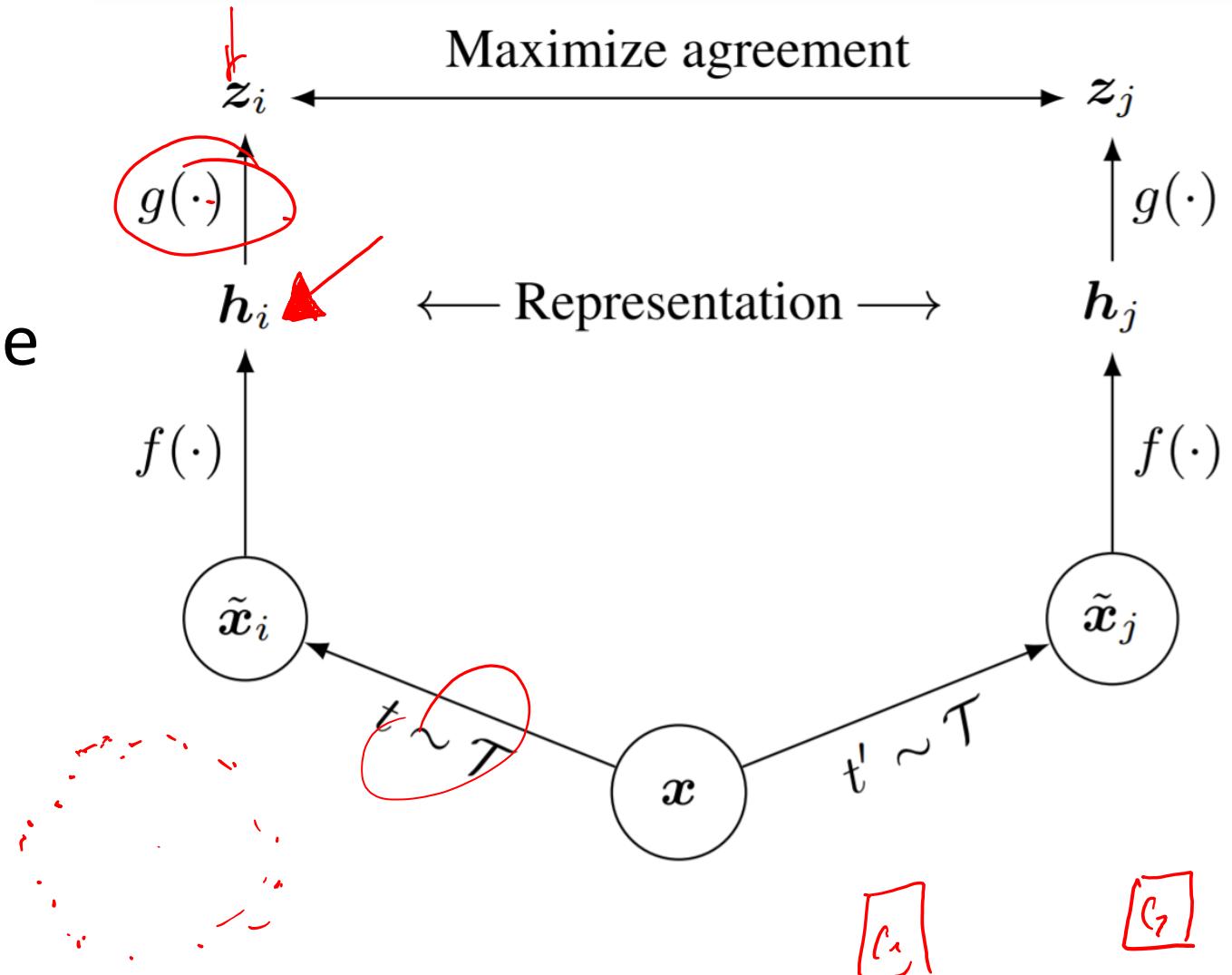


# SimCLR: contrastive losses

Q: what prevents the model from learning  $f(x) = 0$ ?

A: the net is trained to pull the two augmented variants of  $x$  together...

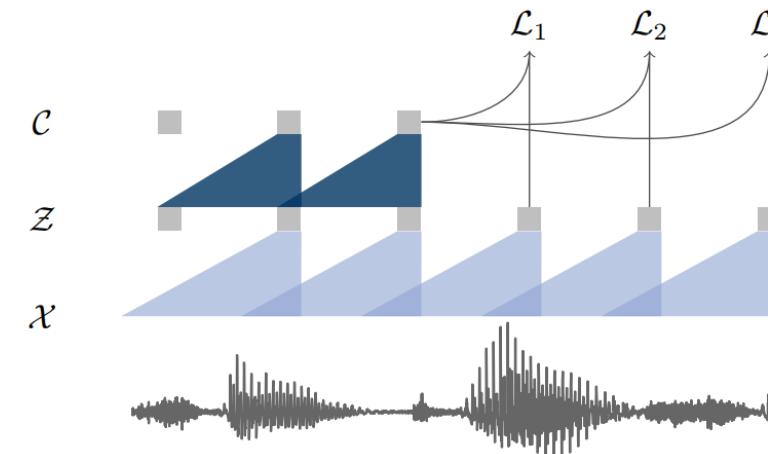
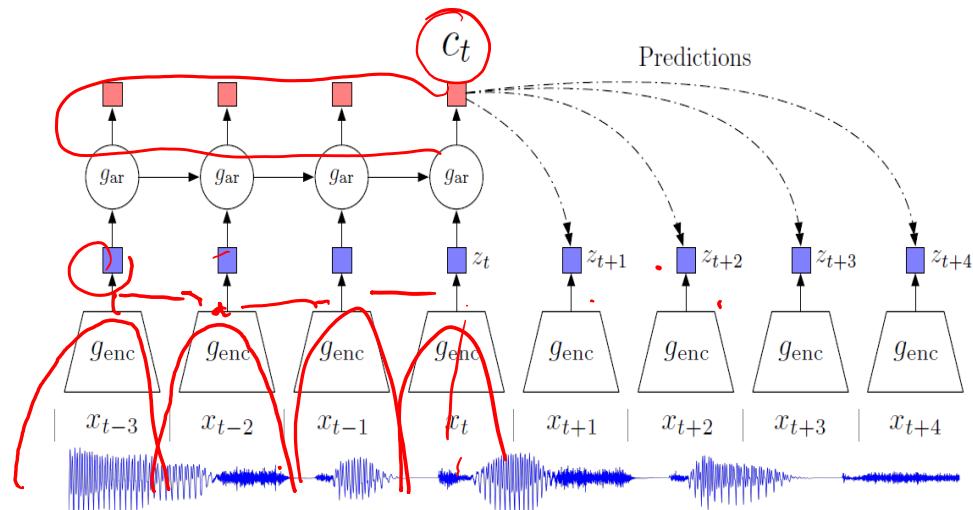
And push them away from other samples (contrastive learning)



# CPC: contrastive training to predict the future

- Encode speech frames into latent vectors
- Combine past frames using  $g_{ar}$
- Predict upcoming frames using a contrastive loss (tell apart the next frame from a randomly sampled one)

$$\begin{cases} X \in \mathbb{R}^{T \times D} \\ \text{1 conv} \\ Z \in \mathbb{R}^{T' \times D'} \end{cases}$$



$$P_2(c_t) \approx \mathbb{I}_{t+2}$$

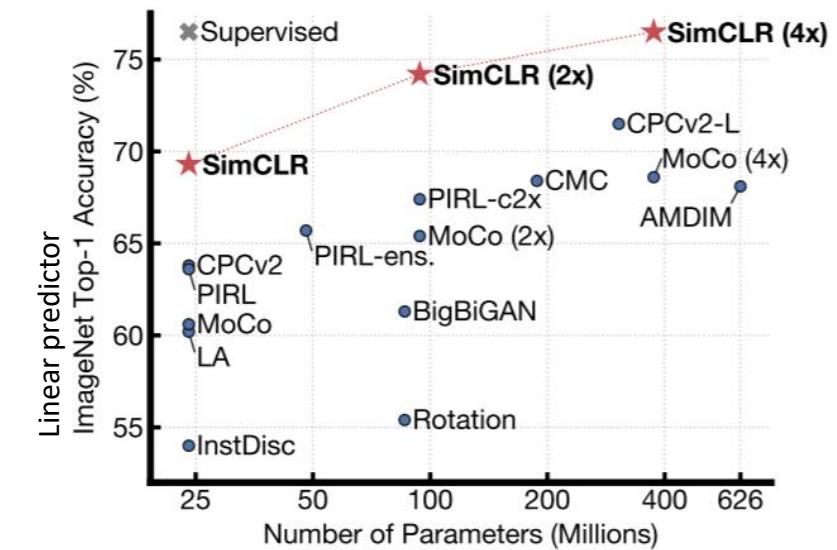
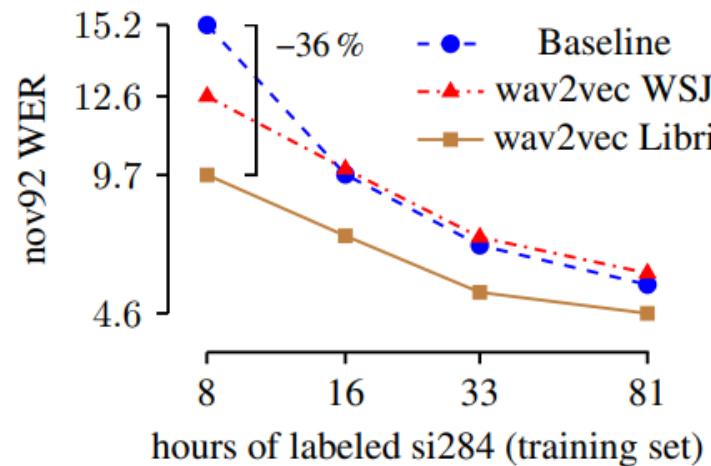
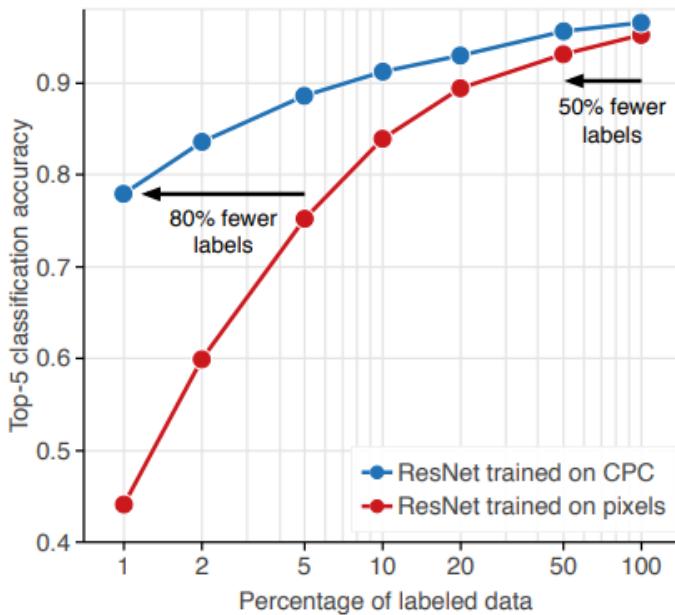
$$P_i(p_t) \neq \mathbb{I}_{t+k} \quad k \neq 2$$

