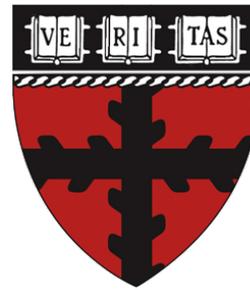
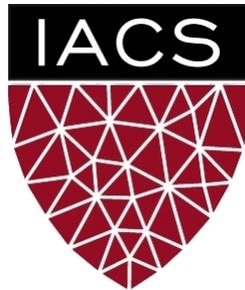


# Lecture 20: Deep Generative Models

CS 109B, STAT 121B, AC 209B, CSE 109B

Mark Glickman and Pavlos Protopapas



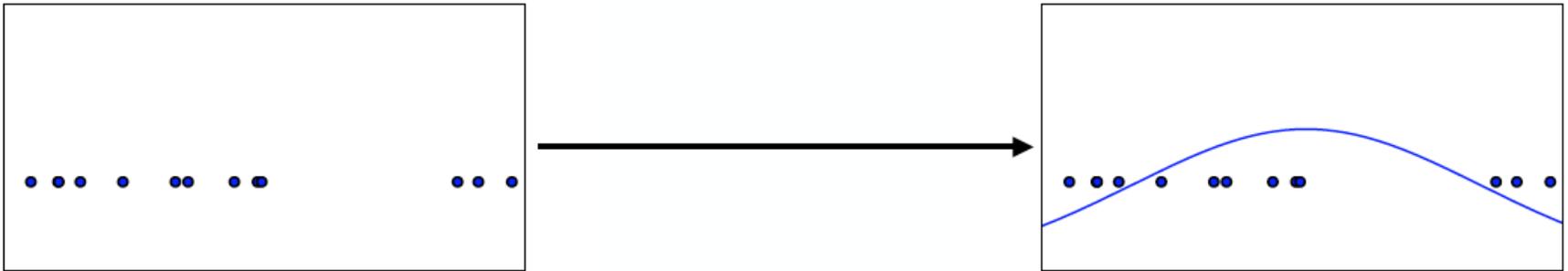
This lecture is based on the following tutorials:

I. Goodfellow, *Generative Adversarial Networks (GANs)*, NIPS 2016 Tutorial

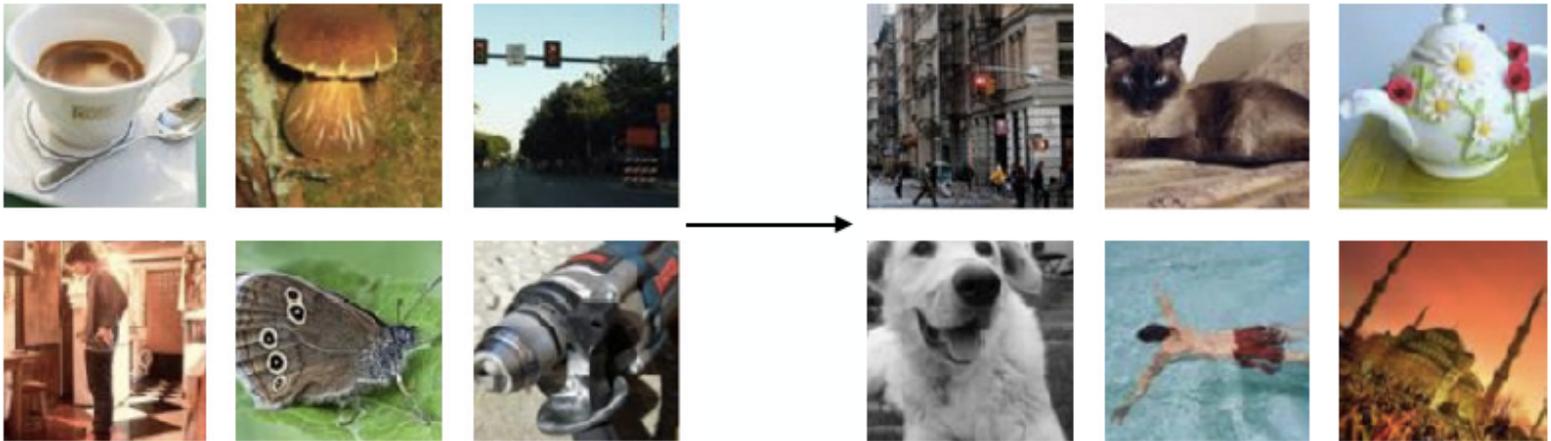
Shenlong Wang, *Deep Generative Models*, [http://www.cs.toronto.edu/~slwang/generative\\_model.pdf](http://www.cs.toronto.edu/~slwang/generative_model.pdf)

