

Deep Unsupervised Learning pt 2

Unsupervised learning goals

① Generative
Models

Autoregressive Models

n-gram LMs

pix CNN

GAN

VAE

ReINFORCE

BERT

Word2vec

② Self-supervised
learning

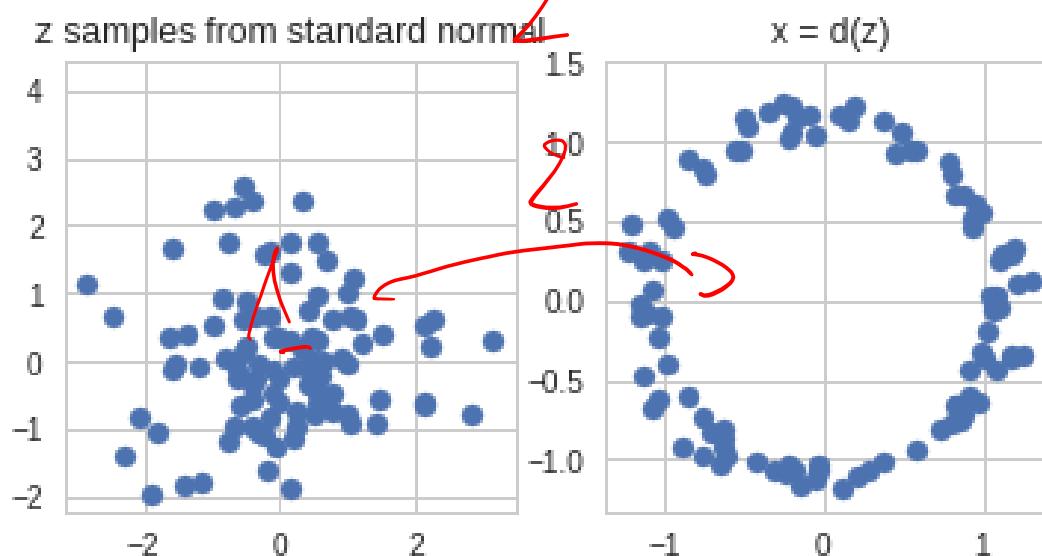
③ Repr learning
models

VAE

Data Generation with Latent Variables: Intuitions

Intuition: to generate something complicated, do:

1. Sample something simple $z \sim \mathcal{N}(0, 1)$
2. Transform it $x = \frac{z}{10} + \frac{z}{\|z\|}$



Latent variable model

Assume a 2 stage data 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

How to train this model?

$$\log \underbrace{p(x)}_{\text{intractable}} = \log \sum_z p(x|z)p(z)$$

$p(x|z)$

This is often intractable!

EM to the rescue!

Want:

$$\max \underbrace{\log p_{\Theta}(x)}_{\text{E}} = \max \log \sum_z p_{\Theta}(x, z) = \max \log \sum_z \underbrace{p_{\Theta}(x|z)p_{\Theta}(z)}_{\text{M}}$$

EM to the rescue!

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

EM to the rescue!

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 $q_{\Phi}(z|x) \neq p_{\Theta}(z|x)$

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 Φ fixed)

$$\max_{\Theta} \sum_z q_{\Phi}(z|x) \log p_{\Theta}(x|z) p_{\Theta}(z)$$

no sum under \log !

EASY !

EM to the rescue!

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 Φ fixed)

Why?

$$\log p(x) = \underbrace{KL(q_{\Phi}(z|x) \parallel p_{\Theta}(z|x))}_{\text{For any } q_{\Phi}(z|x)} + \underbrace{\mathbb{E}_{z \sim q_{\Phi}(z|x)} \left[\log \left(\frac{p_{\Theta}(z,x)}{q_{\Phi}(z|x)} \right) \right]}_{\text{M step maxes}}$$

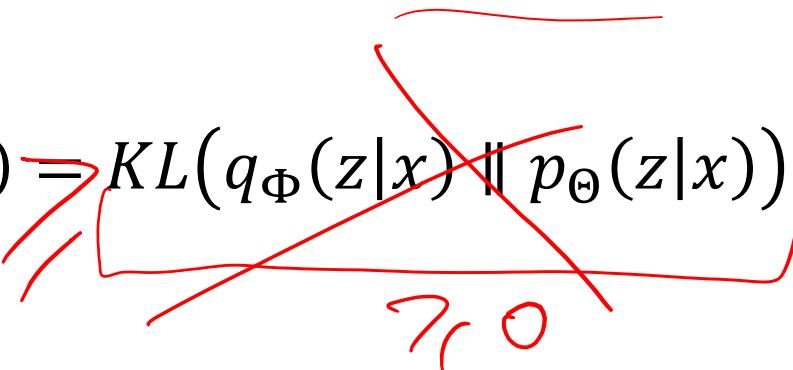
E sets $H_{kl} = 0$

Proof

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

ELBO: A lower bound on $\log p(x)$

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

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


ELBO: A lower bound on $\log p(x)$

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

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

The bound is tight for $p(z|x) = q(z|x)$.

Idea (VAE): Optimize ELBO using gradient techniques jointly over Θ and Φ

ELBO interpretation

ELBO, or evidence lower bound:

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

where:

$\mathbb{E}_{z \sim q(z|x)} [\log p(x|z)]$ reconstruction quality:

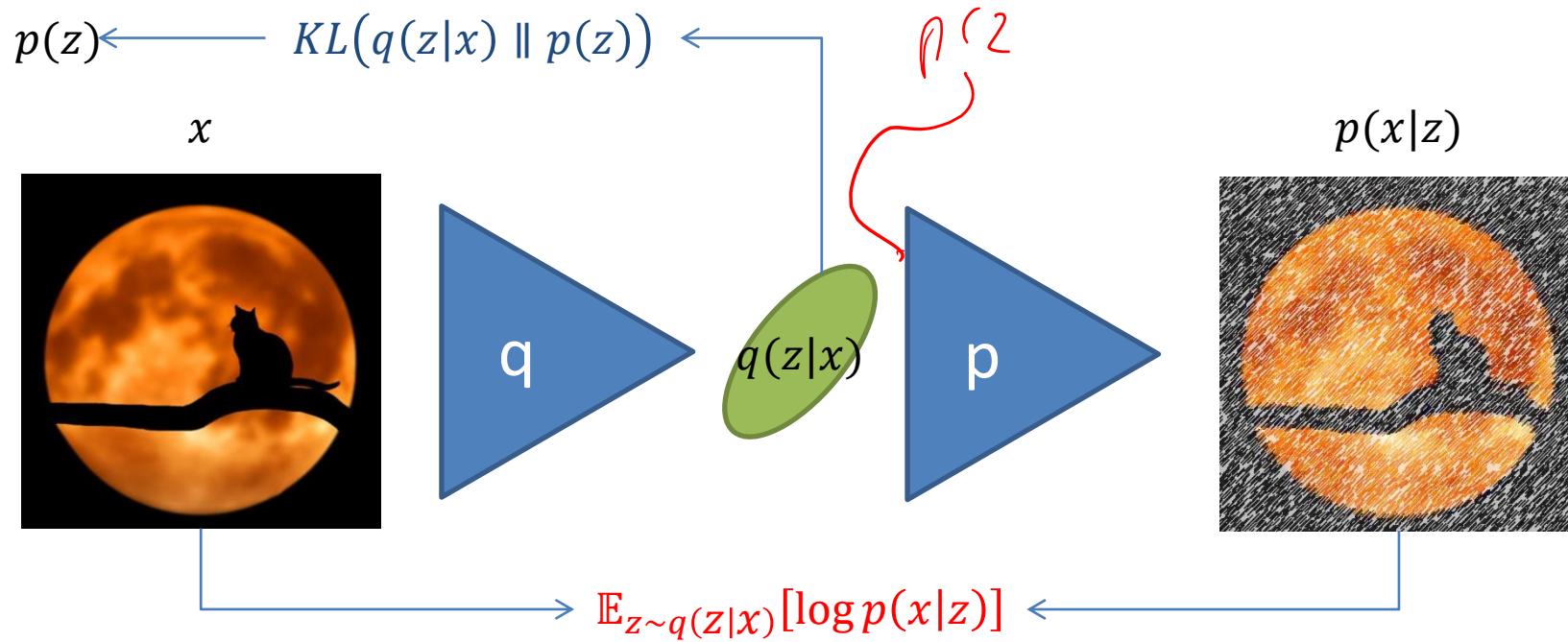
how many nats we need to reconstruct x ,
when someone gives us $q(z|x)$

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

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

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

The Variational Autoencoder



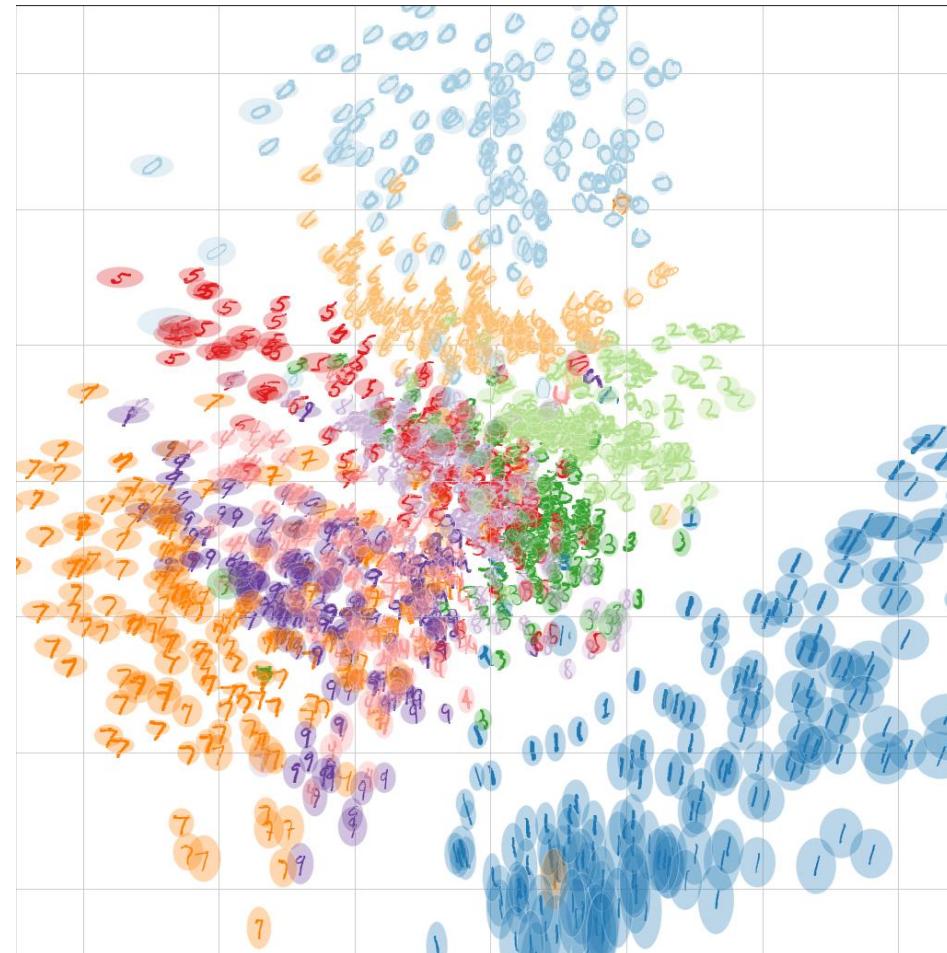
An input x is put through the q network to obtain a distribution over latent code z , $q(z|x)$.

Samples z_1, \dots, z_k are drawn from $q(z|x)$. They k reconstructions $p(x|z_k)$ are computed using the network p .

VAE is an Information Bottleneck

Each sample is represented as a Gaussian

This discards information
(latent representation has low precision)



How to compute the ELBO?

Compute:

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

$\text{KL}(q_{\Phi}(z|x) \| p_{\Theta}(z))$ has closed form for simple $q_{\Phi}(z|x)$:

intractable to compute exactly.

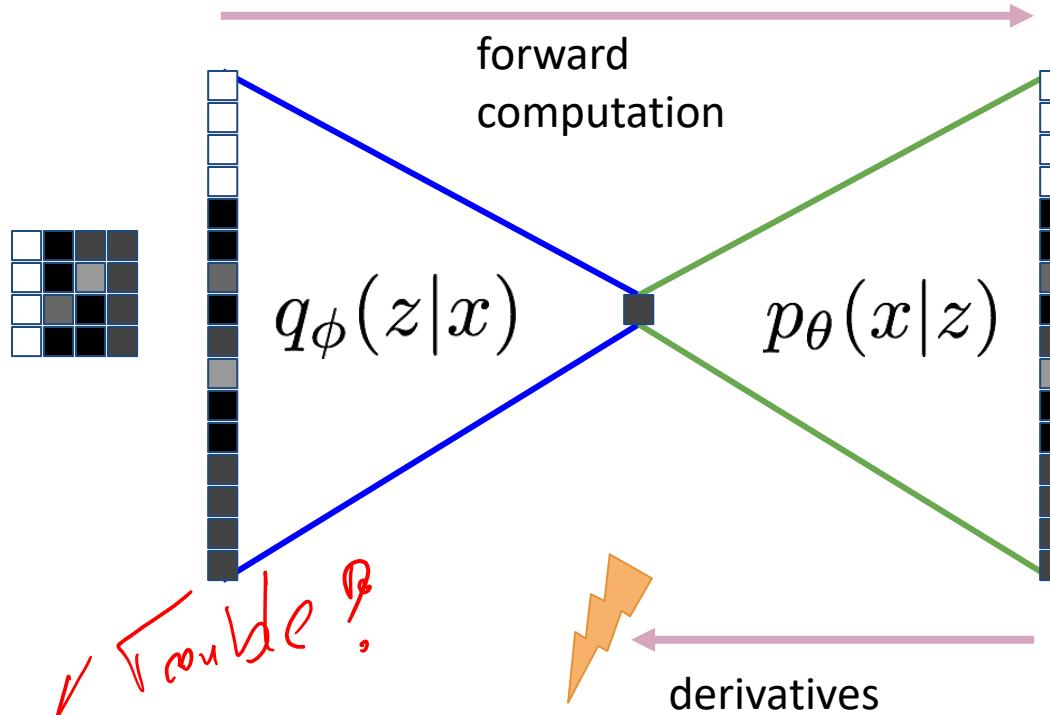
$\mathbb{E}_{z \sim q_{\Phi}(z|x)} [\log p_{\Theta}(x|z)]$ can be approximated using Monte Carlo:

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

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

Monte Carlo
samples

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:

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

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

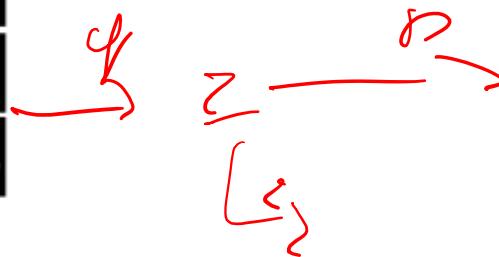
$\overbrace{\quad \quad \quad}^{\text{no dep on } \phi}$ $\overbrace{\quad \quad \quad}^{\epsilon = \mu_z + \sigma_z \cdot \epsilon}$

ϵ is drawn from a fixed distribution.

With ϵ 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

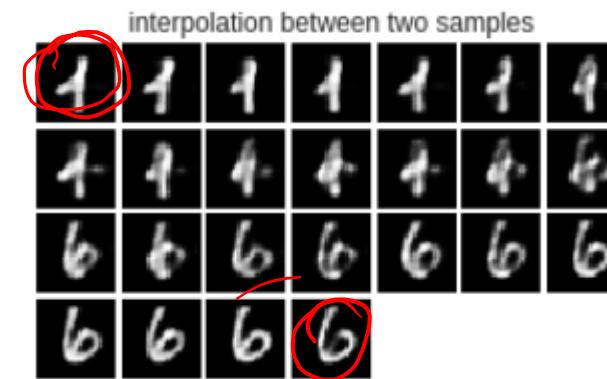
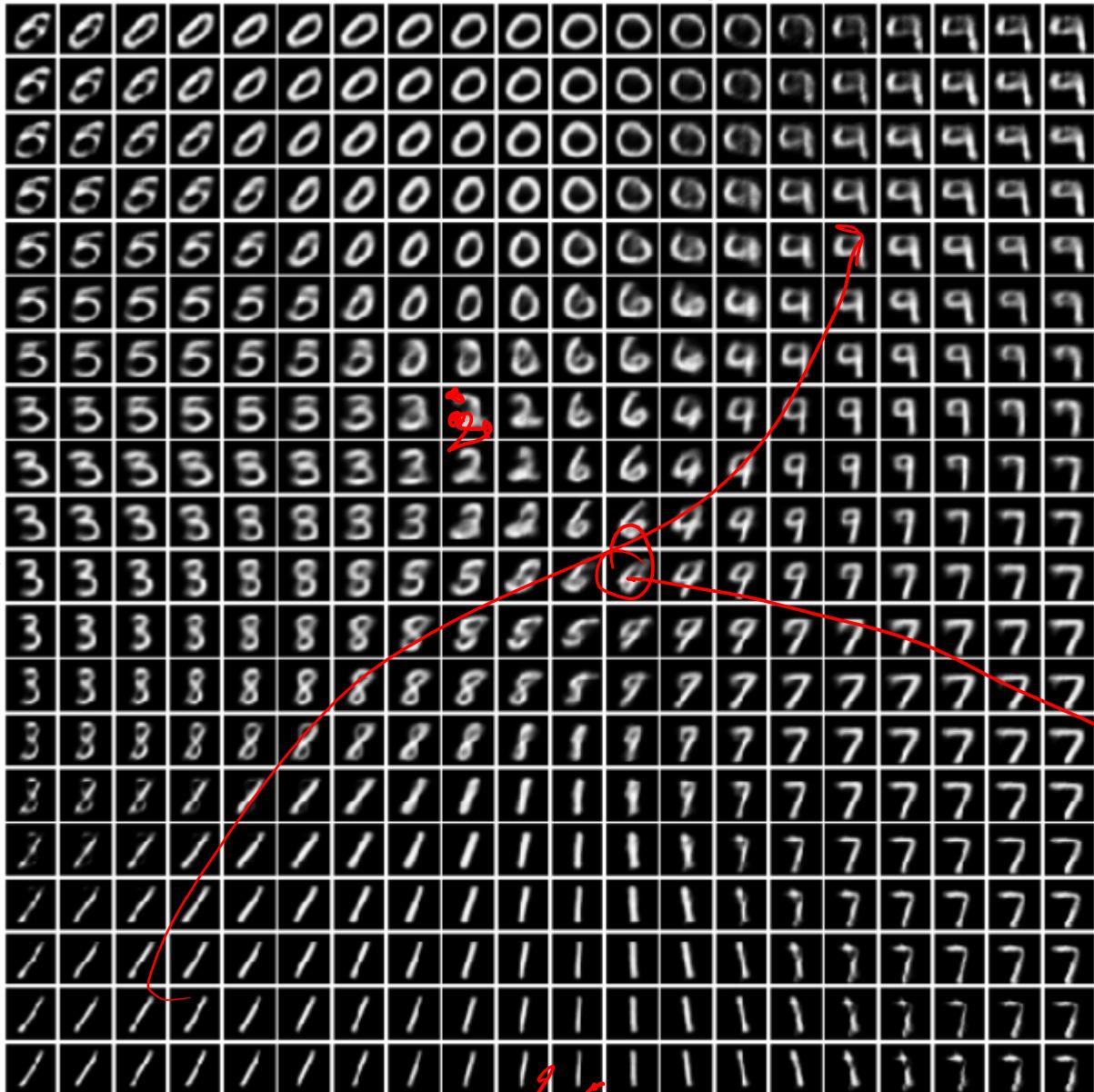
$z >$

