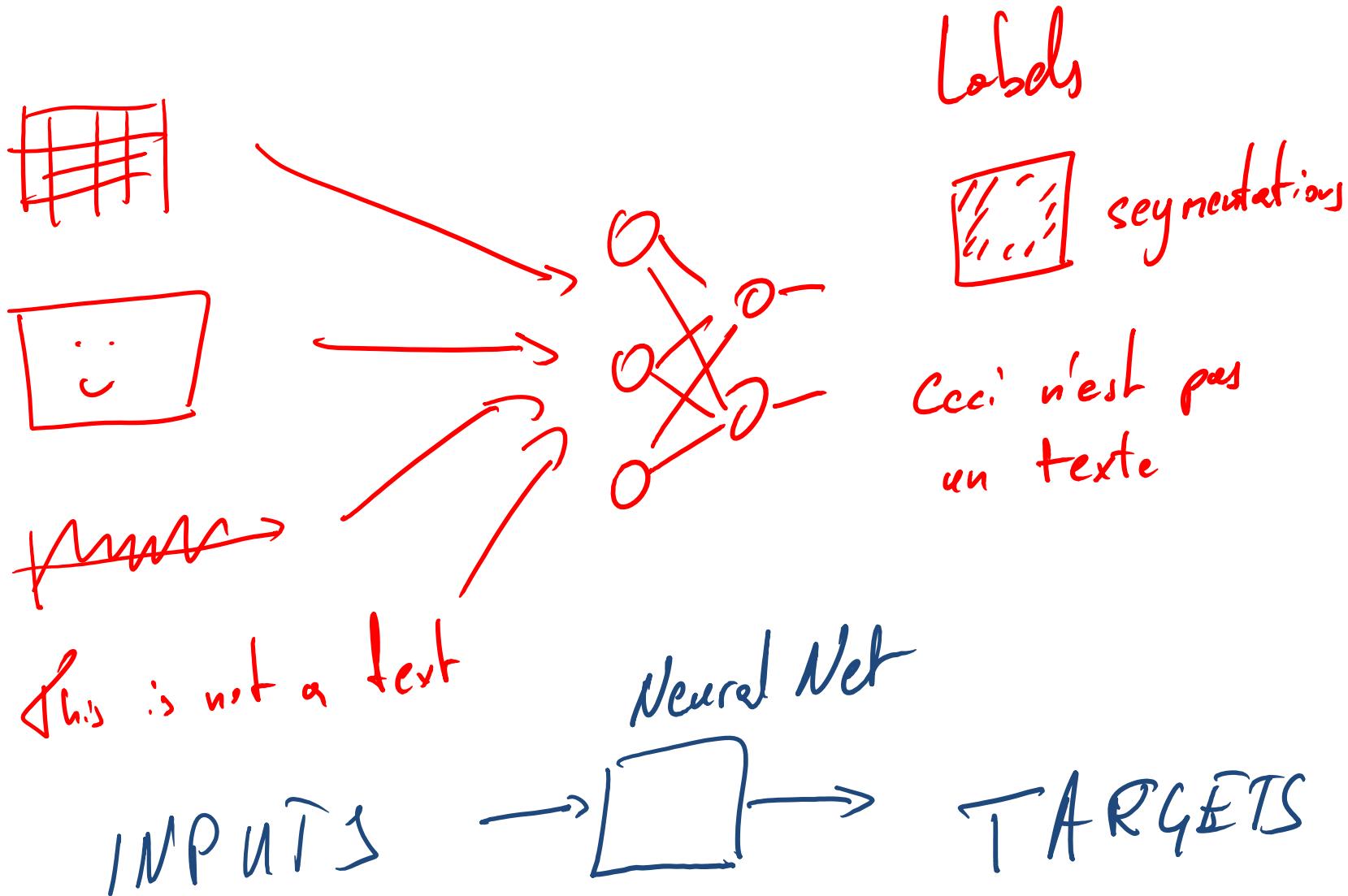


# Deep Unsupervised Learning

Deep Learning @ UWr 2021

Jan Chorowski

# Motivation: Deep Nets excel at supervised tasks



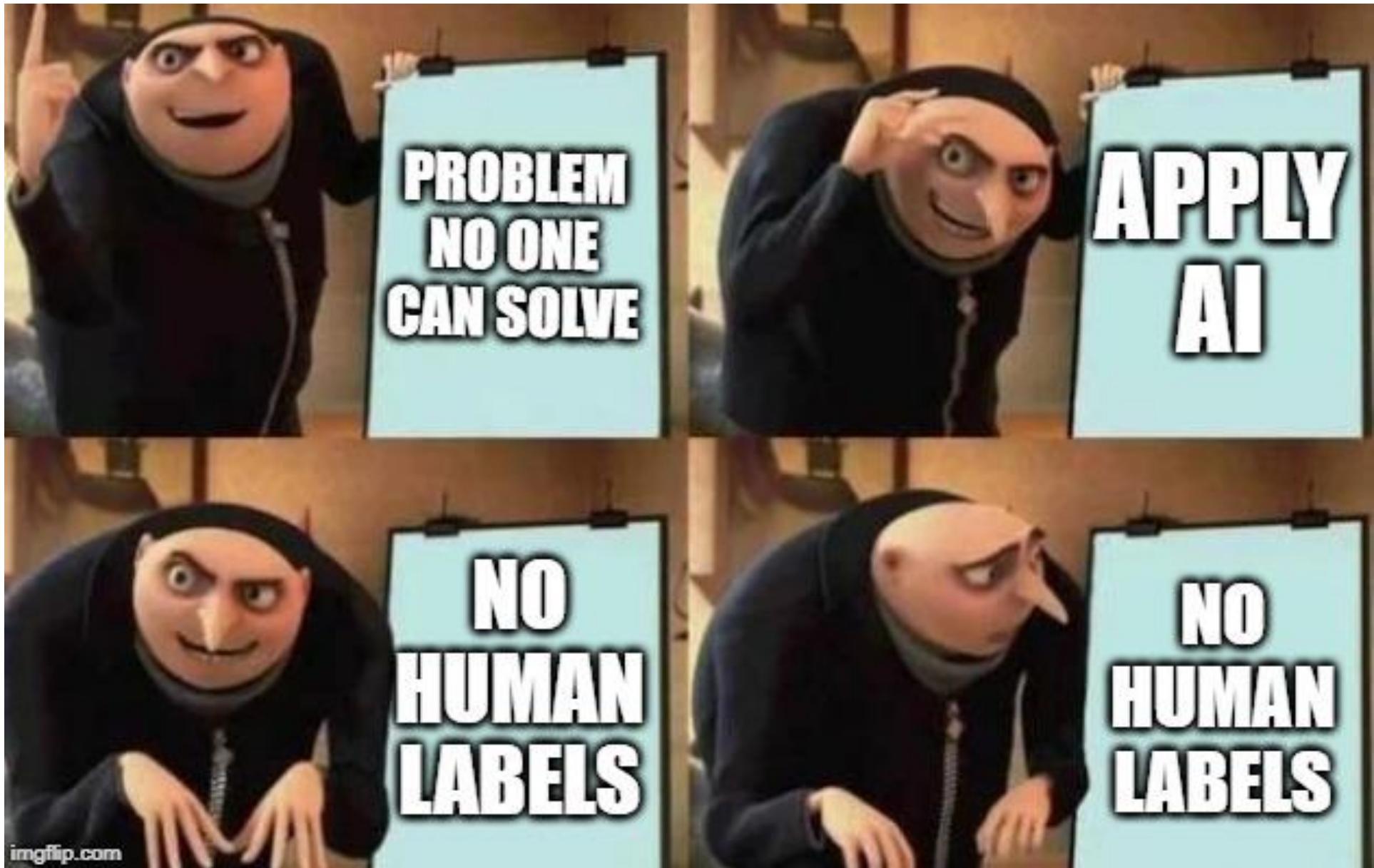
०८८



802 axoott

Scribner





# 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

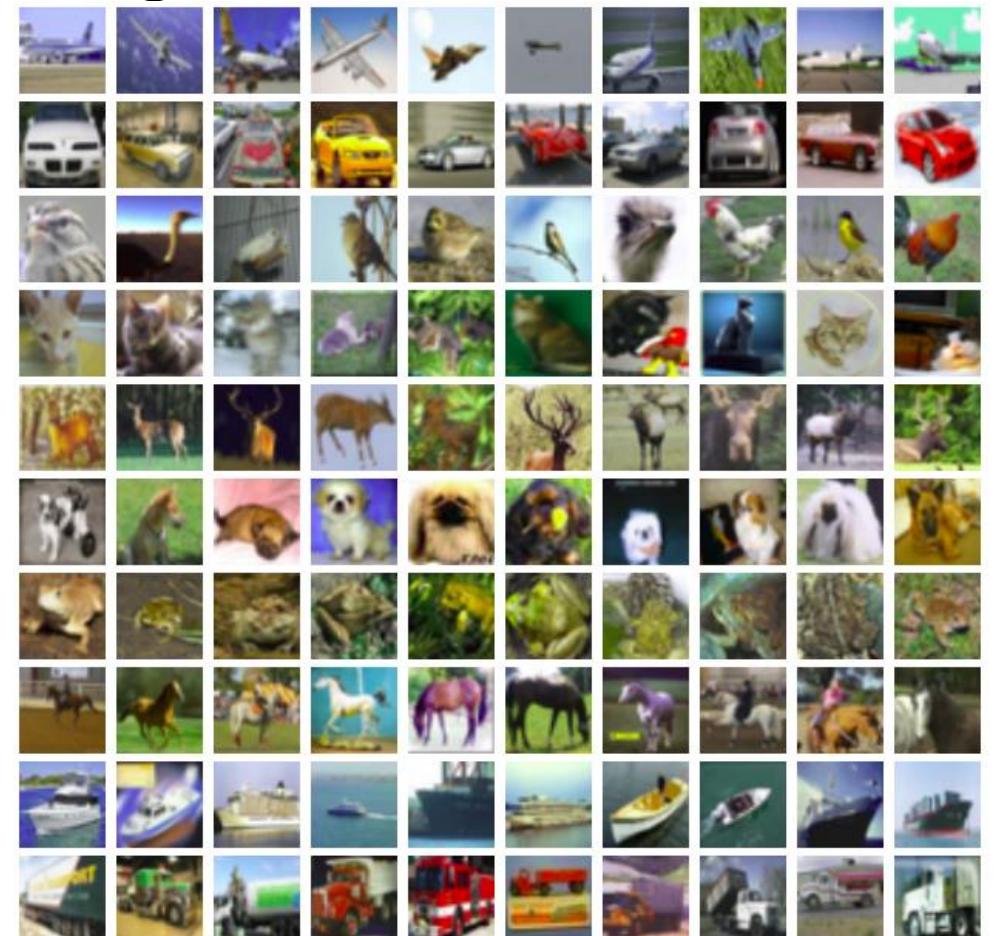
## What losses we can use:

- Log-likelihood
- Learn to produce plausibly looking samples (as judged by an auxiliary network)
- 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!



QQ: what generative models we have  
already seen?

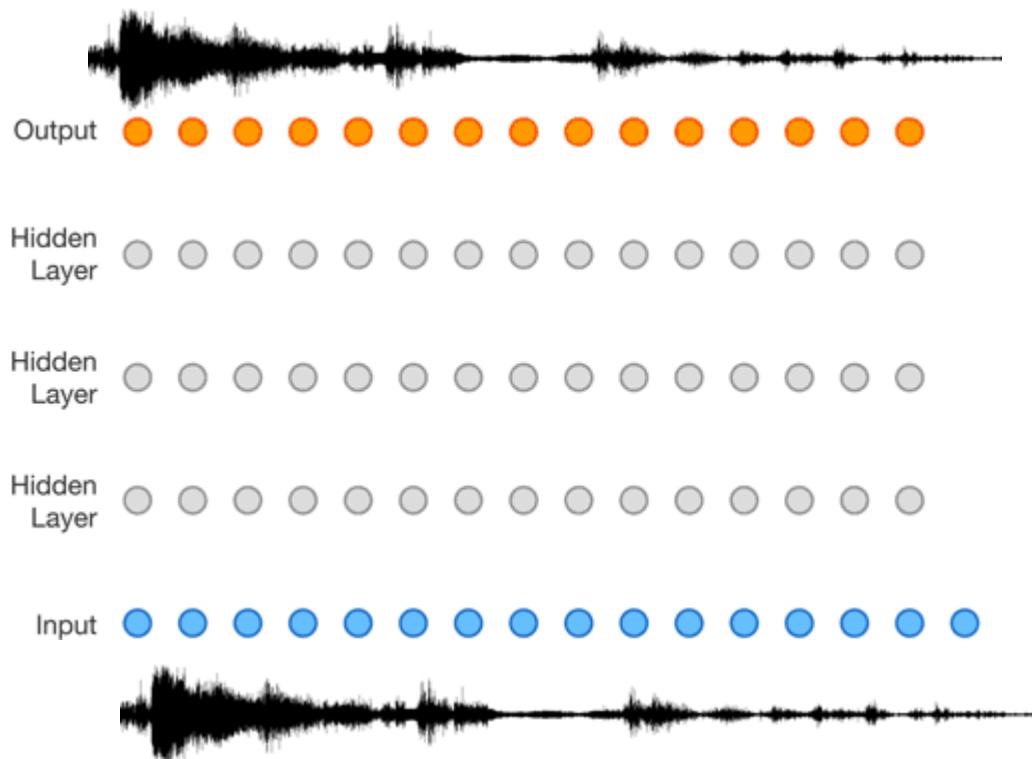
# **AUTOREGRESSIVE MODELS**

# 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

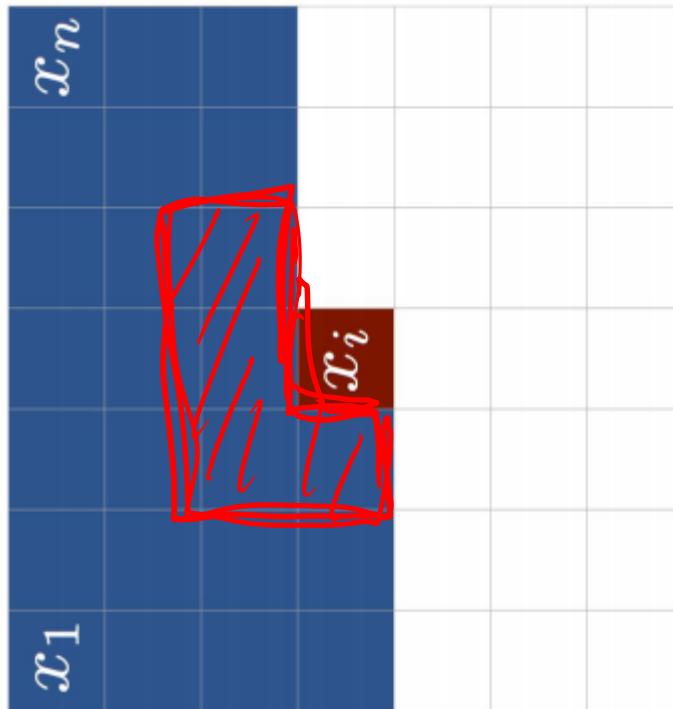
# WaveNet: Autoregressive modeling of speech

- 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.

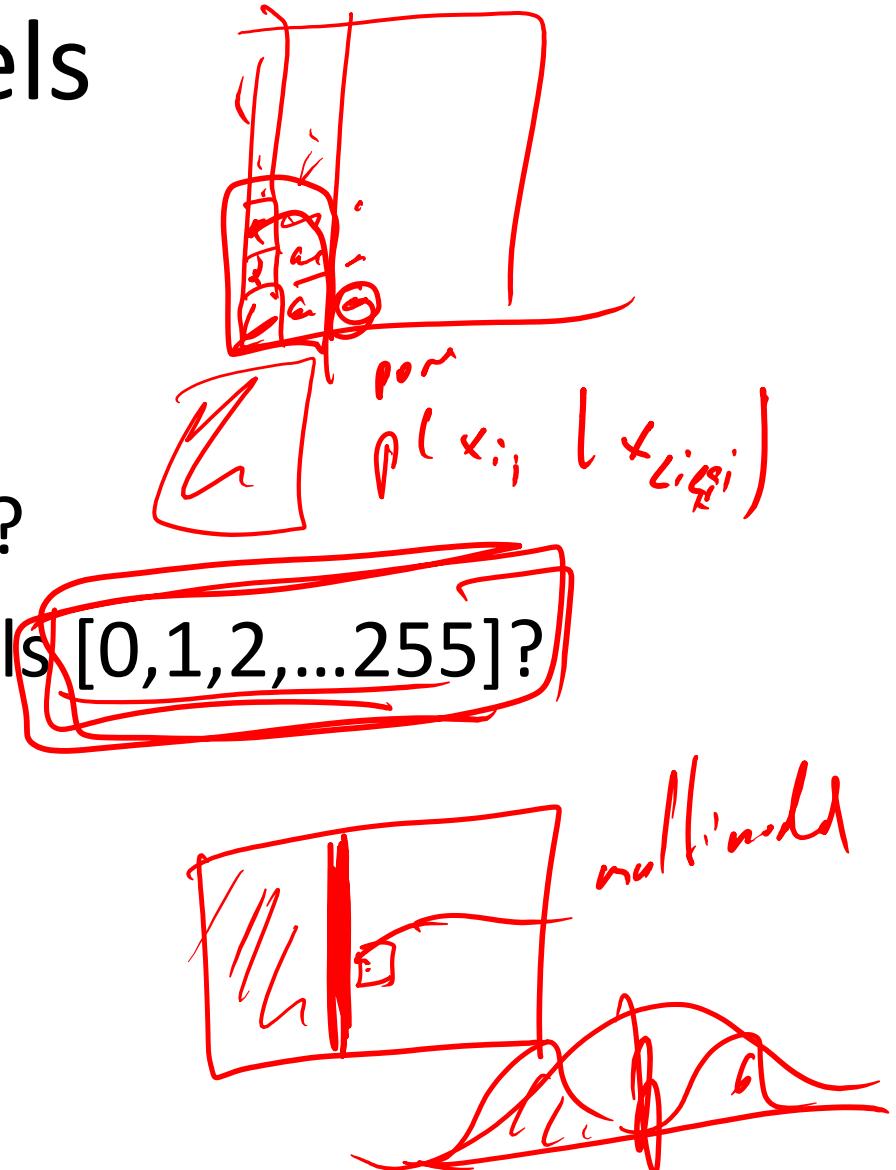
$p(x_i|x_{<i})$  computed using  
convolutional neural nets

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

# Modeling pixels

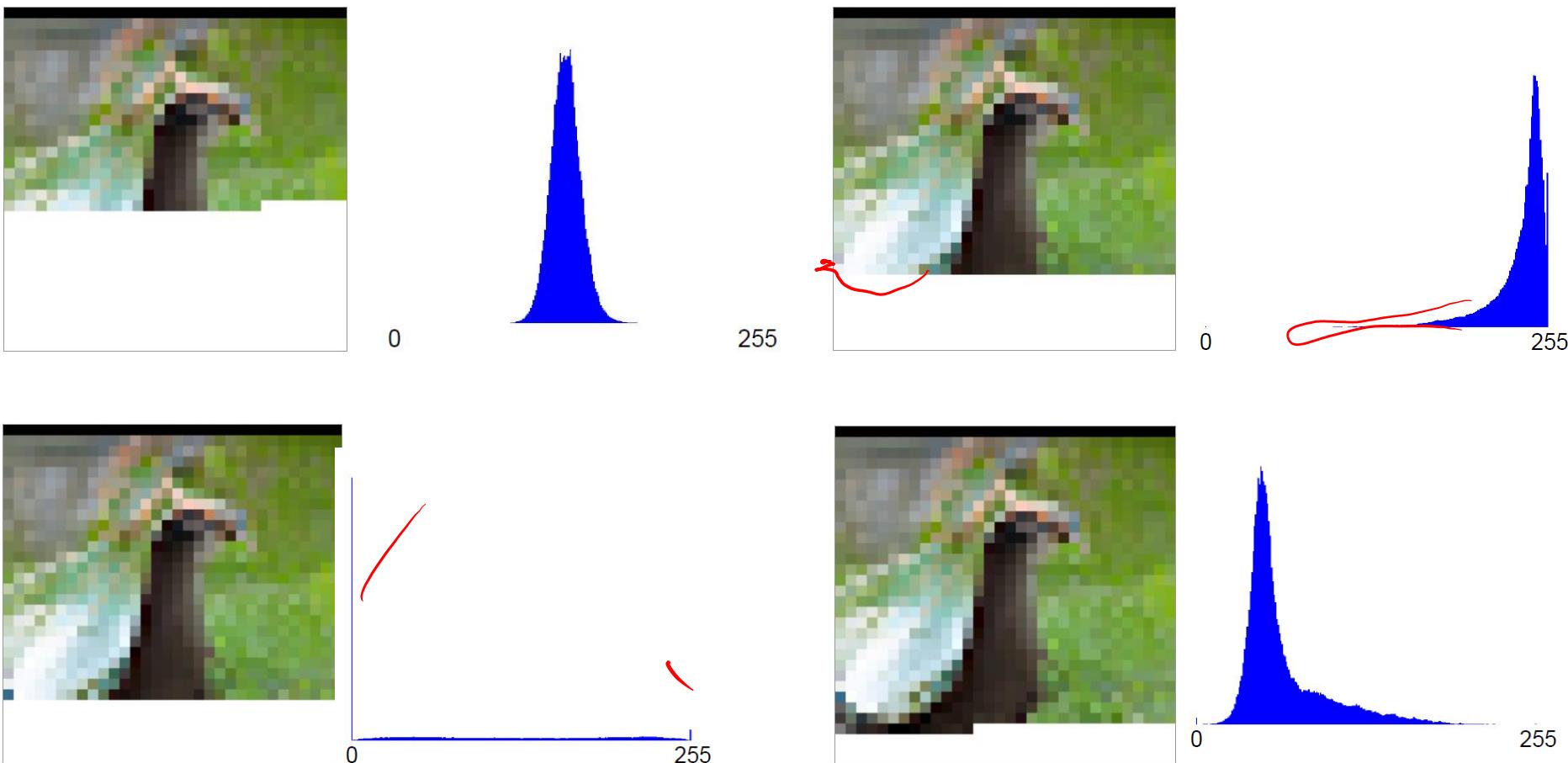
How to model pixel values:

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



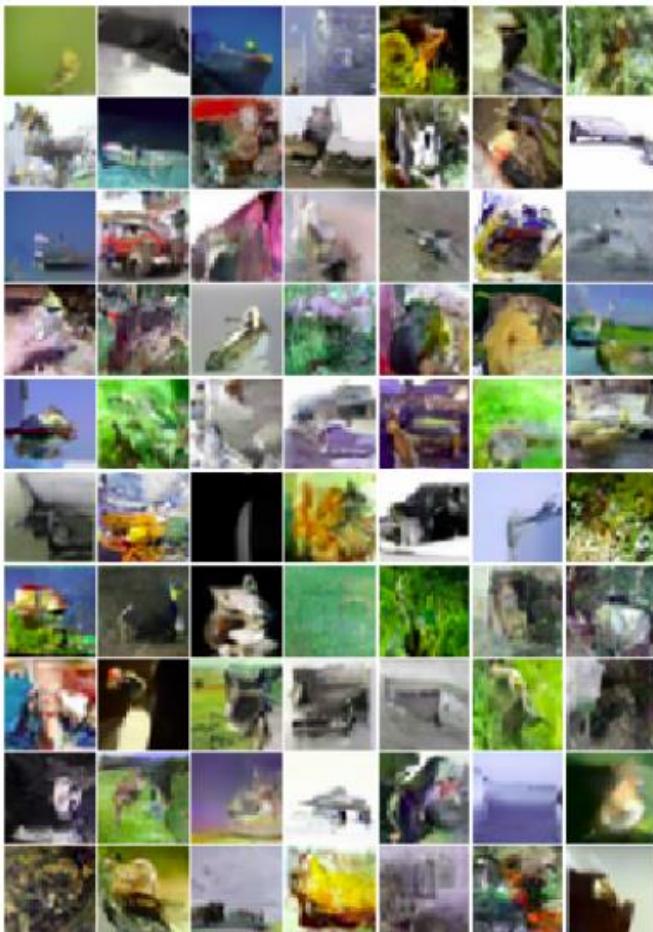
What are the implications?

# Modeling pixel values



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

# PixelCNN samples & completions

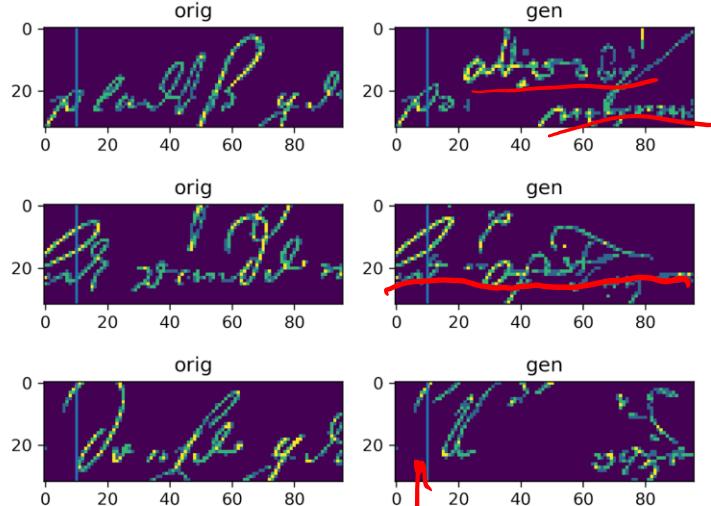


Salimans et al, “A PixelCNN Implementation with Discretized Logistic Mixture Likelihood and Other Modifications”

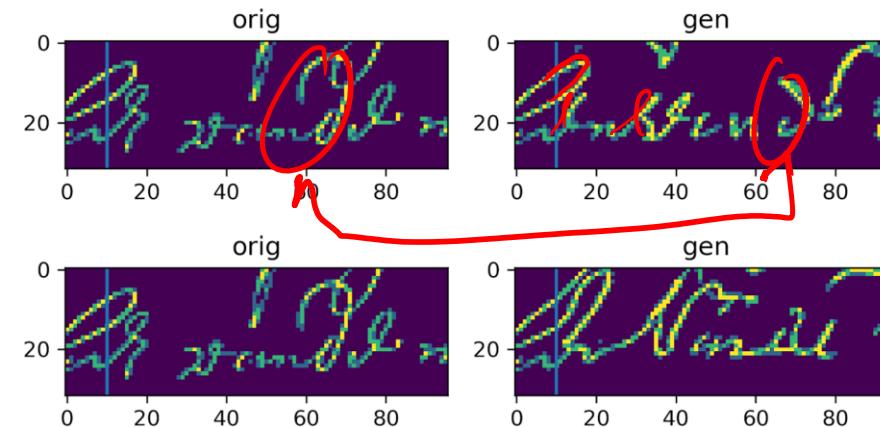
van den Oord, A., et al. “Pixel Recurrent Neural Networks.” ICML (2016).