# Generative Model

Density Estimation



Sample Generation



Training examples

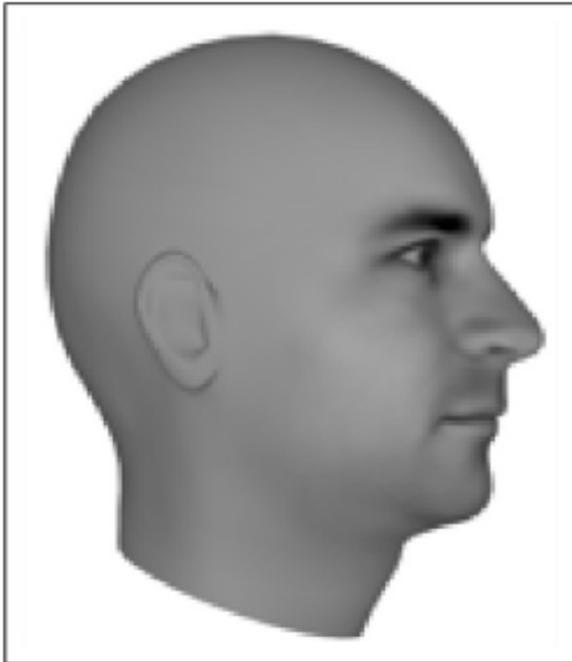
Model samples

# Why generative modeling?

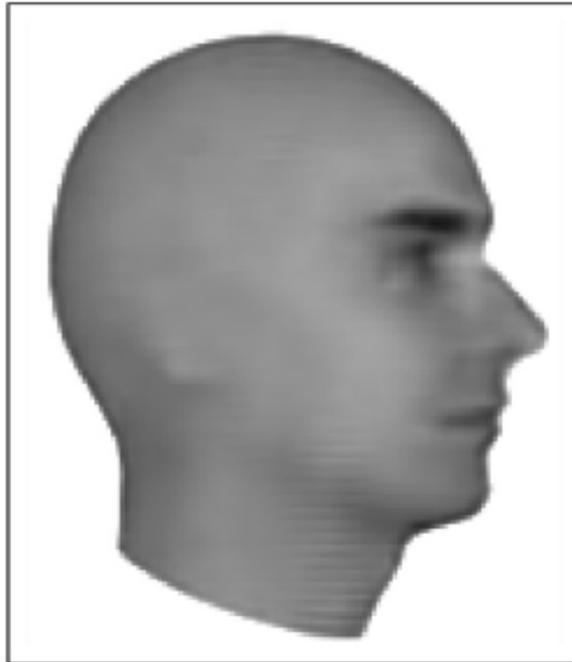
- Modeling high-dim probability distributions
- Denoising damaged samples
- Simulate future events for planning and RL
- Imputing missing data
  - Semi-supervised Learning
- Multi-modal data
  - Multiple correct answers

# Next Video Frame Prediction

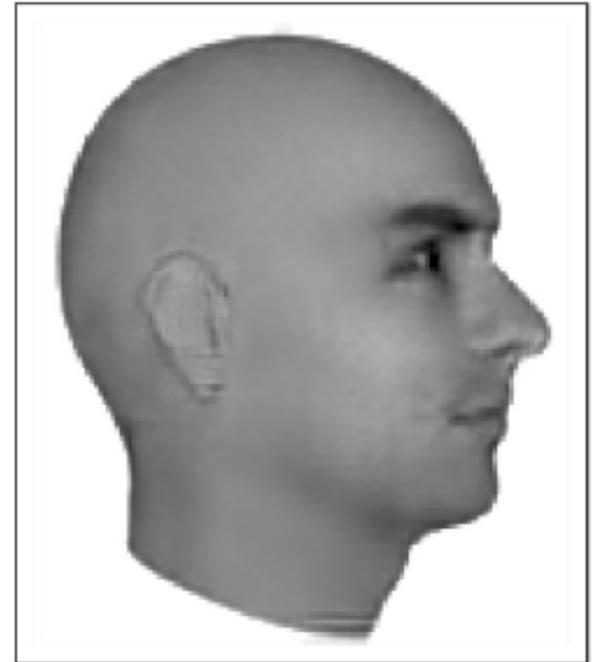
Ground Truth



MSE



Adversarial



# Single Image Super-resolution

- Synthesize a **high-resolution** equivalent from a low-resolution image

original



bicubic  
(21.59dB/0.6423)



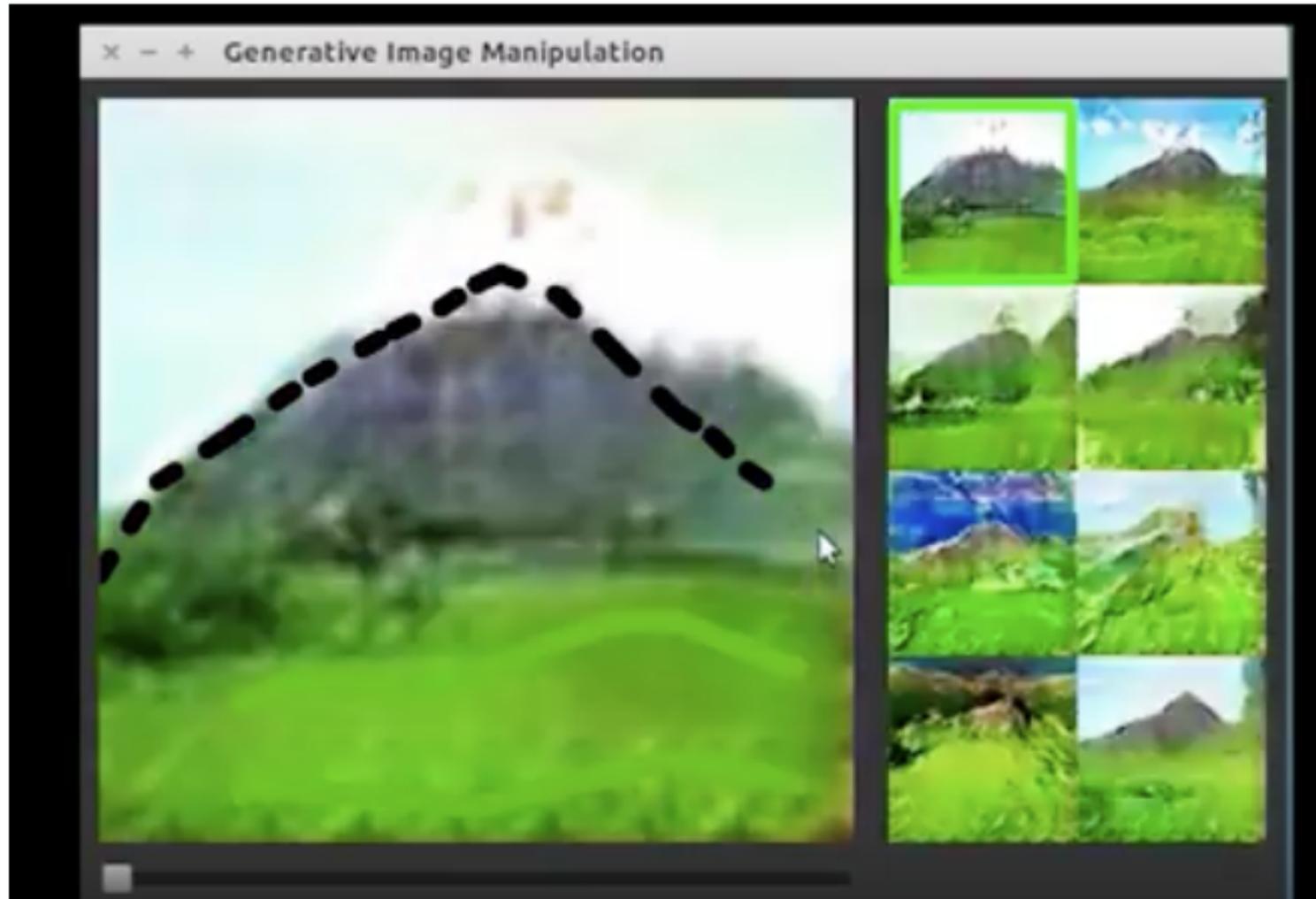
SRResNet  
(23.44dB/0.7777)