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



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

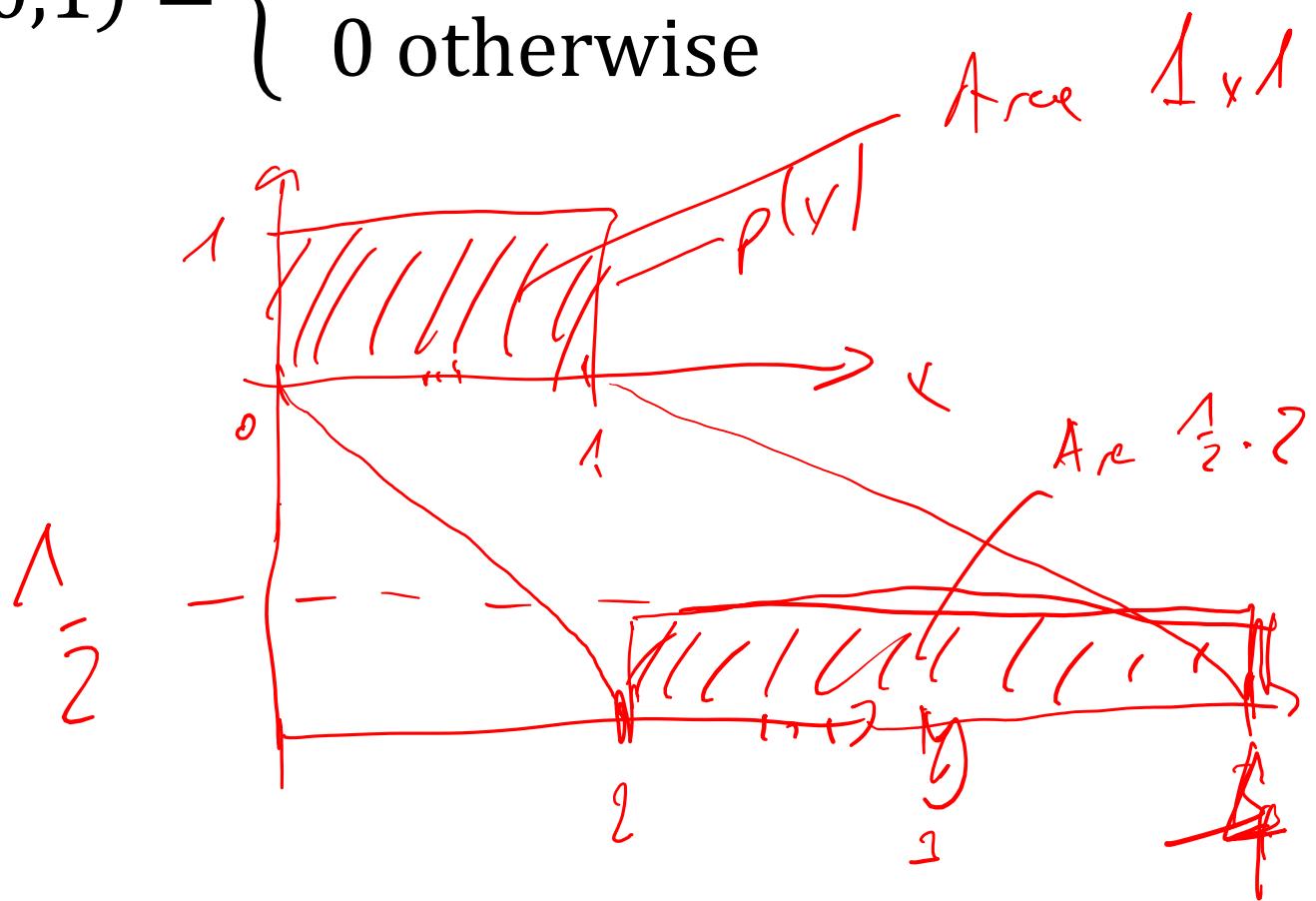
NORMALIZING FLOWS

Exercise

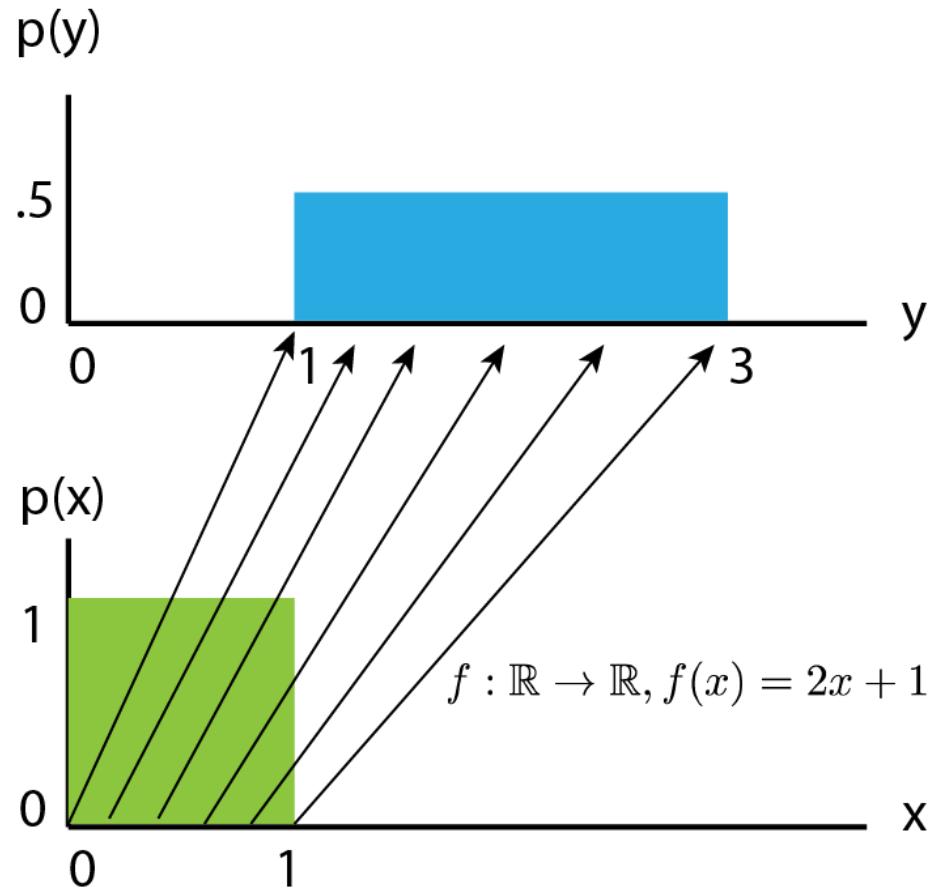
$$p(x) = \text{uniform}(0,1) = \begin{cases} 1 & \text{for } x \in (0,1) \\ 0 & \text{otherwise} \end{cases}$$

$$y = 2 * x + 1$$

$$p(y) = ?$$



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 $y = f(x)$

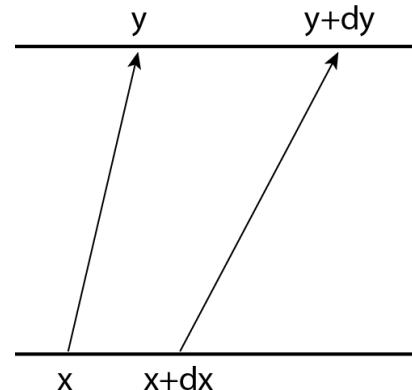
Take a range $(x, x + \Delta x)$

Question:

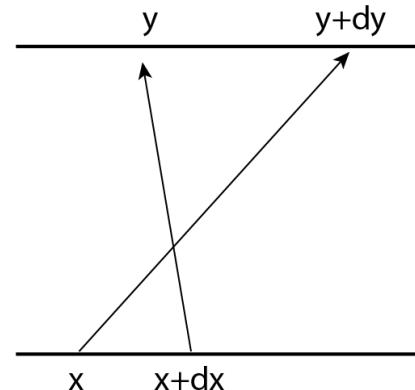
How long is this range?

Δx

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



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



Question:

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

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

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

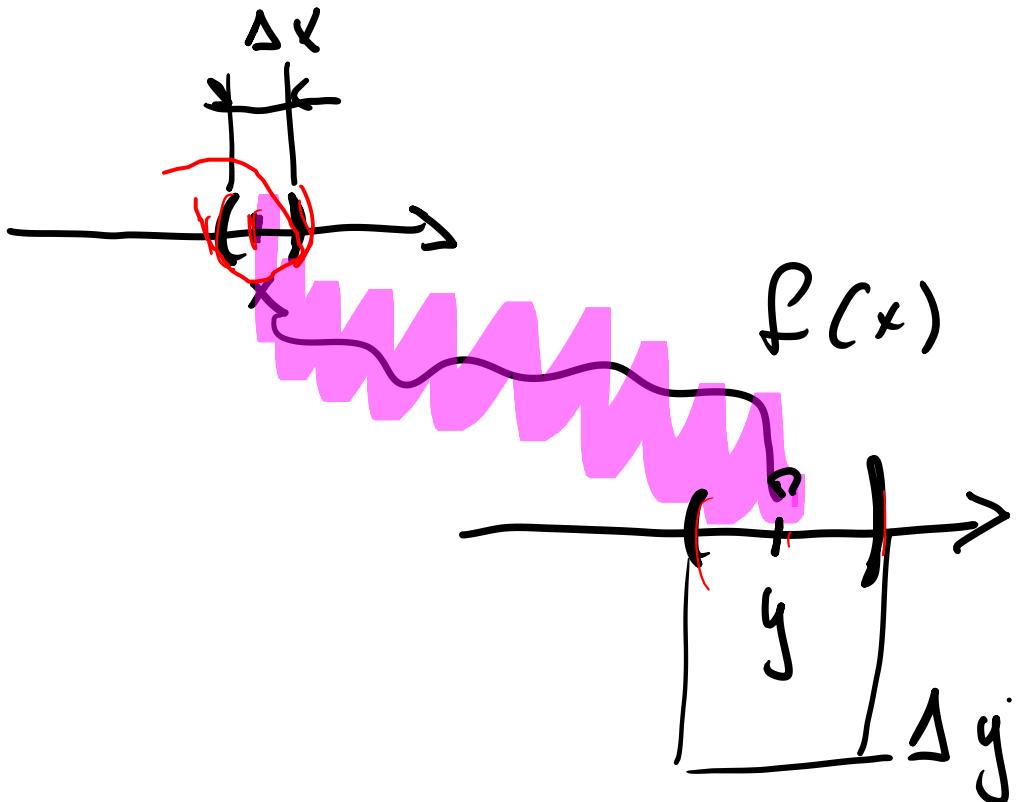
Dot of Jacobian

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:

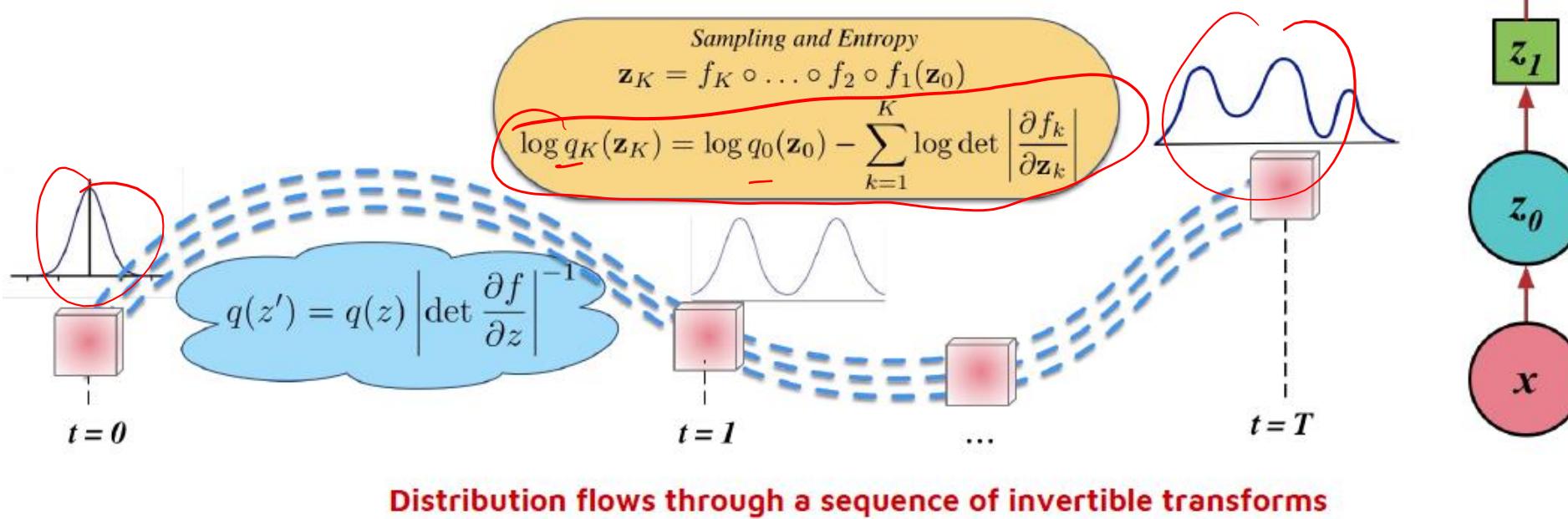
- f is easy to exactly invert
 f and f^{-1} form an exact auto-encoder!
- $\det \frac{\partial f(x)}{x}$ is easy to compute

=> We need special f !

Normalizing Flows

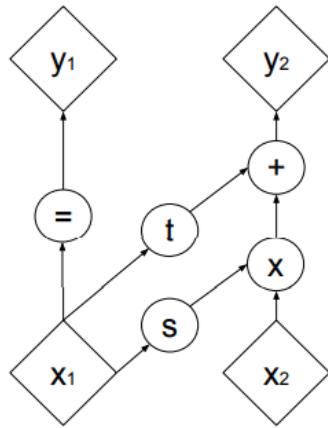
Exploit the rule for change of variables:

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

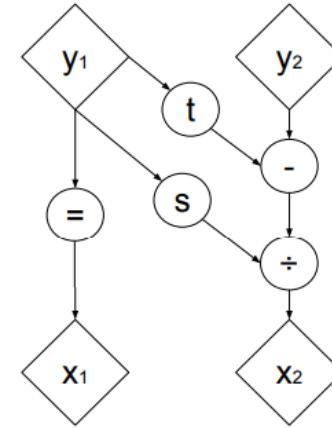


Rezende and Mohamed, 2015

A special form of f



(a) Forward propagation



(b) Inverse propagation

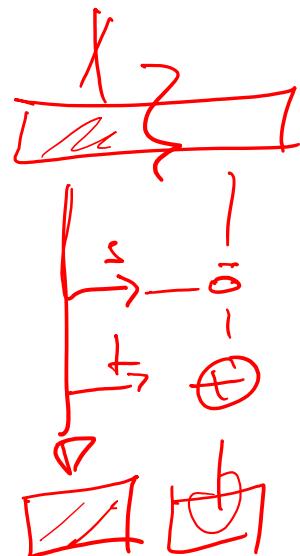
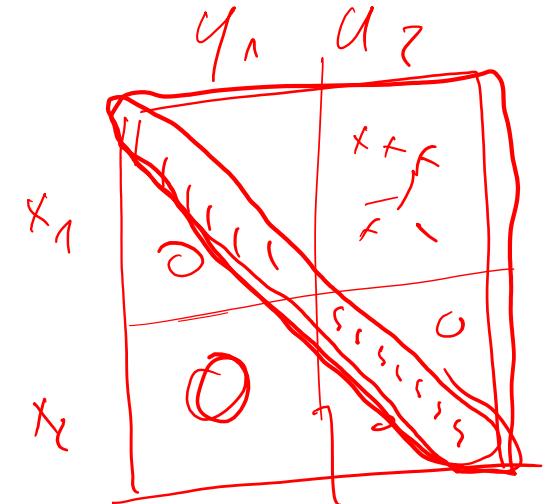
$$x_1, x_2 = \text{split}(x)$$

$$\underline{y_1 = x_1}$$

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

Trivial to invert!

