

# Introductory Overview Lecture

## The Deep Learning Revolution

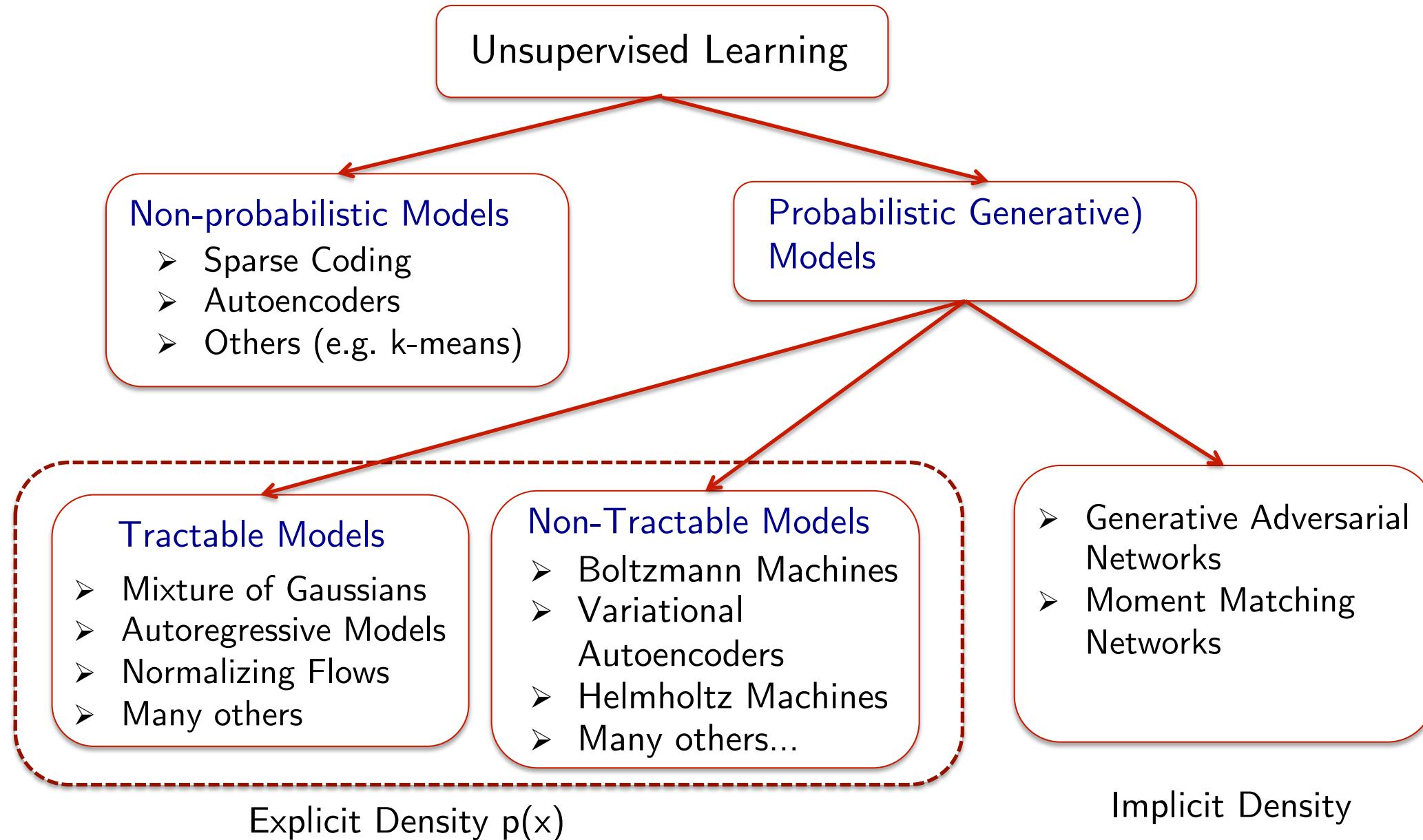
### Part III: Learning Deep Generative Models

Russ Salakhutdinov

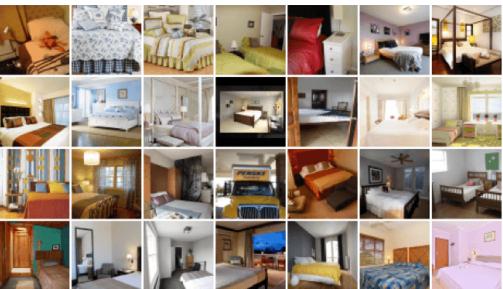
Machine Learning Department  
Carnegie Mellon University  
Canadian Institute for Advanced Research

Carnegie  
Mellon  
University





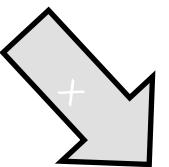
# Statistical Generative Models



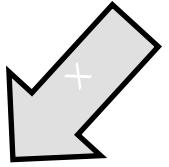
Data

+

Model family, loss function,  
optimization algorithm, etc.



Learning



Prior Knowledge

Image  $x$



A probability  
distribution  
 $p(x)$



probability  $p(x)$

Sampling from  $p(x)$  **generates** new images:



# Statistical Generative Models



Training Data(CelebA)

Model Samples (Karras et.al., 2018)

4 years of progression on Faces



2014

2015

2016

2017

Brundage et al., 2017

# Conditional Generation

- ▶ Conditional generative model  $P(\text{zebra images} | \text{horse images})$



- ▶ Style Transfer



Input Image



Monet



Van Gogh

Zhou et al., Cycle GAN 2017

# Conditional Generation

- ▶ Conditional generative model  $P(\text{zebra images} | \text{horse images})$



- ▶ Failure Case

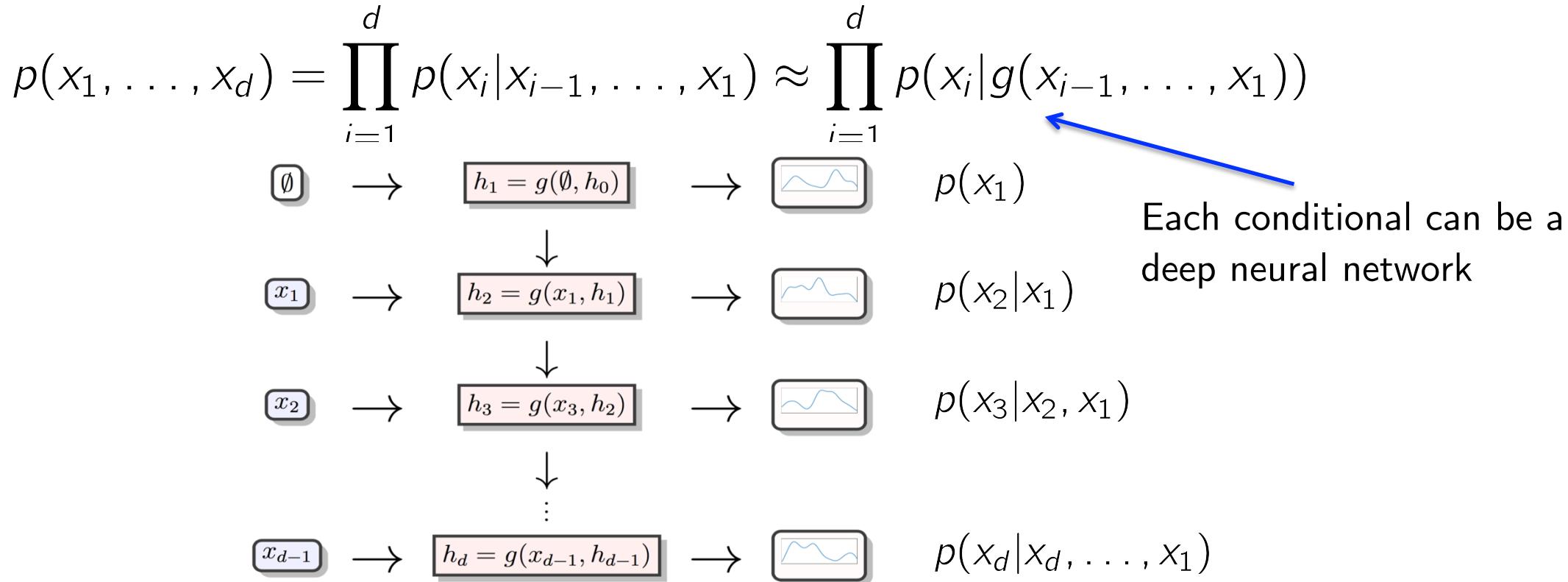


# Talk Roadmap

- ▶ Fully Observed Models
- ▶ Undirected Deep Generative Models
  - ▶ Restricted Boltzmann Machines (RBMs),
  - ▶ Deep Boltzmann Machines (DBMs)
- ▶ Directed Deep Generative Models
  - ▶ Variational Autoencoders (VAEs)
  - ▶ Normalizing Flows
- ▶ Generative Adversarial Networks (GANs)

# Fully Observed Models

- Density Estimation by Autoregression

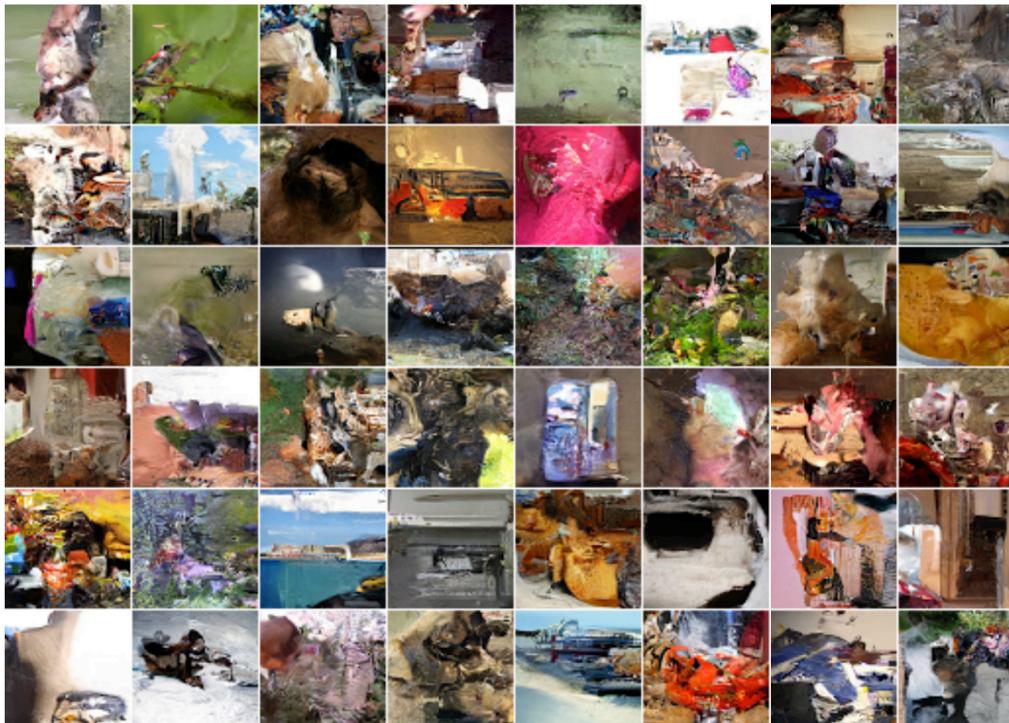


- Ordering of variables is crucial

NADE (Uria 2013), MADE (Germain 2017), MAF (Papamakarios 2017), PixelCNN (van den Oord, et al, 2016)

# Fully Observed Models

- ## ► Density Estimation by Autoregression

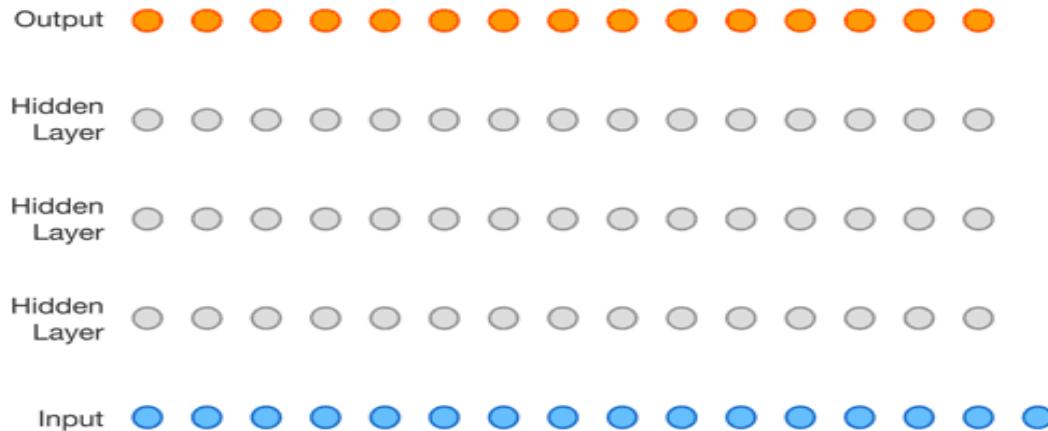


PixelCNN (van den Oord, et al, 2016)

NADE (Uria 2013), MADE (Germain 2017), MAF (Papamakarios 2017), PixelCNN (van den Oord, et al, 2016)