SRGAN  
(20.34dB/0.6562)

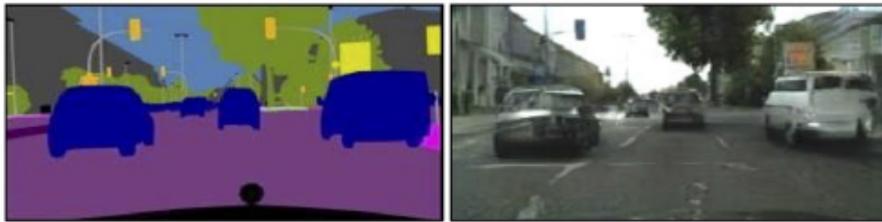


# Interactive Image Generation



# Image-to-Image Translation

Labels to Street Scene



input

output

Aerial to Map

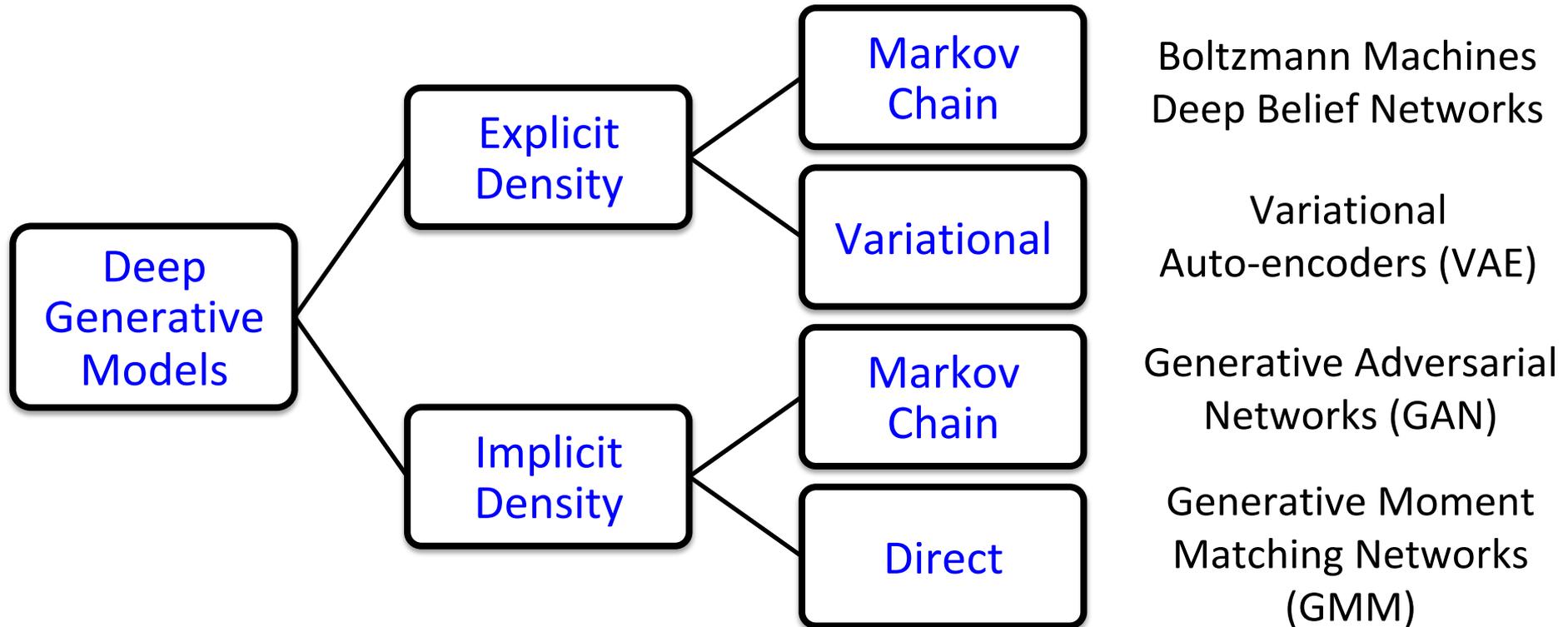


input

output



# Taxonomy of Deep Generative Models



# Outline

- *Explicit Density Models*
  - Boltzmann Machines & Deep Belief Networks (DBN)
  - Variational Auto-encoders (VAE)
- *Implicit Density Models*
  - Generative Adversarial Networks (GAN)

# Boltzmann Machines

$x$ :  $d$ -dimensional binary vector

$$P(x) = \frac{\exp(-E(x))}{Z}$$

$$Z = \sum_x \exp(-E(x))$$

Energy function given by:

$$E(x) = -x^T U x - b^T x$$

# Visible and Hidden Units

$x$ :  $d$ -dimensional binary vector

$$P(v, h) = \frac{\exp(-E(v, h))}{Z}$$

$$Z = \sum_{v, h} \exp(-E(v, h))$$

Energy function given by:

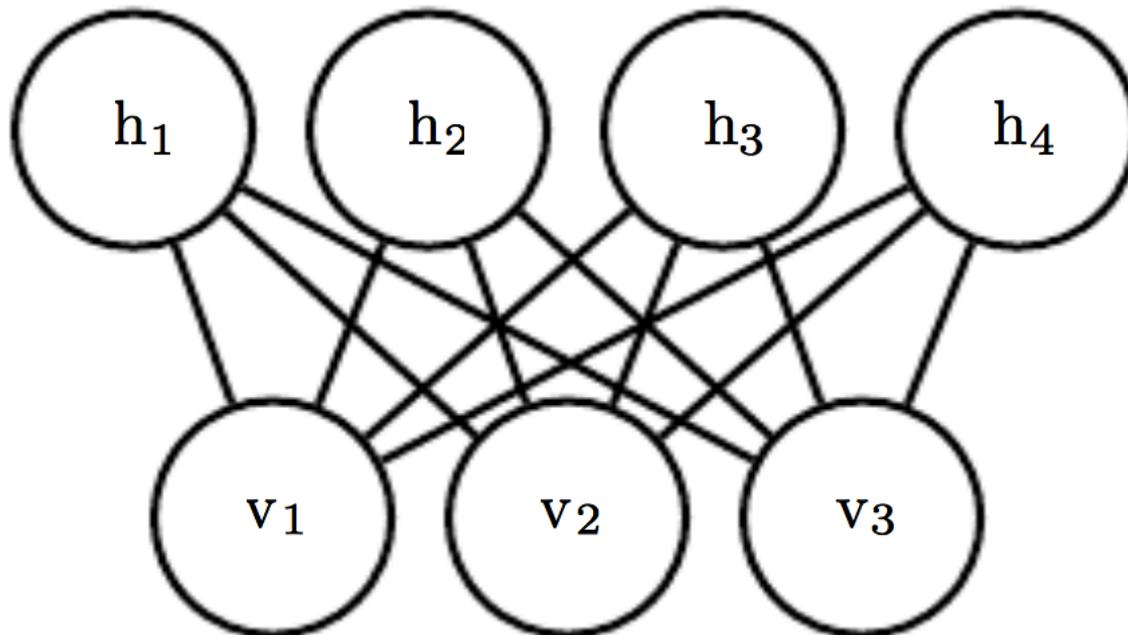
$$E(v, h) = -v^T R v - v^T W h - h^T S h - b^T v - c^T h$$

# Training Boltzmann Machines

- Maximum-likelihood estimation
- Partition function is intractable
- May be estimated with MCMC methods

# Restricted Boltzmann Machines

- No connections among hidden units, and among visible units (bipartite graph)



# Restricted Boltzmann Machines

$x$ :  $d$ -dimensional binary vector

$$P(v, h) = \frac{\exp(-E(v, h))}{Z}$$

$$Z = \sum_{v, h} \exp(-E(v, h))$$

Energy function given by:

$$E(v, h) = -b^T v - c^T h - v^T W h$$

# Training RBMs

- Conditional distributions have simple form

$$P(h_j = 1 | v) = \sigma(c_j + (W^T v)_j)$$

$$P(v_j = 1 | h) = \sigma(b_j + (Wh)_j)$$

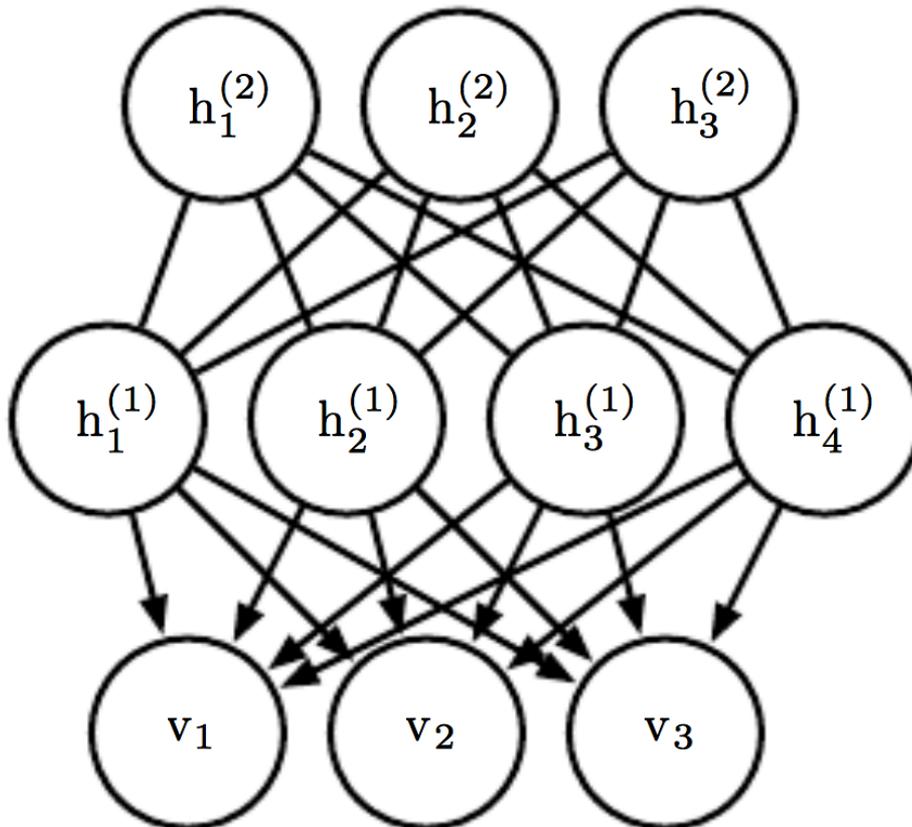
→ Sigmoid

- Efficient MCMC: **block Gibbs sampling**
  - Alternatively sample from  $P(h | v)$  and  $P(v | h)$