# PixelCNN samples

Unconditional



Conditioned on text



Handwriting output system

# 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

state - of - the - art

# 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**

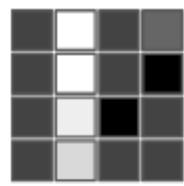
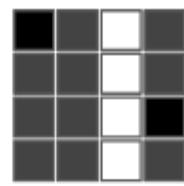
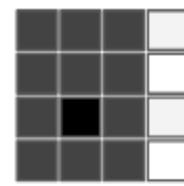
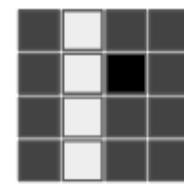
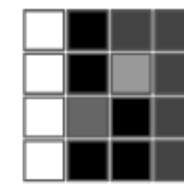
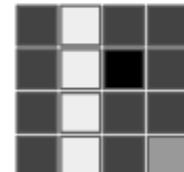
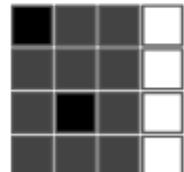
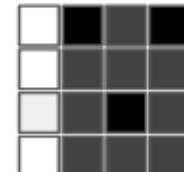
# Ad break

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

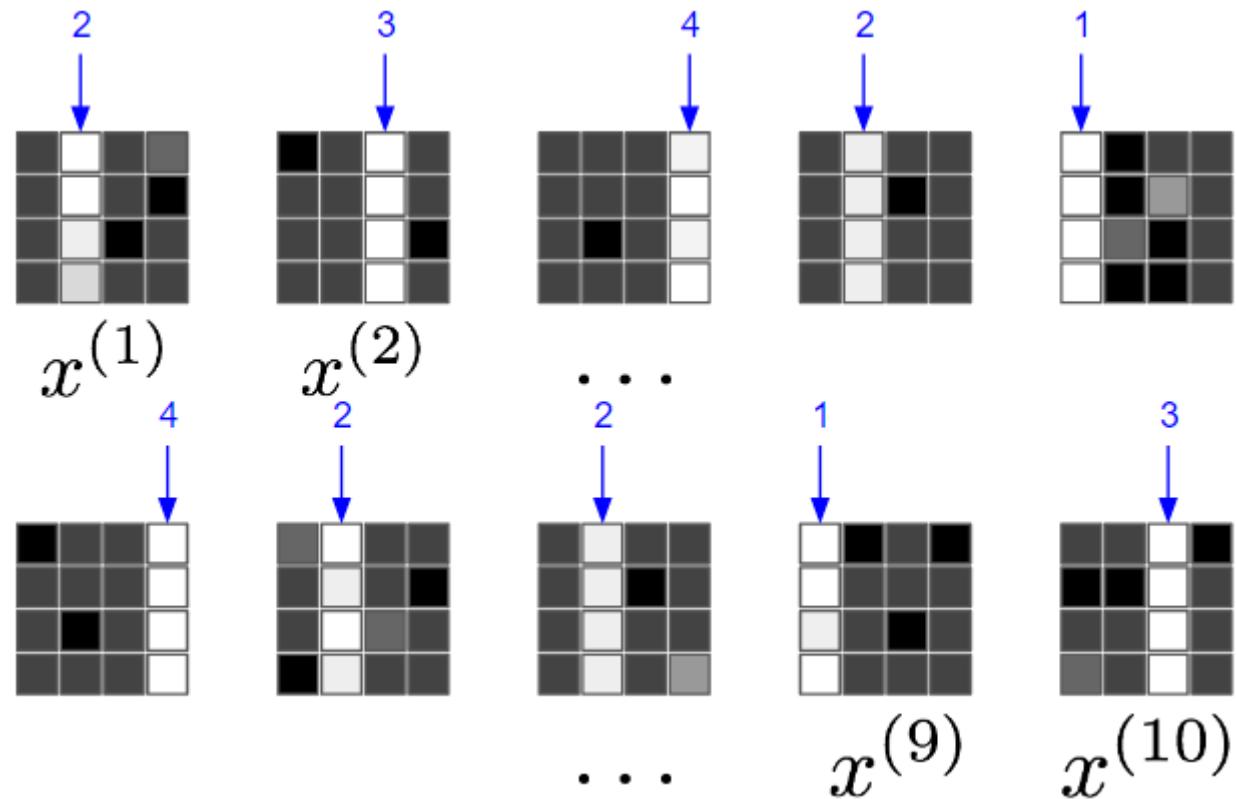
Think about all the properties of this dataset:

 $x^{(1)}$  $x^{(2)}$  $\dots$  $x^{(9)}$  $x^{(10)}$  $\dots$ 

# The rows are correlated

 $x^{(1)}$  $x^{(2)}$  $\dots$  $\dots$  $x^{(9)}$  $x^{(10)}$

There's a simple encoding  
(the LATENT representation)



# Data structure

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

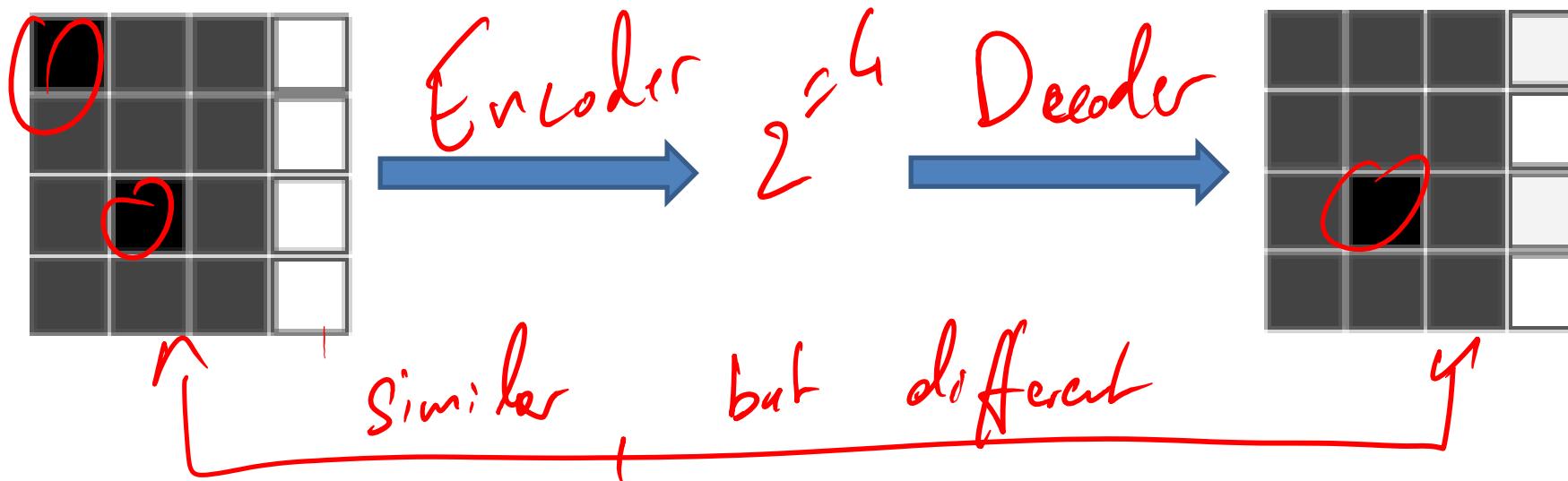
$$z^{(n)} = 1 \text{ or } 2, 3, 4$$

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

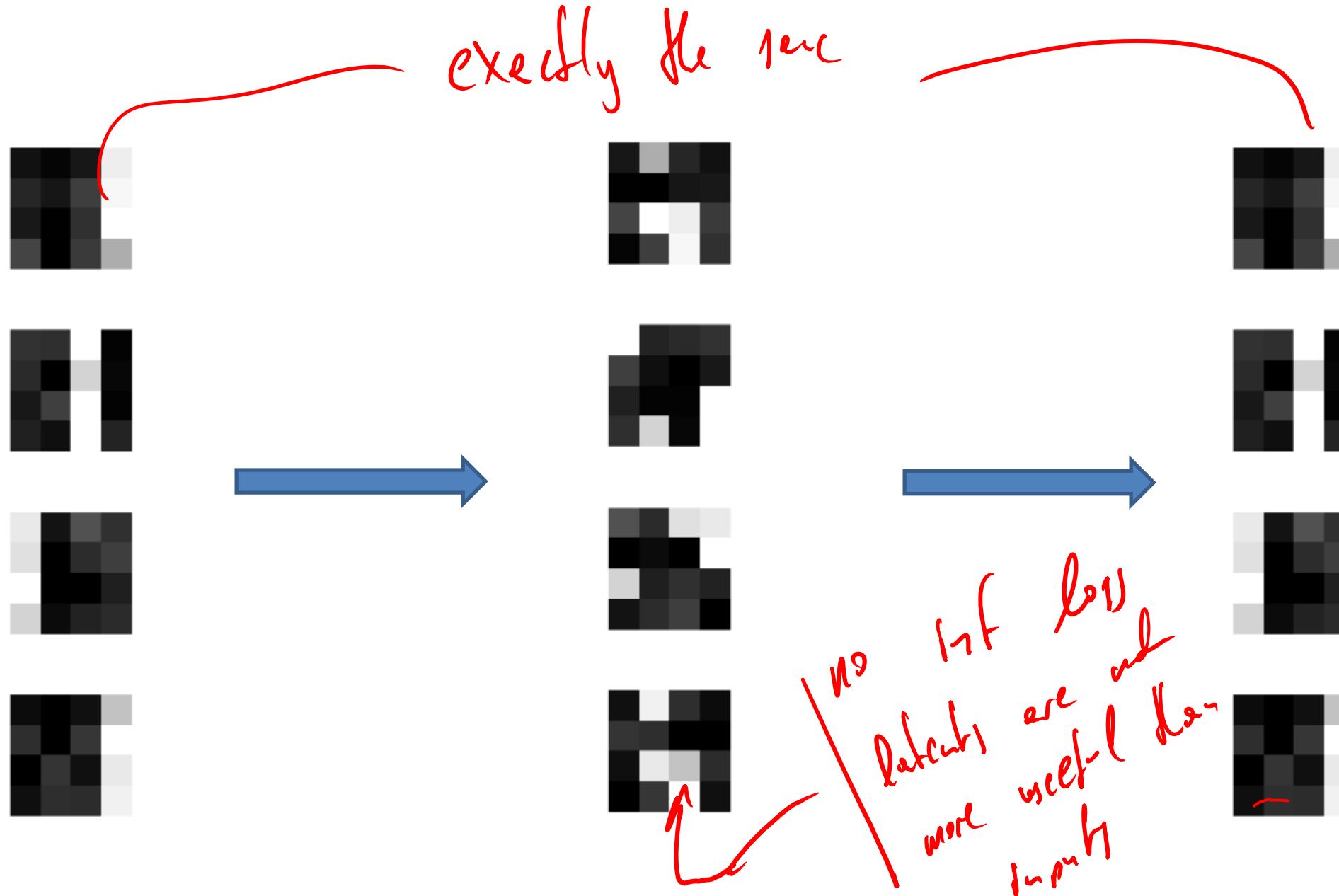
How to train a model to uncover it?

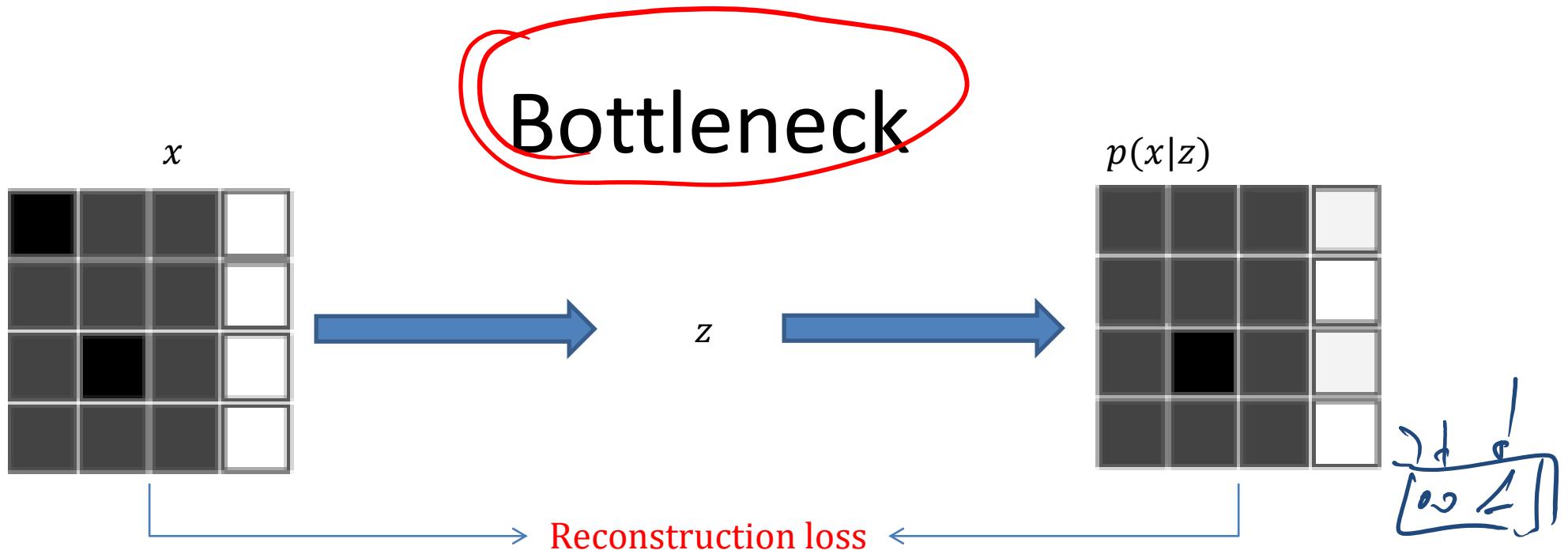
# Autoencoders

SELF



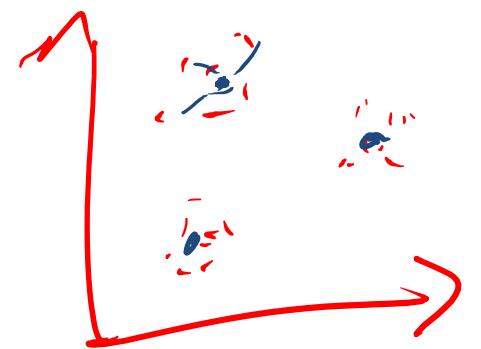
# Perfect Reconstruction === Useful Representations?





Bottleneck prevents "trivial" encodings

- Dimensionality reduction
- Sparsity ( $z$  is mostly 0)
- !



# 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)

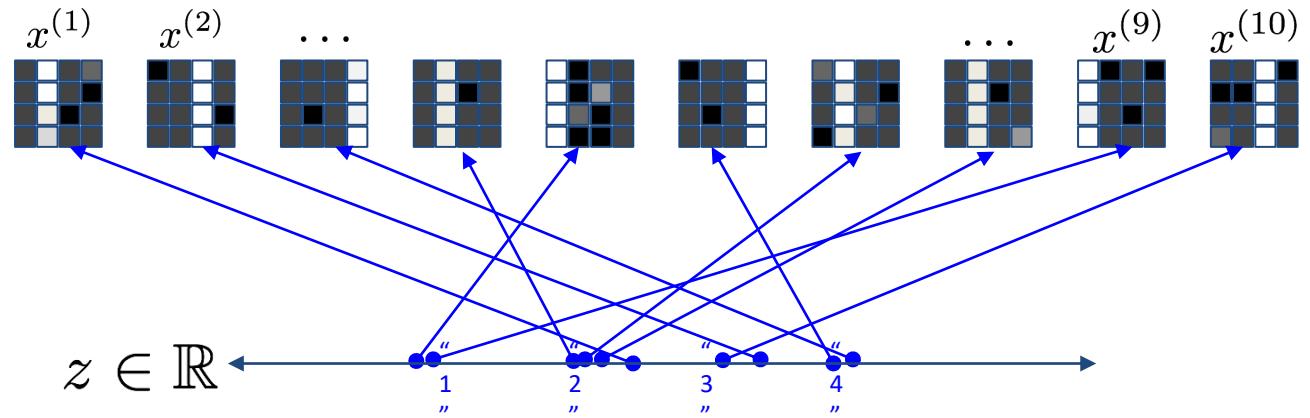
# Simple Auto-encoders summary

- + 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**

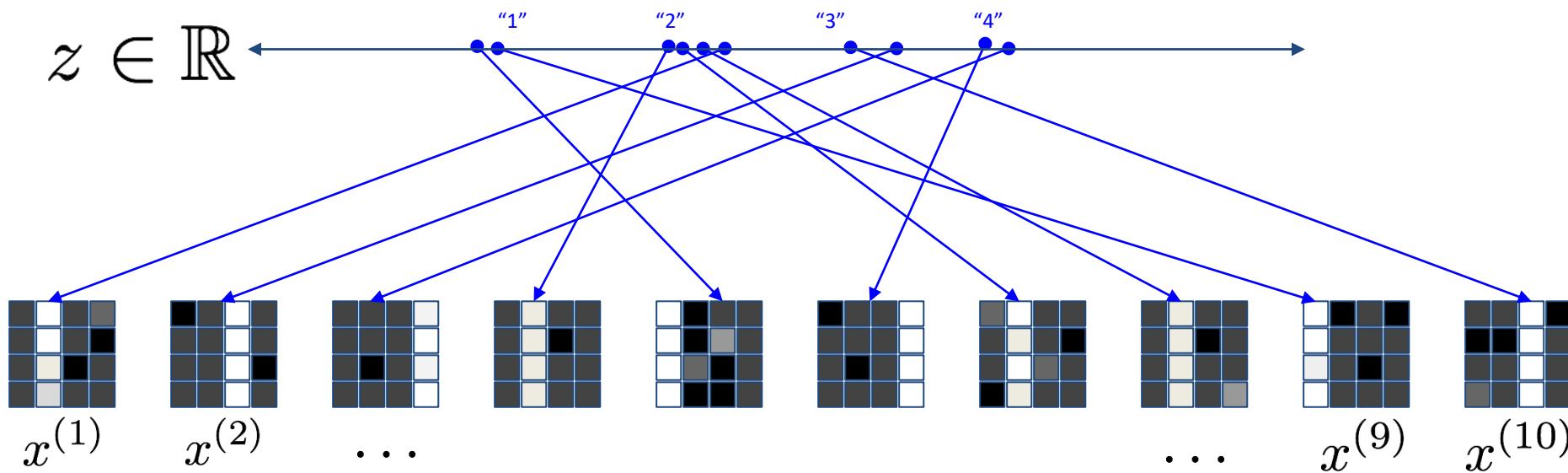
# Intuition

- 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



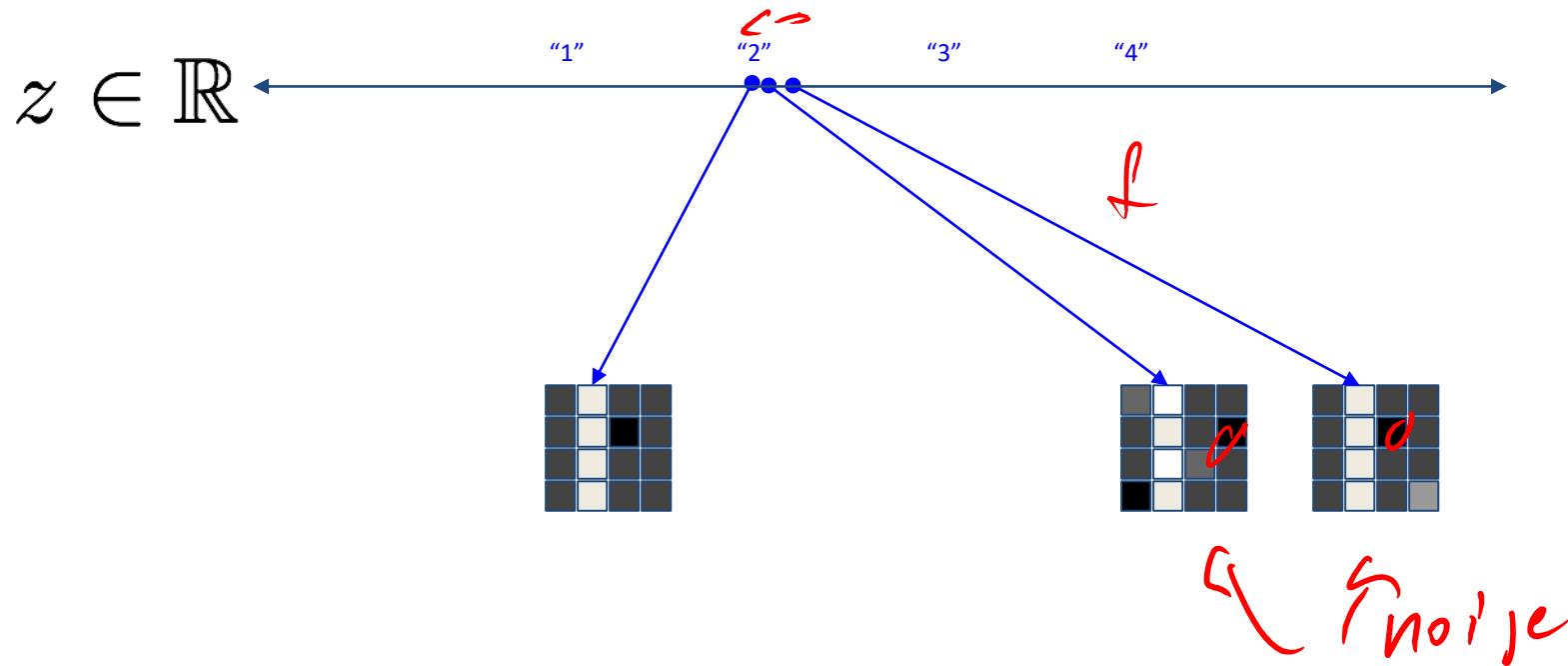
# Data manifold

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

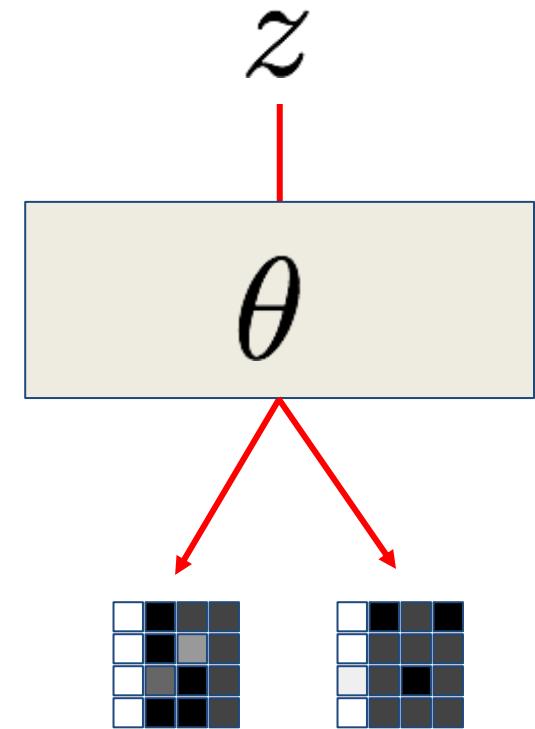
Think how to implement a neural network that takes  $z$  and produces a distribution over  $x$ :

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

Is your input one-dimensional?

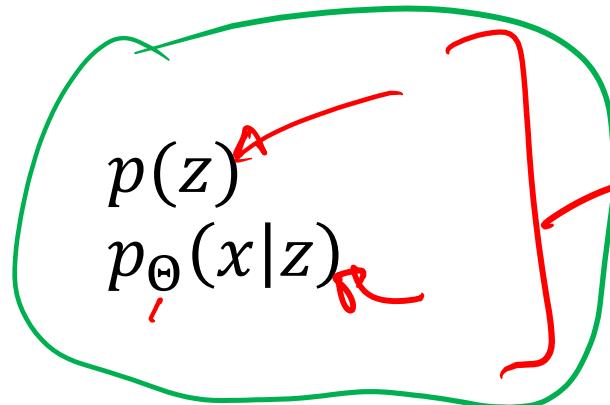
Is your output 16-dimensional?

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



# Inference

Generation:



Generative  
Model

with latents  $z?$ .

Inference:

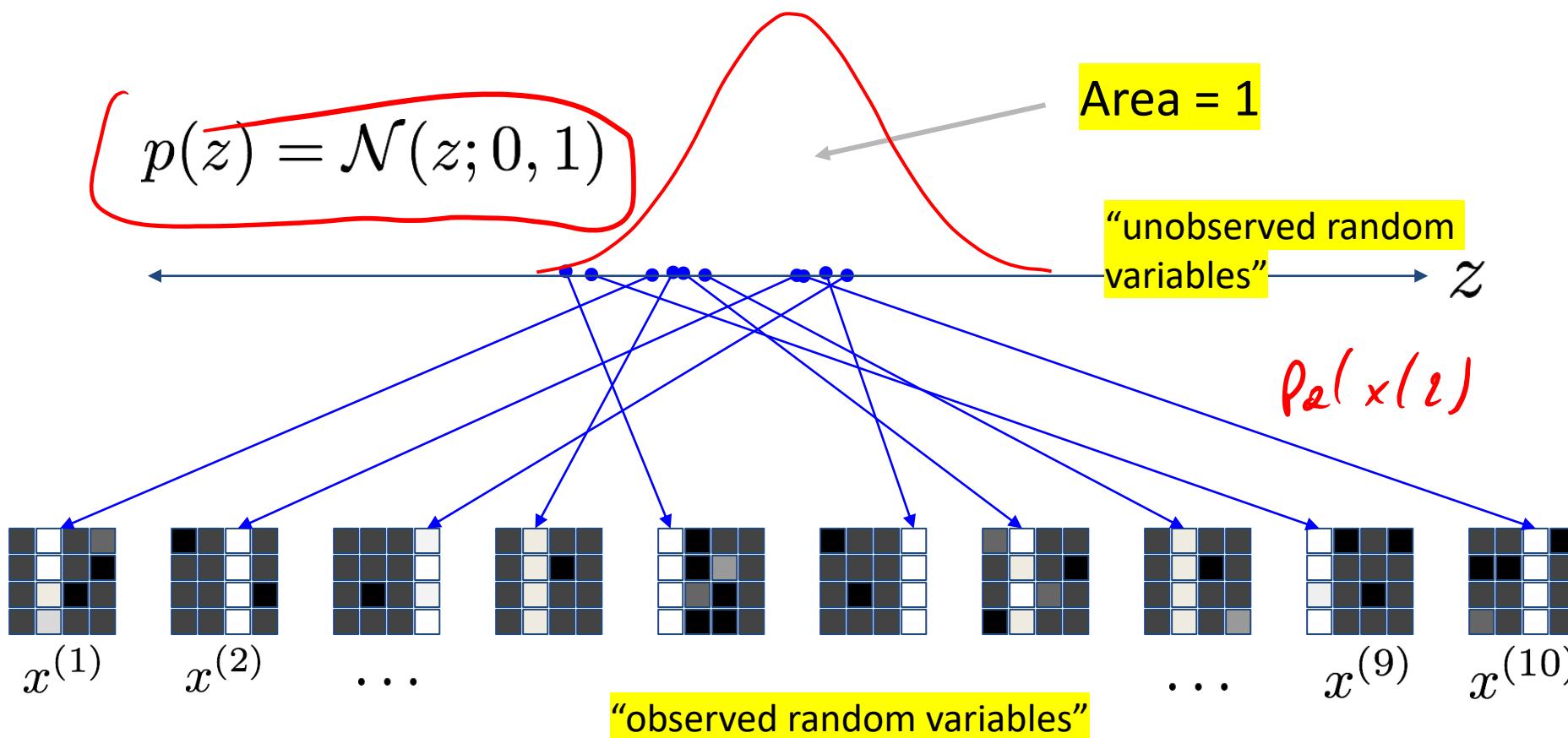
$$p_{\Theta}(z|x) = \text{????}$$

Bayes says:

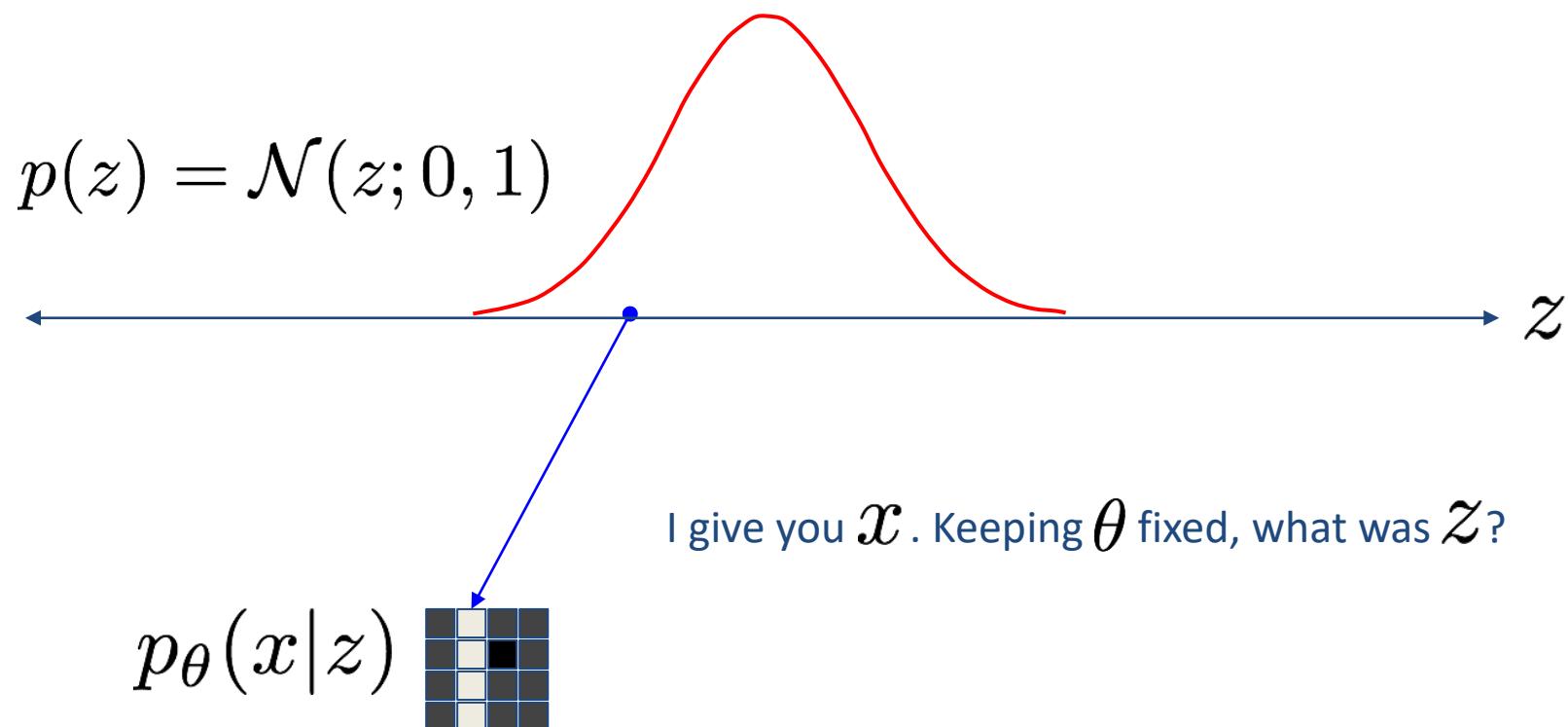
$$p_{\Theta}(z|x) = \frac{p_{\Theta}(x|z)p(z)}{\int dz' p_{\Theta}(x|z')p(z')}$$