$$P_k(c_t) \approx \mathbb{I}_{t+k}$$

$$P_k(c_t) \neq \mathbb{I}_{t+j} \quad j \neq k$$

Oord et al. „Representation Learning with Contrastive Predictive Coding“  
 Schneider et al. „wav2vec: Unsupervised Pre-training for Speech Recognition“

# Self-supervised learning becomes practical



DATA-EFFICIENT IMAGE  
RECOGNITION WITH  
CONTRASTIVE PREDICTIVE  
CODING Olivier J. Henaff,  
Aravind Srinivas, Jeffrey De  
Fauw, Ali Razavi, Carl Doersch,  
S. M. Ali Eslami, Aaron van den  
Oord

WAV2VEC: UNSUPERVISED  
PRE-TRAINING FOR SPEECH  
RECOGNITION Steffen  
Schneider, Alexei Baevski,  
Ronan Collobert, Michael Auli

A Simple Framework for  
Contrastive Learning of Visual  
Representations  
Ting Chen, Simon Kornblith,  
Mohammad Norouzi,  
Geoffrey Hinton

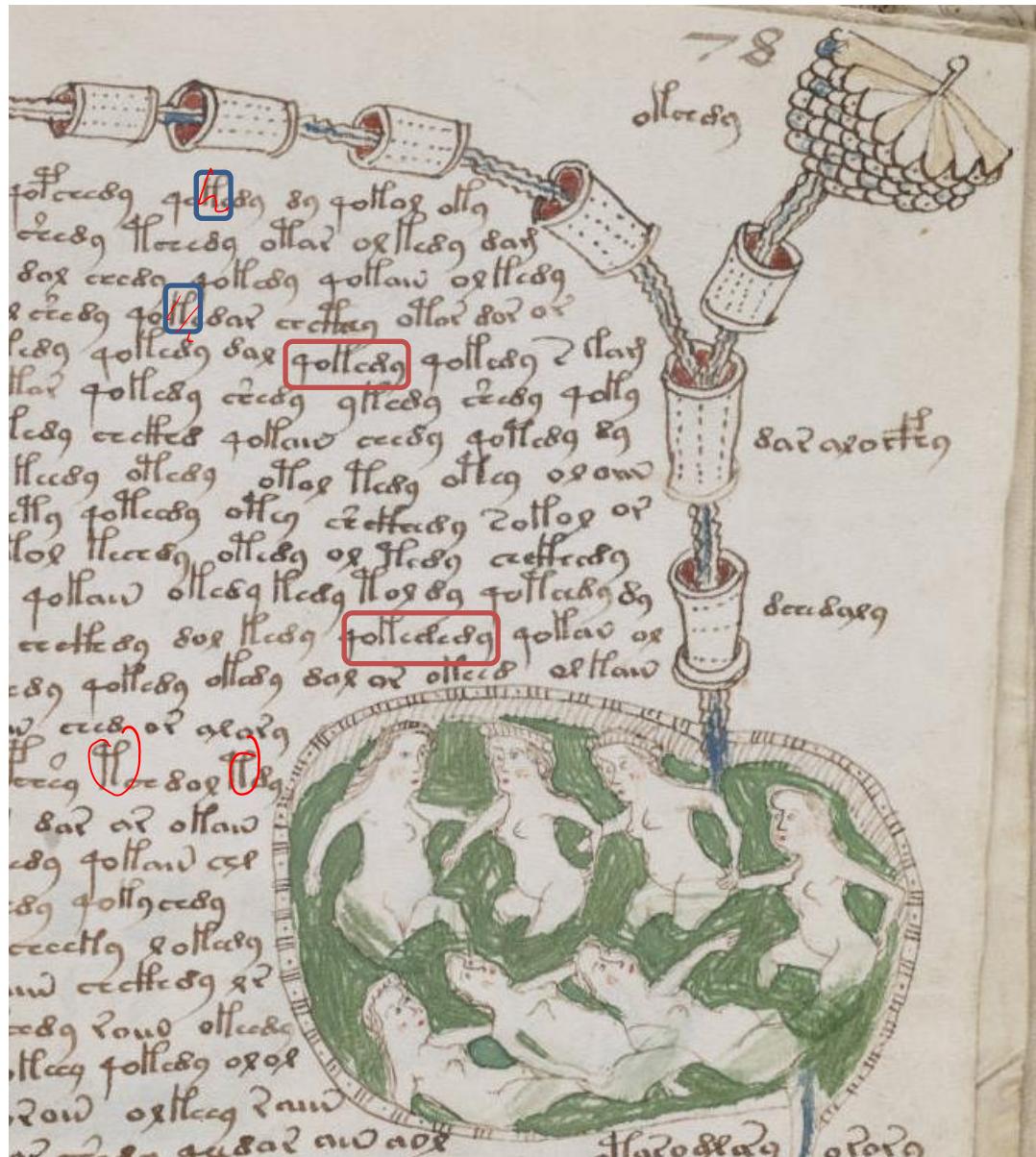
# Summary for unsupervised learning

- We can generate nearly realistic samples with deep generative models!
- We can use unsupervised / self-supervised learning to learn good data representations
- We can approximately control the contents of the latent representation (e.g. separate phone/gender)

# **MY WORK:**

## **VAE + AUTOREGRESSIVE DECODERS**

# Unlabeled data uses

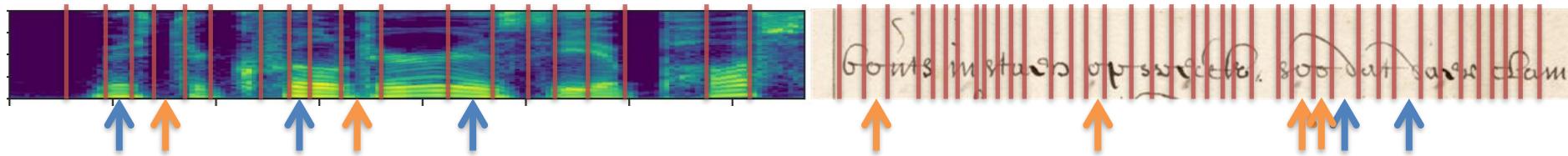


What can we learn  
without labels:

- Features (good data representations)
- Units (characters)
- Word types (sequences of units)
- Their co-occurrence patterns (structure, e.g. bigram stats)

# The goal

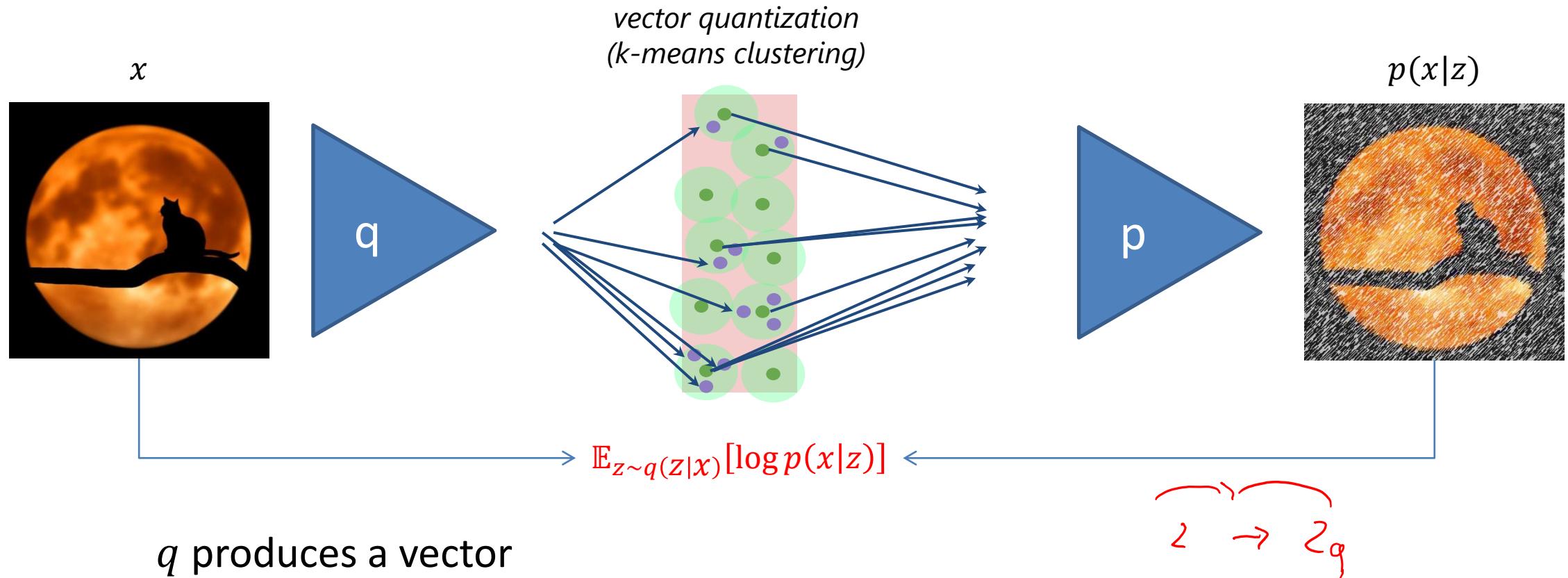
Unsupervised unit discovery  
in speech and handwriting



1. Discover information-bearing units (cluster input segments)
2. Be invariant to other factors of variation (speaker/scribe, style, etc)

Extreme case of representation learning!  
We will combine autoencoders with generative models.