# Deep Belief Networks

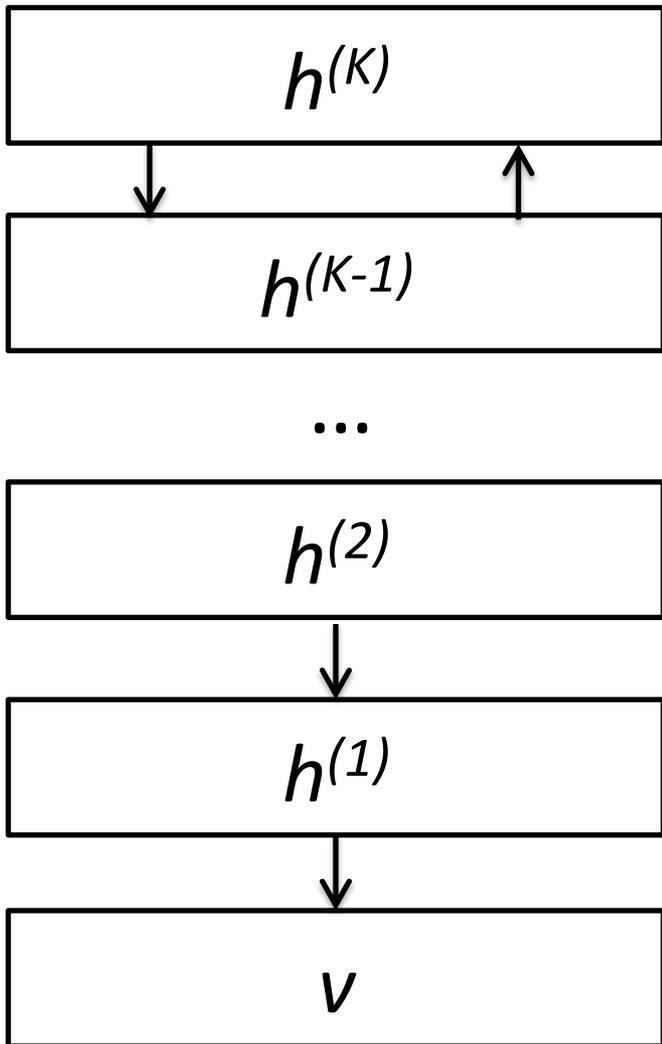
- Multiple hidden layers

Introduction of DBNs in 2006 began the current deep learning renaissance



Undirected connections between last two layers

# Deep Belief Networks



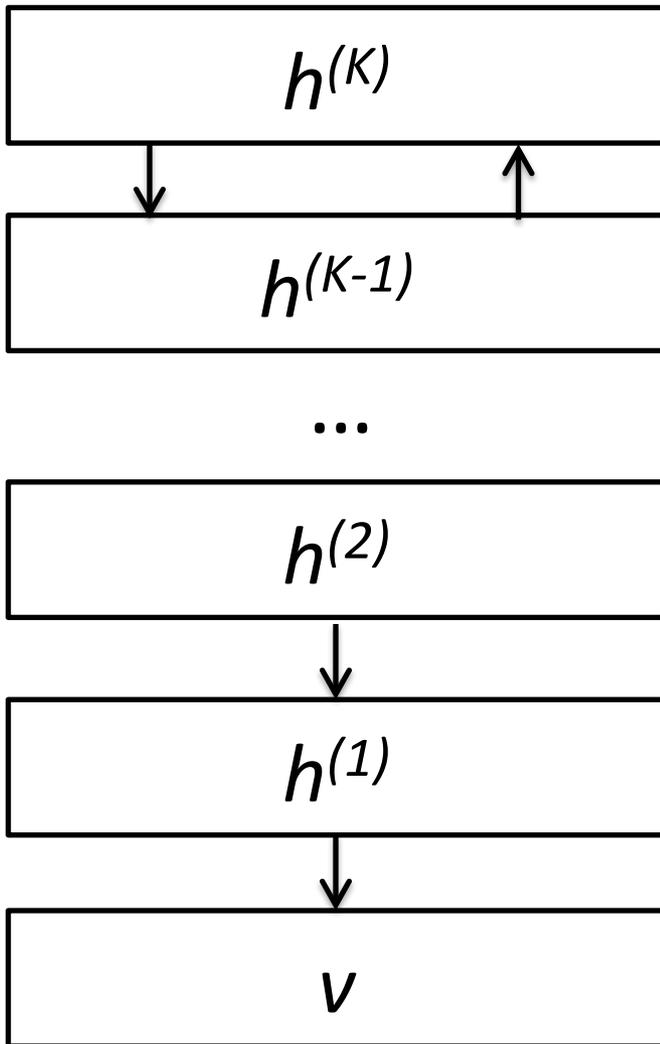
$$P(h^{(K-1)}, h^{(K)}) = \frac{1}{Z} \exp(E(h^{(K-1)}, h^{(K)}))$$