$p_{\Theta}(x)$   
 $p$   
train =  $\exp(\theta)$

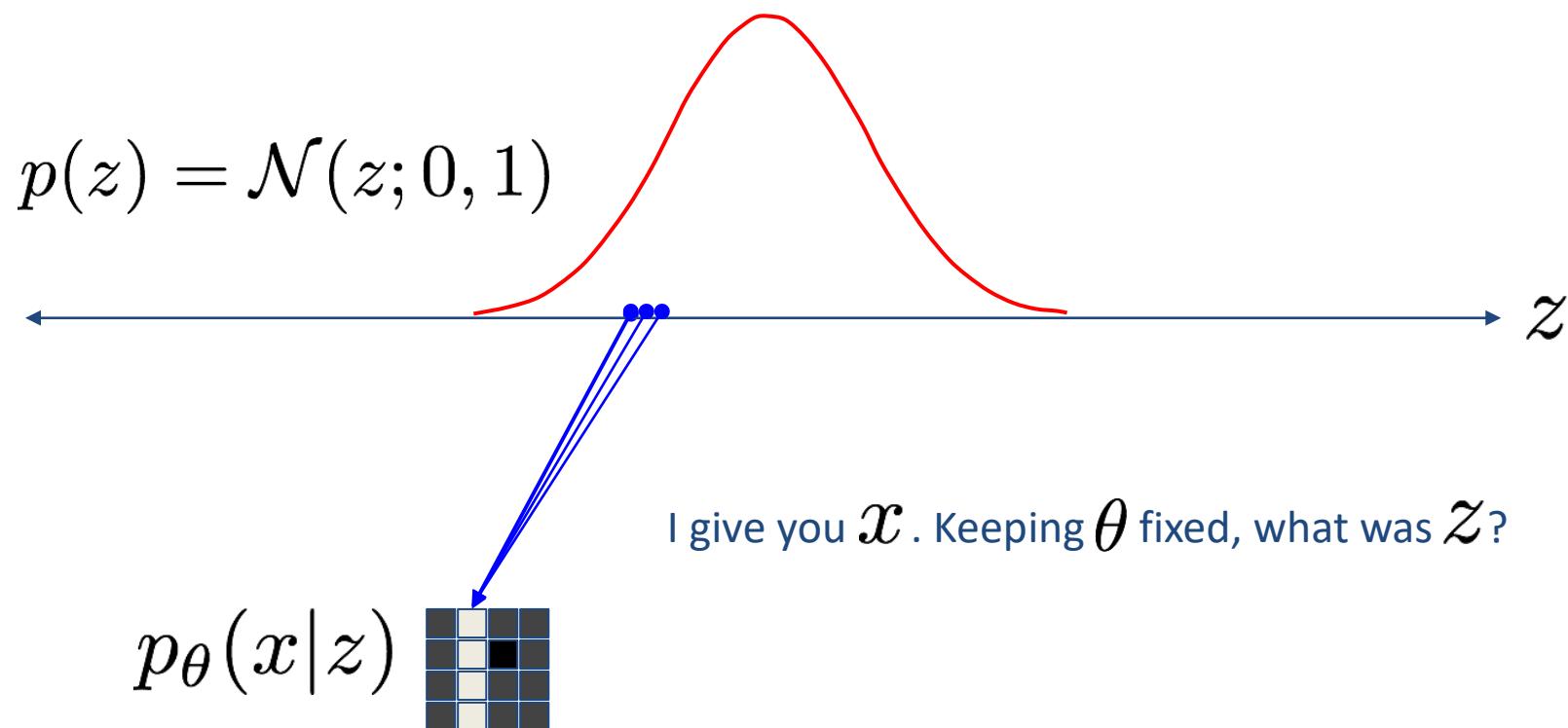
# Inference starts with priors



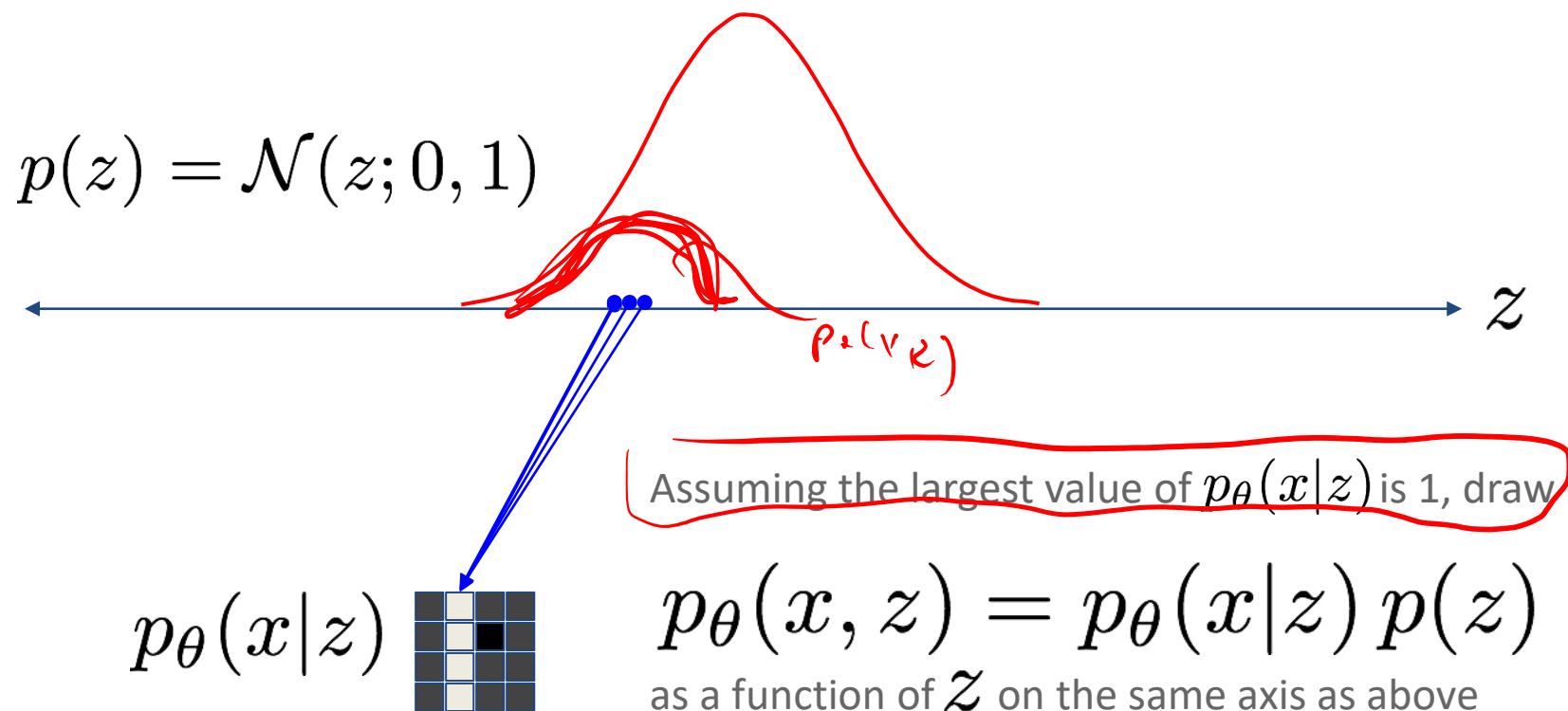
# Inference



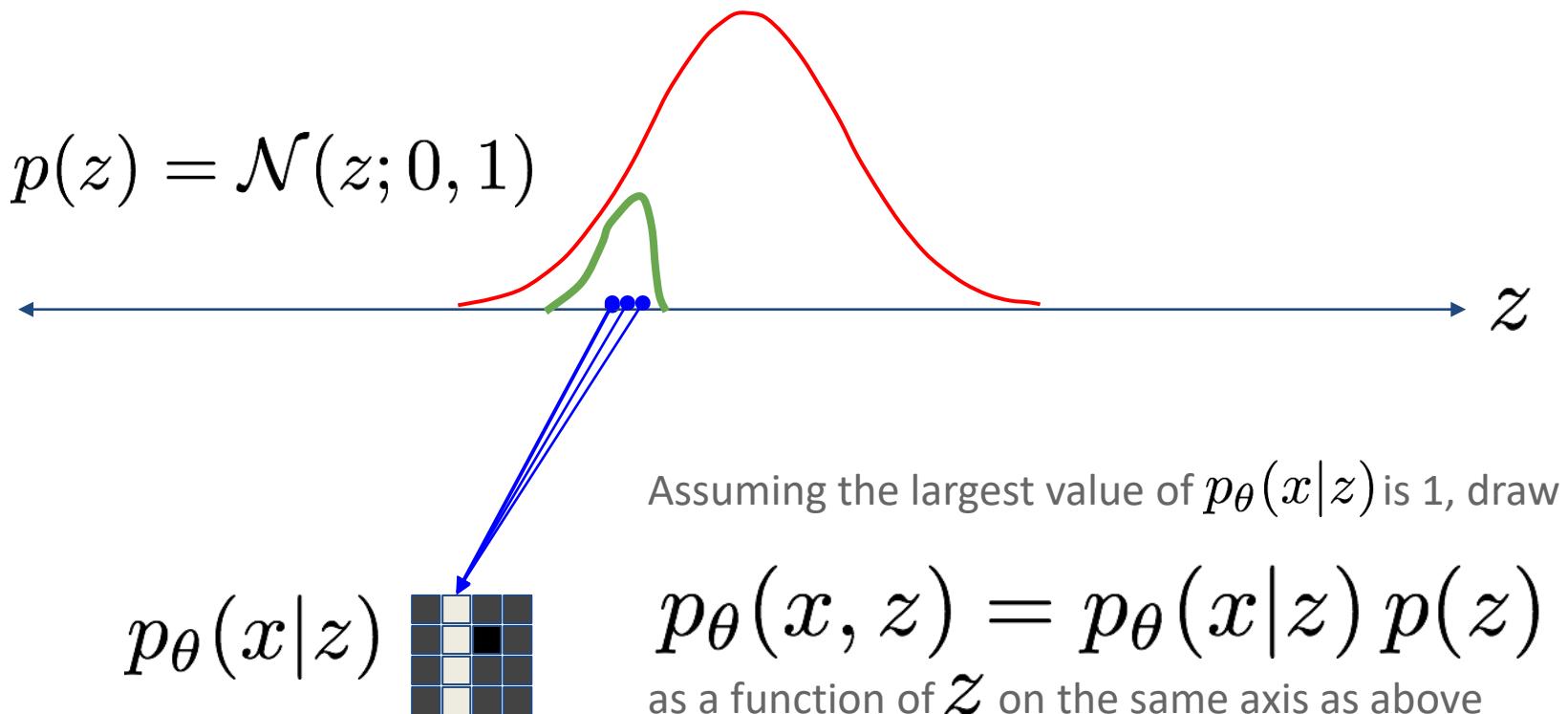
# Inference



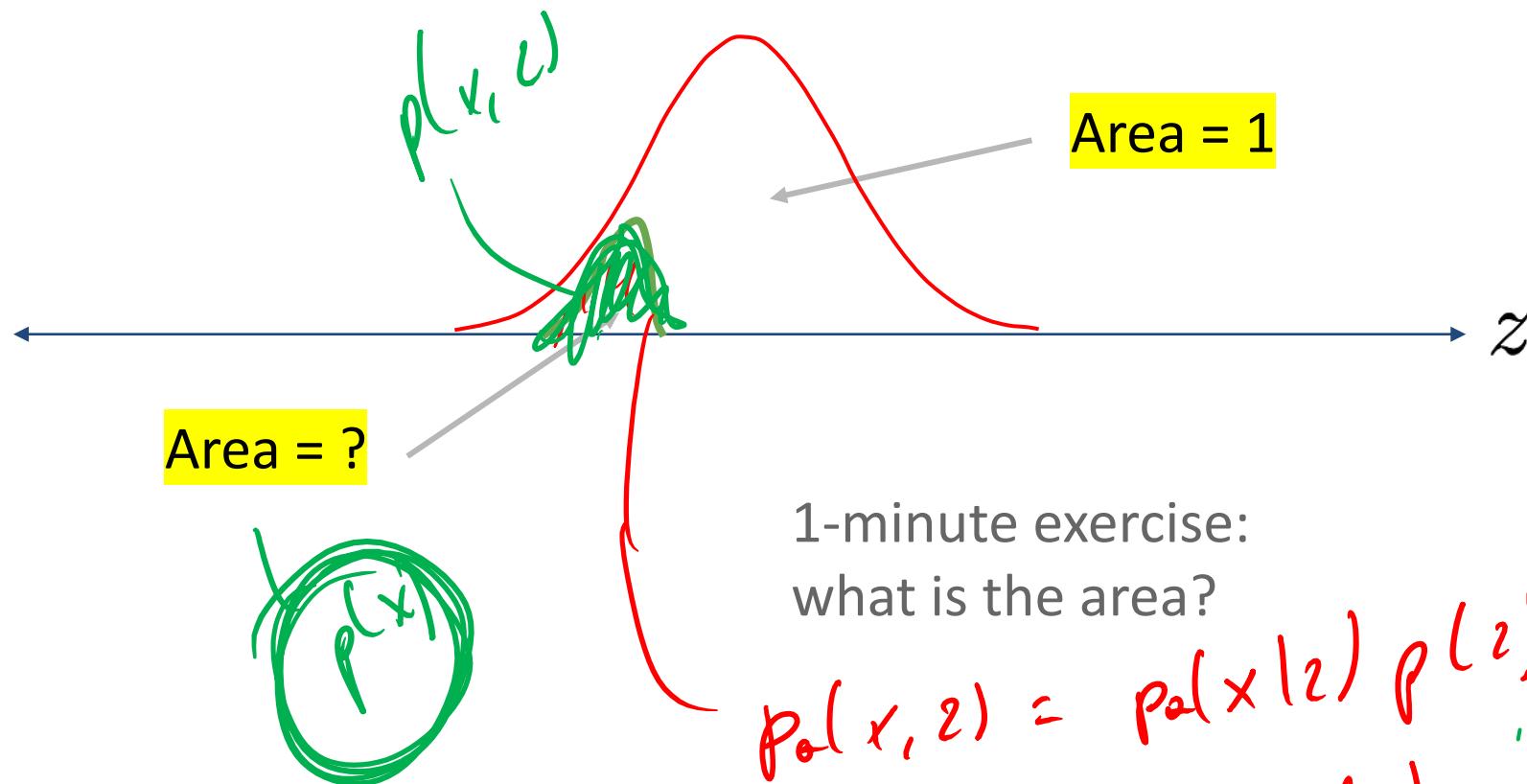
# Exercise



# Joint density (with $x$ observed)



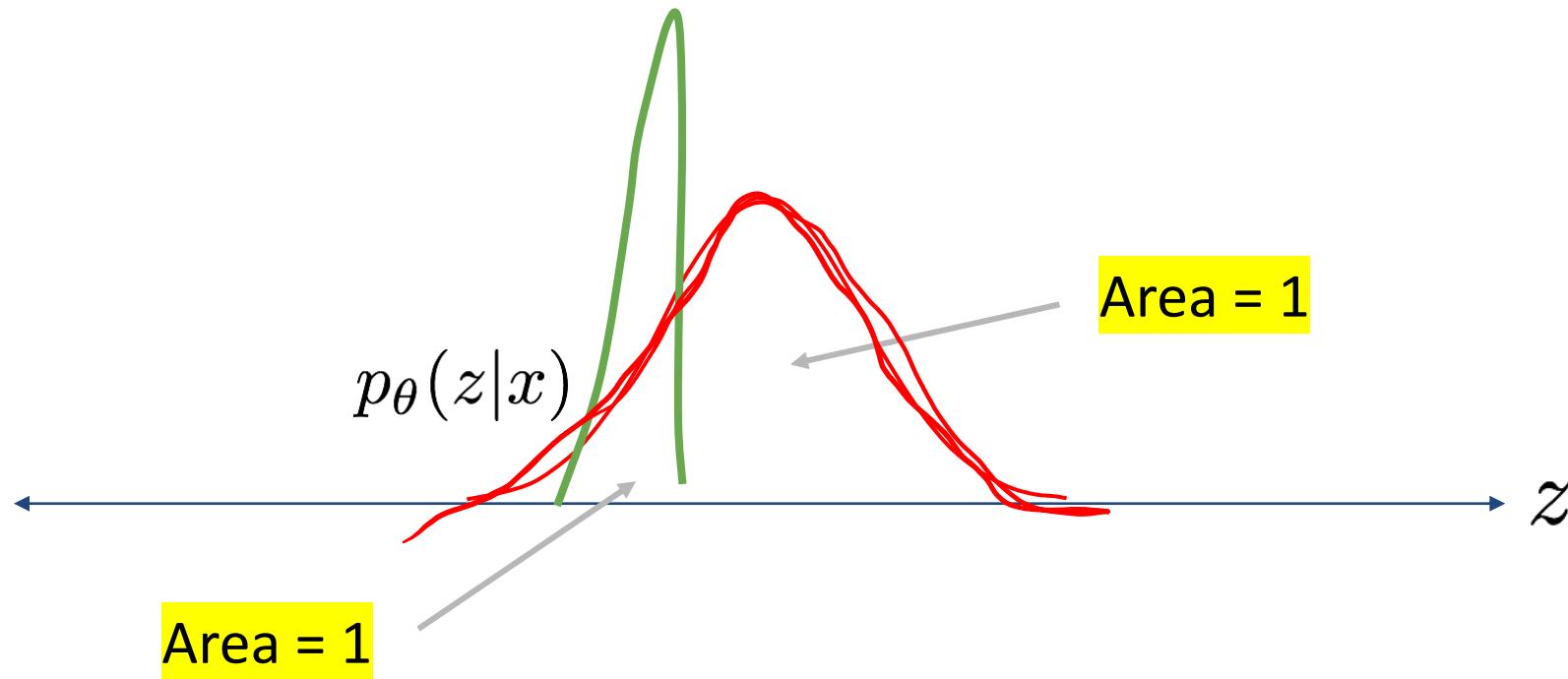
# Joint density (with $x$ observed)



1-minute exercise:  
what is the area?

$$p_0(x, z) = p_0(x|z) p(z)$$
$$\int p_0(x, z) dz = p_0(x) = \text{density in Bayes rule}$$

# Posterior



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

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

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

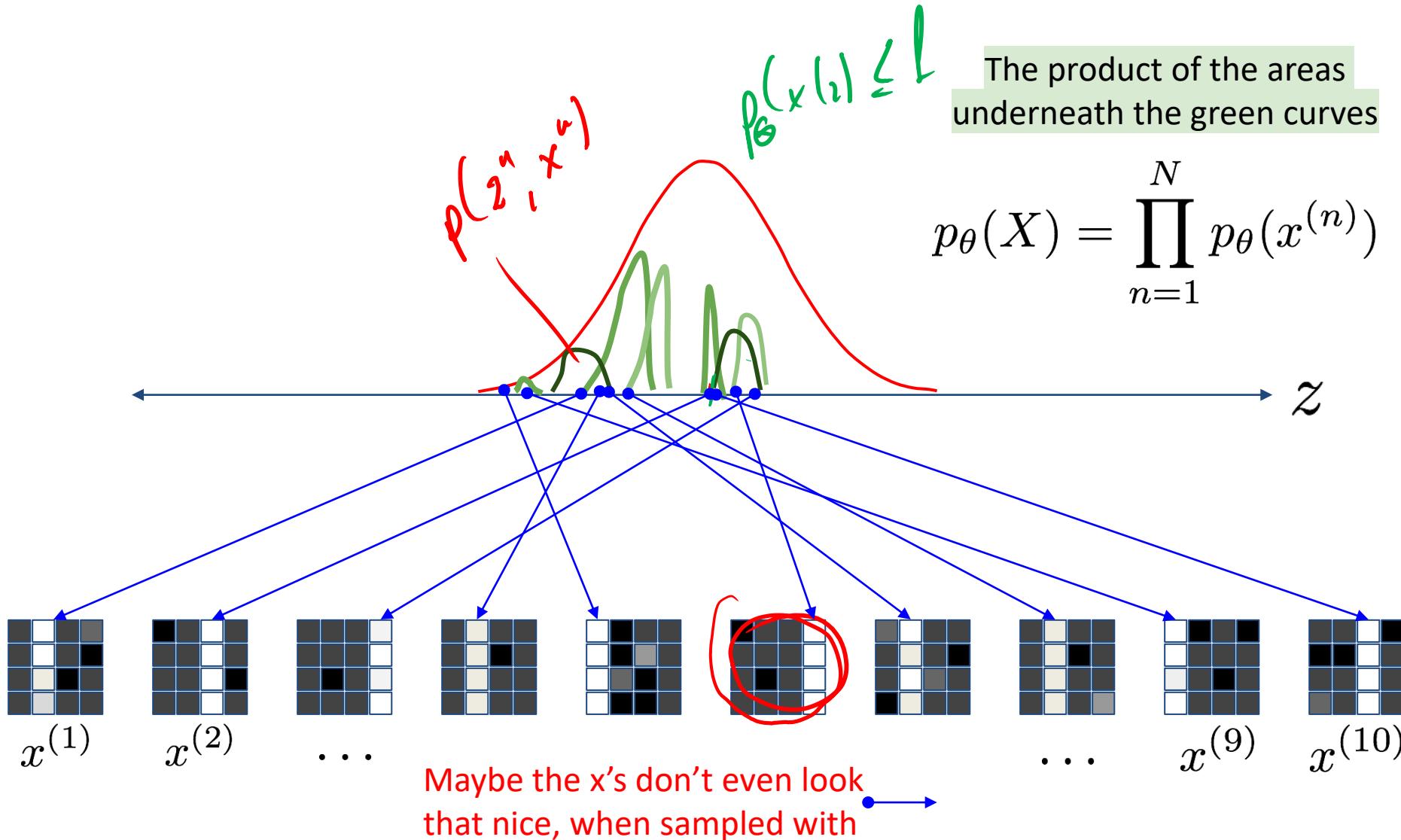
# Evidence of all data points

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

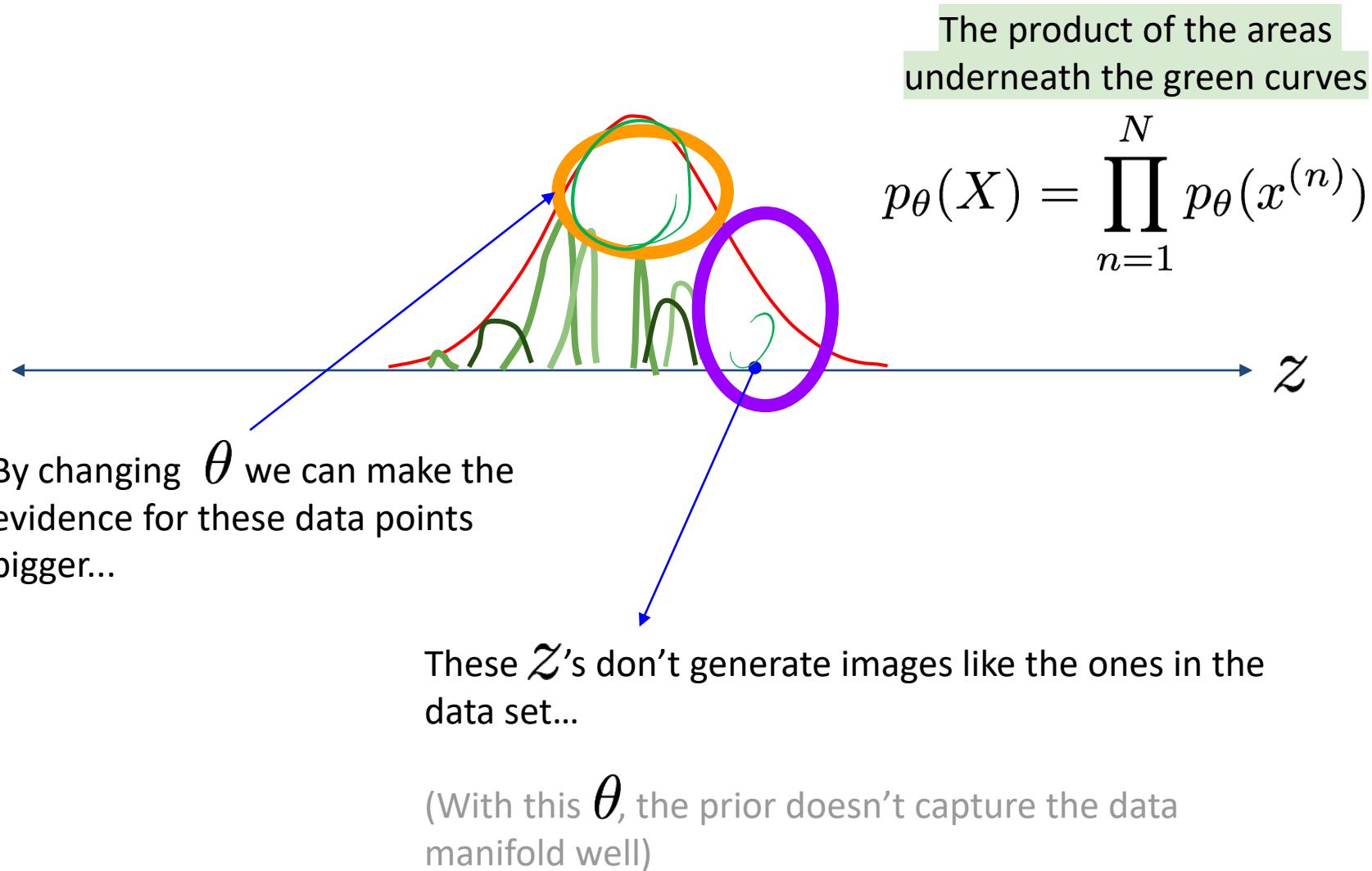


$$\log p_{\theta}(X) = \sum_{n=1}^N \log p_{\theta}(x^{(n)})$$

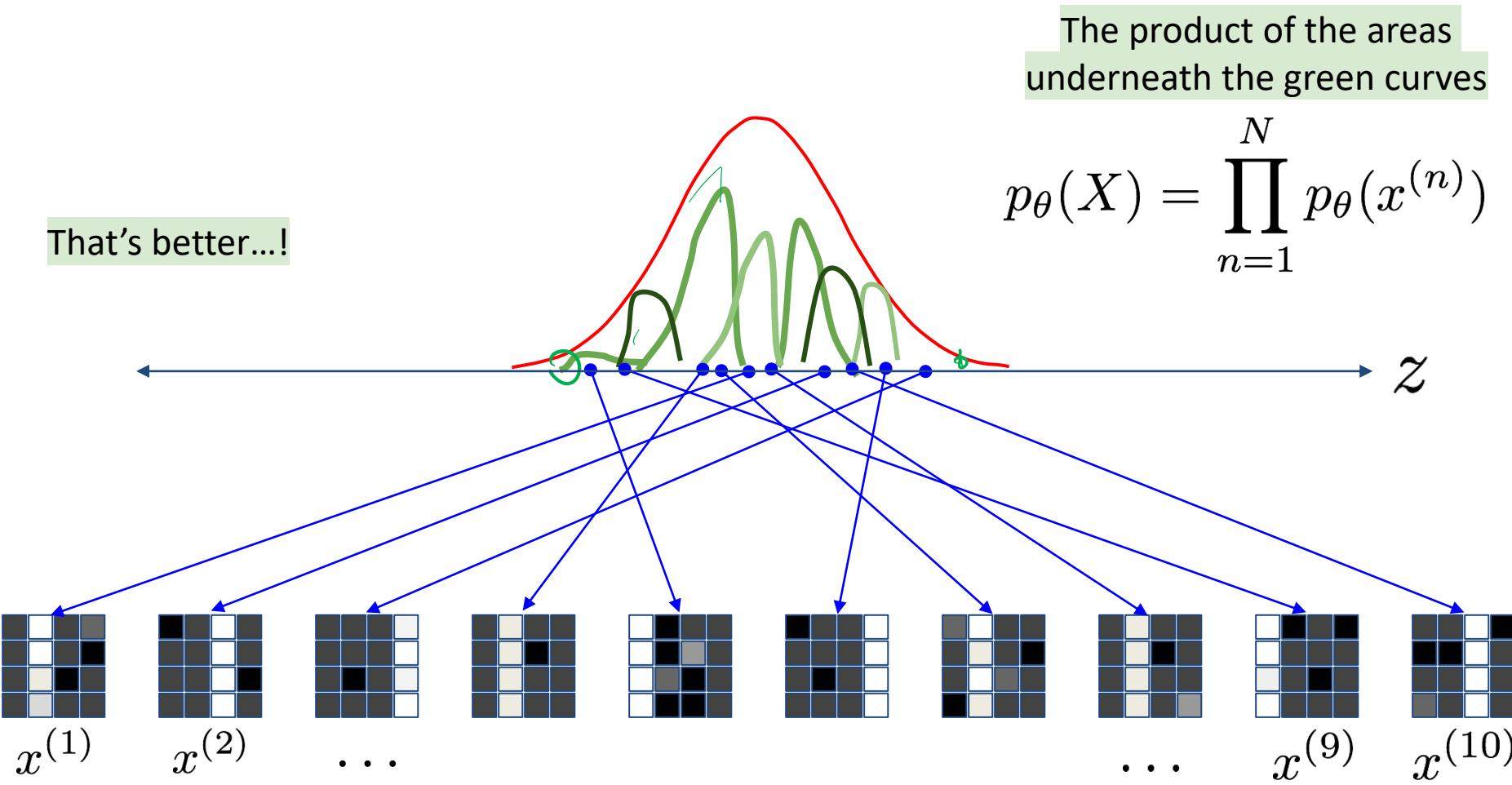
# Evidence for all data points



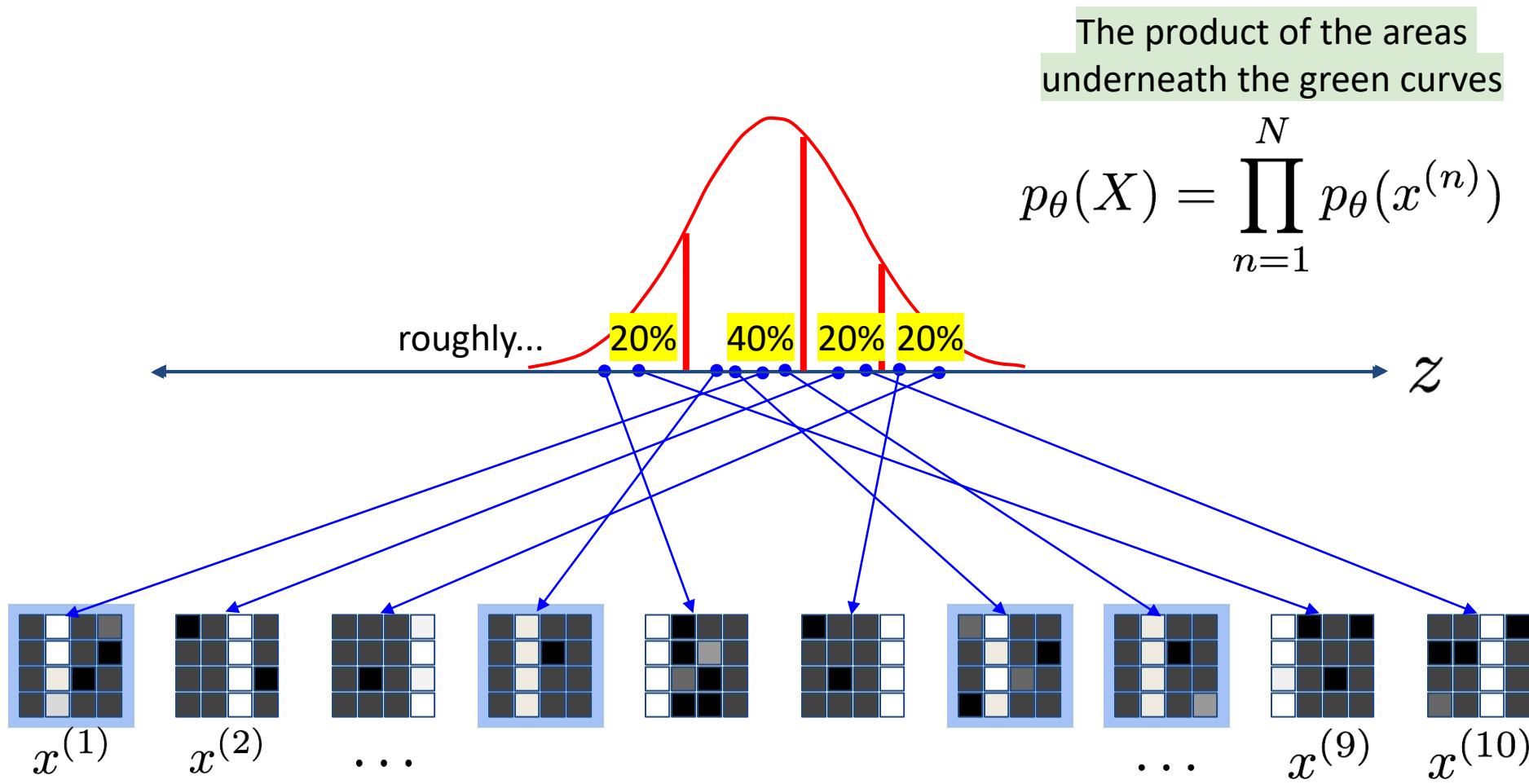
# Maximizing the evidence



# 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  $\log p_\Theta(x)$

$E - d$

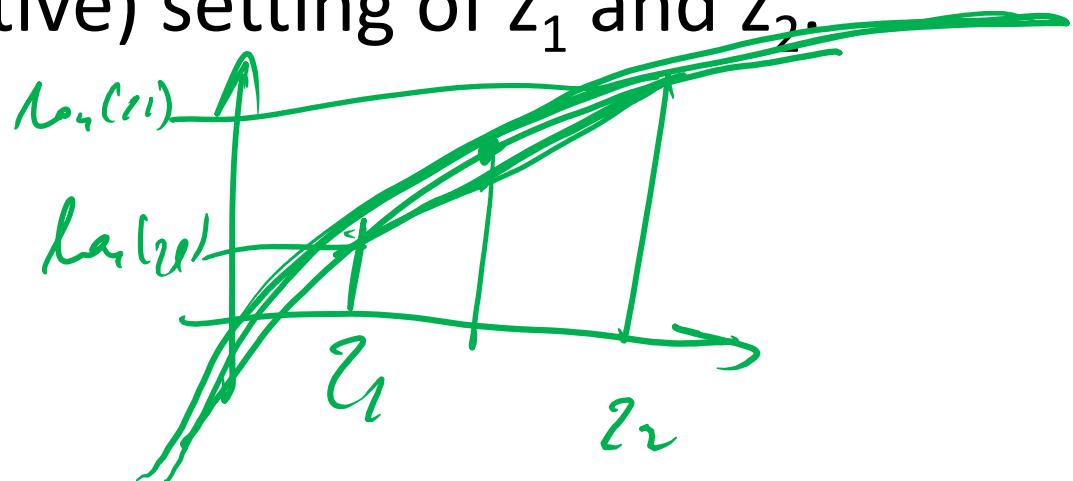
# Exercise

Jensen's inequality

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

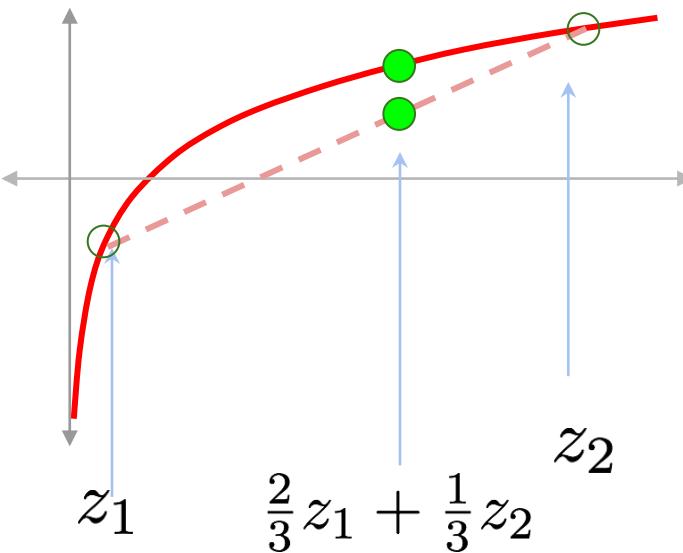
$$\log\left(\frac{2}{3}z_1 + \frac{1}{3}z_2\right) \geq \frac{2}{3}\log z_1 + \frac{1}{3}\log z_2$$