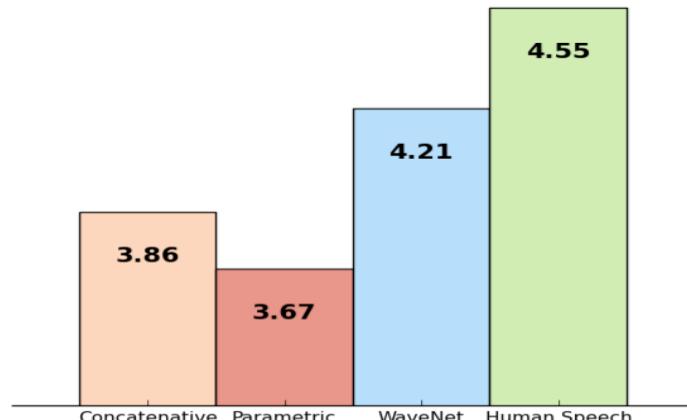
# WaveNet

## ► Generative Model of Speech Signals



1 Second

Quality: Mean Opinion Scores

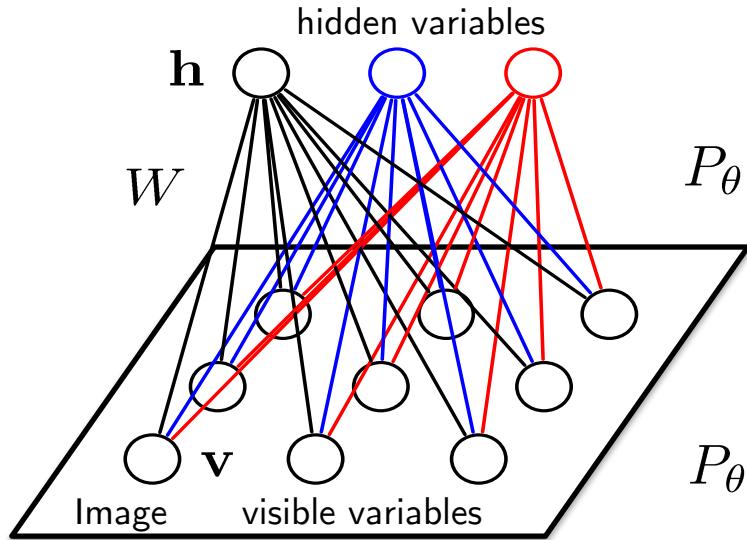


van den Oord et al, 2016

# Talk Roadmap

- ▶ Fully Observed Models
- ▶ Undirected Deep Generative Models
  - ▶ Restricted Boltzmann Machines (RBMs),
  - ▶ Deep Boltzmann Machines (DBMs)
- ▶ Directed Deep Generative Models
  - ▶ Variational Autoencoders (VAEs)
  - ▶ Normalizing Flows
- ▶ Generative Adversarial Networks (GANs)

# Restricted Boltzmann Machines



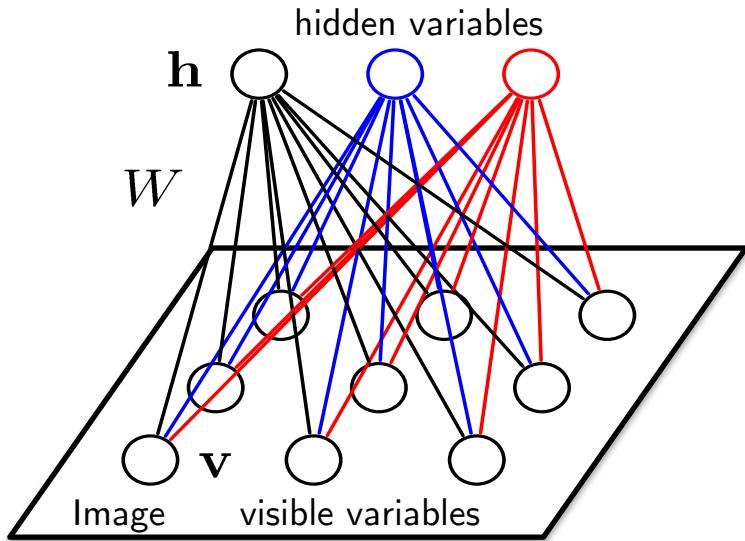
$$P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \exp \left( \sum_{i=1}^D \sum_{j=1}^F W_{ij} v_i h_j + \sum_{i=1}^D v_i b_i + \sum_{j=1}^F h_j a_j \right)$$

$$\theta = \{W, a, b\}$$

$$P_{\theta}(\mathbf{v}|\mathbf{h}) = \prod_{i=1}^D P_{\theta}(v_i|\mathbf{h}) = \prod_{i=1}^D \frac{1}{1 + \exp(-\sum_{j=1}^F W_{ij} v_i h_j - b_i)}$$

- ▶ RBM is a Markov Random Field with
  - ▶ Stochastic binary visible variables  $\mathbf{v} \in \{0, 1\}^D$ .
  - ▶ Stochastic binary hidden variables  $\mathbf{h} \in \{0, 1\}^F$ .
  - ▶ Bipartite connections

# Maximum Likelihood Learning



$$P_\theta(\mathbf{v}) = \frac{P^*(\mathbf{v})}{\mathcal{Z}(\theta)} = \frac{1}{\mathcal{Z}(\theta)} \sum_{\mathbf{h}} \exp \left[ \mathbf{v}^\top W \mathbf{h} + \mathbf{a}^\top \mathbf{h} + \mathbf{b}^\top \mathbf{v} \right]$$

- Maximize log-likelihood objective:

$$L(\theta) = \frac{1}{N} \sum_{n=1}^N \log P_\theta(\mathbf{v}^{(n)})$$

- Derivative of the log-likelihood:

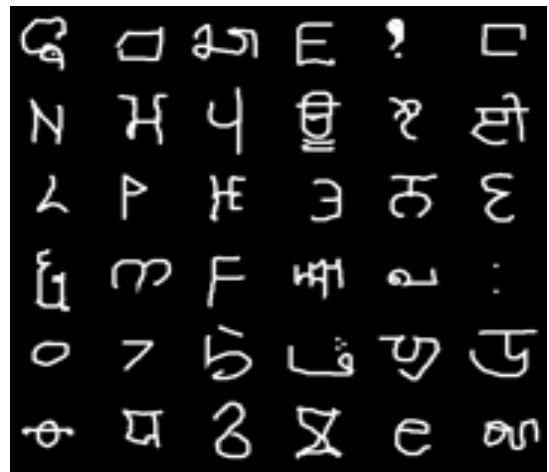
$$\begin{aligned} \frac{\partial L(\theta)}{\partial W_{ij}} &= \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial W_{ij}} \log \left( \sum_{\mathbf{h}} \exp [\mathbf{v}^{(n)\top} W \mathbf{h} + \mathbf{a}^\top \mathbf{h} + \mathbf{b}^\top \mathbf{v}^{(n)}] \right) - \frac{\partial}{\partial W_{ij}} \log \mathcal{Z}(\theta) \\ &= \mathbf{E}_{P_{data}} [v_i h_j] - \underbrace{\mathbf{E}_{P_\theta} [v_i h_j]}_{\text{Difficult to compute: exponentially many configurations}} \end{aligned}$$

$$P_{data}(\mathbf{v}, \mathbf{h}; \theta) = P(\mathbf{h}|\mathbf{v}; \theta)P_{data}(\mathbf{v})$$

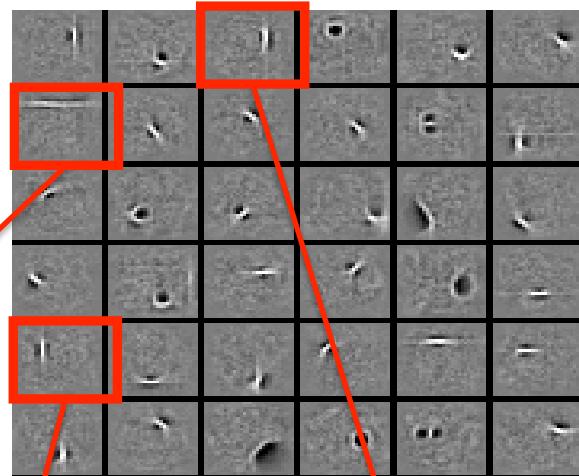
$$P_{data}(\mathbf{v}) = \frac{1}{N} \sum_n \delta(\mathbf{v} - \mathbf{v}^{(n)})$$

# Learning Features

Observed Data  
Subset of 25,000 characters



Learned W: “edges”  
Subset of 1000 features



New Image:

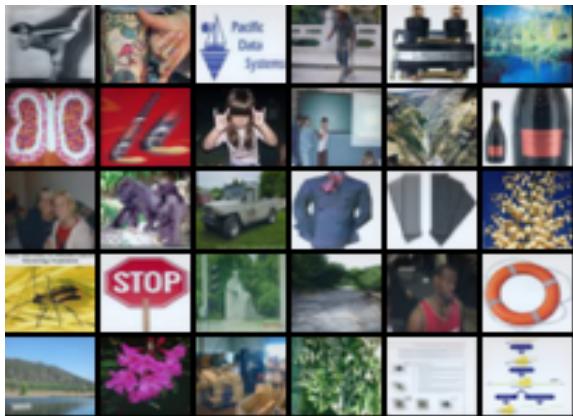
$$\text{Digit} = \sigma\left(0.99 \times \begin{matrix} p(h_7 = 1|v) \\ \downarrow \\ \text{Feature Image} \end{matrix} + 0.97 \times \begin{matrix} p(h_{29} = 1|v) \\ \downarrow \\ \text{Feature Image} \end{matrix} + 0.82 \times \dots\right)$$

$$\sigma(x) = \frac{1}{1+\exp(-x)}$$

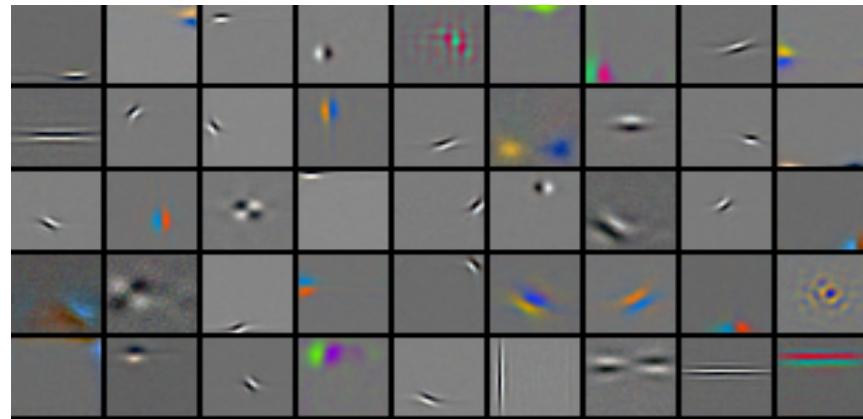
Logistic Function: Suitable for modeling binary images

# RBMs for Real-valued & Count Data

4 million **unlabelled** images



Learned features (out of 10,000)



**REUTERS**   
**AP** Associated Press

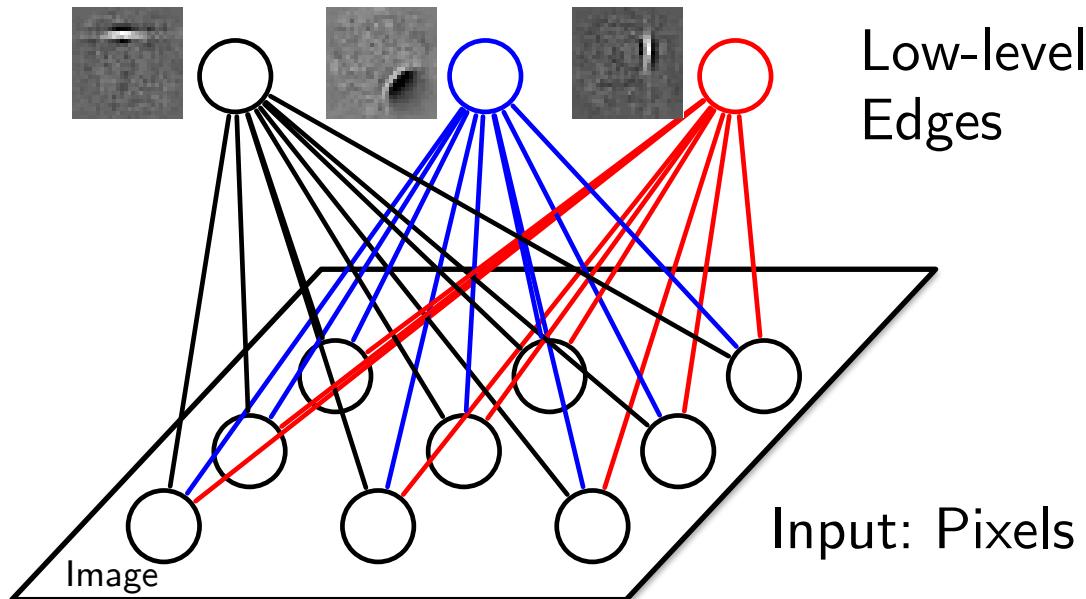
Reuters dataset:  
804,414 **unlabeled**  
newswire stories  
Bag-of-Words



