

# Deep Unsupervised Learning

Russ Salakhutdinov

Machine Learning Department  
Carnegie Mellon University  
Canadian Institute of Advanced Research

Carnegie  
Mellon  
University

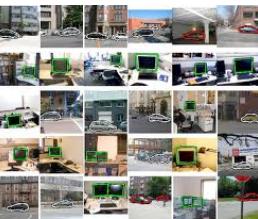


**CIFAR**  
CANADIAN INSTITUTE  
for ADVANCED RESEARCH

# Mining for Structure

Massive increase in both computational power and the amount of data available from web, video cameras, laboratory measurements.

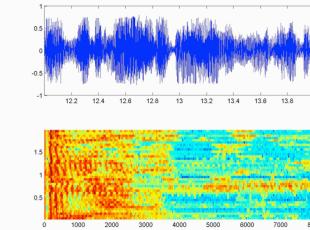
Images & Video



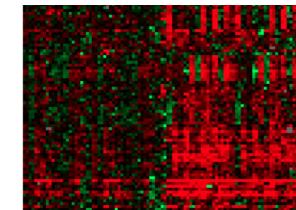
Text & Language



Speech & Audio



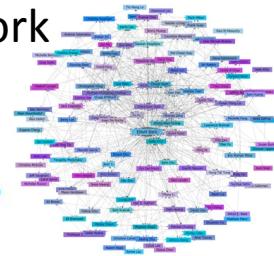
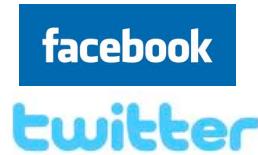
Gene Expression



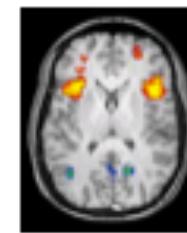
Product Recommendation



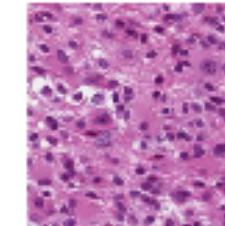
Relational Data/  
Social Network



fMRI



Tumor region



Mostly Unlabeled

- Develop statistical models that can discover underlying structure, cause, or statistical correlation from data in **unsupervised** or **semi-supervised** way.
- Multiple application domains.

# Mining for Structure

Massive increase in both computational power and the amount of data available from web, video cameras, laboratory measurements.

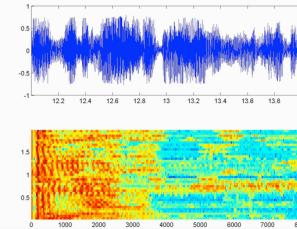
Images & Video



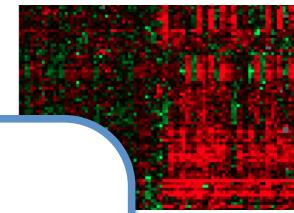
Text & Language



Speech & Audio



Gene Expression



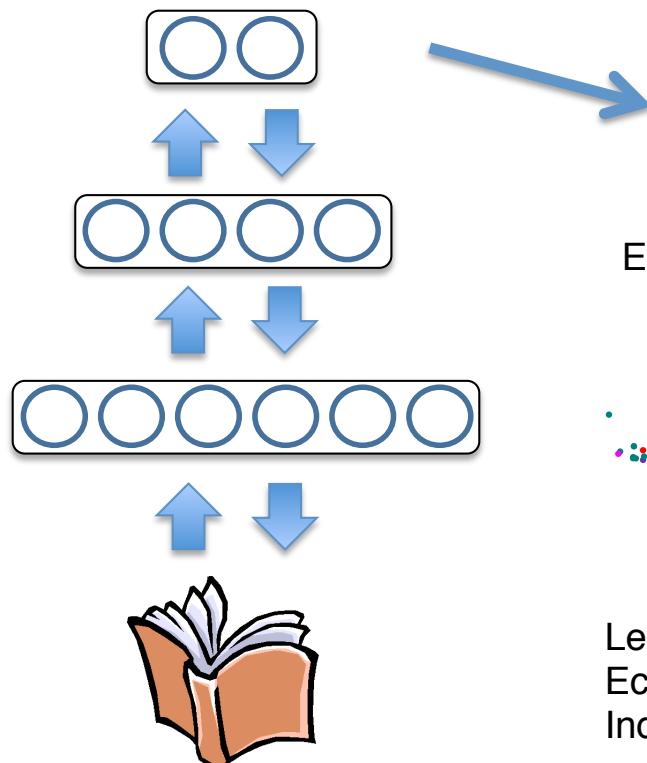
Deep Learning Models that support inferences and discover structure at multiple levels.

Mostly Unlabeled

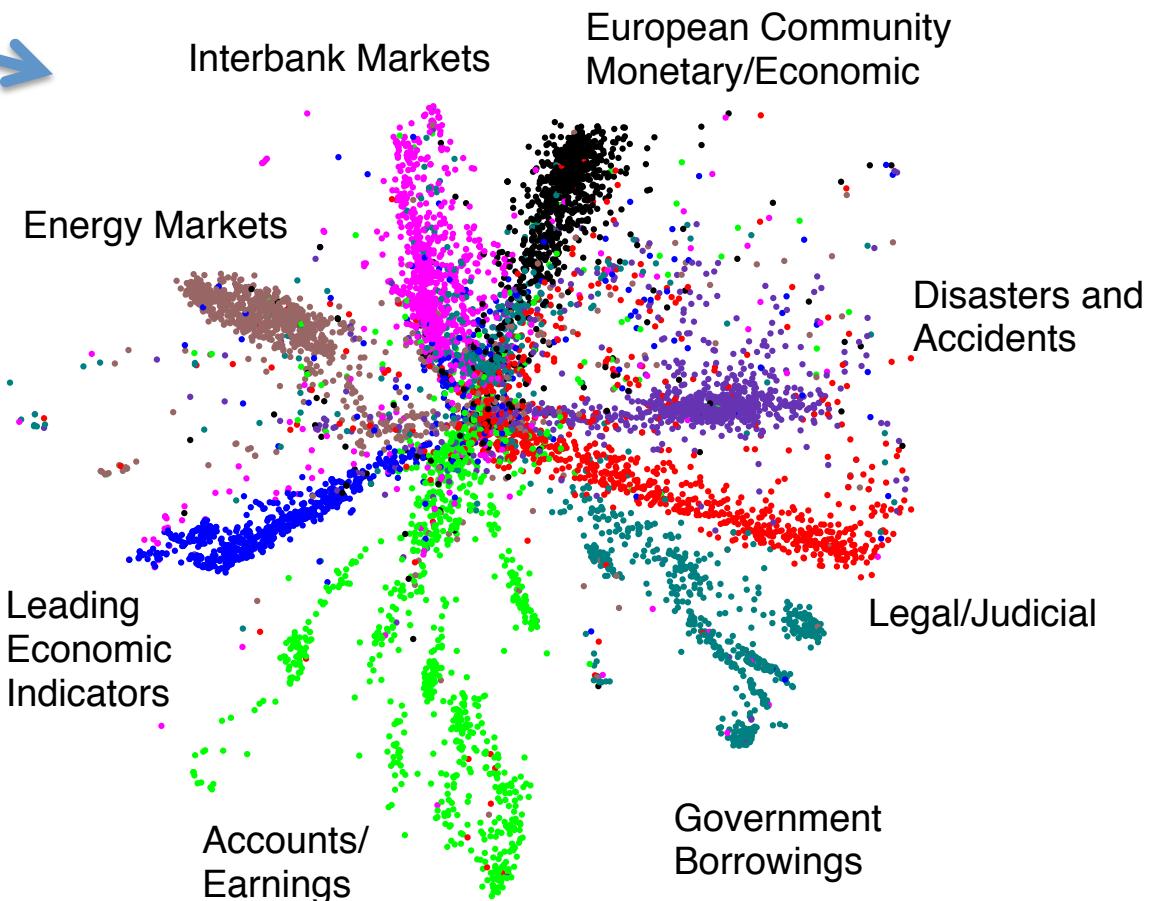
- Develop statistical models that can discover underlying structure, cause, or statistical correlation from data in **unsupervised** or **semi-supervised** way.
- Multiple application domains.

# Deep Autoencoder Model

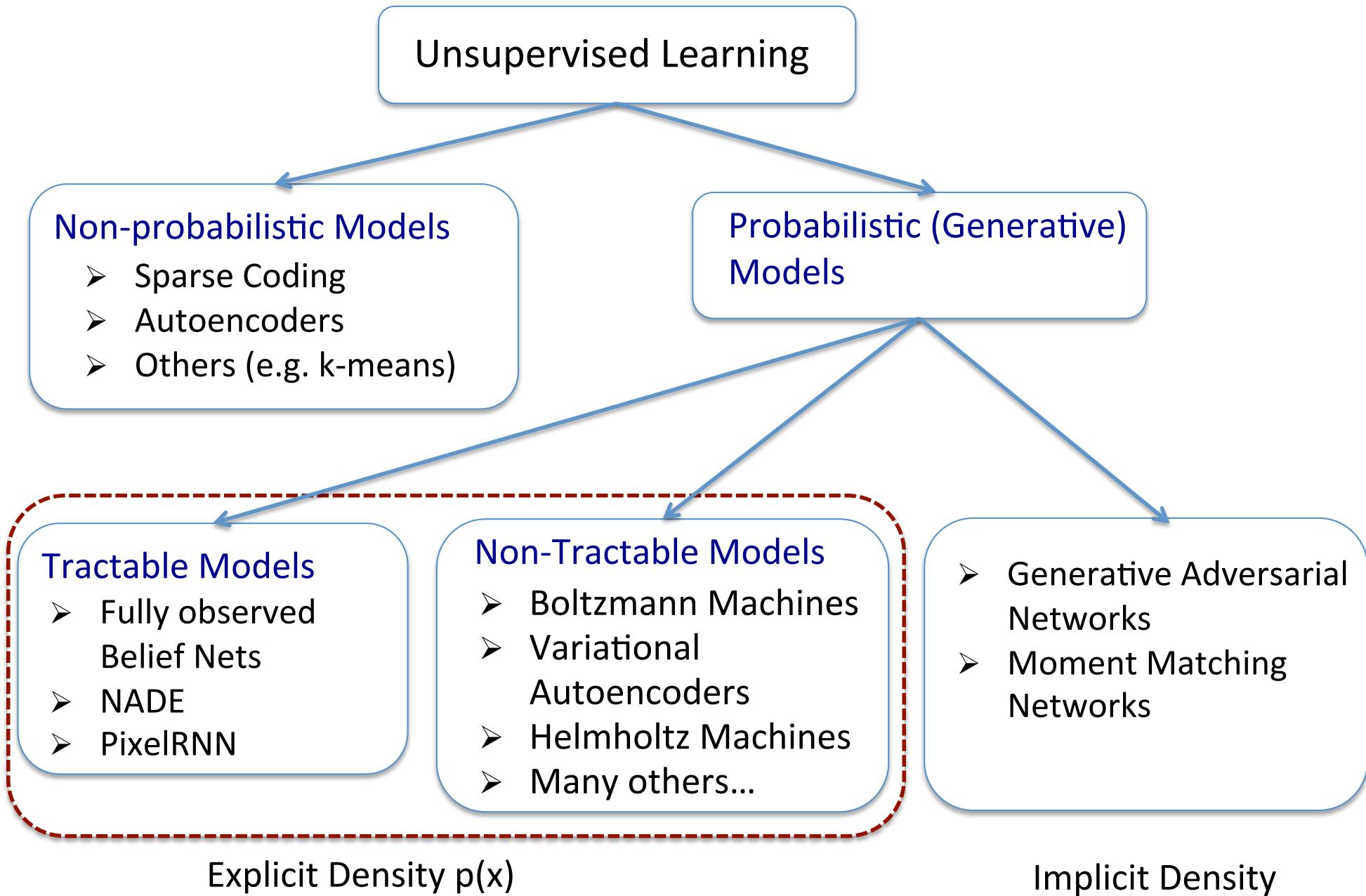
Learned latent code



Reuters dataset: 804,414  
newswire stories: **unsupervised**



(Hinton & Salakhutdinov, Science 2006)



# Talk Roadmap

- Basic Building Blocks:

- Sparse Coding
- Autoencoders

- Deep Generative Models

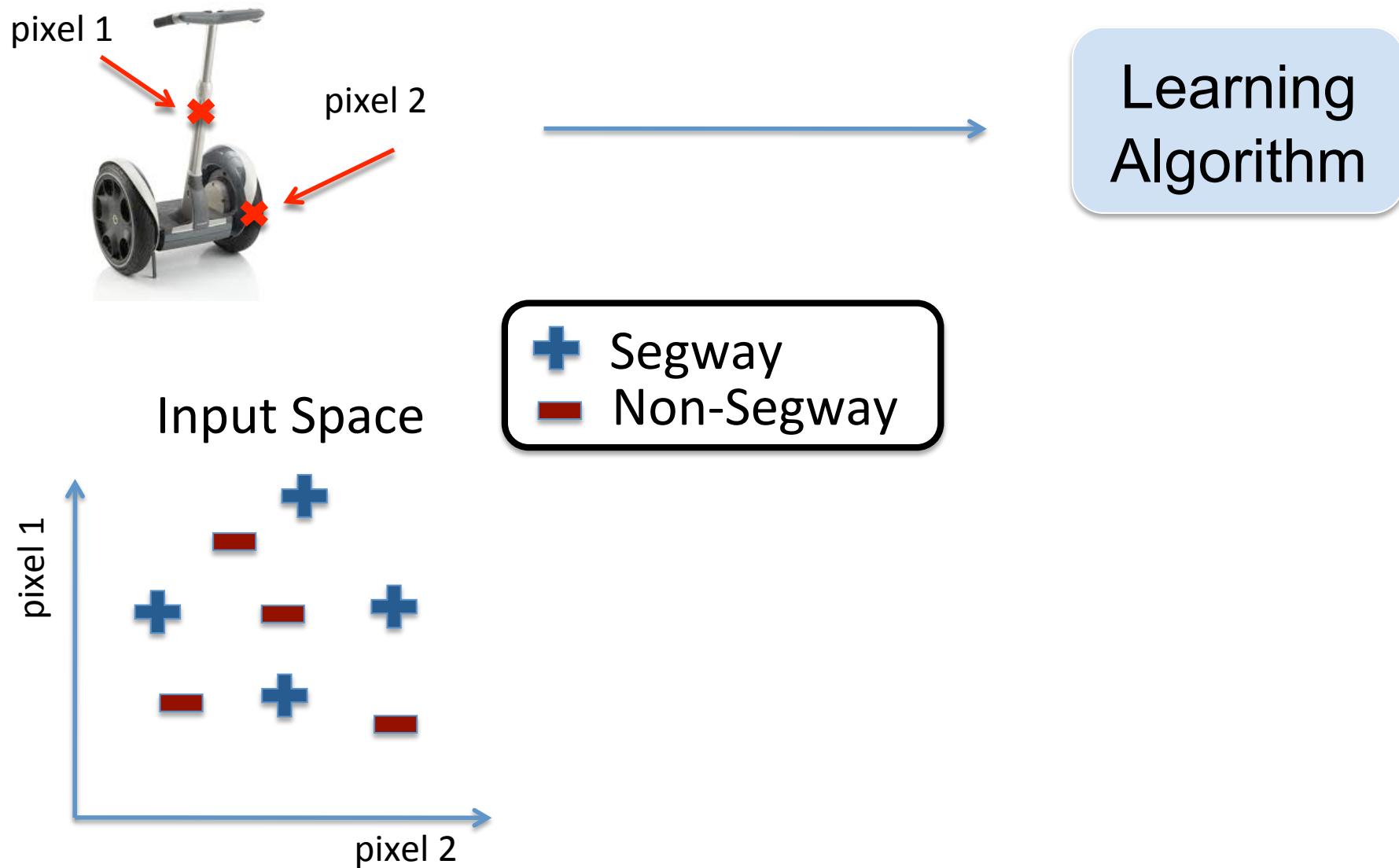
- Restricted Boltzmann Machines
- Deep Boltzmann Machines
- Helmholtz Machines / Variational Autoencoders

- Generative Adversarial Networks

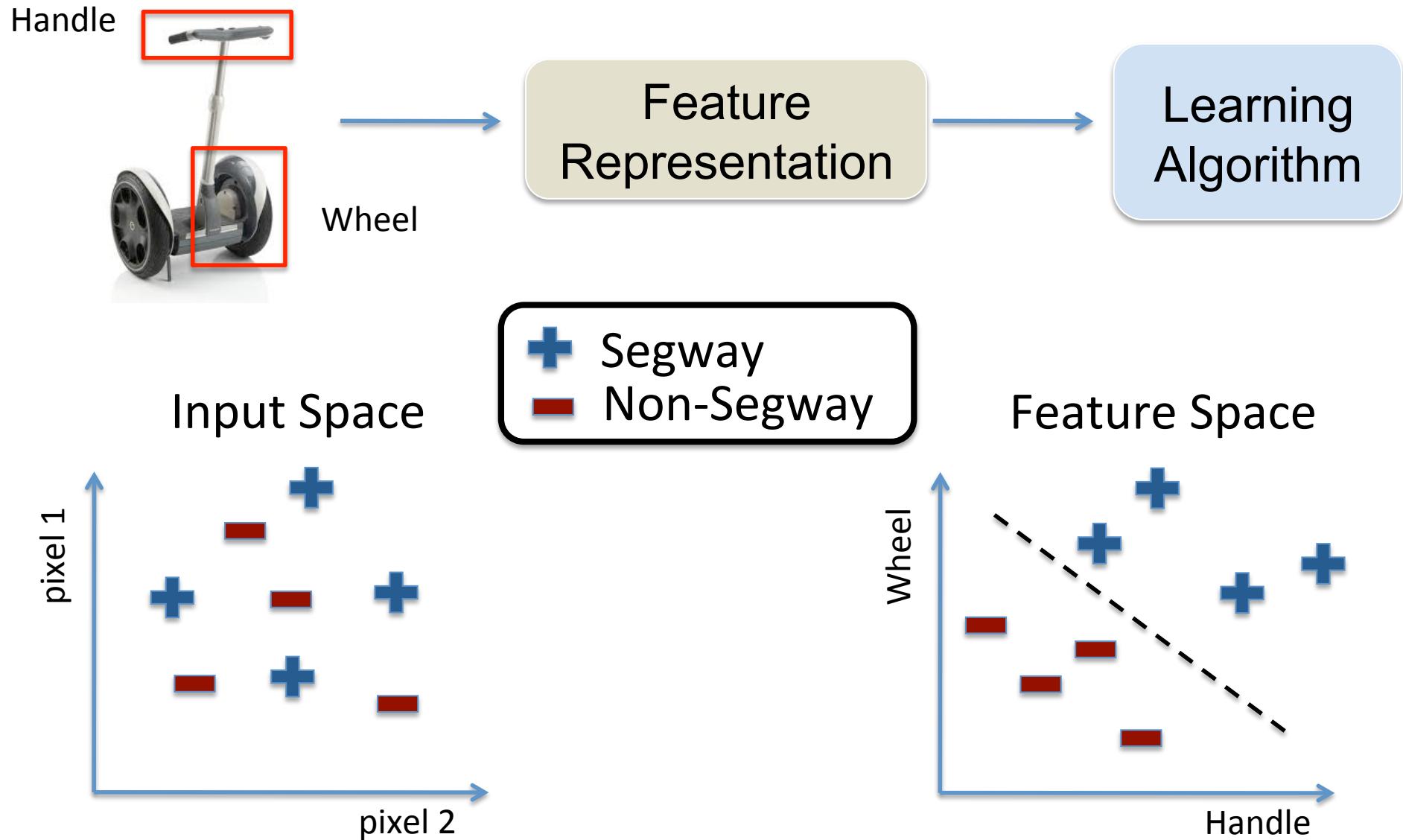
# Talk Roadmap

- Basic Building Blocks:
  - Sparse Coding
  - Autoencoders
- Deep Generative Models
  - Restricted Boltzmann Machines
  - Deep Boltzmann Machines
  - Helmholtz Machines / Variational Autoencoders
- Generative Adversarial Networks

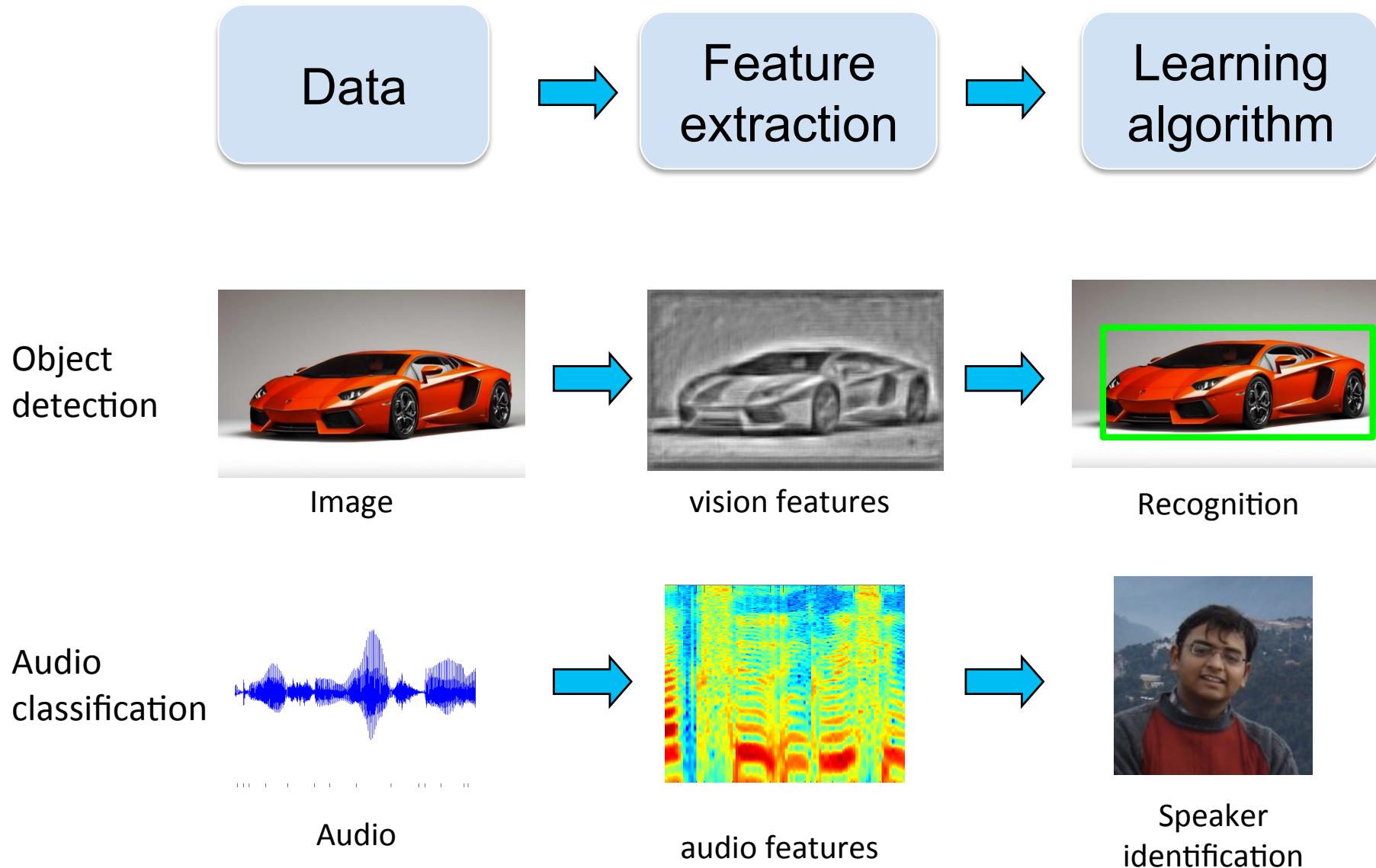
# Learning Feature Representations



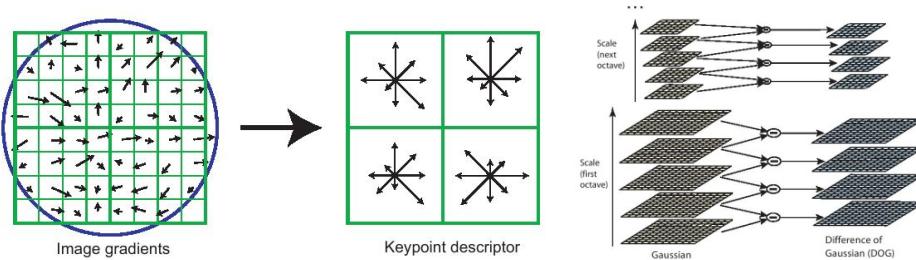
# Learning Feature Representations



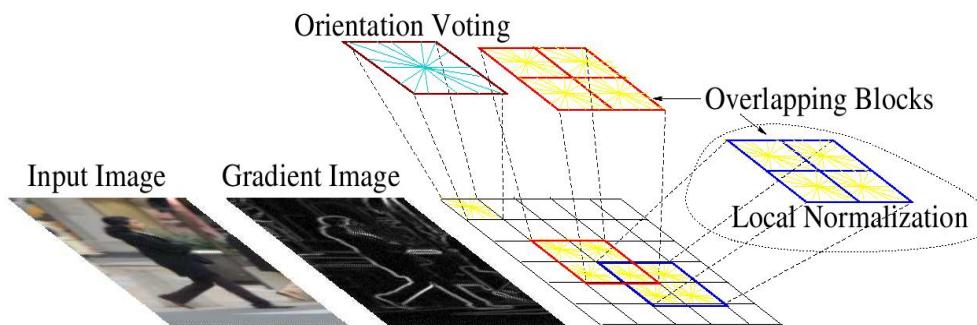
# Traditional Approaches



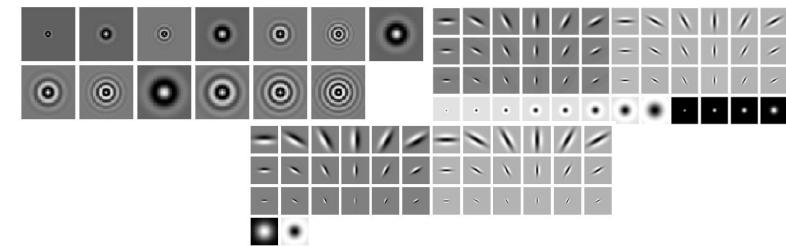
# Computer Vision Features



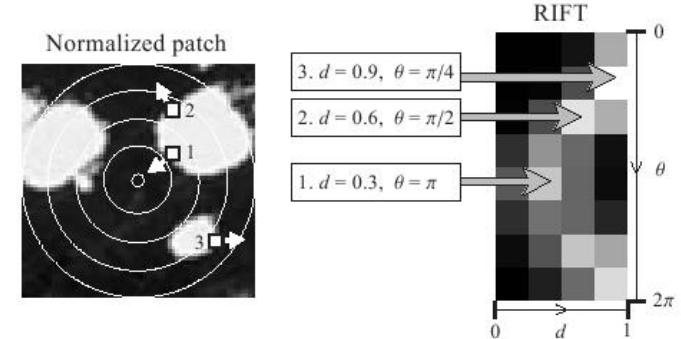
SIFT



HoG

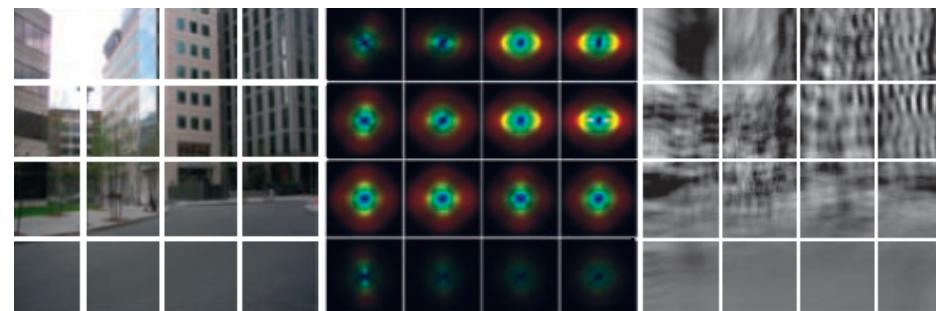


Textons

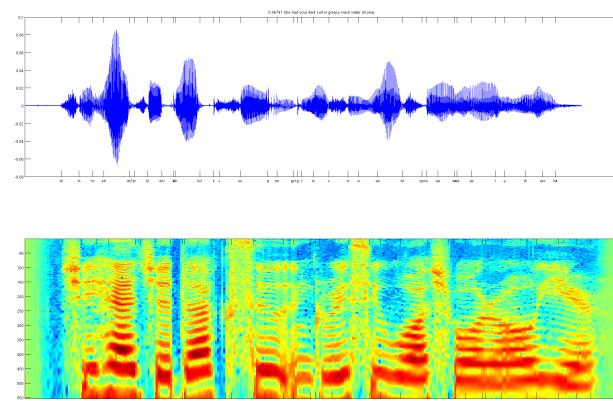


RIFT

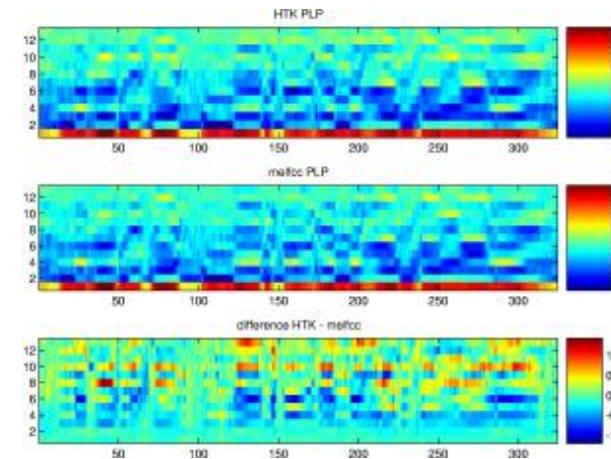
GIST



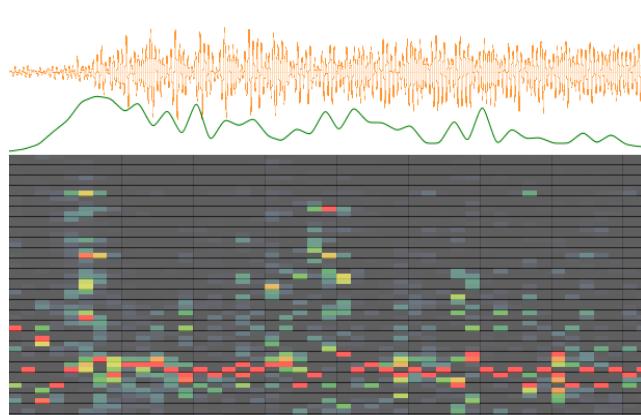
# Audio Features



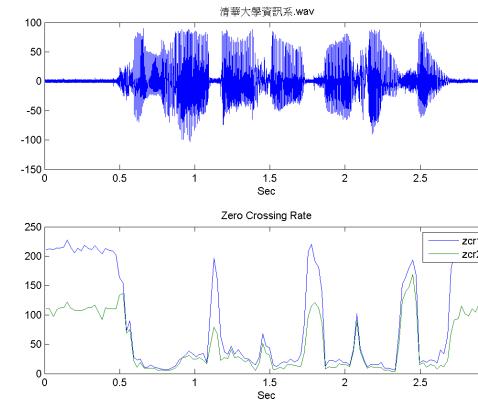
Spectrogram



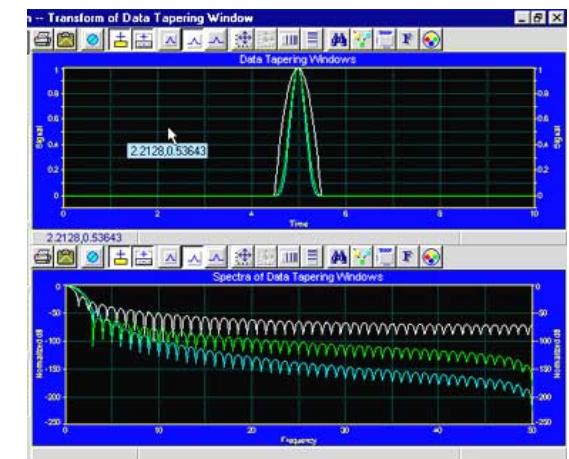
MFCC



Flux

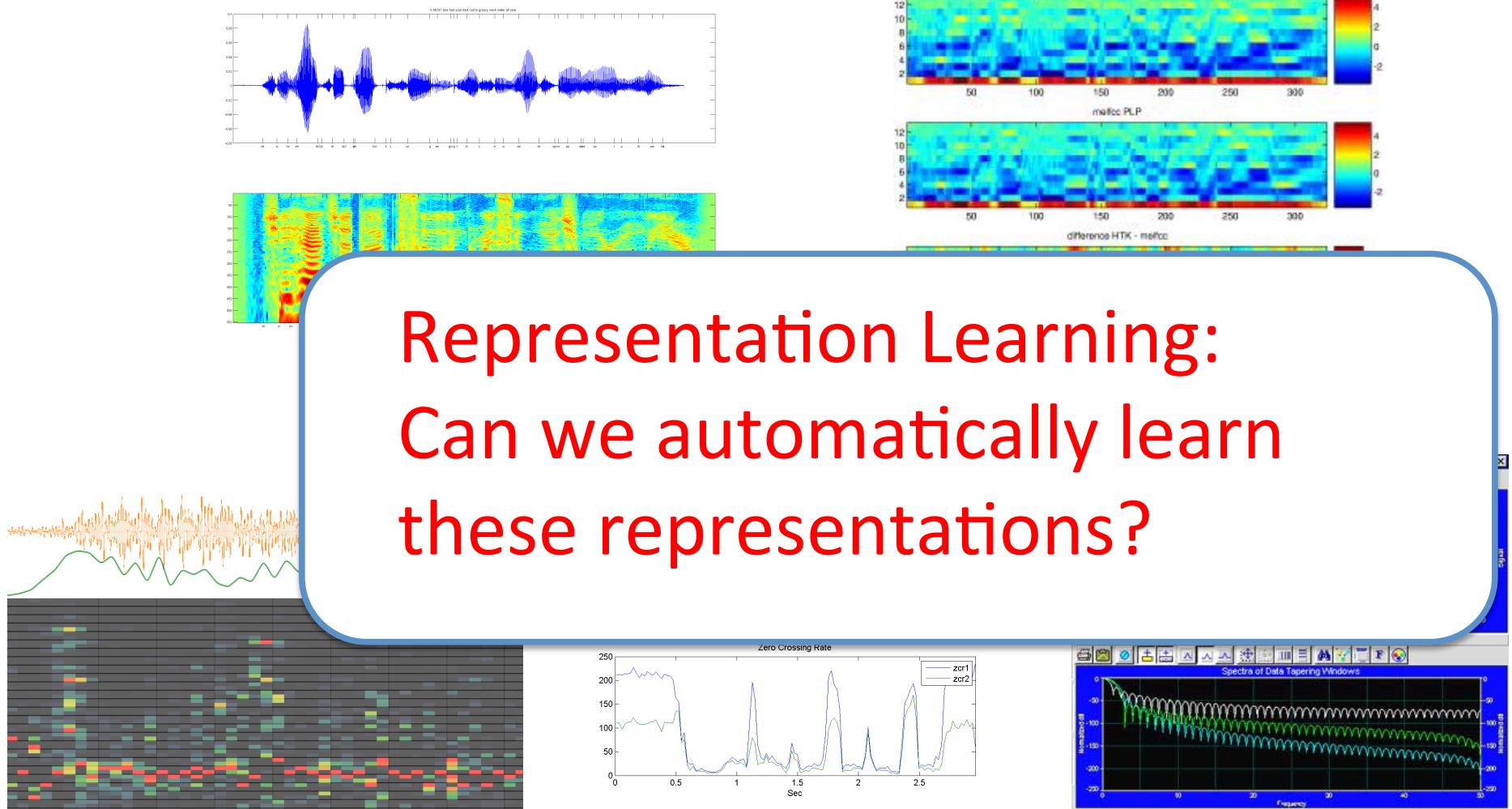


ZCR



Rolloff

# Audio Features



Flux

ZCR

Rolloff

# Sparse Coding

- Sparse coding (Olshausen & Field, 1996). Originally developed to explain early visual processing in the brain (edge detection).
- **Objective:** Given a set of input data vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , learn a dictionary of bases  $\{\phi_1, \phi_2, \dots, \phi_K\}$ , such that:

$$\mathbf{x}_n = \sum_{k=1}^K a_{nk} \phi_k,$$

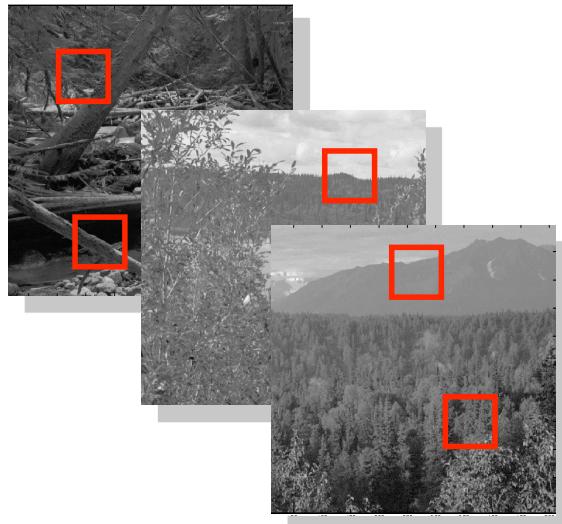
Sparse: mostly zeros



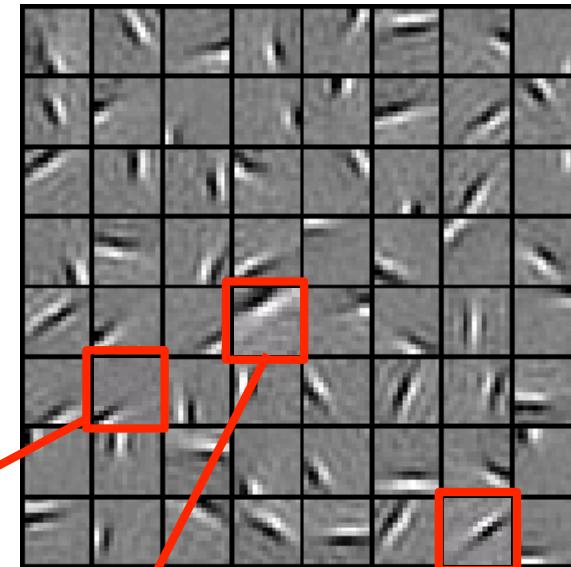
- Each data vector is represented as a sparse linear combination of bases.