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

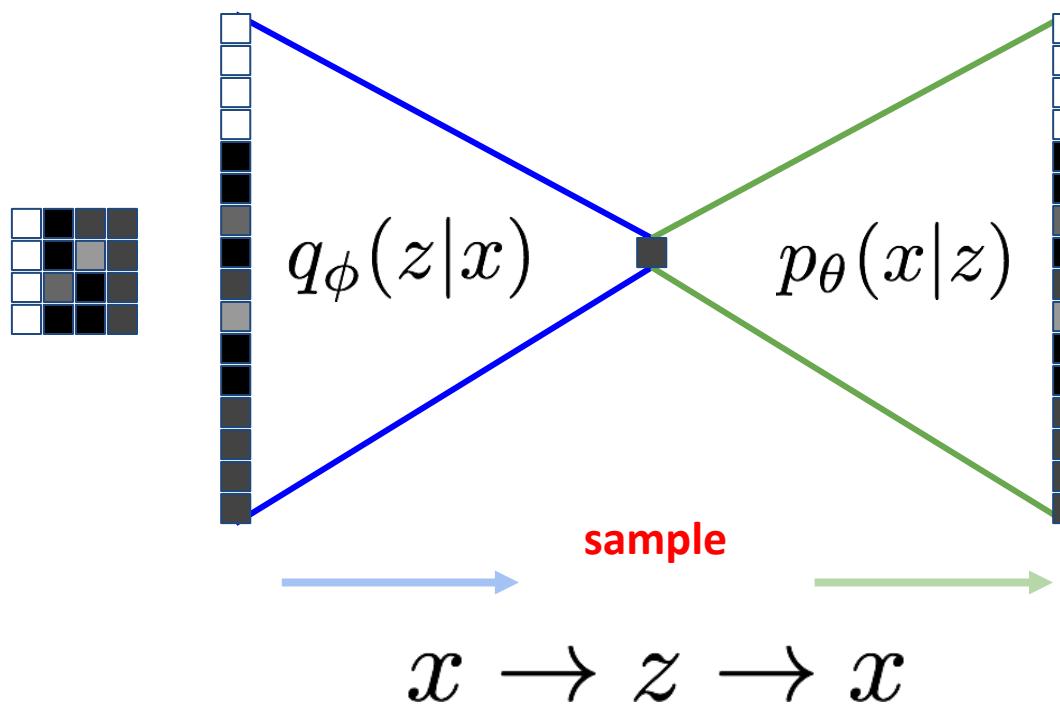
$$\begin{aligned} \log p_{\Theta}(x) &= \log \int dz p_{\Theta}(x, z) = \\ &= \log \int dz q_{\Phi}(z|x) \frac{p_{\Theta}(x, z)}{q_{\Phi}(z|x)} \\ &\geq \int dz q_{\Phi}(z|x) \log \frac{p_{\Theta}(x, z)}{q_{\Phi}(z|x)} \quad \text{works for Jensen} \\ &= \mathbb{E}_{q_{\Phi}(z|x)} \left[ \log \frac{p_{\Theta}(x|z)p(z)}{q_{\Phi}(z|x)} \right] \\ &= \mathbb{E}_{q_{\Phi}(z|x)} [\log p_{\Theta}(x|z)] \\ &\quad - \mathbb{E}_{q_{\Phi}(z|x)} \left[ \log \frac{q_{\Phi}(z|x)}{p(z)} \right] \\ &= \mathbb{E}_{q_{\Phi}(z|x)} [\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z)) \end{aligned}$$

Appl.  $q$   
 $\ll Q(2/\nu) \neq 0$

# ELBO interpretation

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

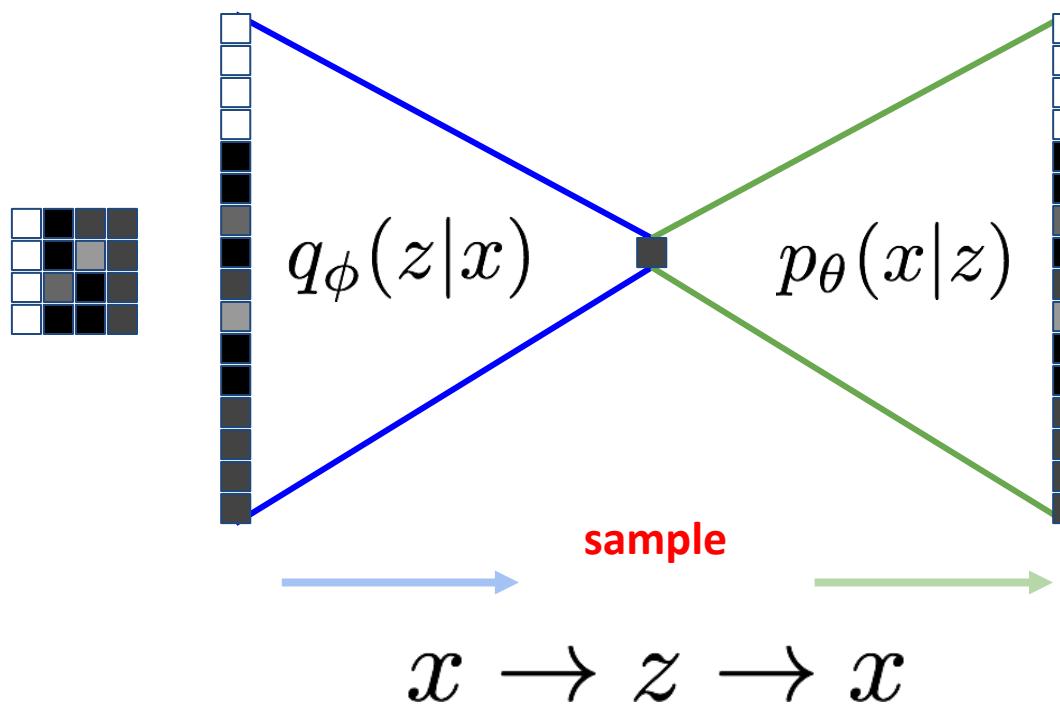
$\mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)]$ : auto-encoding term!



# ELBO interpretation

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

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



Bottleneck

# ELBO interpretation

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:

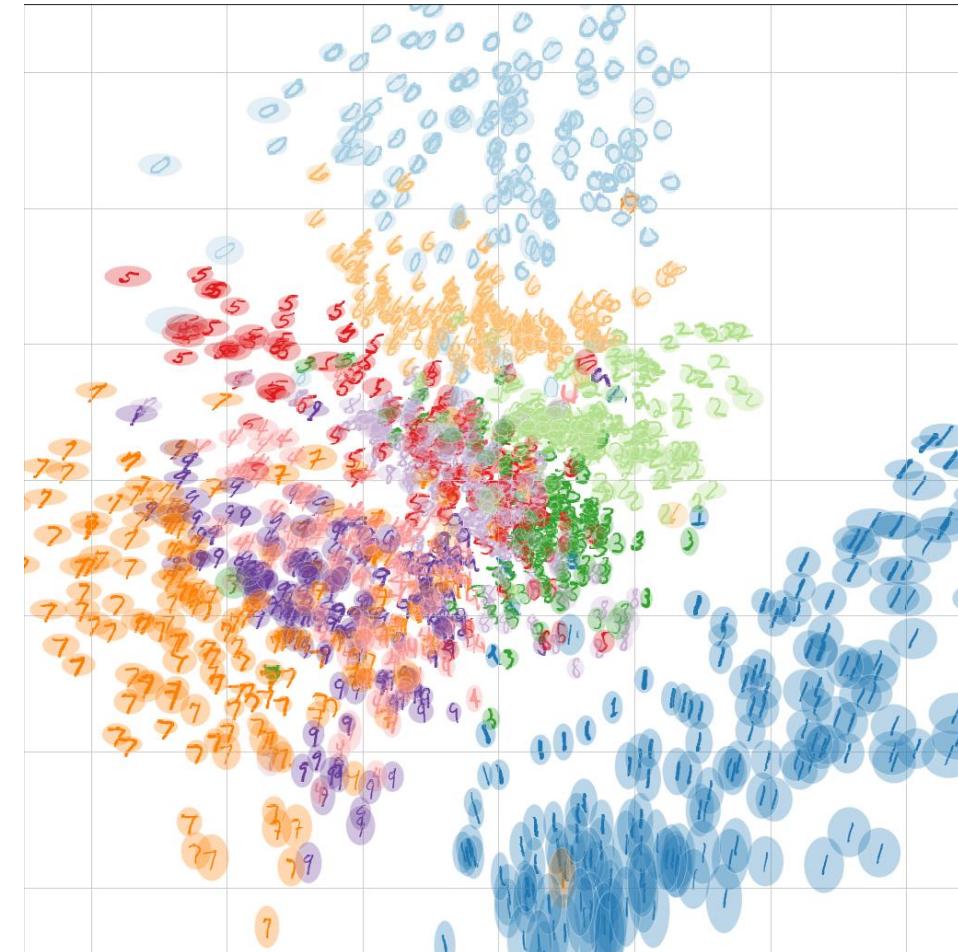
how many nats we transmit about  $x$  in  $q_{\Phi}(z|x)$  rather than  $p(z)$

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:

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

$KL(q_{\Phi}(z|x) \parallel p(z))$  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:

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

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

# ELBO optimization

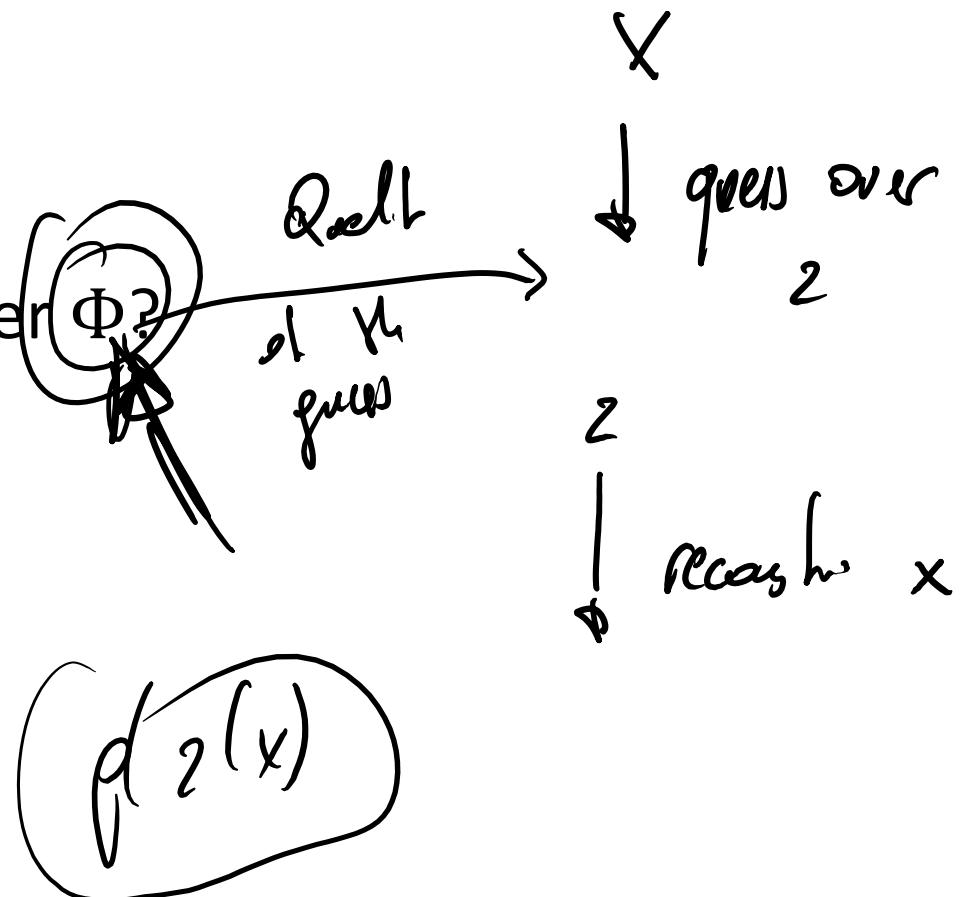
$$\log p_{\Theta}(x) \geq \mathbb{E}_{q_{\Phi}(z|x)} [\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z))$$

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

What does it mean?

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

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

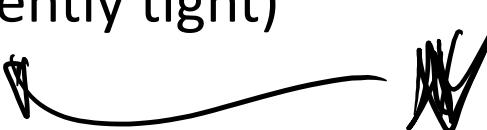
$$\log p_{\Theta}(x) \geq \mathbb{E}_{q_{\Phi}(z|x)}[\log p_{\Theta}(x|z)] - KL(q_{\Phi}(z|x) \parallel p(z))$$

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

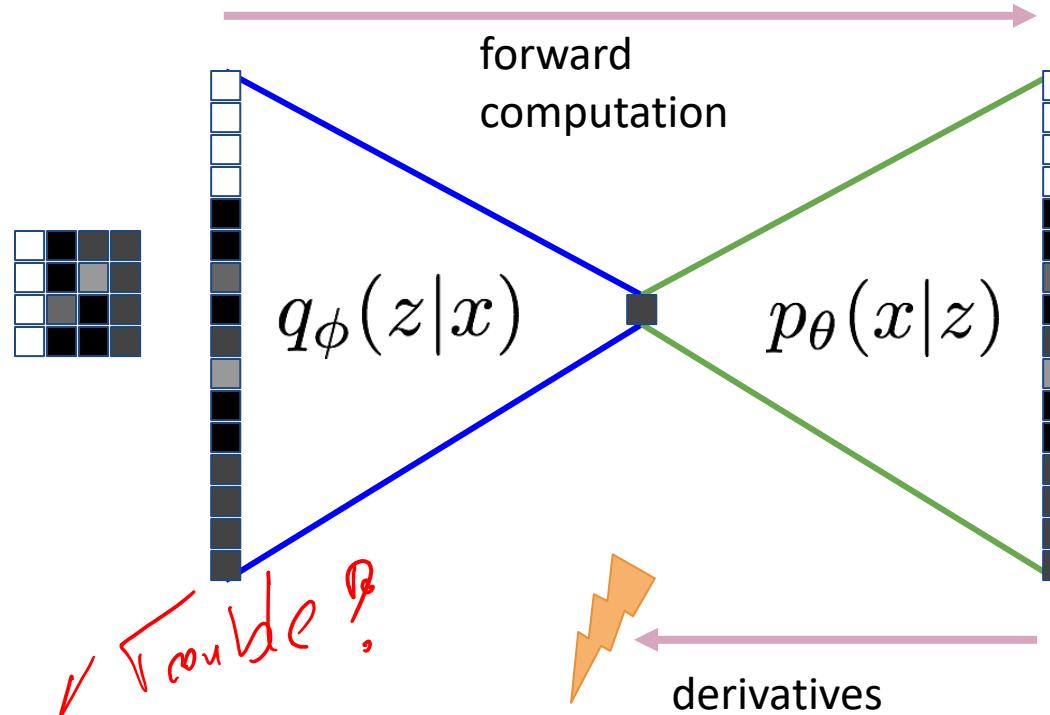
We have shown that for any  $q_\Phi(z|x)$

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

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?



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

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

$\begin{pmatrix} \mu_z + \sigma_z \epsilon \\ \epsilon \end{pmatrix} = g_\phi(x)$

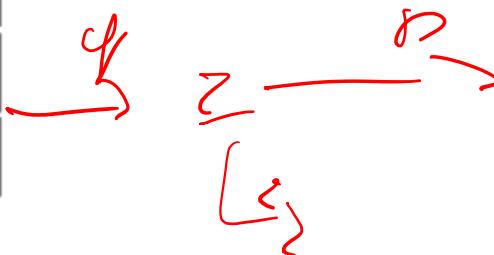
$\xrightarrow{\text{no dep on } \phi} \quad \underbrace{\epsilon = \mu_z + \sigma_z \cdot \epsilon}_{\leftarrow}$

$\epsilon$  is drawn from a fixed distribution.

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

# 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			

~~several  
z >~~

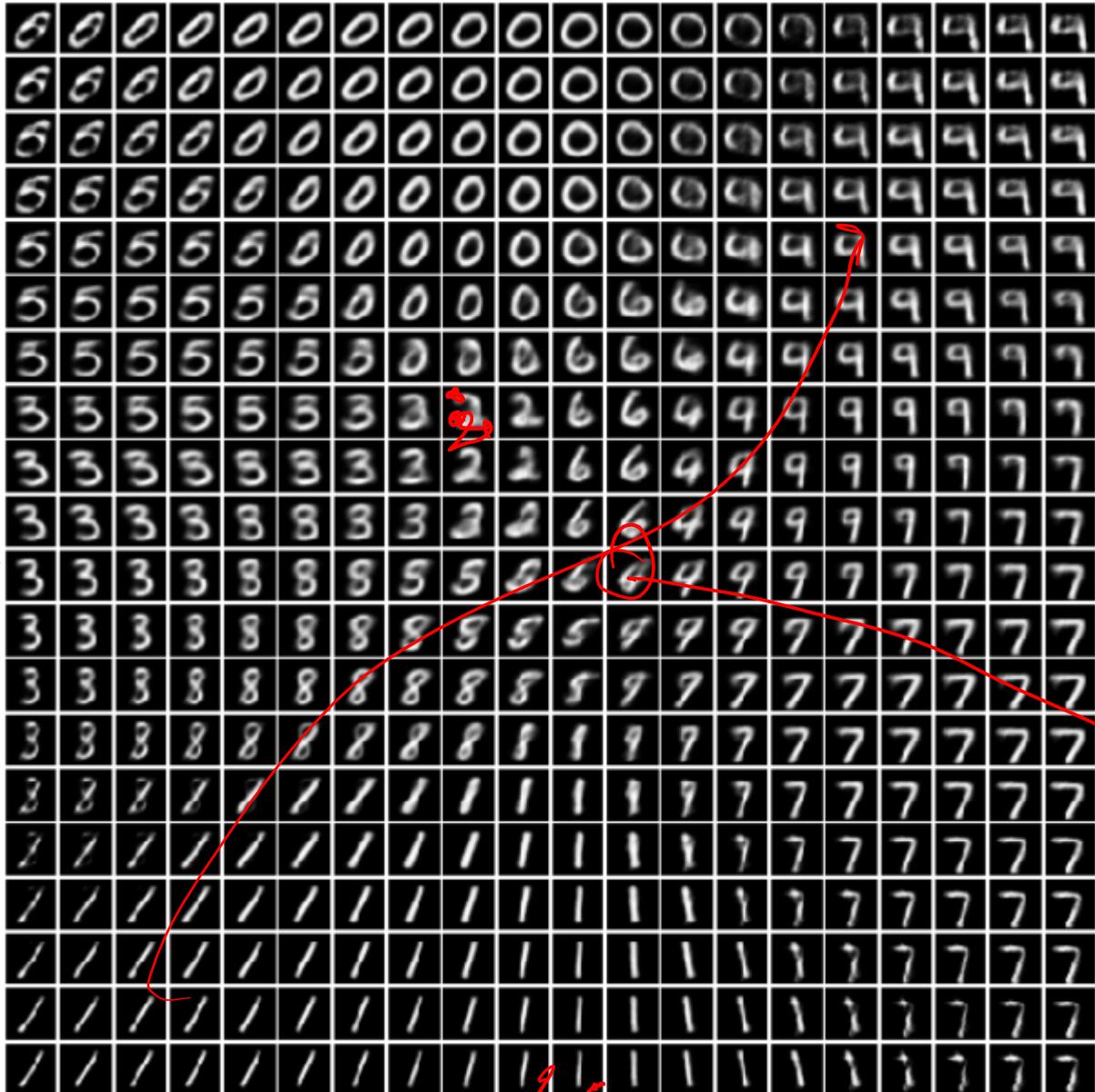
several  
 $z >$

Reconstructions from multiple latent samples

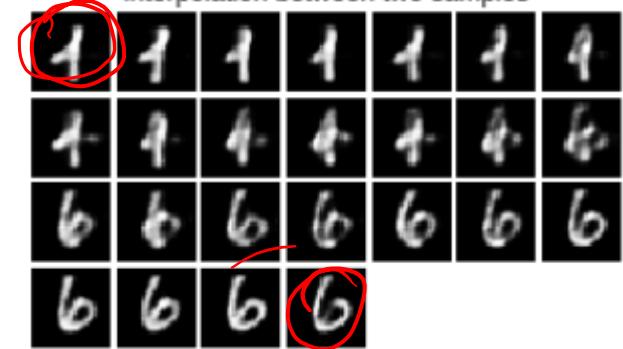
2	2	2	2	2	2
2	2	2	2	2	2
2	2	2	2	2	2
2	2	2	2	2	2

# VAE in action

Reconstructions from a 2D latent space



interpolation between two samples



Dont look  
like a diag.!

# VAE in action

*Gen Samples*

Samples from a 20dimensional VAE

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

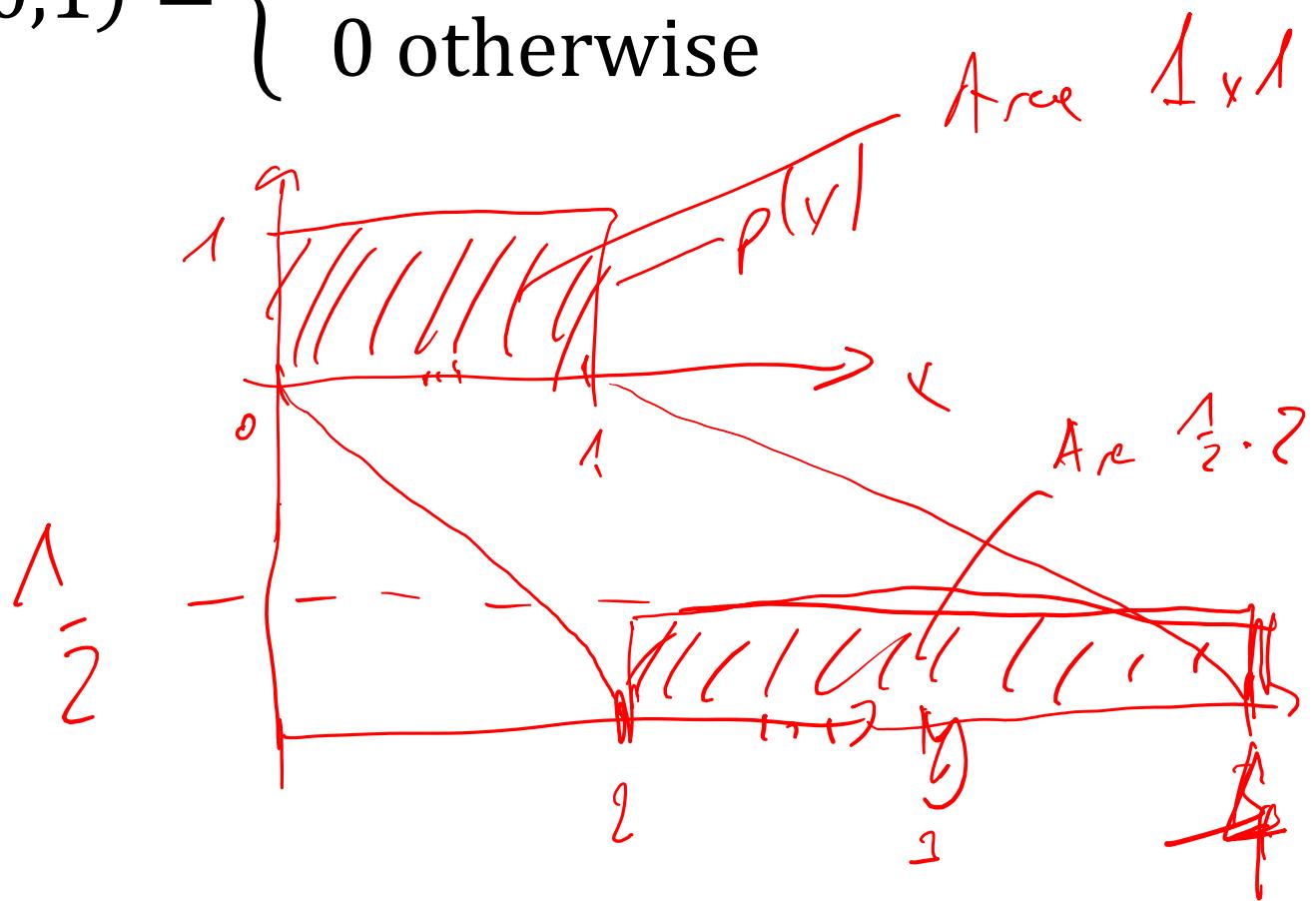
# **NORMALIZING FLOWS**

# Exercise

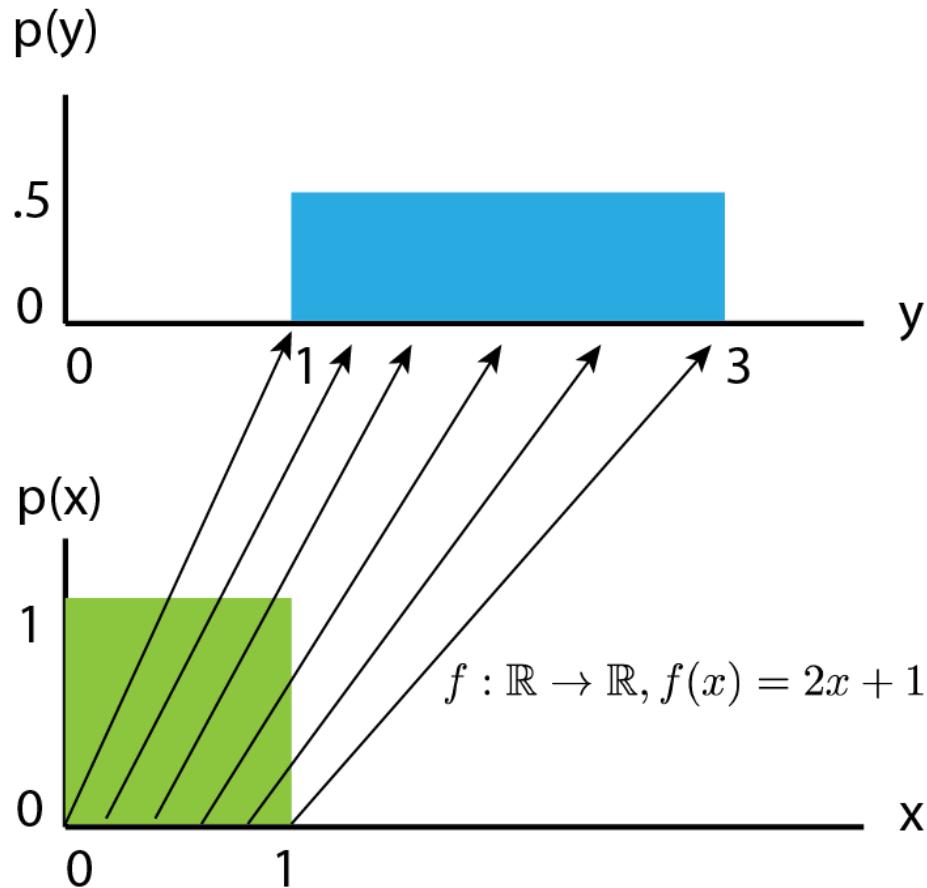
$$p(x) = \text{uniform}(0,1) = \begin{cases} 1 & \text{for } x \in (0,1) \\ 0 & \text{otherwise} \end{cases}$$

$$y = 2 * x + 1$$

$$p(y) = ?$$



# Transforming distributions



$y$  also has a  
uniform  
distribution!

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

# The Jacobian: stretching space

Let  $y = f(x)$

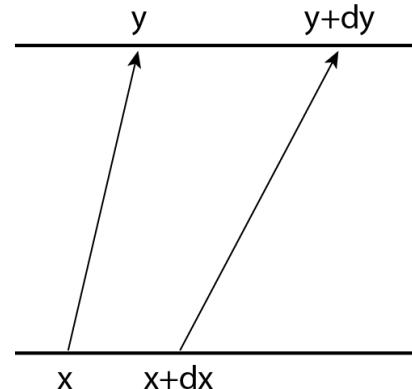
Take a range  $(x, x + \Delta x)$

Question:

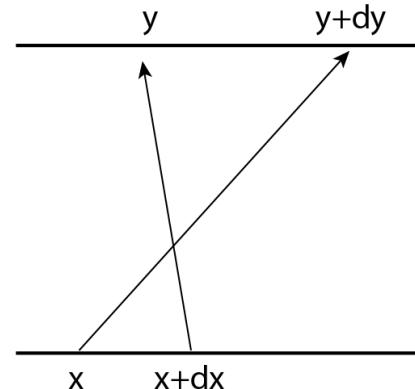
How long is this range?