# The Vector Quantized Autoencoder



$q$  produces a vector

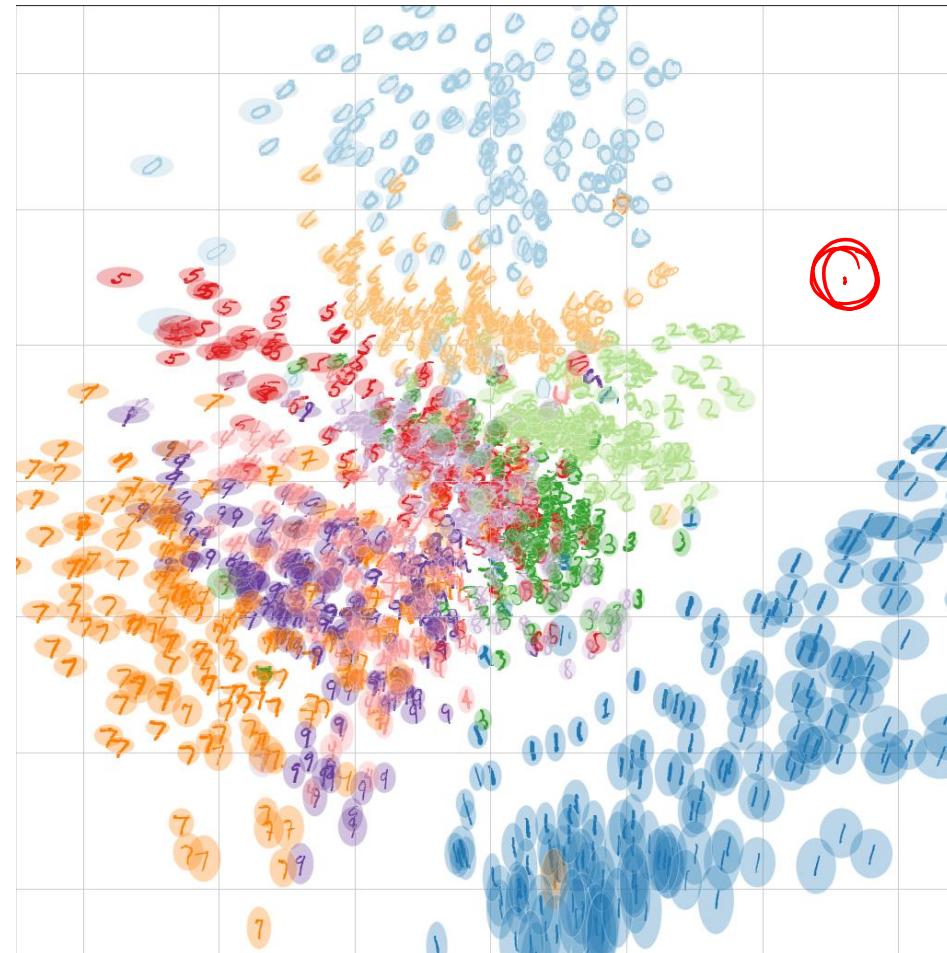
it is replaced by one of K prototypes

Information loss through quantization (rounding)

# Recall: in VAE each $z$ has a volume

Each sample is represented as a Gaussian

This discards information  
(latent representation has low precision)



# VQVAE – deterministic quantization, each z is rounded

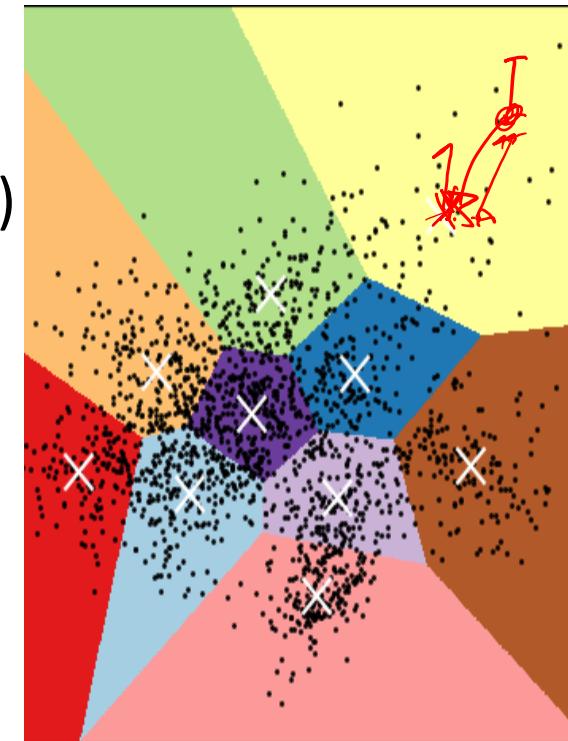
Limit precision of the encoding by quantizing (round each vector to a nearest prototype).

Output can be treated:

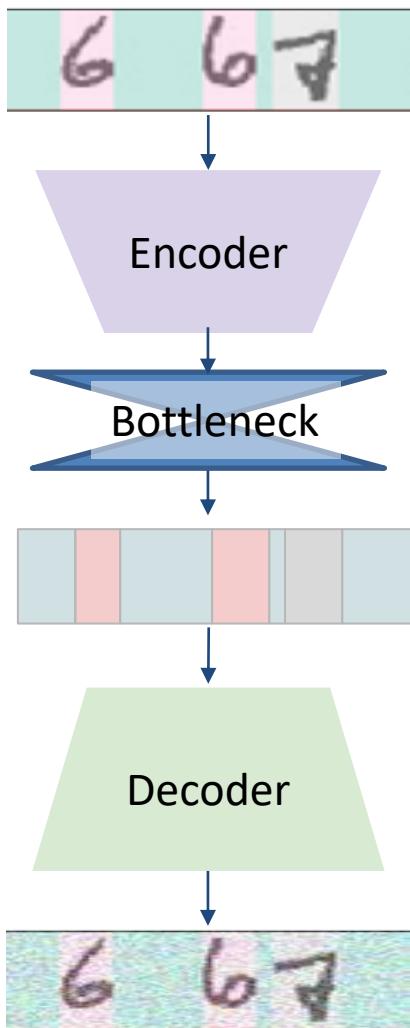
- As a sequence of discrete prototype ids (tokens)
- As a distributed representation (the prototypes themselves)

Train using the straight-through estimator,  
with auxiliary losses:

$$\begin{aligned}\mathcal{L} = & \log p(x | z_q(x)) \\ & + \| \text{sg}(z_e(x)) - e_{q(x)} \|_2^2 + \gamma \| z_e(x) - \text{sg}(e_{q(x)}) \|_2^2\end{aligned}$$



# Autoencoders and missing information



Detailed input:  
handwriting or speech

Encoder with bottleneck removes  
information, enforces segmentation,  
and quantizes the representation

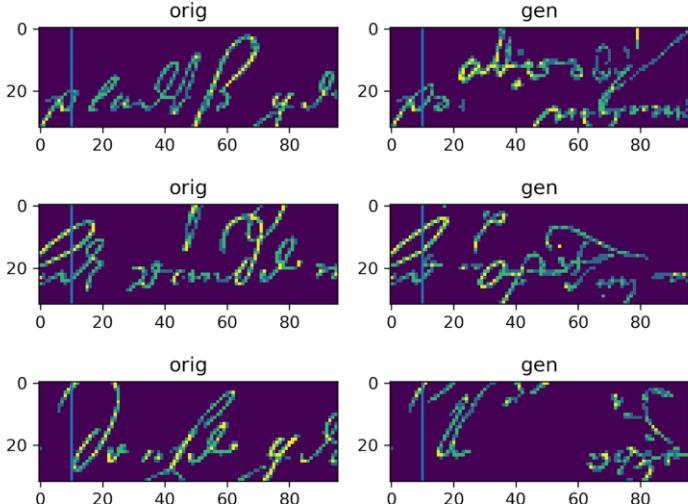
Latent encoding:  
only high level features

Decoder regenerates the inputs.

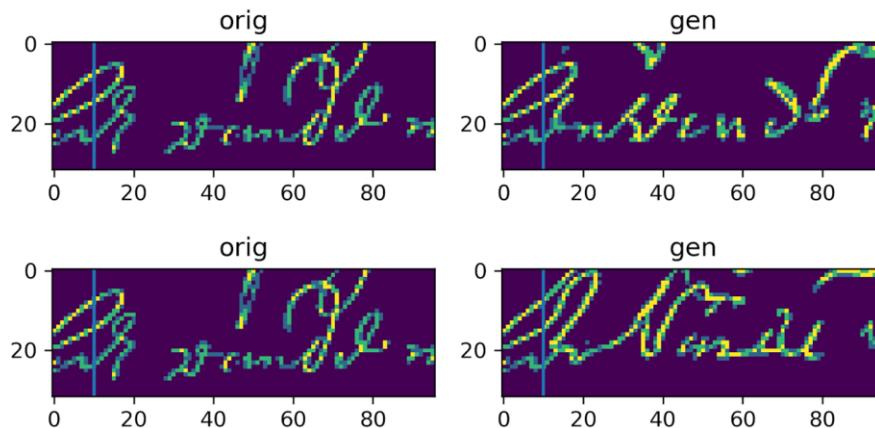
**It learns  $p(x|z)$  and it must fill in all discarded information.**