# Sparse Coding

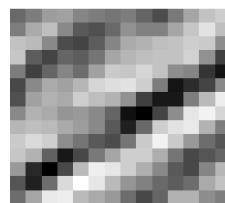
Natural Images



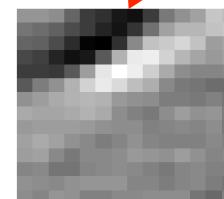
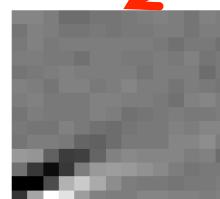
Learned bases: “Edges”



New example



$$x = 0.8 * \phi_{36} + 0.3 * \phi_{42} + 0.5 * \phi_{65}$$



[0, 0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

# Sparse Coding: Training

- Input image patches:  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^D$
- Learn dictionary of bases:  $\phi_1, \phi_2, \dots, \phi_K \in \mathbb{R}^D$

$$\min_{\mathbf{a}, \phi} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \phi_k \right\|_2^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |a_{nk}|$$


Reconstruction error      Sparsity penalty

- Alternating Optimization:
  1. Fix dictionary of bases  $\phi_1, \phi_2, \dots, \phi_K$  and solve for activations  $\mathbf{a}$  (a standard Lasso problem).
  2. Fix activations  $\mathbf{a}$ , optimize the dictionary of bases (convex QP problem).

# Sparse Coding: Testing Time

- Input: a new image patch  $\mathbf{x}^*$ , and  $K$  learned bases  $\phi_1, \phi_2, \dots, \phi_K$
- Output: sparse representation  $\mathbf{a}$  of an image patch  $\mathbf{x}^*$ .

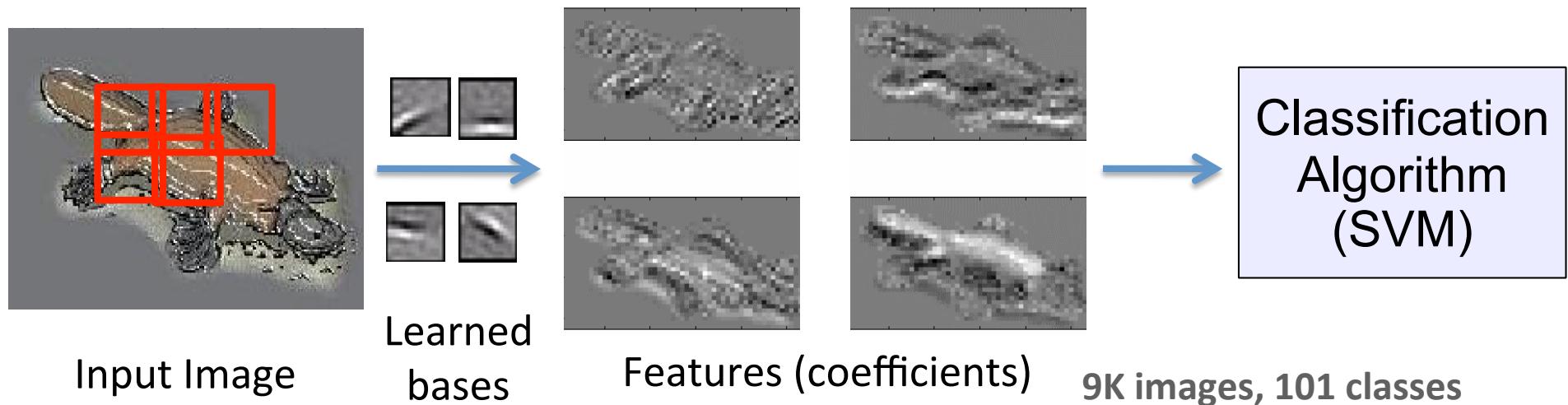
$$\min_{\mathbf{a}} \left\| \mathbf{x}^* - \sum_{k=1}^K a_k \phi_k \right\|_2^2 + \lambda \sum_{k=1}^K |a_k|$$

$$\begin{array}{c} \text{[Image patch]} = 0.8 * \text{[Image patch]} + 0.3 * \text{[Image patch]} + 0.5 * \text{[Image patch]} \\ x^* = 0.8 * \phi_{36} + 0.3 * \phi_{42} + 0.5 * \phi_{65} \end{array}$$

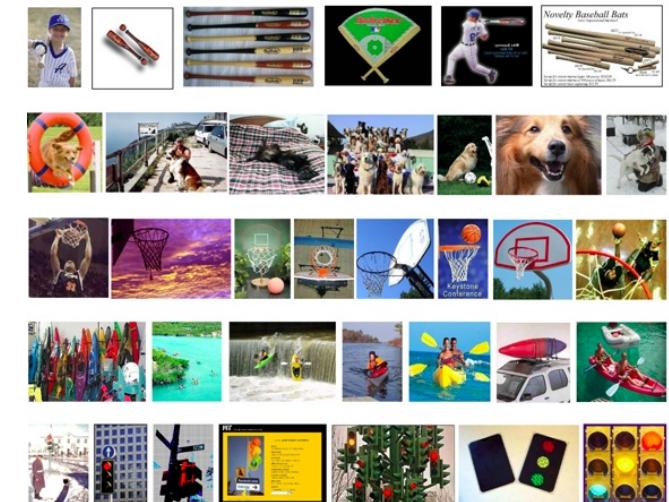
[0, 0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

# Image Classification

Evaluated on Caltech101 object category dataset.

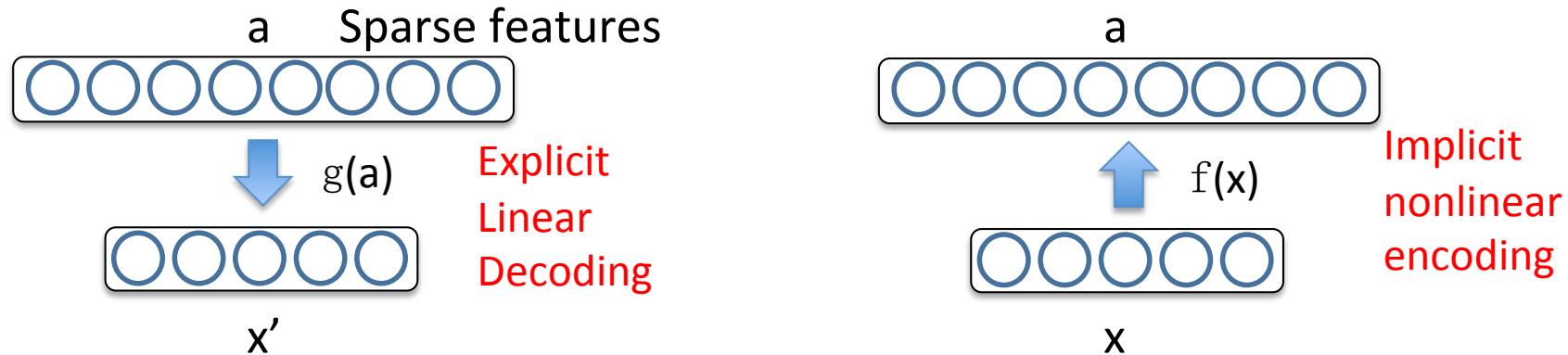


Algorithm	Accuracy
Baseline (Fei-Fei et al., 2004)	16%
PCA	37%
<b>Sparse Coding</b>	<b>47%</b>



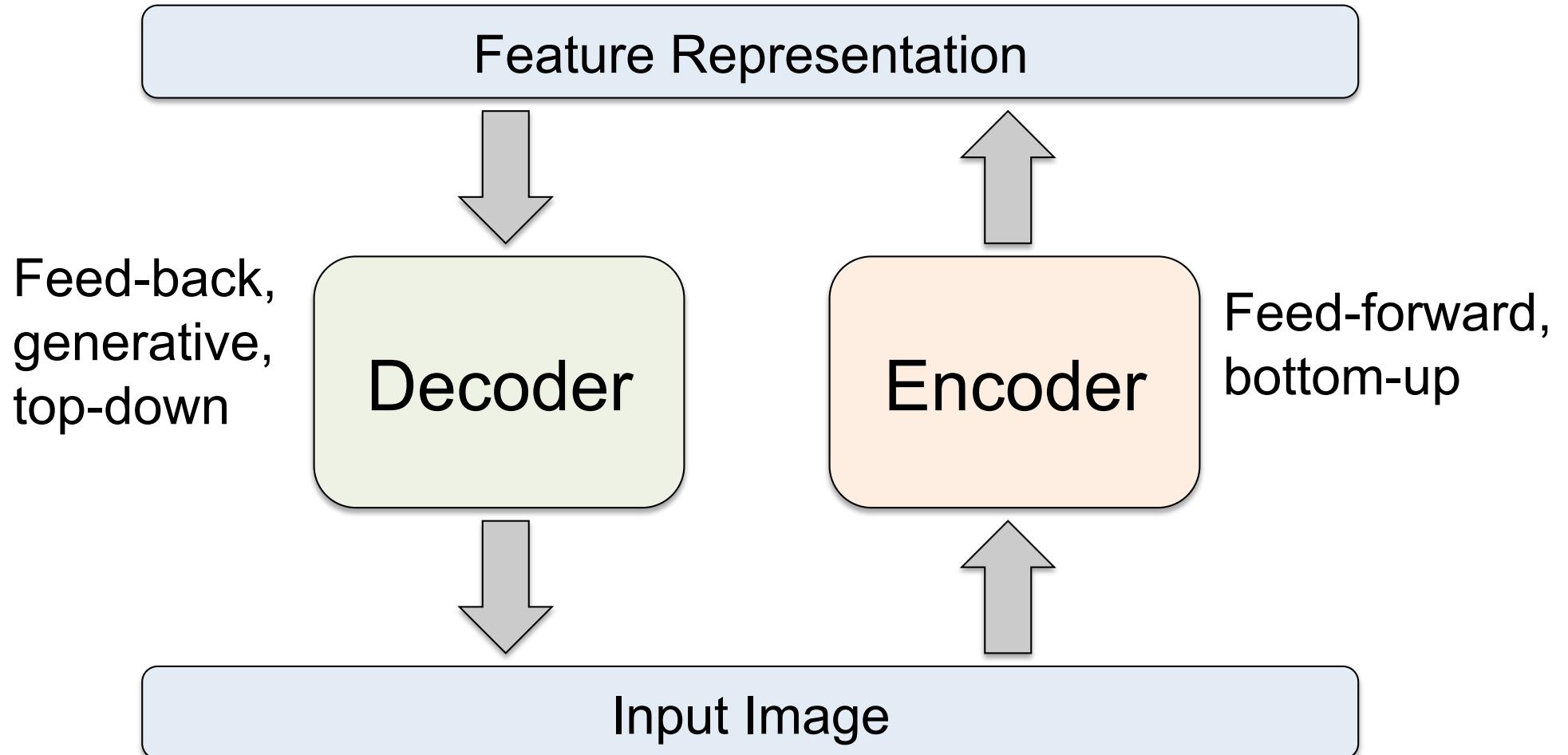
# Interpreting Sparse Coding

$$\min_{\mathbf{a}, \phi} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \phi_k \right\|_2^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |a_{nk}|$$



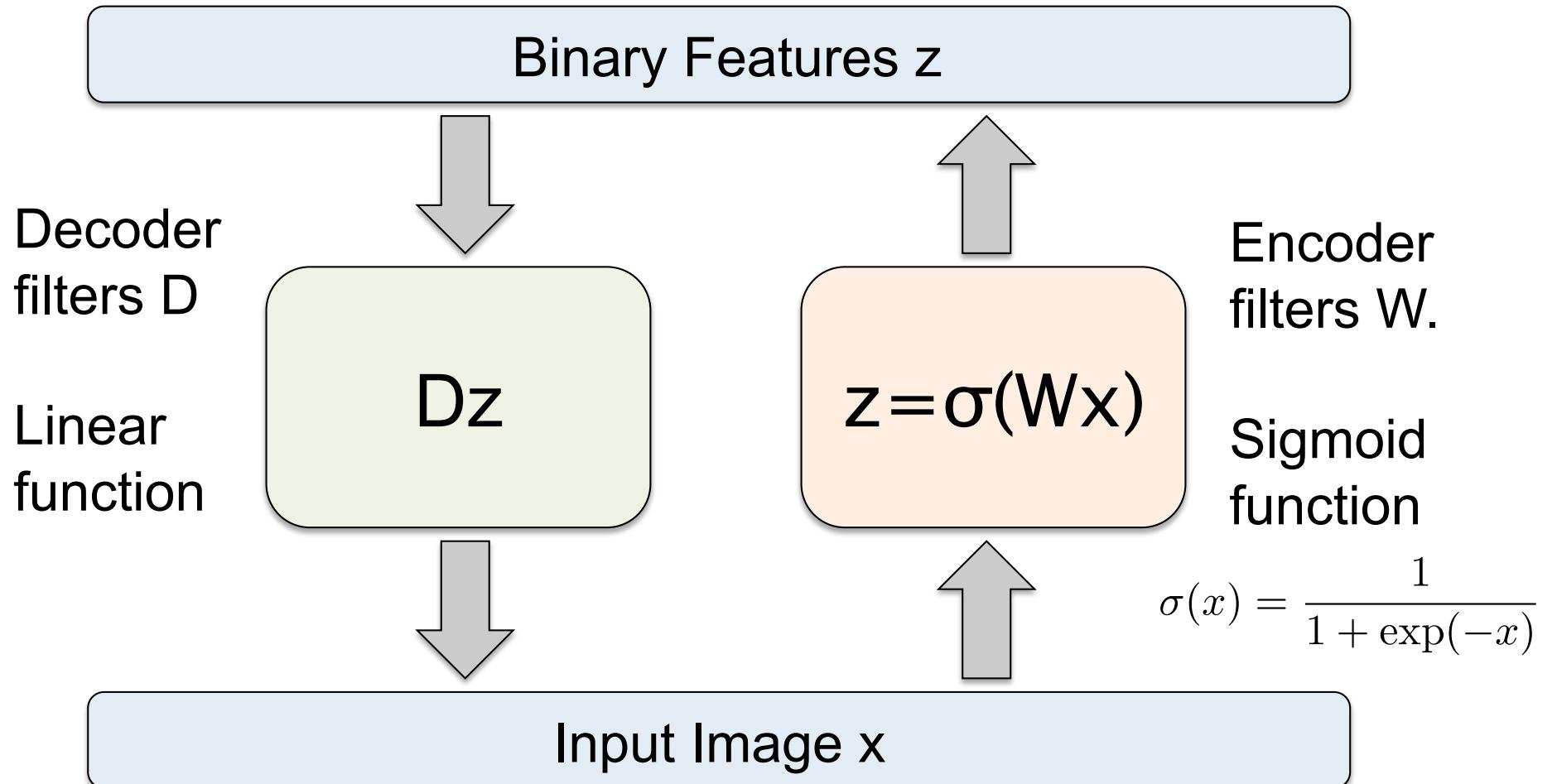
- Sparse, over-complete representation  $\mathbf{a}$ .
- Encoding  $\mathbf{a} = f(\mathbf{x})$  is implicit and nonlinear function of  $\mathbf{x}$ .
- Reconstruction (or decoding)  $\mathbf{x}' = g(\mathbf{a})$  is linear and explicit.

# Autoencoder

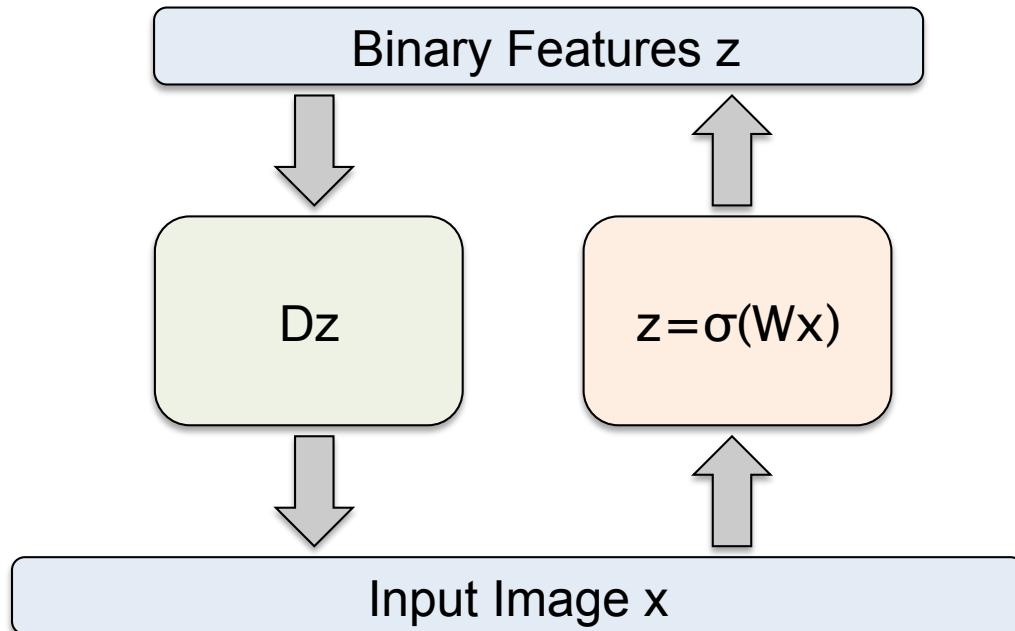


- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity.

# Autoencoder



# Autoencoder



- An autoencoder with D inputs, D outputs, and K hidden units, with K< D.

- Given an input  $x$ , its reconstruction is given by:

$$y_j(\mathbf{x}, W, D) = \underbrace{\sum_{k=1}^K D_{jk} \sigma}_{\text{Decoder}} \left( \underbrace{\sum_{i=1}^D W_{ki} x_i}_{\text{Encoder}} \right), \quad j = 1, \dots, D.$$

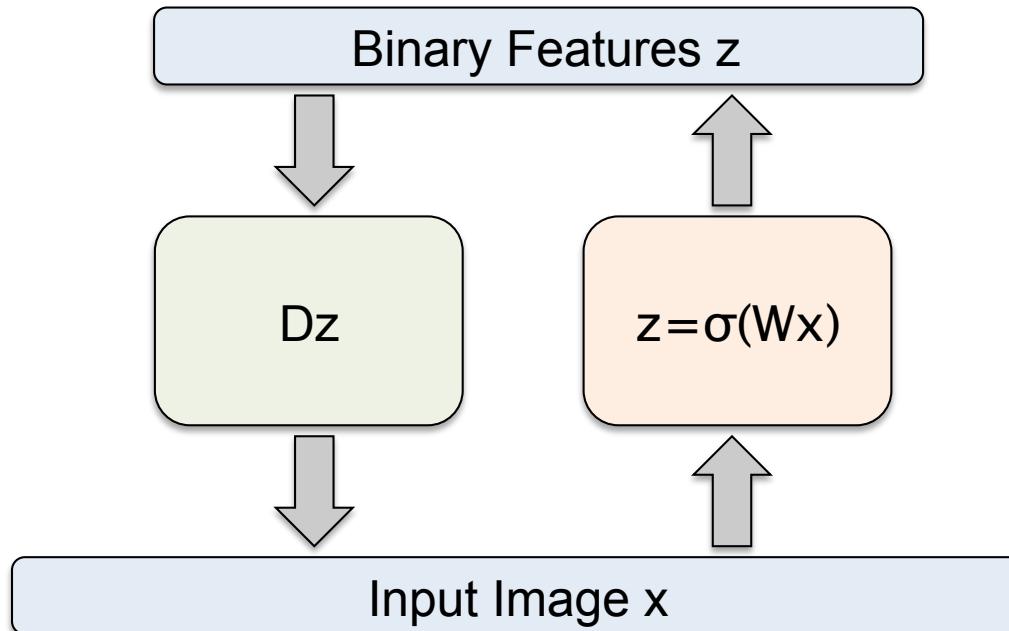
**Decoder**

$$y_j = \sum_{k=1}^K D_{jk} z_k$$

**Encoder**

$$z_k = \sigma \left( \sum_{i=1}^D W_{ki} x_i \right)$$

# Autoencoder

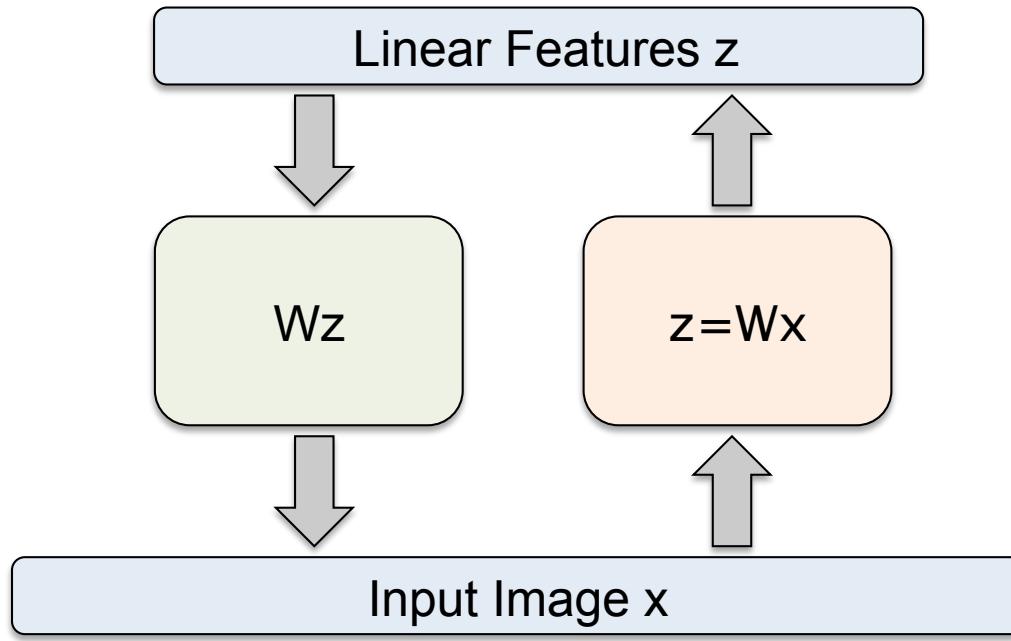


- An autoencoder with D inputs, D outputs, and K hidden units, with K< D.

- We can determine the network parameters W and D by minimizing the reconstruction error:

$$E(W, D) = \frac{1}{2} \sum_{n=1}^N \|y(\mathbf{x}_n, W, D) - \mathbf{x}_n\|^2.$$

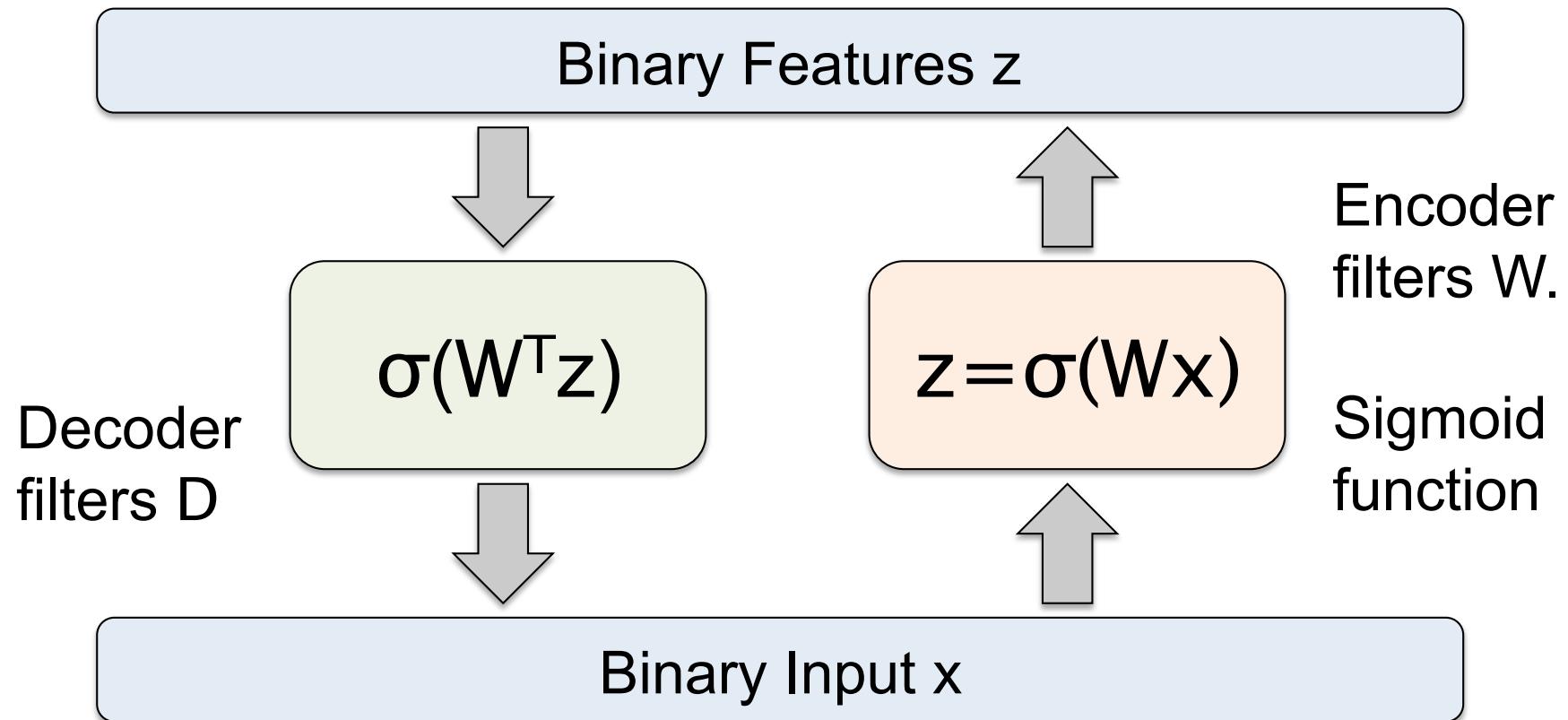
# Autoencoder



- If the hidden and output layers are linear, it will learn hidden units that are a linear function of the data and minimize the squared error.
- The K hidden units will span the same space as the first k principal components. The weight vectors may not be orthogonal.

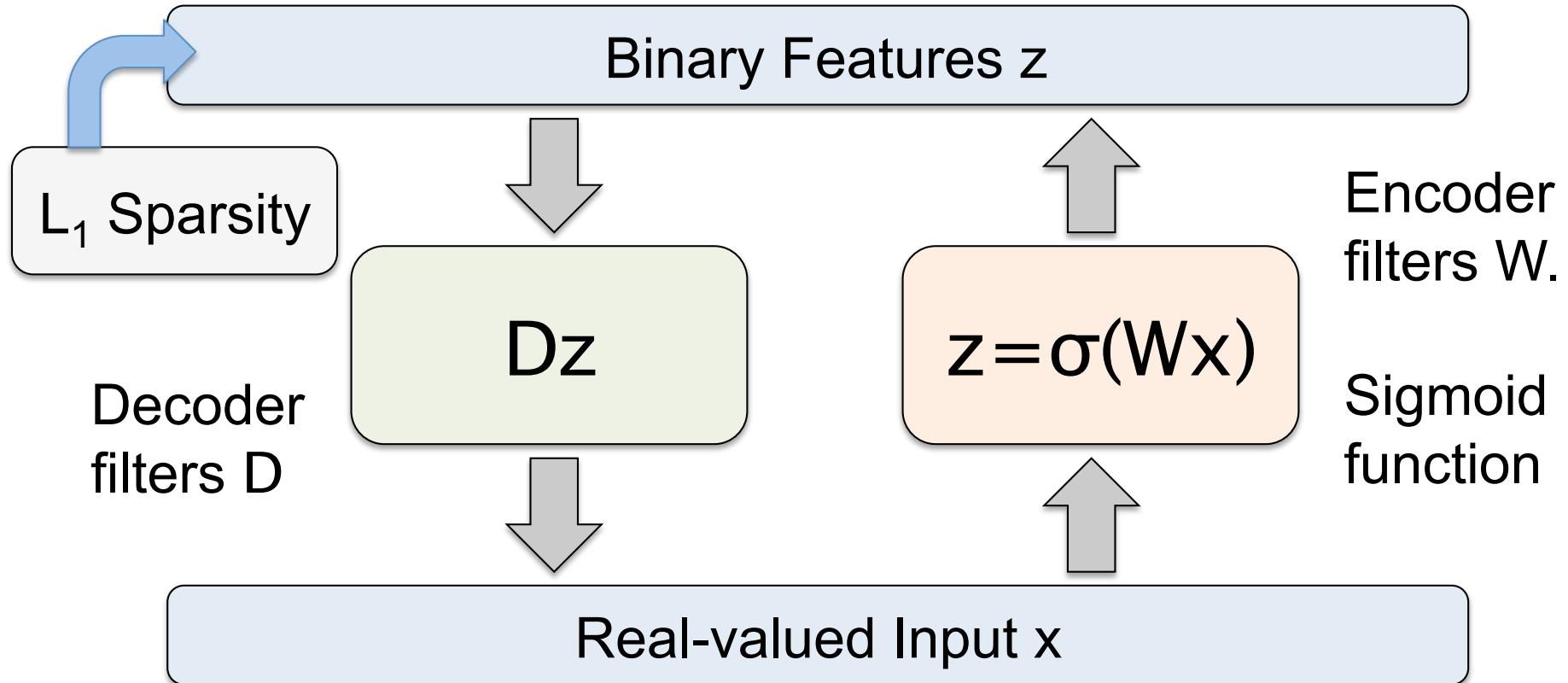
- With nonlinear hidden units, we have a nonlinear generalization of PCA.

# Another Autoencoder Model



- Need additional constraints to avoid learning an identity.
- Relates to Restricted Boltzmann Machines (later).