Learned features: ``topics''

|         |           |          |         |        |
|---------|-----------|----------|---------|--------|
| russian | clinton   | computer | trade   | stock  |
| russia  | house     | system   | country | wall   |
| moscow  | president | product  | import  | street |
| yeltsin | bill      | software | world   | point  |
| soviet  | congress  | develop  | economy | dow    |

# Deep Boltzmann Machines



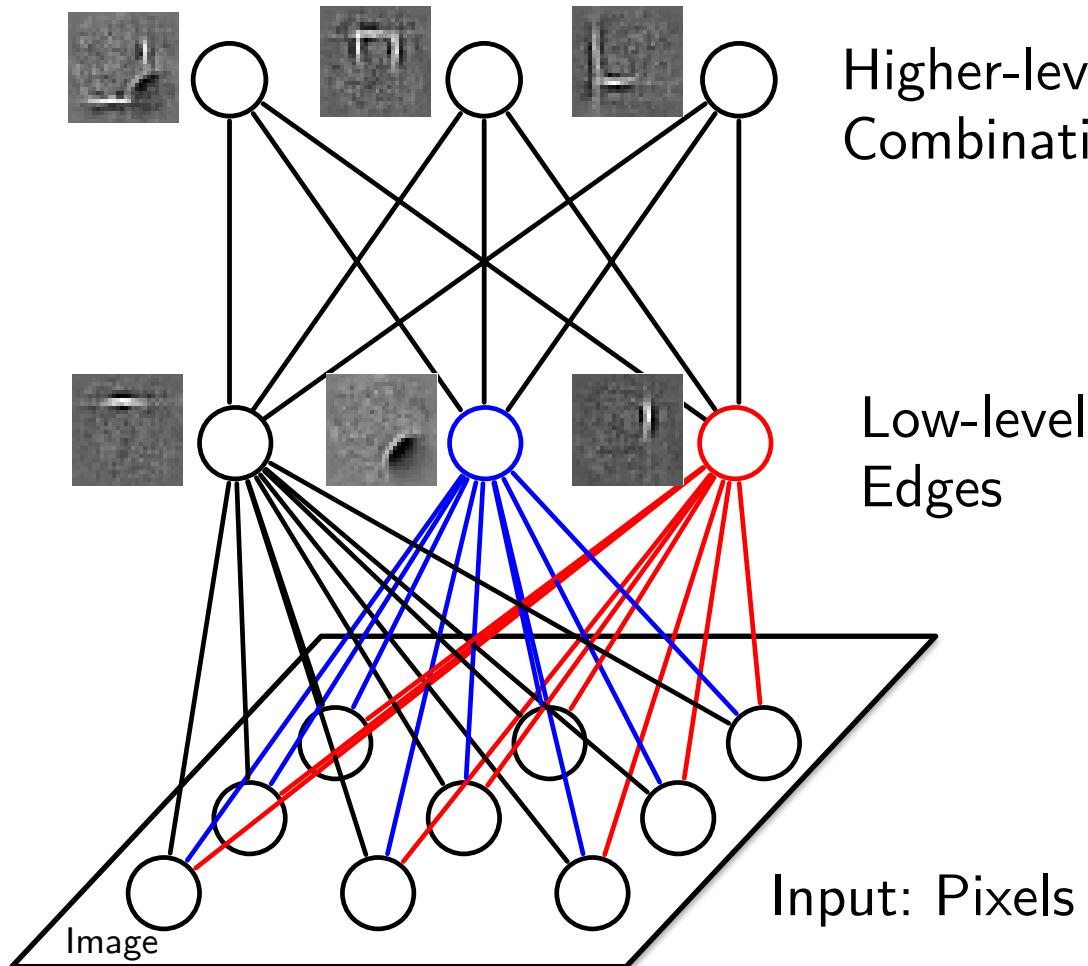
Low-level features:  
Edges

Built from **unlabeled** inputs.

(Salakhutdinov 2008, Salakhutdinov & Hinton 2009)

# Deep Boltzmann Machines

Learn simpler representations,  
then compose more complex ones

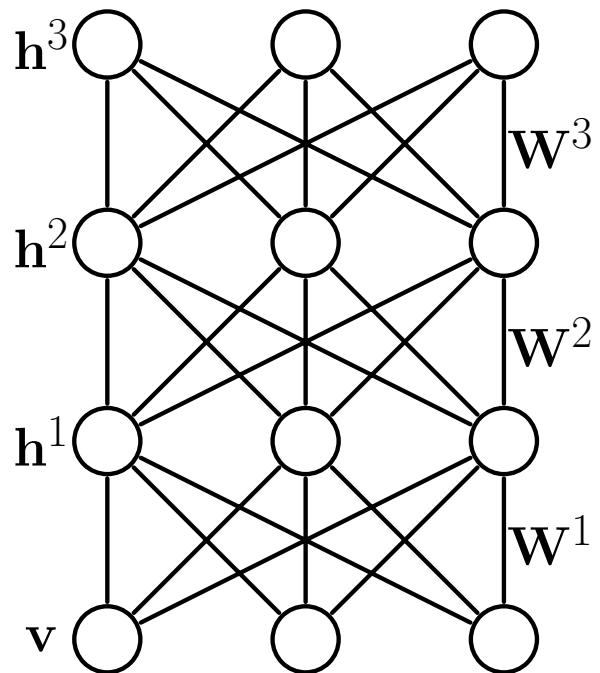


Built from **unlabeled** inputs.

(Salakhutdinov 2008, Salakhutdinov & Hinton 2009)

# Model Formation

$$P_{\theta}(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \frac{1}{\mathcal{Z}(\theta)} \exp \left[ \underbrace{\mathbf{v}^{\top} W^{(1)} \mathbf{h}^{(1)}}_{\text{Same as RBMs}} + \underbrace{\mathbf{h}^{(1)\top} W^{(2)} \mathbf{h}^{(2)}}_{\text{Same as RBMs}} + \underbrace{\mathbf{h}^{(2)\top} W^{(3)} \mathbf{h}^{(3)}}_{\text{Same as RBMs}} \right]$$



$$\theta = \{W^1, W^2, W^3\} \text{ model parameters}$$

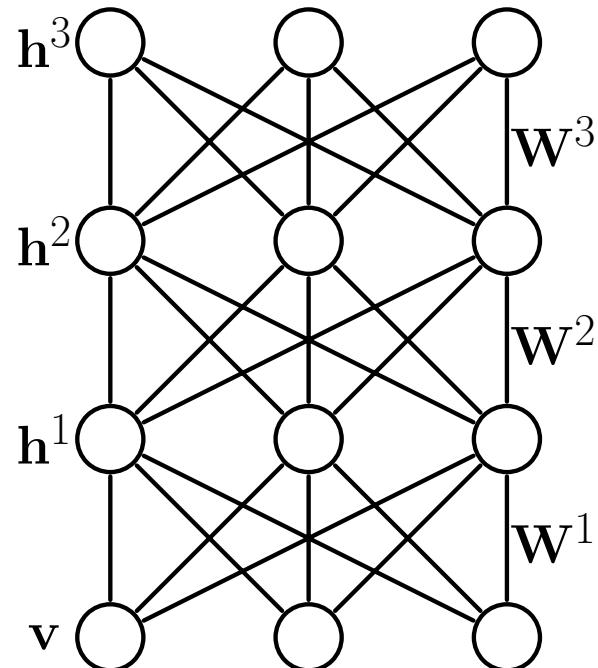
- ▶ Dependencies between hidden variables
- ▶ All connections are undirected
- ▶ Maximum Likelihood Learning:

$$\frac{\partial \log P_{\theta}(\mathbf{v})}{\partial W^1} = \mathbb{E}_{P_{data}} [\mathbf{v} \mathbf{h}^{1\top}] - \mathbb{E}_{P_{\theta}} [\mathbf{v} \mathbf{h}^{1\top}]$$

- ▶ Both expectations are intractable

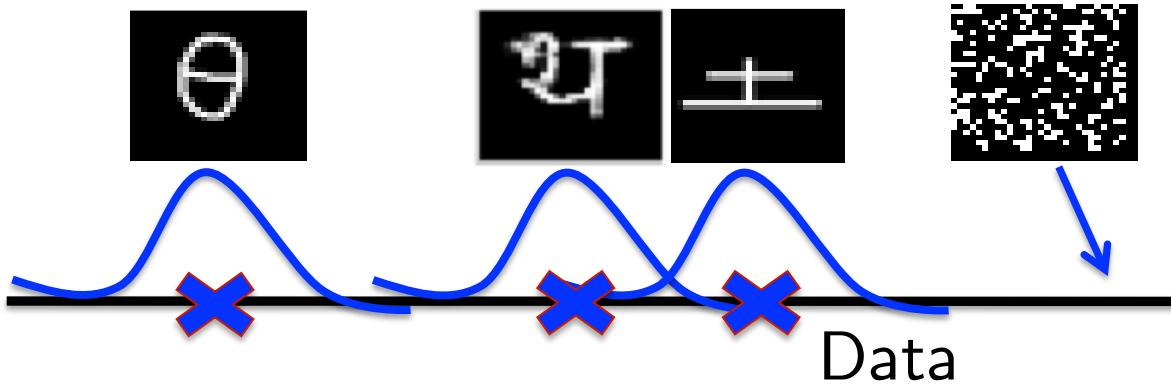
# Approximate Learning

$$P_{\theta}(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \frac{1}{Z(\theta)} \exp \left[ \mathbf{v}^T W^{(1)} \mathbf{h}^{(1)} + \mathbf{h}^{(1)T} W^{(2)} \mathbf{h}^{(2)} + \mathbf{h}^{(2)T} W^{(3)} \mathbf{h}^{(3)} \right]$$



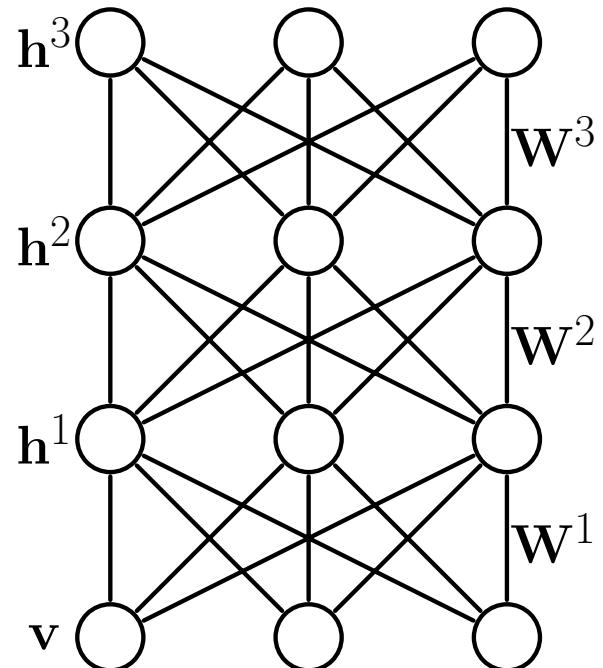
► Maximum Likelihood Learning:

$$\frac{\partial \log P_{\theta}(\mathbf{v})}{\partial W^1} = \mathbb{E}_{P_{data}} [\mathbf{v} \mathbf{h}^1]^T - \mathbb{E}_{P_{\theta}} [\mathbf{v} \mathbf{h}^1]^T$$



# Approximate Learning

$$P_{\theta}(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \frac{1}{\mathcal{Z}(\theta)} \exp \left[ \mathbf{v}^{\top} W^{(1)} \mathbf{h}^{(1)} + \mathbf{h}^{(1)\top} W^{(2)} \mathbf{h}^{(2)} + \mathbf{h}^{(2)\top} W^{(3)} \mathbf{h}^{(3)} \right]$$



► Maximum Likelihood Learning:

$$\frac{\partial \log P_{\theta}(\mathbf{v})}{\partial W^1} = \mathbb{E}_{P_{data}} [\mathbf{v} \mathbf{h}^{1\top}] - \mathbb{E}_{P_{\theta}} [\mathbf{v} \mathbf{h}^{1\top}]$$

Variational  
Inference

Stochastic Approximation  
(MCMC-based)

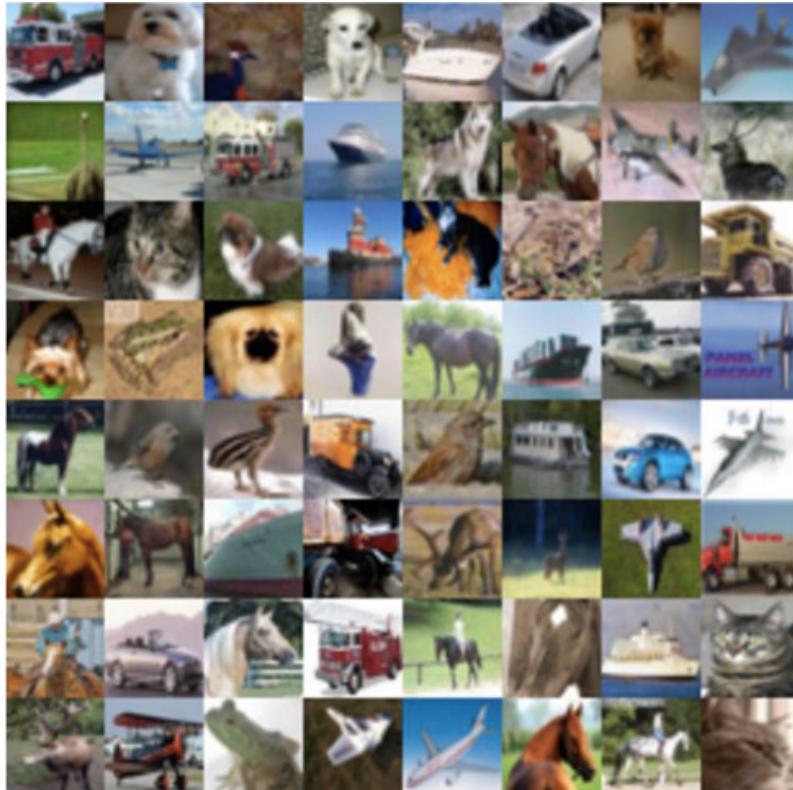
# Good Generative Model?