$\vdots$

$$P(h_i^{(1)} = 1 | h^{(2)}) = \sigma(b_i^{(2)} + w_i^{(2)} \cdot h^{(2)})$$

$$P(v_i = 1 | h^{(1)}) = \sigma(b_i^{(1)} + w_i^{(1)} \cdot h^{(1)})$$

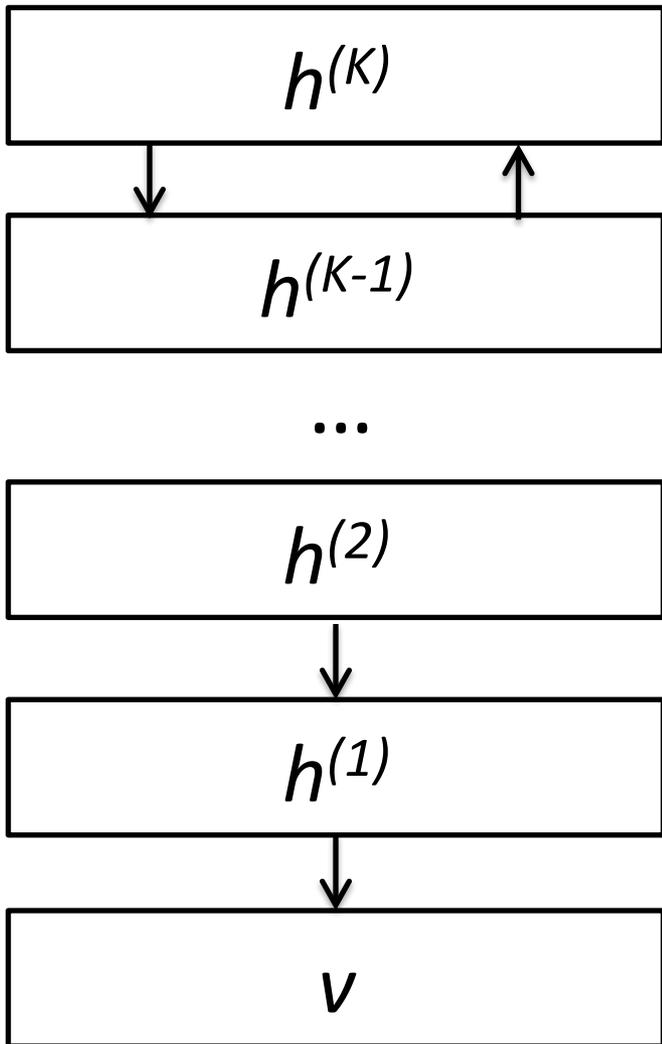
# Layer-wise DBN Training



Train an RBM to  
maximize

$$E_{v \sim p_{data}} \log p_{\theta}(v)$$

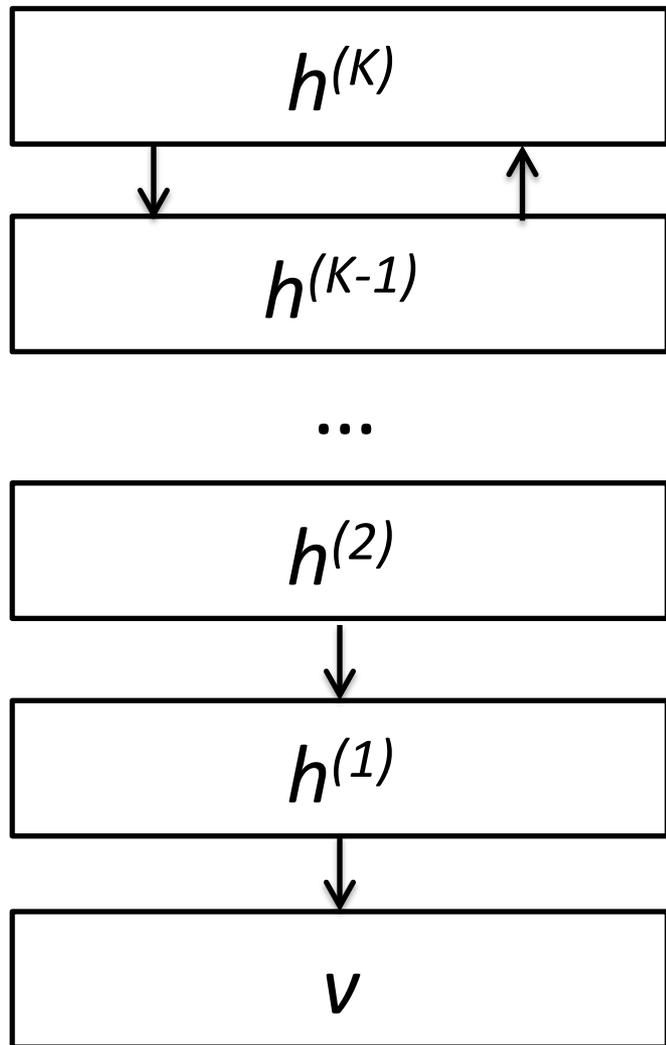
# Layer-wise DBN Training



Train an RBM to  
maximize

$$E_{v \sim p_{data}} E_{h^{(1)} \sim p(h^{(1)}|v)} \log p_{\theta_1}(h^{(1)})$$

# Layer-wise DBN Training



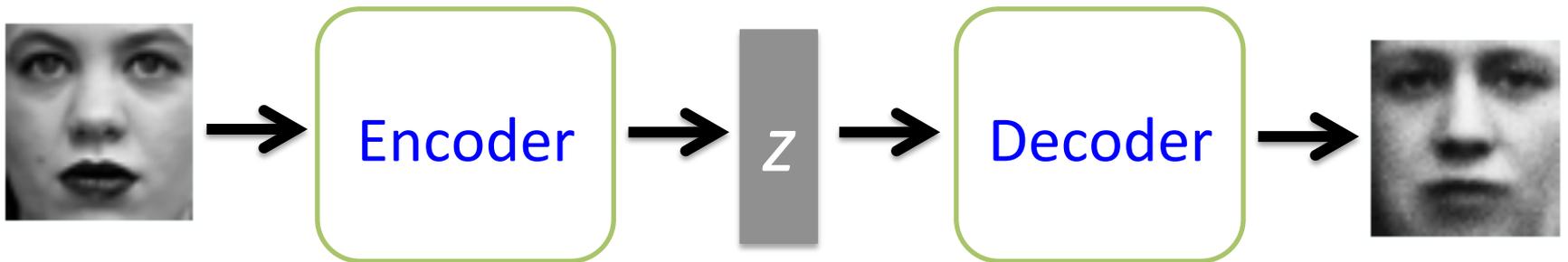
Train an RBM to  
maximize

$$E_{v \sim p_{data}} \cdots E_{h^{(K)} \sim p(h^{(K)} | h^{(K-1)})} \log p_{\theta_K}(h^{(K)})$$

# Outline

- *Explicit Density Models*
  - Boltzmann Machines & Deep Belief Networks (DBN)
  - Variational Auto-encoders (VAE)
- *Implicit Density Models*
  - Generative Adversarial Networks (GAN)

# Recap: Autoencoder



# Variational Autoencoder

- Easier to train using gradient-based methods
- VAE sampling:

Sample from  $P(z)$   
*Standard Gaussian*



Decoder  
 $P(x|z)$

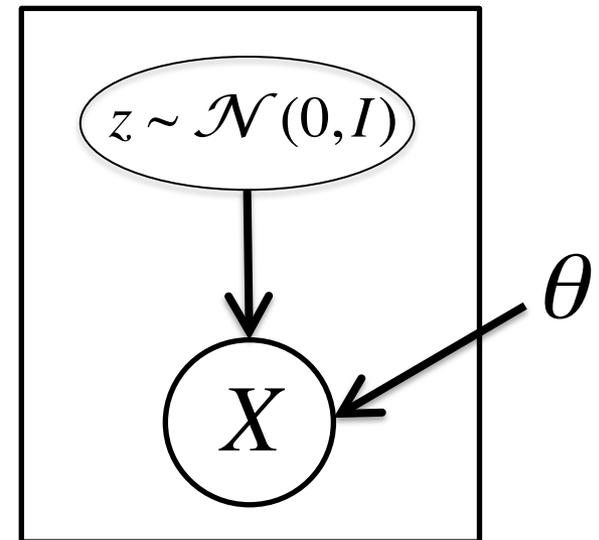


# VAE Likelihood

Neural network

$$p_{\theta}(x) = \int_z p_{\theta}(x|z) p_{\theta}(z) dz$$

Difficult to approximate in high dim through sampling



For most  $z$  values  $p(x|z)$  close to 0

# VAE Likelihood

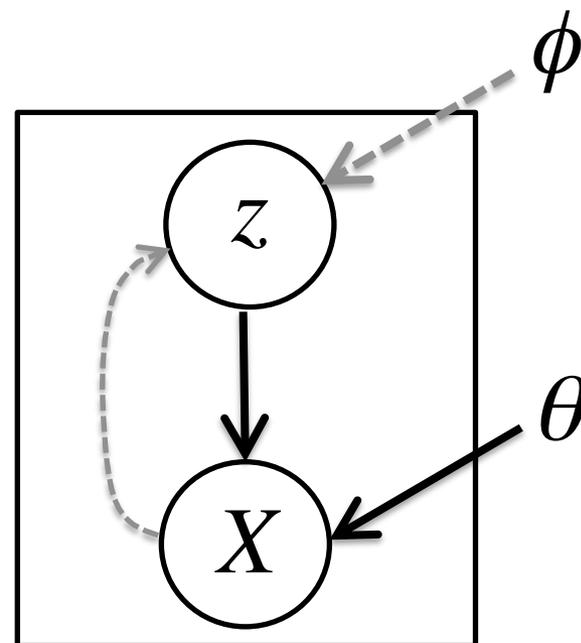
$$p_{\theta}(x) = \int_z p_{\theta}(x|z) q_{\phi}(z|x) dz$$

Another neural net

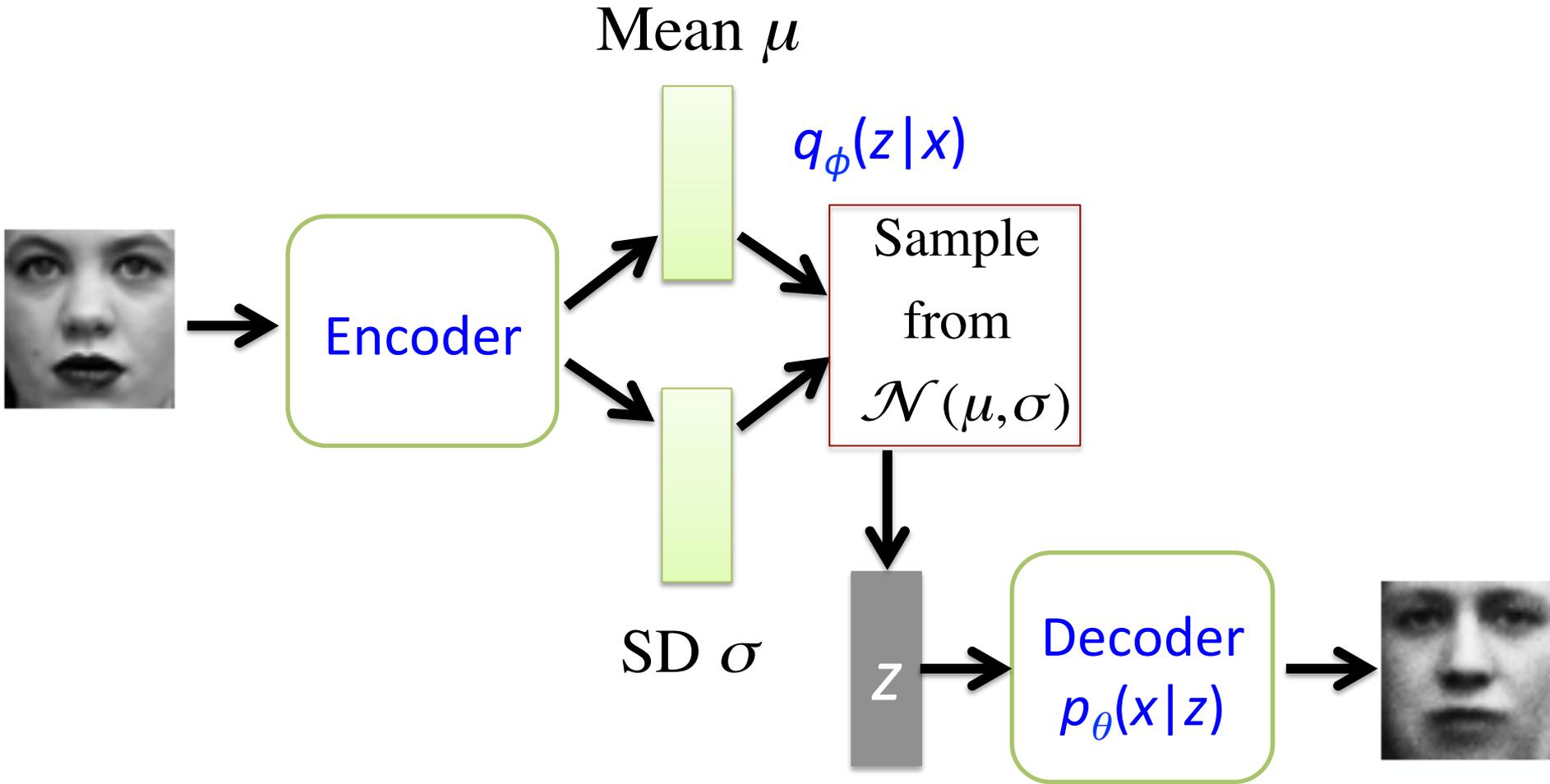


Proposal distribution:

Likely to produce values of  $x$   
for which  $p(x|z)$  is non-zero



# VAE Architecture



# VAE Loss

Reconstruction Loss

$$-\mathbf{E}_{z \sim q_{\phi}(z|x)} \log \left( p_{\theta}(x|z) \right)$$

# VAE Loss

Reconstruction Loss

Proposal distribution should resemble a Gaussian

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + KL(q_\phi(z|x) \| p_\theta(z))$$

# VAE Loss

Reconstruction Loss

Proposal distribution should resemble a Gaussian

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + KL(q_\phi(z|x) \| p_\theta(z))$$

$$\geq -\log p_\theta(x)$$

Variational upper bound  
on loss we care about!