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



Normalizing flows in action

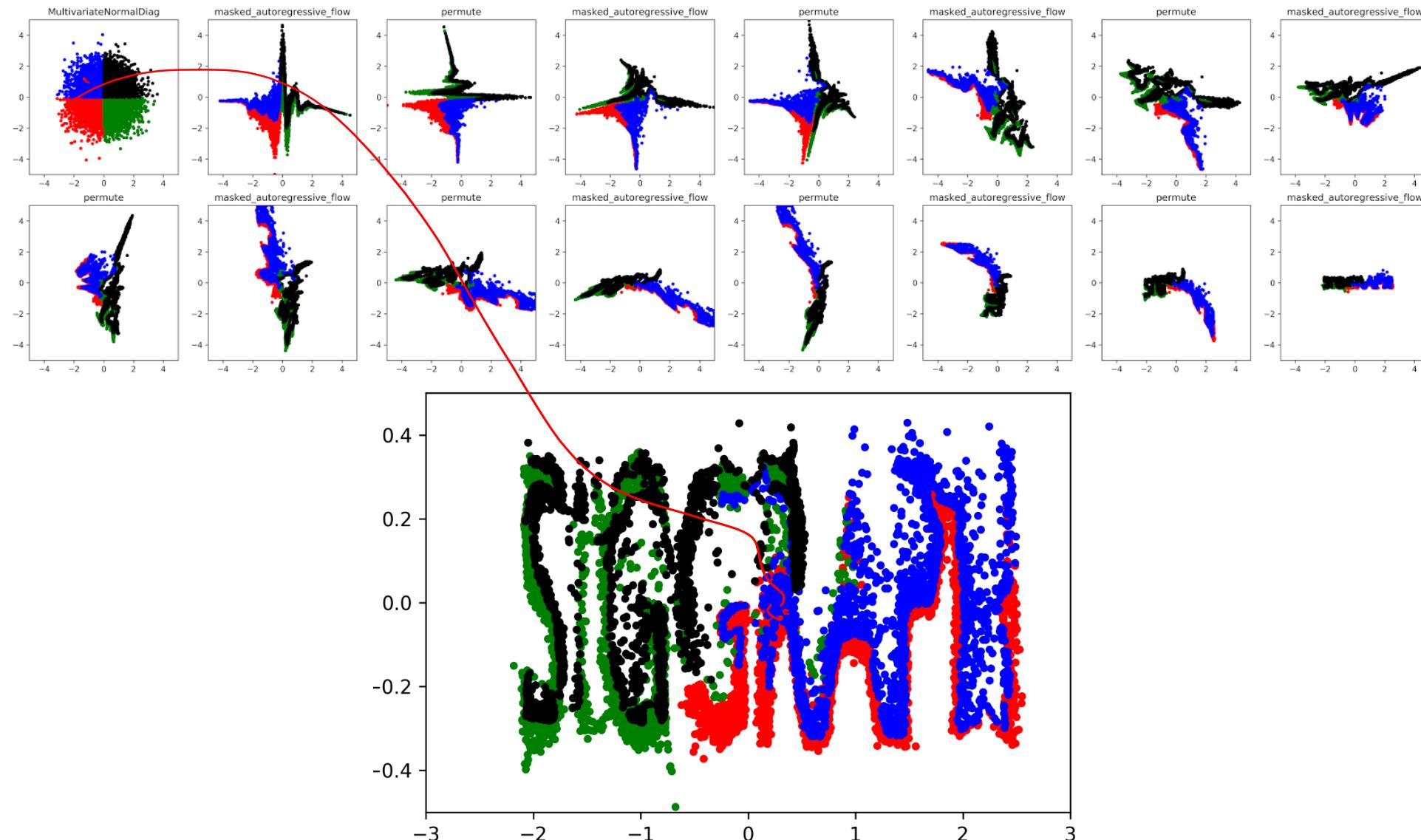
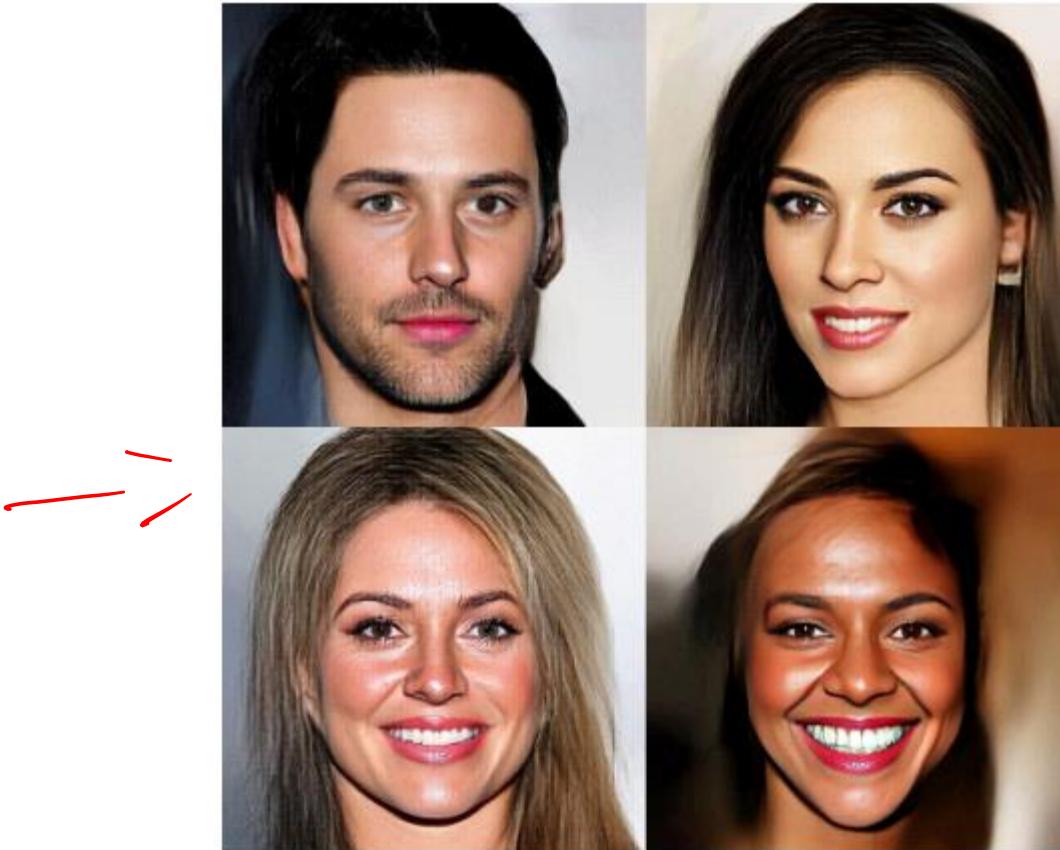
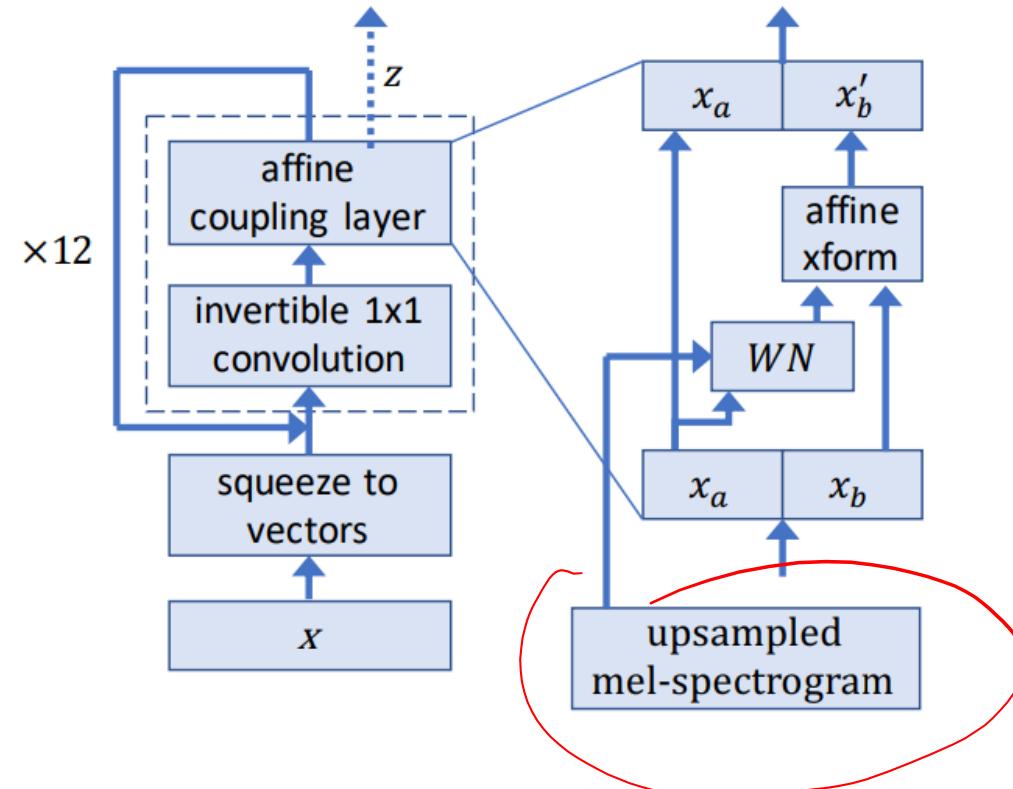


Image credit Eric Jang, <https://blog.evjang.com/2018/01/nf2.html>

Normalizing flows in action



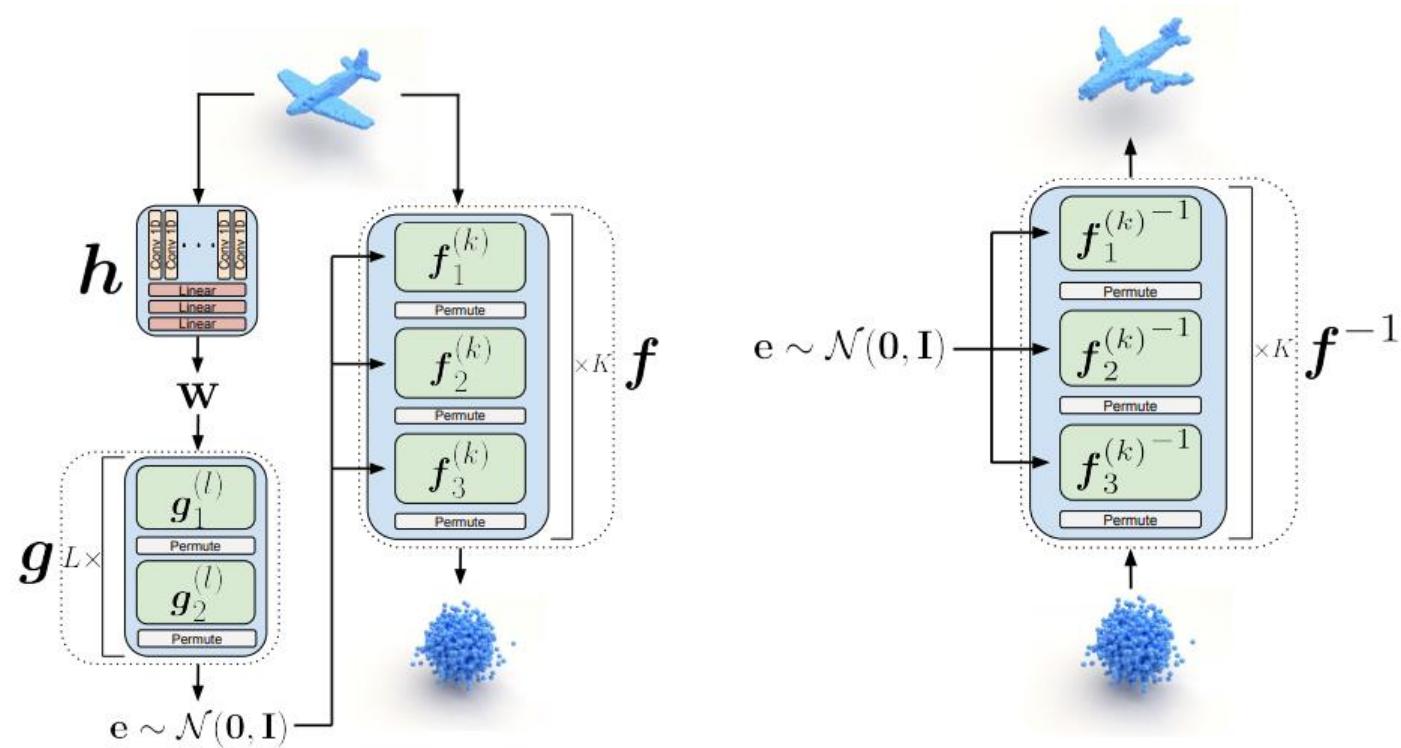
Kingma et al, “GLOW”



Prenger et al, “WAVEGLOW”

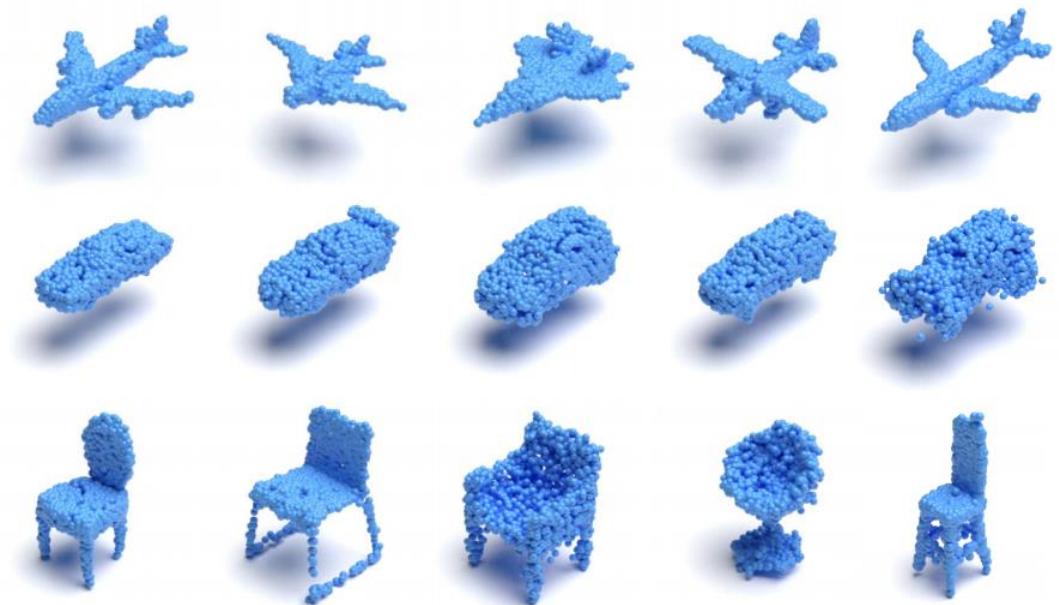
Normalizing Flows in Action

A generative model for point clouds:

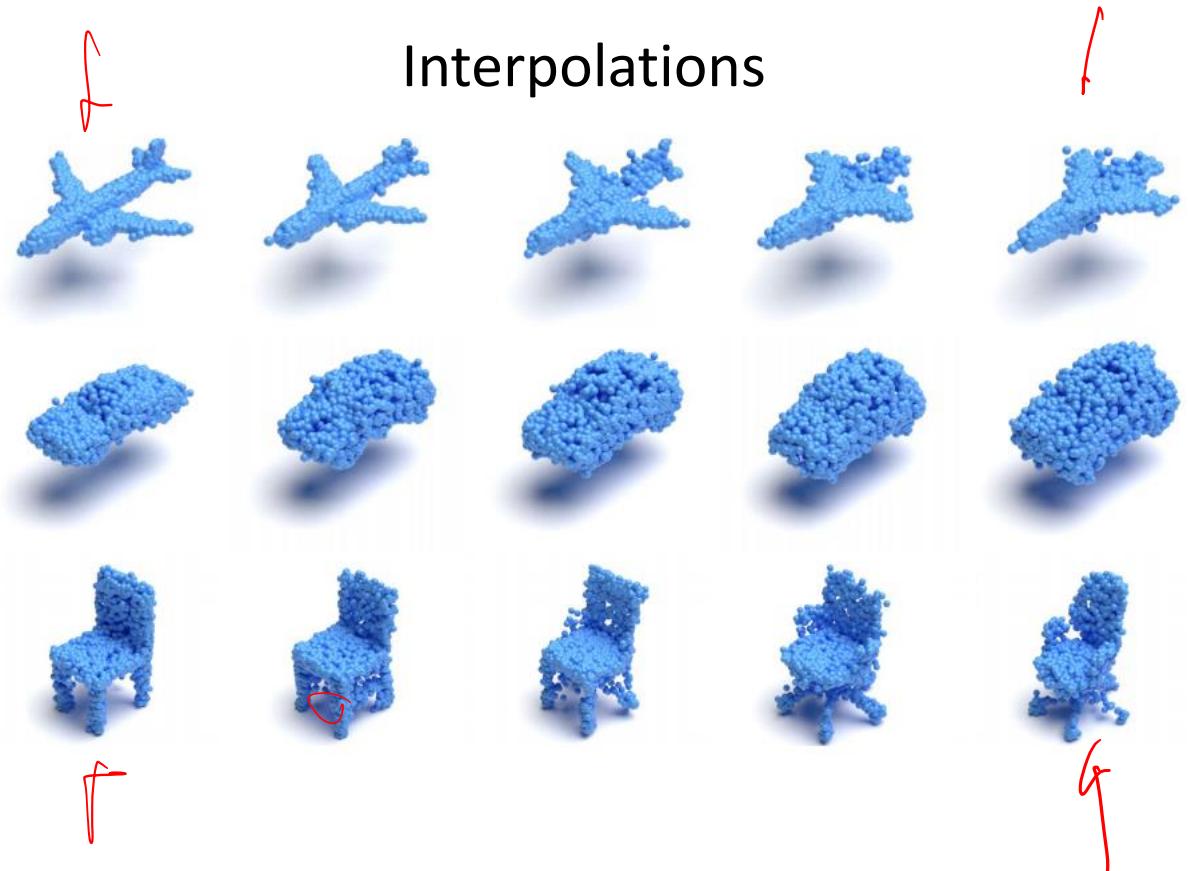


Normalizing Flows in Action

Samples



Interpolations

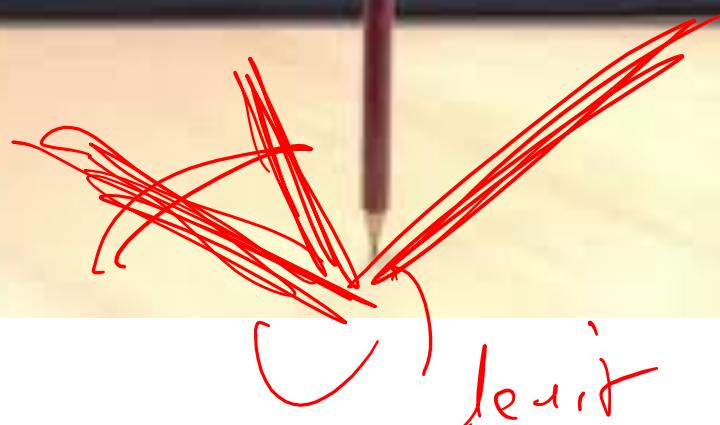
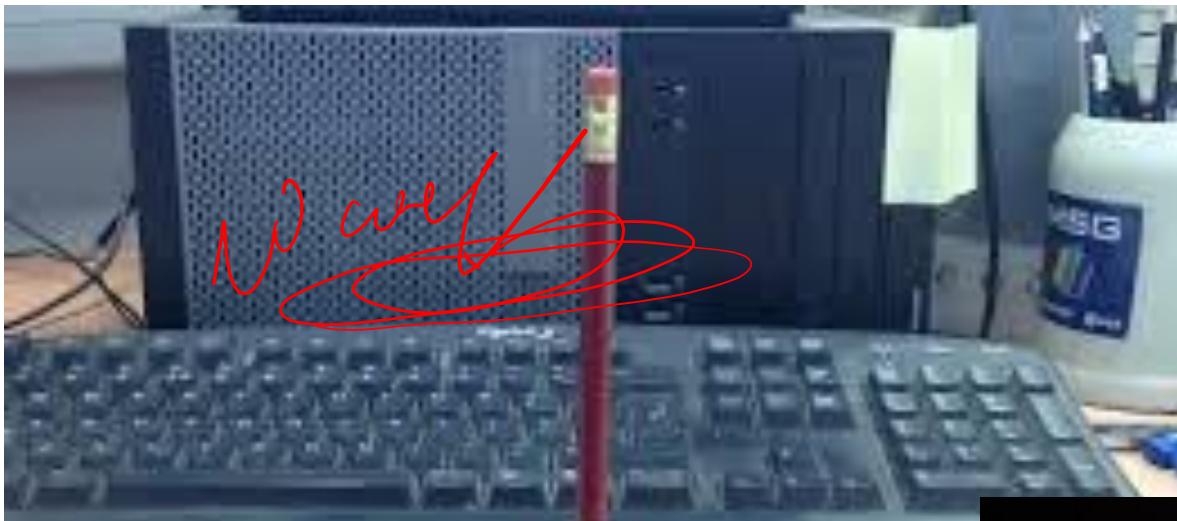