- ▶ Handwritten Characters

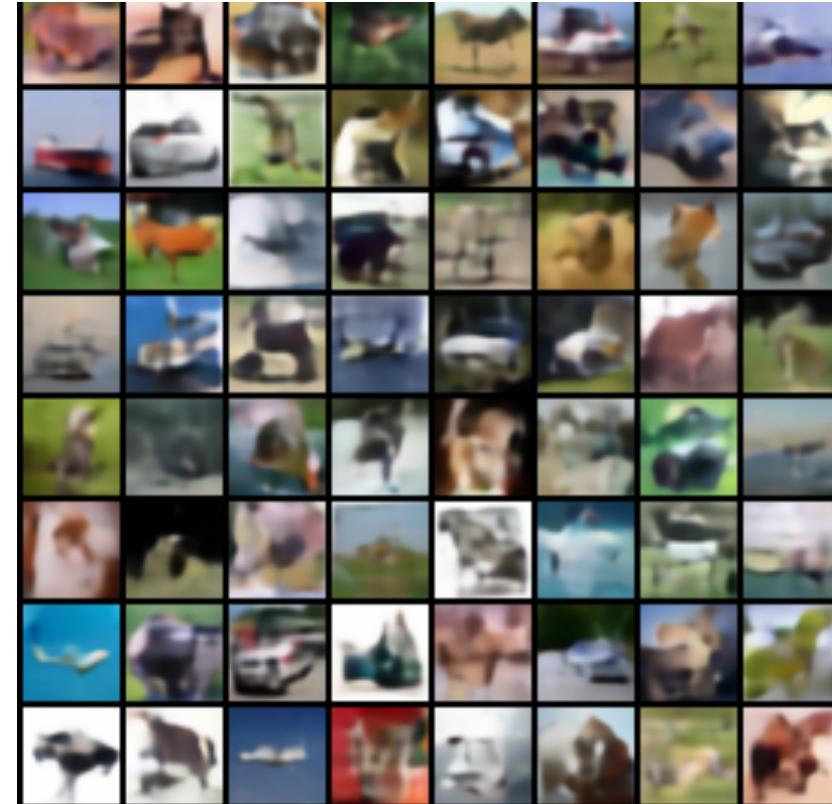


# Good Generative Model?

- CIFAR Dataset



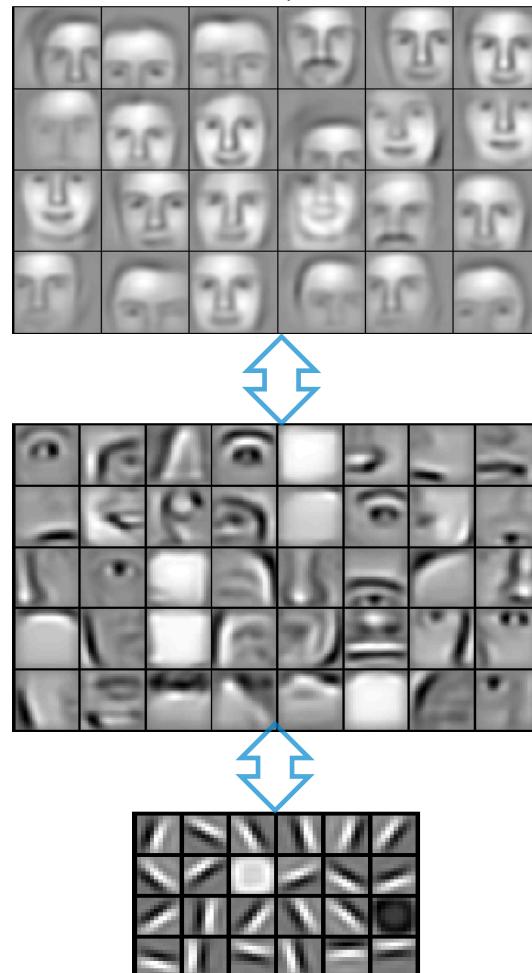
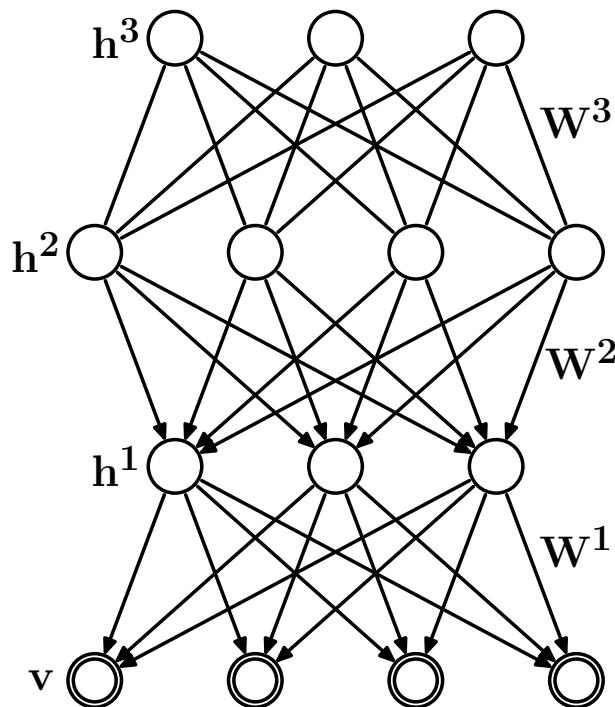
Training



Samples

# Learning Part-Based Representations

Deep Belief Network

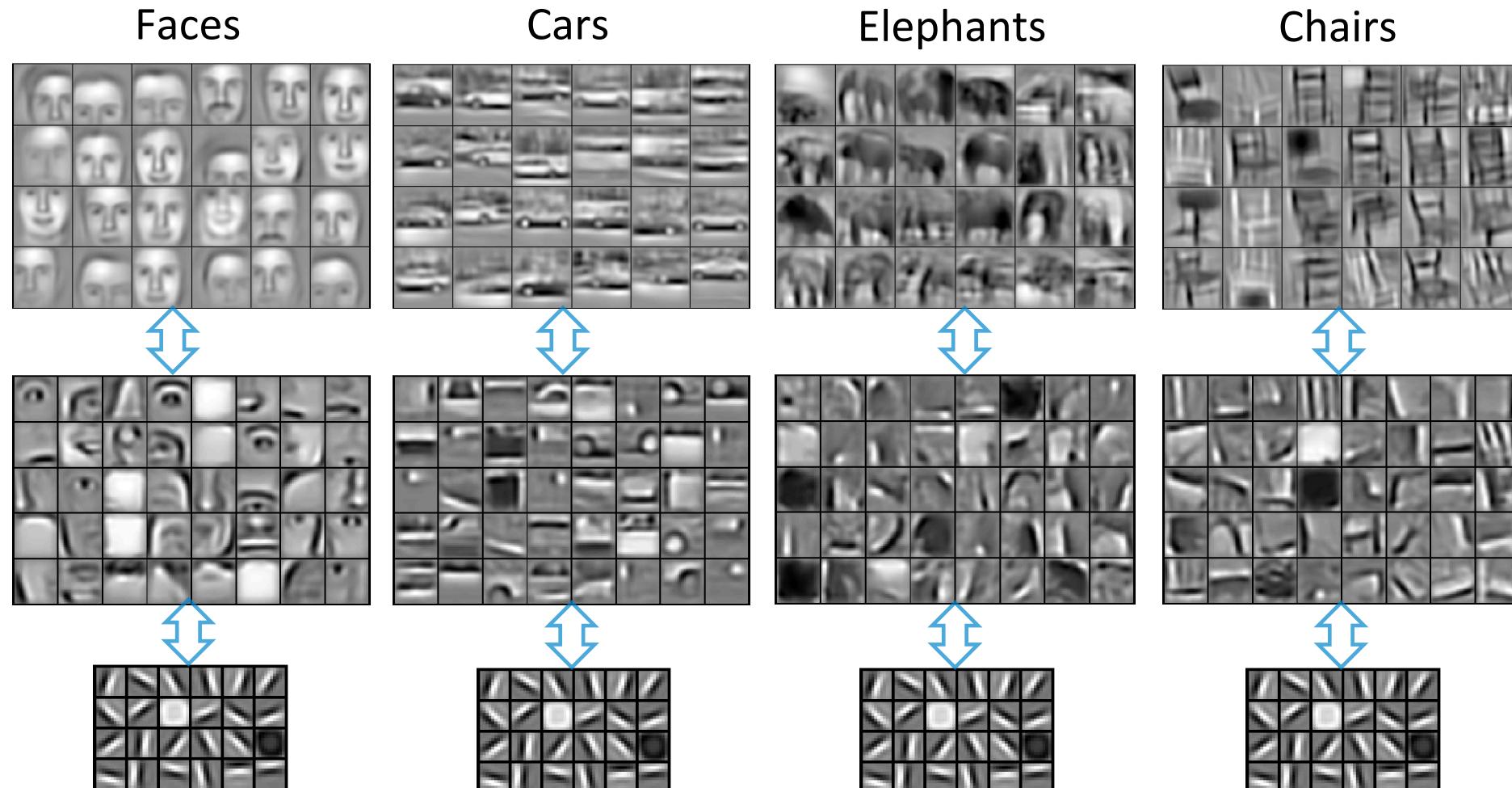


Groups of parts

Object Parts

Trained on face images.

# Learning Part-Based Representations

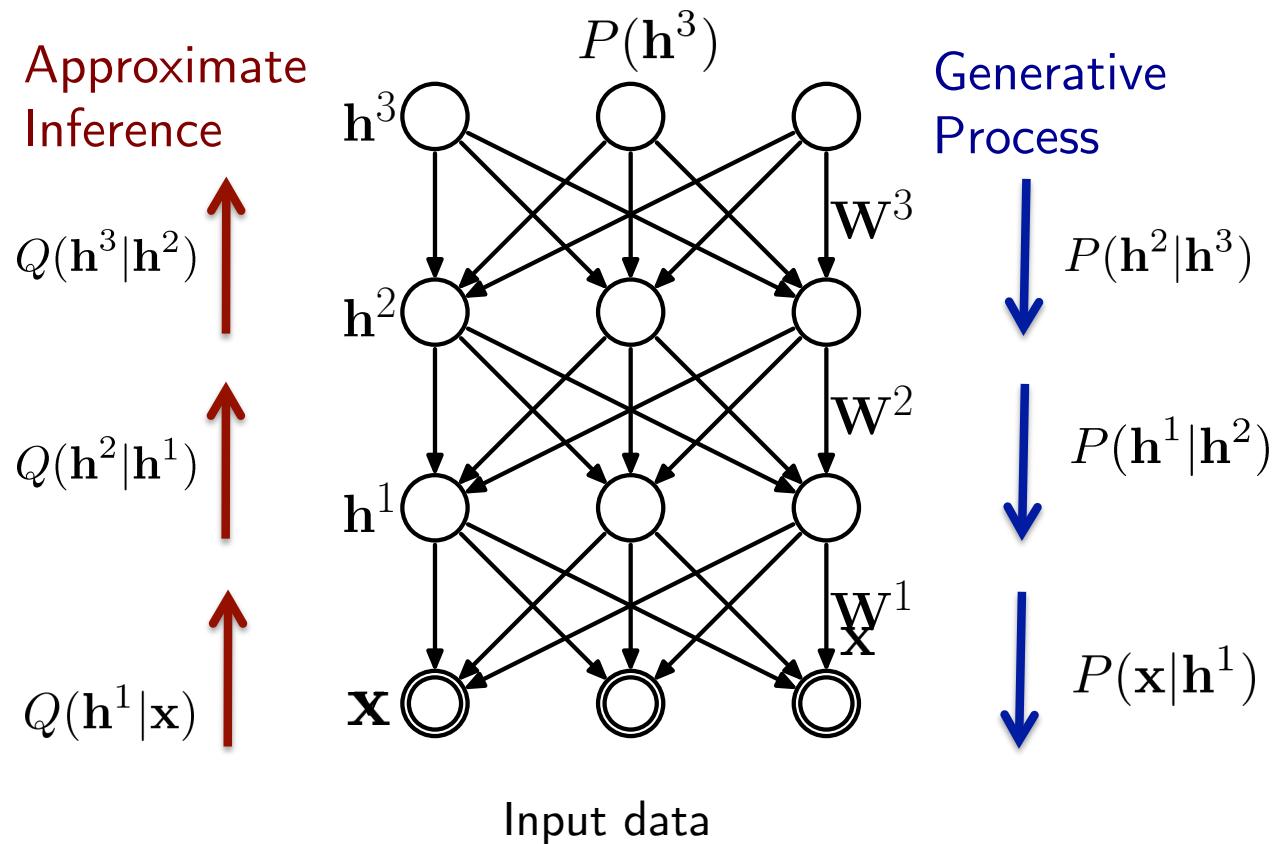


# Talk Roadmap

- ▶ Fully Observed Models
- ▶ Undirected Deep Generative Models
  - ▶ Restricted Boltzmann Machines (RBMs),
  - ▶ Deep Boltzmann Machines (DBMs)
- ▶ Directed Deep Generative Models
  - ▶ Variational Autoencoders (VAEs)
  - ▶ Normalizing Flows
- ▶ Generative Adversarial Networks (GANs)