# Predictive Sparse Decomposition



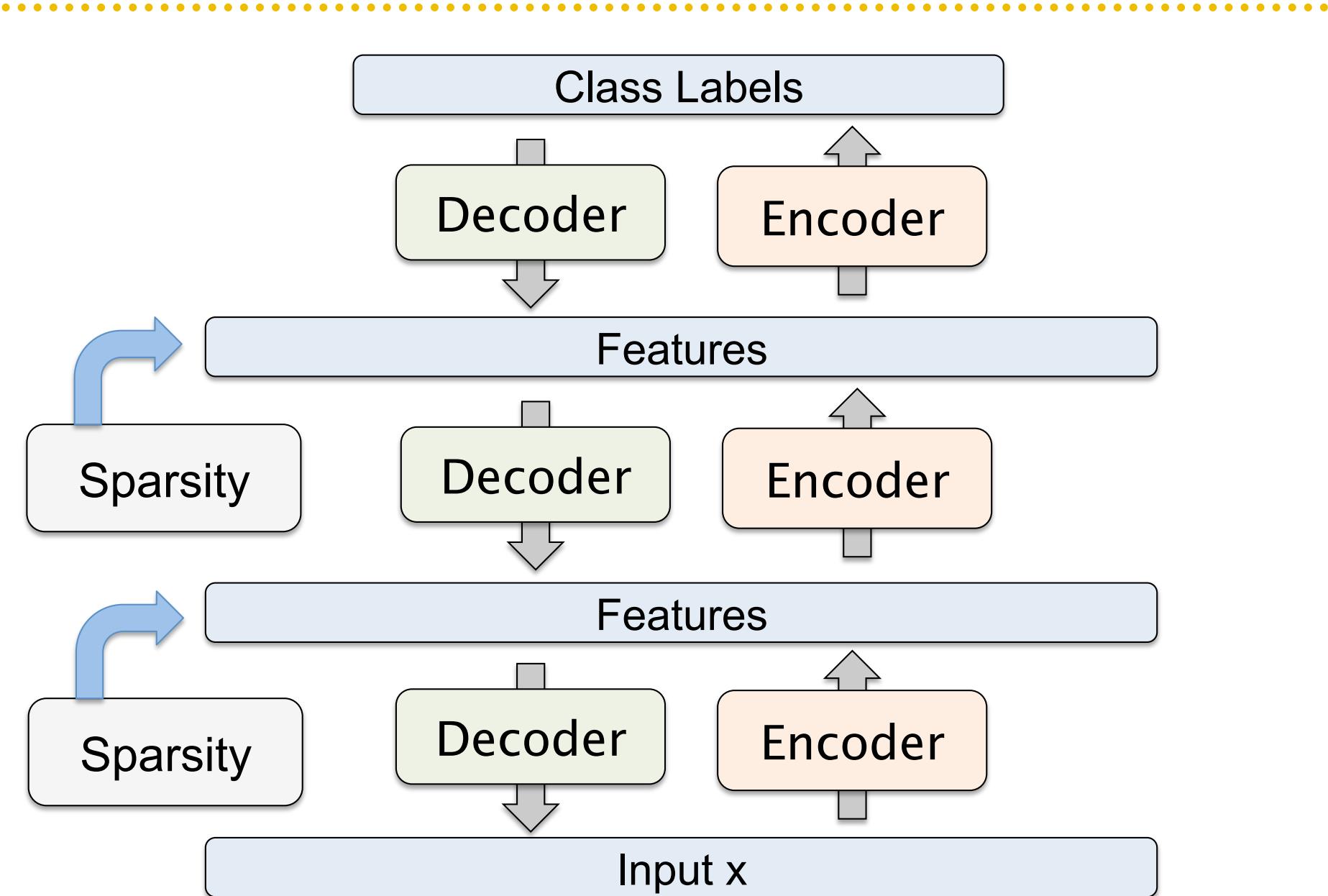
At training time

$$\min_{D, W, z} \|Dz - x\|_2^2 + \lambda|z|_1 + \|\sigma(Wx) - z\|_2^2$$

Decoder

Encoder

# Stacked Autoencoders



The diagram illustrates a Stacked Autoencoder architecture with three layers of features. Each layer consists of an Encoder (orange box) at the top and a Decoder (green box) below it. A blue curved arrow labeled "Sparsity" points from the bottom layer to the middle layer, indicating that the middle layer's features are generated based on the sparsity constraint of the bottom layer's features. The top layer also includes a "Decoder" box. Above the middle layer, there is a "Features" box and a "Class Labels" box. Arrows indicate the flow of data from the bottom layer up through the middle layer to the top layer, and from the bottom layer up through the middle layer to the "Features" and "Class Labels" boxes.

Class Labels

Decoder

Encoder

Features

Sparsity

Decoder

Encoder

Features

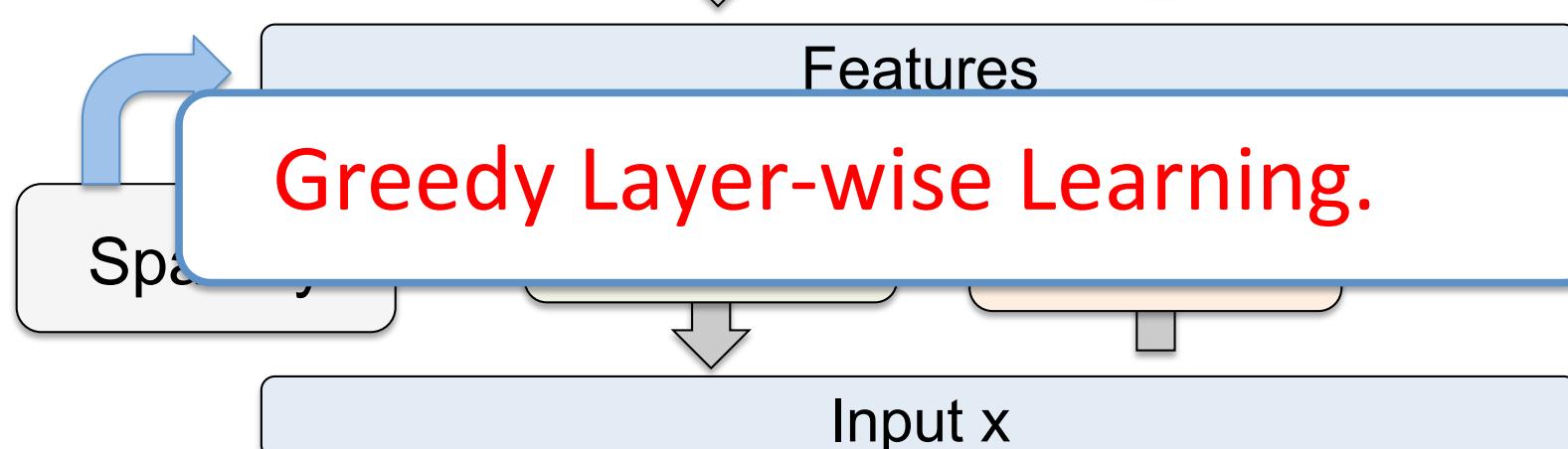
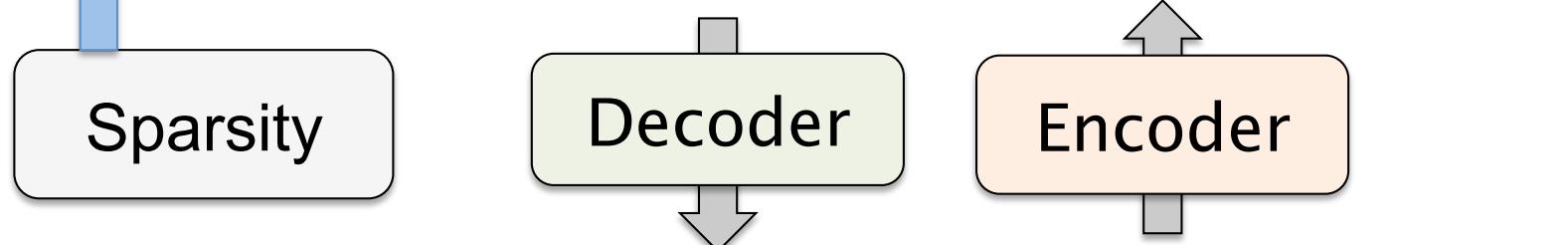
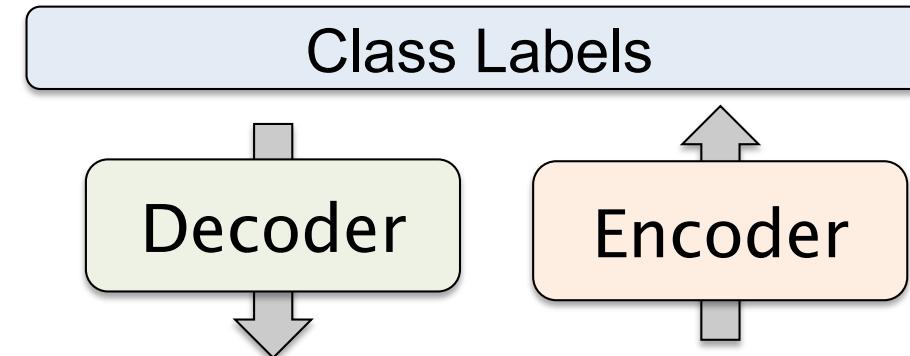
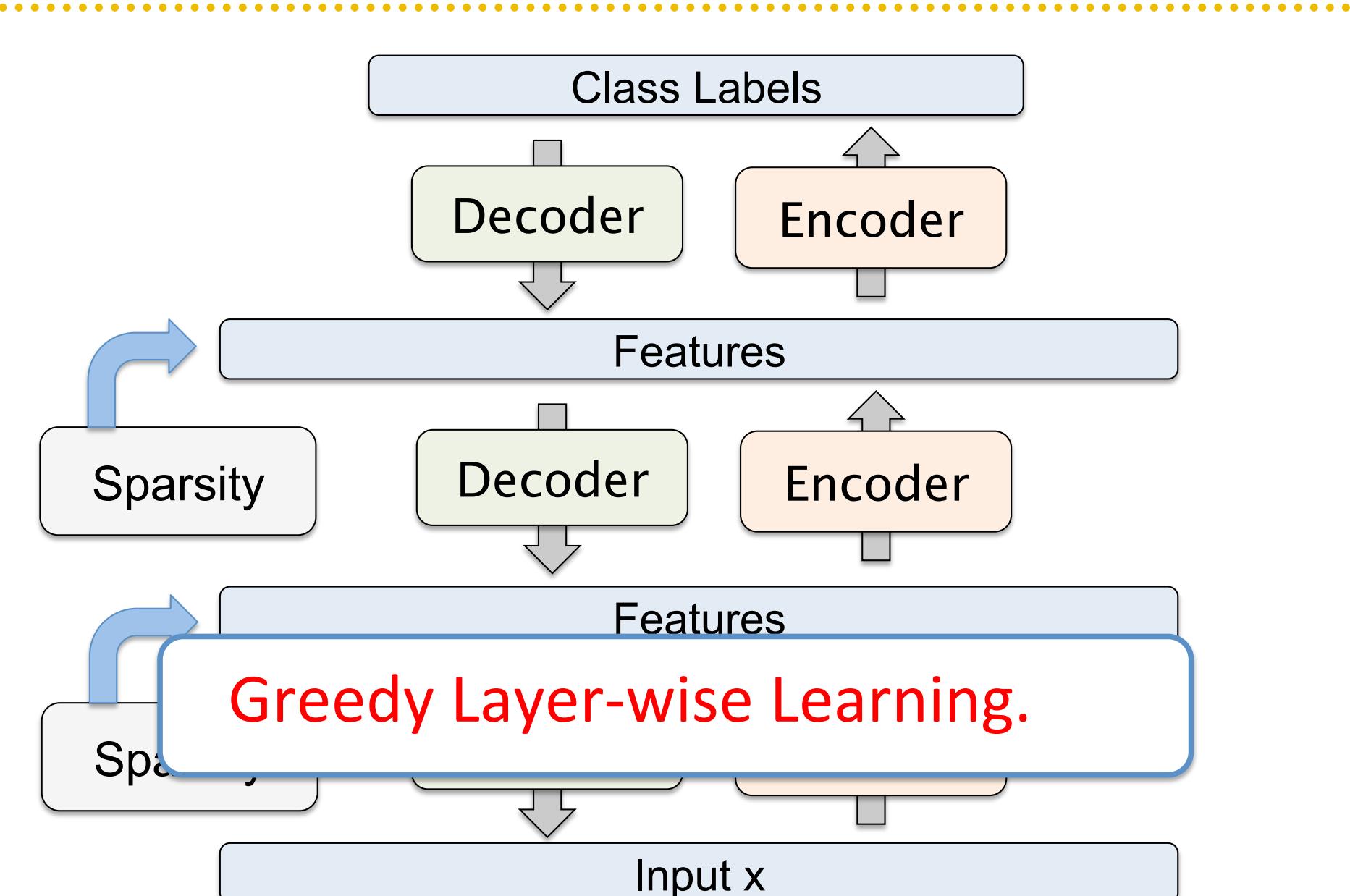
Sparsity

Decoder

Encoder

Input  $x$

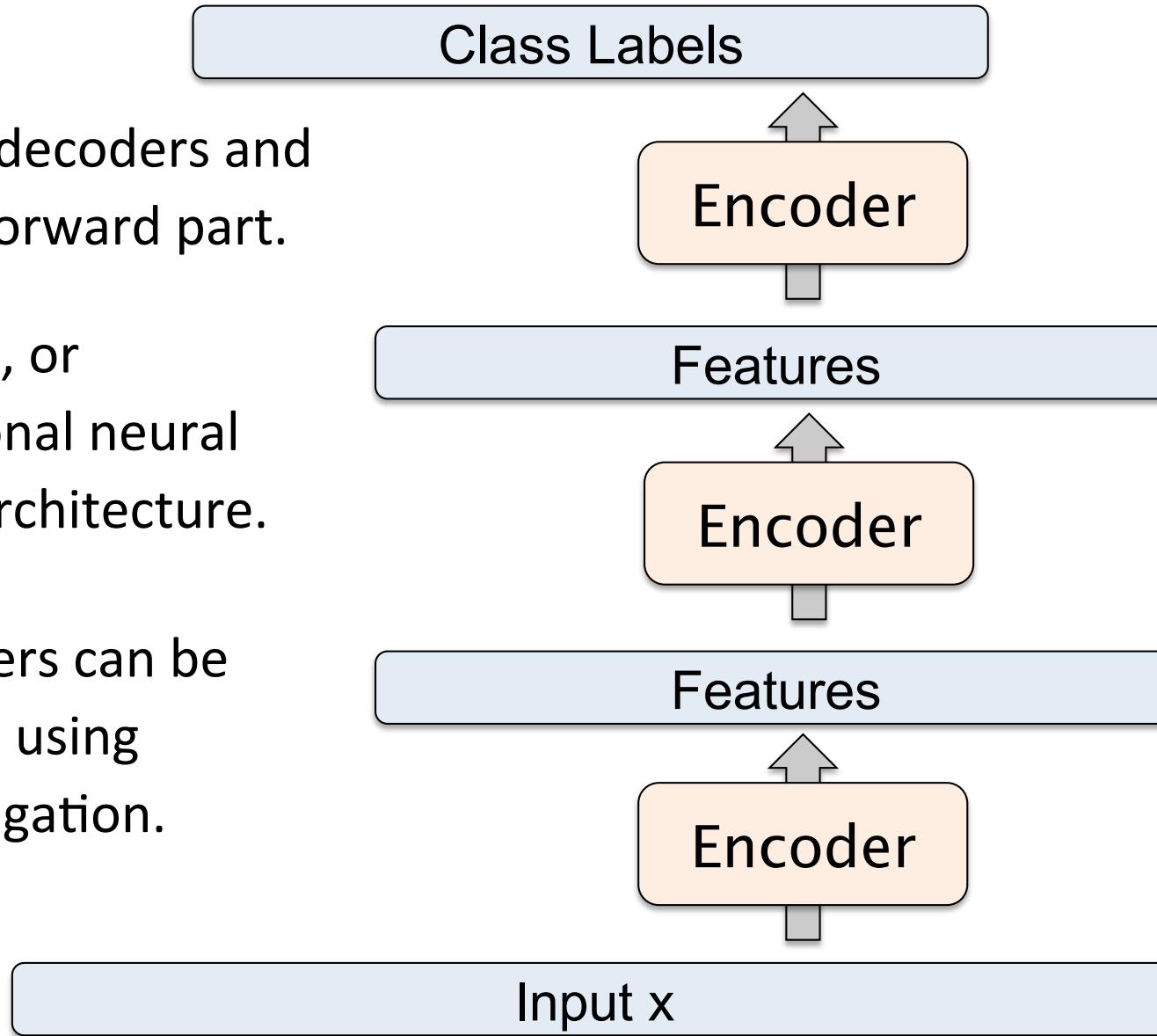
# Stacked Autoencoders



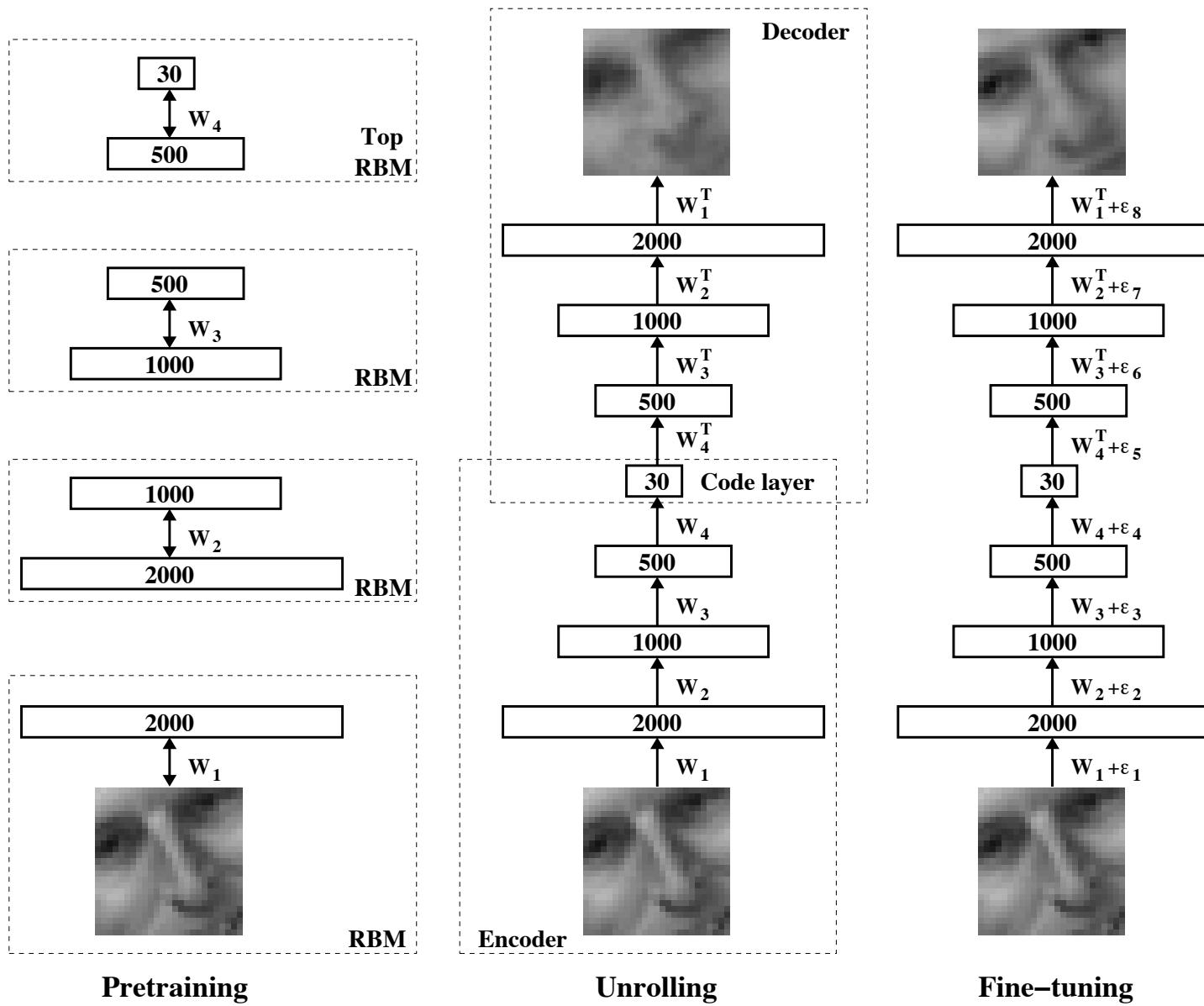
# Stacked Autoencoders

---

- Remove decoders and use feed-forward part.
- Standard, or convolutional neural network architecture.
- Parameters can be fine-tuned using backpropagation.

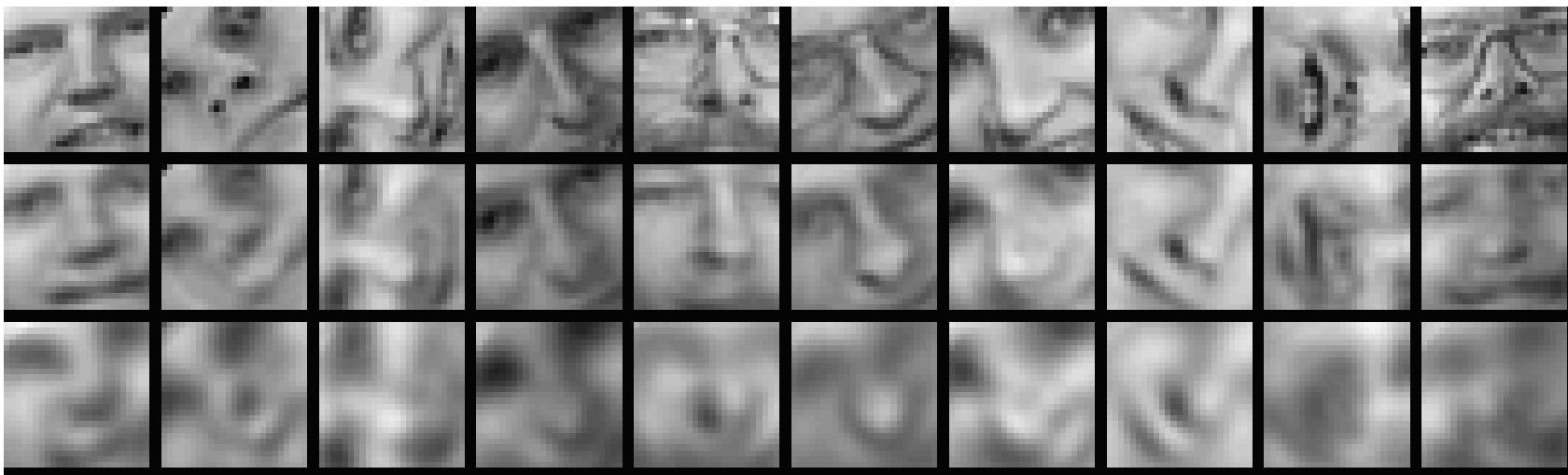


# Deep Autoencoders



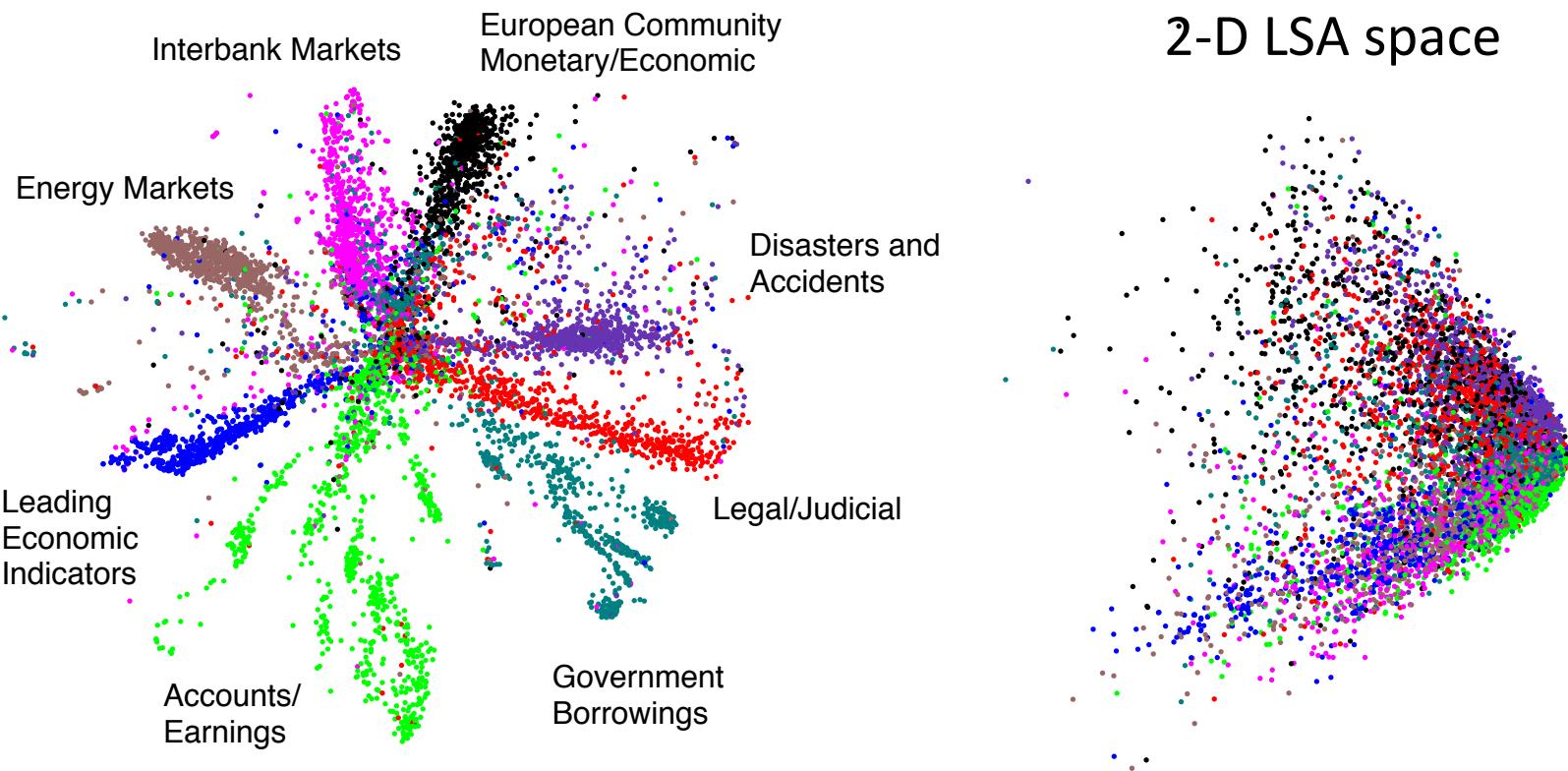
# Deep Autoencoders

- $25 \times 25 - 2000 - 1000 - 500 - 30$  autoencoder to extract 30-D real-valued codes for Olivetti face patches.



- **Top:** Random samples from the test dataset.
- **Middle:** Reconstructions by the 30-dimensional deep autoencoder.
- **Bottom:** Reconstructions by the 30-dimensional PCA.

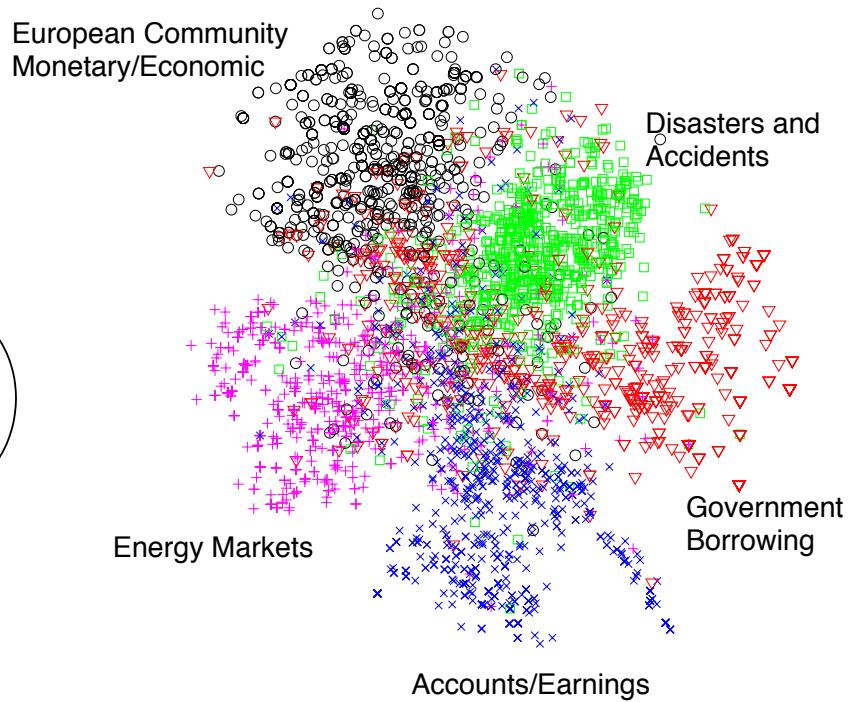
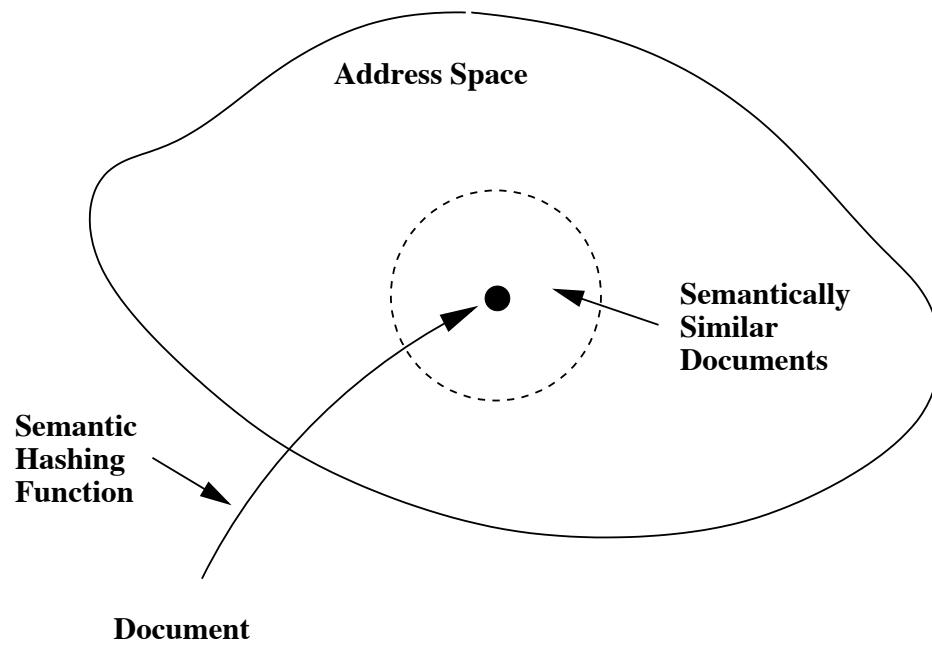
# Information Retrieval



- The Reuters Corpus Volume II contains 804,414 newswire stories (randomly split into **402,207 training** and **402,207 test**).
- “Bag-of-words” representation: each article is represented as a vector containing the counts of the most frequently used 2000 words in the training set.

(Hinton and Salakhutdinov, Science 2006)

# Semantic Hashing

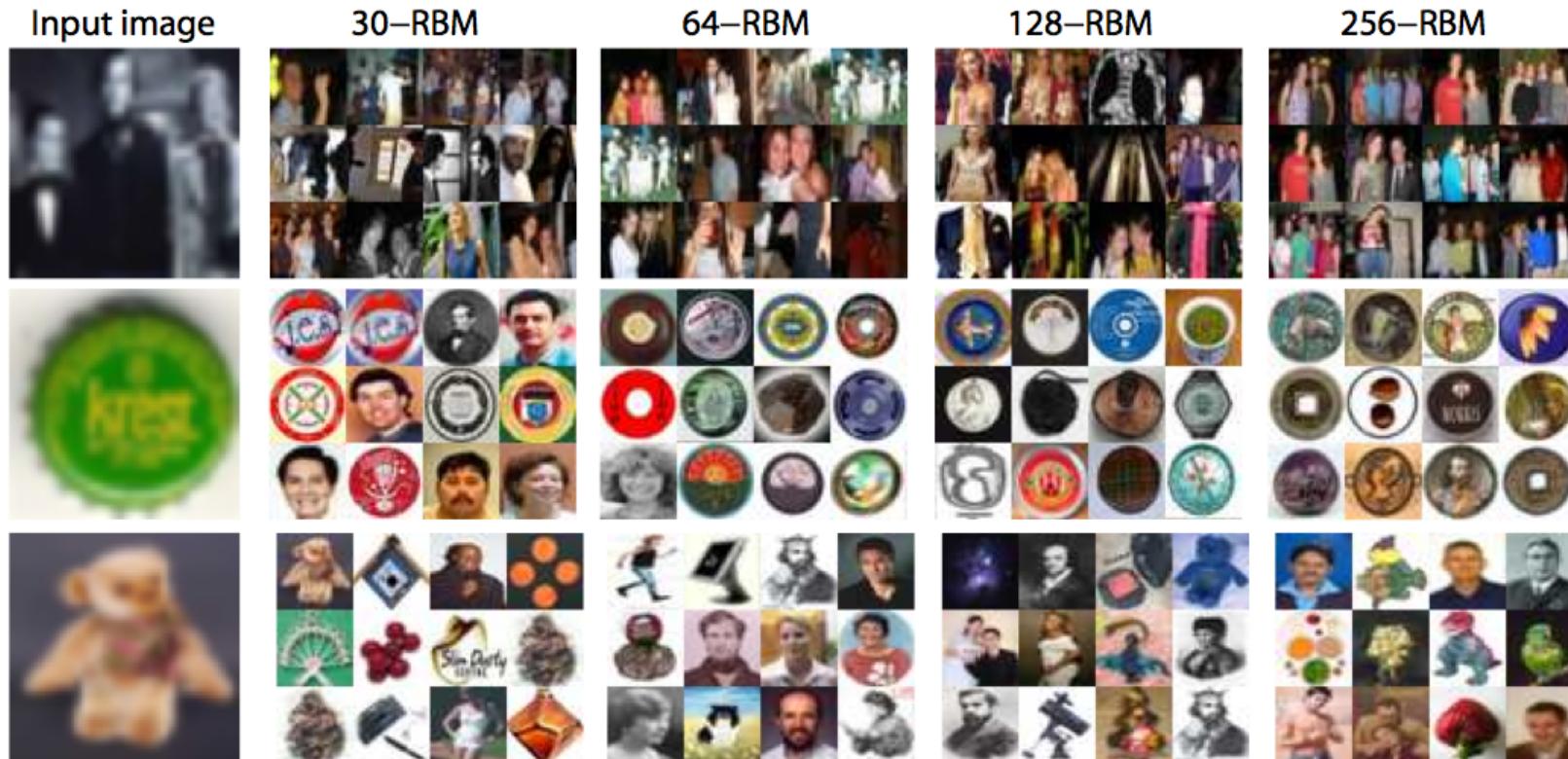


- Learn to map documents into **semantic 20-D binary codes**.
- Retrieve similar documents stored at the nearby addresses **with no search at all**.