# Conditional PixelCNN can regenerate images.

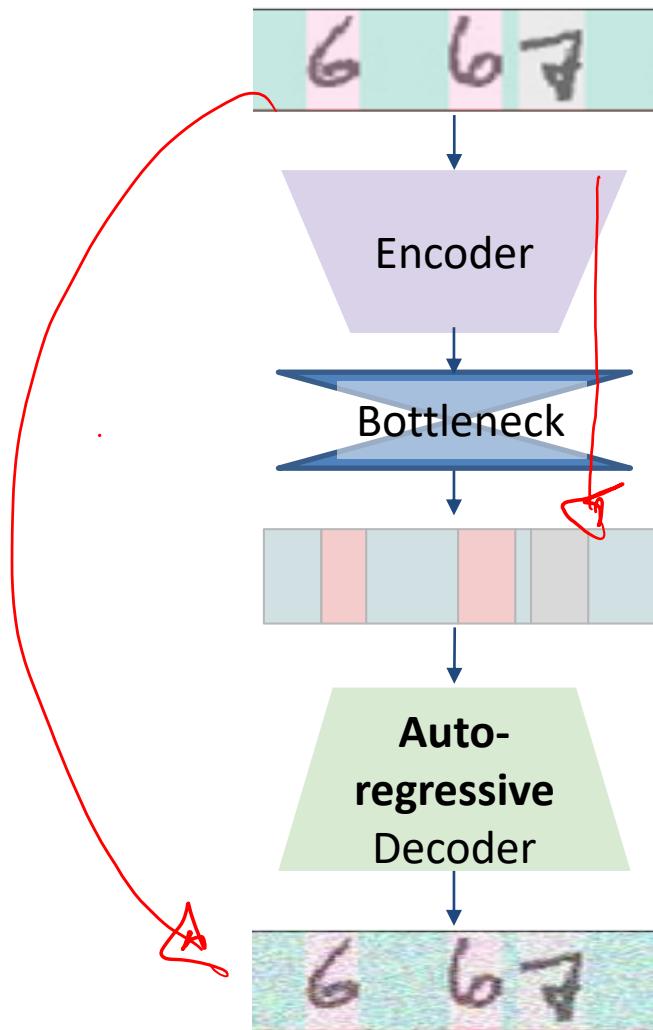
Unconditional



Conditioned on text



# Our model: Autoencoder + autoregressive decoder

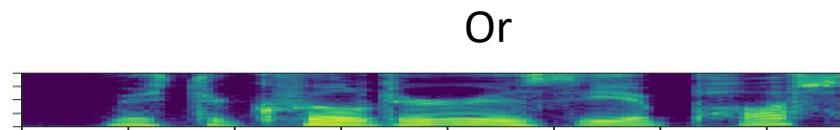
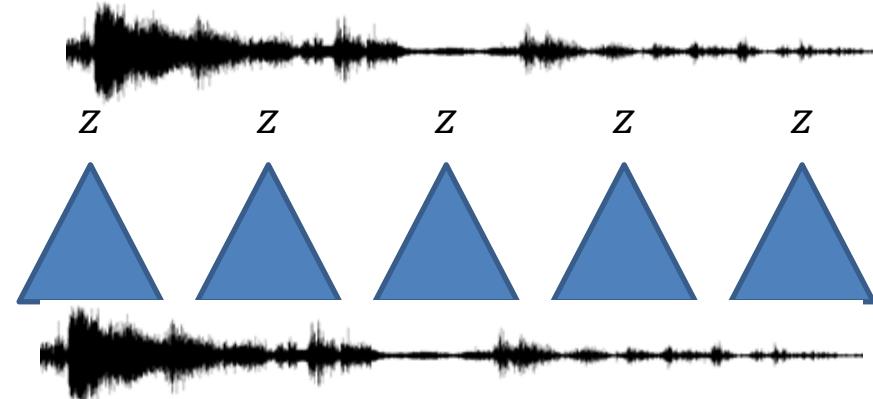
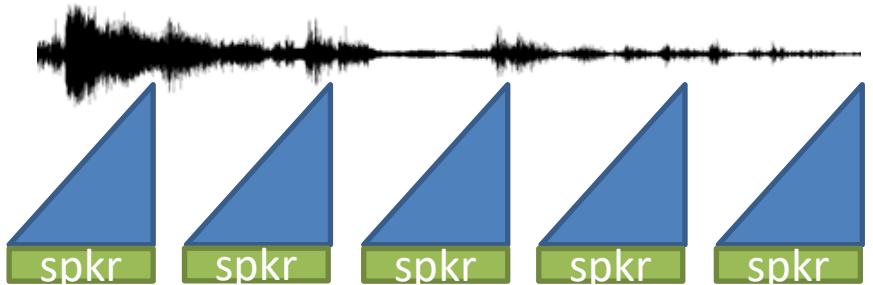


Decoder computes a probability over input images/waveforms conditioned on the encoding.

It combines information from:

- latent representation  
(know what to reconstruct)
- neighboring sounds/pixels  
(recover details)
- other information (e.g. speaker id) ↗  
*Scene*

# Experiments on Speech

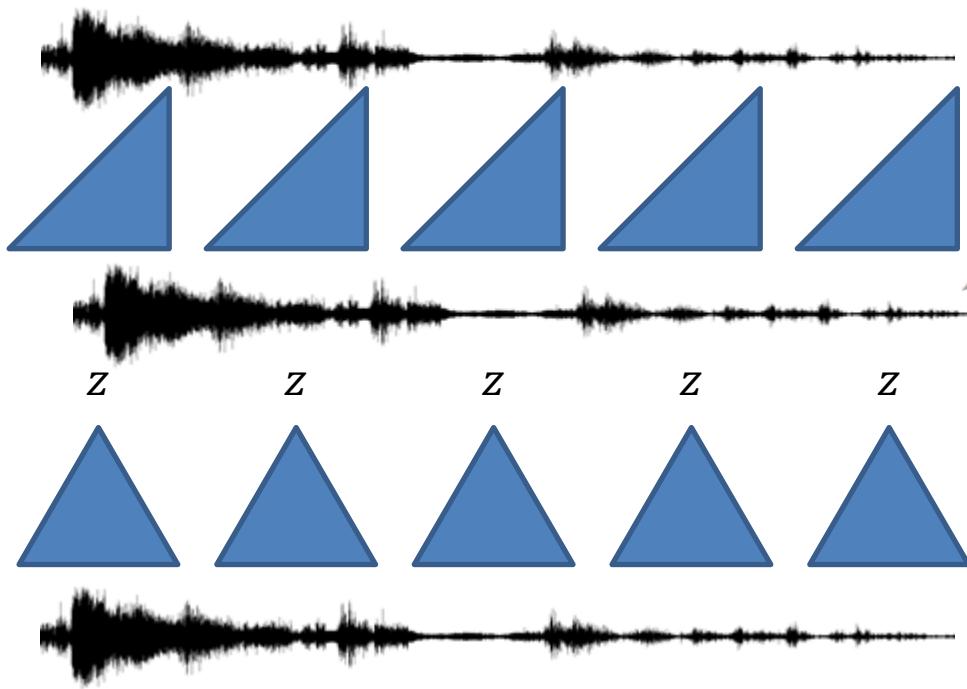


Input:

Waveforms, Mel Filterbanks, MFCCs

# Methodology: probing points

We have inserted probing classifiers at 4 points in the network:



$p_{cond}$ : several  $z$  codes mixed using a convolution.  
The wavenet uses it for conditioning