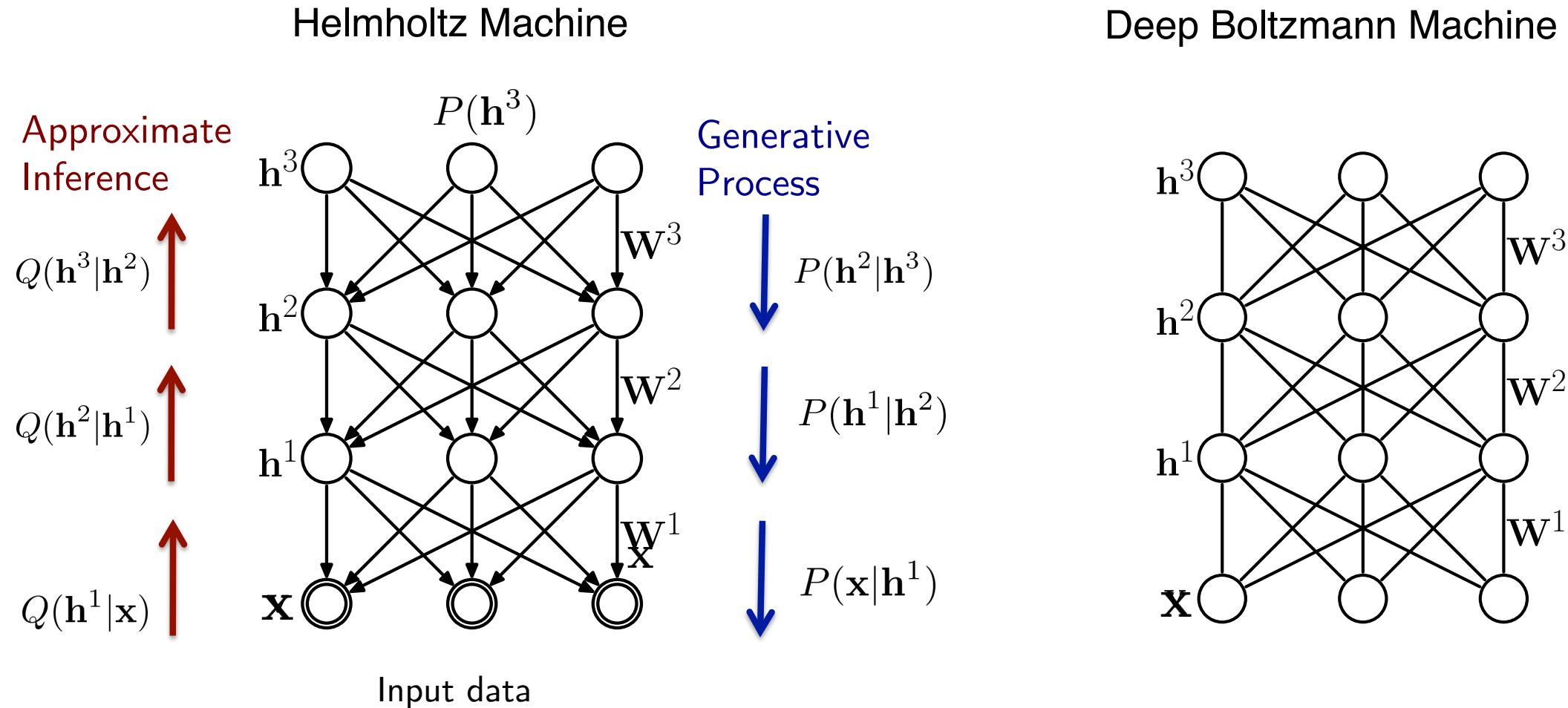
# Helmholtz Machines

- Hinton, G. E., Dayan, P., Frey, B. J. and Neal, R., Science 1995



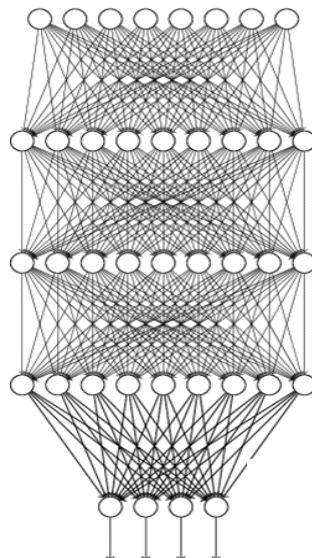
- Kingma & Welling, 2014
- Rezende, Mohamed, Daan, 2014
- Mnih & Gregor, 2014
- Bornschein & Bengio, 2015
- Tang & Salakhutdinov, 2013

# Helmholtz Machines

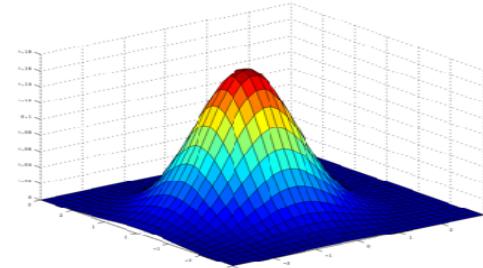


# Deep Directed Generative Models

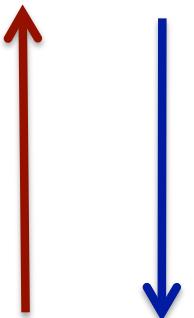
White  
Noise



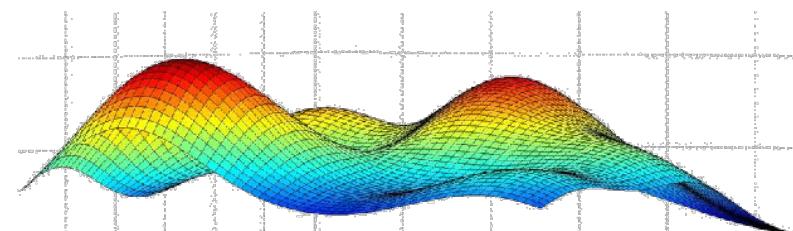
Code Z



- ▶ Recognition
- ▶ Bottom-up
- ▶  $Q(z|x)$

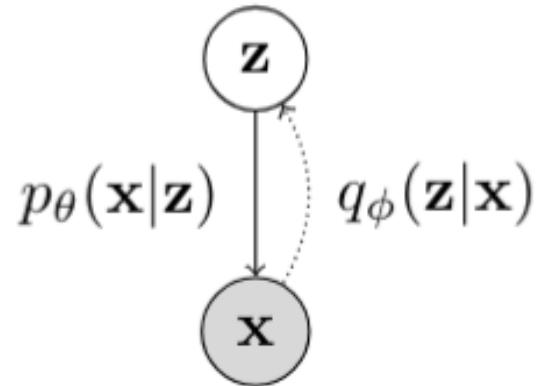


$D_{real}$



- ▶ Generative
- ▶ Top-Down
- ▶  $P(x|z)$

## ► Latent Variable Models

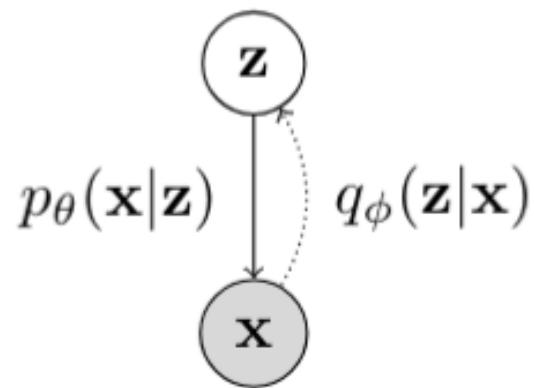


$$\log p_\theta(\mathbf{x}) = \log \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

- ▶ Conditional distributions are parameterized by deep neural networks

# Directed Deep Generative Models

- Directed Latent Variable Models with Inference Network



- Maximum log-likelihood objective

$$\max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x})$$

- Marginal log-likelihood is **intractable**:

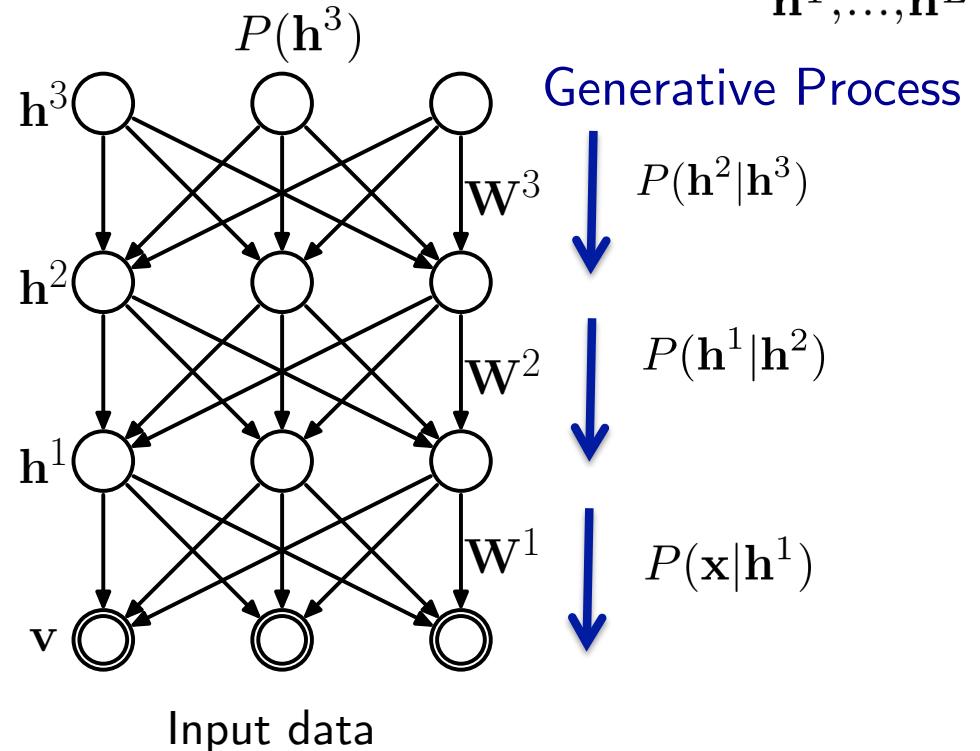
$$\log p_{\theta}(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

- **Key idea:** Approximate true posterior  $p(z|x)$  with a simple, tractable distribution  $q(z|x)$  (inference/recognition network).

# Variational Autoencoders (VAEs)

- The VAE defines a generative process in terms of ancestral sampling through a cascade of hidden stochastic layers:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{\mathbf{h}^1, \dots, \mathbf{h}^L} p(\mathbf{h}^L|\boldsymbol{\theta})p(\mathbf{h}^{L-1}|\mathbf{h}^L, \boldsymbol{\theta}) \cdots p(\mathbf{x}|\mathbf{h}^1, \boldsymbol{\theta})$$

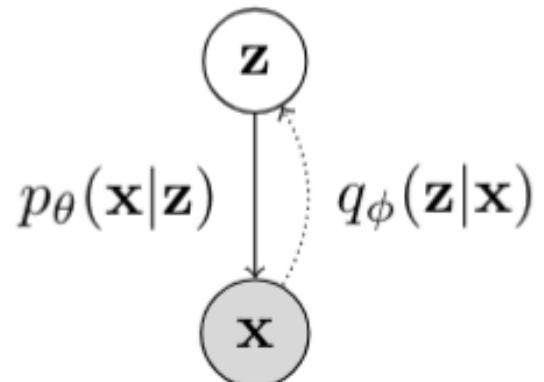
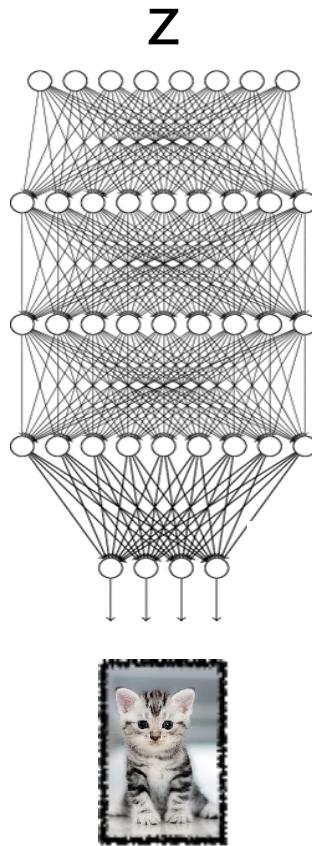


Each conditional term denotes a nonlinear relationship

- $L$  is the number of **stochastic** layers
- Sampling and probability evaluation is tractable for each  $p(\mathbf{h}^\ell|\mathbf{h}^{\ell+1})$

# Variational Autoencoders (VAEs)

- Single stochastic (Gaussian) layer, followed by many deterministic layers



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

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu(\mathbf{z}, \theta), \Sigma(\mathbf{z}, \theta))$$

Deep neural network parameterized by  $\theta$ .  
(Can use different noise models)

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu(\mathbf{x}, \phi), \Sigma(\mathbf{x}, \phi))$$

Deep neural network parameterized by  $\phi$ .