(Salakhutdinov and Hinton, SIGIR 2007)

# Searching Large Image Database using Binary Codes

- Map images into binary codes for fast retrieval.



- Small Codes, Torralba, Fergus, Weiss, CVPR 2008
- Spectral Hashing, Y. Weiss, A. Torralba, R. Fergus, NIPS 2008
- Kulis and Darrell, NIPS 2009, Gong and Lazebnik, CVPR 2011
- Norouzi and Fleet, ICML 2011,

# Talk Roadmap

- Basic Building Blocks:

- Sparse Coding
- Autoencoders

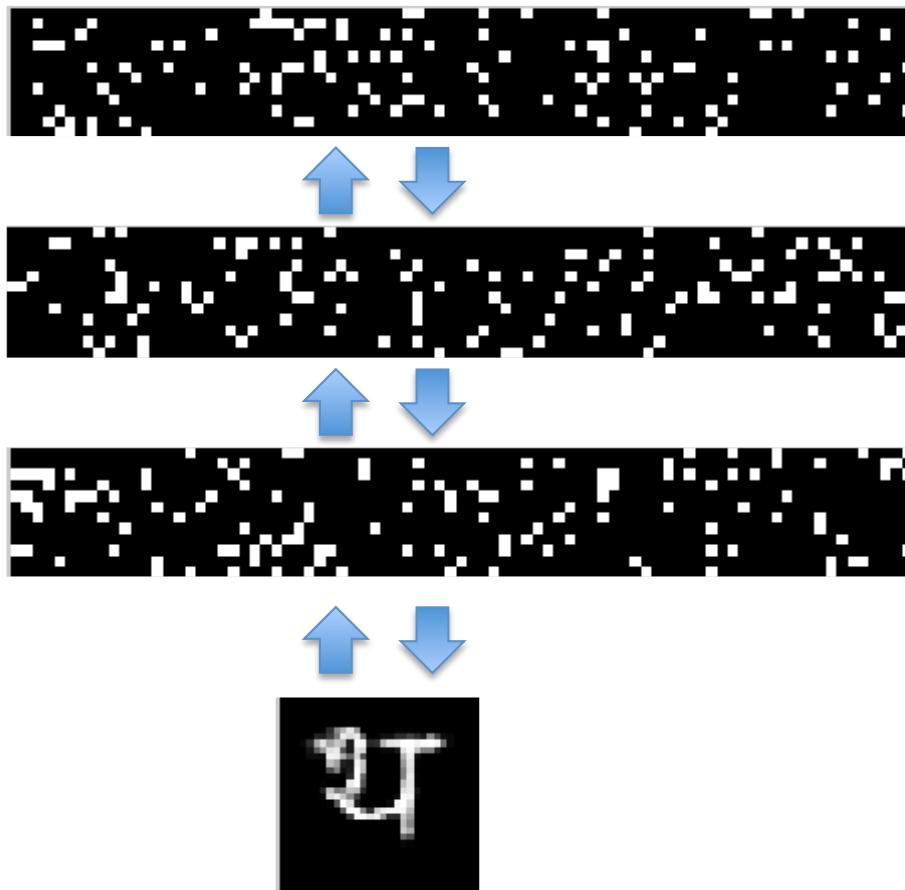
- Deep Generative Models

- Restricted Boltzmann Machines
- Deep Boltzmann Machines
- Helmholtz Machines / Variational Autoencoders

- Generative Adversarial Networks

# Deep Generative Model

Sanskrit



Model P(image)

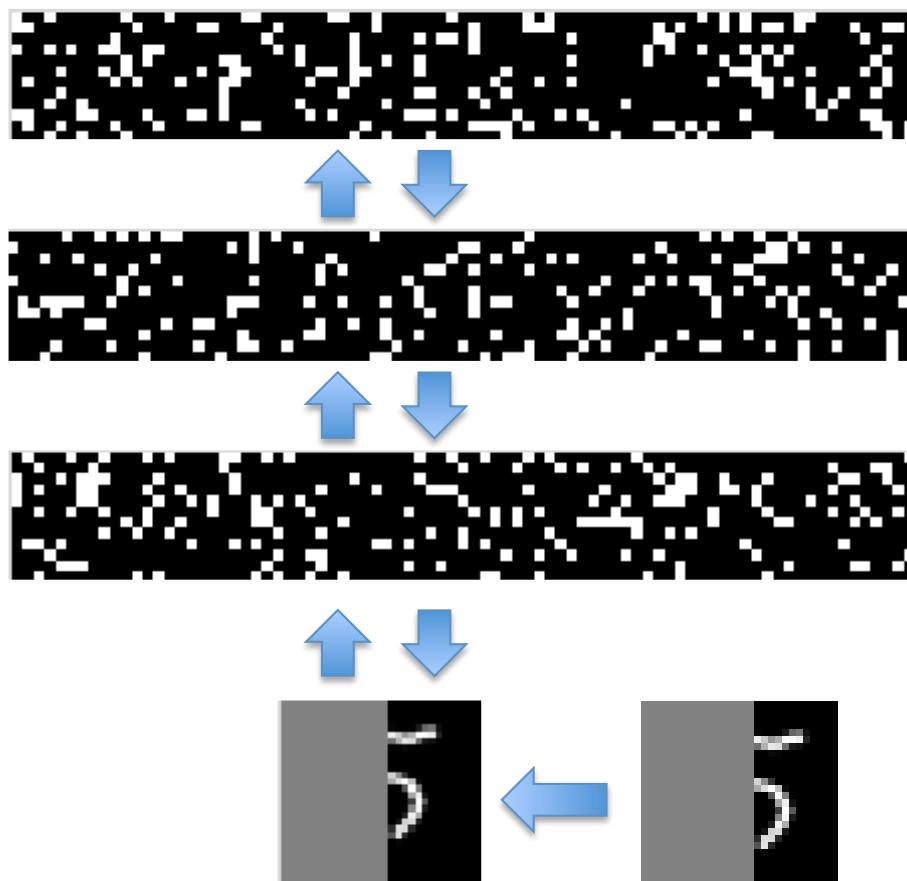
લ ચ થ શ મ છ ણ ણ  
ટ ઢ બ આ લ ઓ ટ ર  
ઝ ઇ લ બ ષ અ ત ઊ  
એ ચ શ ય ક પ ઇ ત્ર

25,000 characters from 50 alphabets around the world.

- 3,000 hidden variables
- 784 observed variables (28 by 28 images)
- About 2 million parameters

Bernoulli Markov Random Field

# Deep Generative Model

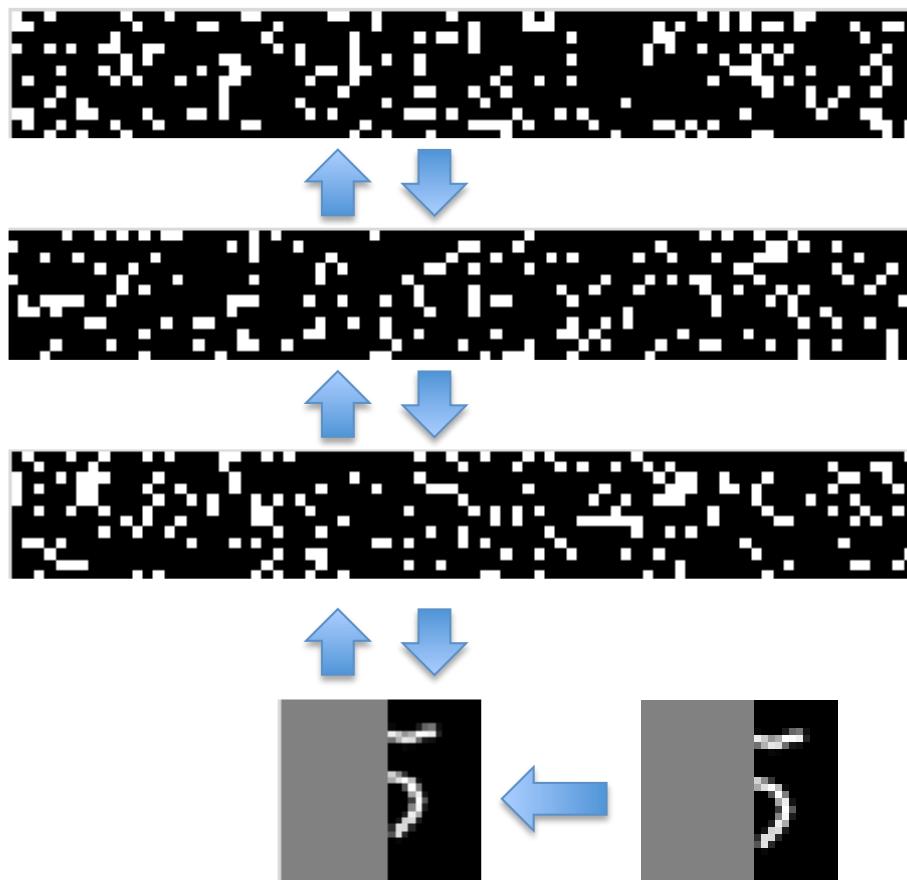


Conditional  
Simulation

$P(\text{image} | \text{partial image})$

Bernoulli Markov Random Field

# Deep Generative Model



$P(\text{image} \mid \text{partial image})$

Conditional  
Simulation

Why so difficult?

28



$2^{28 \times 28}$  possible images!

Bernoulli Markov Random Field

# Fully Observed Models

- Explicitly model conditional probabilities:

$$p_{\text{model}}(\mathbf{x}) = p_{\text{model}}(x_1) \prod_{i=2}^n p_{\text{model}}(x_i \mid x_1, \dots, x_{i-1})$$



Each conditional can be a  
complicated neural network

- A number of successful models, including

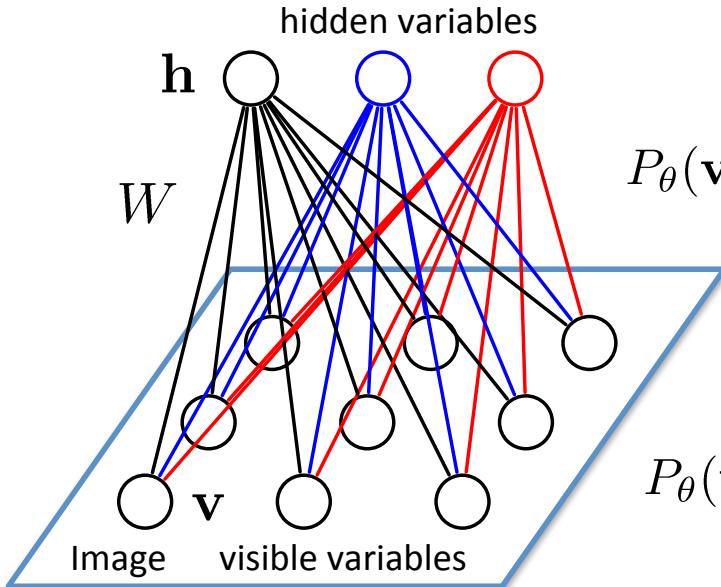
- NADE, RNADE (Larochelle, et.al.  
2001)
- Pixel CNN (van den Ord et. al. 2016)
- Pixel RNN (van den Ord et. al. 2016)



Pixel CNN

# Restricted Boltzmann Machines

Feature Detectors



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

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

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

Pair-wise

Unary

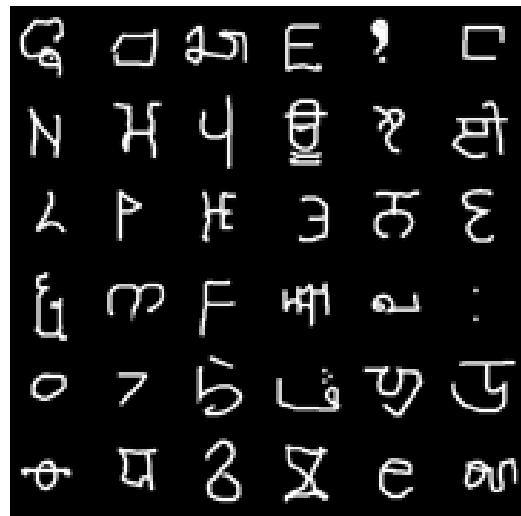
RBM is a Markov Random Field with:

- Stochastic binary visible variables  $\mathbf{v} \in \{0, 1\}^D$ .
- Stochastic binary hidden variables  $\mathbf{h} \in \{0, 1\}^F$ .
- Bipartite connections.

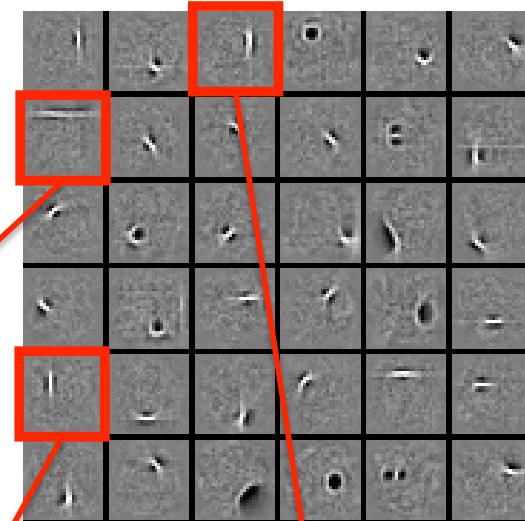
Markov random fields, Boltzmann machines, log-linear models.

# Learning Features

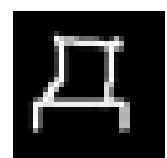
Observed Data  
Subset of 25,000 characters



Learned W: “edges”  
Subset of 1000 features



New Image:  $p(h_7 = 1|v)$

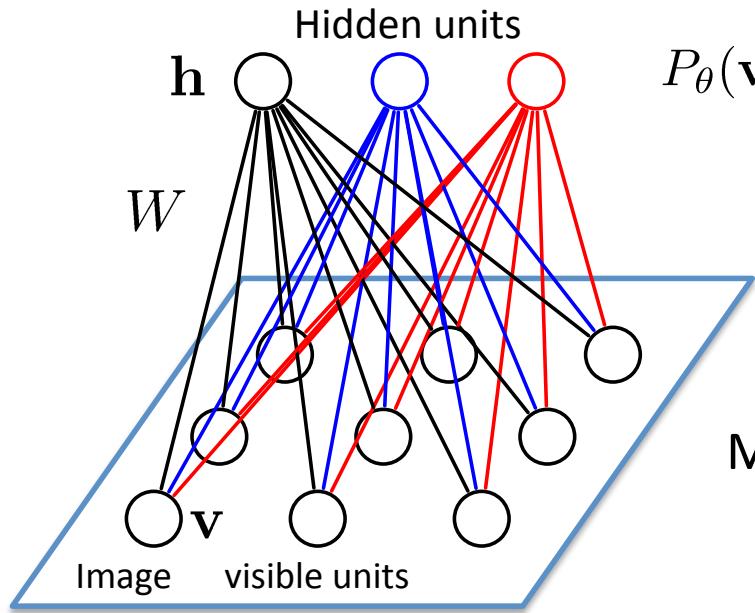


$$= \sigma\left(0.99 \times \begin{array}{c} \text{small image} \\ \text{with vertical stroke} \end{array} + 0.97 \times \begin{array}{c} \text{small image} \\ \text{with horizontal stroke} \end{array} + 0.82 \times \begin{array}{c} \text{small image} \\ \text{with diagonal stroke} \end{array} \dots\right)$$

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

Logistic Function: Suitable for  
modeling binary images

# Model Learning



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

Given a set of *i.i.d.* training examples  $\mathcal{D} = \{\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(N)}\}$ , we want to learn model parameters  $\theta = \{W, a, b\}$ .

Maximize log-likelihood objective:

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

Derivative of the log-likelihood:

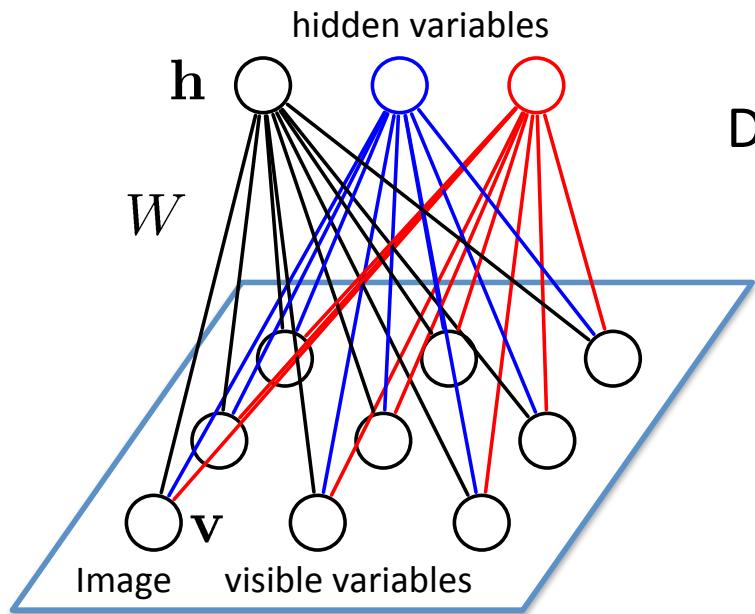
$$\begin{aligned} \frac{\partial L(\theta)}{\partial W_{ij}} &= \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial W_{ij}} \log \left( \sum_{\mathbf{h}} \exp [\mathbf{v}^{(n)\top} W \mathbf{h} + \mathbf{a}^\top \mathbf{h} + \mathbf{b}^\top \mathbf{v}^{(n)}] \right) - \frac{\partial}{\partial W_{ij}} \log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

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

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

Difficult to compute: exponentially many configurations

# Model Learning



Derivative of the log-likelihood:

$$\frac{\partial L(\theta)}{\partial W_{ij}} = \mathbb{E}_{P_{data}}[v_i h_j] - \mathbb{E}_{P_\theta}[v_i h_j]$$

Easy to  
compute exactly

$$\sum_{\mathbf{v}, \mathbf{h}} v_i h_j P_\theta(\mathbf{v}, \mathbf{h})$$

Difficult to compute:  
exponentially many  
configurations.  
Use MCMC

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

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

Approximate maximum likelihood learning

# Approximate Learning

- An approximation to the gradient of the log-likelihood objective:

$$\frac{\partial L(\theta)}{\partial W_{ij}} = \mathbb{E}_{P_{data}}[v_i h_j] - \mathbb{E}_{P_\theta}[v_i h_j]$$

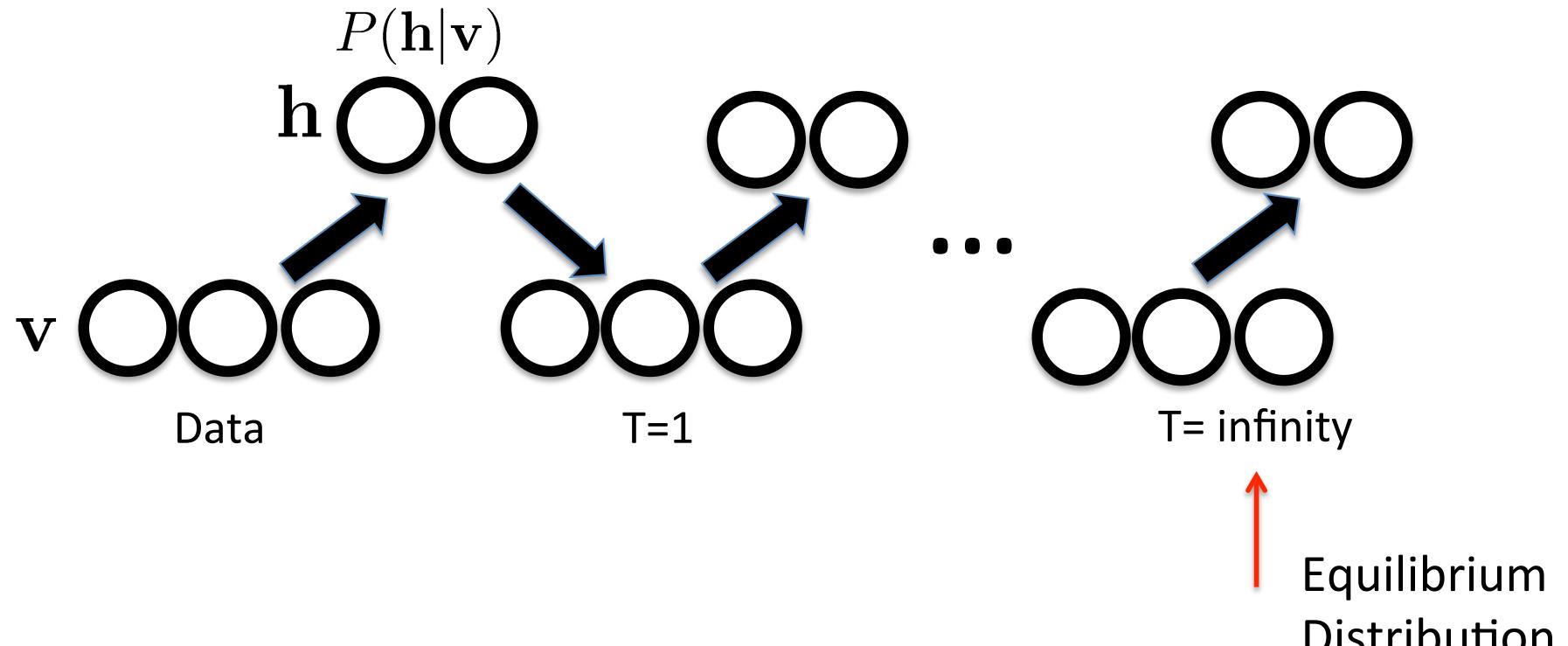
$$\sum_{\mathbf{v}, \mathbf{h}} v_i h_j P_\theta(\mathbf{v}, \mathbf{h})$$

- Replace the average over all possible input configurations by samples.
- Run MCMC chain (Gibbs sampling) starting from the observed examples.

- Initialize  $\mathbf{v}^0 = \mathbf{v}$
- Sample  $\mathbf{h}^0$  from  $P(\mathbf{h} \mid \mathbf{v}^0)$
- For  $t=1:T$ 
  - Sample  $\mathbf{v}^t$  from  $P(\mathbf{v} \mid \mathbf{h}^{t-1})$
  - Sample  $\mathbf{h}^t$  from  $P(\mathbf{h} \mid \mathbf{v}^t)$

# Approximate ML Learning for RBMs

Run Markov chain (alternating Gibbs Sampling):

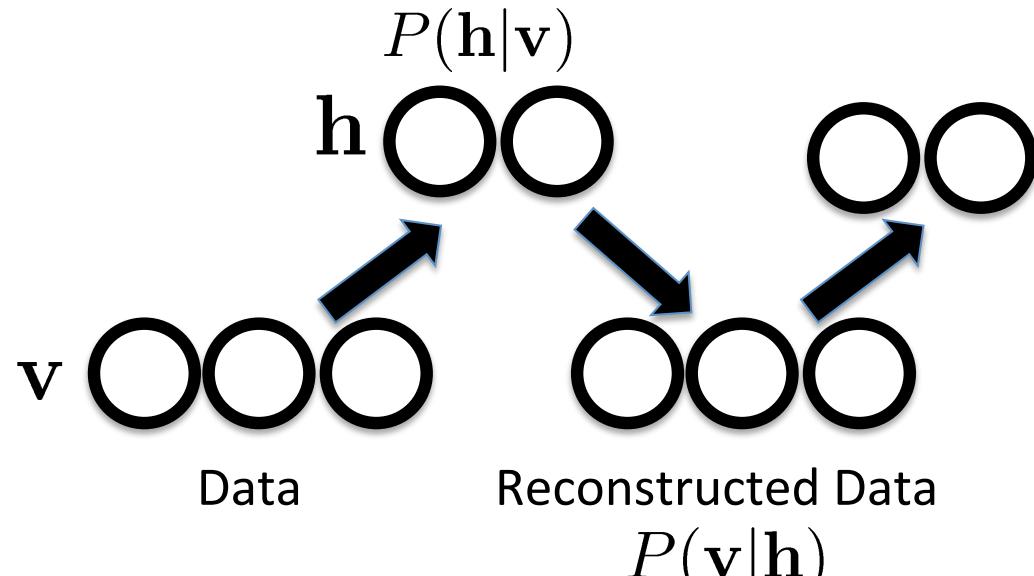


$$P(\mathbf{h}|\mathbf{v}) = \prod_j P(h_j|\mathbf{v}) \quad P(h_j = 1|\mathbf{v}) = \frac{1}{1 + \exp(-\sum_i W_{ij} v_i - a_j)}$$

$$P(\mathbf{v}|\mathbf{h}) = \prod_i P(v_i|\mathbf{h}) \quad P(v_i = 1|\mathbf{h}) = \frac{1}{1 + \exp(-\sum_j W_{ij} h_j - b_i)}$$

# Contrastive Divergence

A quick way to learn RBM:



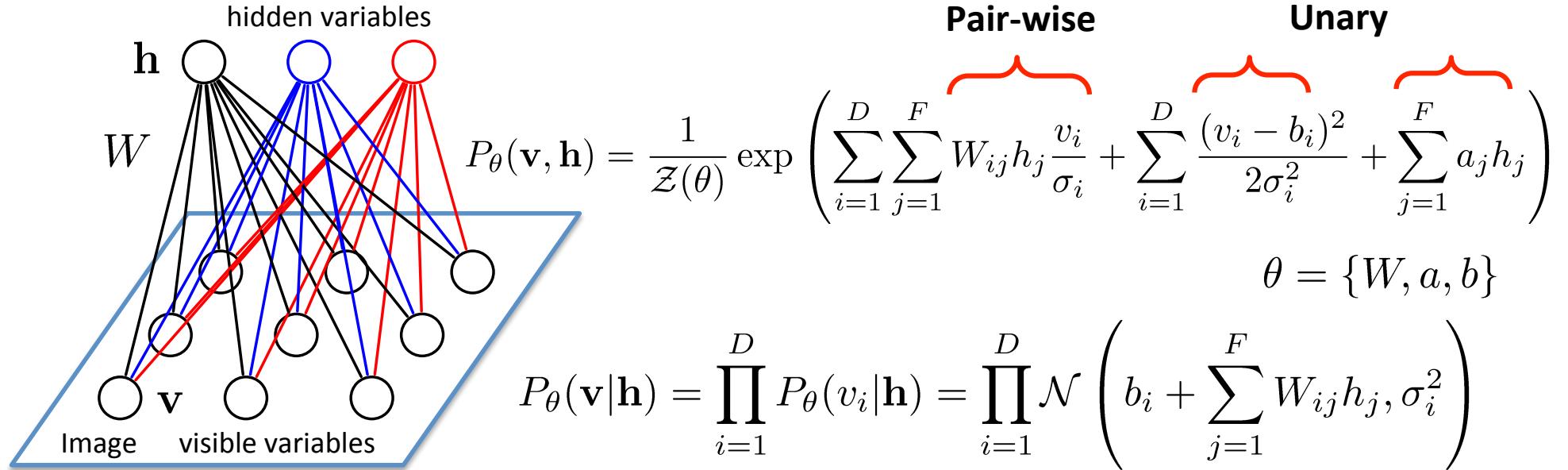
- Start with a training vector on the visible units.
- Update all the hidden units in parallel.
- Update all the visible units in parallel to get a “reconstruction”.
- Update the hidden units again.

Update model parameters:

$$\Delta W_{ij} = \text{E}_{P_{data}}[v_i h_j] - \text{E}_{P_1}[v_i h_j]$$

Implementation: ~10 lines of Matlab code.

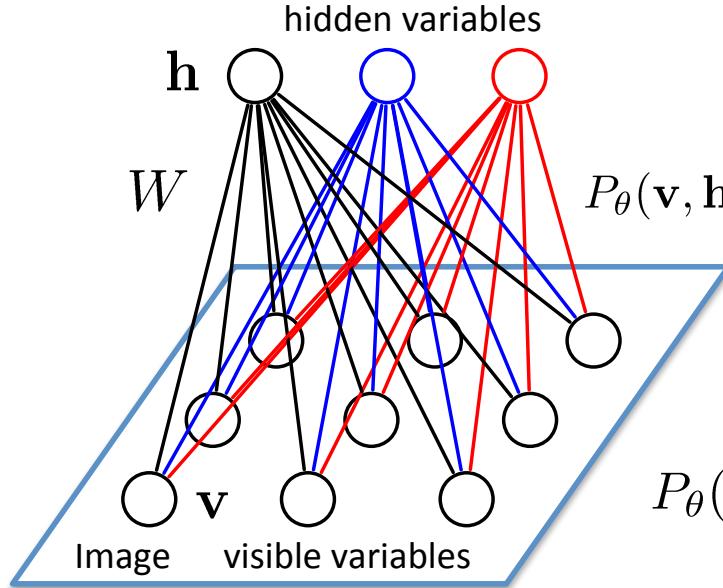
# RBM<sup>s</sup> for Real-valued Data



Gaussian-Bernoulli RBM:

- Stochastic real-valued visible variables  $\mathbf{v} \in \mathbb{R}^D$ .
- Stochastic binary hidden variables  $\mathbf{h} \in \{0, 1\}^F$ .
- Bipartite connections.

# RBM<sup>s</sup> for Real-valued Data



$$P_\theta(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \exp \left( \sum_{i=1}^D \sum_{j=1}^F W_{ij} h_j \frac{v_i}{\sigma_i} + \sum_{i=1}^D \frac{(v_i - b_i)^2}{2\sigma_i^2} + \sum_{j=1}^F a_j h_j \right)$$

**Pair-wise**      **Unary**

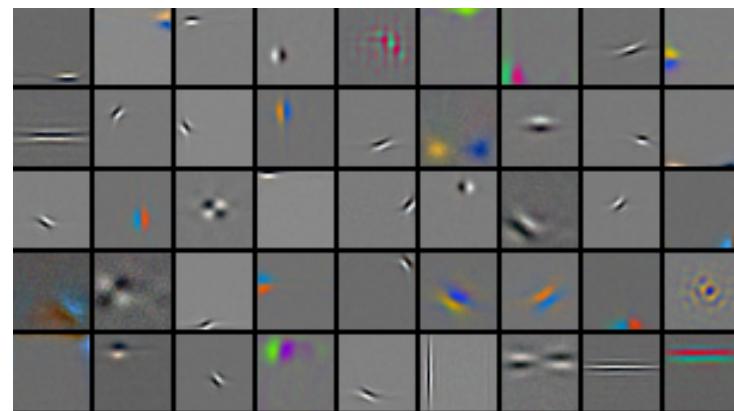
$$\theta = \{W, a, b\}$$

$$P_\theta(\mathbf{v}|\mathbf{h}) = \prod_{i=1}^D P_\theta(v_i|\mathbf{h}) = \prod_{i=1}^D \mathcal{N} \left( b_i + \sum_{j=1}^F W_{ij} h_j, \sigma_i^2 \right)$$

4 million **unlabelled** images



Learned features (out of 10,000)

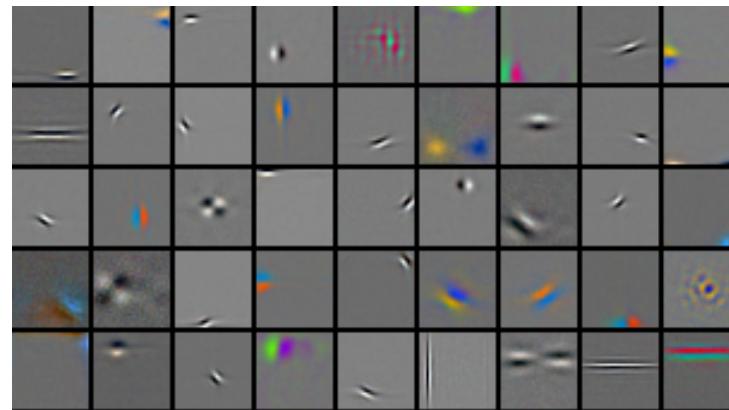


# RBM<sup>s</sup> for Real-valued Data

4 million **unlabelled** images

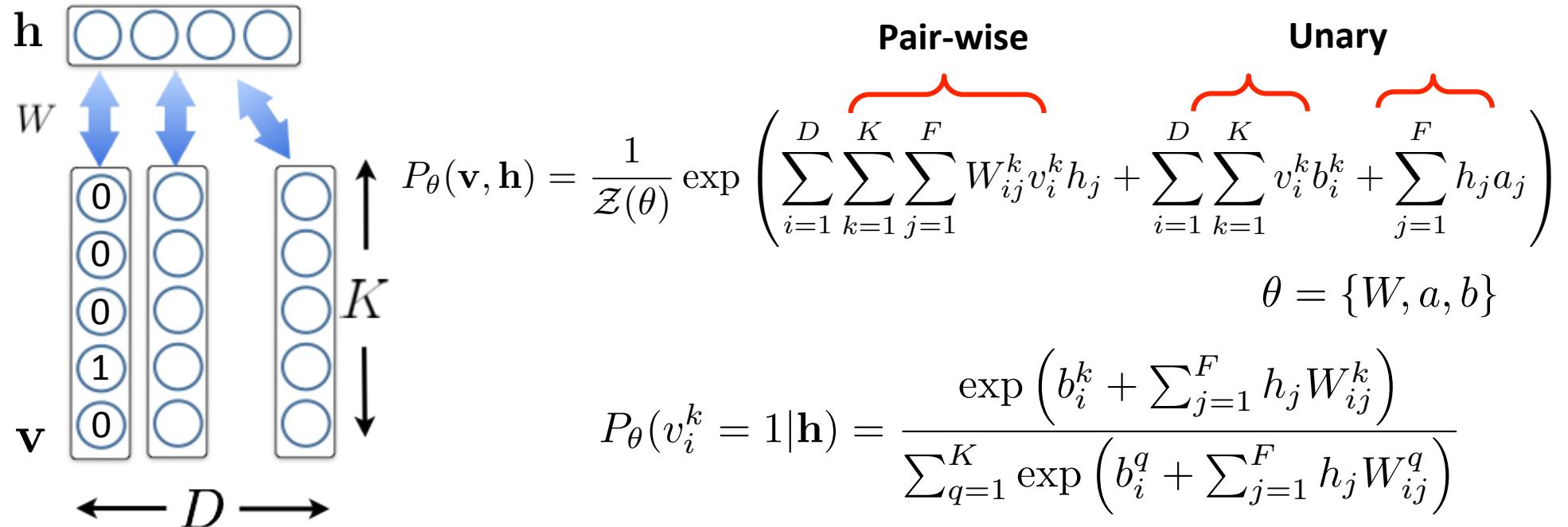


Learned features (out of 10,000)



$$\text{New Image} \quad p(h_7 = 1|v) \quad p(h_{29} = 1|v)$$
$$= 0.9 * \text{[feature image]} + 0.8 * \text{[feature image]} + 0.6 * \text{[feature image]} \dots$$

# RBMs for Word Counts

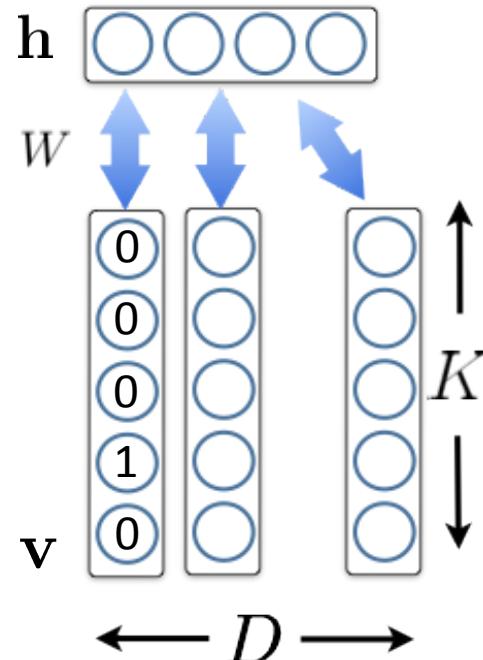


Replicated Softmax Model: undirected topic model:

- Stochastic 1-of- $K$  visible variables.
- Stochastic binary hidden variables  $\mathbf{h} \in \{0, 1\}^F$ .
- Bipartite connections.