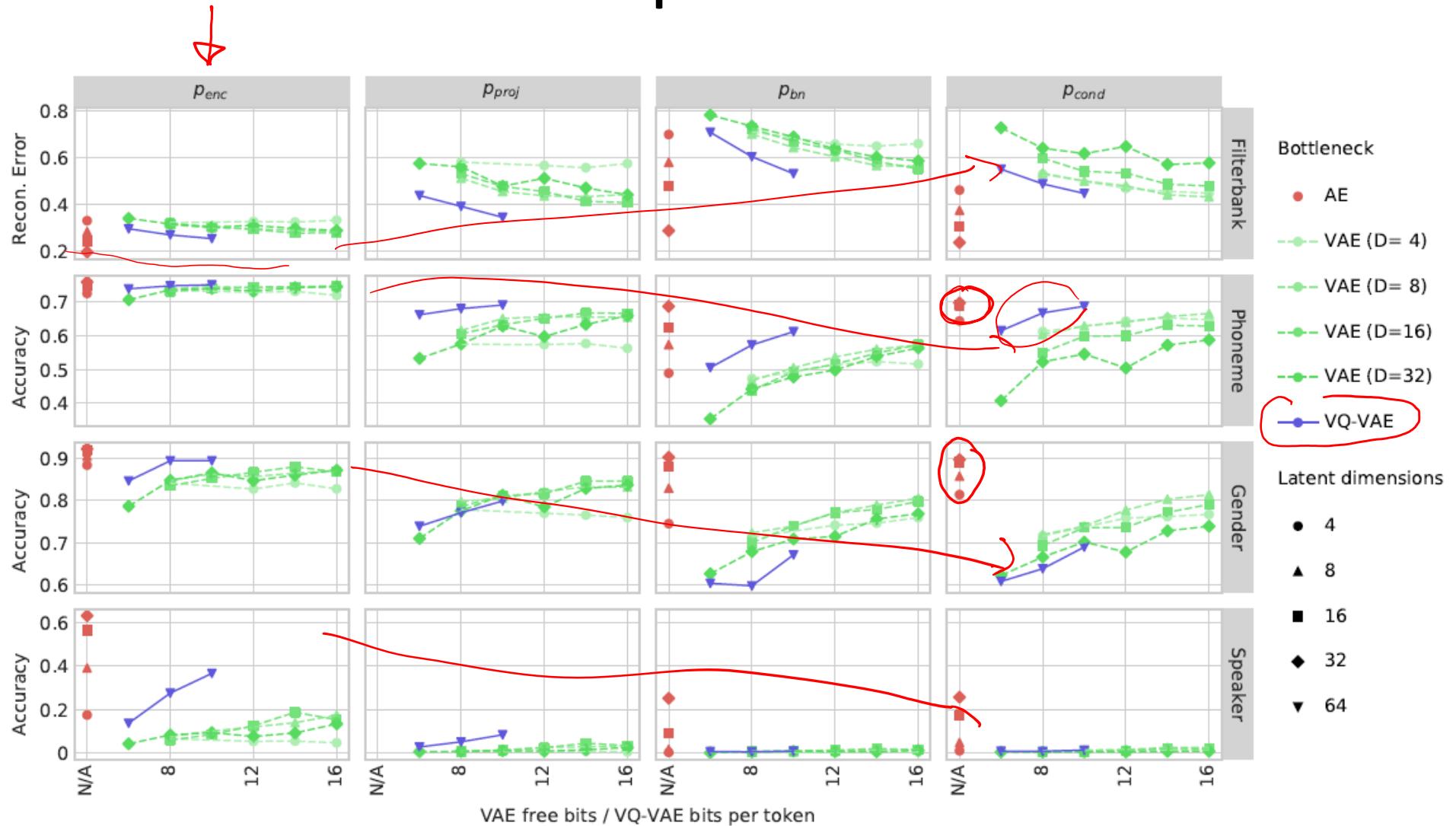


$p_{bn}$ : the latent codes

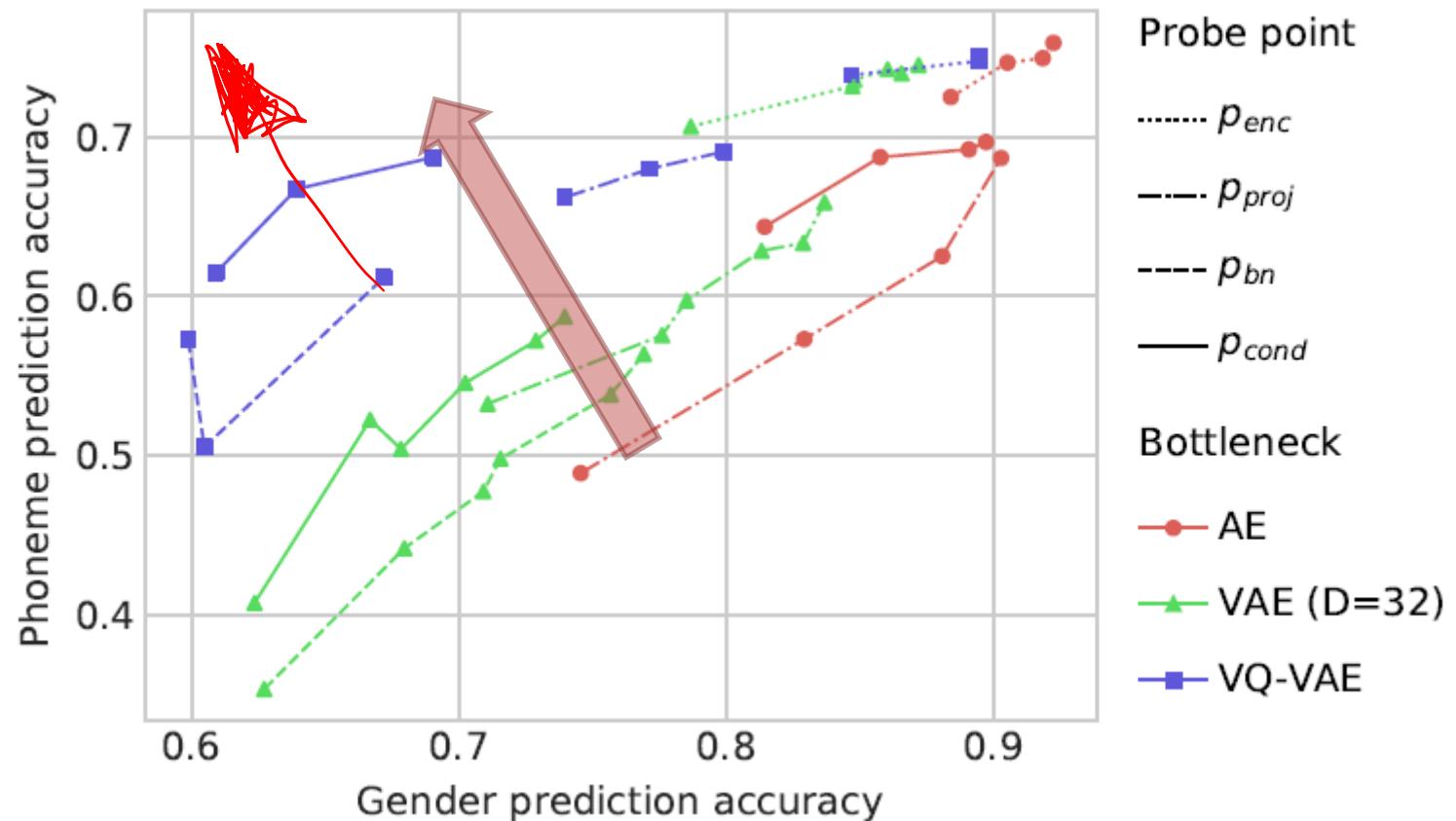
$p_{proj}$ : low dimensional representation input to the bottleneck layer

$p_{enc}$ : high dimensional representation coming out of the encoder

# What information is captured in the latent codes?



# Phonemes vs Gender tradeoff



# Performance on ZeroSpeech 2017 unit discovery

Model	Within-speaker									Across-speaker								
	English (45h)			French (24h)			Mandarin (2.4h)			English (45h)			French (24h)			Mandarin (2.4h)		
	1s	10s	2m	1s	10s	2m	1s	10s	2m	1s	10s	2m	1s	10s	2m	1s	10s	2m
Unsupervised baseline	12.0	12.1	12.1	12.5	12.6	12.6	11.5	11.5	11.5	23.4	23.4	23.4	25.2	25.5	25.2	21.3	21.3	21.3
Supervised topline	6.5	5.3	5.1	8.0	6.8	6.8	9.5	4.2	4.0	8.6	6.9	6.7	10.6	9.1	8.9	12.0	5.7	5.1
VQ-VAE (per lang, $p_{\text{cond}}$ )	<b>5.6</b>	<b>5.5</b>	<b>5.5</b>	<b>7.3</b>	<b>7.5</b>	<b>7.5</b>	11.2	10.7	10.8	<b>8.1</b>	<b>8.0</b>	<b>8.0</b>	<b>11.0</b>	<b>10.8</b>	<b>11.1</b>	12.2	11.7	11.9
Heck et al. [57]	6.9	6.2	6.0	9.7	8.7	8.4	<b>8.8</b>	<b>7.9</b>	<b>7.8</b>	10.1	8.7	8.5	13.6	11.7	11.3	<b>8.8</b>	<b>7.4</b>	<b>7.3</b>
Chen et al. [58]	8.5	7.3	7.2	11.2	9.4	9.4	10.5	8.7	8.5	12.7	11.0	10.8	17.0	14.5	14.1	11.9	10.3	10.1
Ansari et al. [59]	7.7	6.8	N/A	10.4	N/A	8.8	10.4	9.3	9.1	13.2	12.0	N/A	17.2	N/A	15.4	13.0	12.2	12.3
Yuan et al. [60]	9.0	7.1	7.0	11.9	9.5	9.5	11.1	8.5	8.2	14.0	11.9	11.7	18.6	15.5	14.9	12.7	10.8	10.7

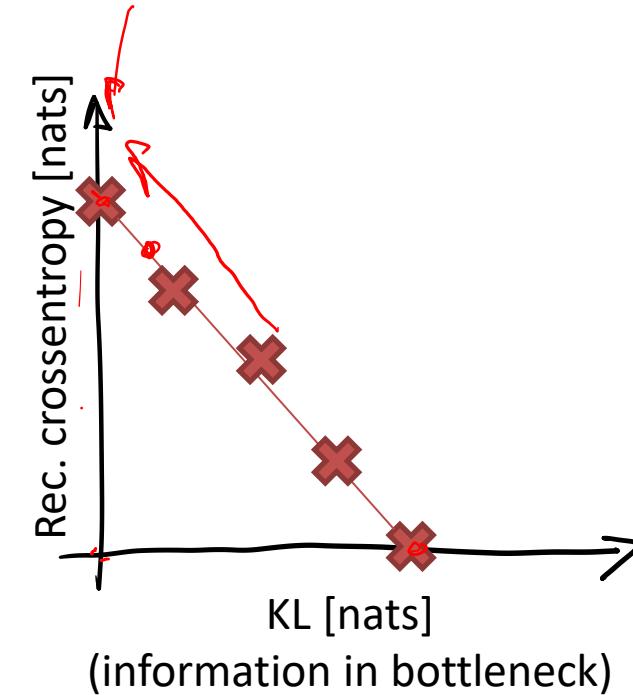
SOTA results in unsupervised phoneme discrimination Fr and EN ZeroSpeech challenge.

Mandarin shows limitation of the method:

- Too little training data (only 2.4h unsup. speech)
- Tonal information is discarded.