# Variational Bound

- VAE is trained to maximize the **variational lower bound**:

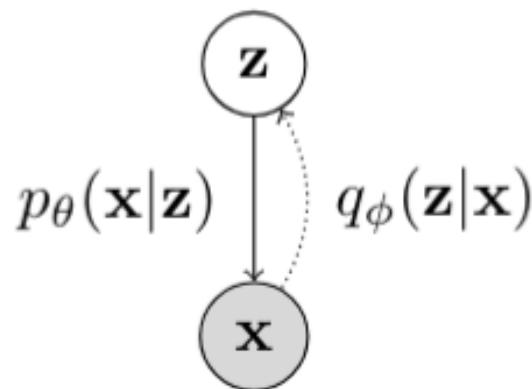
$$\begin{aligned}
 \log p_\theta(\mathbf{x}) &= \log \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \log \int q_\phi(\mathbf{z}|\mathbf{x}) \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\
 &= \log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\
 &\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \left[ \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\
 &= \log p_\theta(\mathbf{x}) - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x})) = \mathcal{L}(\mathbf{x})
 \end{aligned}$$

Tightness Condition:  
 $q_\phi(\mathbf{z}|\mathbf{x}) = p_\theta(\mathbf{z}|\mathbf{x})$

- Trading off the data log-likelihood and the KL divergence from the true posterior
- Hard to optimize the variational bound with respect to the q recognition network (high variance)
- Key idea of Kingma and Welling is to use reparameterization trick

# Reparameterization

- ▶ Assume that the recognition distribution is Gaussian:



$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu(\mathbf{x}, \phi), \Sigma(\mathbf{x}, \phi))$$

- ▶ Alternatively, we can express this in term of auxiliary variable:

$$\mathbf{z}(\epsilon, \mathbf{x}, \phi) = \Sigma(\mathbf{x}, \phi)^{1/2}\epsilon + \mu(\mathbf{x}, \phi), \quad \epsilon \sim \mathcal{N}(0, I)$$

- ▶ The recognition distribution can be expressed as a deterministic mapping

$$\underbrace{\mathbf{z}(\epsilon, \mathbf{x}, \phi)}$$

- ▶ Distribution of  $\epsilon$  does not depend on  $\phi$

Deterministic Encoder

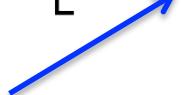
# Computing Gradients

- The gradients of the variational bound w.r.t the recognition (similar w.r.t the generative) parameters:

$$\begin{aligned}
 \nabla_{\phi} \mathcal{L}(\mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\
 &= \nabla_{\phi} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z}(\epsilon, \mathbf{x}, \phi))}{q_{\phi}(\mathbf{z}(\epsilon, \mathbf{x}, \phi)|\mathbf{x})} \right] \\
 &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ \nabla_{\phi} \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z}(\epsilon, \mathbf{x}, \phi))}{q_{\phi}(\mathbf{z}(\epsilon, \mathbf{x}, \phi)|\mathbf{x})} \right]
 \end{aligned}$$

Gradients can be computed by backprop

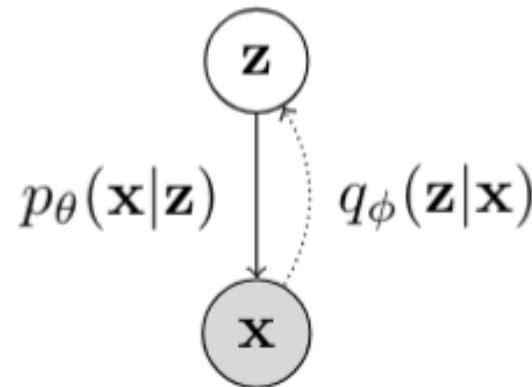
The mapping  $\mathbf{z}$  is a deterministic neural net for fixed  $\epsilon$



# Importance Weighted Autoencoder

- ▶ Improve VAE by using the following k-sample importance weighting of the log-likelihood:

$$\mathcal{L}_k(\mathbf{x}) = \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{1}{k} \sum_{i=1}^k \frac{p_\theta(\mathbf{x}, \mathbf{z}_i)}{q_\phi(\mathbf{z}_i|\mathbf{x})} \right]$$



$$= \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{1}{k} \sum_{i=1}^k w_i \right]$$

- ▶ where multiple  $z$  are sampled from the recognition network.

unnormalized  
importance weights

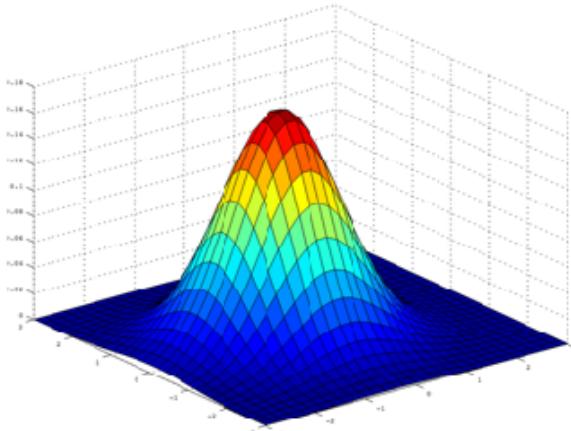
- ▶ Can improve the tightness of the bound.

# Talk Roadmap

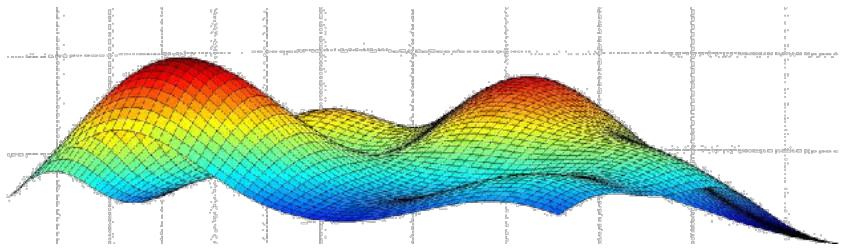
- ▶ Fully Observed Models
- ▶ Undirected Deep Generative Models
  - ▶ Restricted Boltzmann Machines (RBMs),
  - ▶ Deep Boltzmann Machines (DBMs)
- ▶ Directed Deep Generative Models
  - ▶ Variational Autoencoders (VAEs)
  - ▶ Normalizing Flows
- ▶ Generative Adversarial Networks (GANs)

# Generative Adversarial Networks (GAN)

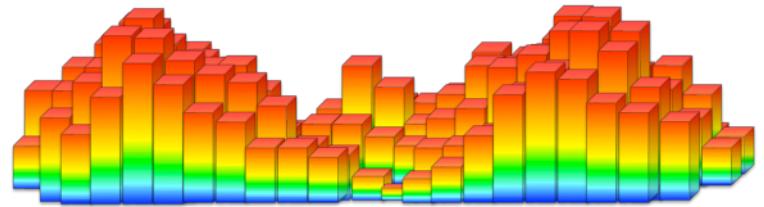
- ▶ Implicit generative model for an unknown target density  $p(x)$
- ▶ Converts sample from a known noise density  $p_Z(z)$  to the target  $p(x)$



Noise density  $p_Z(z)$  over space  $\mathcal{Z}$



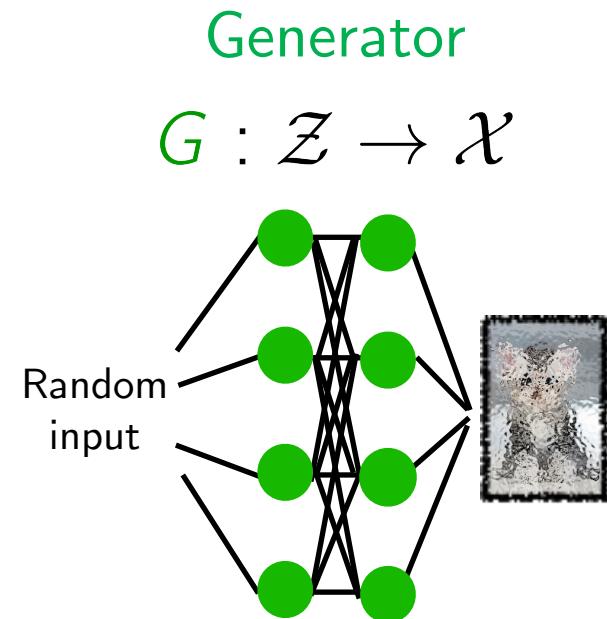
Unknown target density  $p(x)$  of data over domain  $\mathcal{X}$ , e.g.  $\mathbb{R}^{32 \times 32}$



Distribution of generated samples should follow target density  $p(x)$

# GAN Formulation

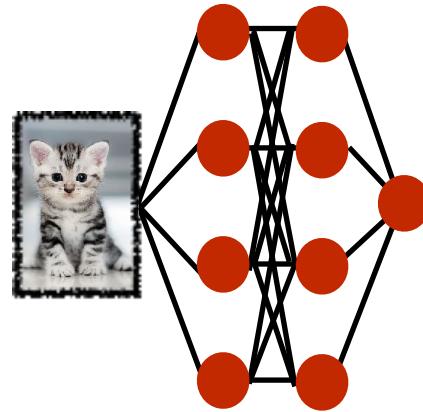
- GAN consists of two components



Goal: Produce samples  
indistinguishable from true data

Discriminator

$$D : \mathcal{X} \rightarrow \mathbb{R}$$

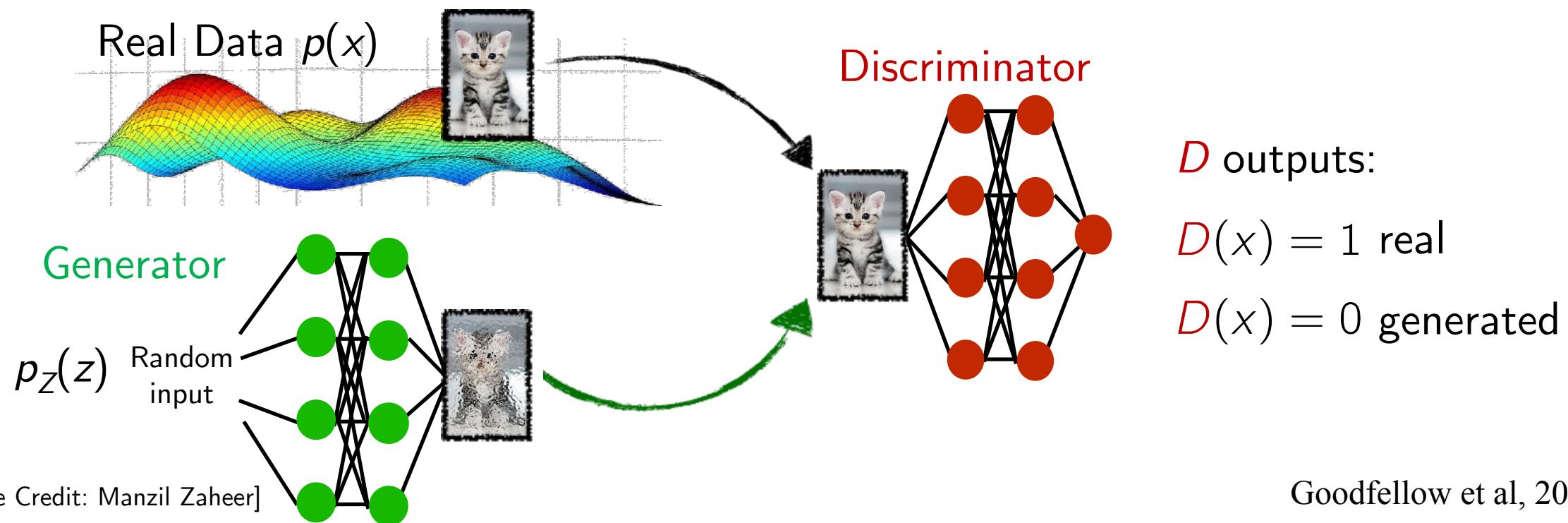


Goal: Distinguish  
true and generated  
data apart

# GAN Formulation: Discriminator

- Discriminator's objective: Tell real and generated data apart like a classifier