$\Delta x$

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



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



Question:

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

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

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

Dot of

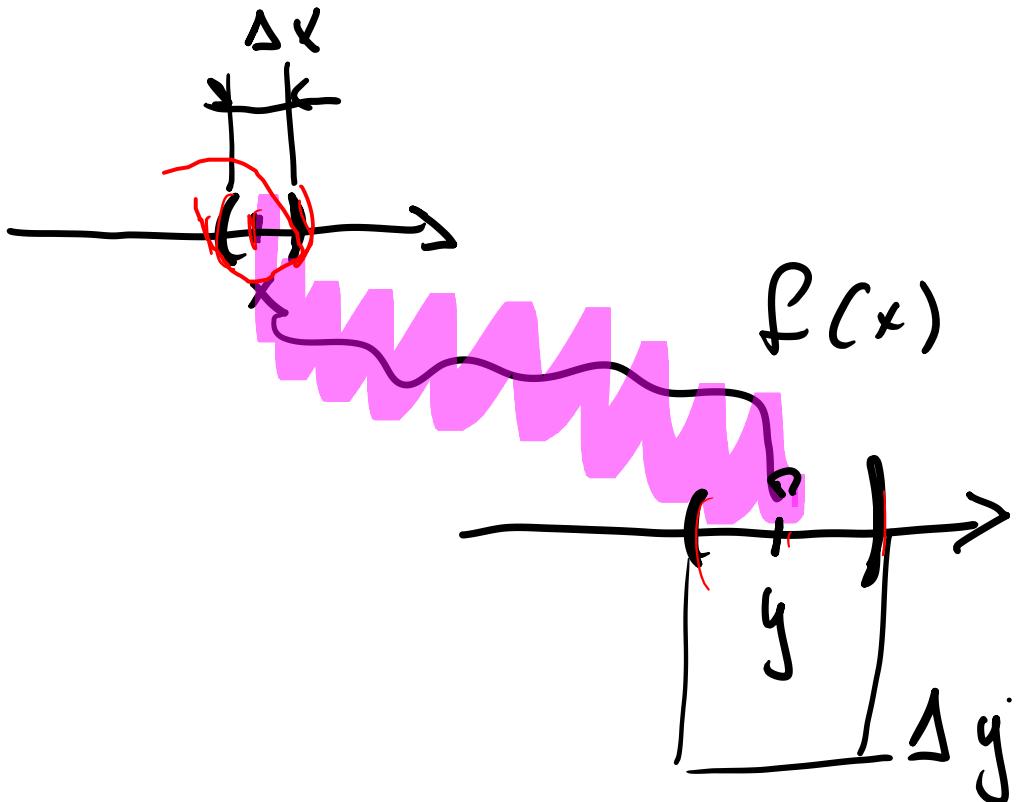
Jacobian

# Change of variables

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

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

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



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:

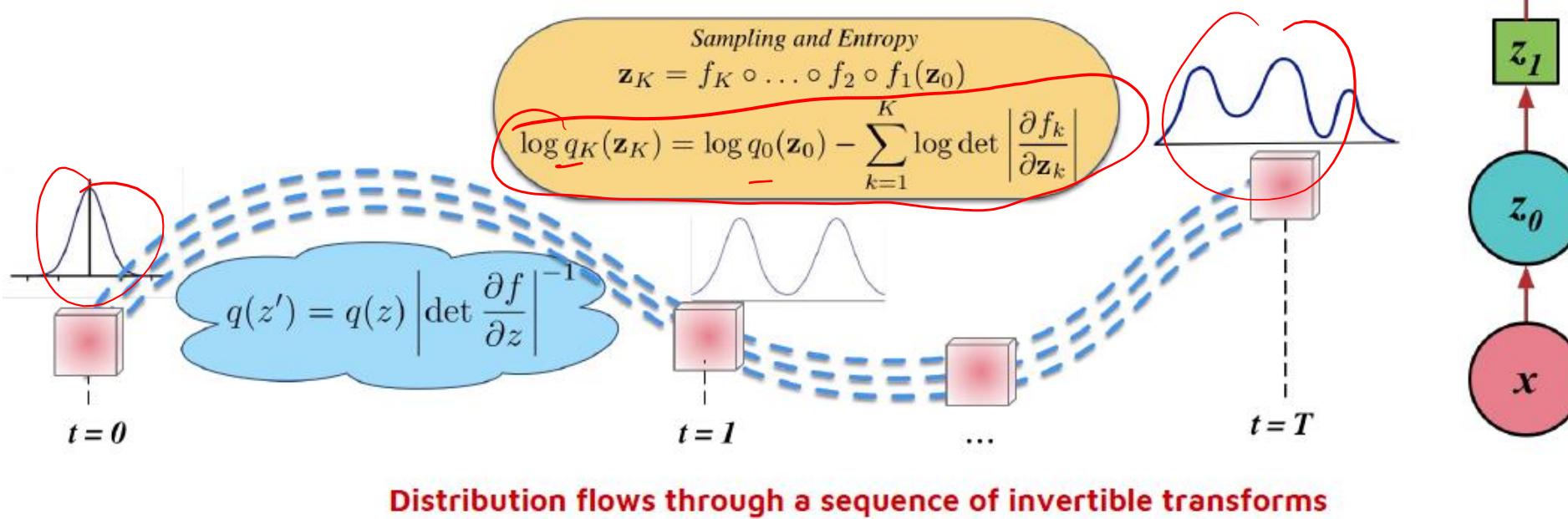
- $f$  is easy to exactly invert  
 $f$  and  $f^{-1}$  form an exact auto-encoder!
- $\det \frac{\partial f(x)}{x}$  is easy to compute

=> We need special  $f$ !

# Normalizing Flows

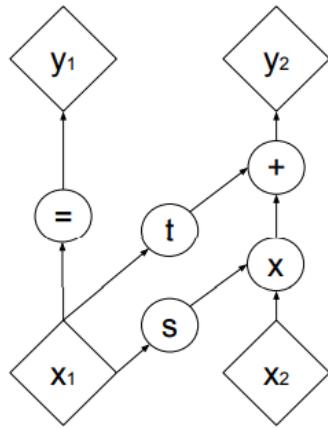
Exploit the rule for change of variables:

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

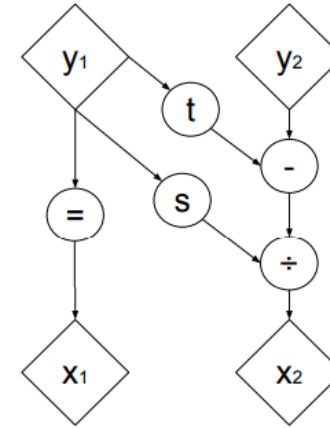


Rezende and Mohamed, 2015

# A special form of $f$



(a) Forward propagation



(b) Inverse propagation

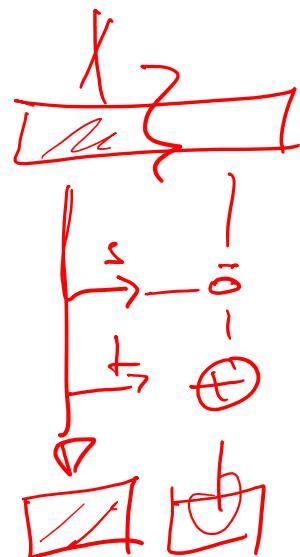
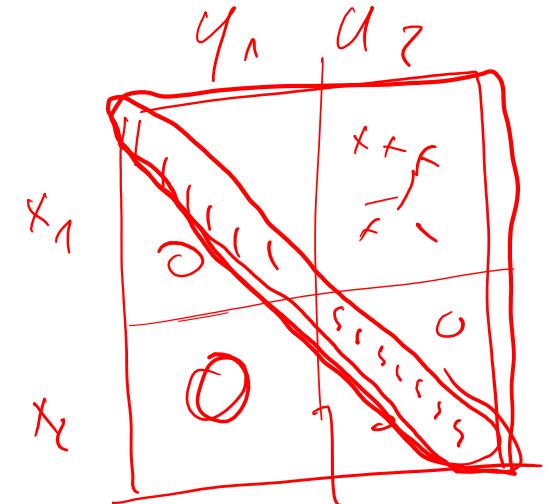
$$x_1, x_2 = \text{split}(x)$$

$$\underline{y_1 = x_1}$$

$$y_2 = x_2 s(x_1) + t(x_1)$$

Trivial to invert!

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



# Normalizing flows in action

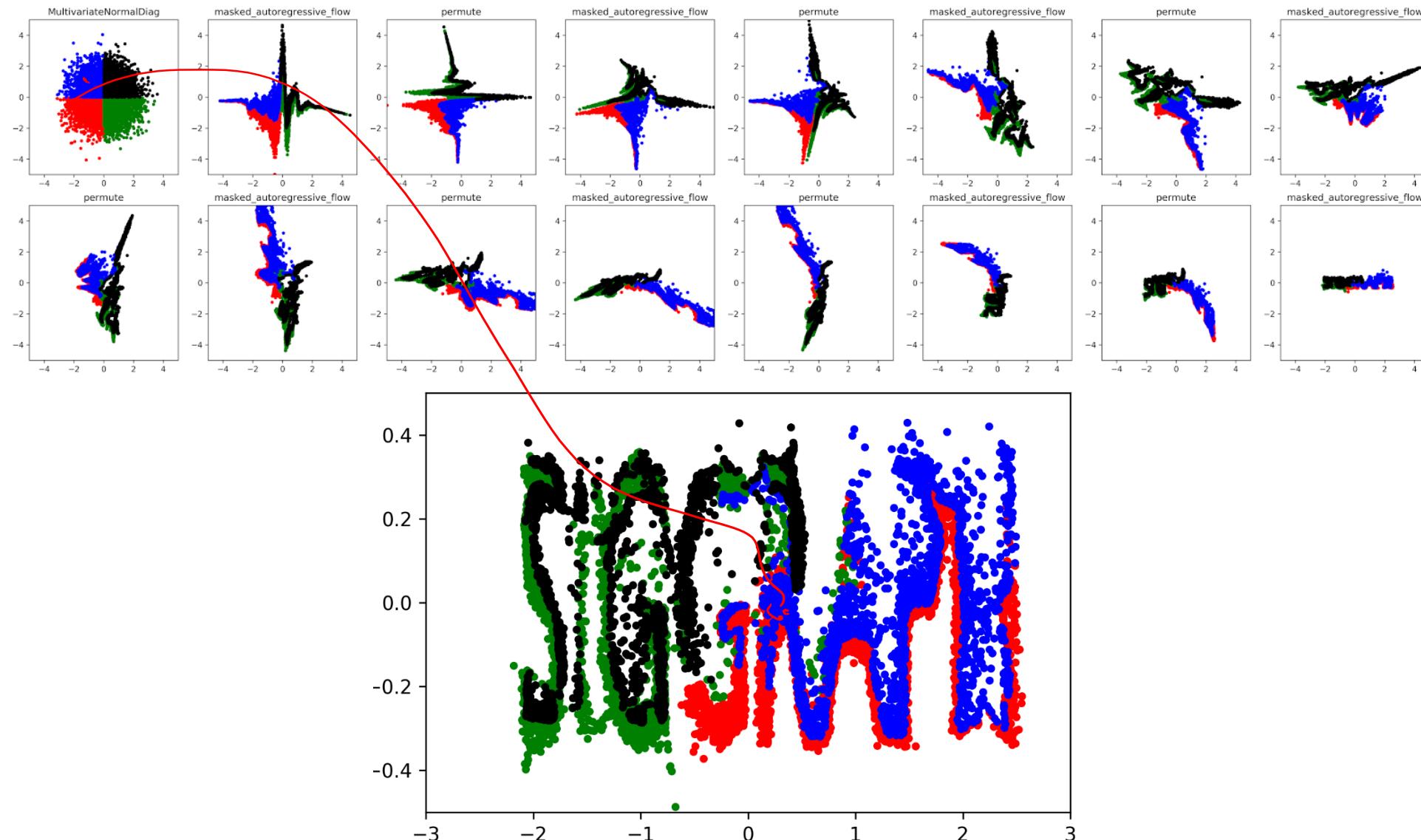
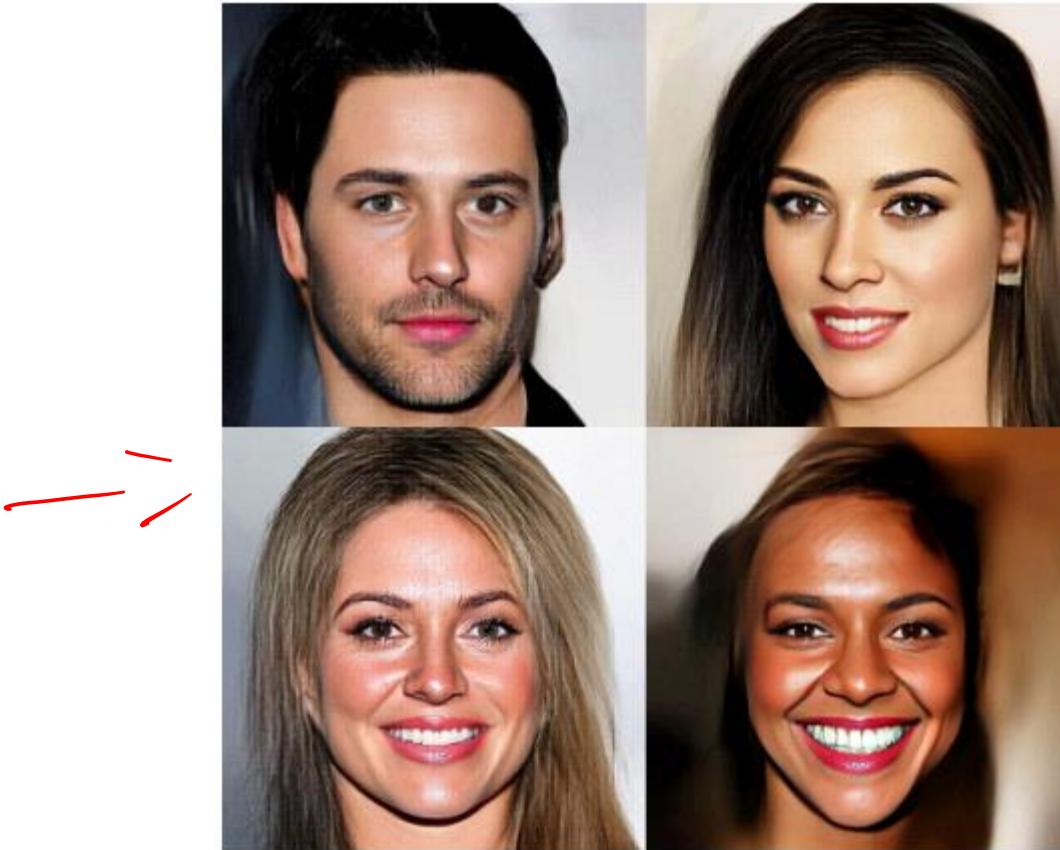
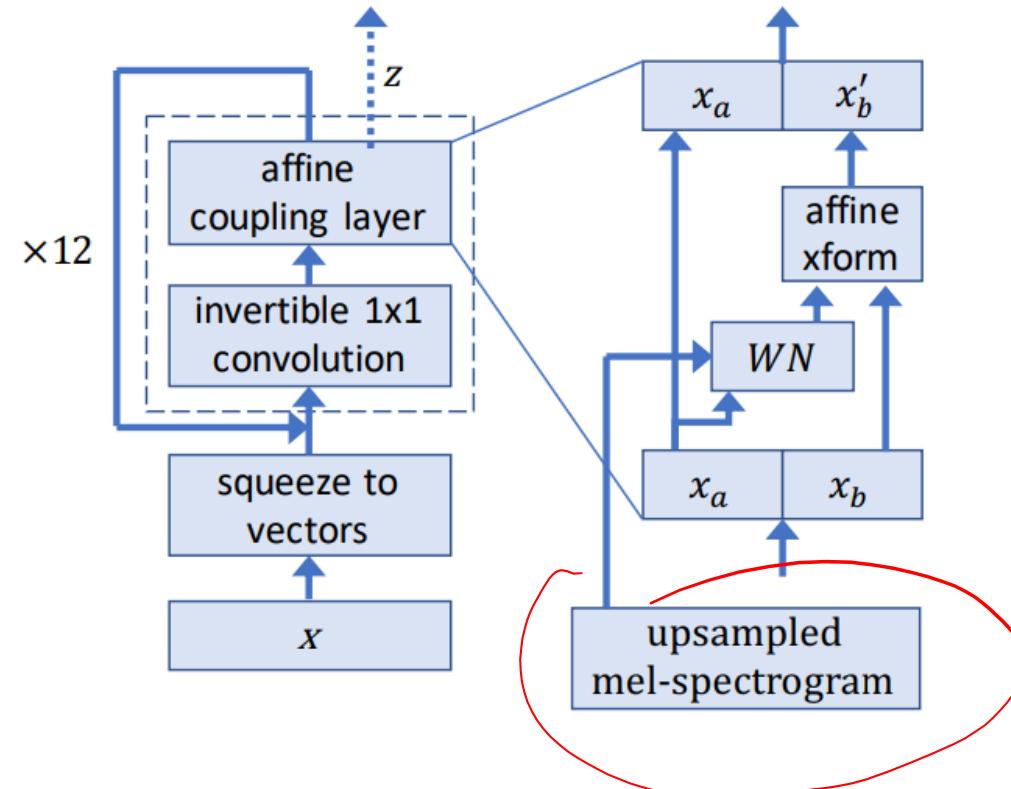


Image credit Eric Jang, <https://blog.evjang.com/2018/01/nf2.html>

# Normalizing flows in action



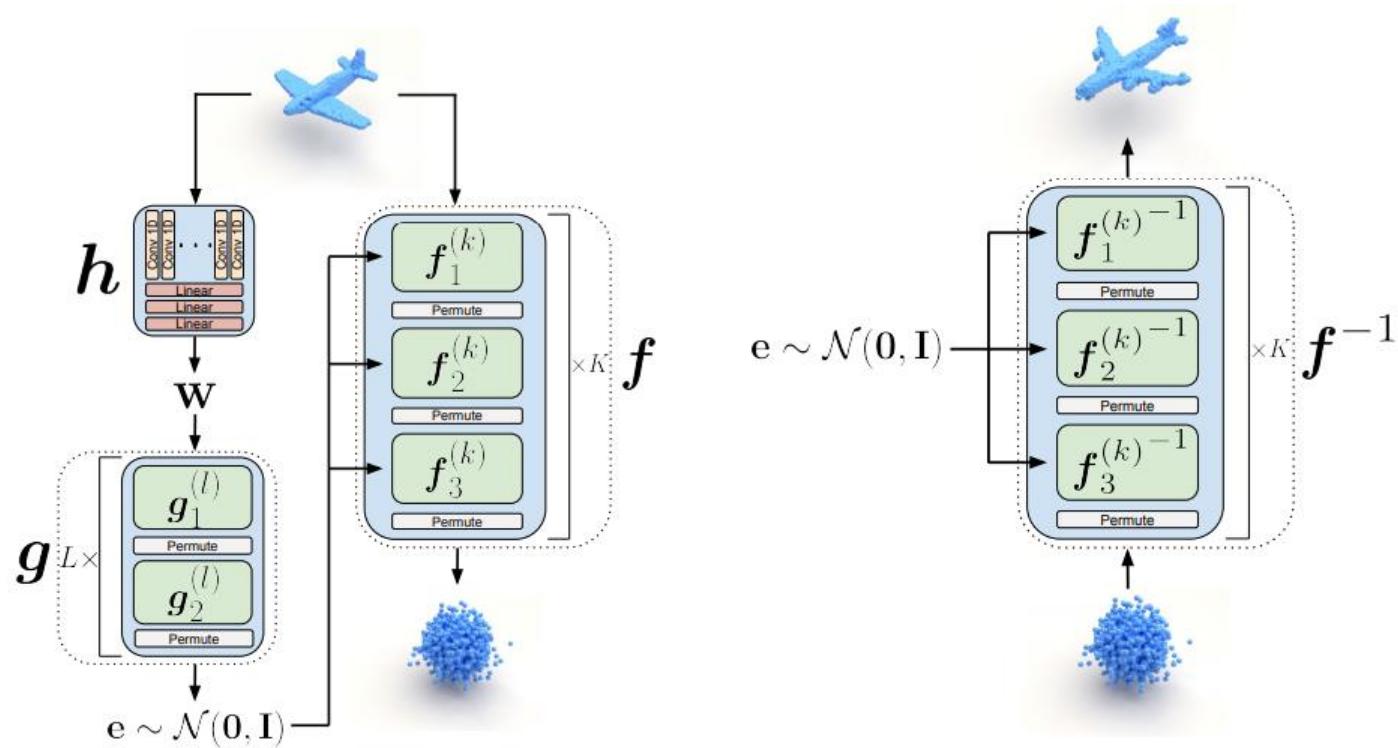
Kingma et al, “GLOW”



Prenger et al, “WAVEGLOW”

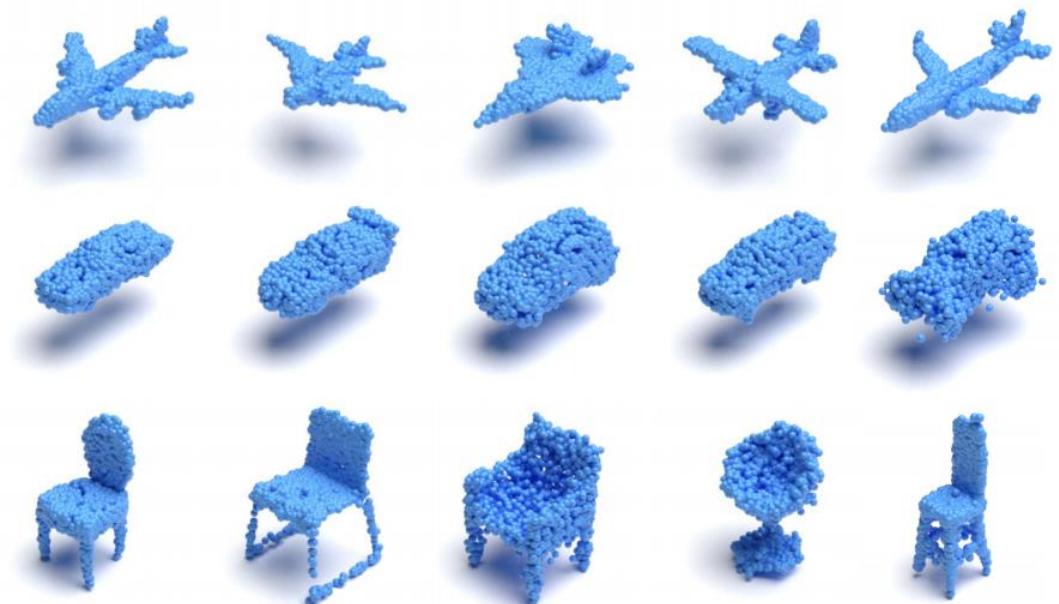
# Normalizing Flows in Action

A generative model for point clouds:

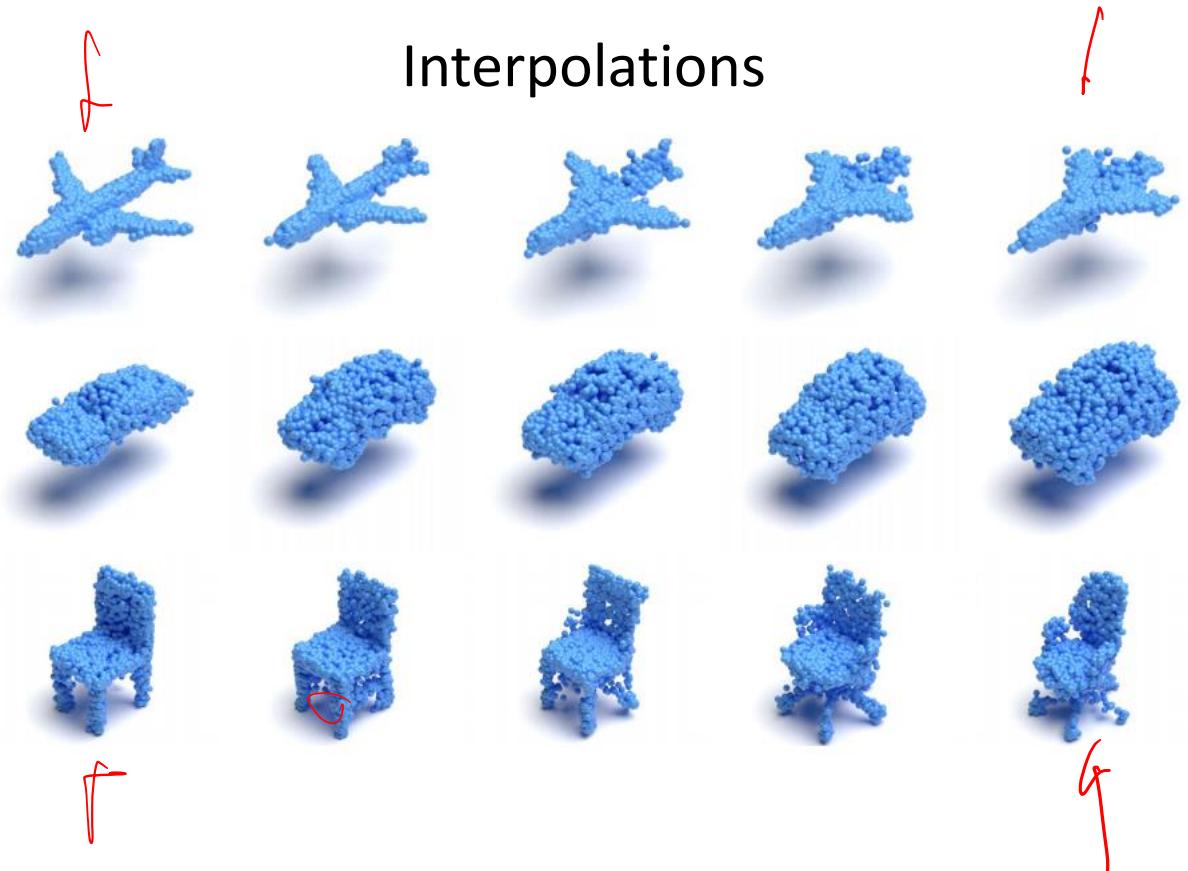


# Normalizing Flows in Action

Samples



Interpolations

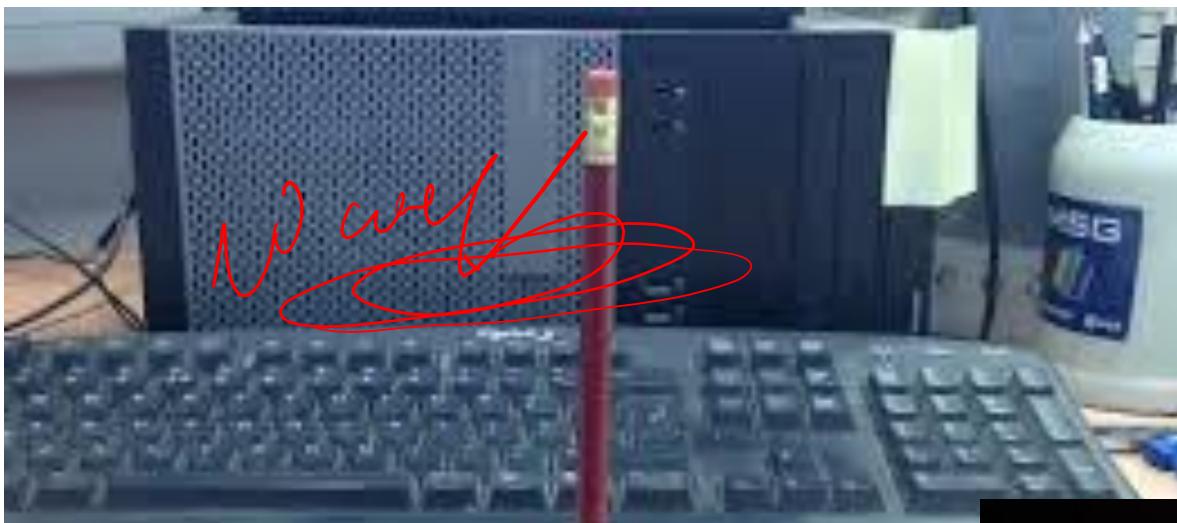


**GAN**

# What will happen next?

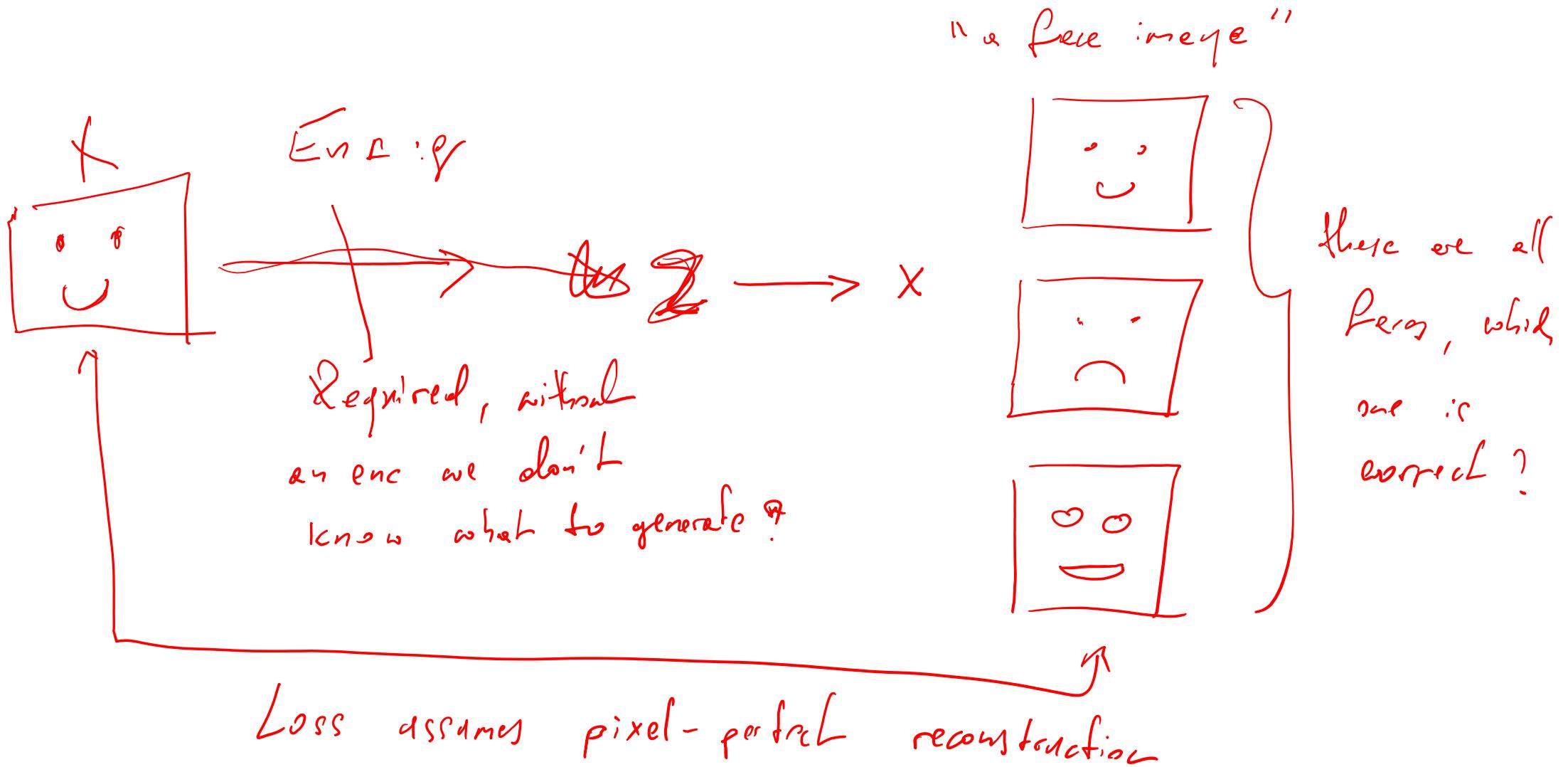
$f_0$

$p(f_2 | f_0, f_1)$



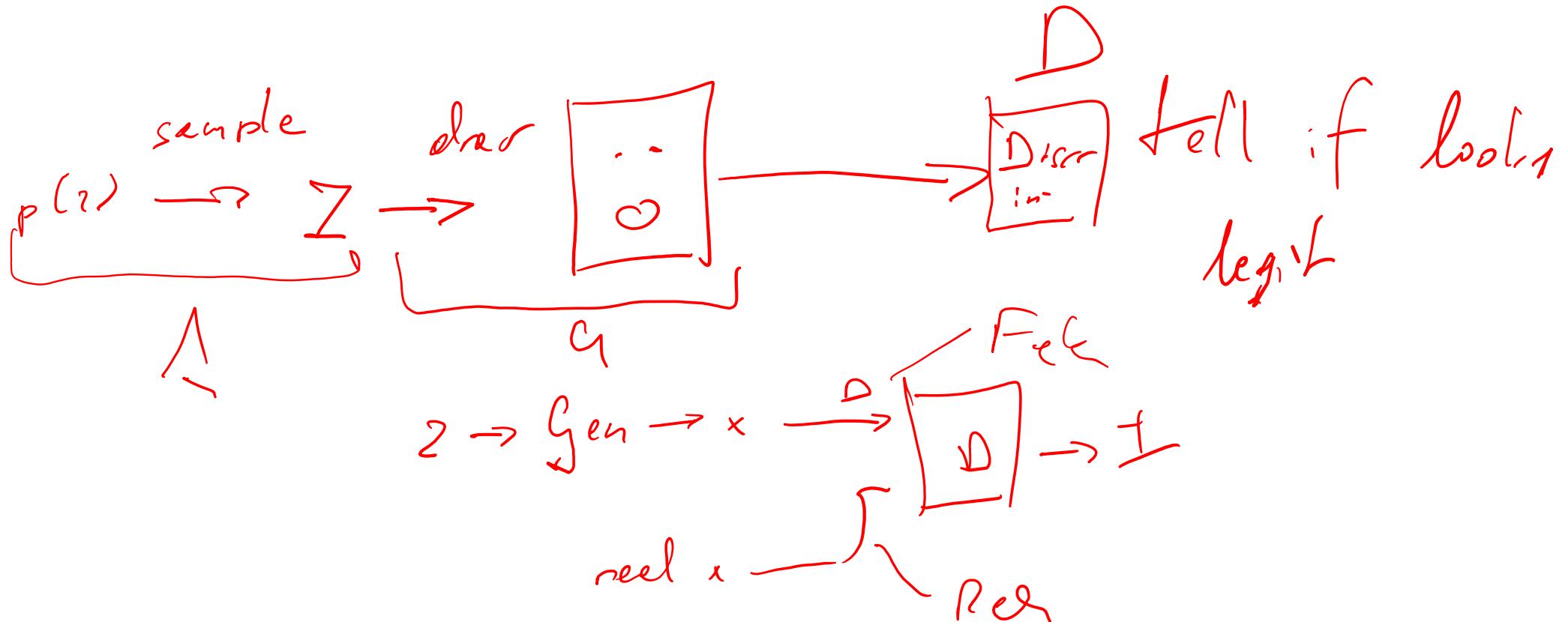
shutterstock.com • 767078788

# Trouble with autoencoders

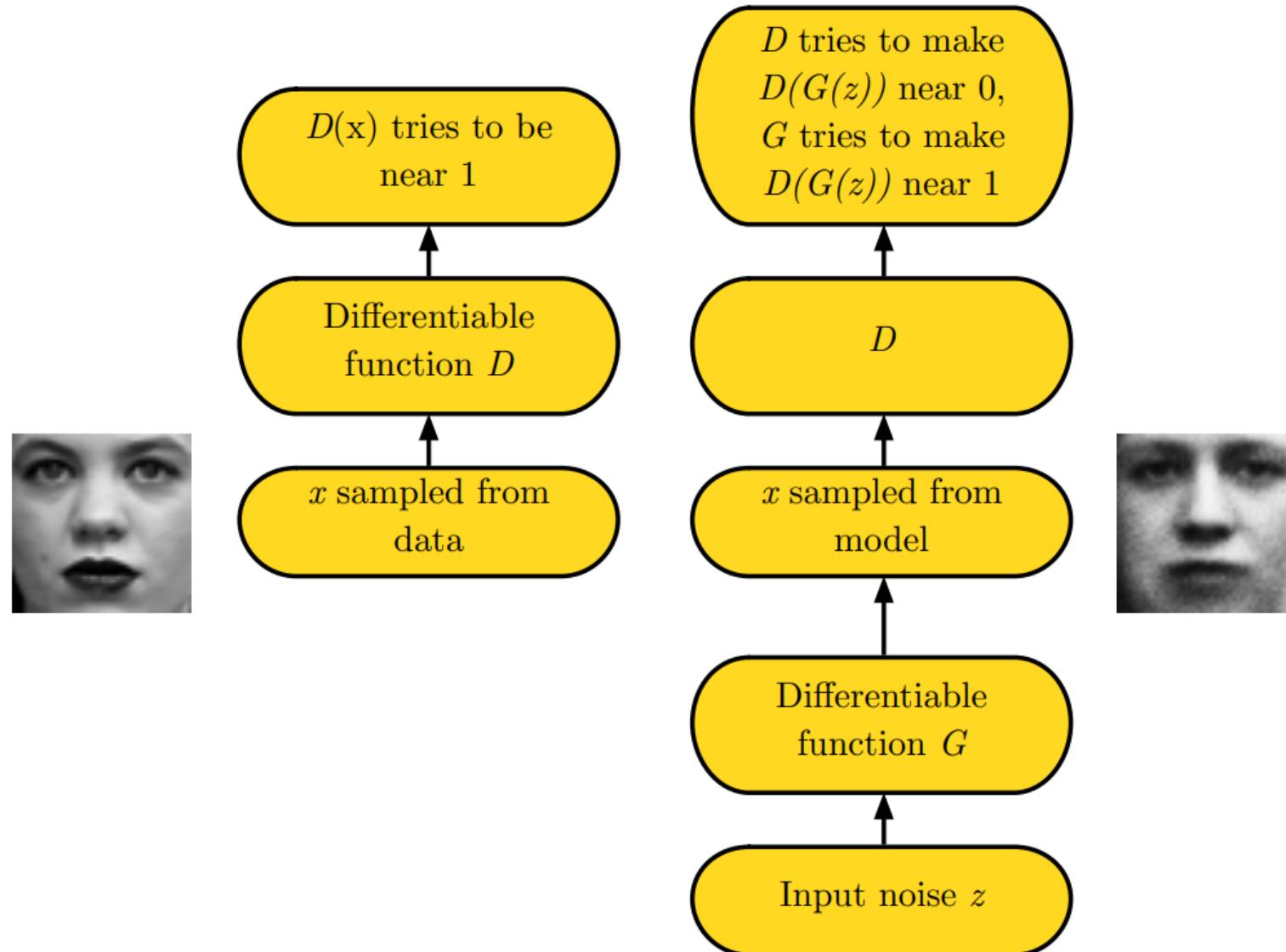


# GAN idea

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



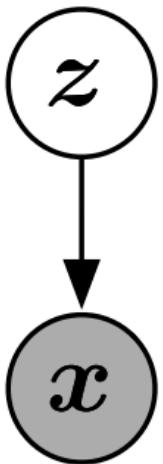
# Adversarial Nets Framework



(Goodfellow 2016)

# Generator Network

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



- 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-Info*

- 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

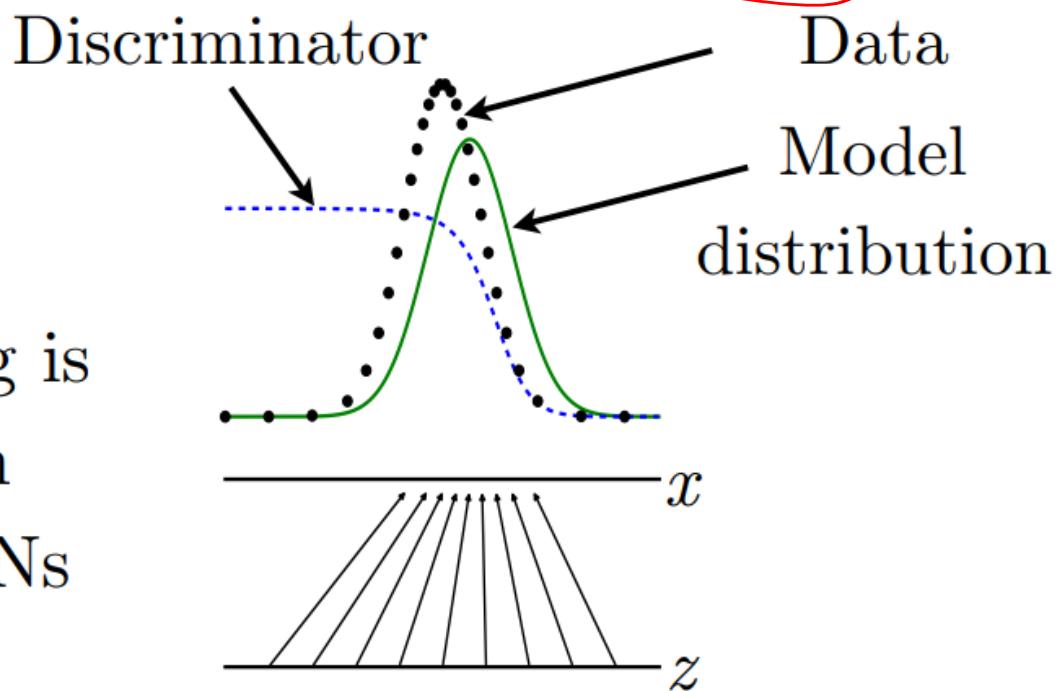
$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$
$$J^{(G)} = -J^{(D)}$$

- 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?

# Discriminator Strategy

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

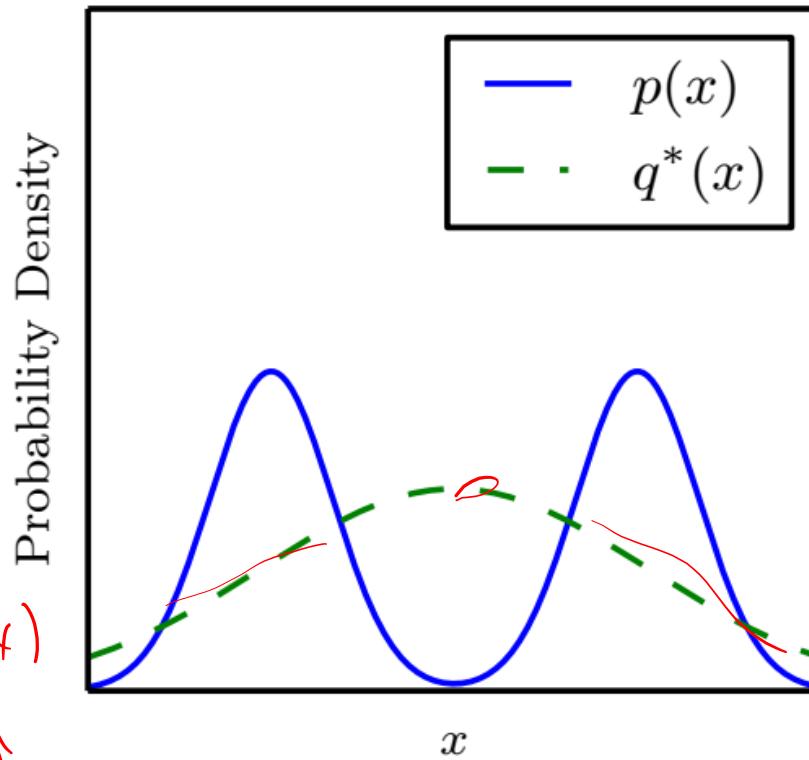
$$D(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}$$



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

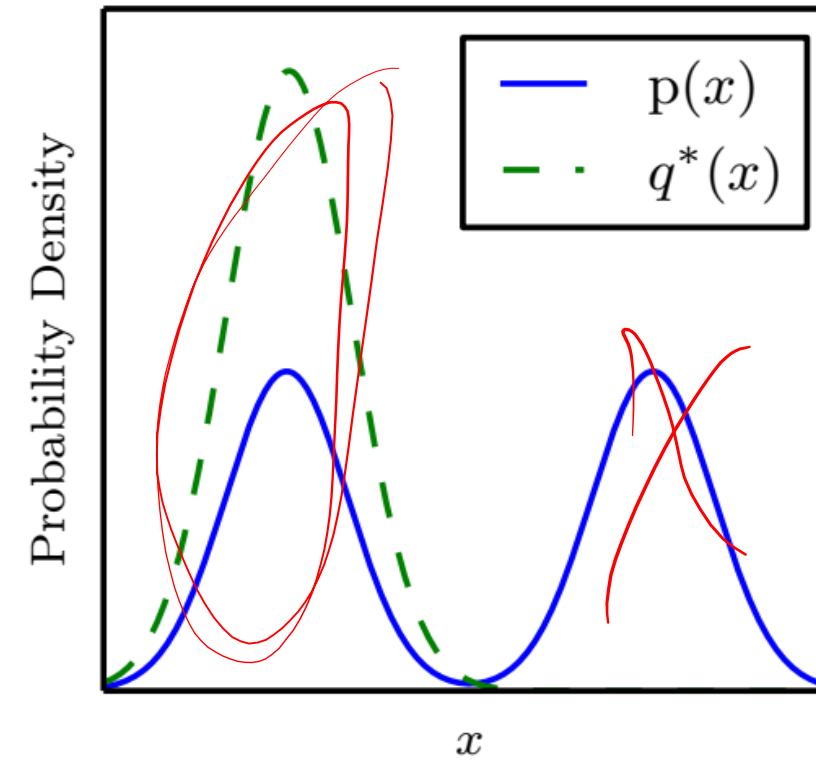
# Is the divergence important?

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



Maximum likelihood  
VAE *new in*  
(Goodfellow et al 2016)

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



Reverse KL  
CAU

# GAN Trick

Non-saturating loss

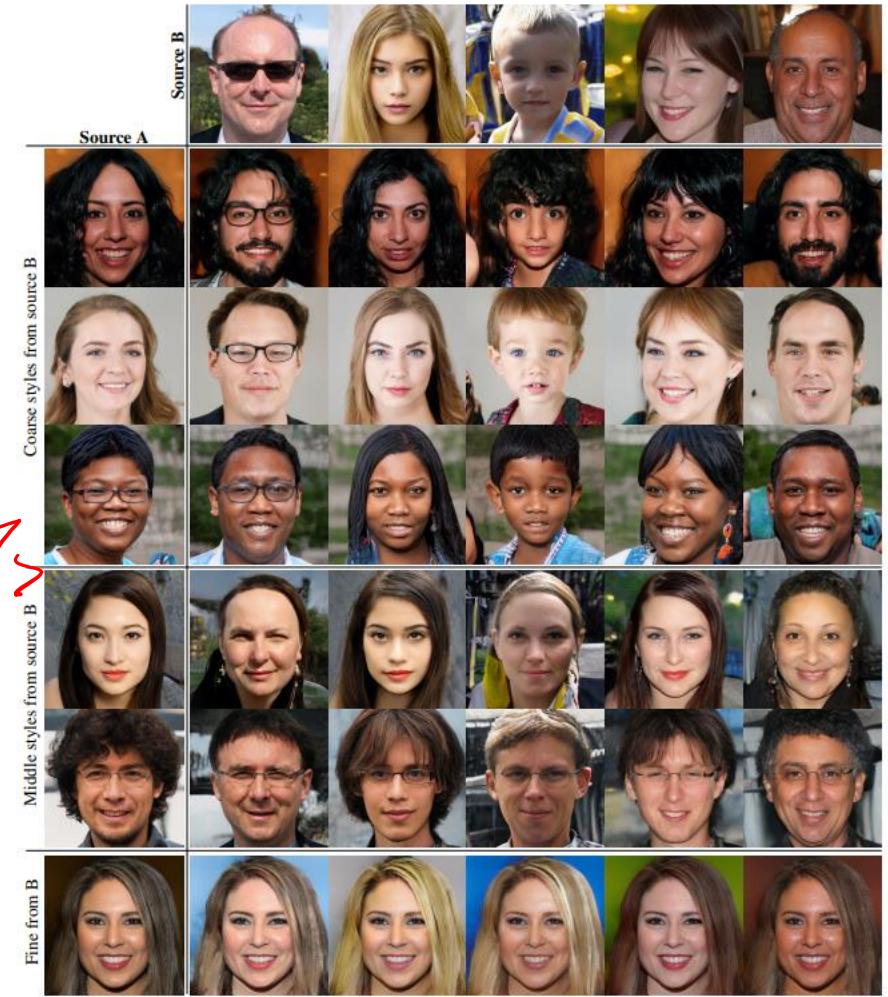
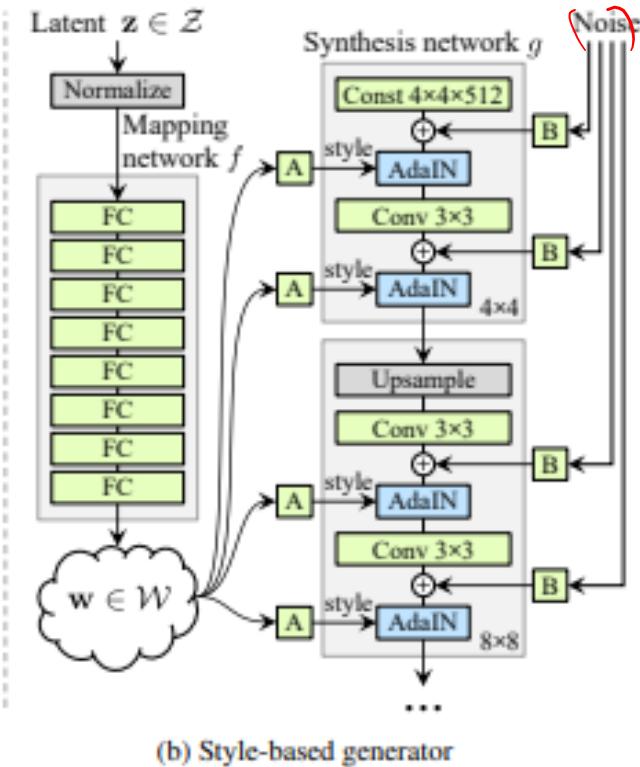
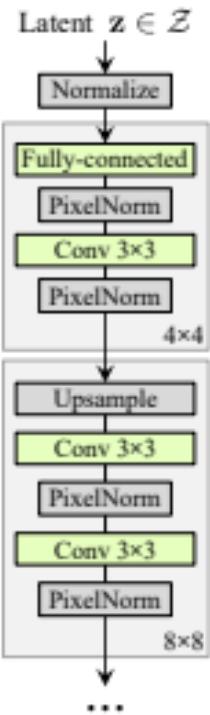
$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

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

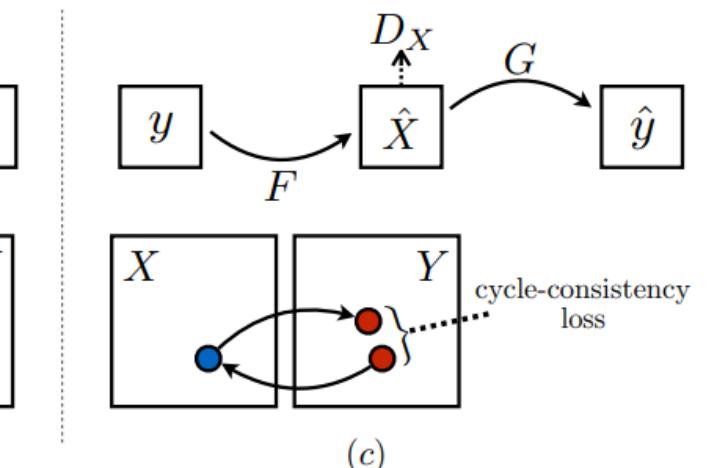
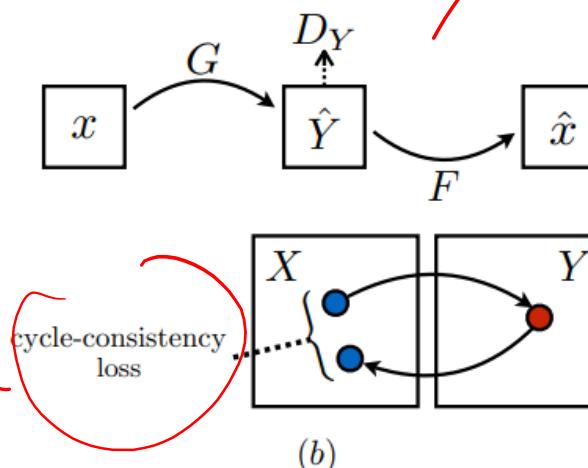
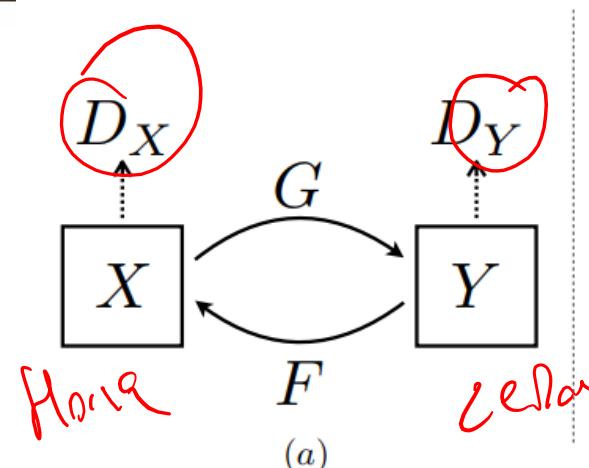
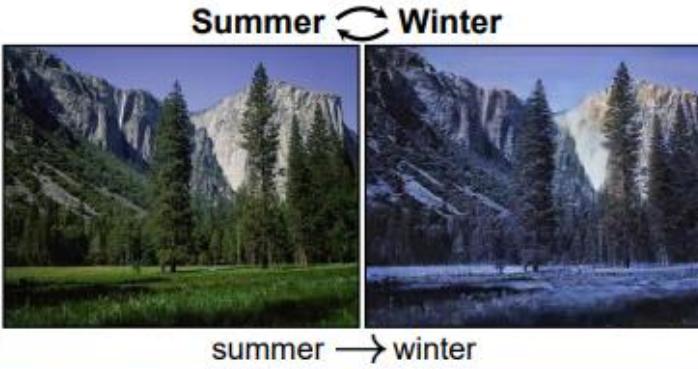
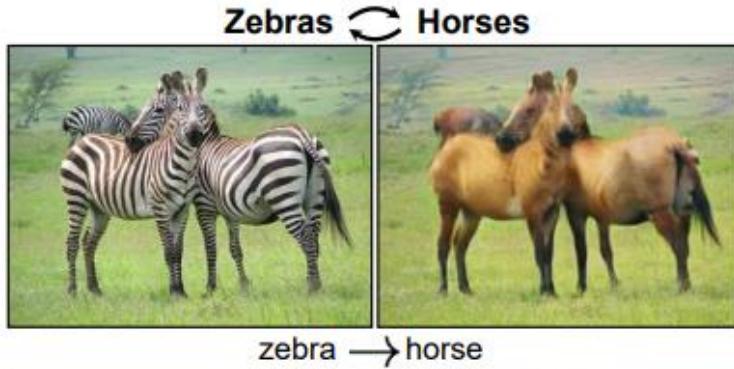
$$J^G \doteq -J^D - \cancel{\frac{1}{2} \mathbb{E}_{\mathbf{z}} \left( \log \cancel{D}(G(\mathbf{z})) \right)}$$

# GAN Uses: image generation

2 min



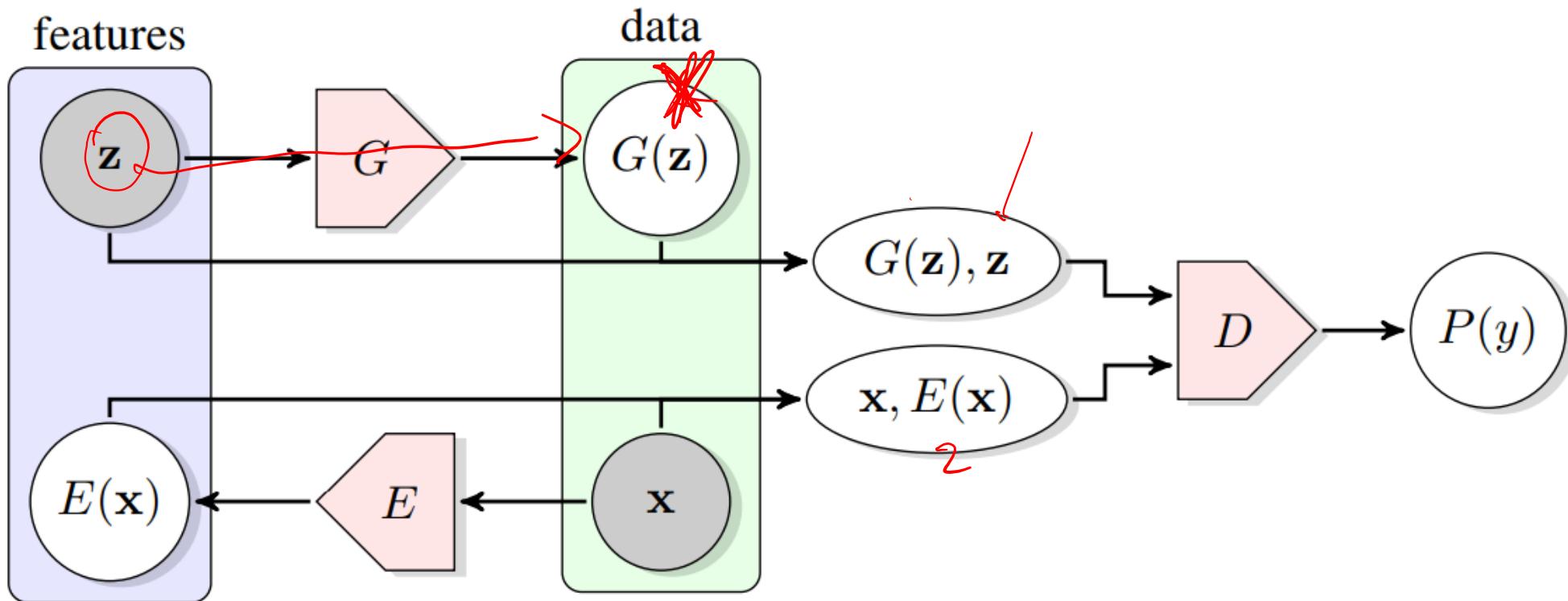
# GAN Uses: unpaired image-image translation



Cycle GAN

# 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>

$\rightarrow c \rightarrow$  important latent dim

$z \rightarrow$  noise

GAN

$$\min_G \max_D V_I(D, G) = V(D, G) - \lambda I(c; G(z, c))$$

MT between  $c$   $G(c, z)$

$c_1$

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

(a) Varying  $c_1$  on InfoGAN (Digit type)

(b) Varying  $c_1$  on regular GAN (No clear meaning)

1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1
8 8 8 8 8 8 8 8	8 8 8 8 8 8 8 8
3 3 3 3 3 3 3 3	3 3 3 3 3 3 3 3
9 9 9 9 9 9 9 9	9 9 9 9 9 9 9 9
5 5 5 5 5 5 5 5	5 5 5 5 5 5 5 5

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

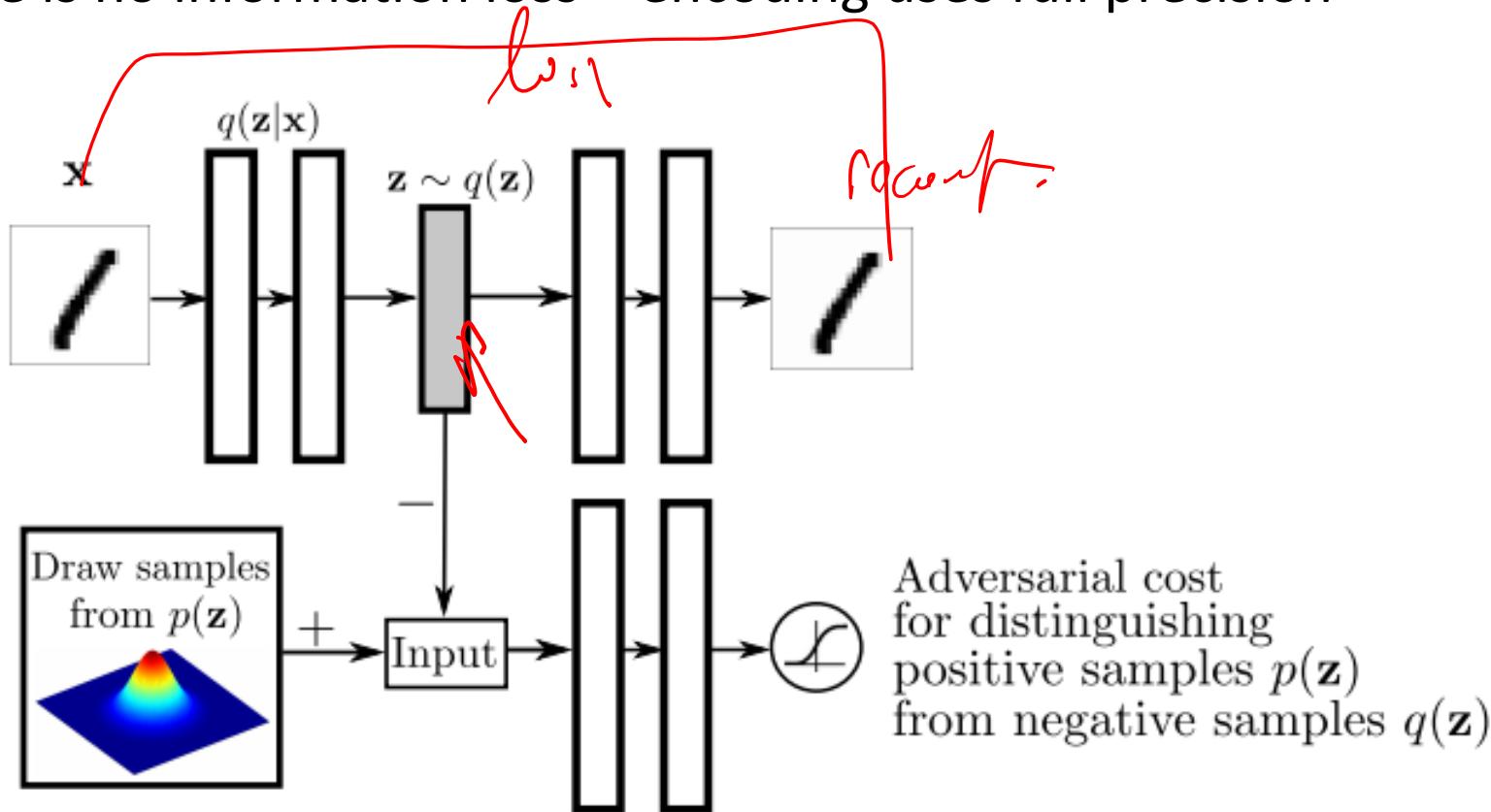
(d) Varying  $c_3$  from -2 to 2 on InfoGAN (Width)

# 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

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



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

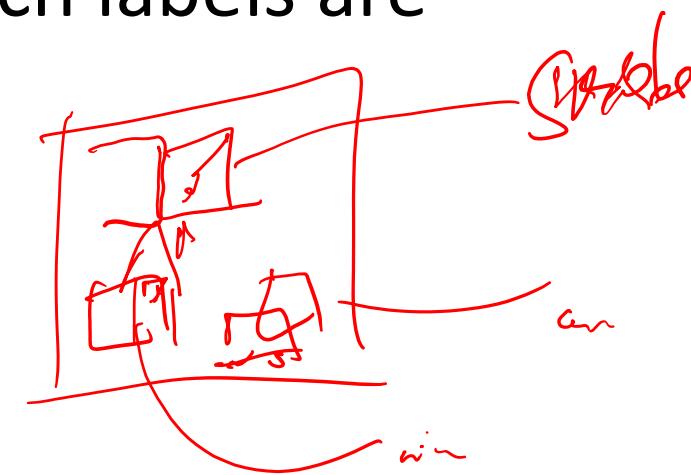
# Core Idea

Solve a task requiring data understanding, for which labels are easy to get.

E.g. cut image into parts, predict their location.