# VAE + autoregressive models: the danger of collapsing latents

- Purely Autoregressive models:  
SOTA unconditional log-likelihoods
- Autoencoder:  
information passed through bottleneck  
increases conditional log-likelihood
- Latent priors help us move along this line
- Too strong priors push the model to collapse
  - (ignored latents, unconditional behavior)
  - Hacky solutions: stop optimizing KL term  
(free bits), make it a hyperparameter (VQVAE)

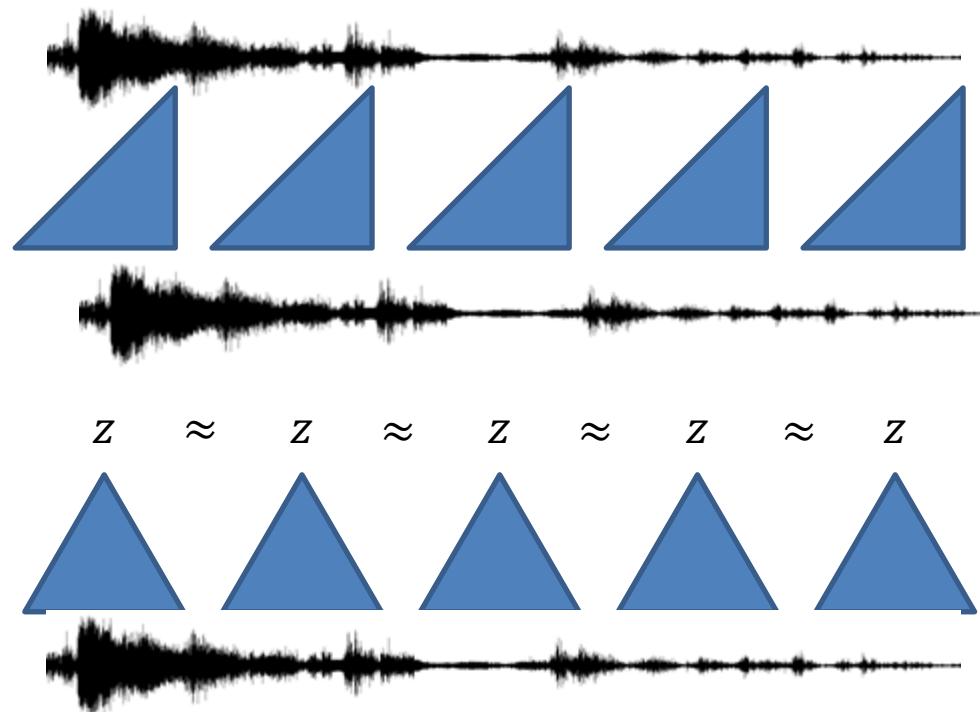


# Naïve latent representation smoothing

Apply a penalty on pairwise differences

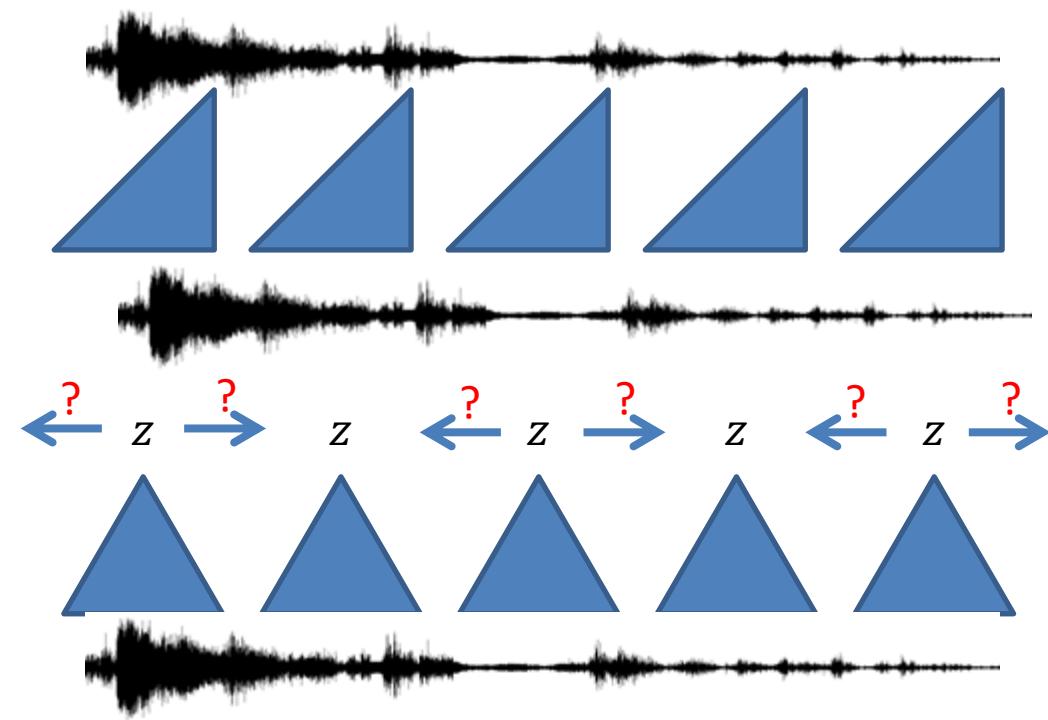
$$\sum_t (z_t - z_{t-1})^2$$

Trivially minimized when  $z_1 = z_2 = z_3 = \dots = z_T$

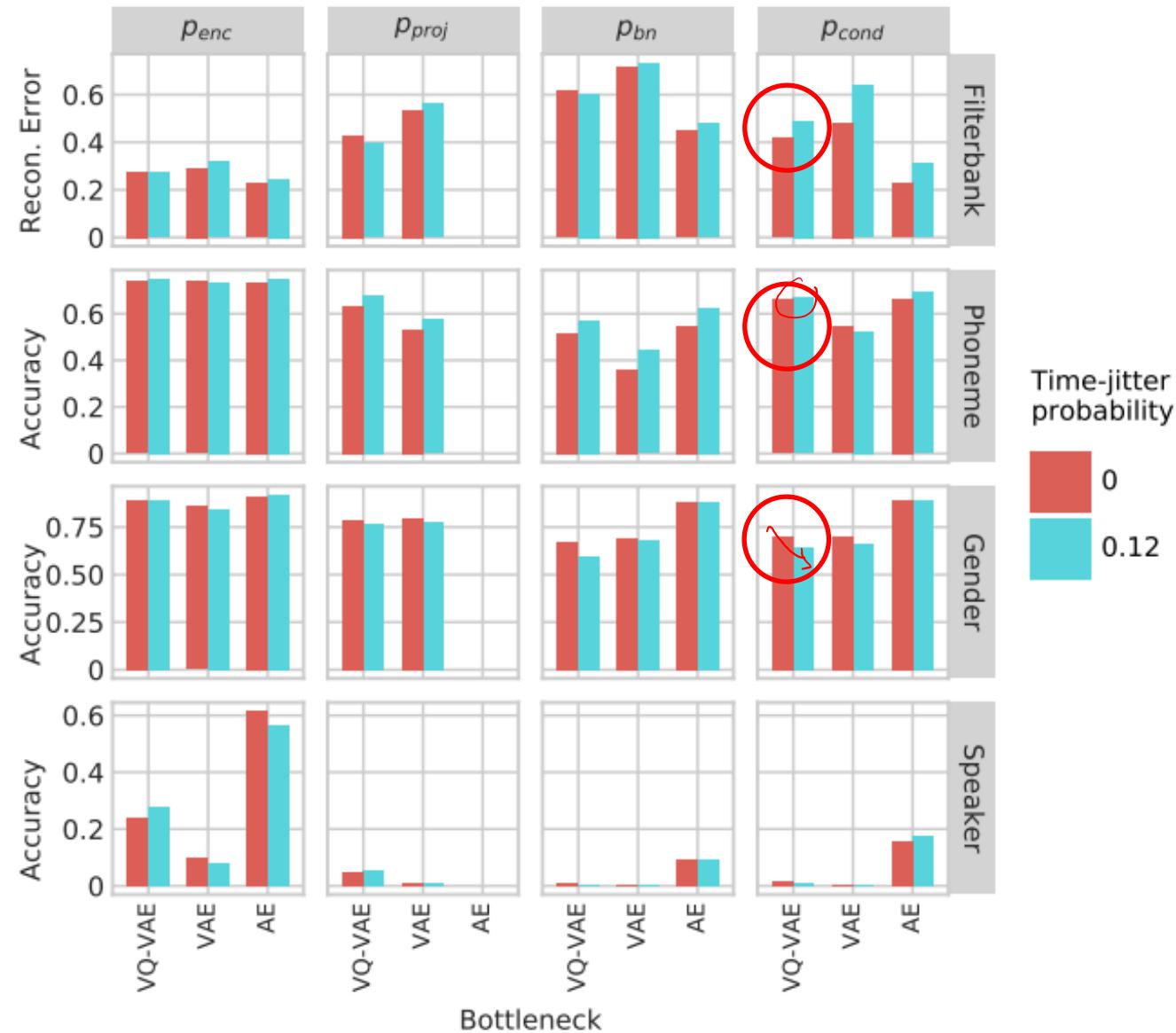


# Randomized time jitter

Don't penalize, **randomize!**

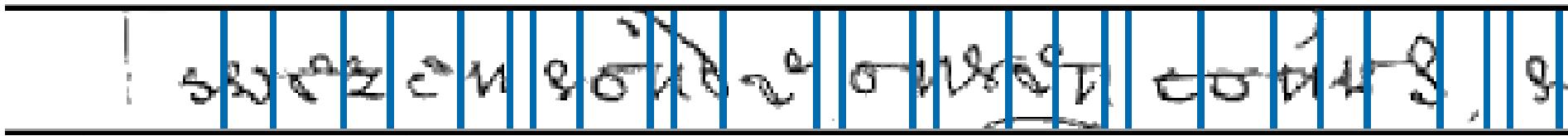


# Randomized time jitter results



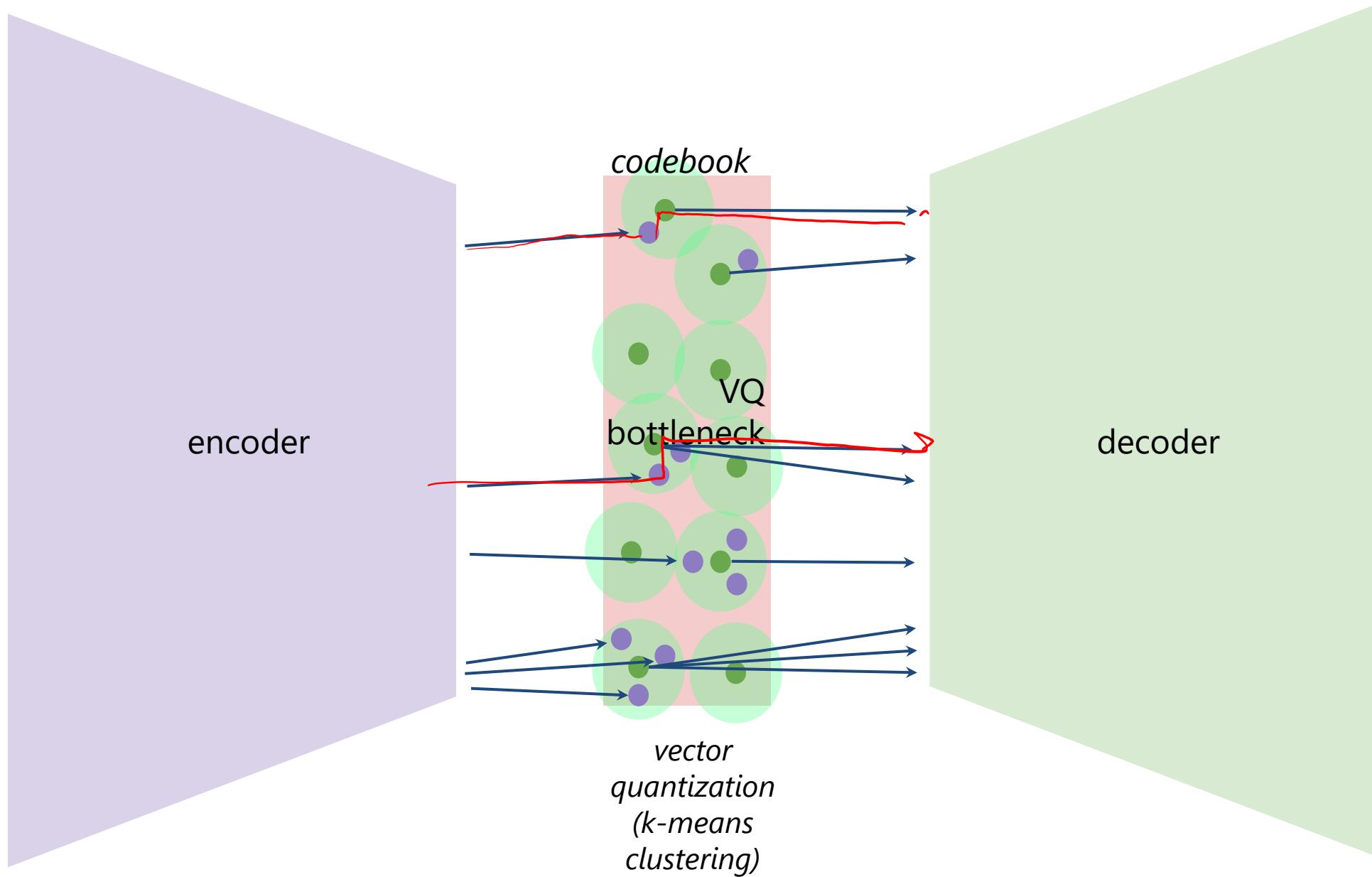
# Learning to segment

We can delineate individual characters



Desired latent representation  
piece-wise constant, with abrupt changes!

# VQ-VAE



# Learning piecewise-constant representations

Quantization criterion:

$$\min_{q_1 \dots q_T} \sum_{t=1}^T \|x_t - q_t\|$$

*encoder outputs*      *VQ tokens*

subject to:  $\sum_{t=1}^T [q_{t-1} \neq q_t] = K$

*expected # of characters*

A red oval highlights the constraint part of the equation:  $\sum_{t=1}^T [q_{t-1} \neq q_t] = K$ . Red arrows point from the labels 'encoder outputs', 'VQ tokens', and 'expected # of characters' to the terms  $\|x_t - q_t\|$ ,  $q_t$ , and  $K$  respectively.

During training: enforce as a minibatch constraint

- easy to specify as a hyperparameter