GAN

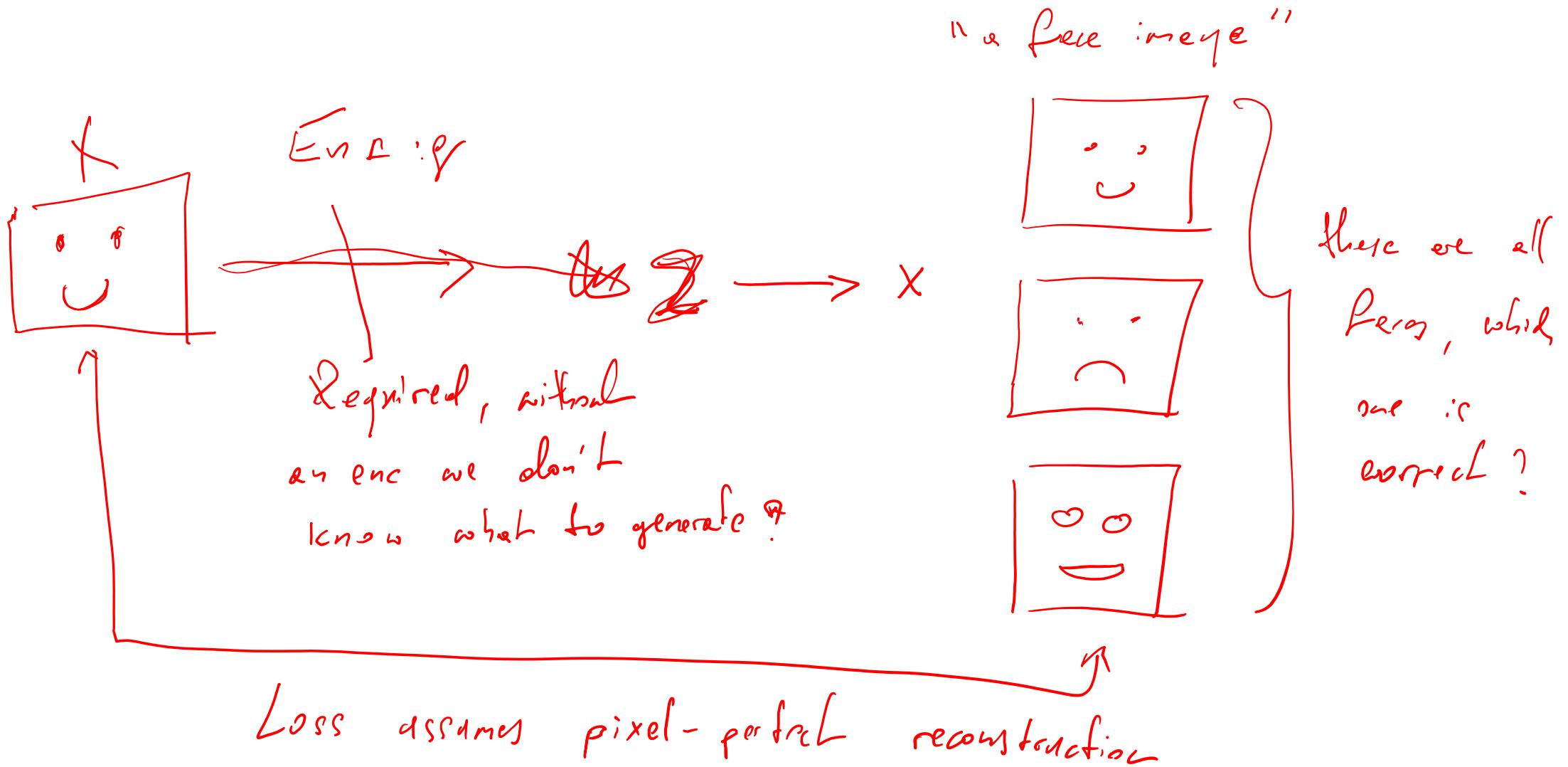
What will happen next?

f_0

$p(f_2 | f_0, f_1)$

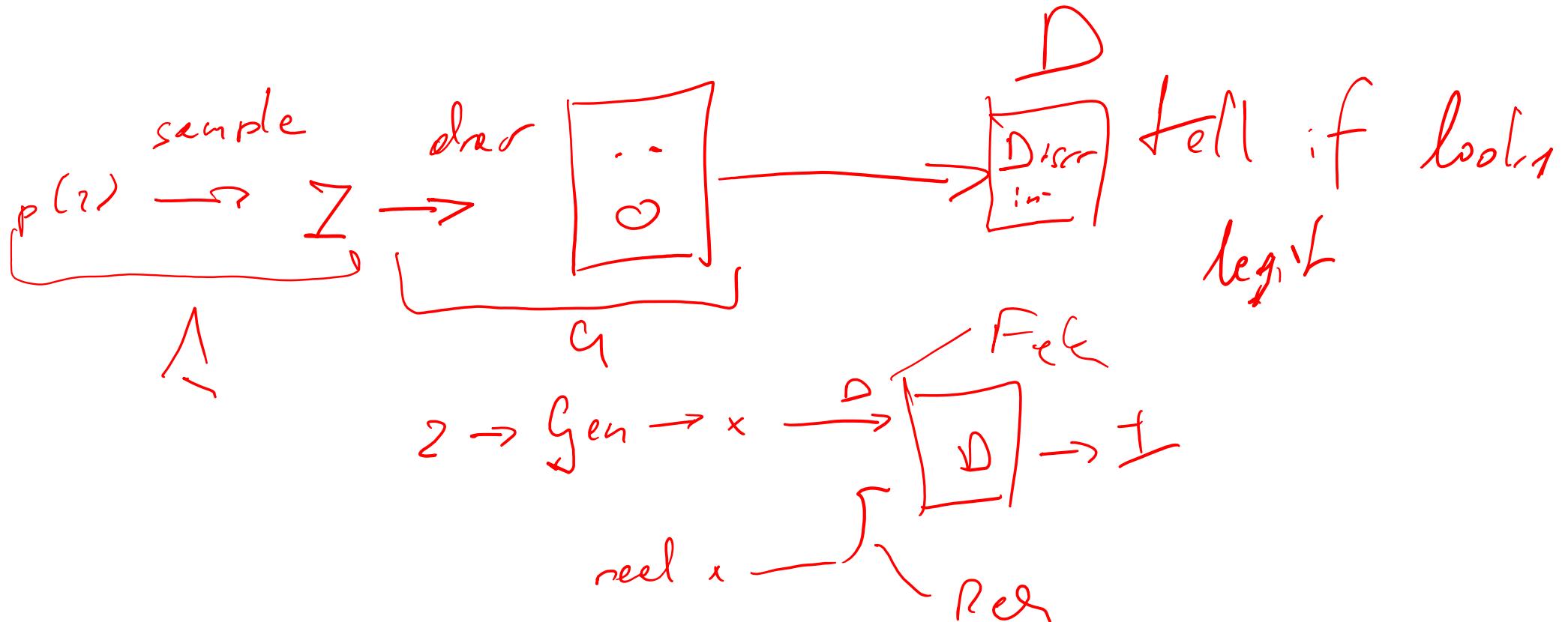


Trouble with autoencoders

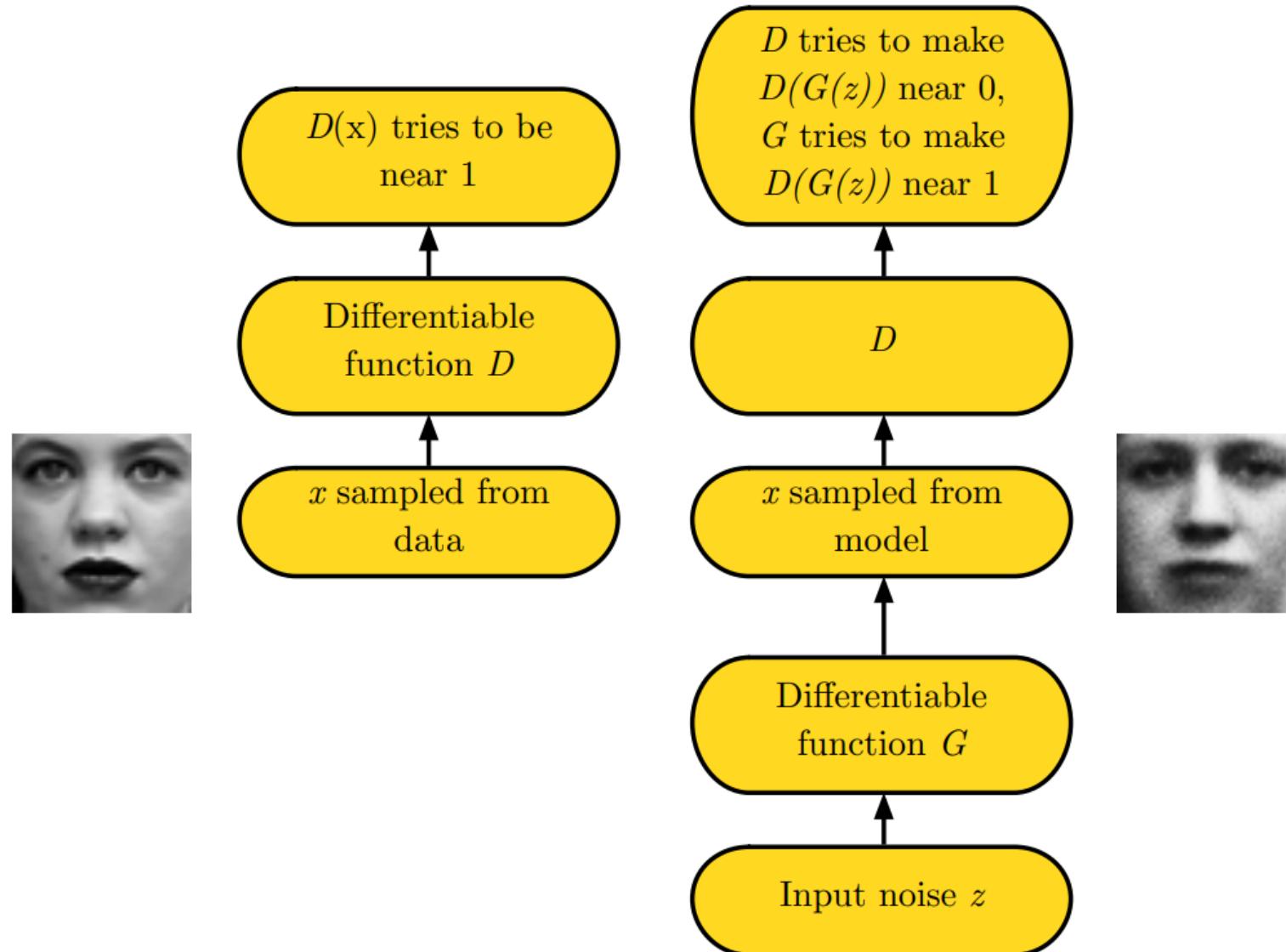


GAN idea

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



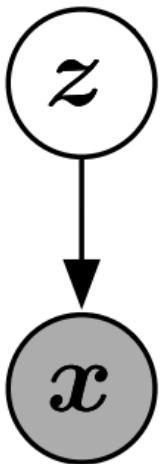
Adversarial Nets Framework



(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

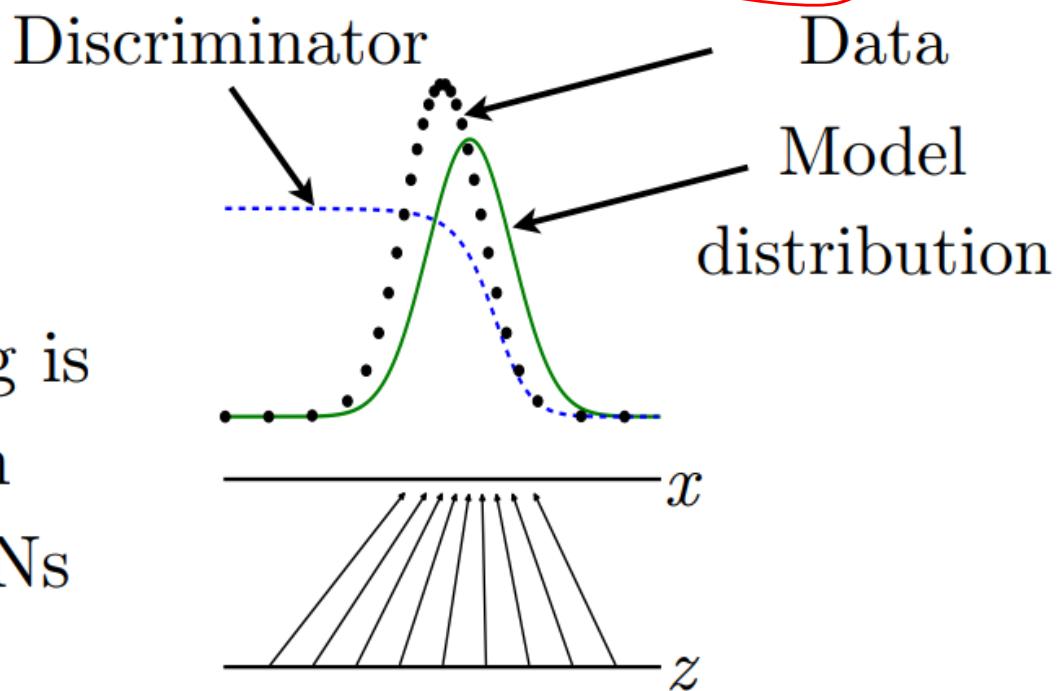
$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$
$$J^{(G)} = -J^{(D)}$$

- What is the solution to $D(x)$ in terms of p_{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

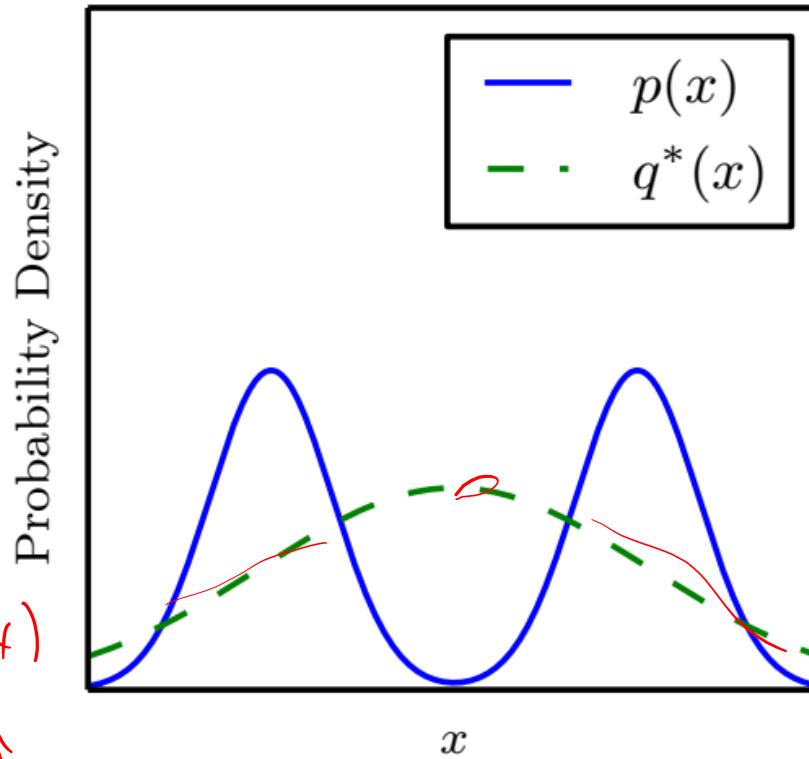
$$D(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}$$



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

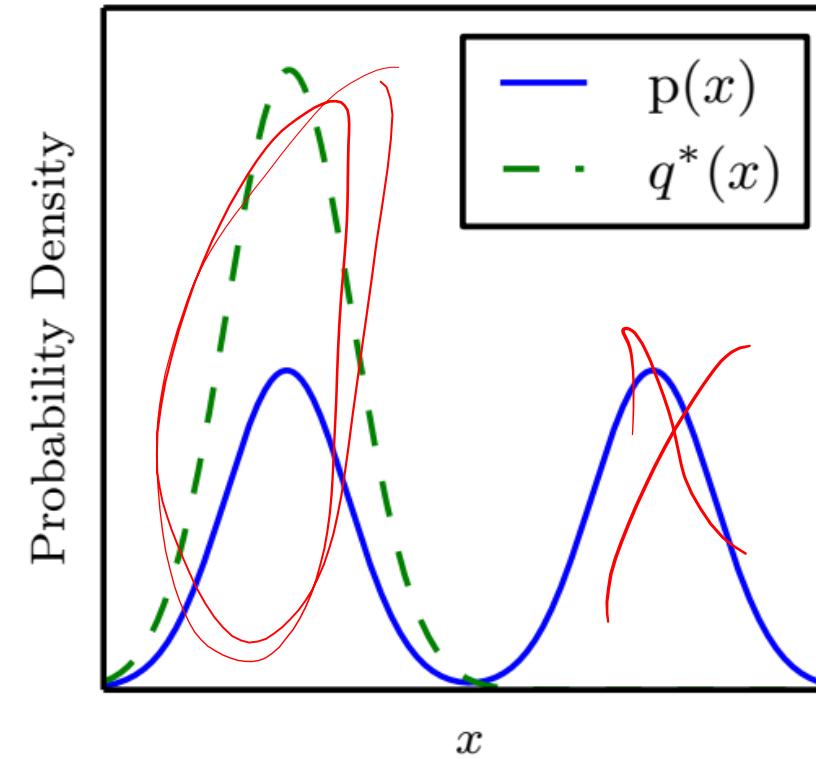
Is the divergence important?