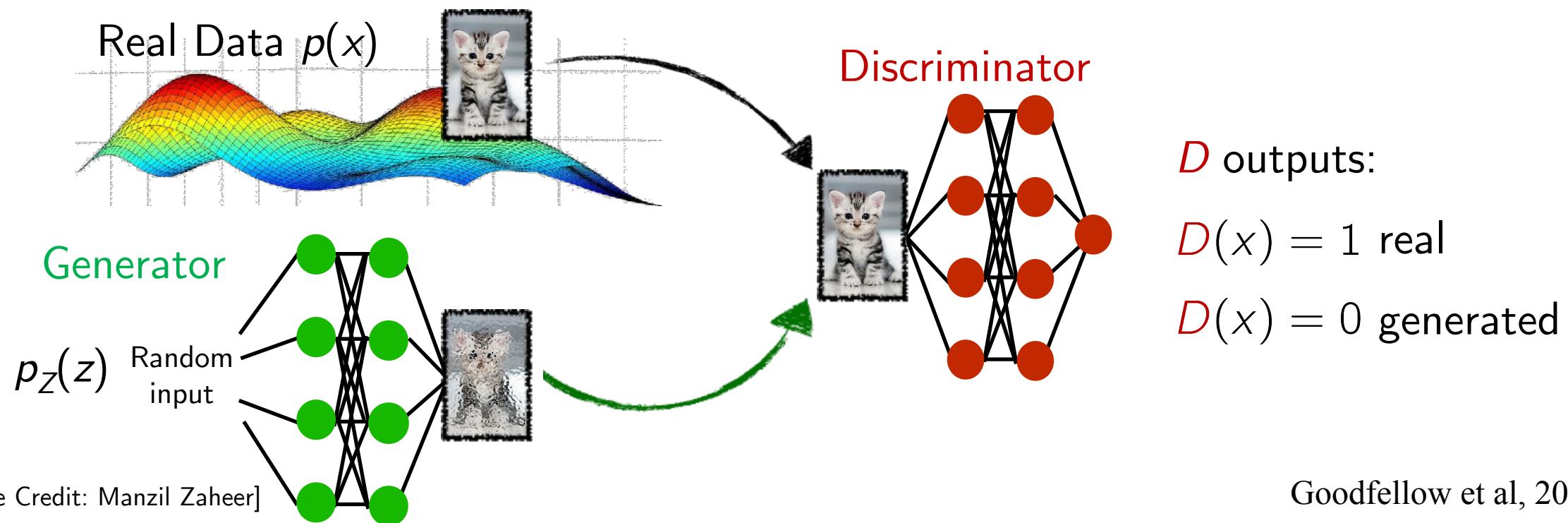
$$\max_D \mathbb{E}_{x \sim p} [\log D(x)] + \mathbb{E}_{z \sim p_Z} [\log (1 - D(G(z)))]$$



# GAN Formulation: Generator

- Generator's objective: Fool the best **discriminator**

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



# GAN Formulation: Optimization

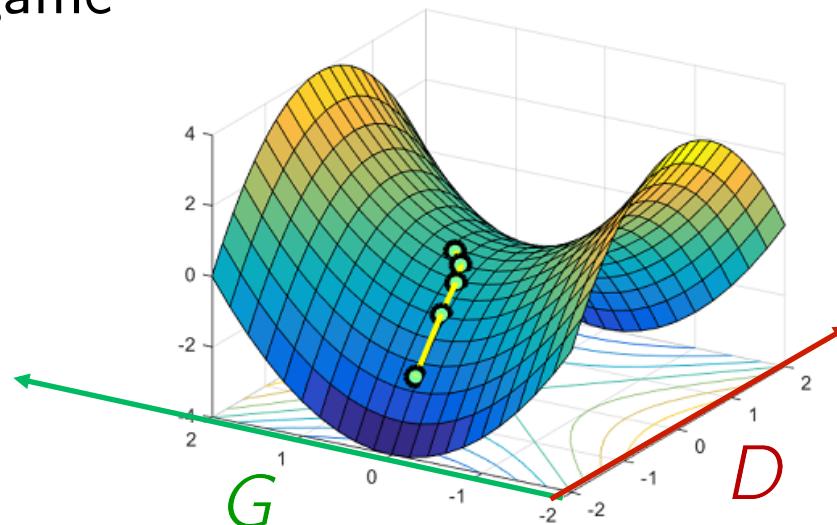
- ▶ Overall GAN optimization

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

- ▶ The generator-discriminator are iteratively updated using SGD to find “equilibrium” of a “min-max objective” like a game

$$G \leftarrow G - \eta_G \nabla_G V(G, D)$$

$$D \leftarrow D - \eta_D \nabla_D V(G, D)$$



# Wasserstein GAN

- ▶ WGAN optimization

$$\min_G \max_D W(G, D) = \mathbb{E}_{x \sim p} [D(x)] - \mathbb{E}_{z \sim p_Z} [D(G(z))]$$

- ▶ Difference in expected output on real vs. generated images
  - ▶ Generator attempts to drive objective  $\approx 0$
- ▶ More stable optimization

Compare to training DBMs

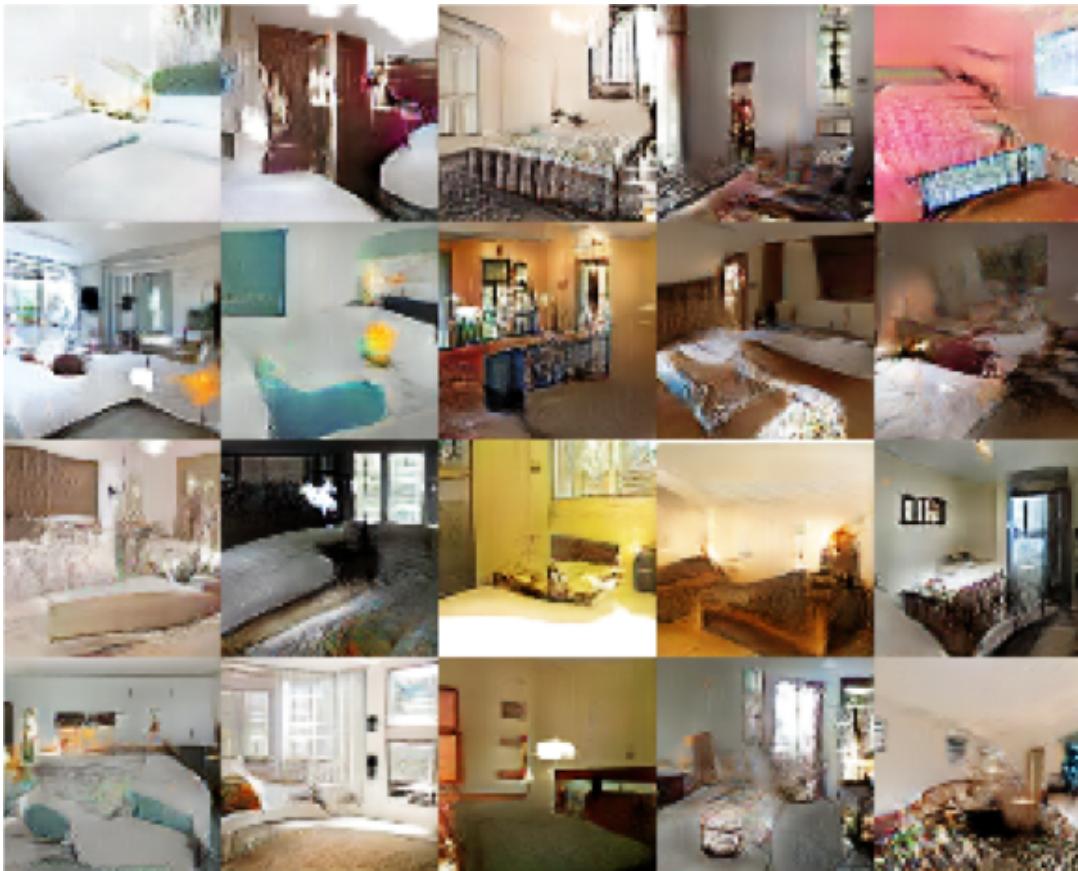
$$\frac{\partial \log P_\theta(\mathbf{v})}{\partial W^1} = \mathbb{E}_{P_{data}} [\mathbf{v} \mathbf{h}^1] - \mathbb{E}_{P_\theta} [\mathbf{v} \mathbf{h}^1]$$

$D$  outputs:

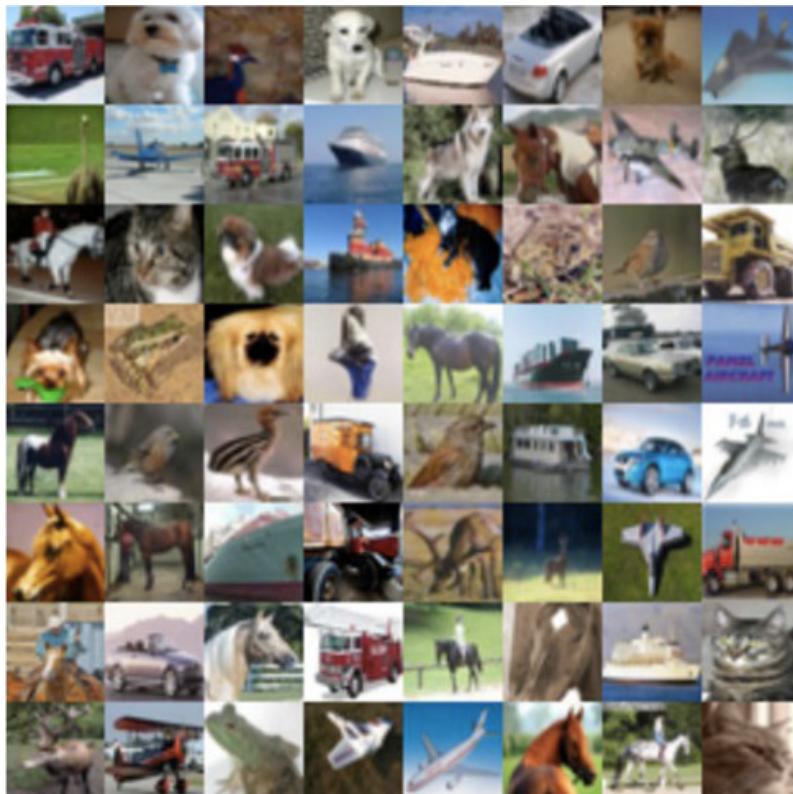
$$D(x) = 1 \text{ real}$$

$$D(x) = 0 \text{ generated}$$

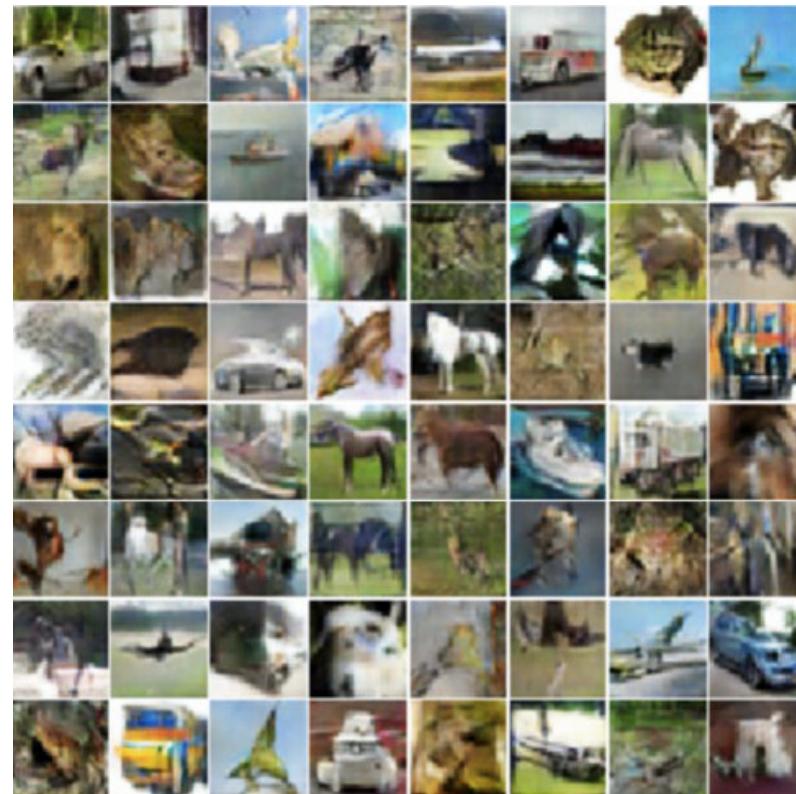
# LSUN Bedroom: Samples



# CIFAR Dataset

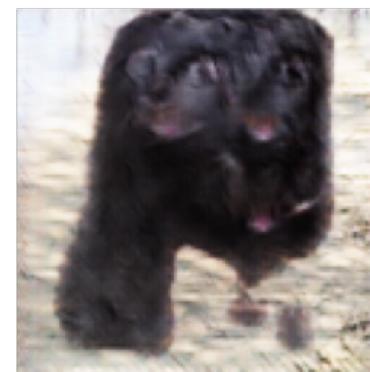
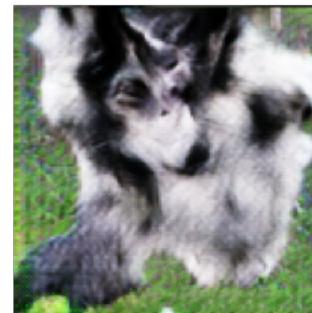
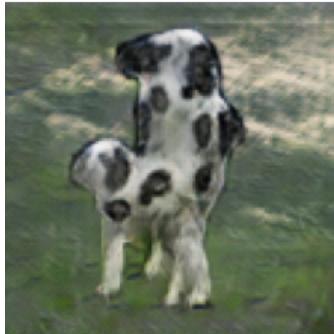
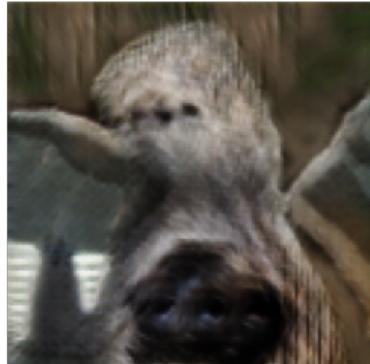
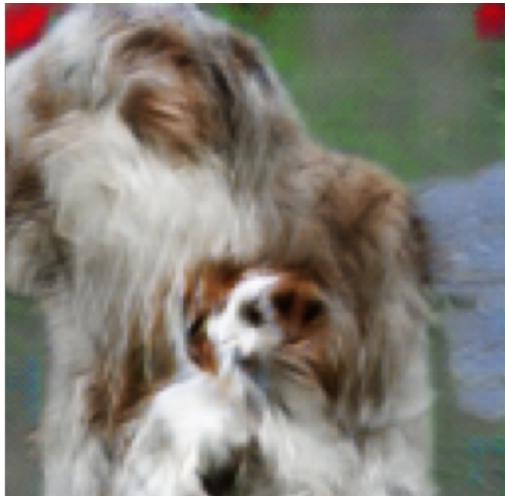


Training



Samples

# ImageNet: Cherry-Picked Samples

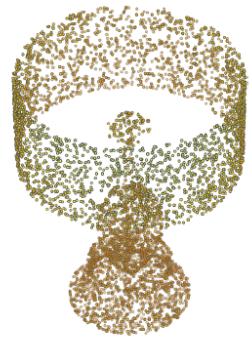


- ▶ Open Question: How can we quantitatively evaluate these models!