(Salakhutdinov & Hinton, NIPS 2010, Srivastava & Salakhutdinov, NIPS 2012)

# RBMs for Word Counts



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

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

$$P_{\theta}(v_i^k = 1 | \mathbf{h}) = \frac{\exp(b_i^k + \sum_{j=1}^F h_j W_{ij}^k)}{\sum_{q=1}^K \exp(b_i^q + \sum_{j=1}^F h_j W_{ij}^q)}$$



REUTERS

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

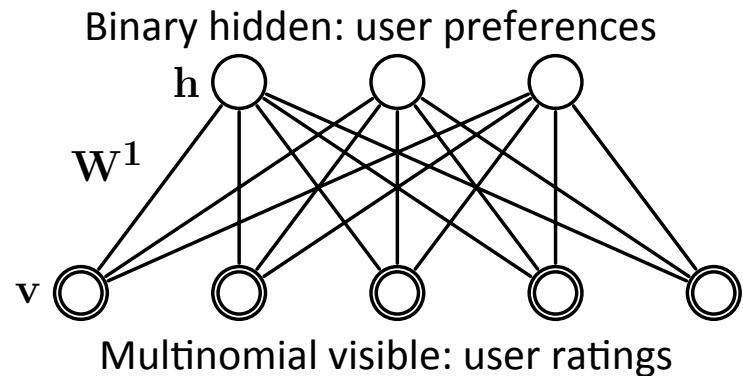
# RBM<sub>s</sub> for Word Counts

One-step reconstruction from the Replicated Softmax model.

<b>Input</b>	<b>Reconstruction</b>
chocolate, cake	cake, chocolate, sweets, dessert, cupcake, food, sugar, cream, birthday
nyc	nyc, newyork, brooklyn, queens, gothamist, manhattan, subway, streetart
dog	dog, puppy, perro, dogs, pet, filmshots, tongue, pets, nose, animal
flower, high, 花	flower, 花, high, japan, sakura, 日本, blossom, tokyo, lily, cherry
girl, rain, station, norway	norway, station, rain, girl, oslo, train, umbrella, wet, railway, weather
fun, life, children	children, fun, life, kids, child, playing, boys, kid, play, love
forest, blur	forest, blur, woods, motion, trees, movement, path, trail, green, focus
españa, agua, granada	españa, agua, spain, granada, water, andalucía, naturaleza, galicia, nieve

# Collaborative Filtering

$$P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{\mathcal{Z}(\theta)} \exp \left( \sum_{ijk} W_{ij}^k v_i^k h_j + \sum_{ik} b_i^k v_i^k + \sum_j a_j h_j \right)$$



Netflix dataset:

480,189 users

17,770 movies

Over 100 million ratings



Learned features: ``genre''

Fahrenheit 9/11  
Bowling for Columbine  
The People vs. Larry Flynt  
Canadian Bacon  
La Dolce Vita

Independence Day  
The Day After Tomorrow  
Con Air  
Men in Black II  
Men in Black

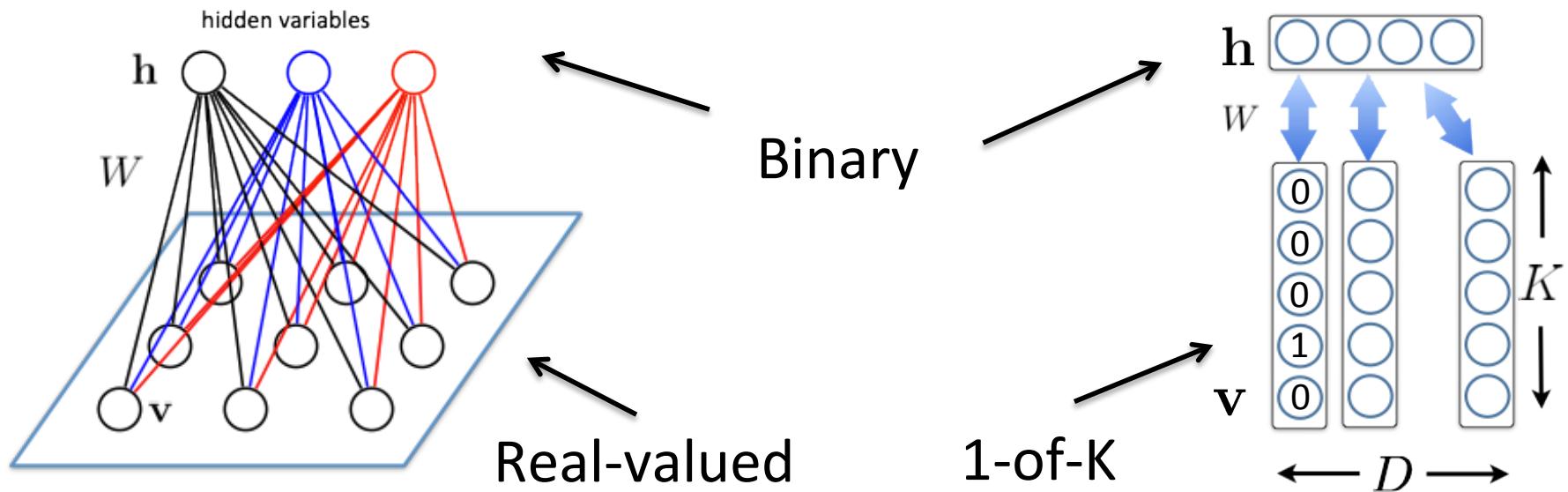
Friday the 13th  
The Texas Chainsaw Massacre  
Children of the Corn  
Child's Play  
The Return of Michael Myers

Scary Movie  
Naked Gun  
Hot Shots!  
American Pie  
Police Academy

(Salakhutdinov, Mnih, Hinton, ICML 2007)

# Different Data Modalities

- Binary/Gaussian/Softmax RBMs: All have binary hidden variables but use them to model different kinds of data.



- It is easy to infer the states of the hidden variables:

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

# Product of Experts

The joint distribution is given by:

$$P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \exp \left( \sum_{ij} W_{ij} v_i h_j + \sum_i b_i v_i + \sum_j a_j h_j \right)$$

Marginalizing over hidden variables:

$$P_{\theta}(\mathbf{v}) = \sum_{\mathbf{h}} P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \prod_i \exp(b_i v_i) \prod_j \left( 1 + \exp(a_j + \sum_i W_{ij} v_i) \right)$$

government	clinton	bribery	mafia	stock	...
authority	house	corruption	business	wall	
power	president	dishonesty	gang	street	
empire	bill	corrupt	mob	point	
federation	congress	fraud	insider	dow	

**Product of Experts**

Silvio Berlusconi

Topics “government”, “corruption” and “mafia” can combine to give very high probability to a word “Silvio Berlusconi”.

# Product of Experts

The joint distribution is given by:

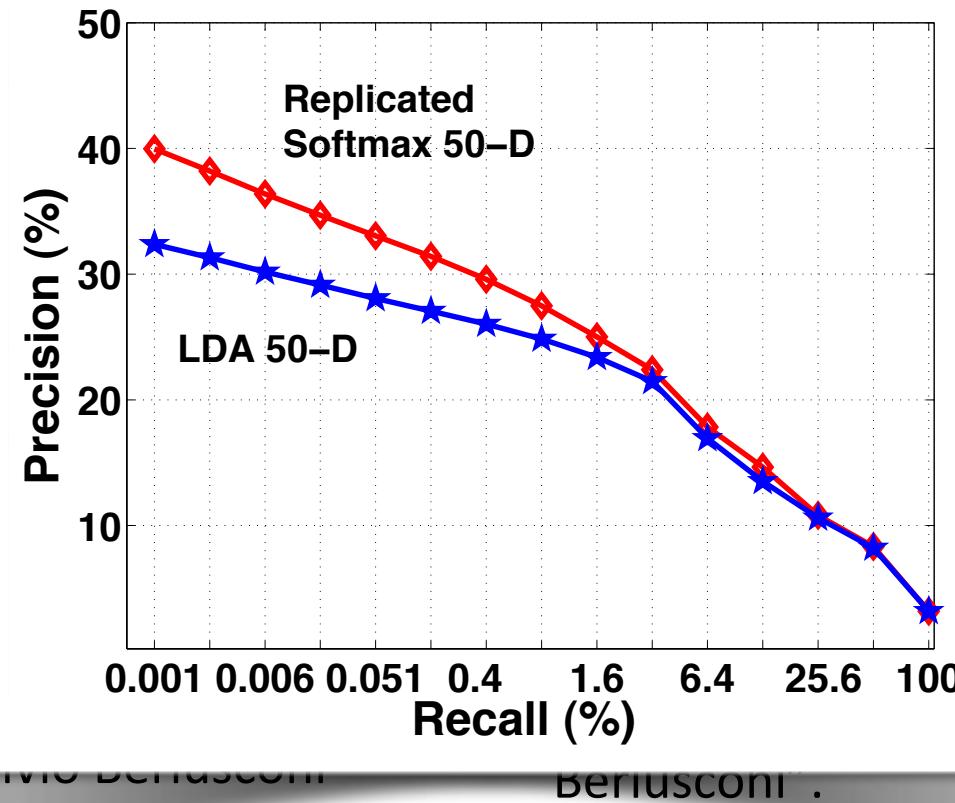
$$P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \exp \left( \sum_{ij} W_{ij} v_i h_j + \sum_i b_i v_i + \sum_j a_j h_j \right)$$

Marginalizing over  $\mathbf{h}$ :

$$P_{\theta}(\mathbf{v}) = \sum_{\mathbf{h}}$$

government  
authority  
power  
empire  
federation

clint  
hou  
pres  
bill  
congr



Product of Experts

$$W_{ij} v_i)$$

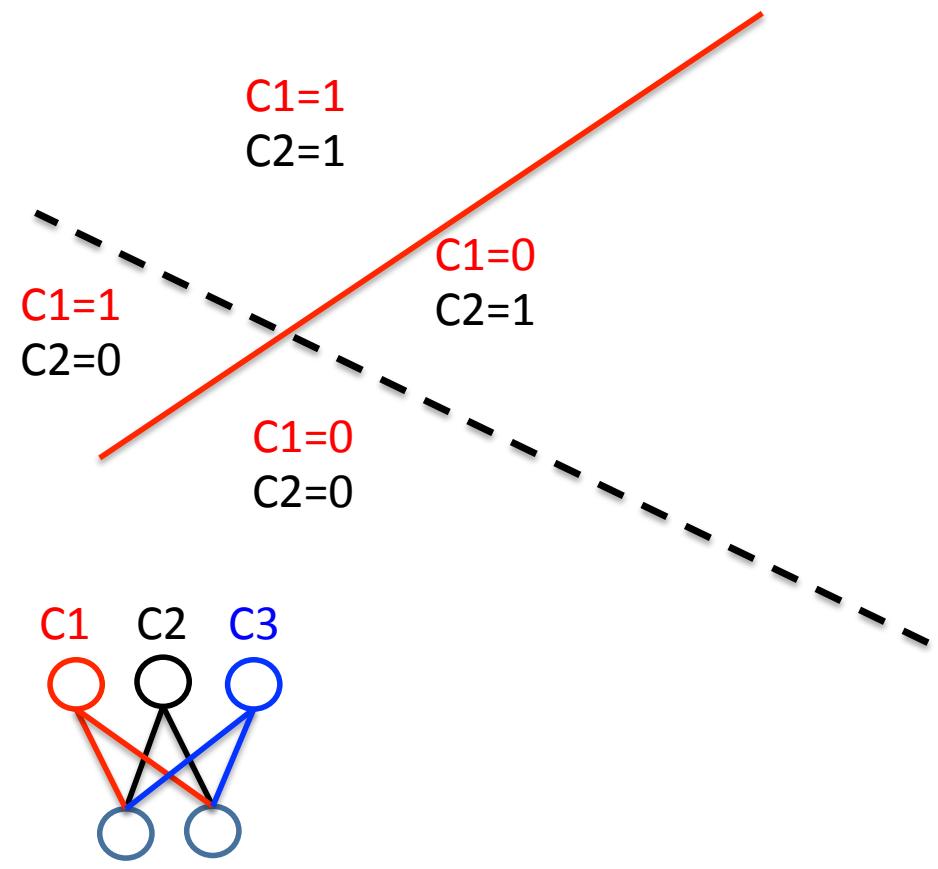
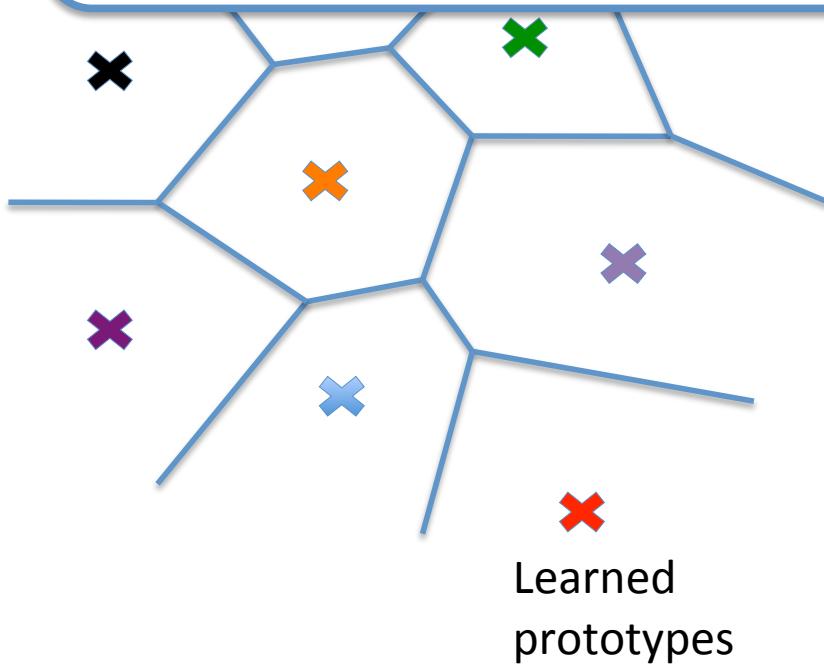
, "corruption"  
bine to give very  
word "Silvio

# Local vs. Distributed Representations

- Clustering, Nearest Neighbors, RBF SVM, local density estimators

- RBMs, Factor models, PCA, Sparse Coding, Deep models

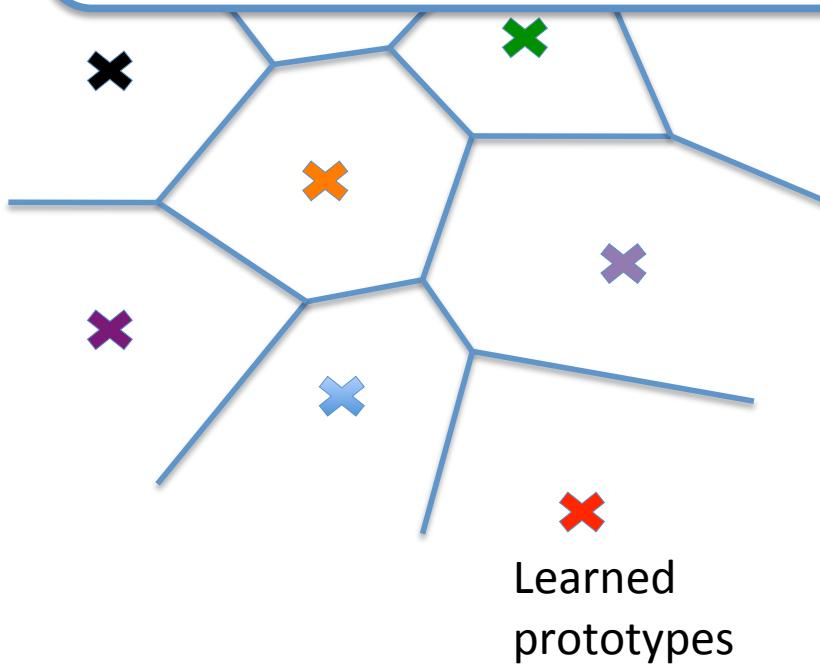
- Parameters for each region.
- # of regions is linear with # of parameters.



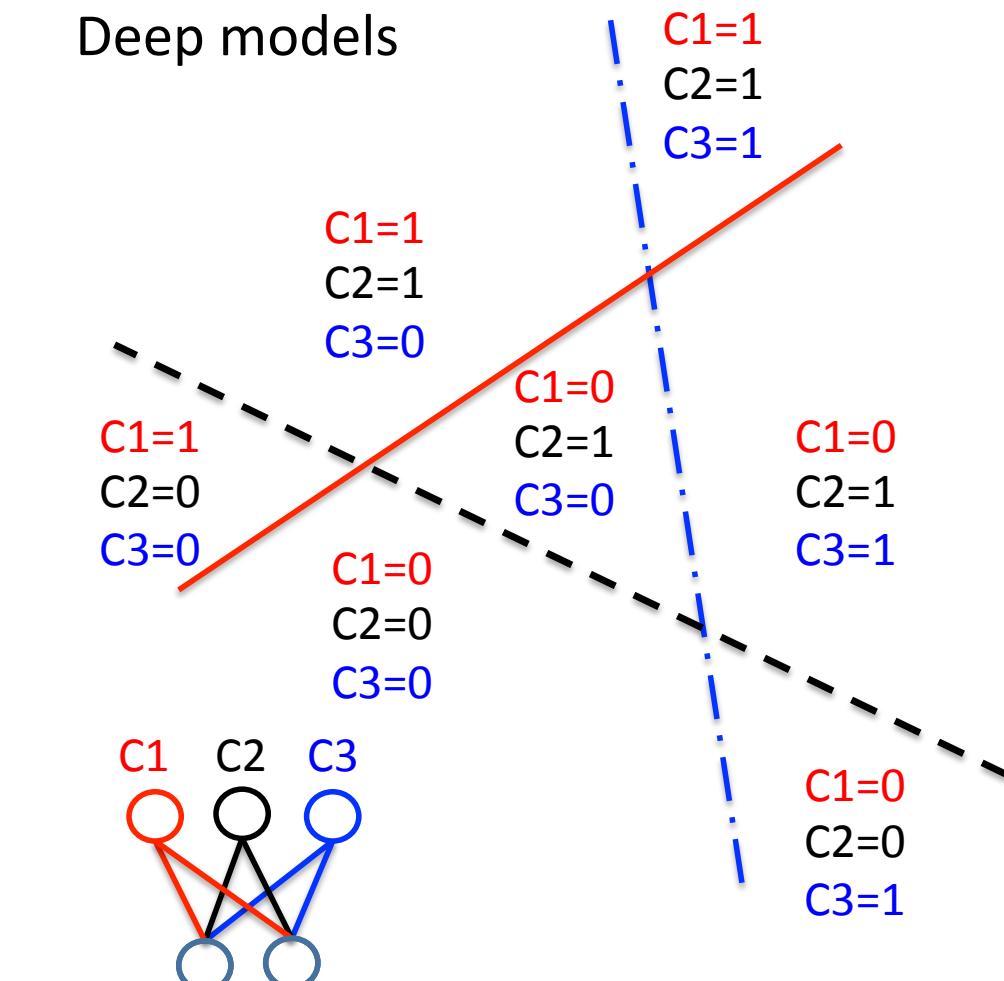
# Local vs. Distributed Representations

- Clustering, Nearest Neighbors, RBF SVM, local density estimators

- Parameters for each region.
- # of regions is linear with # of parameters.



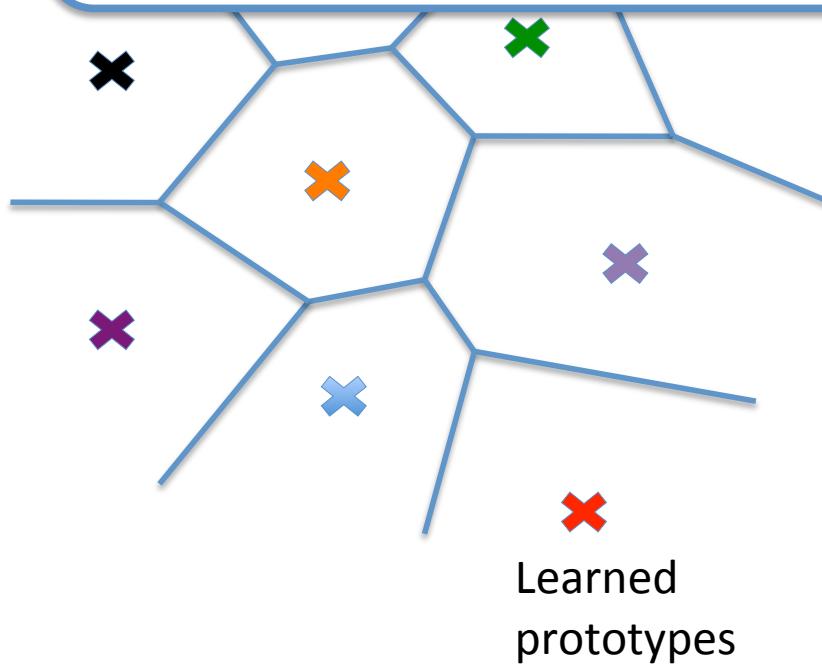
- RBMs, Factor models, PCA, Sparse Coding, Deep models



# Local vs. Distributed Representations

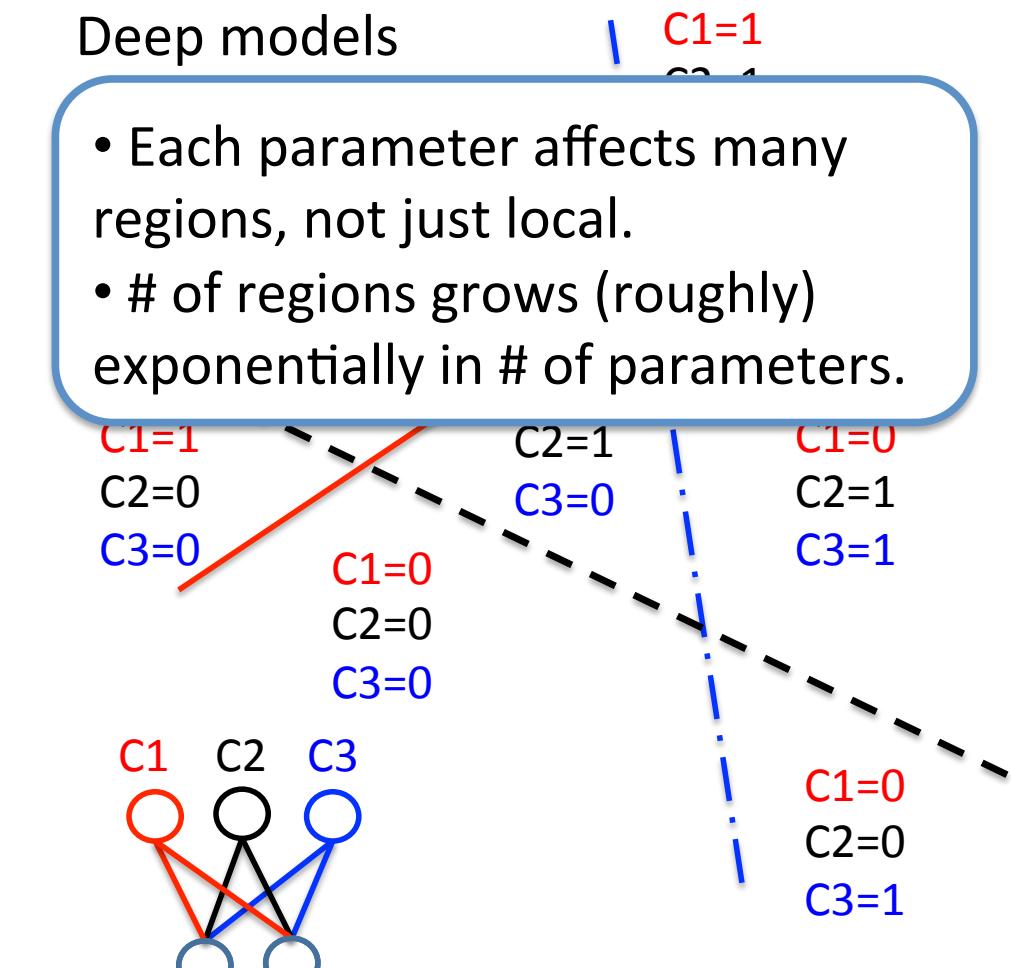
- Clustering, Nearest Neighbors, RBF SVM, local density estimators

- Parameters for each region.
- # of regions is linear with # of parameters.



- RBMs, Factor models, PCA, Sparse Coding, Deep models

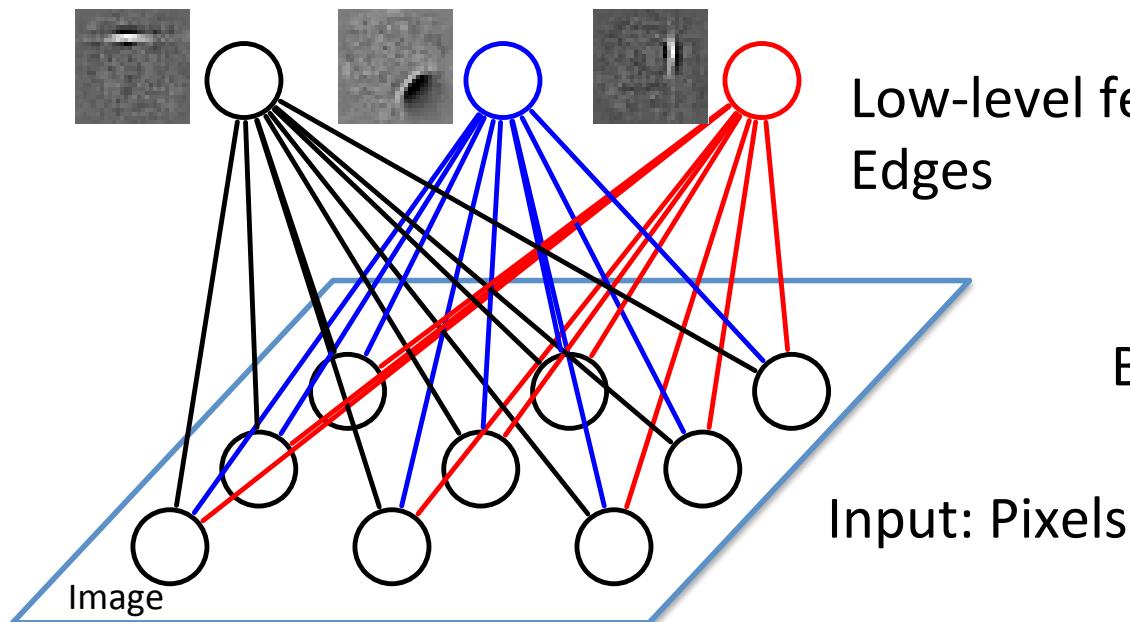
- Each parameter affects many regions, not just local.
- # of regions grows (roughly) exponentially in # of parameters.



# Talk Roadmap

- Basic Building Blocks (non-probabilistic models):
  - Sparse Coding
  - Autoencoders
- Deep Generative Models
  - Restricted Boltzmann Machines
  - Deep Boltzmann Machines
  - Helmholtz Machines / Variational Autoencoders
- Generative Adversarial Networks

# Deep Boltzmann Machines

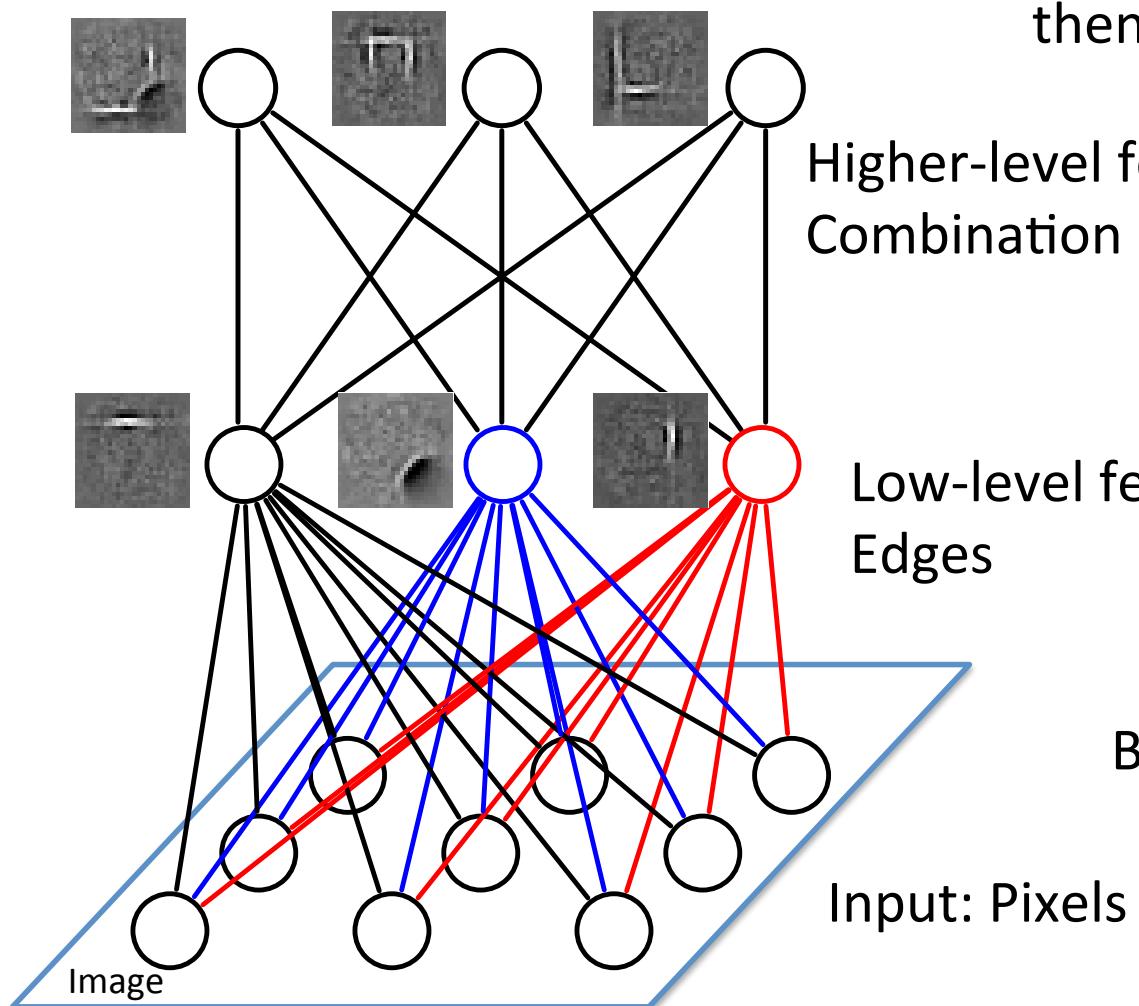


Built from **unlabeled** inputs.

Input: Pixels

(Salakhutdinov 2008, Salakhutdinov & Hinton 2012)

# Deep Boltzmann Machines



Learn simpler representations,  
then compose more complex ones

Higher-level features:  
Combination of edges

Low-level features:  
Edges

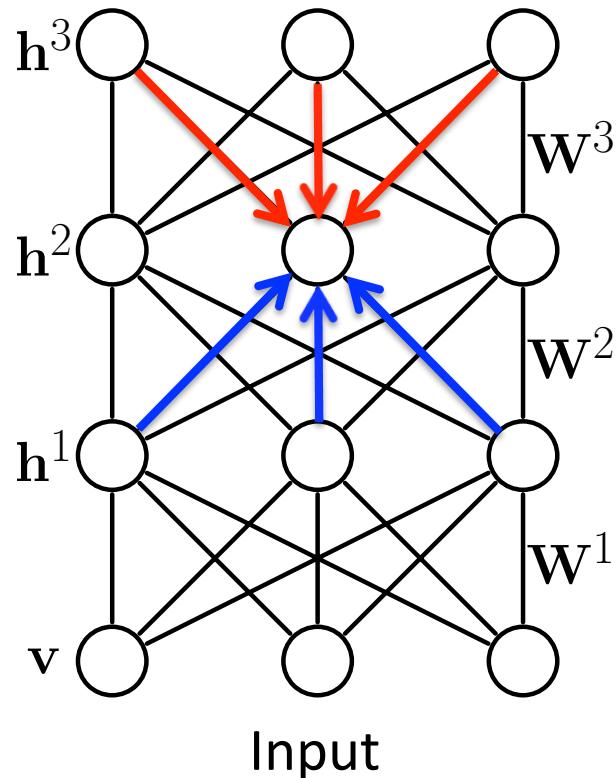
Built from **unlabeled** inputs.