# Training VAE

- Apply stochastic gradient descent
- Sampling step not differentiable
- Use a re-parameterization trick
  - Move sampling to input layer, so that the sampling step is independent of the model

# Boltzmann Machines vs. VAE

## *Pros:*

- VAE is easier to train
- Does not require MCMC sampling
- Has theoretical backing
- Applicable to wider range of models

## *Cons:*

- Samples from VAE tend to be **blurry**



6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
4	4	4	4	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	2	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	2	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	2	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	2	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	2	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	2	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	2	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	2	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	2	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	2	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	2	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	2	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4

# Outline

- *Explicit Density Models*
  - Boltzmann Machines & Deep Belief Networks (DBN)
  - Variational Auto-encoders (VAE)
- *Implicit Density Models*
  - Generative Adversarial Networks (GAN)

# Generative Adversarial Networks

- No Markov chains needed
- No variational approximations needed
- Often regarded as producing better examples

# Two Player Game

## Generator $G$

Seeks to create  
samples from  $p_{data}$

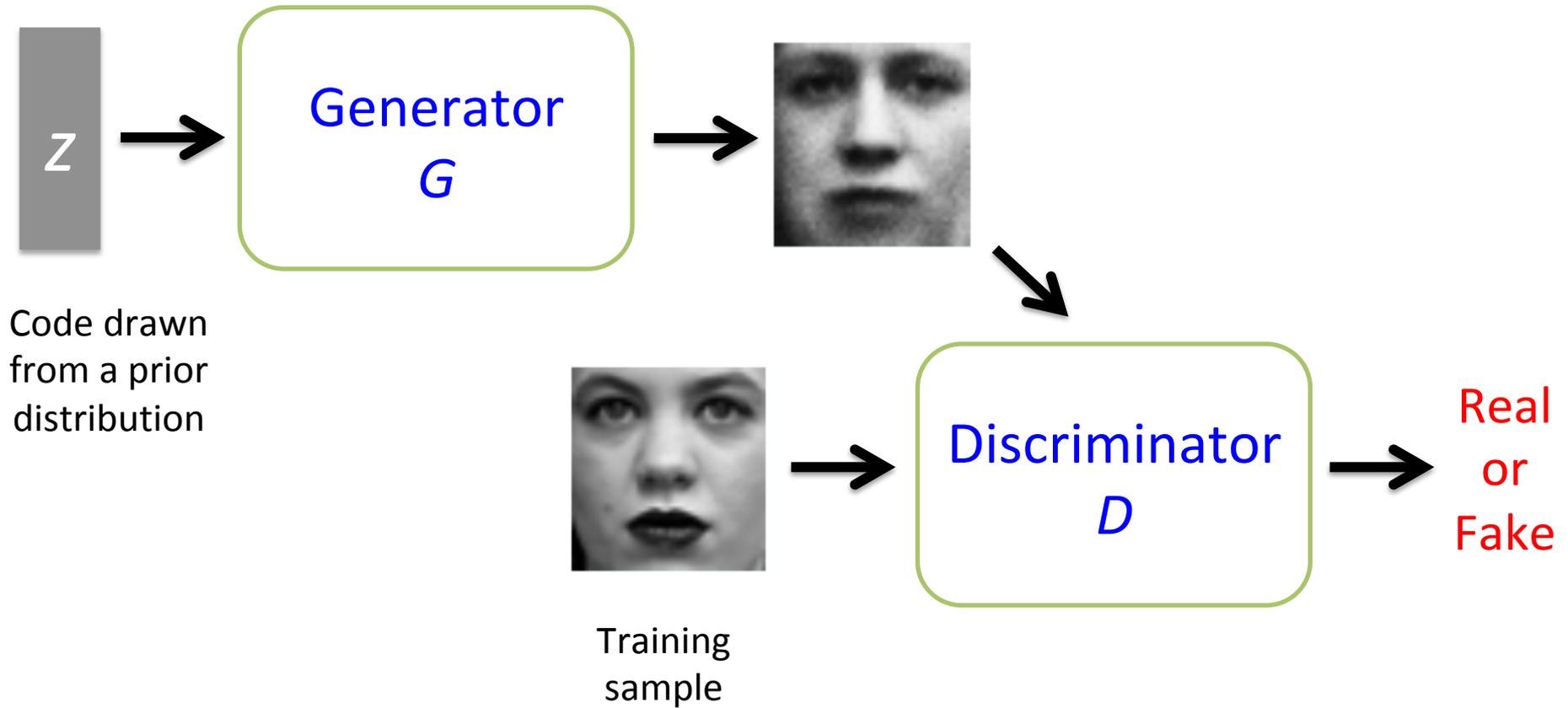
Seeks to trick discriminator  
into believing that its samples  
are genuine

## Discriminator $D$

Classifies samples  
as real or fake

Seeks to not get tricked  
by the generator

# GAN Overview



# Min-max Cost

Discriminator seeks to  
**predict 1** on training  
samples

*Discriminator* wants to  
**predict 0** on samples  
generated by  $G$

$$V(\theta^D, \theta^G) = \mathbf{E}_{x \sim P_{data}} \log D(x) + \mathbf{E}_z \log(1 - D(G(z)))$$

Generator wants  $D$  to not distinguish between  
original and generated samples!

$$J^{(G)} = V(\theta^{(D)}, \theta^{(G)})$$

$$J^{(D)} = -V(\theta^{(D)}, \theta^{(G)})$$

# Min-max Cost

$$\min_{\theta^G} \max_{\theta^D} J(\theta^G, \theta^D)$$

- Equilibrium is a saddle-point
- Training is difficult in practice

# Training GANs

- Sample mini-batch of training images  $x$ , and generator codes  $z$
- Update  $G$  using back-prop
- Update  $D$  using back-prop

*Optional:* Run  $k$  steps of one player for every step of the other player