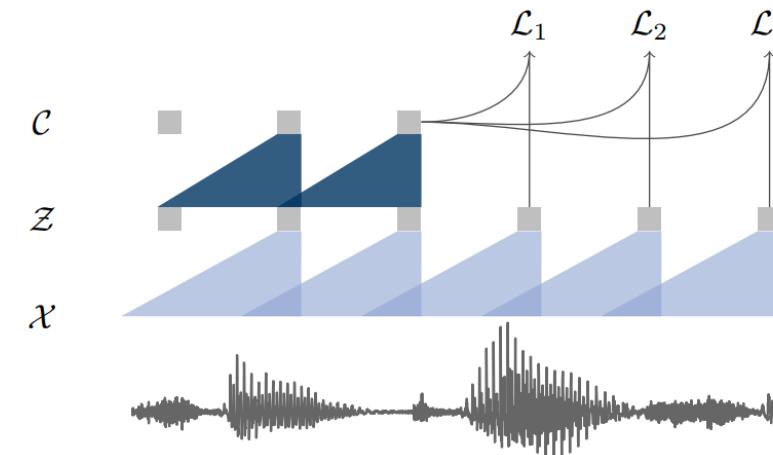
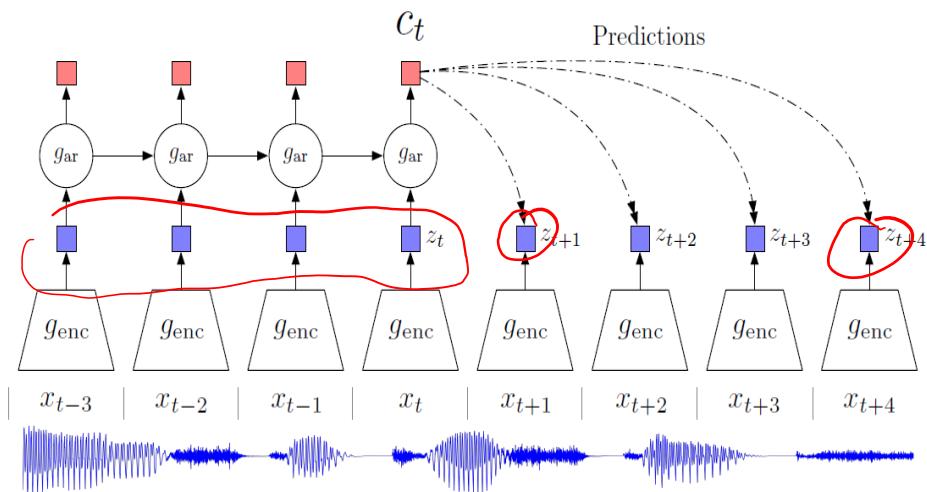
Or cut video into parts, predict time direction

Or cut sentence into words, predict neighbors (quick, what was the name of this model??)



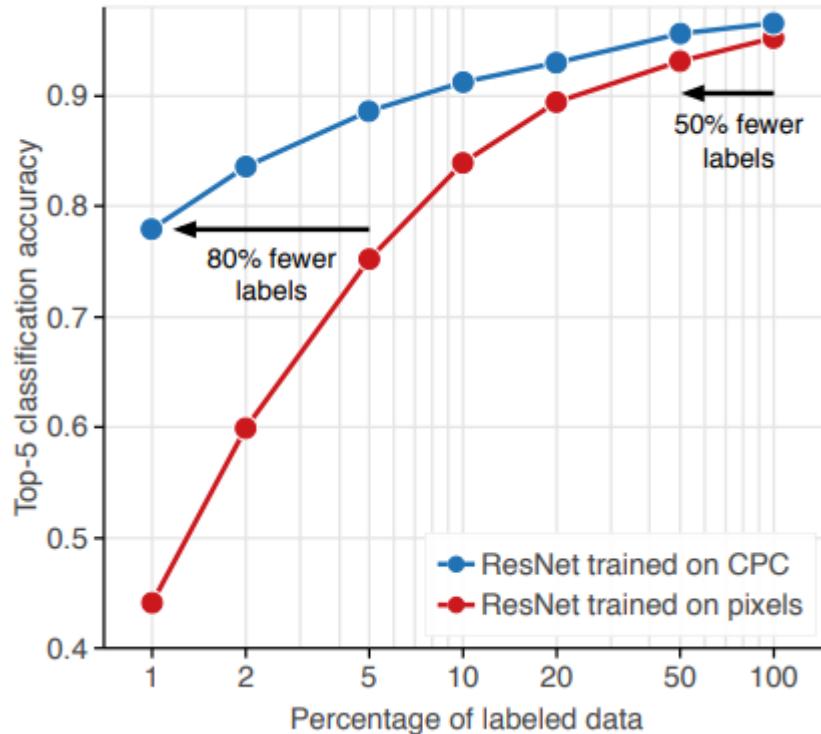
# Decoder-less approach: CPC

Contrastive coding learns representations that can tell a frame from a randomly sampled one

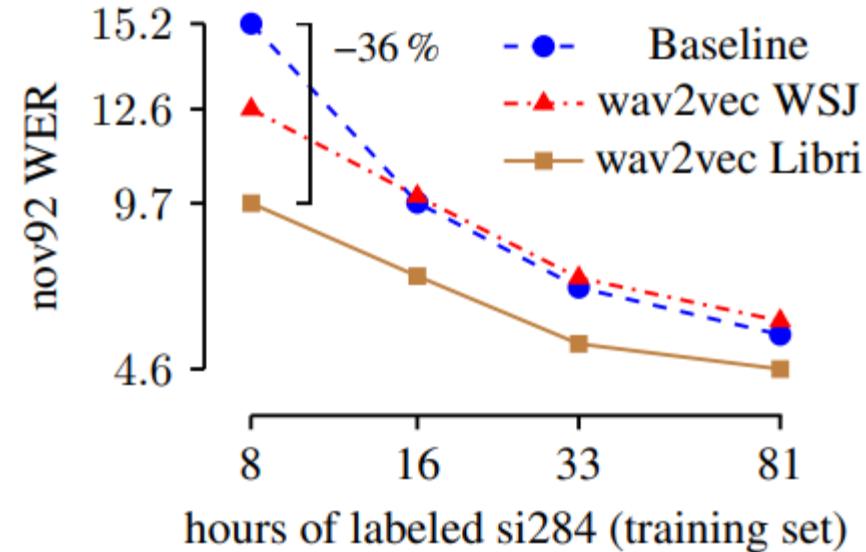


Oord et al. „Representation Learning with Contrastive Predictive Coding“  
Schneider et al. „wav2vec: Unsupervised Pre-training for Speech Recognition“

# CPC 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

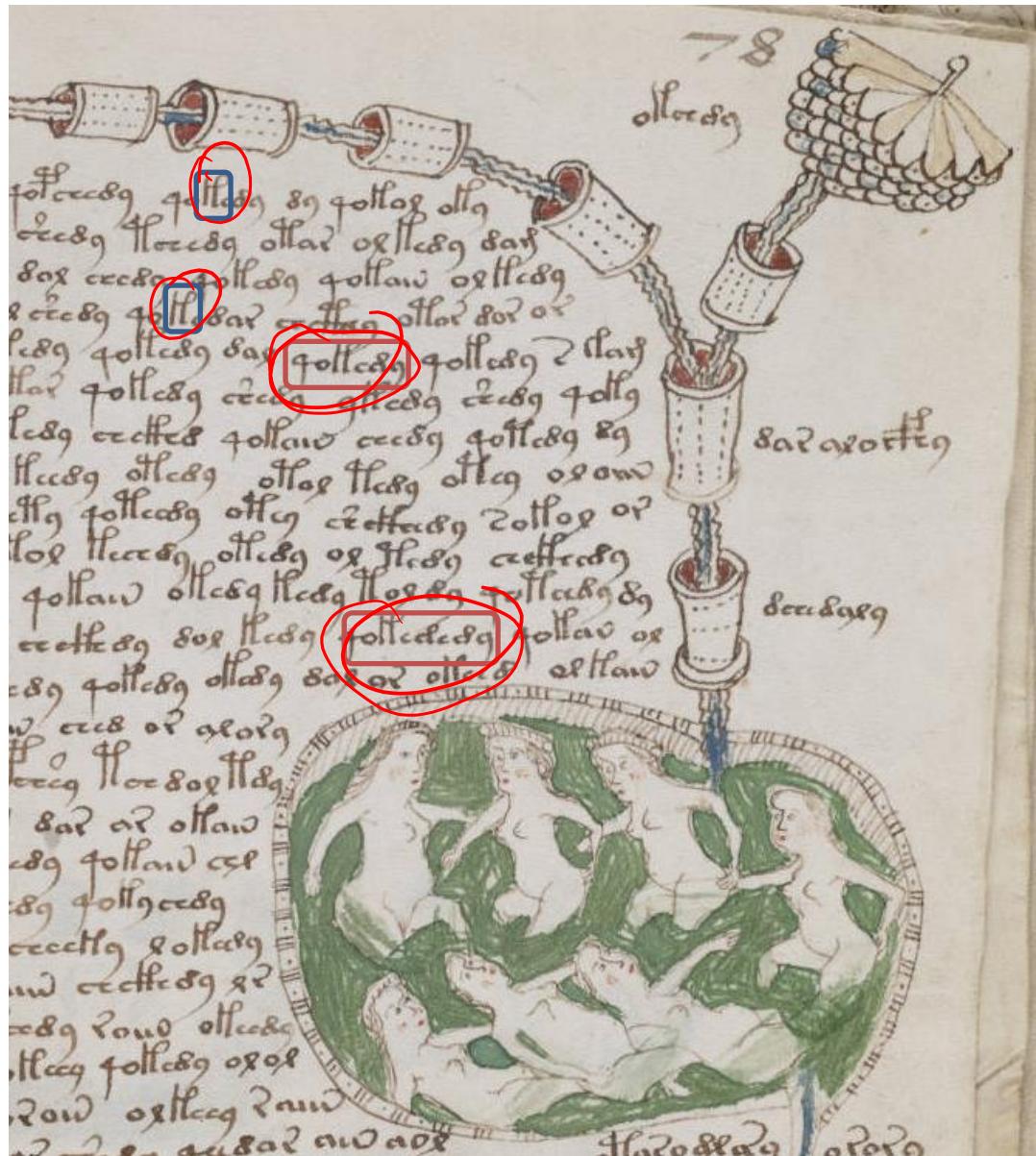
# Summary for unsupervised learning

- 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

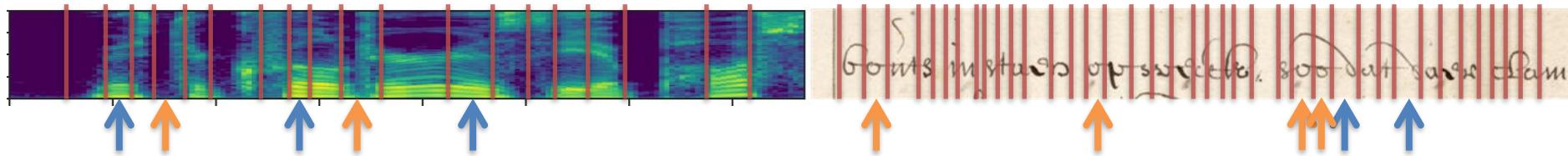


What can we learn  
without labels:

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

# Our goal

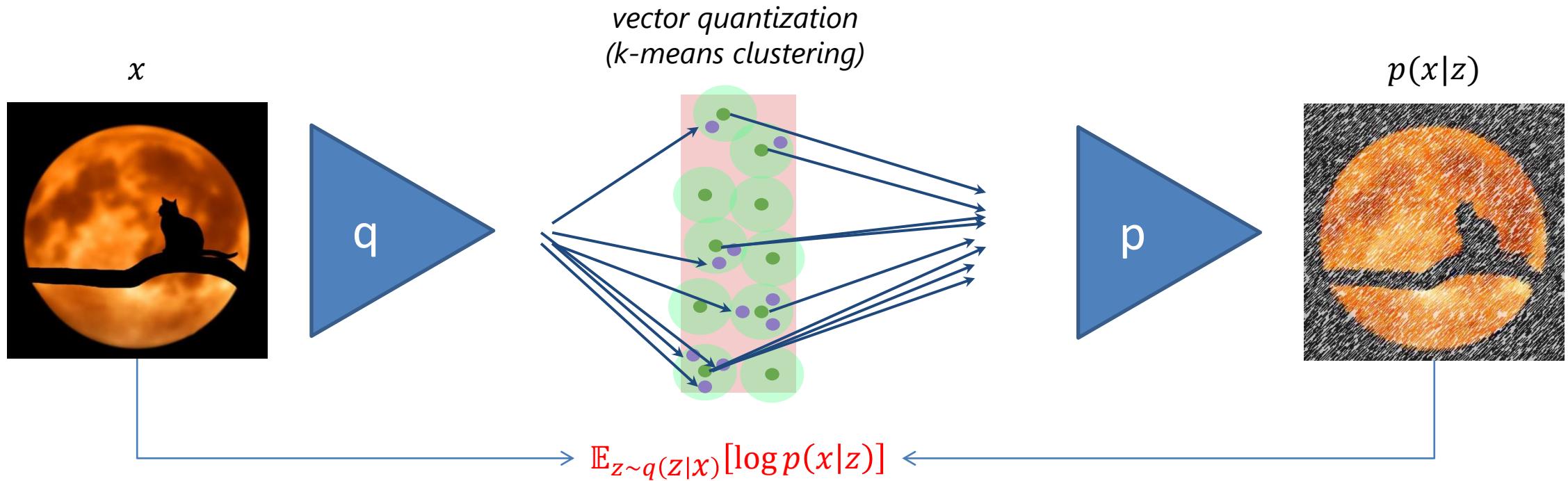
Unsupervised unit discovery  
in speech and handwriting



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

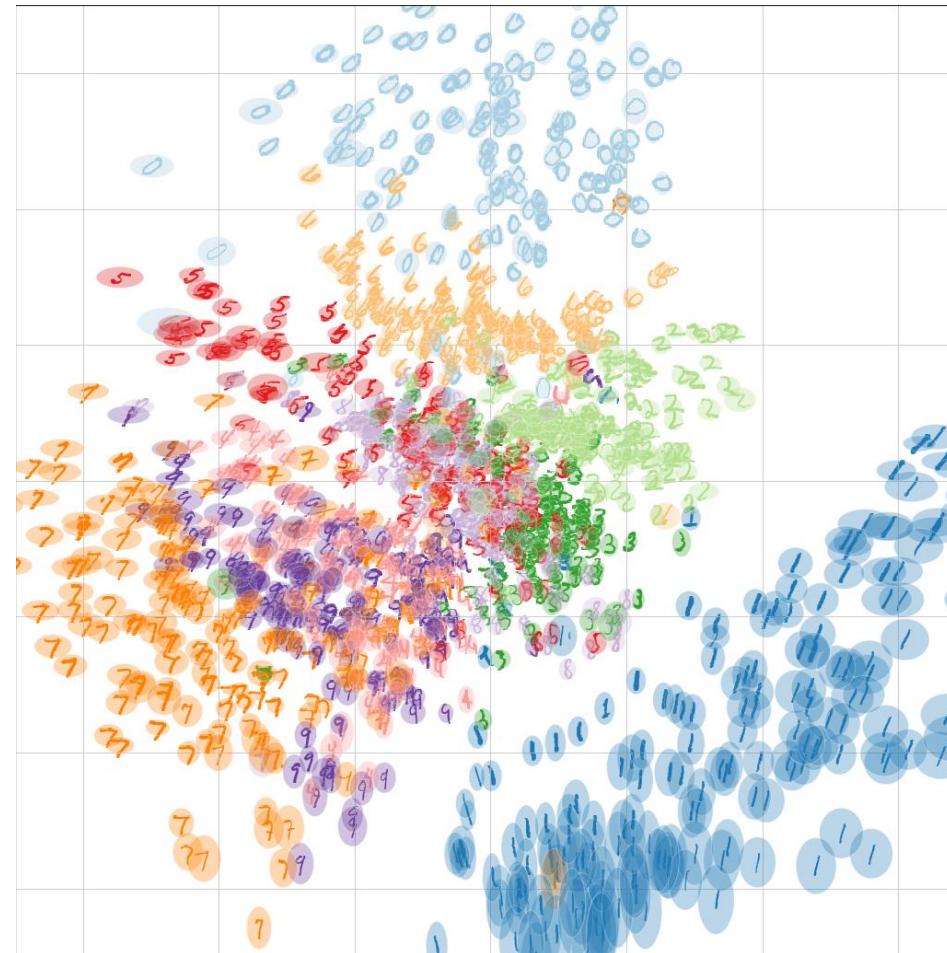
it is replaced by one of K prototypes

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

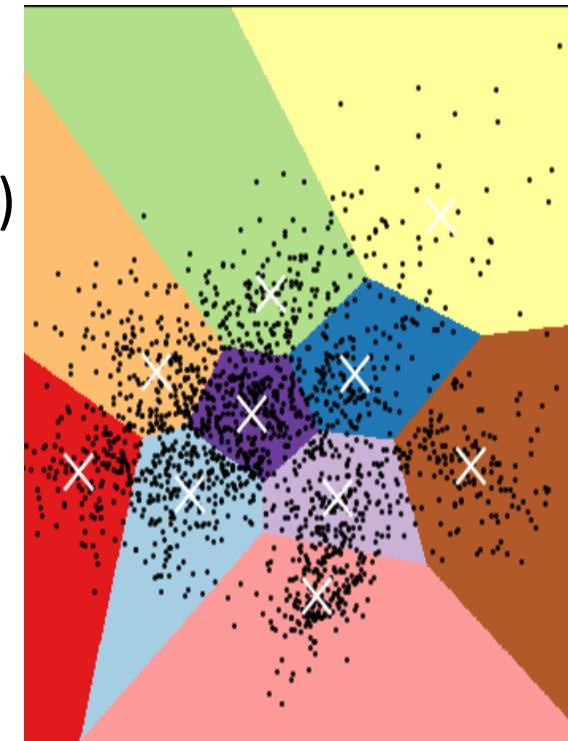
Limit precision of the encoding by quantizing (round each vector to a nearest prototype).

Output can be treated:

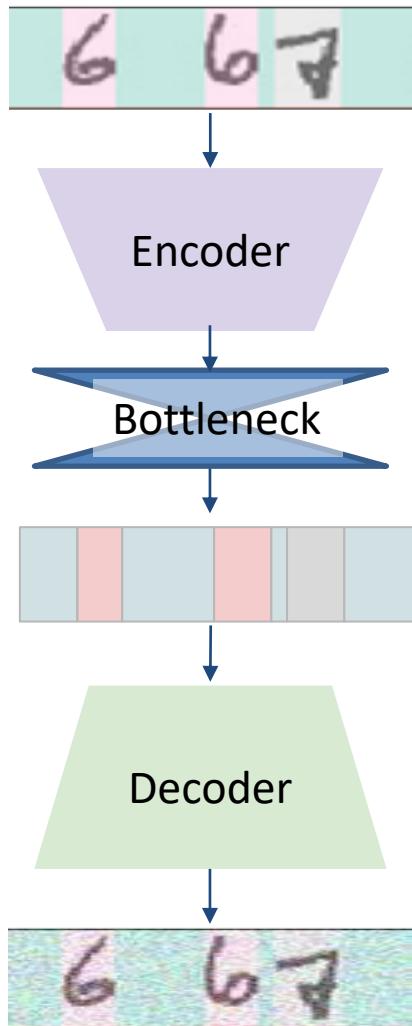
- As a sequence of discrete prototype ids (tokens)
- As a distributed representation (the prototypes themselves)

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

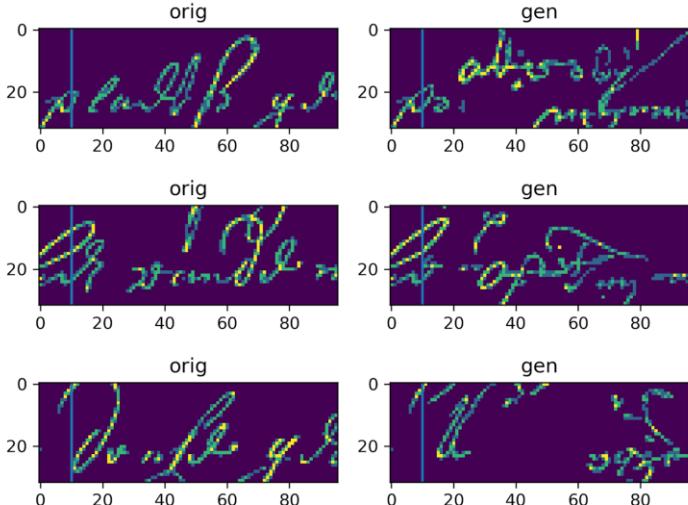
Latent encoding:  
only high level features

Decoder regenerates the inputs.

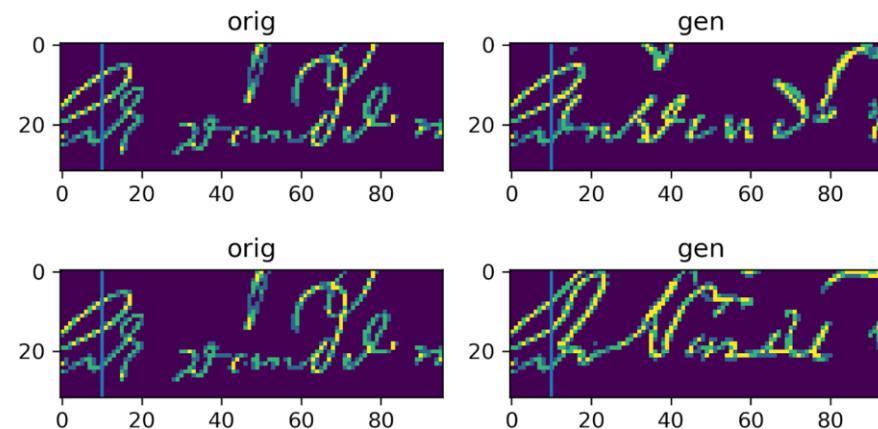
**It learns  $p(x|z)$  and it must fill in all discarded information.**

# Conditional PixelCNN can regenerate images.

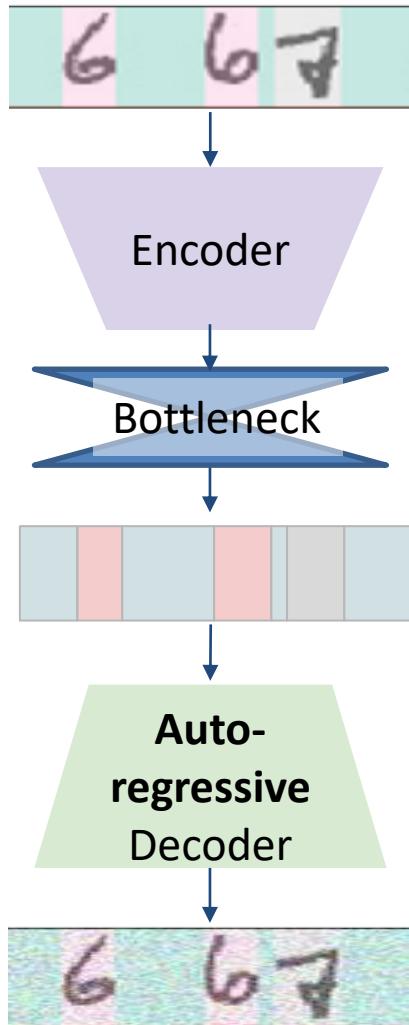
Unconditional



Conditioned on text



# Our model: Autoencoder + autoregressive decoder

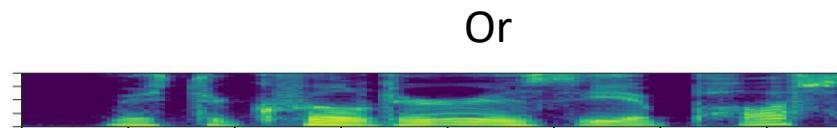
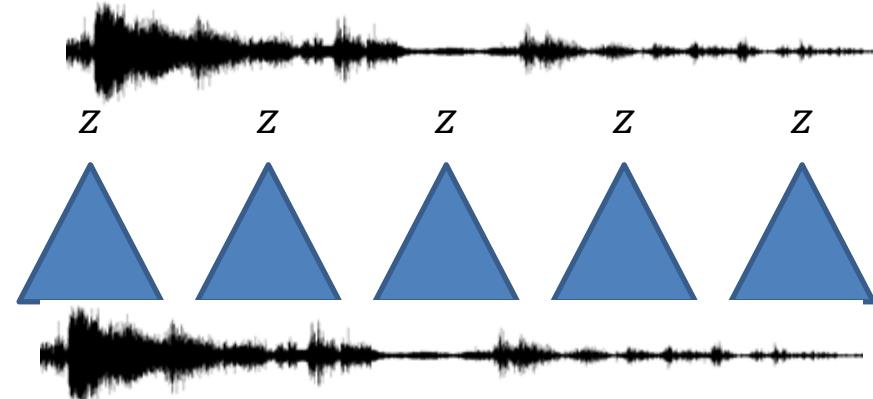
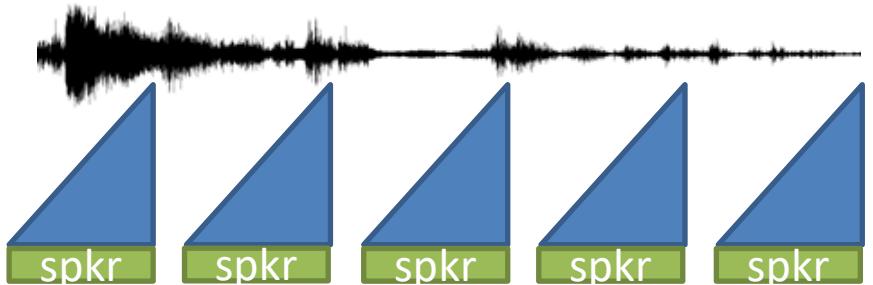


Decoder computes a probability over input images/waveforms conditioned on the encoding.

It combines information from:

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

# Experiments on Speech



Or

WaveNet decoder conditioned on:

- latents extracted at 24Hz-50Hz
- speaker

3 bottleneck evaluated:

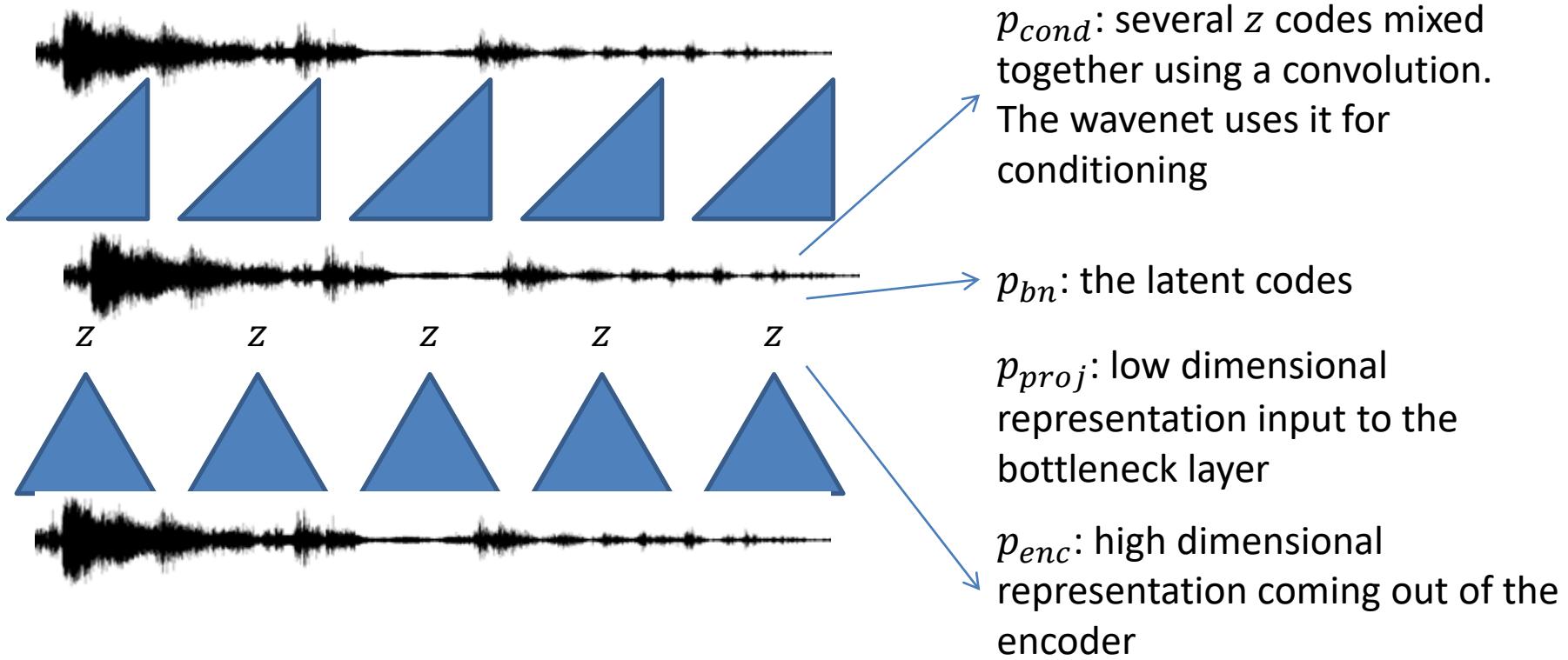
- Dimensionality reduction, max 32 bits/dim
- VAE,  $KL(q(z|x) \parallel \mathcal{N}(0,1))$  nats (bits)
- VQVAE with  $K$  protos:  $\log_2 K$  bits

Input:

Waveforms, Mel Filterbanks, MFCCs

# Representation probing points

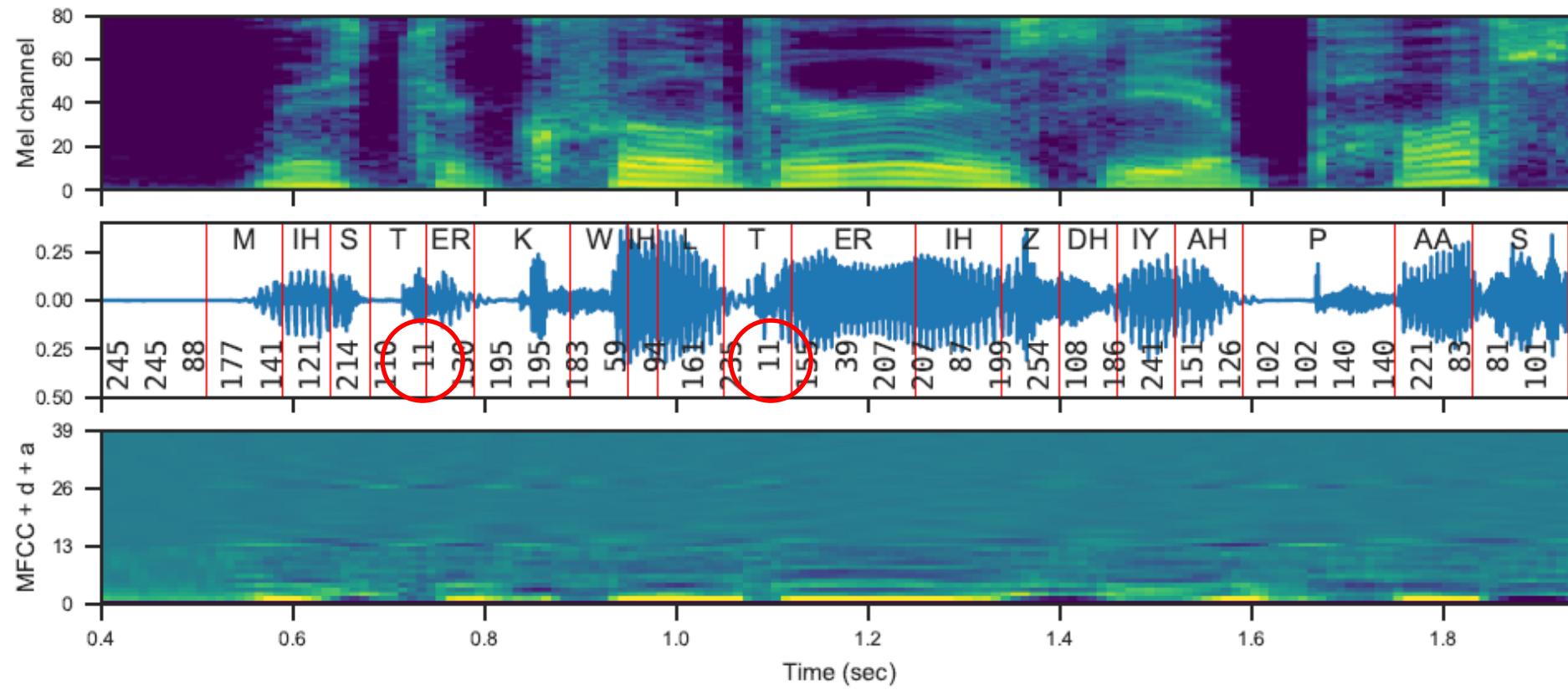
We have inserted probing classifiers at 4 points in the network:



# Experimental Questions #1

- What information is captured in the latent codes/probing points?
- What is the role of the bottleneck layer?
- How good is the representation on downstream tasks?
- What design choices affect it?