Input: Pixels

(Salakhutdinov 2008, Salakhutdinov & Hinton 2012)

# Model Formulation

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



Same as RBMs

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

- Dependencies between hidden variables.
- All connections are undirected.
- Bottom-up and Top-down:

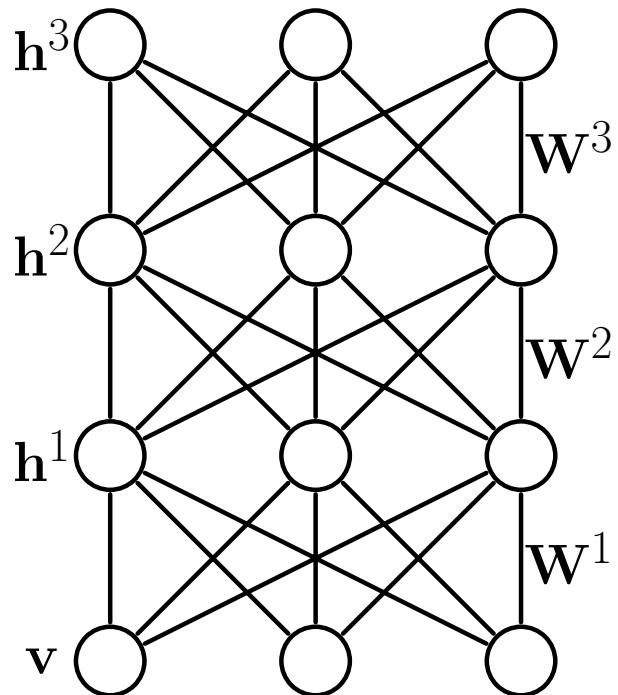
$$P(h_j^2 = 1 | \mathbf{h}^1, \mathbf{h}^3) = \sigma \left( \sum_k W_{kj}^3 h_k^3 + \sum_m W_{mj}^2 h_m^1 \right)$$

Top-down                      Bottom-up

- Hidden variables are dependent even when **conditioned on the input**.

# Approximate Learning

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



(Approximate) Maximum Likelihood:

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

- Both expectations are intractable!

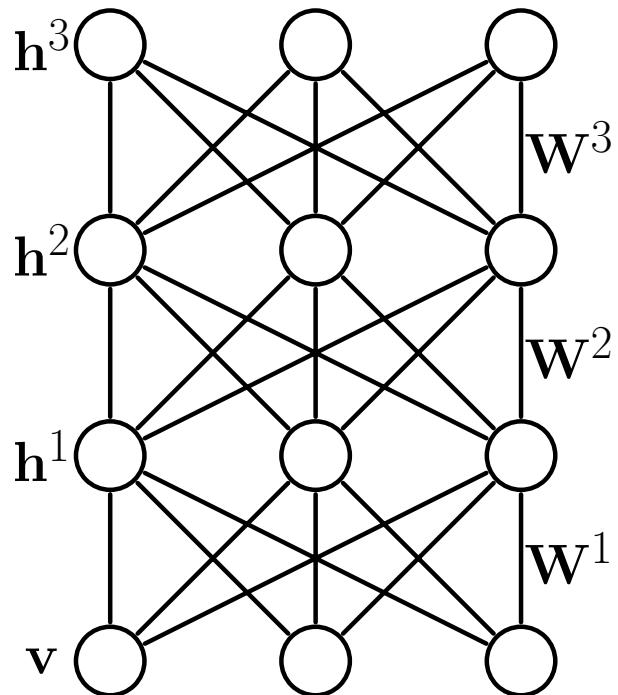
$$P_{data}(\mathbf{v}, \mathbf{h}^1) = P_{\theta}(\mathbf{h}^1 | \mathbf{v}) P_{data}(\mathbf{v})$$

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

Not factorial any more!

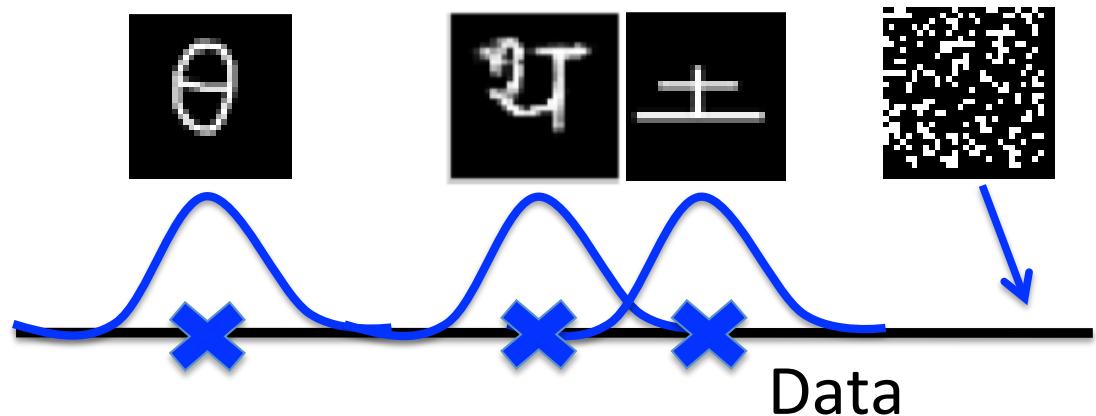
# Approximate Learning

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



(Approximate) Maximum Likelihood:

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



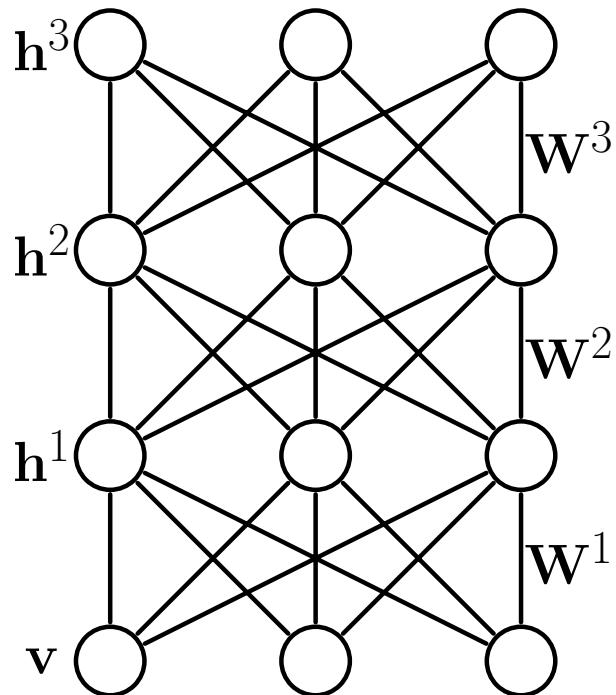
$$P_{data}(\mathbf{v}, \mathbf{h}^1) = P_{\theta}(\mathbf{h}^1 | \mathbf{v}) P_{data}(\mathbf{v})$$

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

Not factorial any more!

# Approximate Learning

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



(Approximate) Maximum Likelihood:

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

Variational  
Inference

Stochastic  
Approximation  
(MCMC-based)

$$P_{data}(\mathbf{v}, \mathbf{h}^1) = P_\theta(\mathbf{h}^1 | \mathbf{v}) P_{data}(\mathbf{v})$$

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

Not factorial any more!

# Good Generative Model?

Handwritten Characters

# Good Generative Model?

## Handwritten Characters

手 書 か ら で は る と  
た ち て ま す く な い  
し う か い て は じ ま  
し か せ ば あ は い  
し か ま す そ う す て  
し か く き て は じ ま  
し か ま す そ う す て  
し か く き て は じ ま  
し か ま す そ う す て

手 書 ま め ; は じ ま  
た ち て ま す く な い  
し う か い て は じ ま  
し か せ ば あ は い  
し か ま す そ う す て  
し か く き て は じ ま  
し か ま す そ う す て  
し か く き て は じ ま  
し か ま す そ う す て

# Good Generative Model?

Handwritten Characters

Simulated

Real Data

# Good Generative Model?

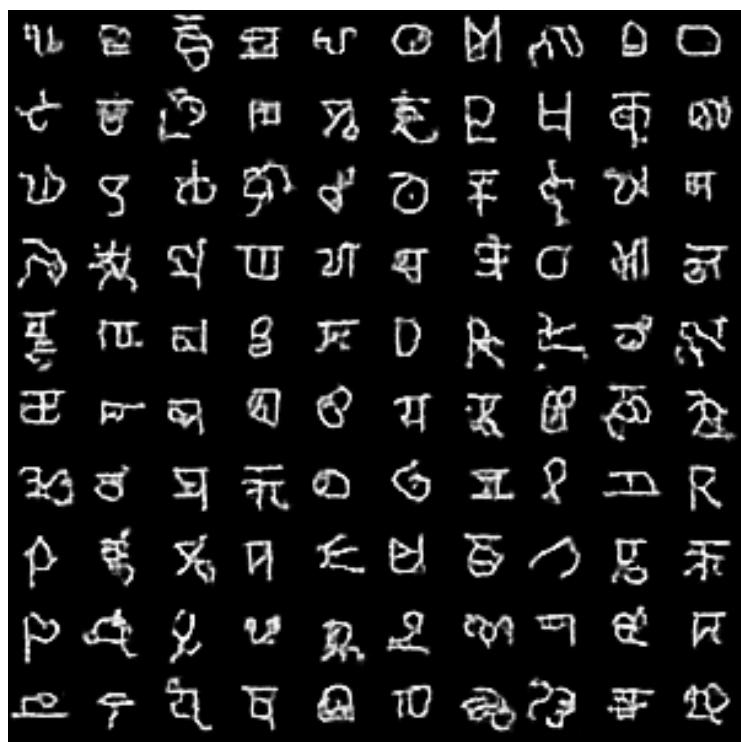
Handwritten Characters

Real Data

Simulated

# Good Generative Model?

# Handwritten Characters



# Handwriting Recognition

MNIST Dataset

60,000 examples of 10 digits

Learning Algorithm	Error
Logistic regression	12.0%
K-NN	3.09%
Neural Net (Platt 2005)	1.53%
SVM (Decoste et.al. 2002)	1.40%
Deep Autoencoder (Bengio et. al. 2007)	1.40%
Deep Belief Net (Hinton et. al. 2006)	1.20%
<b>DBM</b>	<b>0.95%</b>

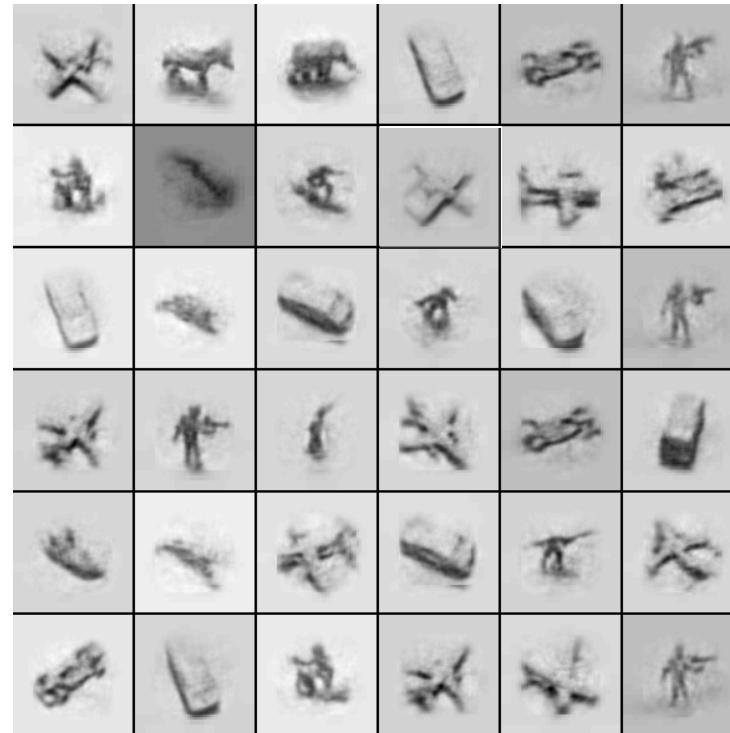
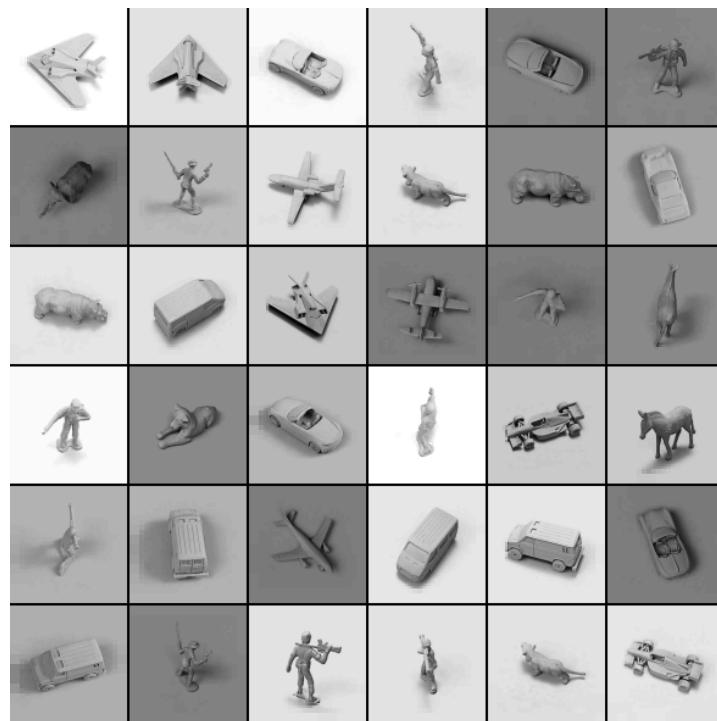
Optical Character Recognition

42,152 examples of 26 English letters

Learning Algorithm	Error
Logistic regression	22.14%
K-NN	18.92%
Neural Net	14.62%
SVM (Larochelle et.al. 2009)	9.70%
Deep Autoencoder (Bengio et. al. 2007)	10.05%
Deep Belief Net (Larochelle et. al. 2009)	9.68%
<b>DBM</b>	<b>8.40%</b>

Permutation-invariant version.

# Generative Model of 3-D Objects

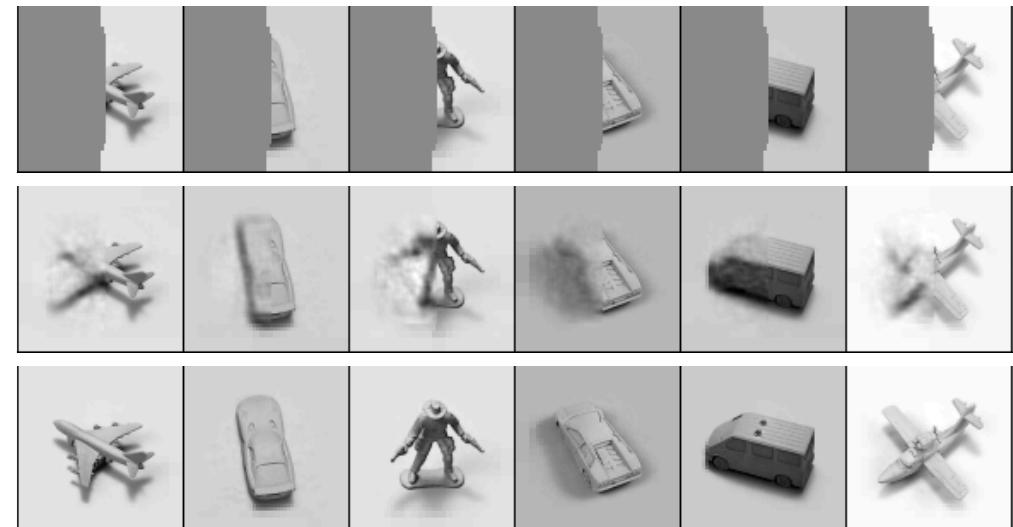


24,000 examples, 5 object categories, 5 different objects within each category, 6 lightning conditions, 9 elevations, 18 azimuths.

# 3-D Object Recognition

Learning Algorithm	Error
Logistic regression	22.5%
K-NN (LeCun 2004)	18.92%
SVM (Bengio & LeCun 2007)	11.6%
Deep Belief Net (Nair & Hinton 2009)	9.0%
<b>DBM</b>	<b>7.2%</b>

Pattern Completion

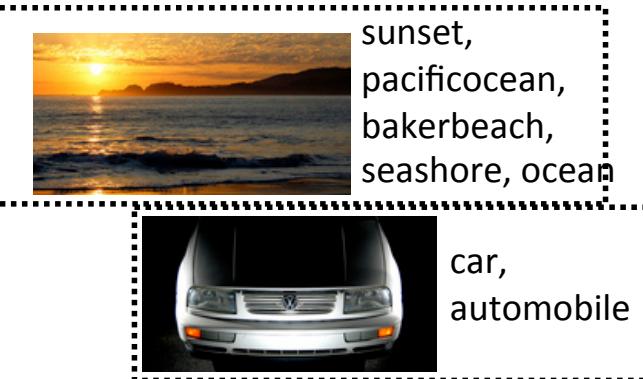


Permutation-invariant version.

Where else can we use  
generative models?

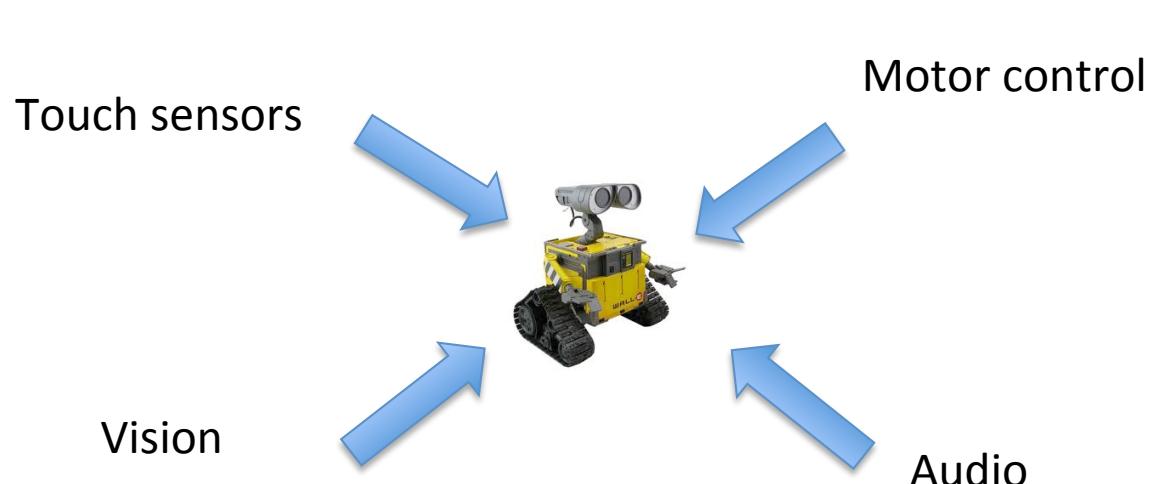
# Data – Collection of Modalities

- Multimedia content on the web - image + text + audio.



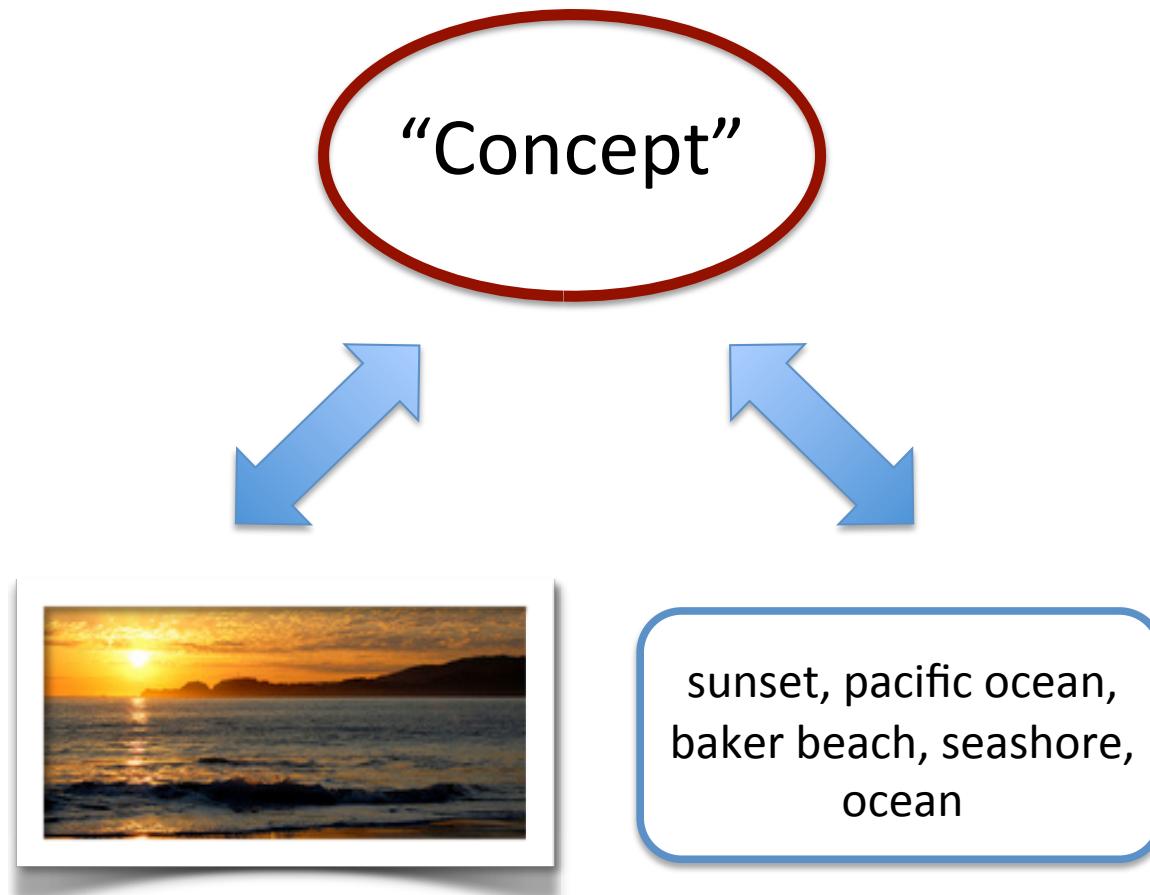
- Product recommendation systems.

- Robotics applications.



# Shared Concept

“Modality-free” representation



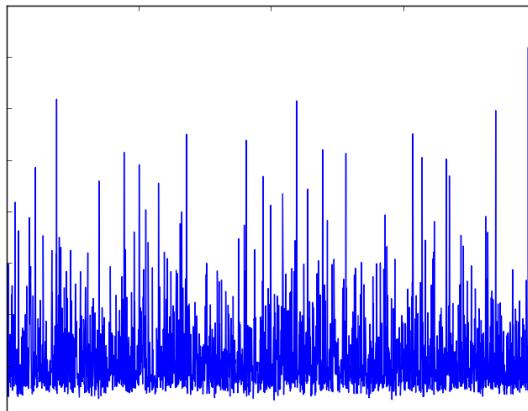
“Modality-full” representation

# Challenges - I

Image



Dense



Text

sunset, pacific ocean,  
baker beach, seashore,  
ocean



Very different input representations

- Images – real-valued, dense
- Text – discrete, sparse

Difficult to learn cross-modal features from low-level representations.

# Challenges - II

Image



Text

pentax, k10d,  
pentaxda50200,  
kangarooisland, sa,  
australiansealion

Noisy and missing data



mickikrimmel,  
mickipedia,  
headshot



< no text>



unseulpixel,  
naturey

# Challenges - II

## Image



## Text

pentax, k10d,  
pentaxda50200,  
kangarooisland, sa,  
australiansealion

## Text generated by the model

beach, sea, surf, strand,  
shore, wave, seascape,  
sand, ocean, waves



mickikrimmel,  
mickipedia,  
headshot

portrait, girl, woman, lady,  
blonde, pretty, gorgeous,  
expression, model



< no text>

night, notte, traffic, light,  
lights, parking, darkness,  
lowlight, nacht, glow

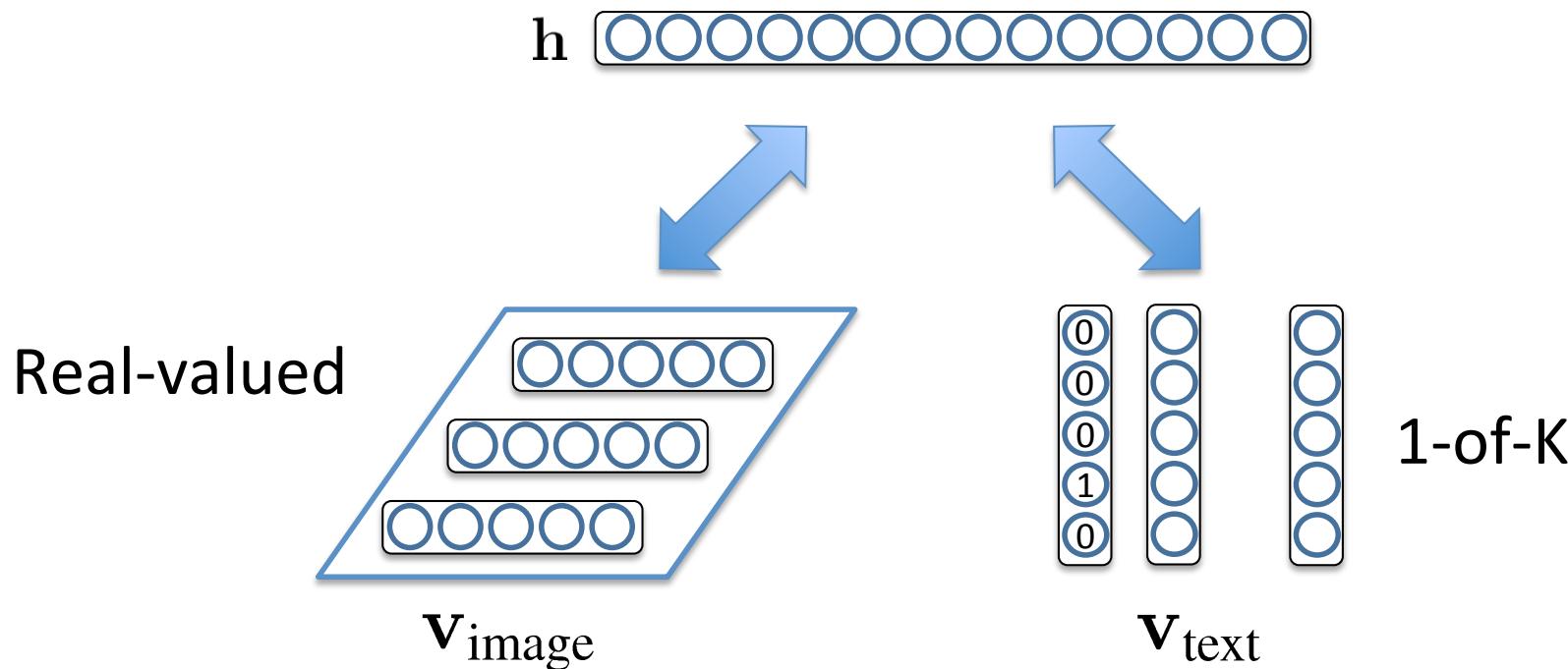


unseulpixel,  
naturey

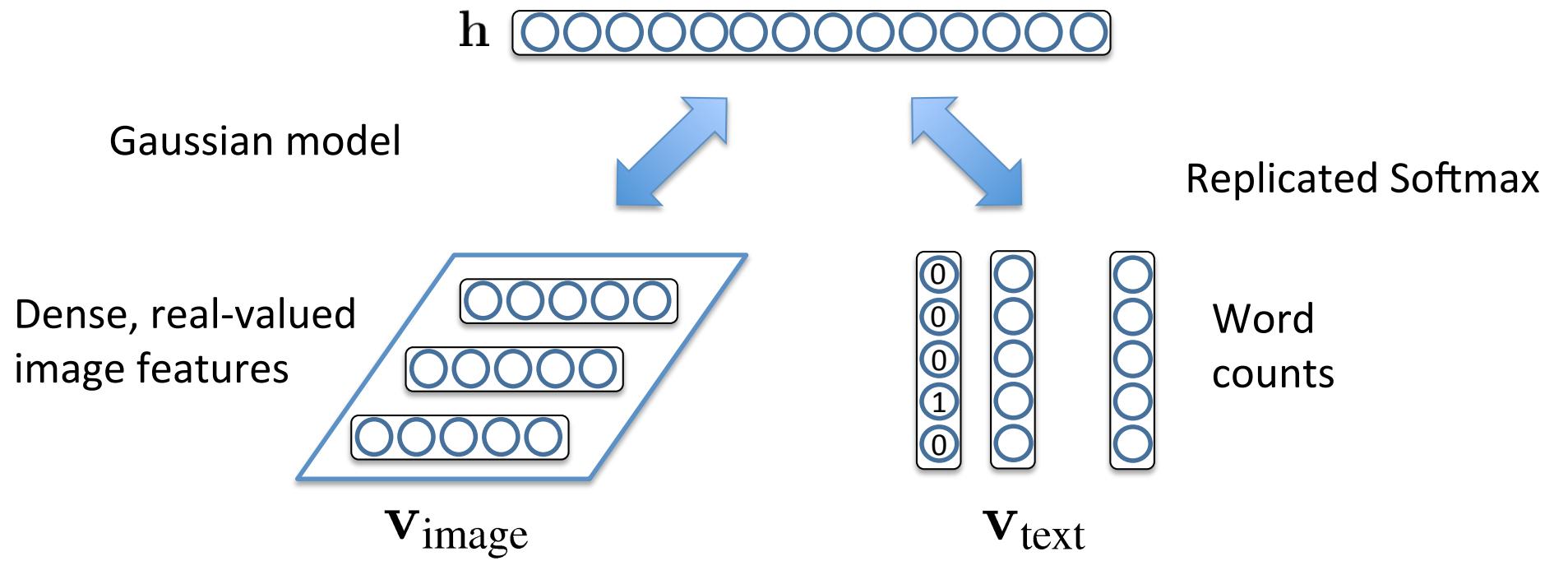
fall, autumn, trees, leaves,  
foliage, forest, woods,  
branches, path

# A Simple Multimodal Model

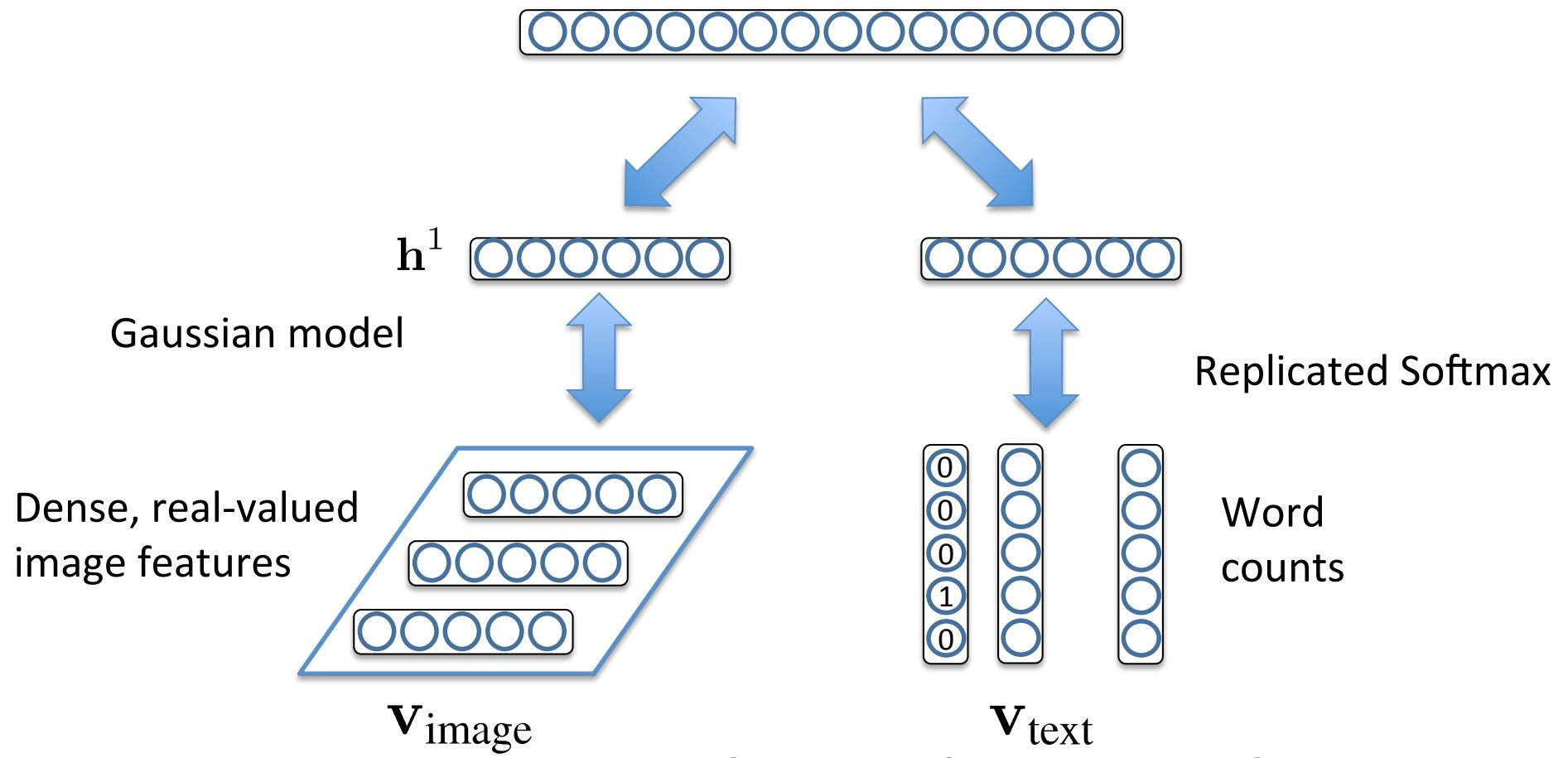
- Use a joint binary hidden layer.
- **Problem:** Inputs have very different statistical properties.
- Difficult to learn cross-modal features.