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



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

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



Reverse KL
CVAU

GAN Trick

Non-saturating loss

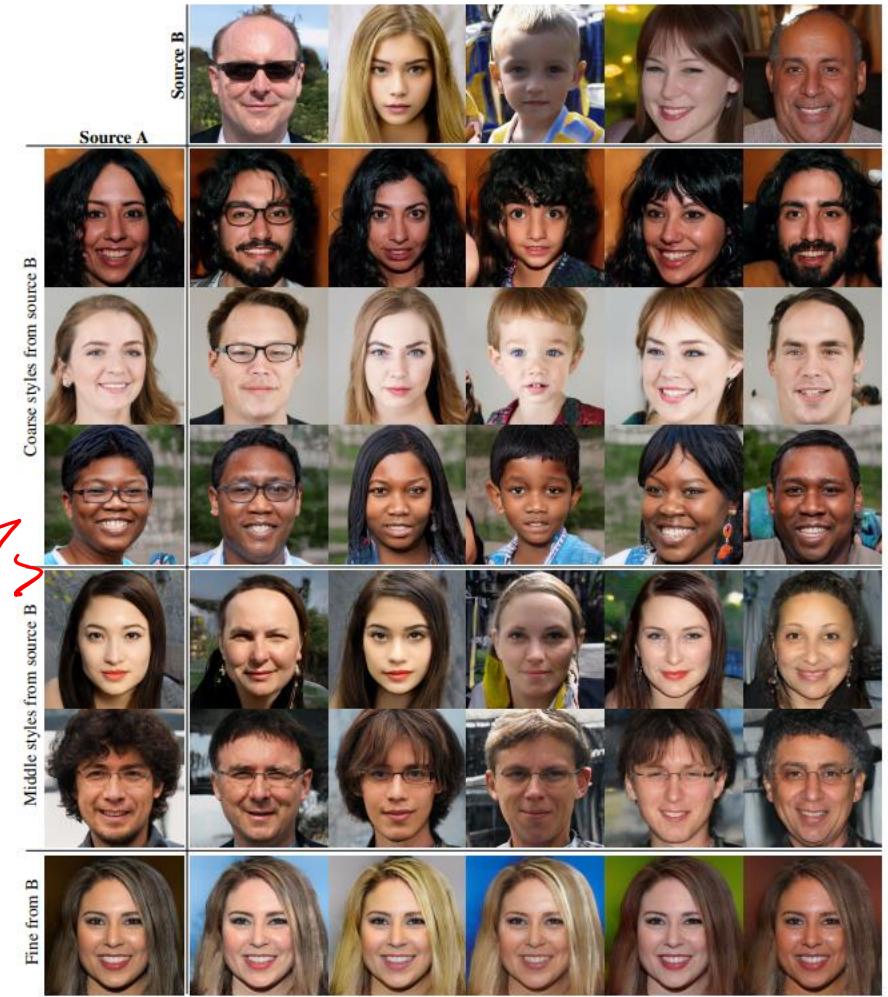
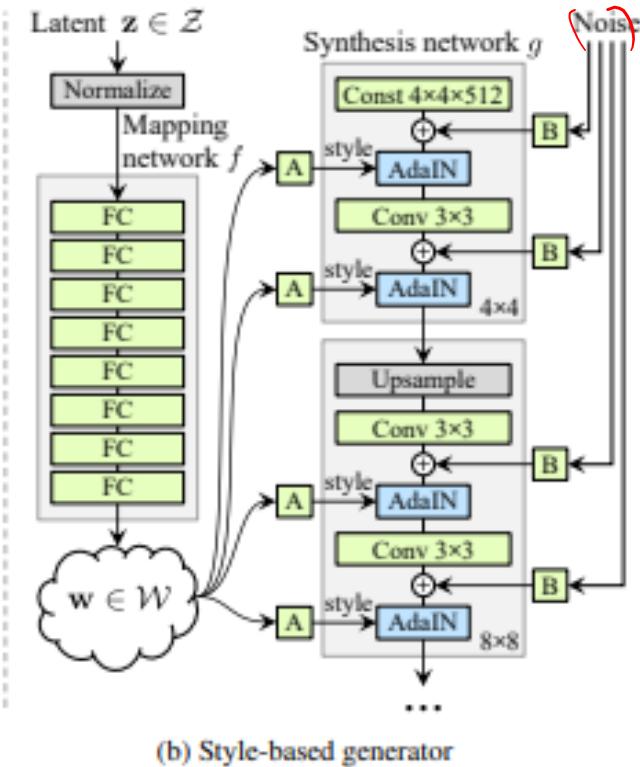
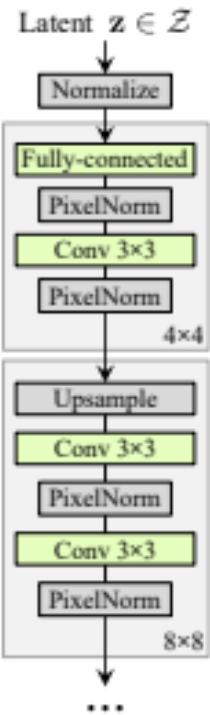
$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

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

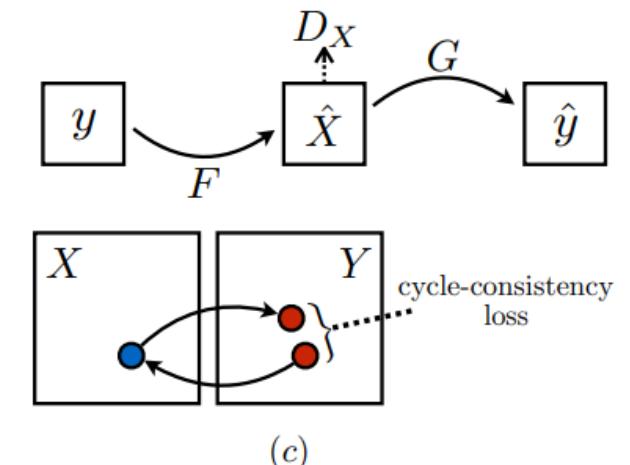
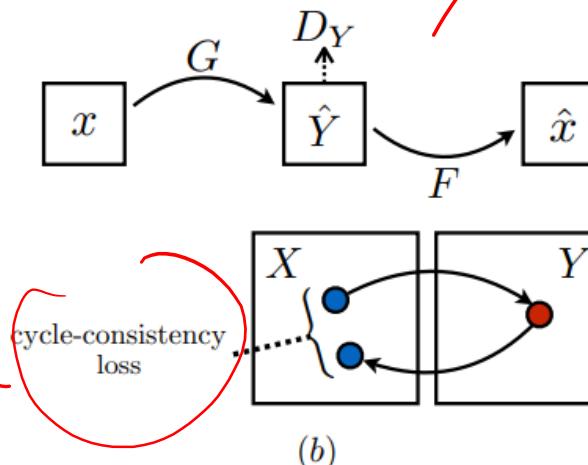
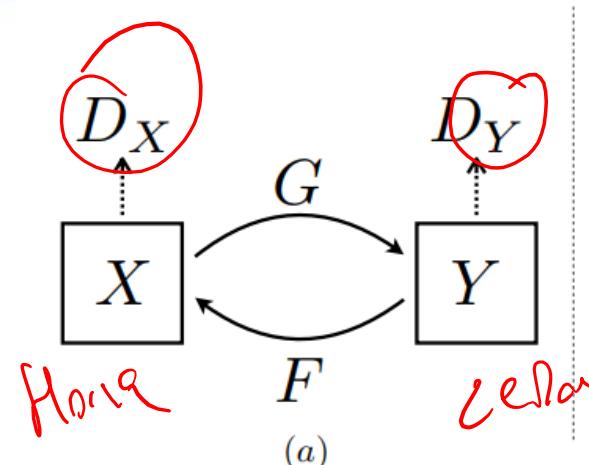
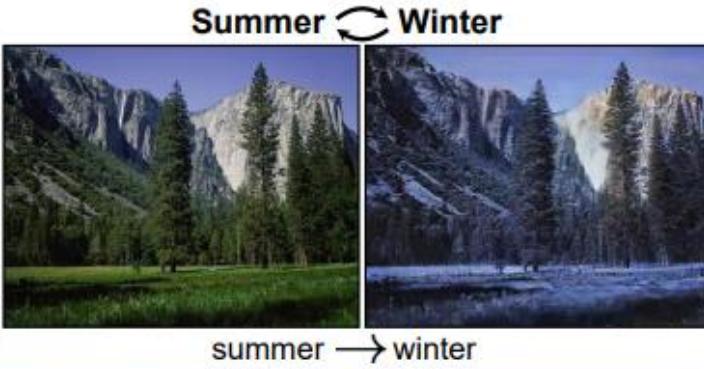
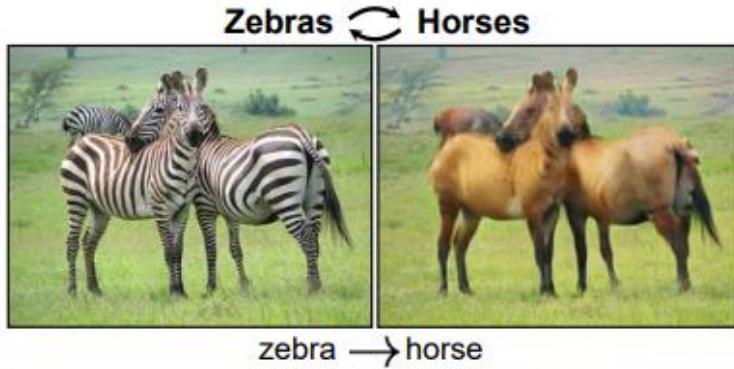
$$J^G \doteq -J^D - \cancel{\frac{1}{2} \mathbb{E}_{\mathbf{z}} \left(\log \cancel{D}(G(\mathbf{z})) \right)}$$

GAN Uses: image generation

2 min



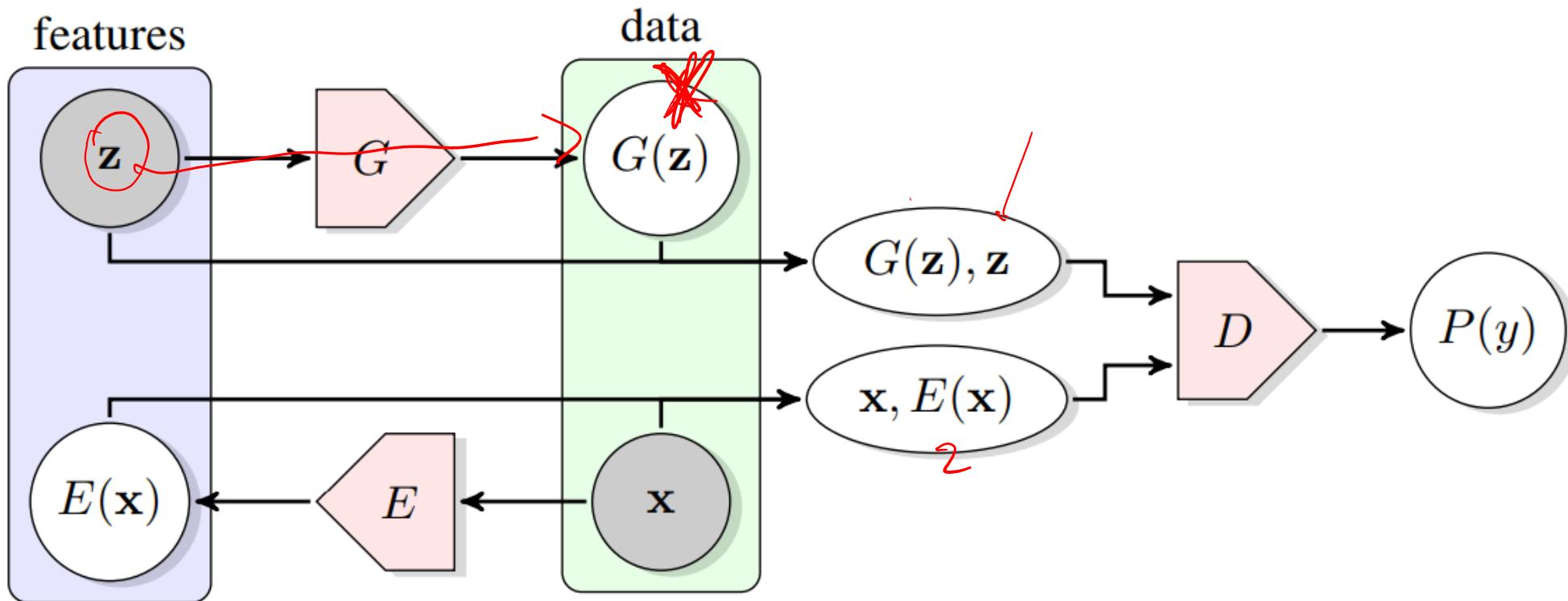
GAN Uses: unpaired image-image translation



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	7	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	1	1	1	1
8	8	8	8	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	3	3	3	3
9	9	9	9	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	5	5	5	5

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

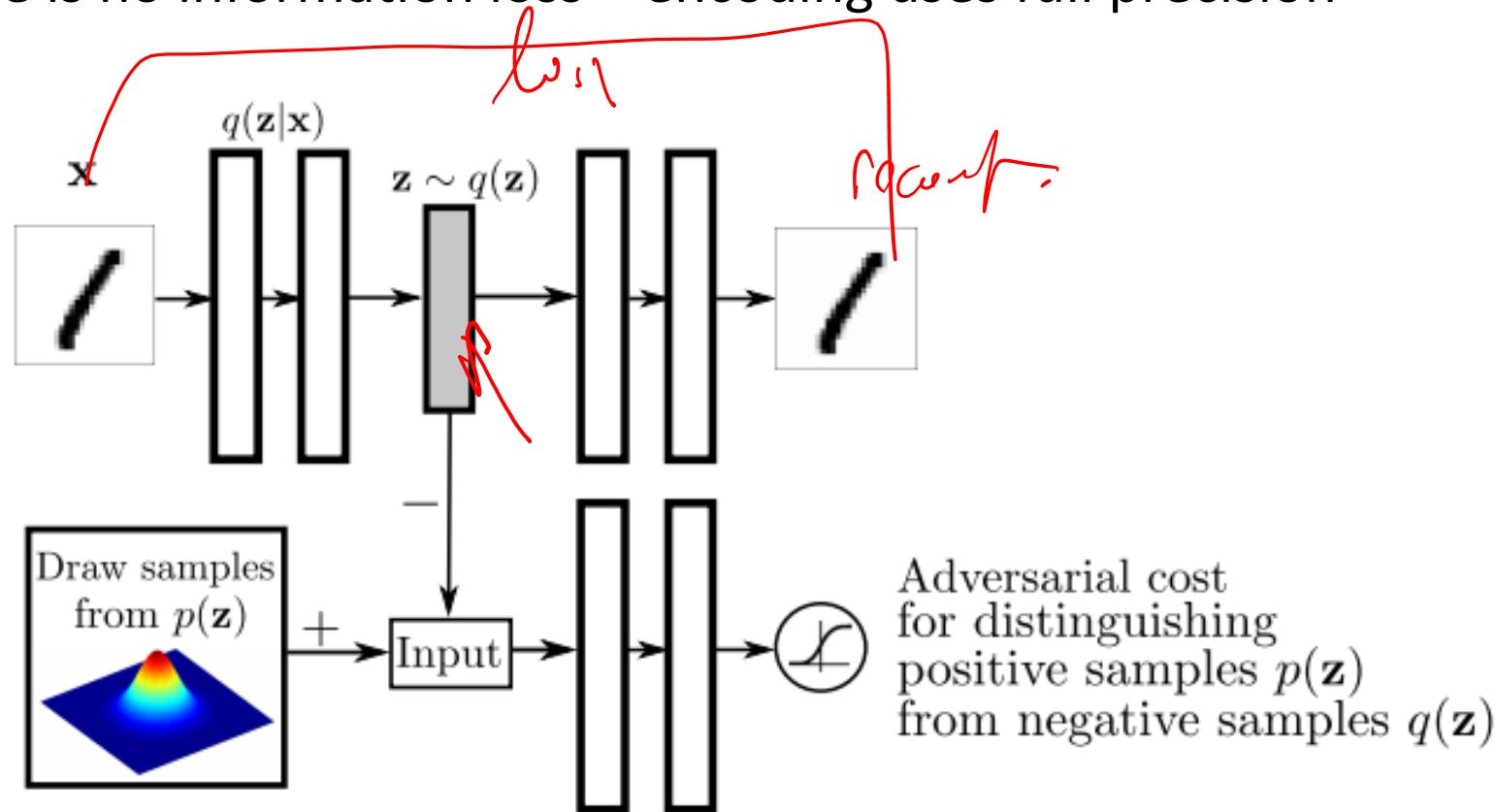
(d) Varying c_3 from -2 to 2 on InfoGAN (Width)

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

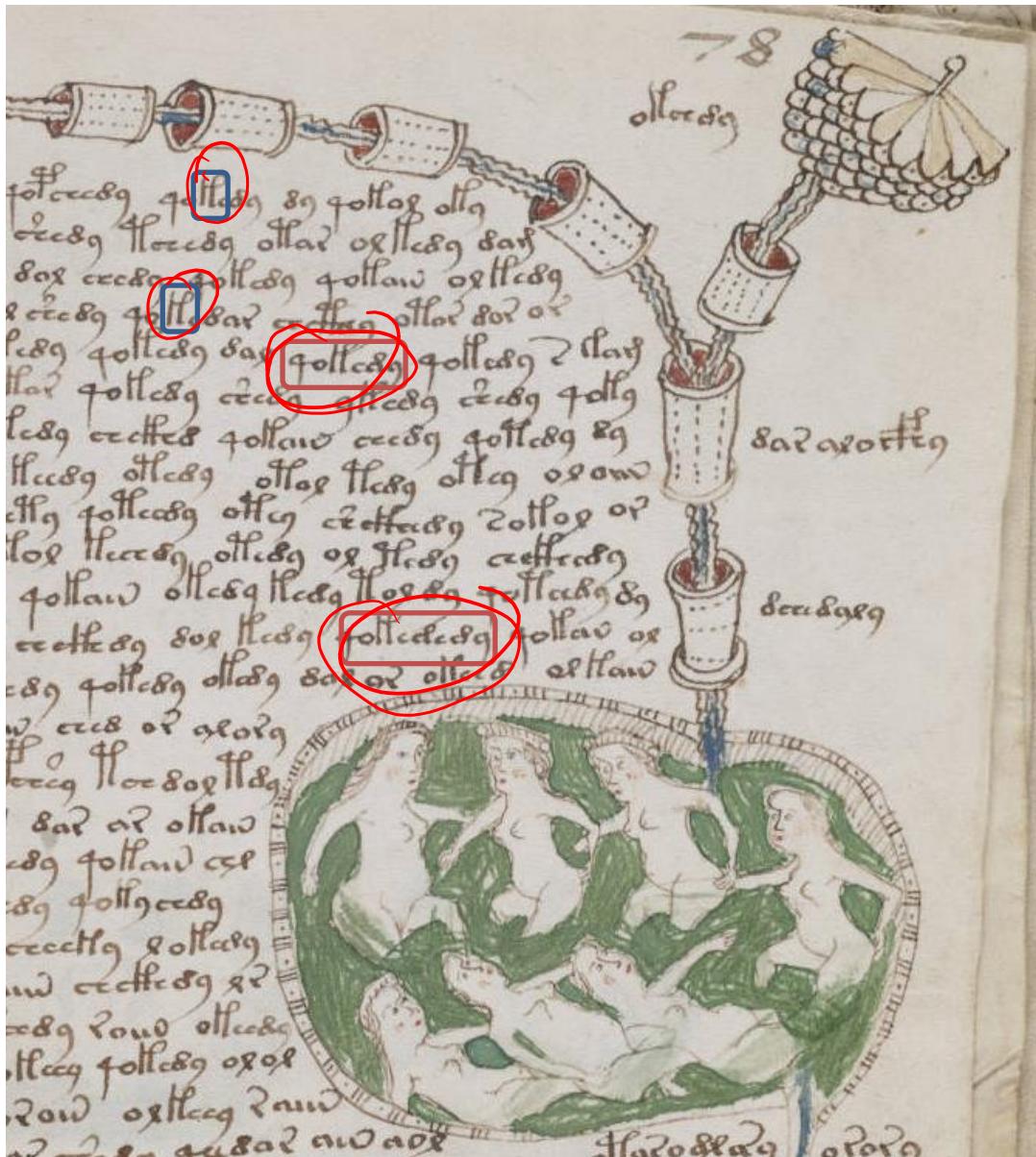
Unlike VAE, there is no information loss – encoding uses full precision



MY WORK:

VAE + AUTOREGRESSIVE DECODERS

Unlabeled data uses

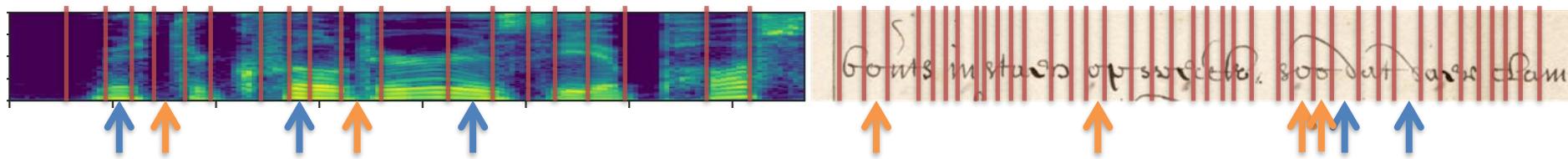


What can we learn without labels:

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

Our goal

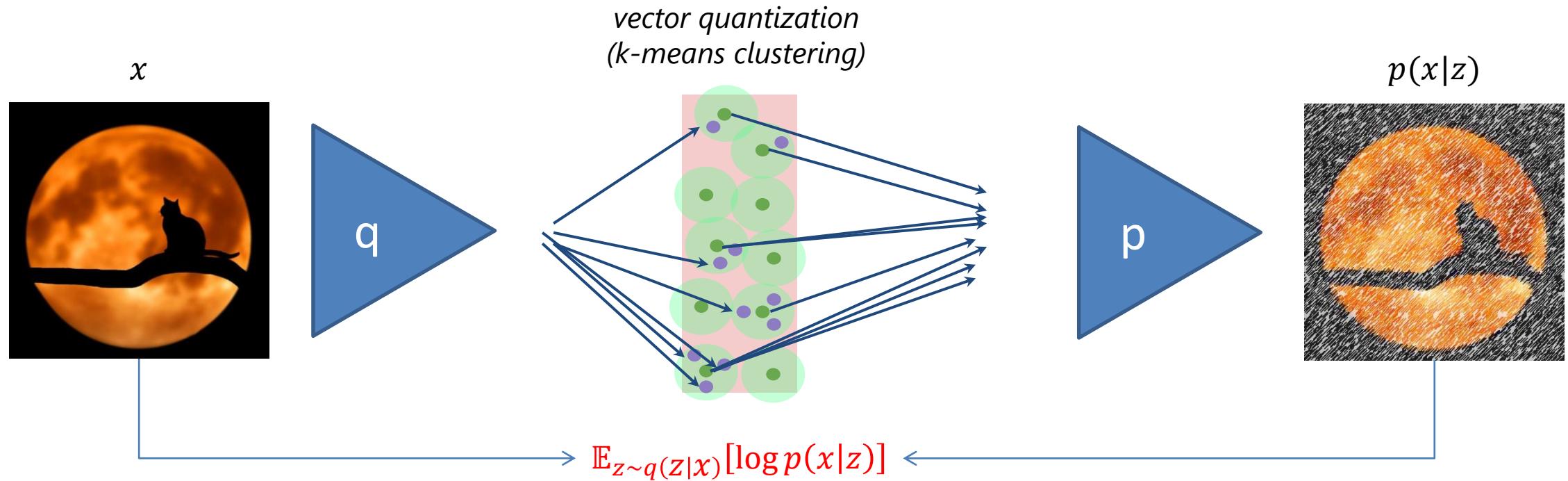
Unsupervised unit discovery
in speech and handwriting



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

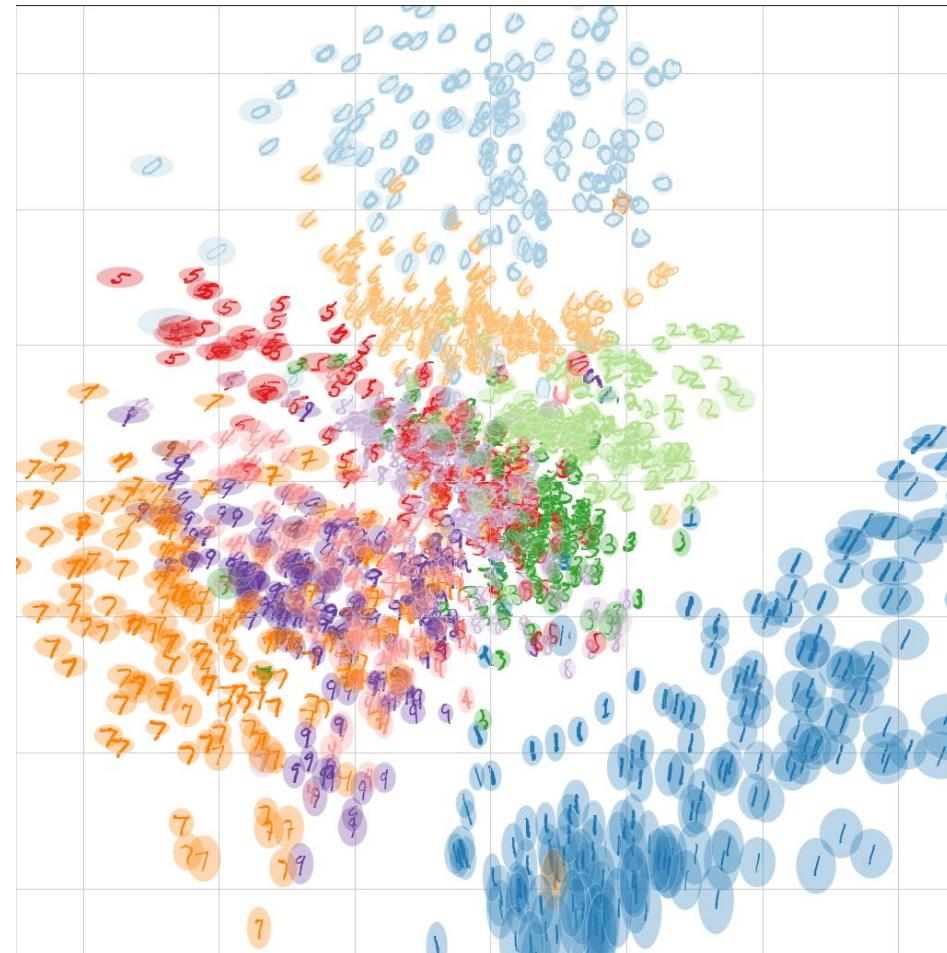
it is replaced by one of K prototypes

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

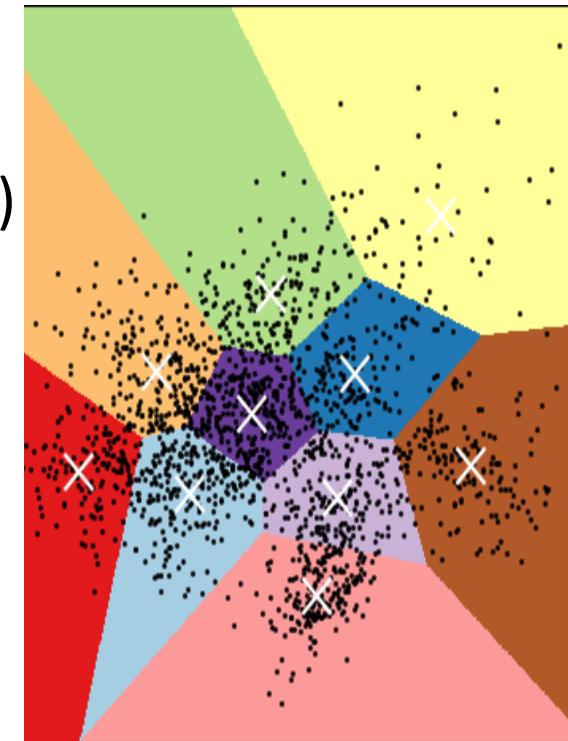
Limit precision of the encoding by quantizing (round each vector to a nearest prototype).

Output can be treated:

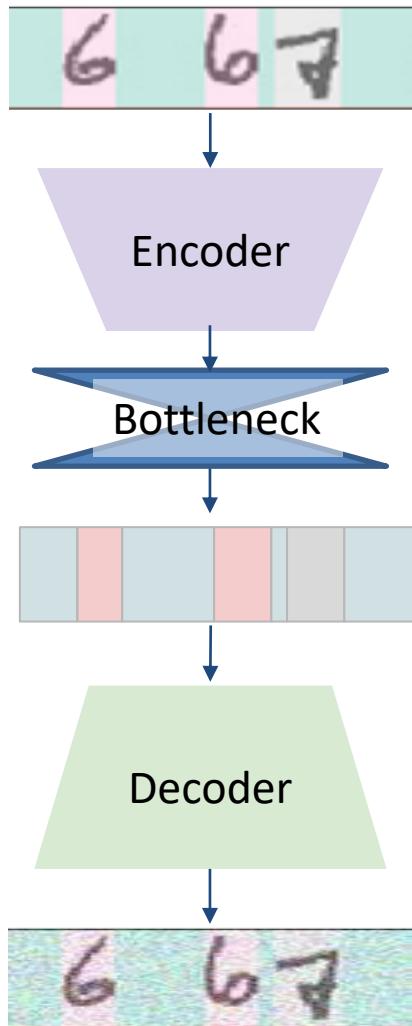
- As a sequence of discrete prototype ids (tokens)
- As a distributed representation (the prototypes themselves)

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

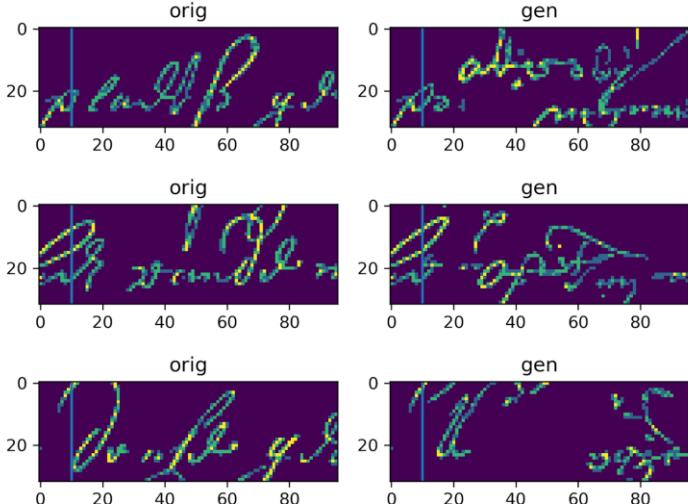
Latent encoding:
only high level features

Decoder regenerates the inputs.

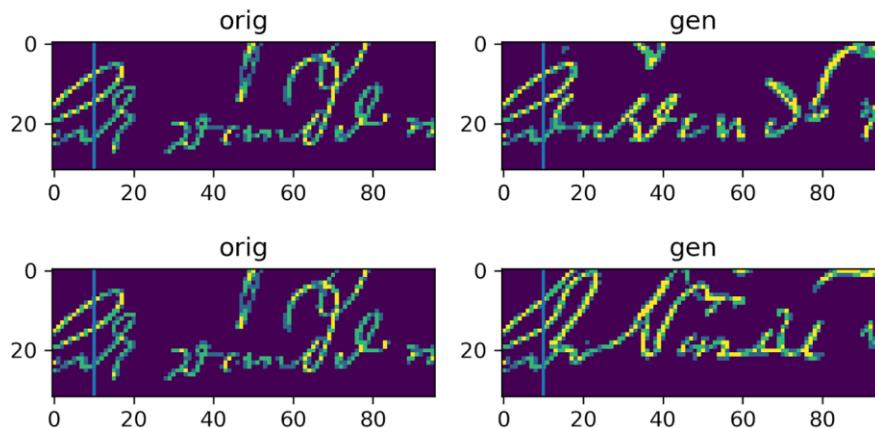
It learns $p(x|z)$ and it must fill in all discarded information.

Conditional PixelCNN can regenerate images.

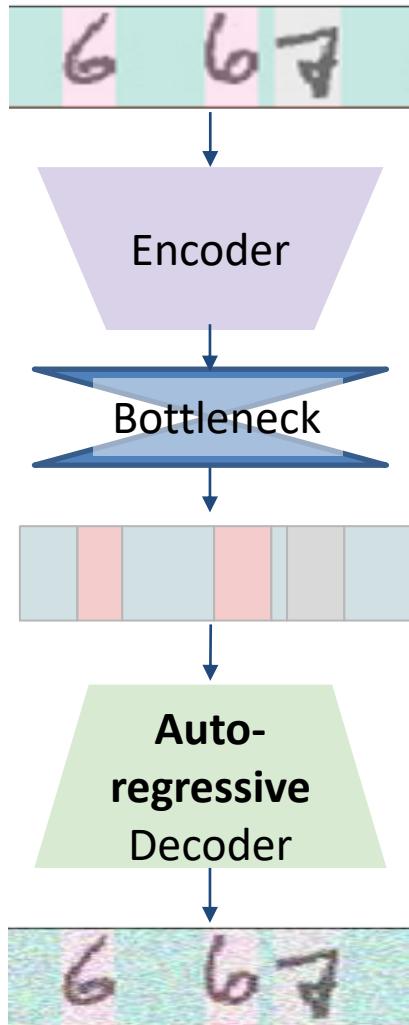
Unconditional



Conditioned on text



Our model: Autoencoder + autoregressive decoder

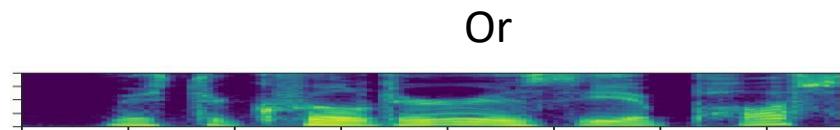
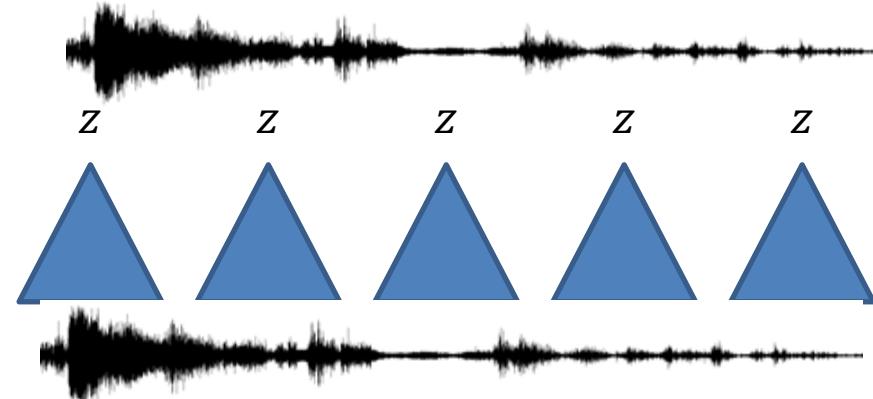
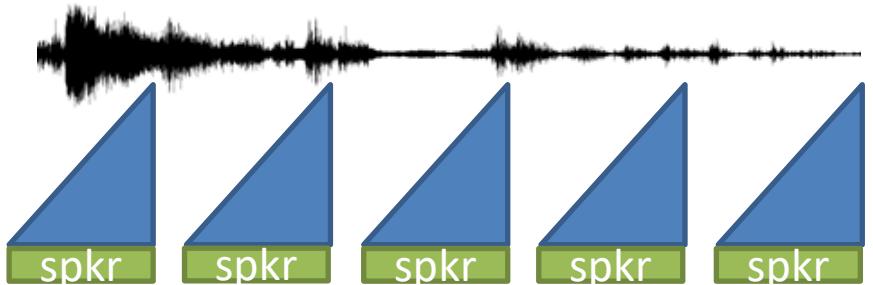


Decoder computes a probability over input images/waveforms conditioned on the encoding.

It combines information from:

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

Experiments on Speech

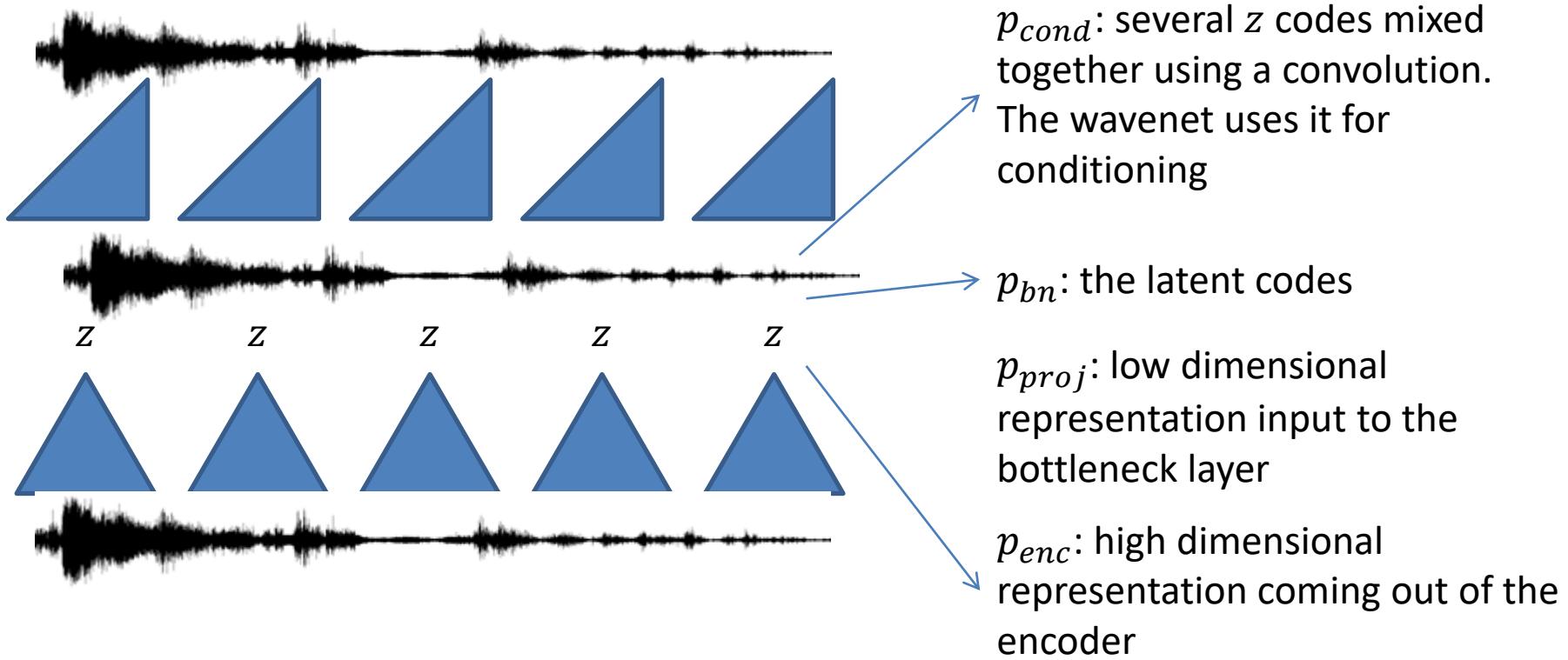


Input:

Waveforms, Mel Filterbanks, MFCCs

Representation probing points

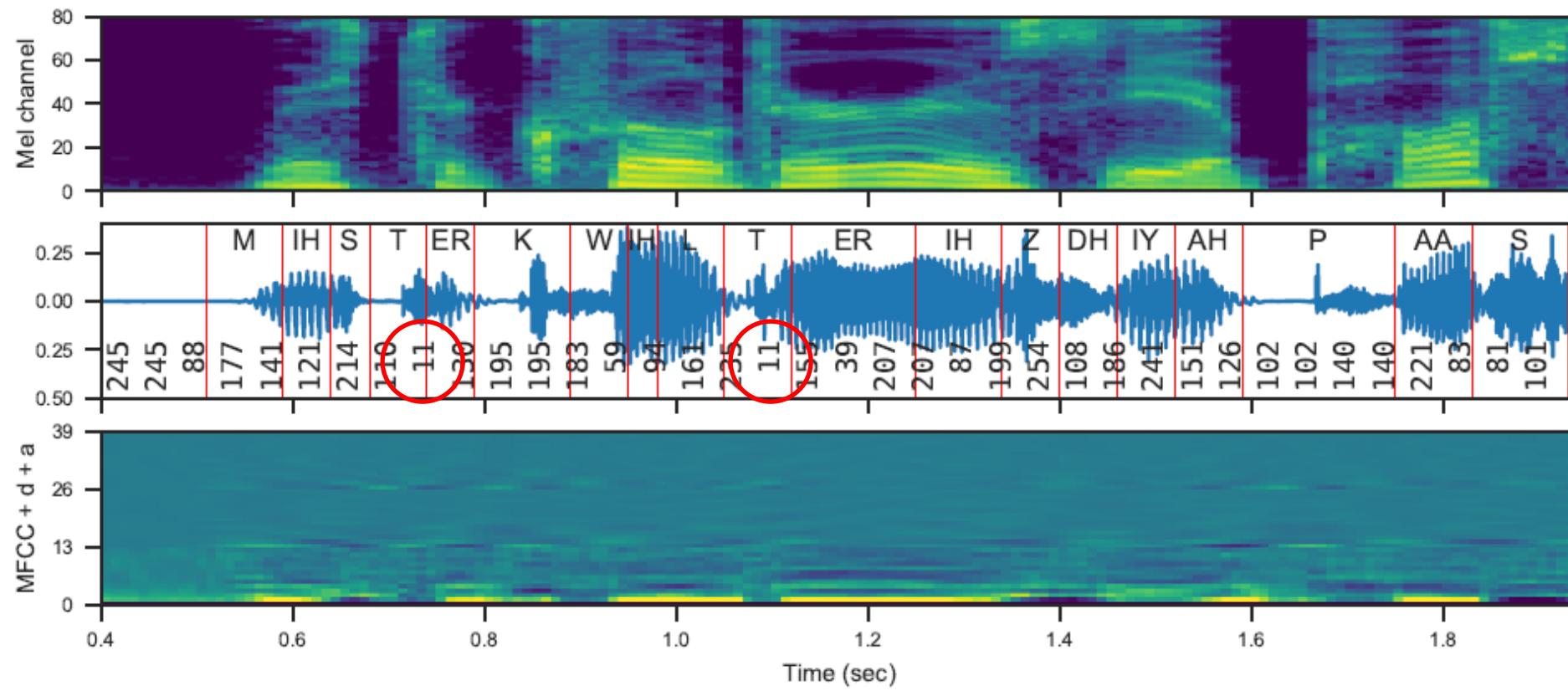
We have inserted probing classifiers at 4 points in the network:



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

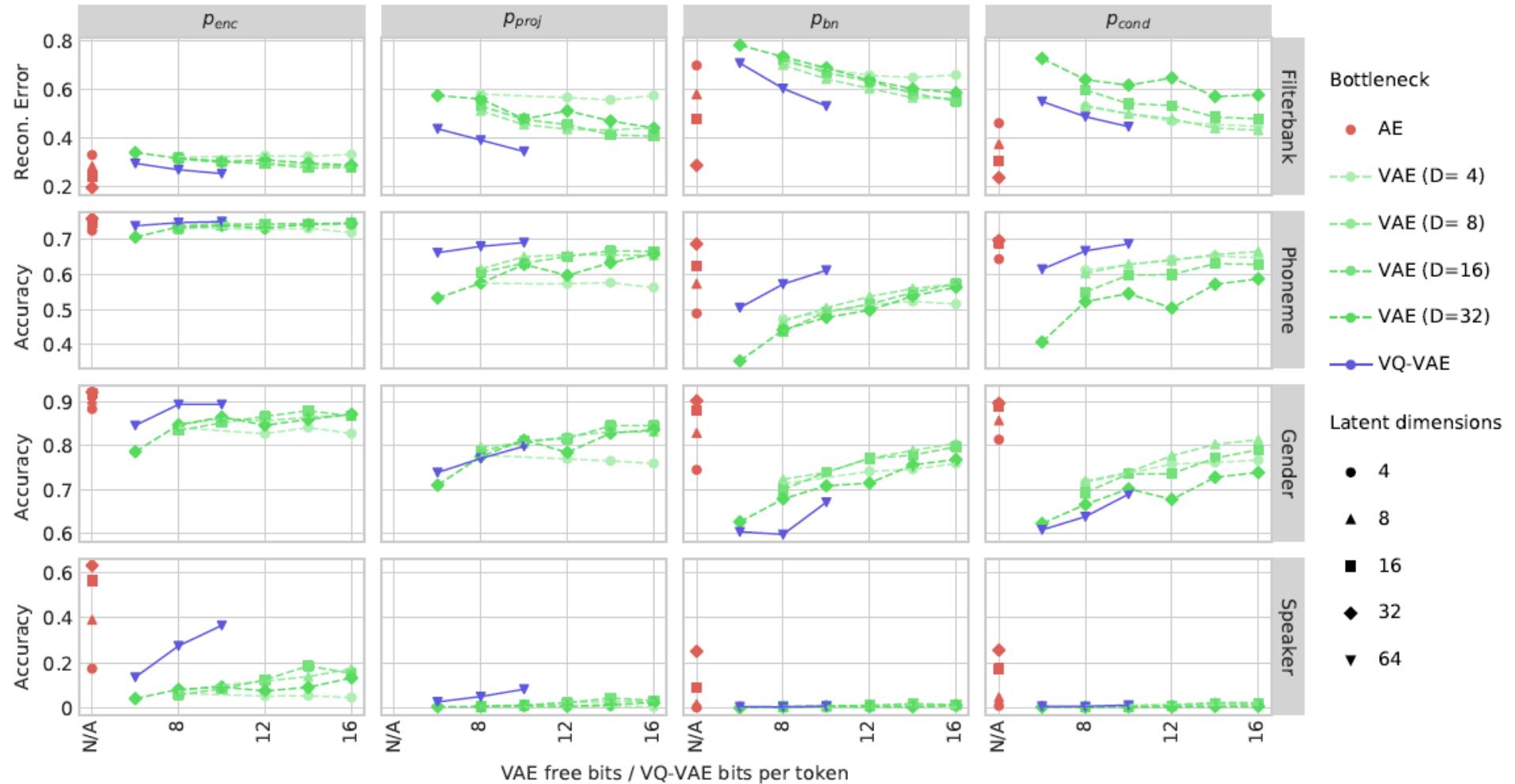


What information is captured in the latent codes?

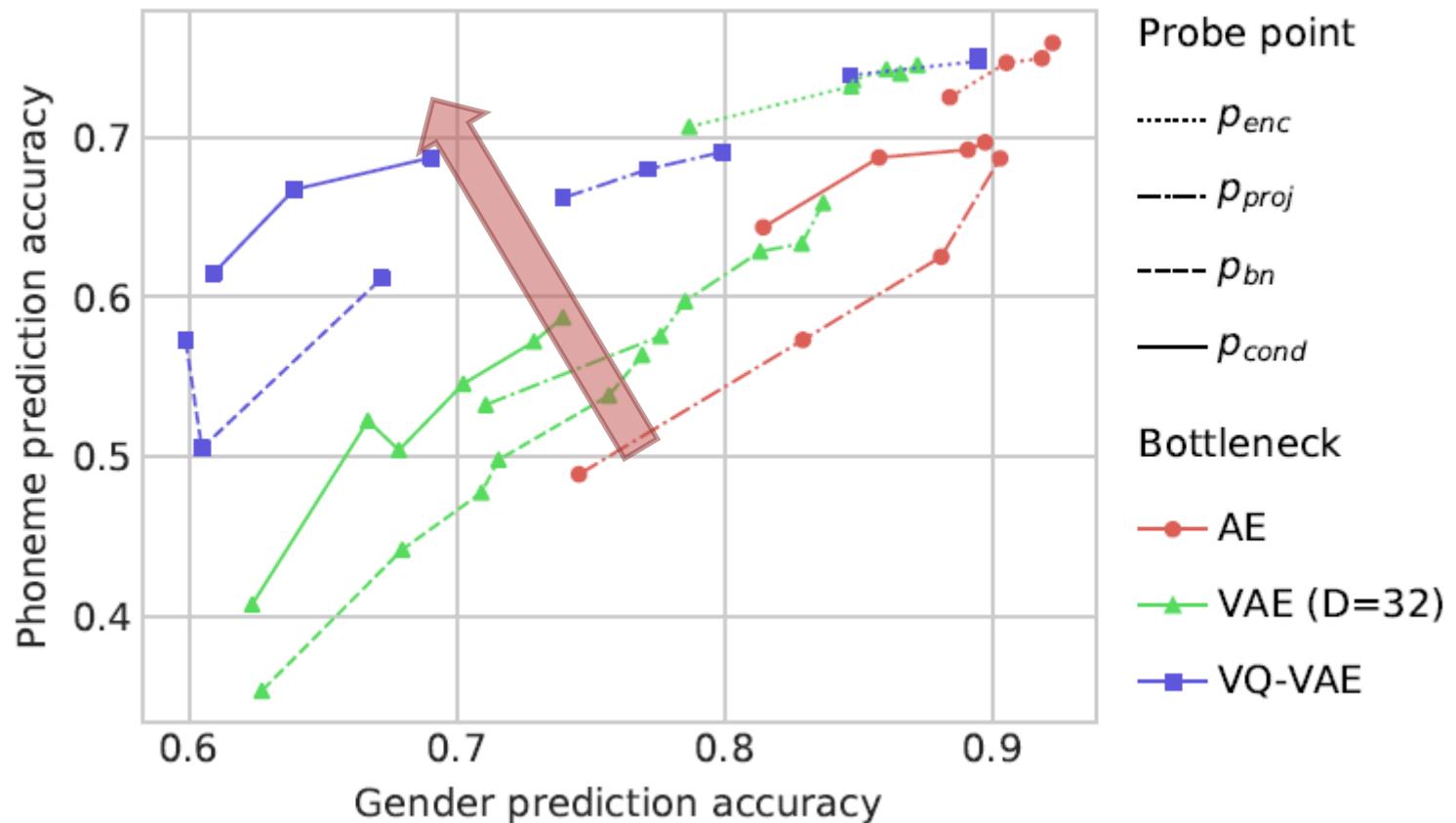
For each probing point, we have trained predictors for:

- 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_{cond})	5.6	5.5	5.5	7.3	7.5	7.5	11.2	10.7	10.8	8.1	8.0	8.0	11.0	10.8	11.1	12.2	11.7	11.9
Heck et al. [57]	6.9	6.2	6.0	9.7	8.7	8.4	8.8	7.9	7.8	10.1	8.7	8.5	13.6	11.7	11.3	8.8	7.4	7.3
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_{cond})	5.6	8.0
VQ-VAE (per lang, MFCC, p_{bn})	6.2	8.8
VQ-VAE (per lang, MFCC, p_{proj})	5.9	9.0
VQ-VAE (all lang, MFCC, p_{cond})	5.8	8.6
VQ-VAE (all lang, MFCC, p_{bn})	6.3	9.2
VQ-VAE (all lang, MFCC, p_{proj})	5.8	9.3
VQ-VAE (all lang, fbank, p_{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

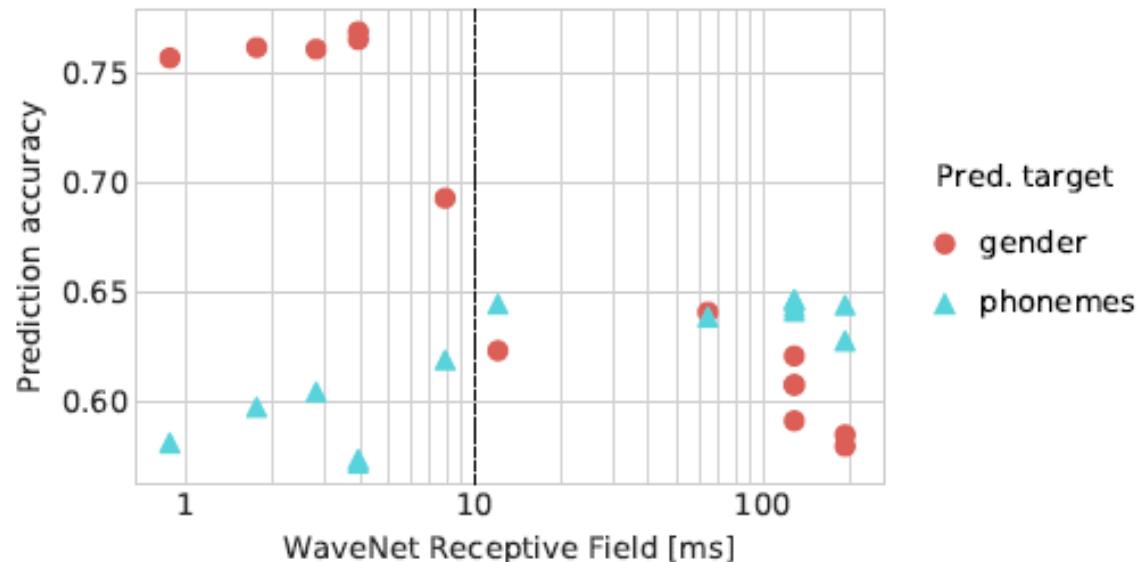
Mandarin	Within spkr.	Across spkr.
VQ-VAE (per lang, MFCC, p_{cond})	11.2	12.2
VQ-VAE (per lang, MFCC, p_{bn})	10.8	11.9
VQ-VAE (per lang, MFCC, p_{proj})	9.9	11.0
VQ-VAE (all lang, MFCC, p_{cond})	9.2	10.3
VQ-VAE (all lang, MFCC, p_{bn})	9.0	9.9
VQ-VAE (all lang, MFCC, p_{proj})	7.4	8.6
VQ-VAE (all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

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

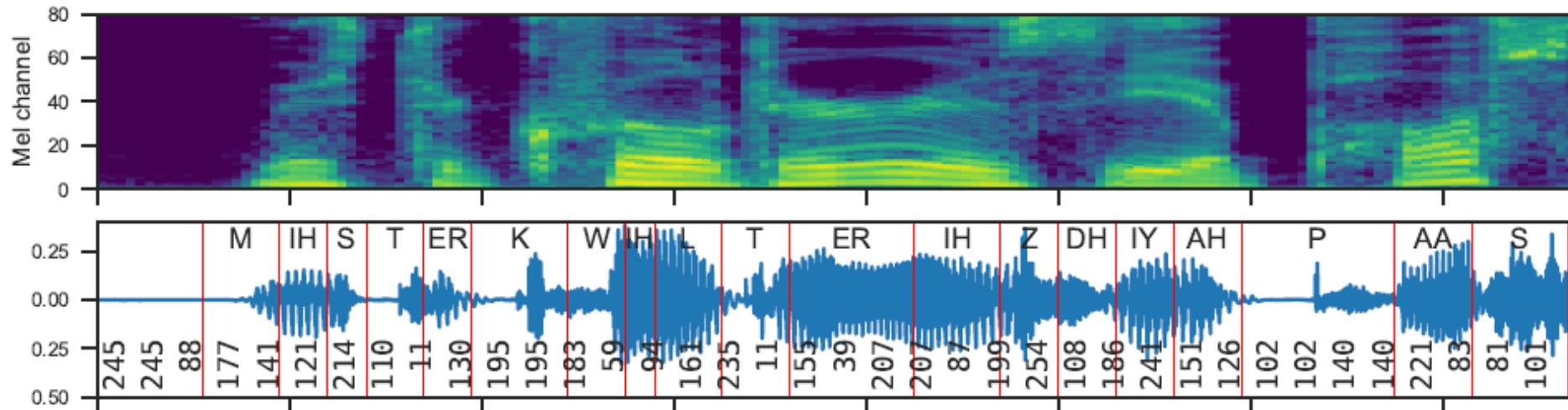
What impacts the representation?