- “large minibatch” statistics are close-ish to full data  
(our batches have many small chunks)
- fast with greedy merging algorithm  
(could be exact with Viterbi)

During inference: ?

# Learning piecewise-constant representations

Quantization criterion:

$$\min \sum_{t=1}^T \|x_t - q_t\| \quad \text{subject to: } \sum_{t=1}^T [q_{t-1} \neq q_t] = K$$

Reformulate as:

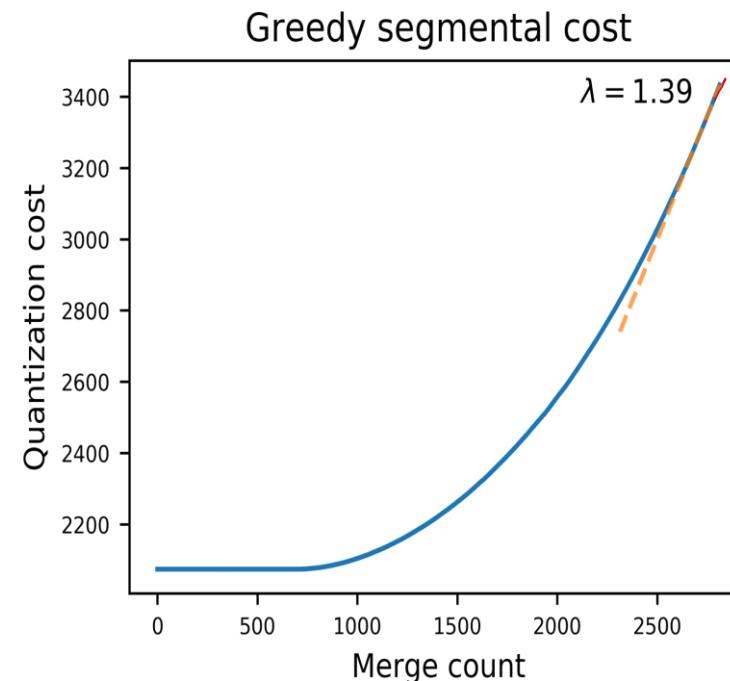
$$\min \sum_{t=1}^T \|x_t - q_t\| + \lambda \left( \sum_{t=1}^T [q_{t-1} \neq q_t] - K \right)$$

During training:

- . estimate  $\lambda$

During inference:

- . use  $\lambda$  to balance the costs



# Dual formulation and Markovian dynamics

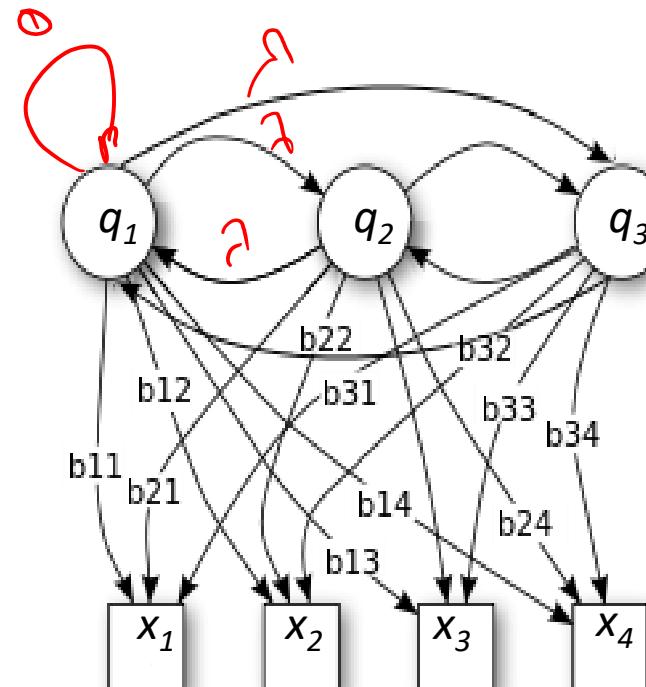
Dual formulation resembles HMM cost:

$$\min \sum_{t=1}^T \|x_t - q_t\| + \lambda \left( \sum_{t=1}^T [q_{t-1} \neq q_t] - K \right)$$

*state change cost*

*encoder outputs*

Observed states: encoder outputs  
Hidden states: sequence of VQ tokens

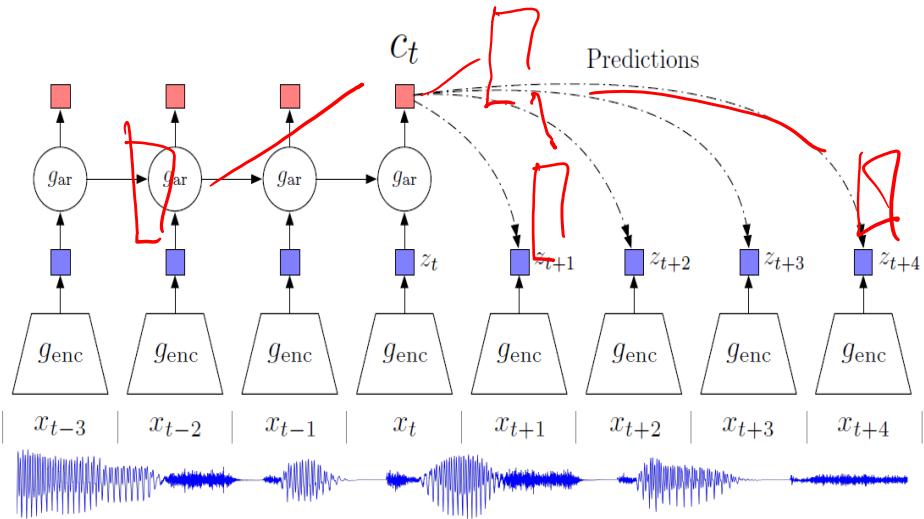


*Image source: Wikipedia*

# Piecewise-constant results: Scribblelens

<i>Dataset</i>	<i>Model</i>	<i>Downstream CER</i>	<i>MLP Accuracy @Bottleneck</i>	<i>Token -&gt; Char Accuracy</i>	<i>Rand score</i>	<i>Mutual information</i>
<i>Single scribe Tasman</i>	<i>Unsupervised base</i>	47%	54%	34%	0.15	0.15
	<b><i>Piecewise-const VQ bottleneck</i></b>	<b>30%</b>	<b>65%</b>	<b>57%</b>	<b>0.23</b>	<b>0.36</b>
<i>All scribes</i>	<i>Unsupervised base</i>	60%	42%	30%	0.13	0.10
	<b><i>Piecewise-const VQ bottleneck</i></b>	<b>51%</b>	<b>49%</b>	<b>39%</b>	<b>0.18</b>	<b>0.20</b>

# Contrastive Coding



At step  $t$  make  $K$  independent predictions  $p_t^k$  for  $K$  next latent vectors.