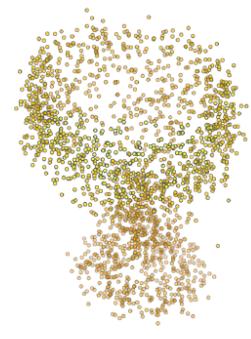
Slide Credit: Ian Goodfellow

# Modelling Point Cloud Data

Data



AAE



PC-GAN

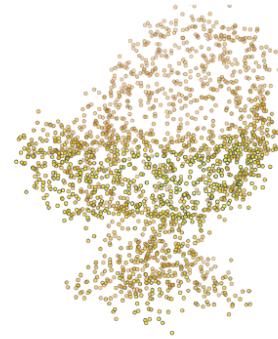


(a) Lamp

Data



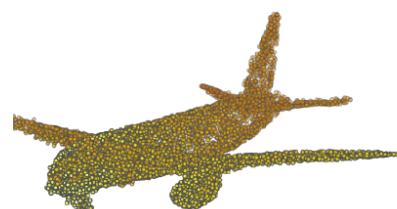
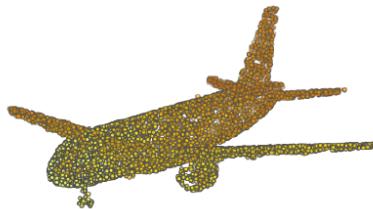
AAE



PC-GAN



(b) Chair

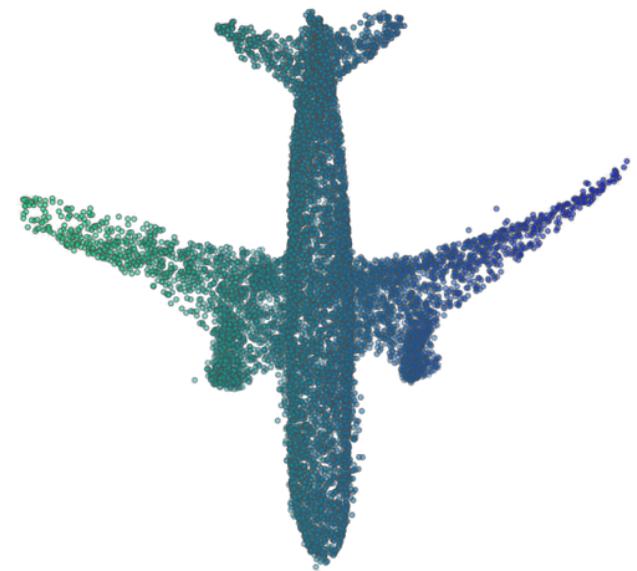
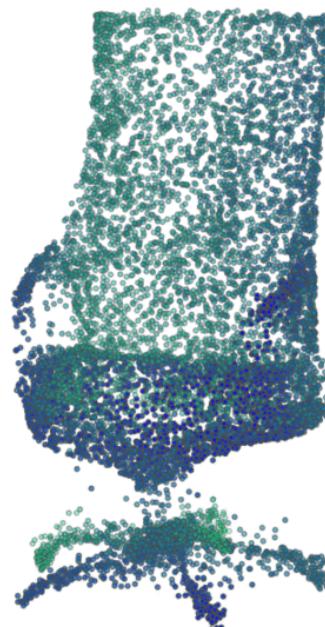
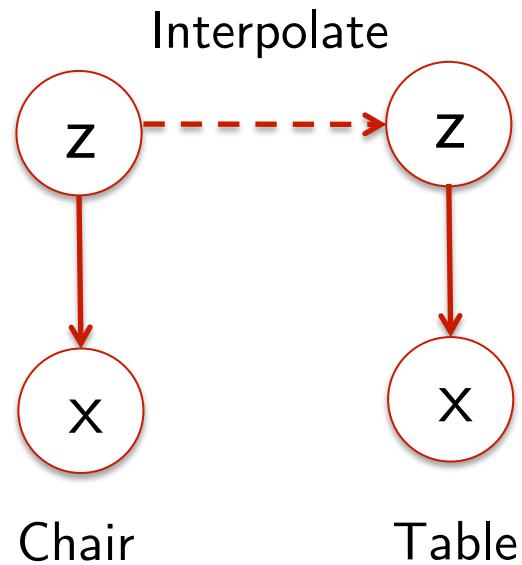


(c) Plane



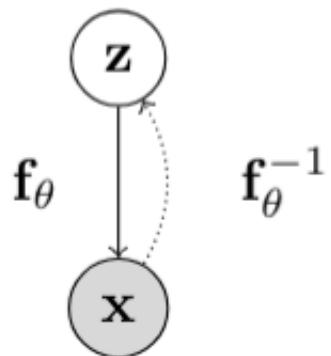
(d) Guitar

# Interpolation in Latent Space



# Normalizing Flows

- Directed Latent Variable Invertible models



- The mapping between  $x$  and  $z$  is deterministic and invertible:

$$\mathbf{x} = \mathbf{f}_\theta(\mathbf{z})$$

$$\mathbf{z} = \mathbf{f}_\theta^{-1}(\mathbf{x})$$

- Use change-of-variables to relate densities between  $z$  and  $x$

$$p_X(\mathbf{x}; \theta) = p_Z(\mathbf{z}) \left| \det \frac{\partial \mathbf{f}_\theta^{-1}(\mathbf{x})}{\partial \mathbf{X}} \right|_{\mathbf{X}=\mathbf{x}}$$

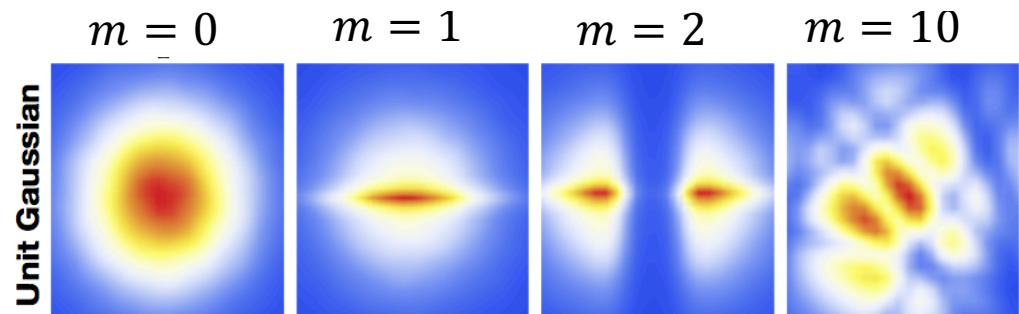
# Normalizing Flows

- Invertible transformations can be composed:

$$\mathbf{x} = \mathbf{f}_\theta^M \circ \dots \circ \mathbf{f}_\theta^1(\mathbf{z}^0); \quad p_X(\mathbf{x}; \theta) = p_{Z^0}(\mathbf{z}^0) \prod_{m=1}^M \left| \det \frac{\partial (\mathbf{f}_\theta^m)^{-1}}{\partial Z^m} \right|_{Z^m=\mathbf{z}^m}$$

- Planar Flows

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}g(\mathbf{w}^\top \mathbf{z} + b)$$



Rezende and Mohamed, 2016

# Normalizing Flows

- ▶ Maximum log-likelihood objective

$$\max_{\theta} \log p_X(\mathcal{D}; \theta) = \sum_{\mathbf{x} \in \mathcal{D}} \left( \log p_Z(\mathbf{z}) - \log \left| \det \frac{\partial (\mathbf{f}_{\theta})^{-1}}{\partial X} \right|_{X=\mathbf{x}} \right)$$

- ▶ Exact log-likelihood evaluation via inverse transformations
- ▶ Sampling from the model

$$\mathbf{z} \sim p_Z(\mathbf{z}), \quad \mathbf{x} = \mathbf{f}_{\theta}(\mathbf{z})$$

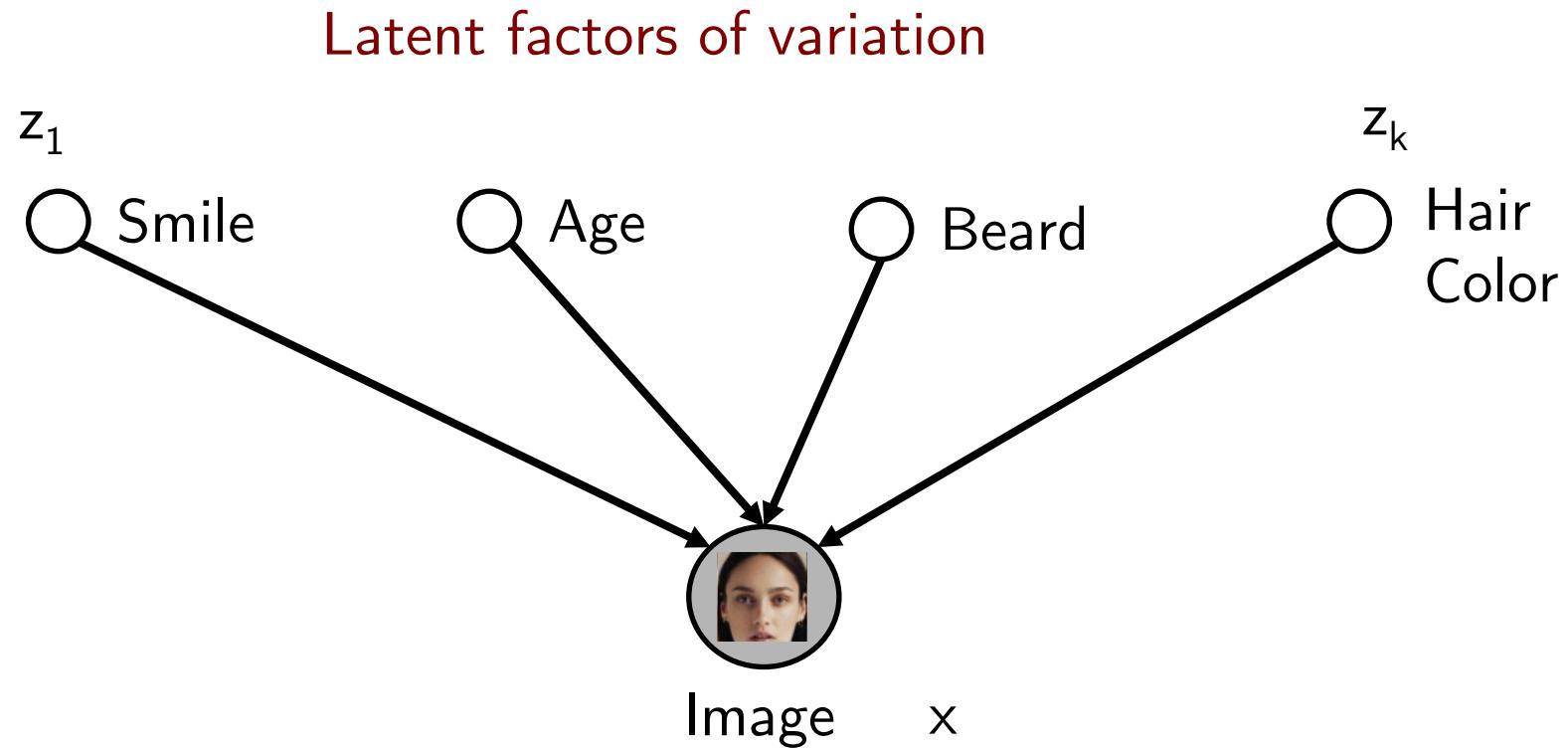
- ▶ Inference over the latent representations:

$$\mathbf{z} = \mathbf{f}_{\theta}^{-1}(\mathbf{x})$$

# Example: GLOW

- Generative Flow with Invertible 1x1 Convolutions

<https://blog.openai.com/glow/>



# Example: GLOW

Input



Smile



Add Beard



Increase Age



<https://blog.openai.com/>



Remove Beard



Decrease Age

Thank you