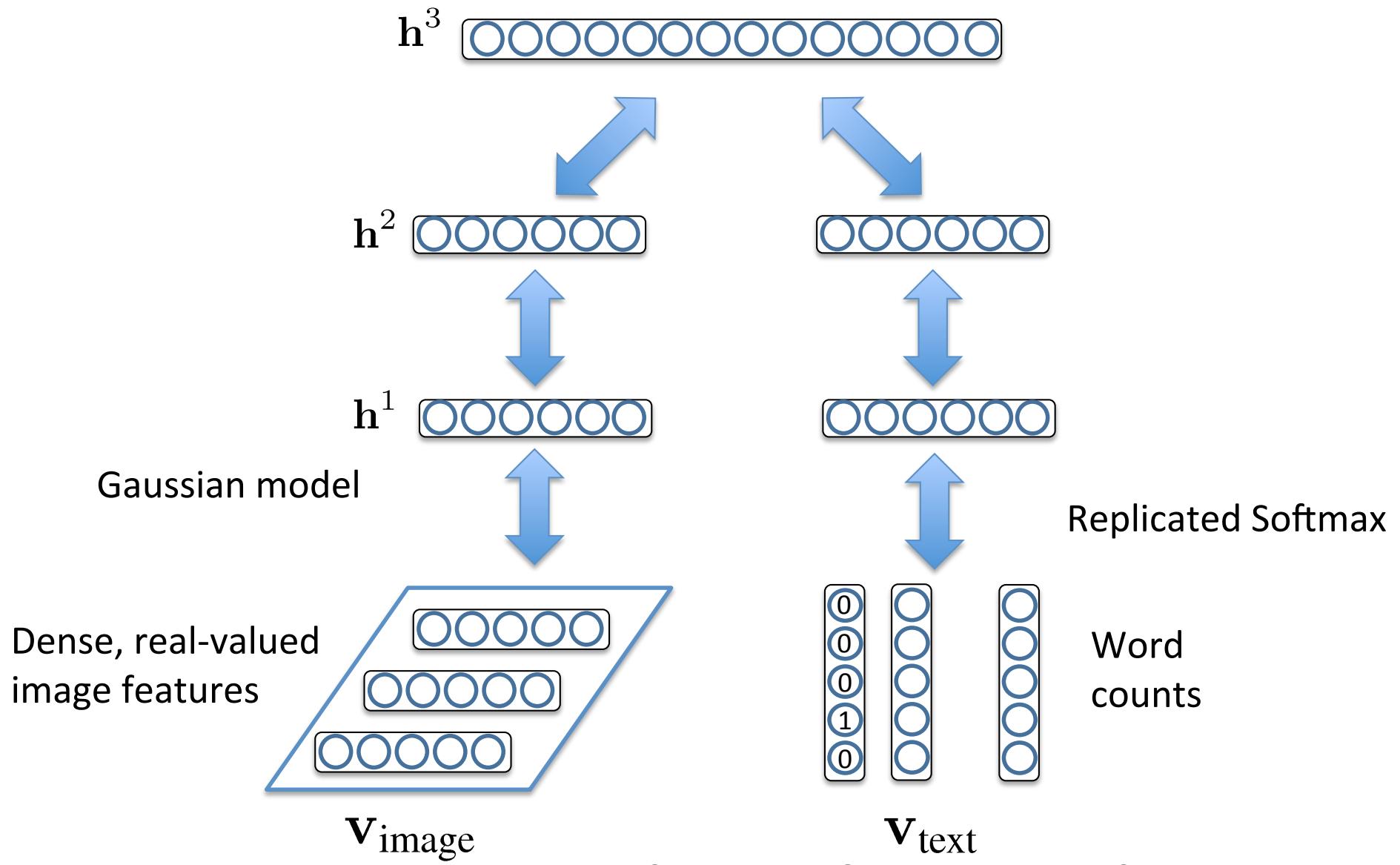
# Multimodal DBM



# Multimodal DBM

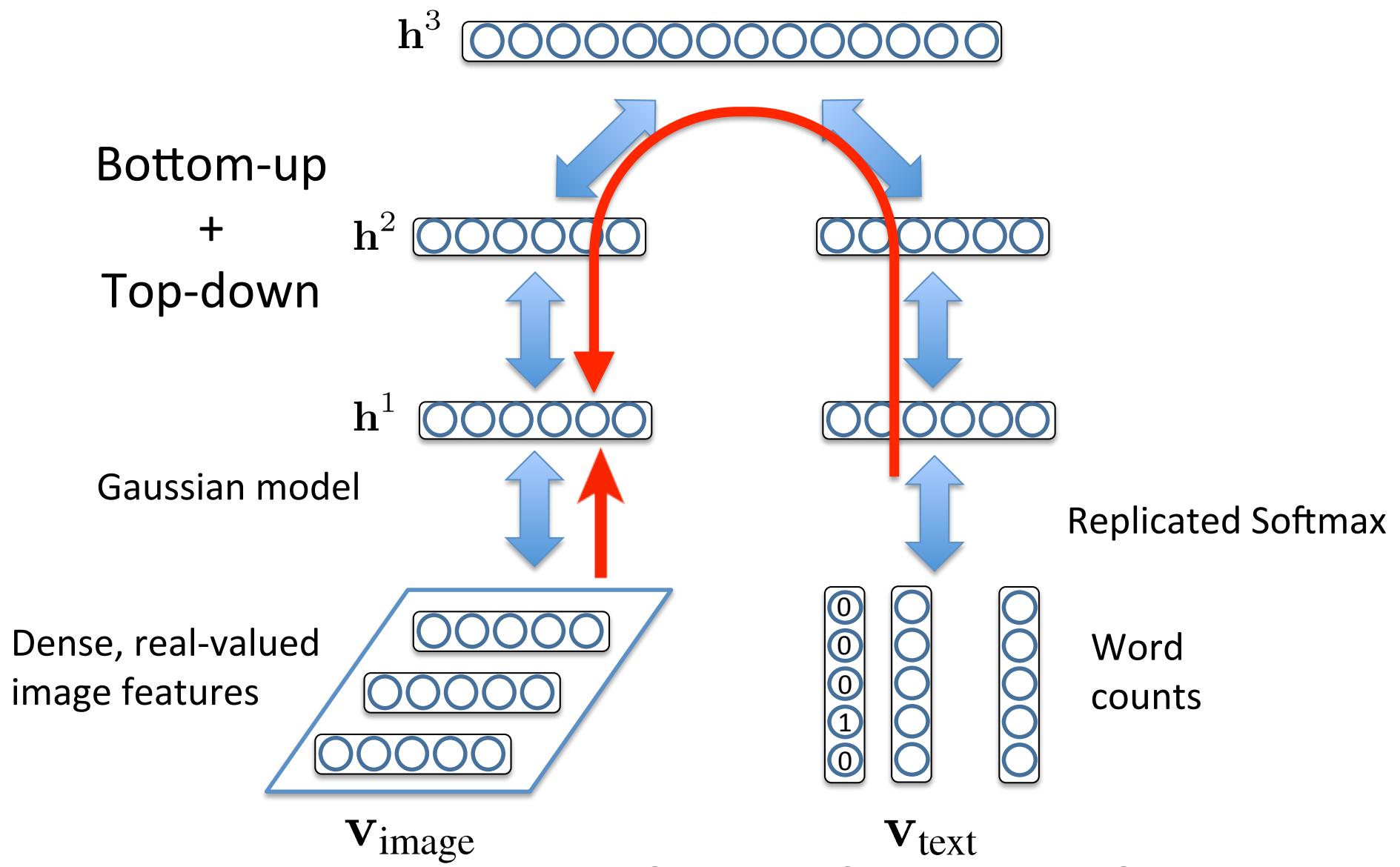


# Multimodal DBM



(Srivastava & Salakhutdinov, NIPS 2012, JMLR 2014)

# Multimodal DBM



(Srivastava & Salakhutdinov, NIPS 2012, JMLR 2014)

# Text Generated from Images

Given



Generated

dog, cat, pet, kitten,  
puppy, ginger, tongue,  
kitty, dogs, furry



sea, france, boat, mer,  
beach, river, bretagne,  
plage, brittany



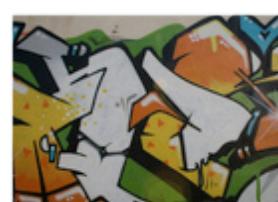
portrait, child, kid,  
ritratto, kids, children,  
boy, cute, boys, italy

Given



Generated

insect, butterfly, insects,  
bug, butterflies,  
lepidoptera



graffiti, streetart, stencil,  
sticker, urbanart, graff,  
sanfrancisco



canada, nature,  
sunrise, ontario, fog,  
mist, bc, morning

# Text Generated from Images

Given



Generated

portrait, women, army, soldier,  
mother, postcard, soldiers

Given

  
A photograph of a white heron with long legs and a long beak, standing on a wire mesh fence that spans across a body of blue water. The heron is facing towards the left of the frame.

Generated

obama, barackobama, election,  
politics, president, hope, change,  
sanfrancisco, convention, rally



Generated

water, glass, beer, bottle,  
drink, wine, bubbles, splash,  
drops, drop

# Images from Text

Given

water, red,  
sunset

Retrieved



nature, flower,  
red, green



blue, green,  
yellow, colors



chocolate, cake



# MIR-Flickr Dataset

- 1 million images along with user-assigned tags.



sculpture, beauty, stone



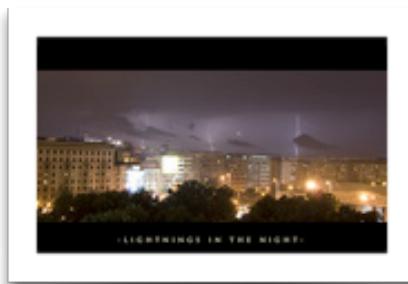
d80



nikon, abigfave, goldstaraward, d80, nikond80



food, cupcake, vegan



anawesomeshot, theperfectphotographer, flash, damniwishidtakenthat, spiritofphotography



nikon, green, light, photoshop, apple, d70



white, yellow, abstract, lines, bus, graphic



sky, geotagged, reflection, cielo, bilbao, reflejo

Huiskes et. al.

# Results

- Logistic regression on top-level representation.
- Multimodal Inputs

Mean Average Precision



Labeled  
25K  
examples

Learning Algorithm	MAP	Precision@50
Random	0.124	0.124
LDA [Huiskes et. al.]	0.492	0.754
SVM [Huiskes et. al.]	0.475	0.758
DBM-Labelled	0.526	0.791

# Results

- Logistic regression on top-level representation.
  - Multimodal Inputs

Learning Algorithm	MAP	Precision@50
Random	0.124	0.124
LDA [Huiskes et. al.]	0.492	0.754
SVM [Huiskes et. al.]	0.475	0.758
DBM-Labelled	0.526	0.791
Deep Belief Net	0.638	0.867
Autoencoder	0.638	0.875
DBM	0.641	0.873

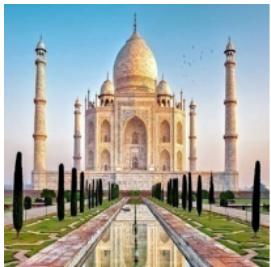
# Mean Average Precision

Labeled  
25K  
examples

+ 1 Million  
unlabelled

# Multimodal Linguistic Regularities

Nearest Images



- day + night =



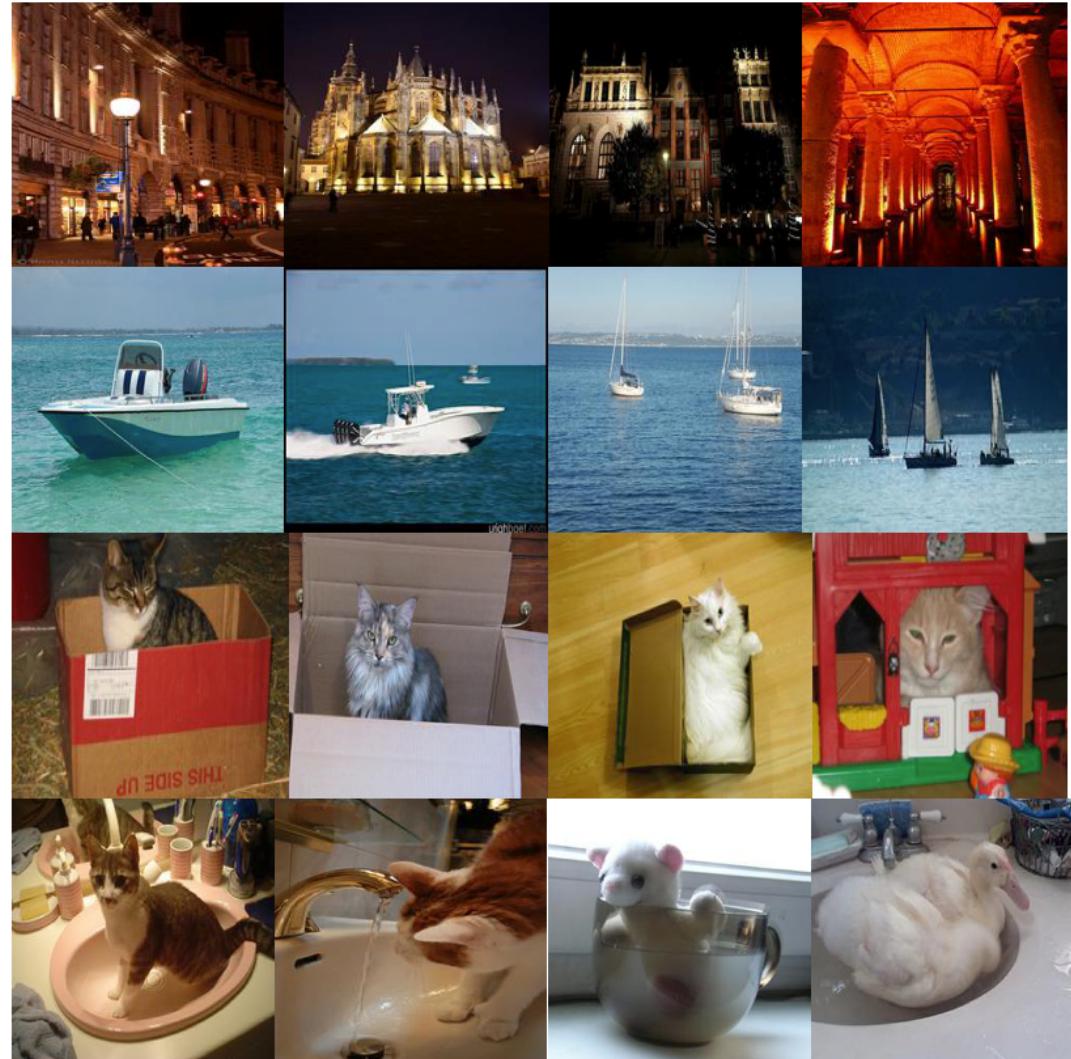
- flying + sailing =



- bowl + box =



- box + bowl =



(Kiros, Salakhutdinov, Zemel, TACL 2015)

# Talk Roadmap

- Basic Building Blocks:

- Sparse Coding
- Autoencoders

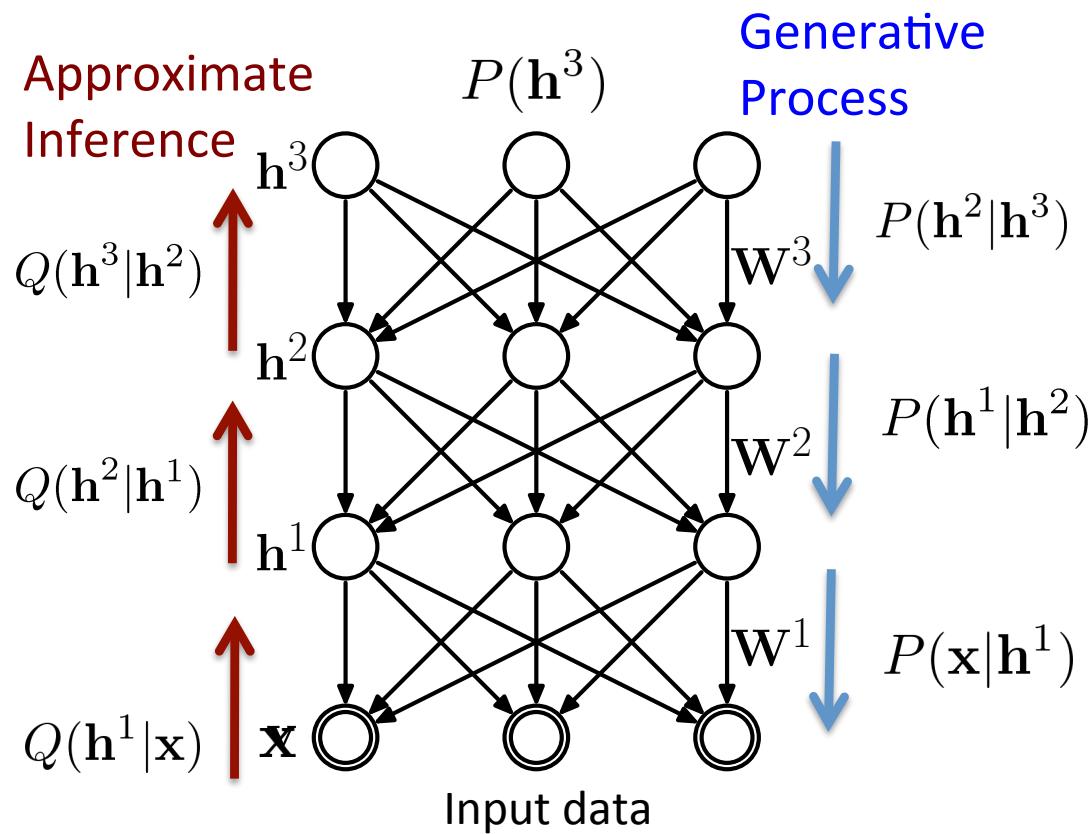
- Deep Generative Models

- Restricted Boltzmann Machines
- Deep Boltzmann Machines
- Helmholtz Machines / Variational Autoencoders

- Generative Adversarial Networks

# Helmholtz Machines

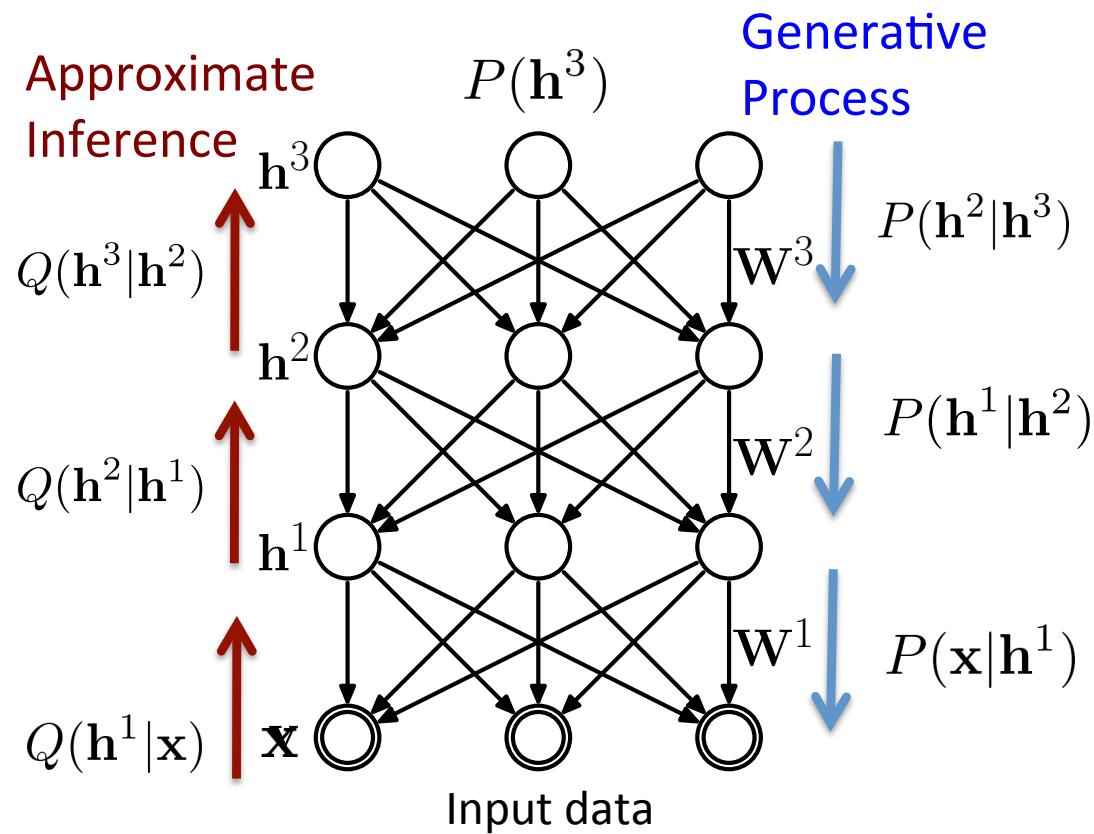
- Hinton, G. E., Dayan, P., Frey, B. J. and Neal, R., Science 1995



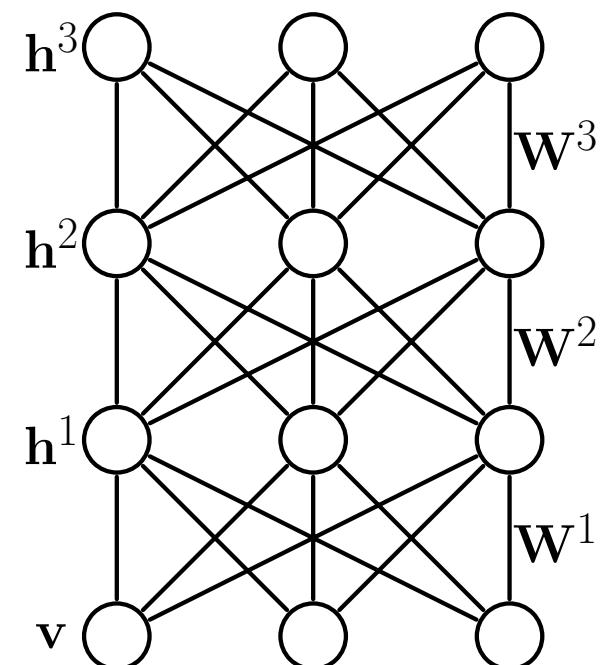
- Kingma & Welling, 2014
- Rezende, Mohamed, Daan, 2014
- Mnih & Gregor, 2014
- Bornschein & Bengio, 2015
- Tang & Salakhutdinov, 2013

# Helmholtz Machines vs. DBMs

Helmholtz Machine



Deep Boltzmann Machine



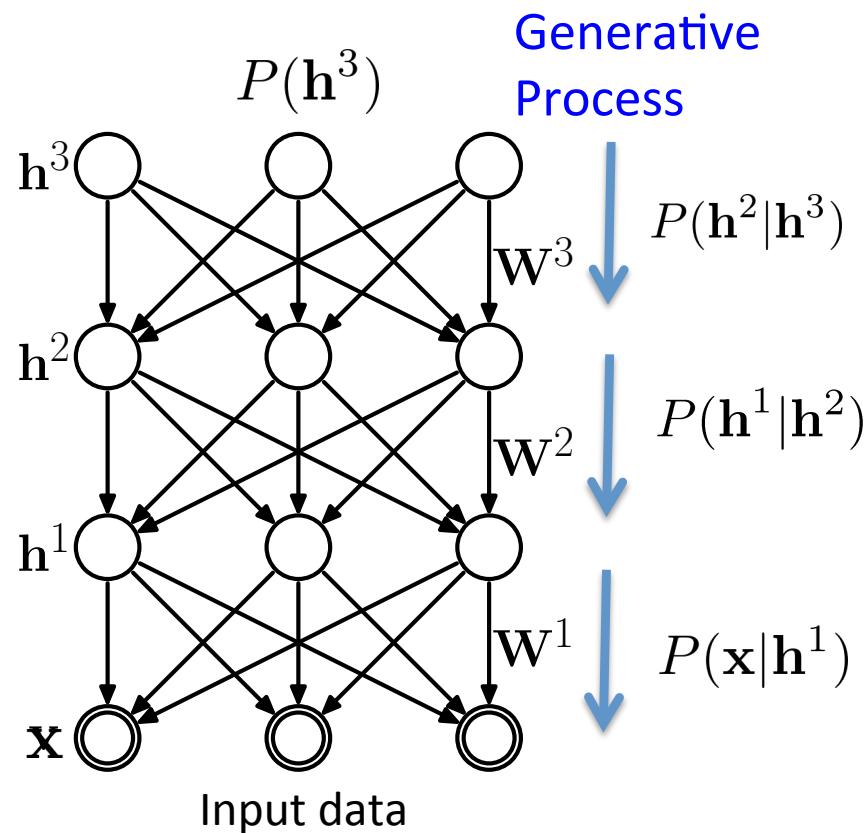
# Variational Autoencoders (VAEs)

- 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 term may denote a complicated nonlinear relationship



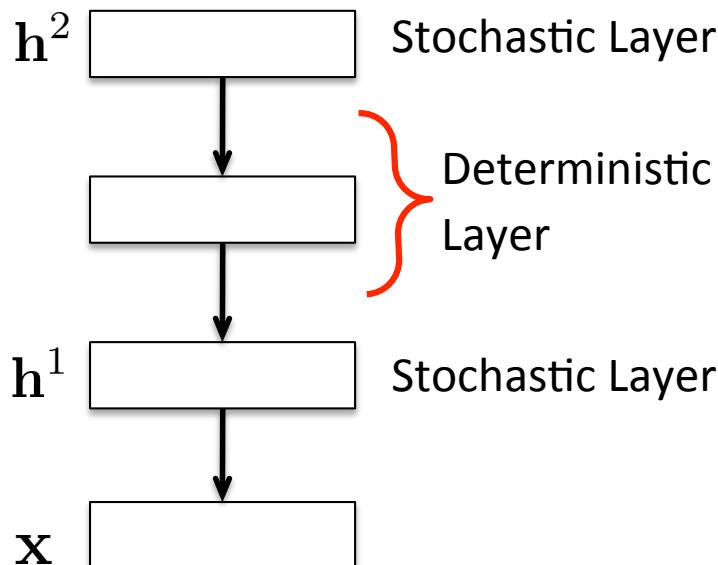
- $\boldsymbol{\theta}$  denotes parameters of VAE.
- $L$  is the number of **stochastic** layers.
- Sampling and probability evaluation is tractable for each  $p(\mathbf{h}^\ell|\mathbf{h}^{\ell+1})$ .

# VAE: Example

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

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{\mathbf{h}^1, \mathbf{h}^2} p(\mathbf{h}^2|\boldsymbol{\theta})p(\mathbf{h}^1|\mathbf{h}^2, \boldsymbol{\theta})p(\mathbf{x}|\mathbf{h}^1, \boldsymbol{\theta})$$

This term denotes a one-layer neural net.



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

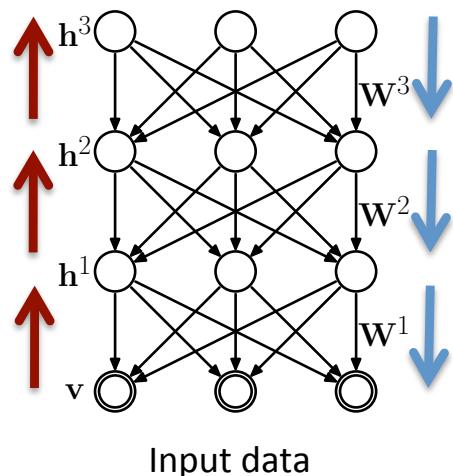
# Variational Bound

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

$$\log p(\mathbf{x}) = \log \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] \geq \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] = \mathcal{L}(\mathbf{x})$$

$$\mathcal{L}(\mathbf{x}) = \log p(\mathbf{x}) - D_{KL}(q(\mathbf{h}|\mathbf{x}))||p(\mathbf{h}|\mathbf{x}))$$

- Trading off the data log-likelihood and the KL divergence from the true posterior.



- Hard to optimize the variational bound with respect to the recognition network (high-variance).
- Key idea of Kingma and Welling is to use reparameterization trick.

# Reparameterization Trick

- Assume that the recognition distribution is Gaussian:

$$q(\mathbf{h}^\ell | \mathbf{h}^{\ell-1}, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta}), \boldsymbol{\Sigma}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta}))$$

with mean and covariance computed from the state of the hidden units at the previous layer.

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

$$\boldsymbol{\epsilon}^\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{h}^\ell (\boldsymbol{\epsilon}^\ell, \mathbf{h}^{\ell-1}, \boldsymbol{\theta}) = \boldsymbol{\Sigma}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta})^{1/2} \boldsymbol{\epsilon}^\ell + \boldsymbol{\mu}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta})$$

# Reparameterization Trick

- Assume that the recognition distribution is Gaussian:

$$q(\mathbf{h}^\ell | \mathbf{h}^{\ell-1}, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta}), \boldsymbol{\Sigma}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta}))$$

- Or

$$\boldsymbol{\epsilon}^\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{h}^\ell (\boldsymbol{\epsilon}^\ell, \mathbf{h}^{\ell-1}, \boldsymbol{\theta}) = \boldsymbol{\Sigma}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta})^{1/2} \boldsymbol{\epsilon}^\ell + \boldsymbol{\mu}(\mathbf{h}^{\ell-1}, \boldsymbol{\theta})$$

- The recognition distribution  $q(\mathbf{h}^\ell | \mathbf{h}^{\ell-1}, \boldsymbol{\theta})$  can be expressed in terms of a deterministic mapping:

$$\underbrace{\mathbf{h}(\boldsymbol{\epsilon}, \mathbf{x}, \boldsymbol{\theta})}_{\text{Deterministic Encoder}}, \quad \text{with} \quad \boldsymbol{\epsilon} = \underbrace{(\boldsymbol{\epsilon}^1, \dots, \boldsymbol{\epsilon}^L)}_{\text{Distribution of } \boldsymbol{\epsilon}}$$

Deterministic  
Encoder

Distribution of  $\boldsymbol{\epsilon}$   
does not depend on  $\boldsymbol{\theta}$

# Computing the Gradients

- The gradient w.r.t the parameters: both recognition and generative:

$$\nabla_{\theta} \mathbb{E}_{\mathbf{h} \sim q(\mathbf{h}|\mathbf{x}, \theta)} \left[ \log \frac{p(\mathbf{x}, \mathbf{h}|\theta)}{q(\mathbf{h}|\mathbf{x}, \theta)} \right]$$

Autoencoder  
↓

$$= \nabla_{\theta} \mathbb{E}_{\epsilon^1, \dots, \epsilon^L \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h}(\epsilon, \mathbf{x}, \theta)|\theta)}{q(\mathbf{h}(\epsilon, \mathbf{x}, \theta)|\mathbf{x}, \theta)} \right]$$

$$= \mathbb{E}_{\epsilon^1, \dots, \epsilon^L \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \nabla_{\theta} \log \frac{p(\mathbf{x}, \mathbf{h}(\epsilon, \mathbf{x}, \theta)|\theta)}{q(\mathbf{h}(\epsilon, \mathbf{x}, \theta)|\mathbf{x}, \theta)} \right]$$



Gradients can be computed by backprop

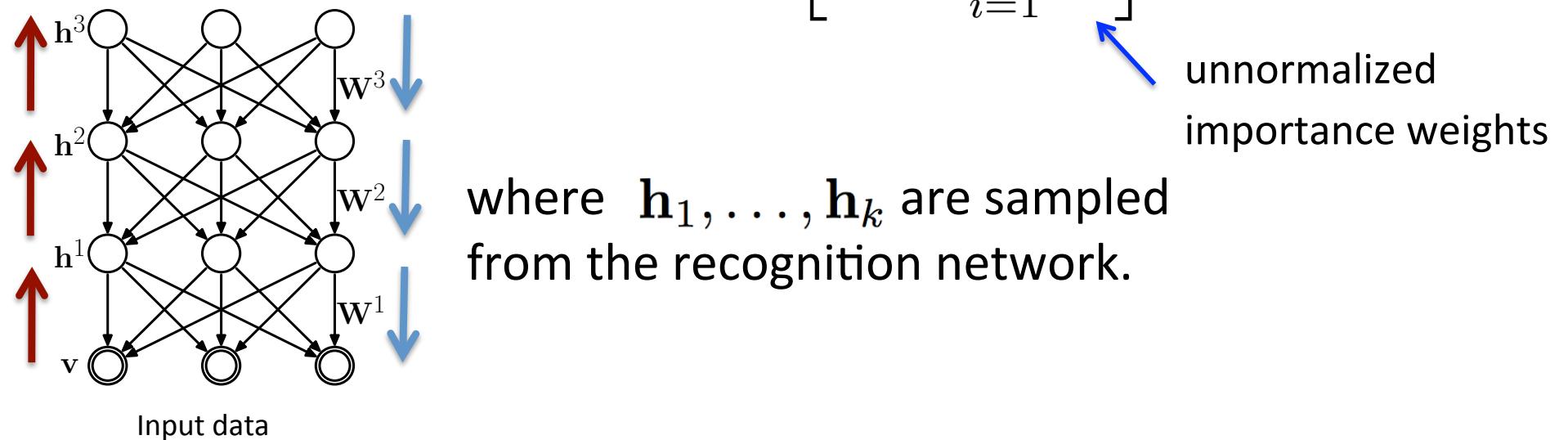
The mapping  $\mathbf{h}$  is a deterministic neural net for fixed  $\epsilon$ .