Score them using negative sampling

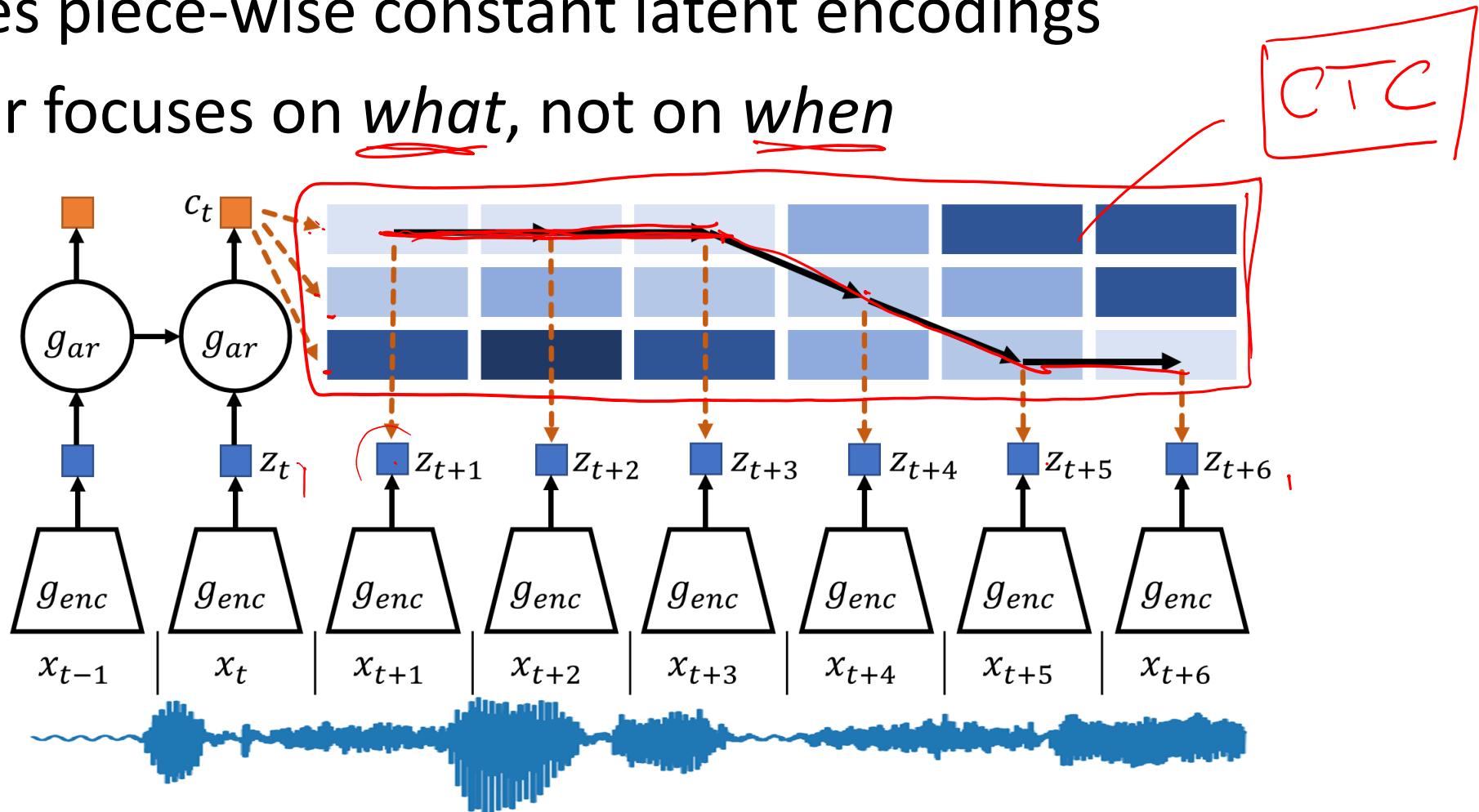
$$s_t^k = \frac{e^{<p_t^k, z_{t+k}>}}{e^{<p_t^k, z_{t+k}>} + \sum_n e^{<p_t^k, z_n>}}$$

Where  $z_n$  are some randomly selected latent codes.

# Aligned CPC ([arxiv.org/abs/2104.11946](https://arxiv.org/abs/2104.11946))

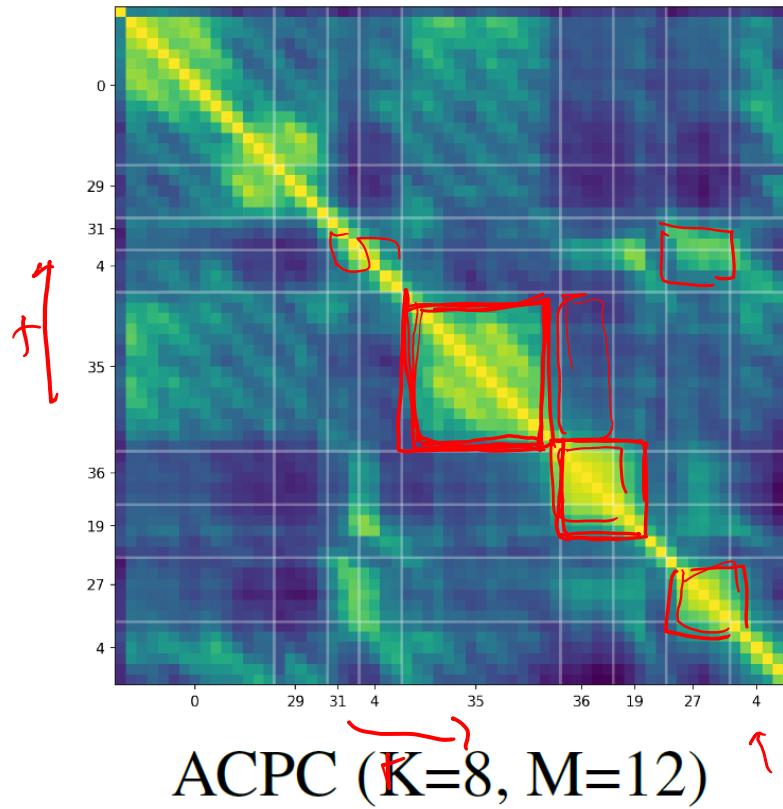
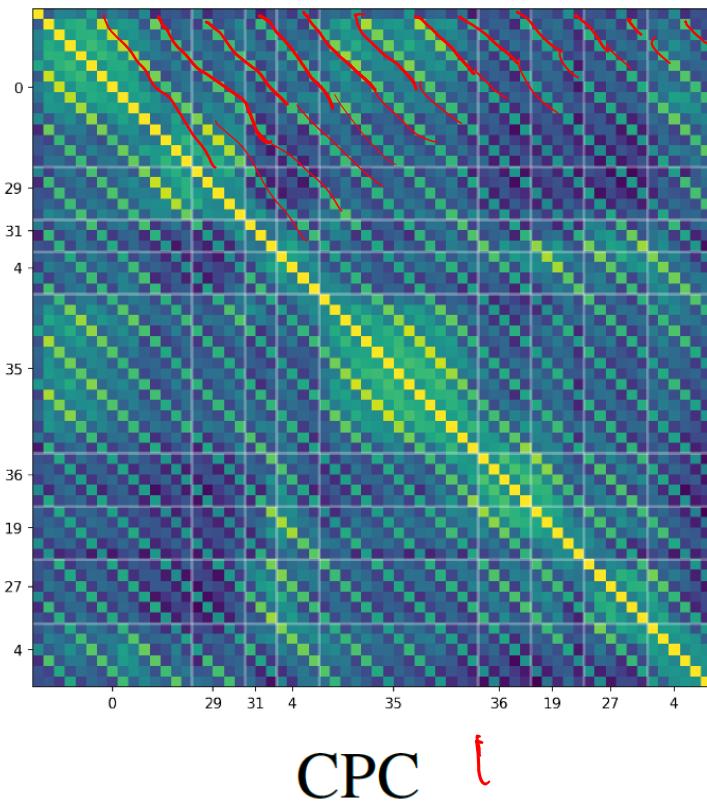
Core idea: make few predictions, align to many frames:

- Promotes piece-wise constant latent encodings
- Predictor focuses on *what*, not on *when*



# ACPC: more stable representations

Vanilla CPC exploits conv strides to learn relative position  
Phoneme boundaries more clearly visible in ACPC



# ACPC: better ABX and phone recognition scores

ABX Scores (lower is better)

Dev	CPC (ZeroSpeech [10])		CPC (our baseline)		ACPC M=12 K=8	
	Within	Across	Within	Across	Within	Across
clean	6.18%	8.02%	6.68%	8.39%	<b>5.37%</b>	<b>7.09%</b>
other	8.46%	13.59%	9.03%	13.87%	<b>7.46%</b>	<b>12.60%</b>

Framewise phoneme prediction  
(higher is better)

Model	Step time	Encoded $z_t$		Contextual $c_t$	
		Train	Val	Train	Val
CPC	2.6ms	51.5%	51.2%	67.8%	67.5%
ACPC M=12 K=10	2.4ms	51.8%	51.3%	68.8%	68.4%
ACPC M=12 K=8	2.1ms	52.0%	51.9%	<b>69.7%</b>	68.6%
ACPC M=12 K=6	1.8ms	<b>52.2%</b>	<b>52.1%</b>	69.2%	<b>68.8%</b>
ACPC M=12 K=4	<b>1.5ms</b>	52.1%	51.9%	69.0%	68.7%
ACPC M=16 K=10	2.4ms	52.1%	51.9%	69.1%	<b>68.8%</b>
ACPC M=20 K=12	2.6ms	52.1%	51.9%	68.0%	67.8%

# Summary

- We can design latent variables models that work well on sequential data
- We can enforce priors on the latent space, by
  - Applying hidden dynamics
    - during training, force a fixed number of code changes
    - during testing, use the equivalent HMM formulation
    - Constraint-penalty duality: very generic technique, can also use it to e.g. enforce unit unigram probabilities.
  - Aligning model predictions to latent encodings