# VQVAE Latent representation

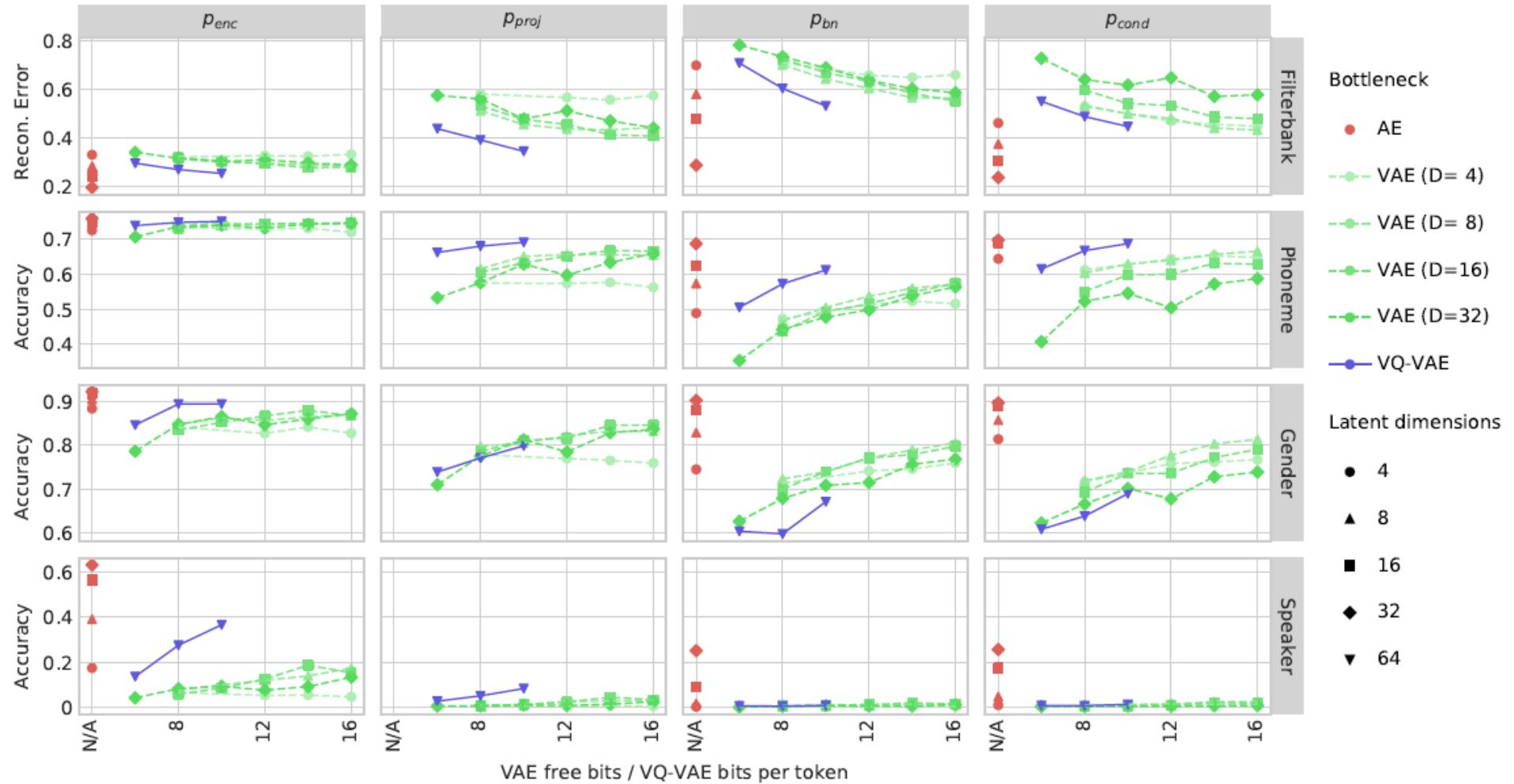


# What information is captured in the latent codes?

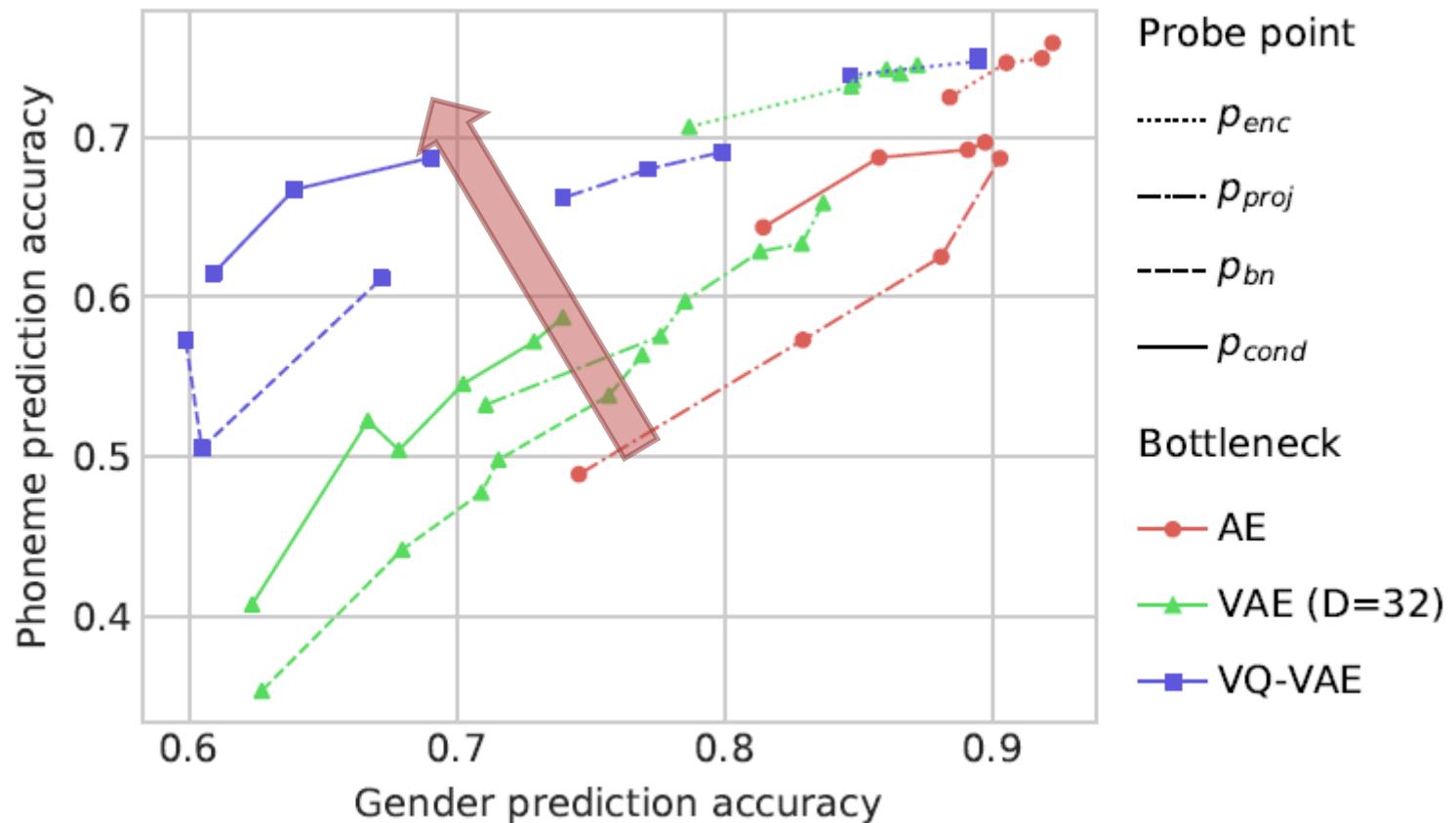
For each probing point, we have trained predictors for:

- Framewise phoneme prediction
- Speaker prediction
- Gender prediction
- Mel Filterbank reconstruction

# Results



# 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.

# English: VQVAE bottleneck adds speaker invariance

English	Within spkr.	Across spkr.
VQ-VAE (per lang, MFCC, $p_{\text{cond}}$ )	<b>5.6</b>	<b>8.0</b>
VQ-VAE (per lang, MFCC, $p_{\text{bn}}$ )	6.2	8.8
VQ-VAE (per lang, MFCC, $p_{\text{proj}}$ )	5.9	9.0
VQ-VAE (all lang, MFCC, $p_{\text{cond}}$ )	5.8	8.6
VQ-VAE (all lang, MFCC, $p_{\text{bn}}$ )	6.3	9.2
VQ-VAE (all lang, MFCC, $p_{\text{proj}}$ )	5.8	9.3
VQ-VAE (all lang, fbank, $p_{\text{proj}}$ )	6.0	10.1

The quantization discards speaker info, improving across-speaker results  
MFCCs slightly better than FBanks

# Mandarin: VQVAE bottleneck discards phone information

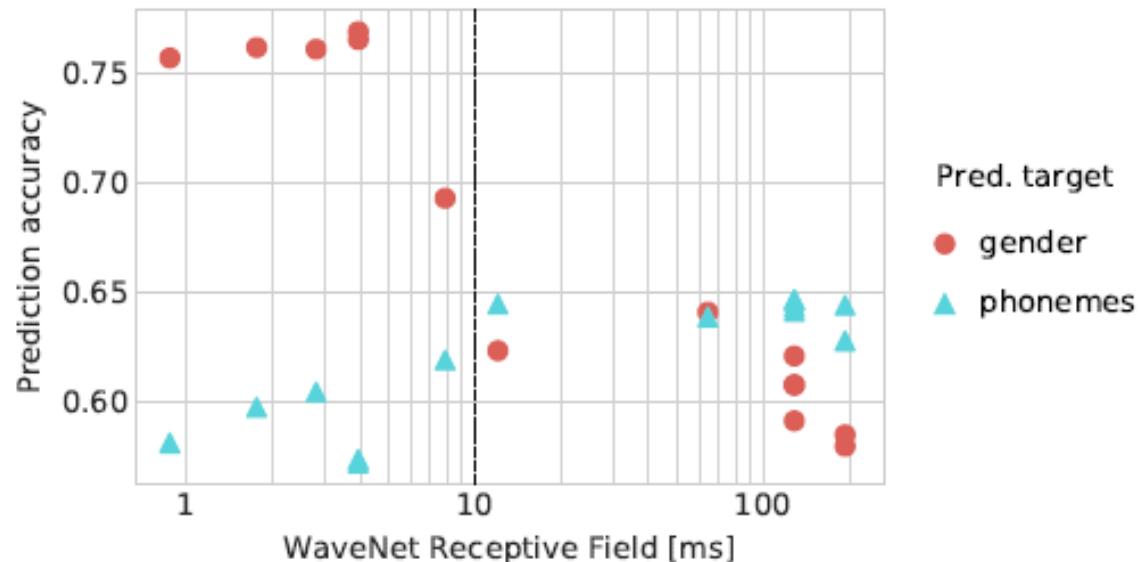
Mandarin	Within spkr.	Across spkr.
VQ-VAE (per lang, MFCC, $p_{\text{cond}}$ )	11.2	12.2
VQ-VAE (per lang, MFCC, $p_{\text{bn}}$ )	10.8	11.9
VQ-VAE (per lang, MFCC, $p_{\text{proj}}$ )	9.9	11.0
VQ-VAE (all lang, MFCC, $p_{\text{cond}}$ )	9.2	10.3
VQ-VAE (all lang, MFCC, $p_{\text{bn}}$ )	9.0	9.9
VQ-VAE (all lang, MFCC, $p_{\text{proj}}$ )	7.4	8.6
VQ-VAE (all lang, fbank, $p_{\text{proj}}$ )	6.8	7.8

The quantization discards too much (tone insensitivity?)  
MFCCs worse than FBanks

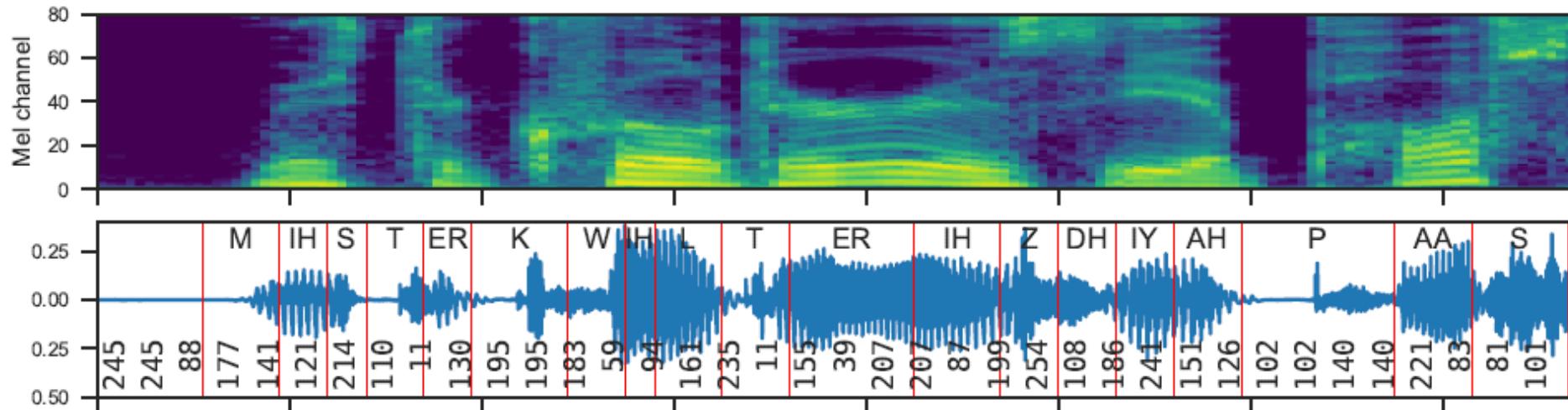
# What impacts the representation?

Implicit time constant of the model:

- Input field of view of the encoder – optimum close to 0.3s
- WaveNet field of view - needs at minimum 10ms



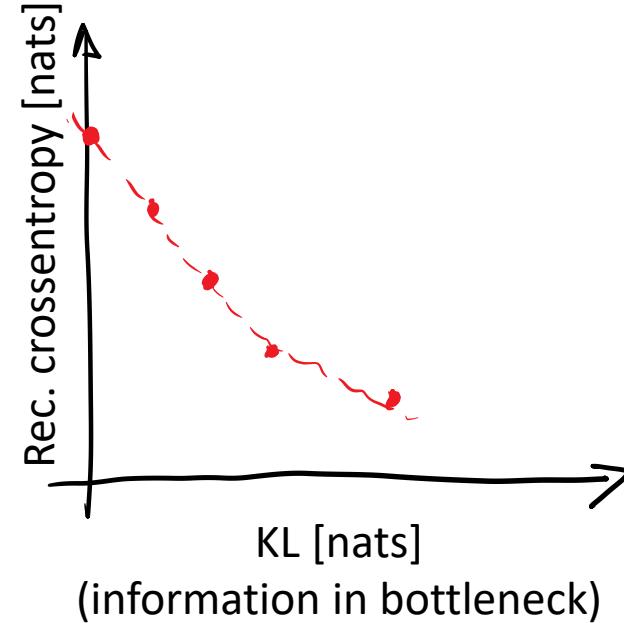
# Experimental Questions #2



- Can we regularize the latent representation?
- How to promote latent code stability/segmentation?

# VAE + autoregressive models: latent collapse danger

- Purely Autoregressive models: SOTA log-likelihoods
- Conditioning on latents:  
information passed through bottleneck  
lower reconstruction x-entropy
- In standard VAE model actively tries to
  - reduce information in the latents
  - maximally use autoregressive information=> Collapse: latents are not used!
- Solution: stop optimizing KL term  
(free bits), make it a hyperparam (VQVAE)



# Priors on latents promote collapse

Consider slow features:

L1 penalty on the difference of neighboring codes

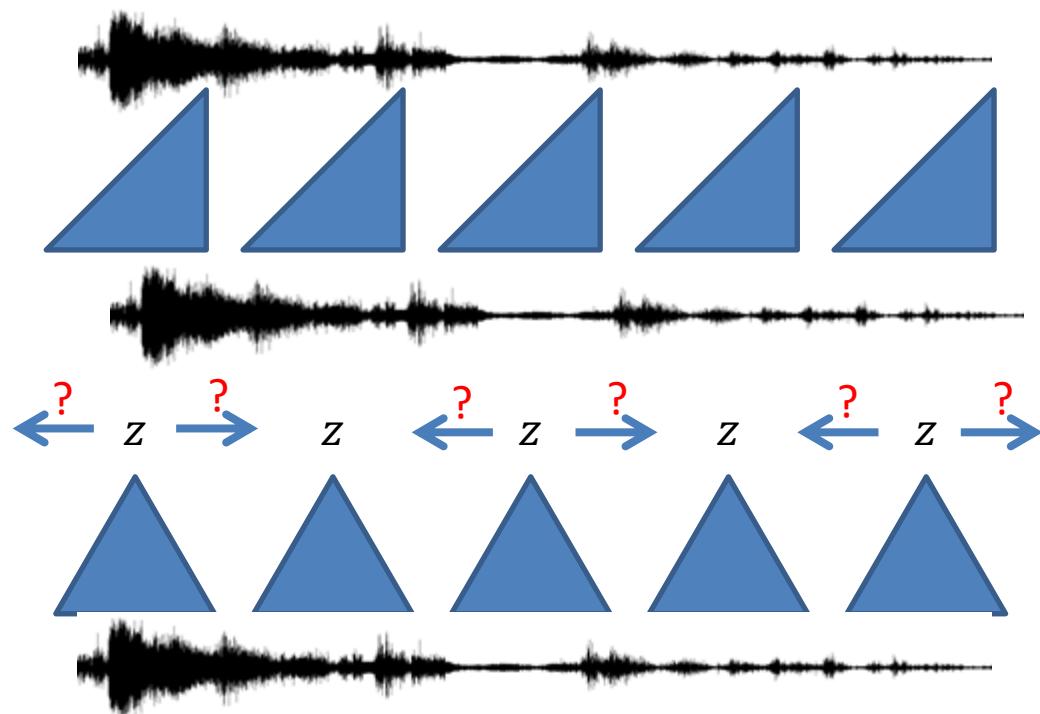
Promotes a constant encoding

Leads to collapse

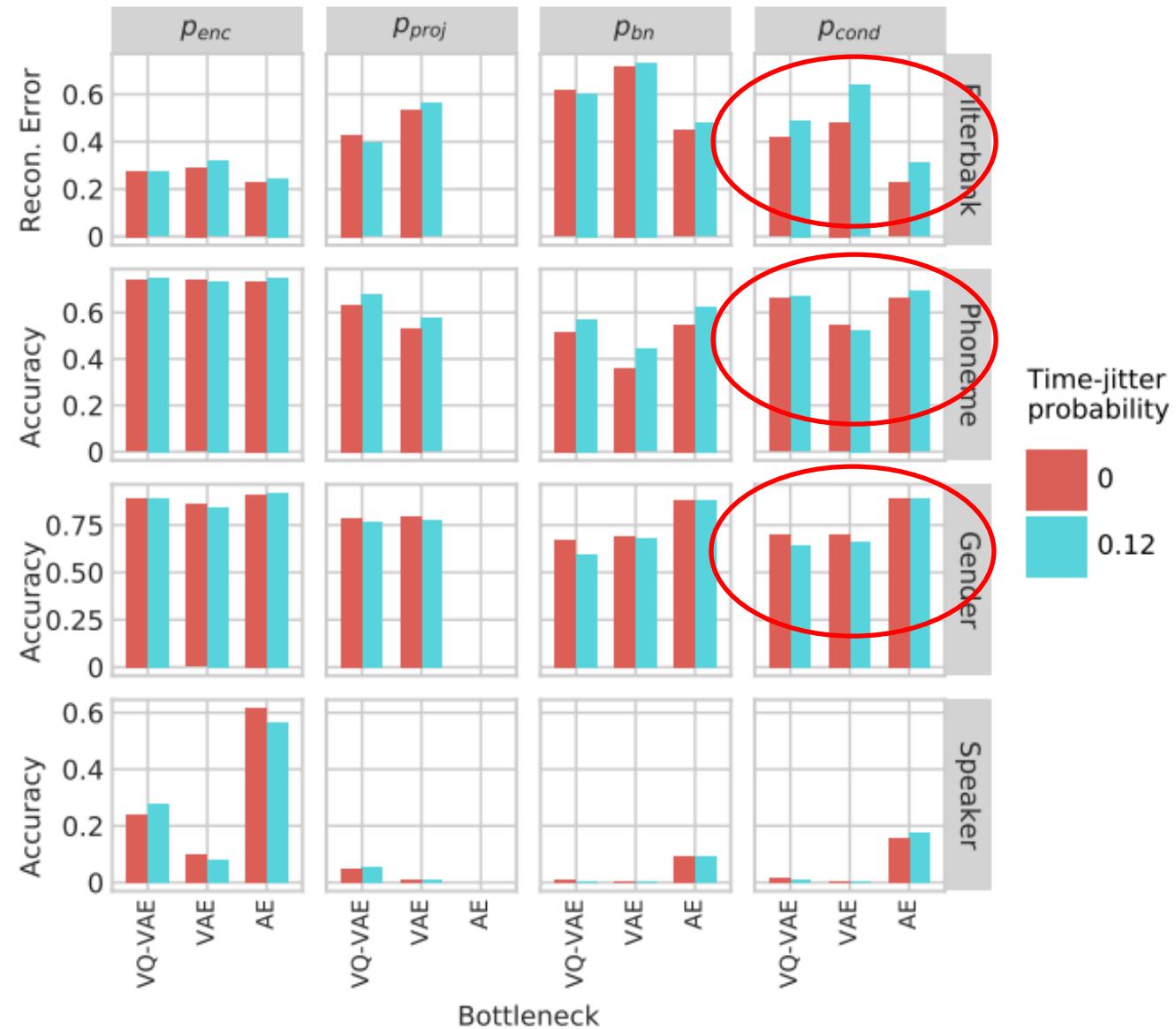


# Randomized time jitter

Rather than putting a penalty on changes of the latent  $z$  vectors, add time jitter to them.  
This forces the model to have a more stable representation over time.



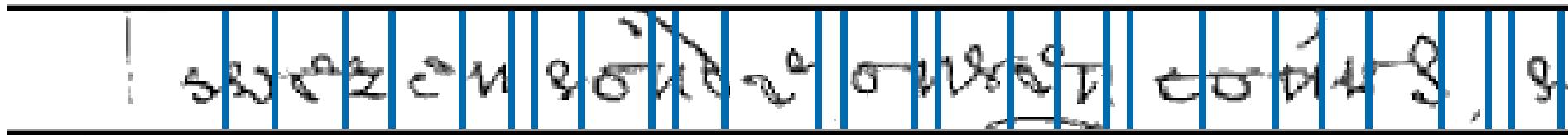
# Randomized time jitter results



# Segmenting auto-encoder

We extract a latent vector every 12 pixels

How to extract at every character?

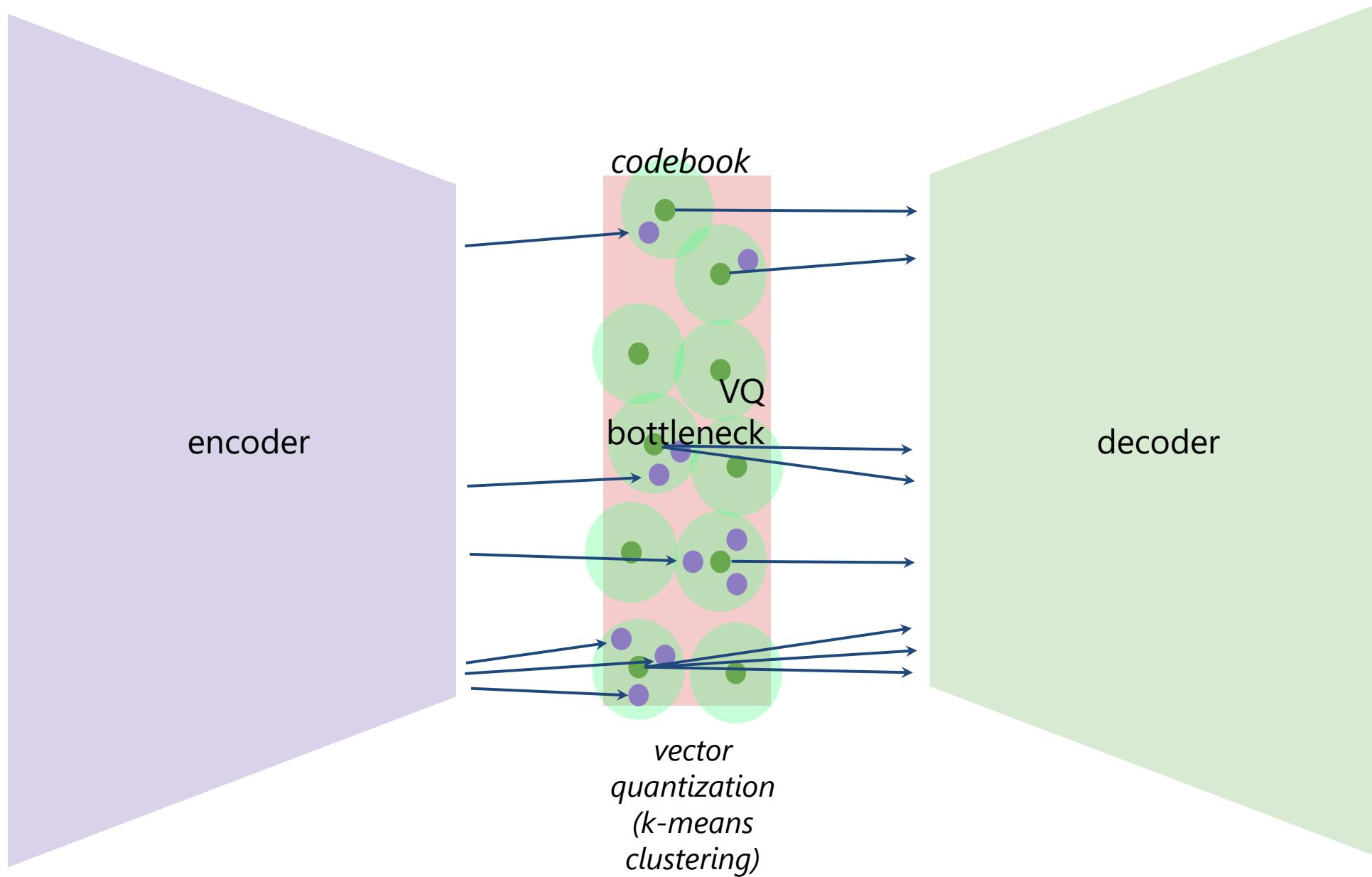


Enforce domain knowledge:

- timing
- average unit/state duration
- state and transition frequency modeling

Deep Learning value: bring hard constraints into smooth, differentiable models.

# VQ-VAE



# Learning piecewise-constant representations

Quantization criterion:

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

*encoder outputs*

*VQ tokens*

*expected  
# of  
characters*

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:

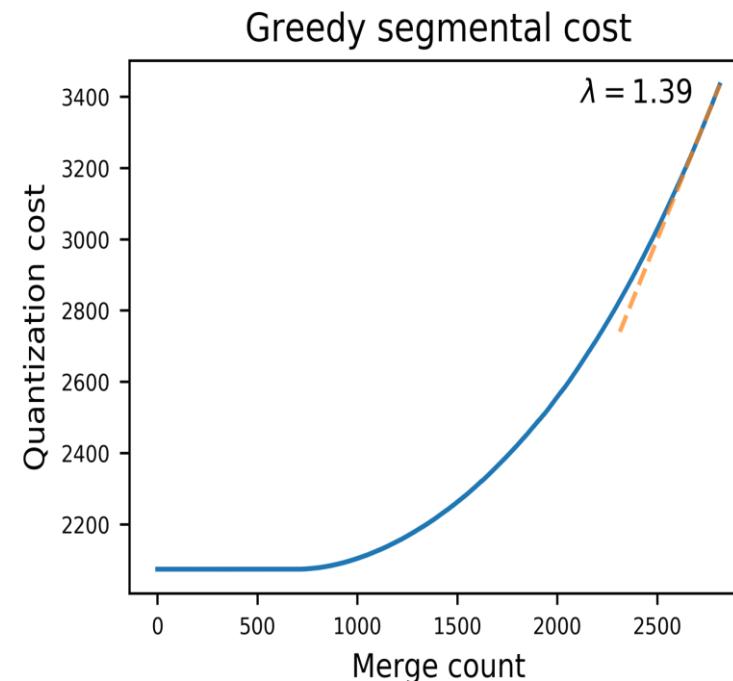
$$\min \sum_{t=1}^T \|x_t - q_t\| + \lambda \left( \sum_{t=1}^T [q_{t-1} \neq q_t] - K \right)$$

During training:

- . estimate  $\lambda$

During inference:

- . use  $\lambda$  to balance the costs



# Dual formulation and Markovian dynamics

Dual formulation resembles HMM cost:

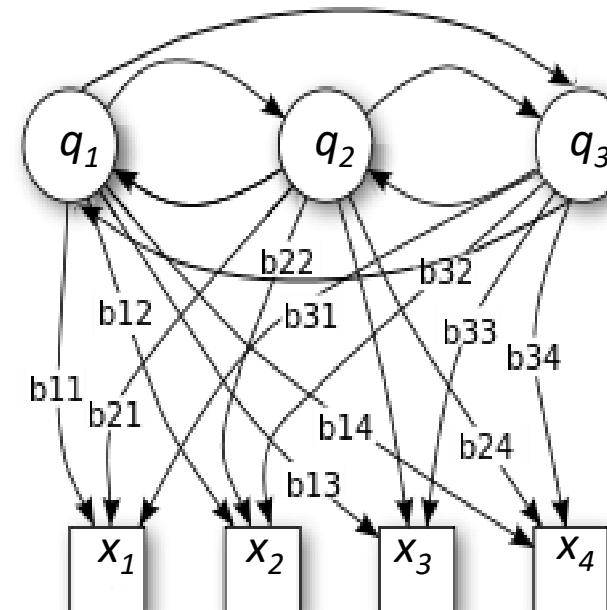
$$\min \sum_{t=1}^T \|x_t - q_t\| + \lambda \left( \sum_{t=1}^T [q_{t-1} \neq q_t] - K \right)$$

*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
	<i>Piecewise-const VQ bottleneck</i>	<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
	<i>Piecewise-const VQ bottleneck</i>	<b>51%</b>	<b>49%</b>	<b>39%</b>	<b>0.18</b>	<b>0.20</b>