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



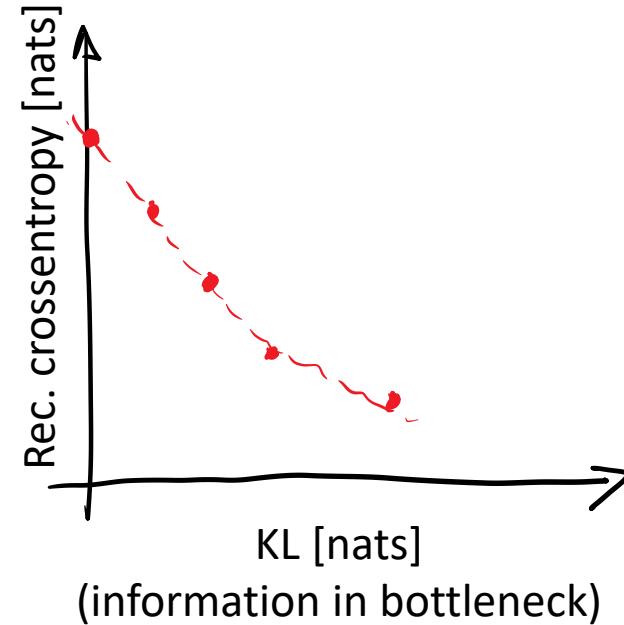
Experimental Questions #2



- 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

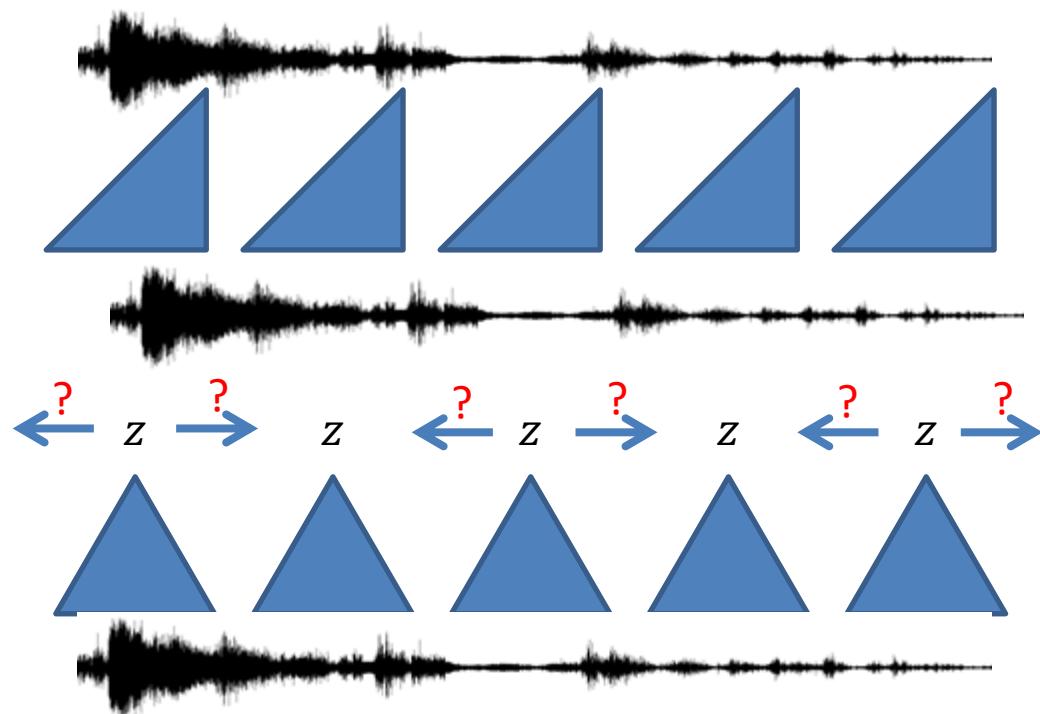
Promotes a constant encoding

Leads to collapse

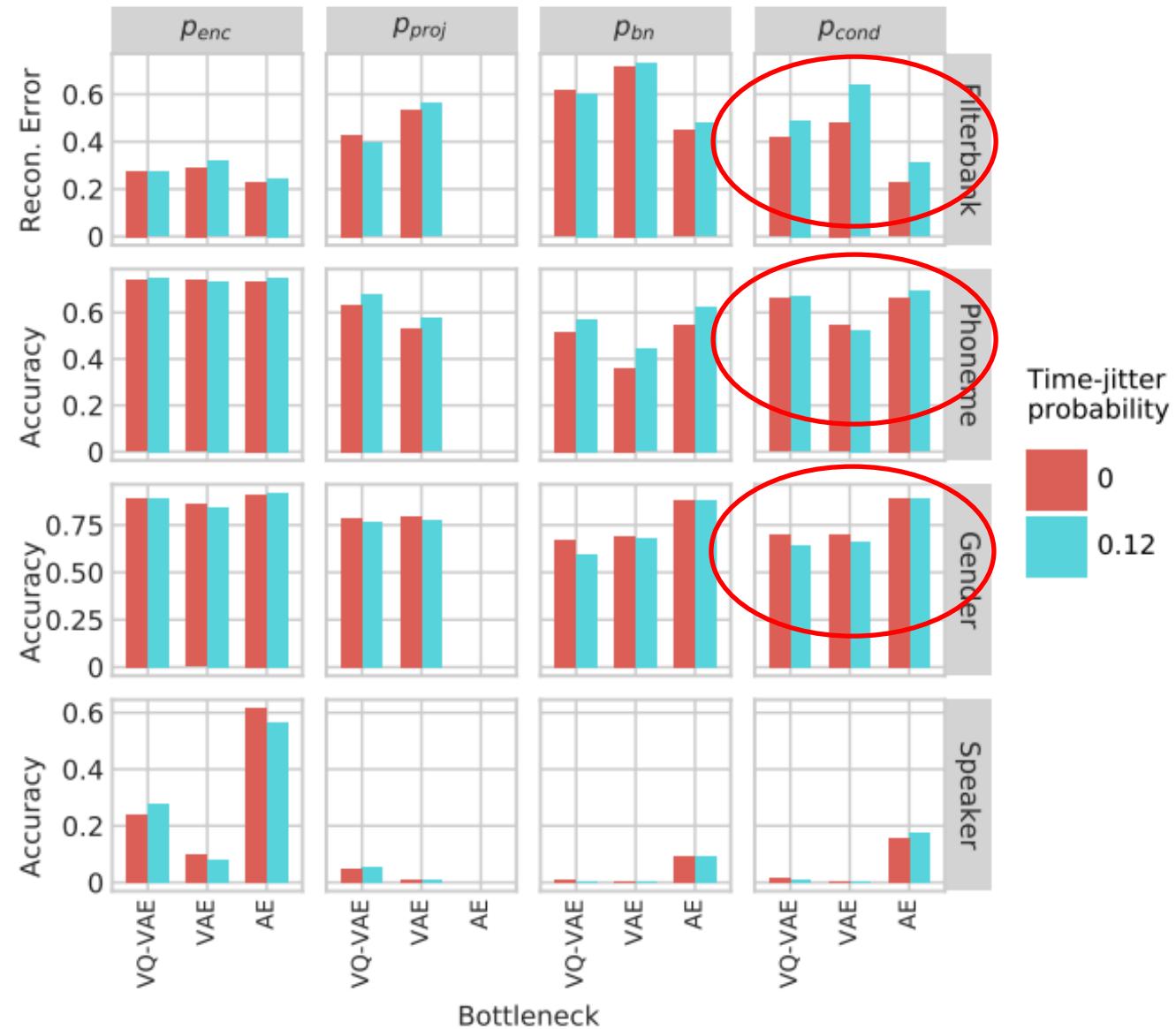


Randomized time jitter

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.



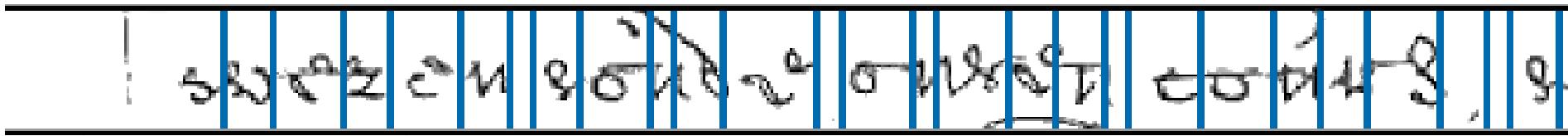
Randomized time jitter results



Segmenting auto-encoder

We extract a latent vector every 12 pixels

How to extract at every character?

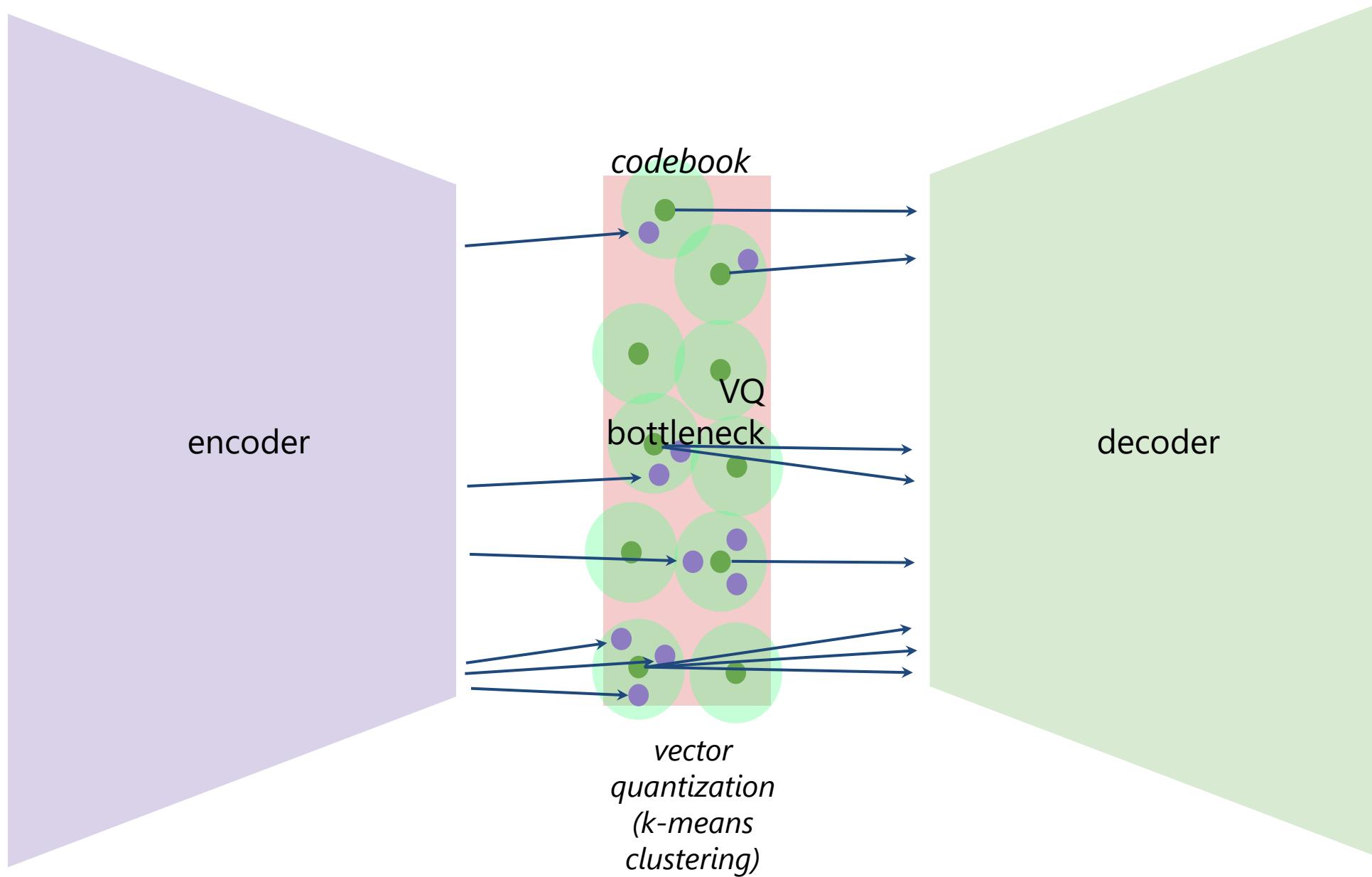


Enforce domain knowledge:

- 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:

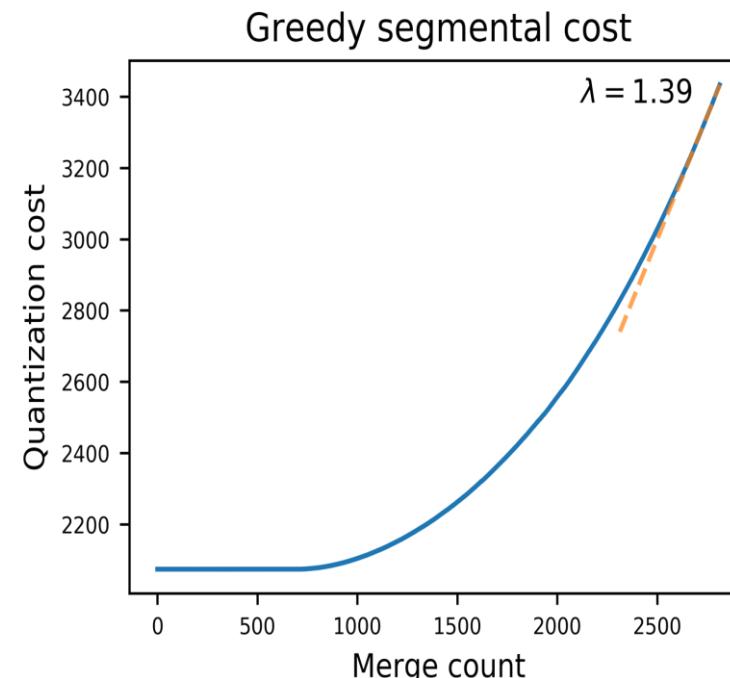
$$\min \sum_{t=1}^T \|x_t - q_t\| + \lambda \left(\sum_{t=1}^T [q_{t-1} \neq q_t] - K \right)$$

During training:

- . estimate λ

During inference:

- . use λ 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

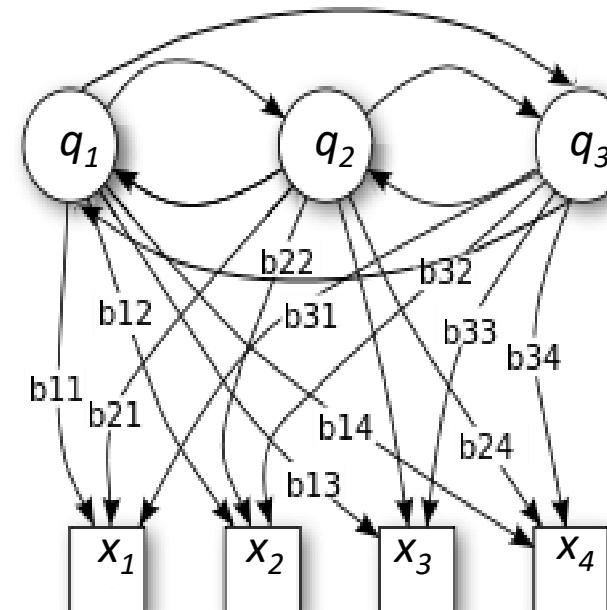


Image source: Wikipedia

Piecewise-constant results: Scribblelens

<i>Dataset</i>	<i>Model</i>	<i>Downstream CER</i>	<i>MLP Accuracy @Bottleneck</i>	<i>Token -> 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>	30%	65%	57%	0.23	0.36
<i>All scribes</i>	<i>Unsupervised base</i>	60%	42%	30%	0.13	0.10
	<i>Piecewise-const VQ bottleneck</i>	51%	49%	39%	0.18	0.20

THE FUTURE OF UNSUPERVISED: SELF-SUPERVISED LEARNING

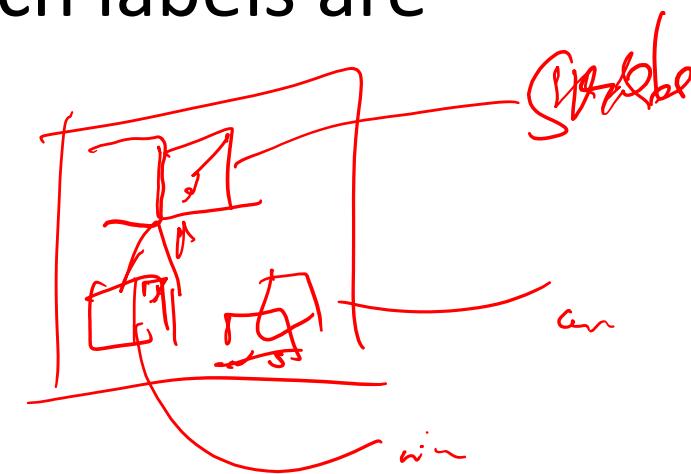
Core Idea

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

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

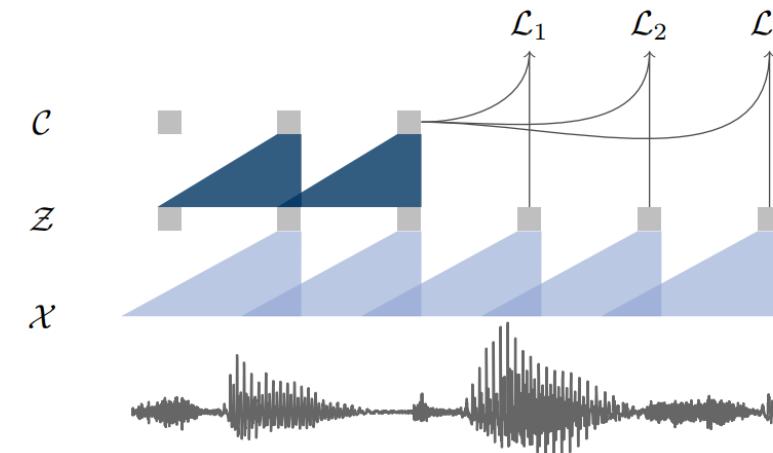
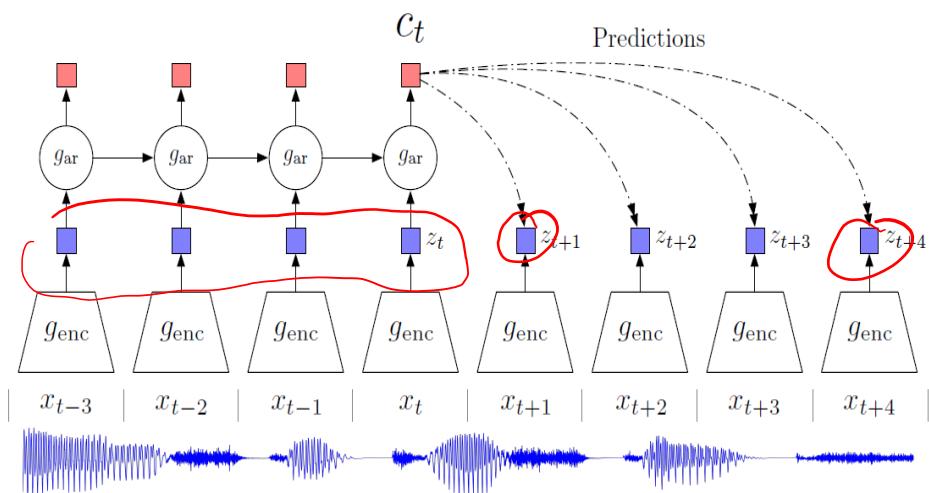
Or cut video into parts, predict time direction

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



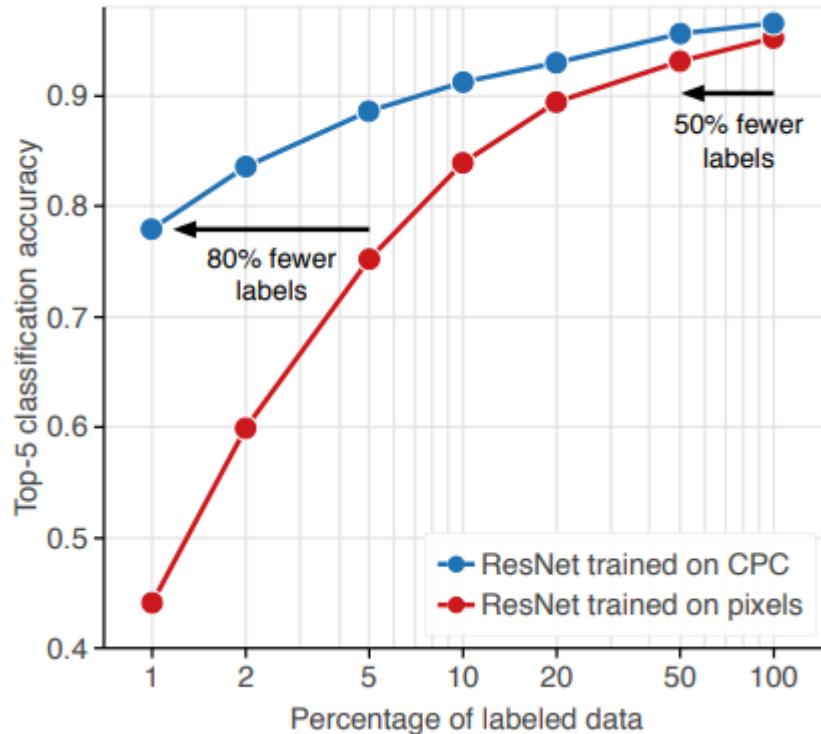
Decoder-less approach: CPC

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

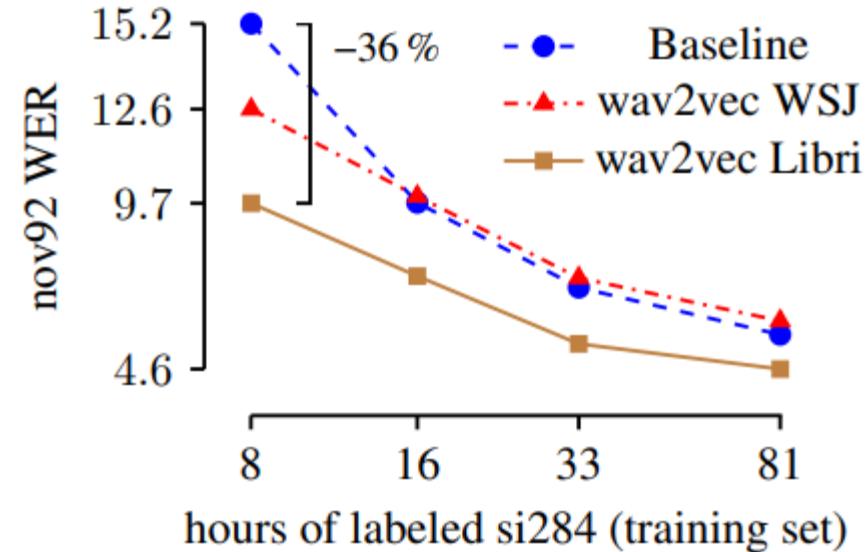


Oord et al. „Representation Learning with Contrastive Predictive Coding“
Schneider et al. „wav2vec: Unsupervised Pre-training for Speech Recognition“

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)

End of semester announcements

- Grades!
Make sure you have everything graded, if not ping me on discord -> I'll be grading A6 next week
- Project presentations
They are mandatory, and you know it since a month!