# Importance Weighted Autoencoders

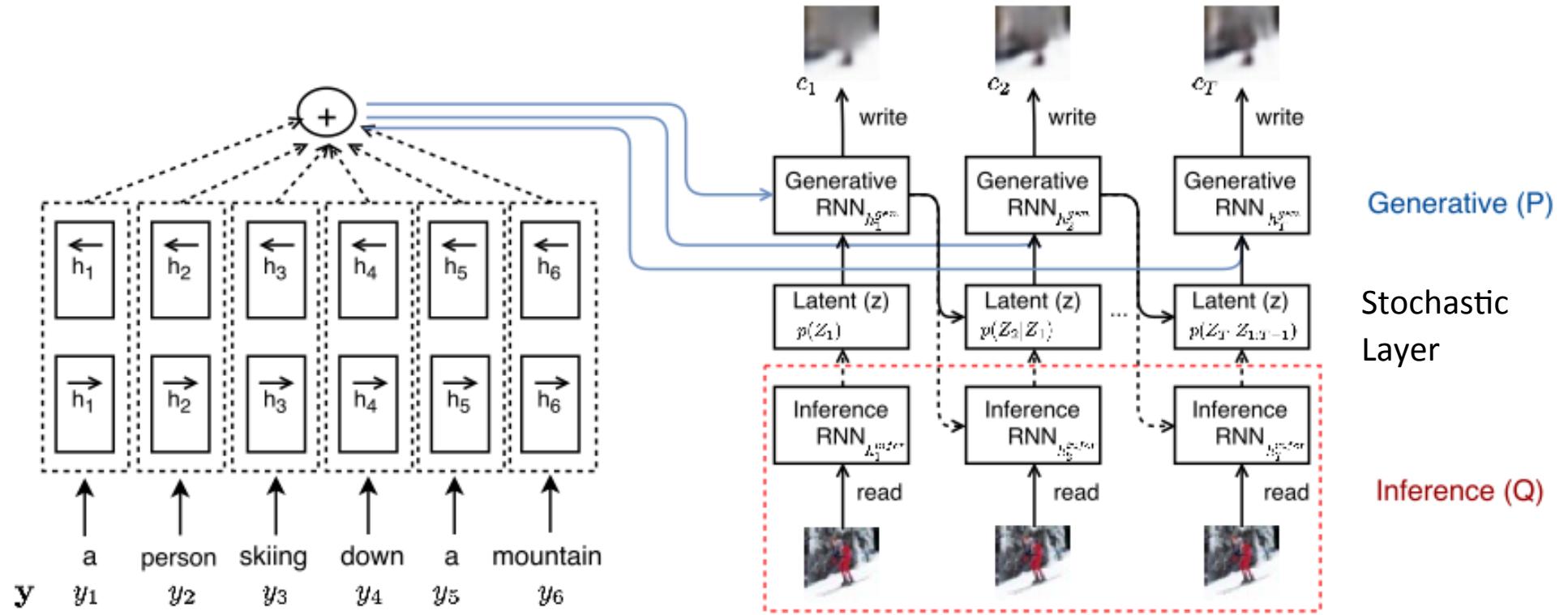
- Can improve VAE by using following k-sample importance weighting of the log-likelihood:

$$\mathcal{L}_k(\mathbf{x}) = \mathbb{E}_{\mathbf{h}_1, \dots, \mathbf{h}_k \sim q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{1}{k} \sum_{i=1}^k \frac{p(\mathbf{x}, \mathbf{h}_i)}{q(\mathbf{h}_i|\mathbf{x})} \right]$$

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



# Generating Images from Captions

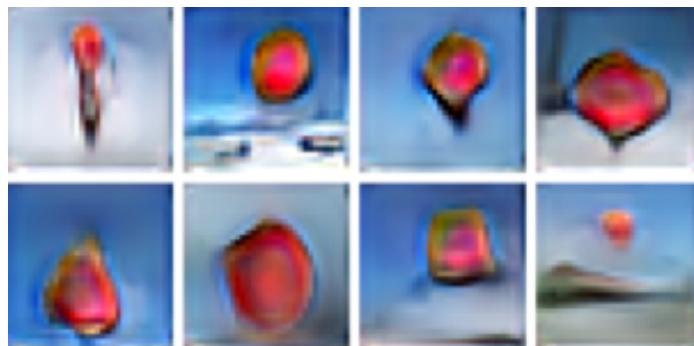


- **Generative Model:** Stochastic Recurrent Network, chained sequence of Variational Autoencoders, with a single stochastic layer.
- **Recognition Model:** Deterministic Recurrent Network.

# Motivating Example

- Can we generate images from natural language descriptions?

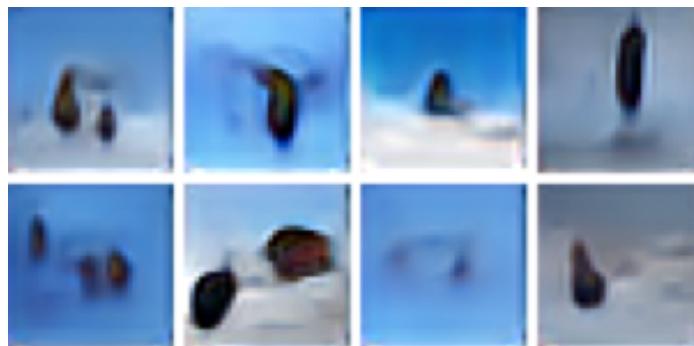
A **stop sign** is flying in blue skies



A **pale yellow school bus** is flying in blue skies



A **herd of elephants** is flying in blue skies



A **large commercial airplane** is flying in blue skies



(Mansimov, Parisotto, Ba, Salakhutdinov, 2015)

# Flipping Colors

A **yellow school bus** parked in the parking lot



A **red school bus** parked in the parking lot



A **green school bus** parked in the parking lot



A **blue school bus** parked in the parking lot



(Mansimov, Parisotto, Ba, Salakhutdinov, 2015)

# Qualitative Comparison

*A group of people walk on a beach with surf boards*

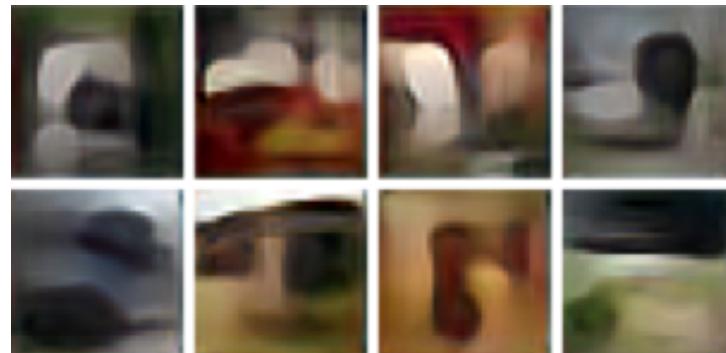
Our Model



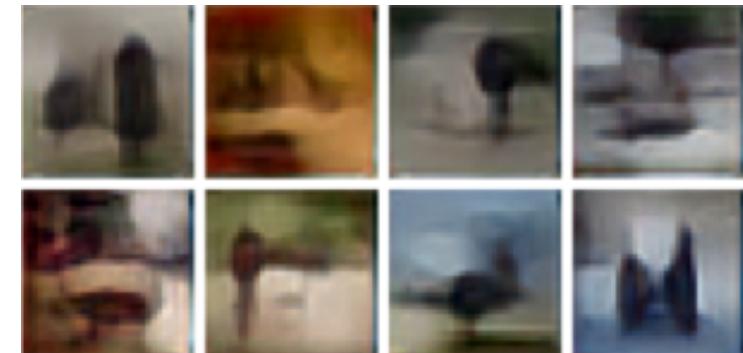
LAPGAN (Denton et. al. 2015)



Conv-Deconv VAE

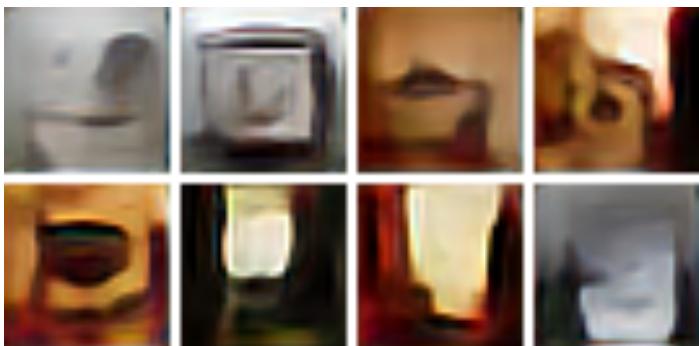


Fully Connected VAE



# Novel Scene Compositions

A toilet seat sits open in the bathroom



A toilet seat sits open in the grass field



Ask Google?



# Neural Story Telling

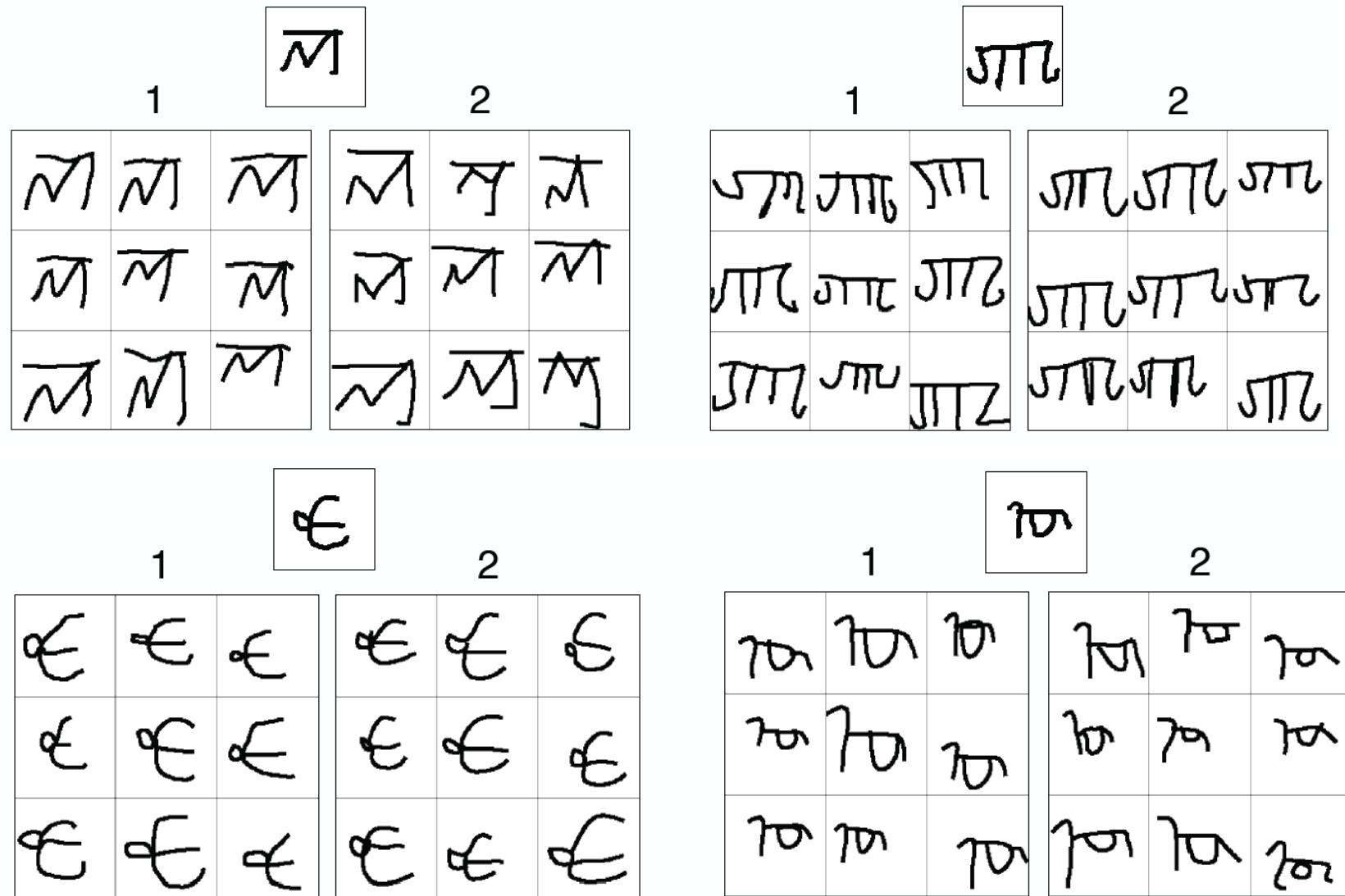


**Sample from the Generative Model  
(recurrent neural network):**

We were barely able to catch the breeze at the beach, and it felt as if someone stepped out of my mind.

She was in love with him for the first time in months, so she had no intention of escaping. The sun had risen from the ocean, making her feel more alive than normal . She is beautiful, but the truth is that I do not know what to do. The sun was just starting to fade away, leaving people scattered around the Atlantic Ocean.

# One-Shot Learning: Humans vs. Machines



(Lake, Salakhutdinov, Tenenbaum, Science, 2015)

# Talk Roadmap

- Basic Building Blocks:

- Sparse Coding
- Autoencoders

- Deep Generative Models

- Restricted Boltzmann Machines
- Deep Boltzmann Machines
- Helmholtz Machines / Variational Autoencoders

- Generative Adversarial Networks

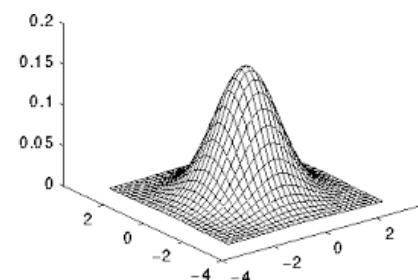
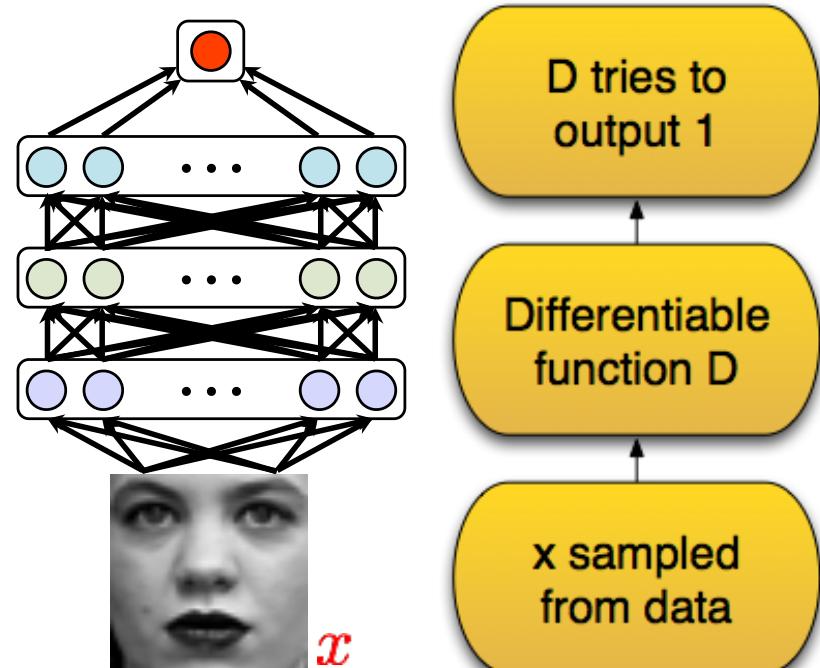
# Generative Adversarial Networks

- There is no explicit definition of the density for  $p(x)$  – Only need to be able to sample from it.
- No variational learning, no maximum-likelihood estimation, no MCMC. How?
- By playing a game!

# Generative Adversarial Networks

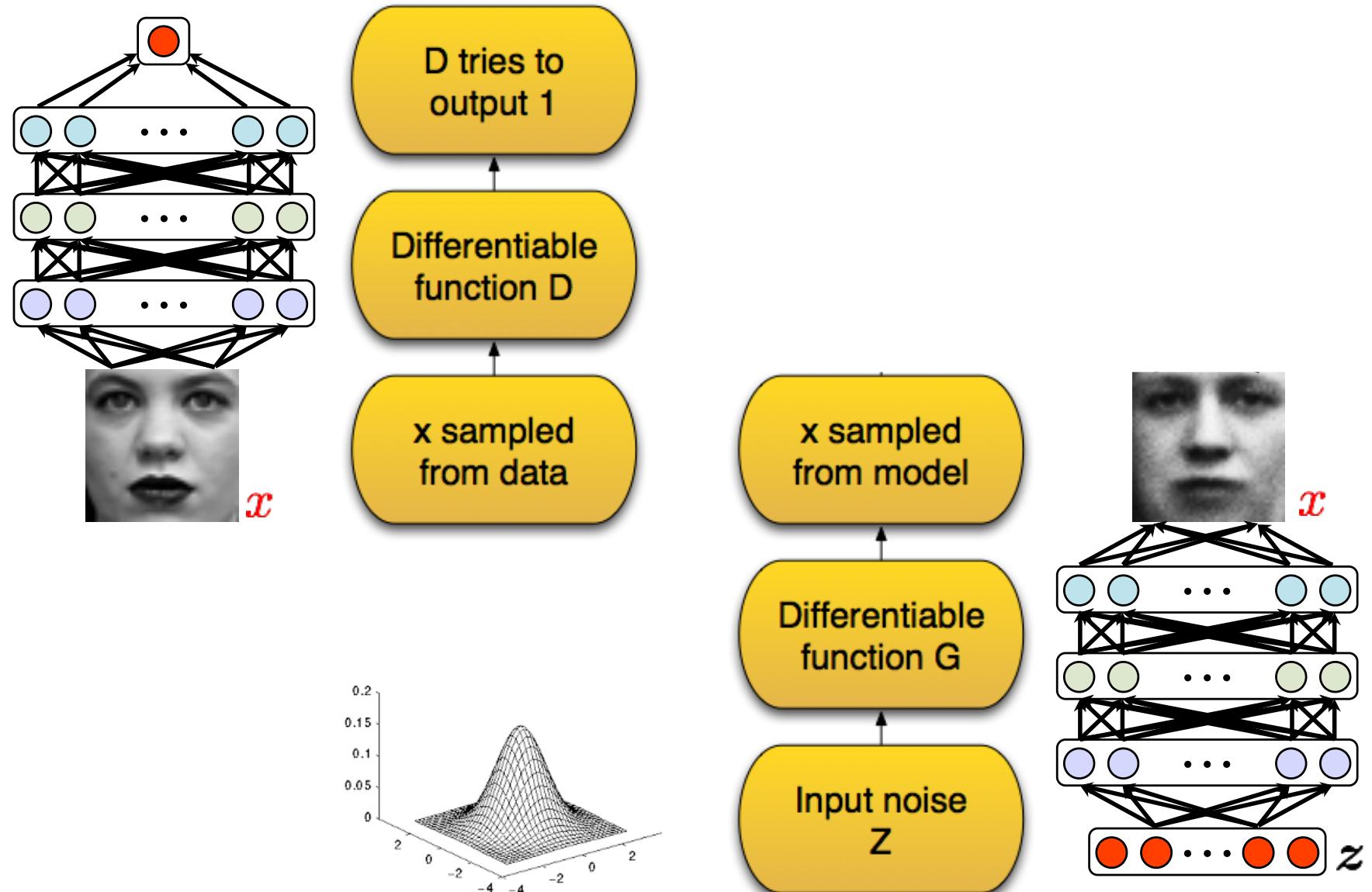
- Set up a game between two players:
  - Discriminator D
  - Generator G
- **Discriminator D** tries to discriminate between:
  - A sample from the data distribution.
  - And a sample from the generator G.
- The **Generator G** attempts to “fool” D by generating samples that are hard for D to distinguish from the real data.

# Generative Adversarial Networks



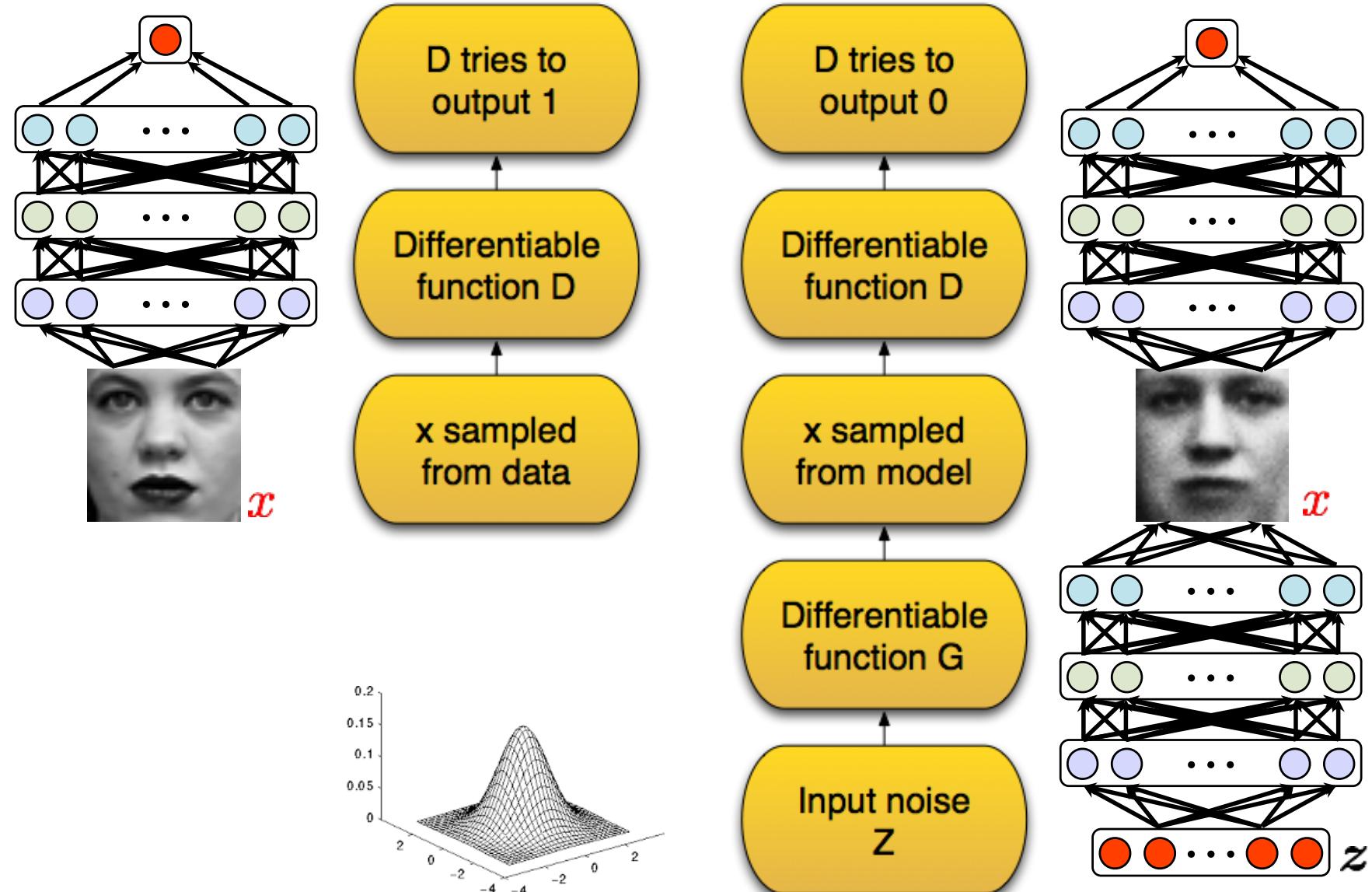
Slide Credit: Ian Goodfellow

# Generative Adversarial Networks



Slide Credit: Ian Goodfellow

# Generative Adversarial Networks



Slide Credit: Ian Goodfellow

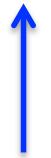
# Generative Adversarial Networks

- Minimax value function

Generator: generate samples  
that D would classify as real



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



Discriminator:  
Pushes up



Discriminator: Classify  
data as being real



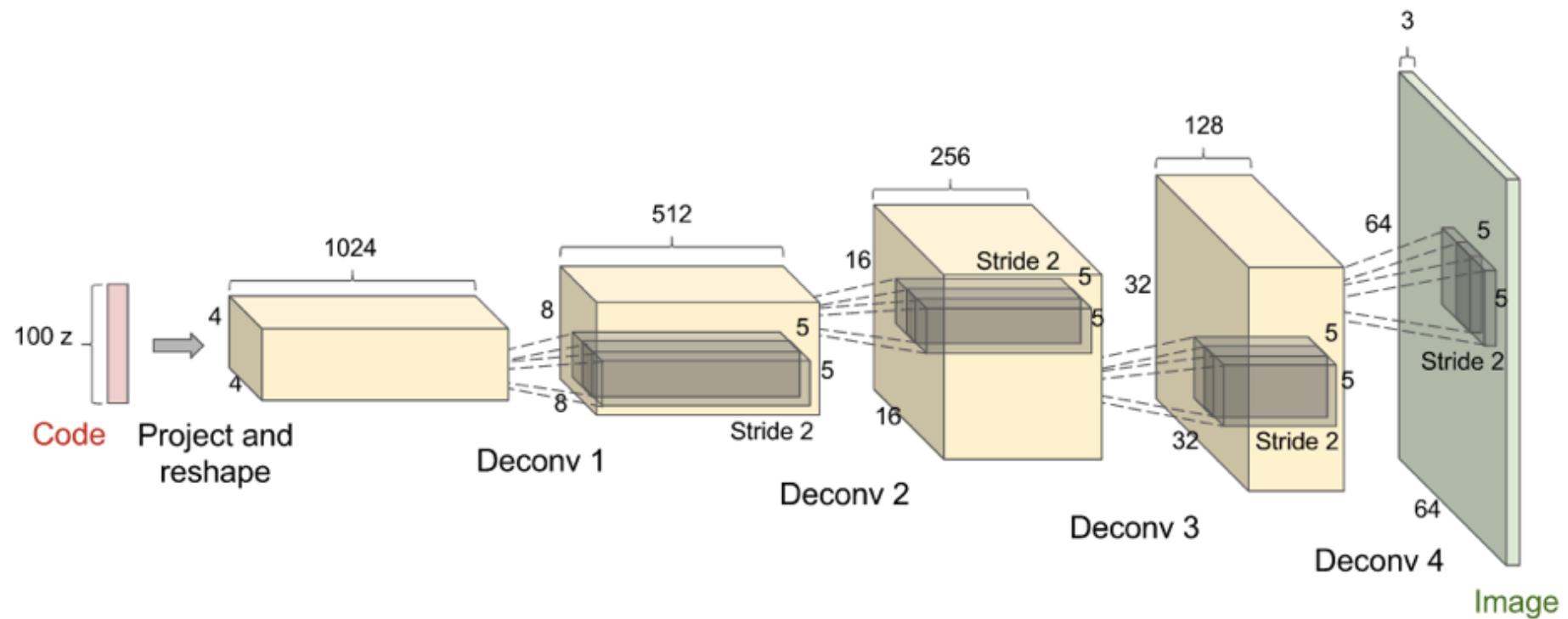
Discriminator: Classify  
generator samples as  
being fake

Generator:  
Pushes down

- Optimal strategy for Discriminator is:

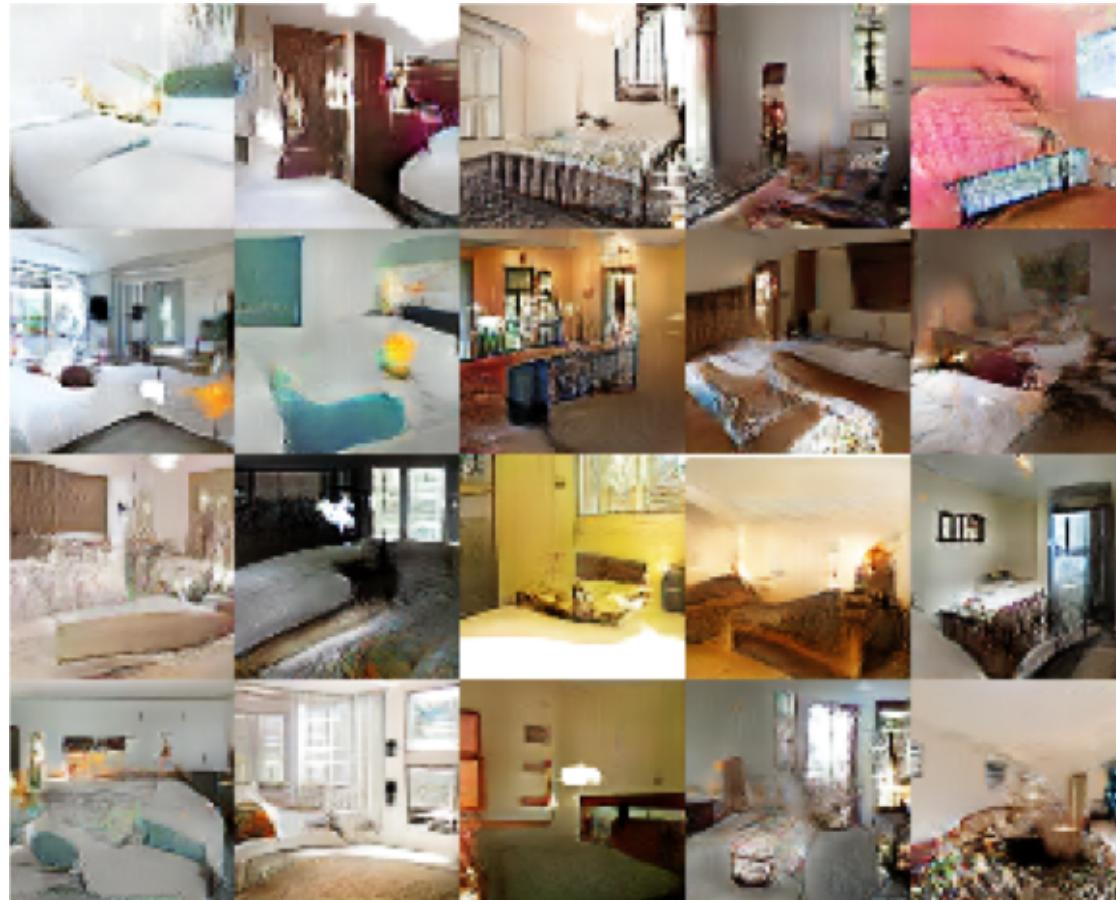
$$D(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\text{model}}(x)}$$

# DCGAN Architecture



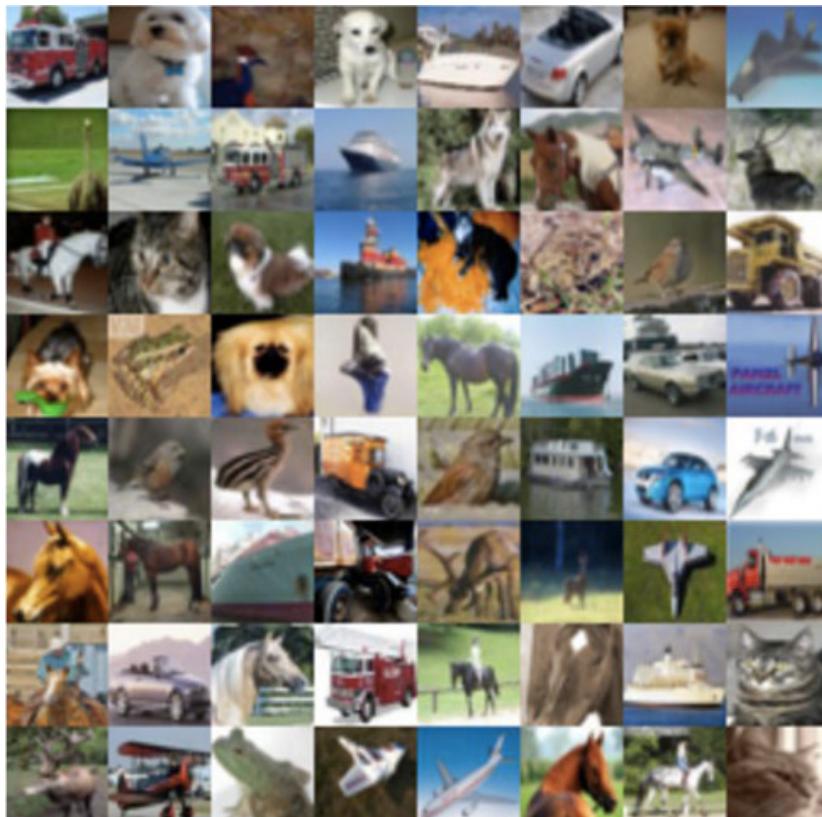
(Radford et al 2015)

# LSUN Bedrooms: Samples

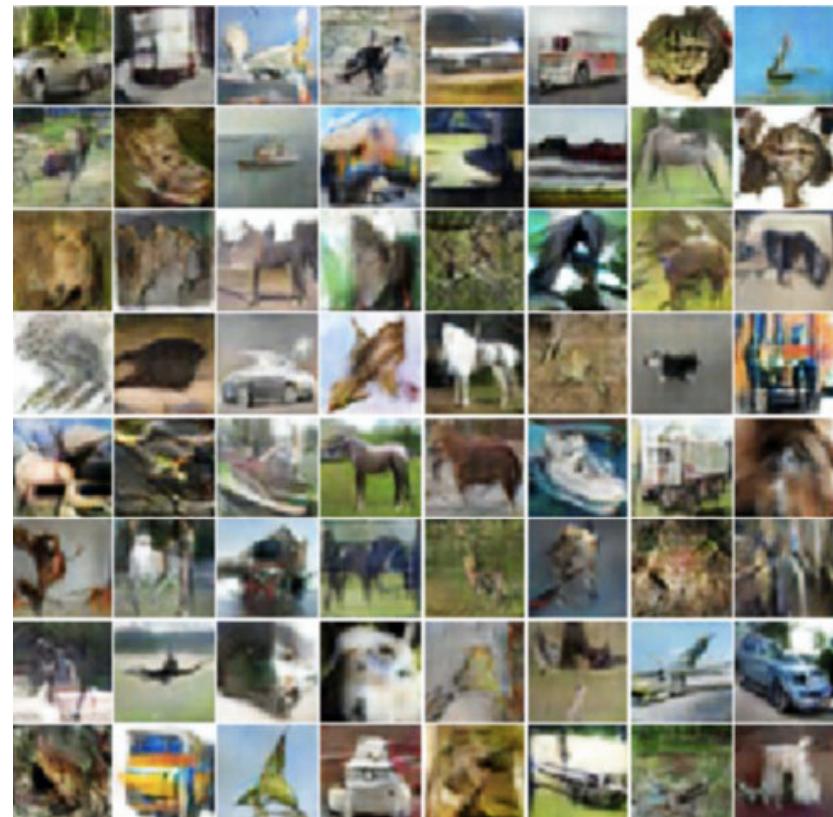


(Radford et al 2015)

# CIFAR



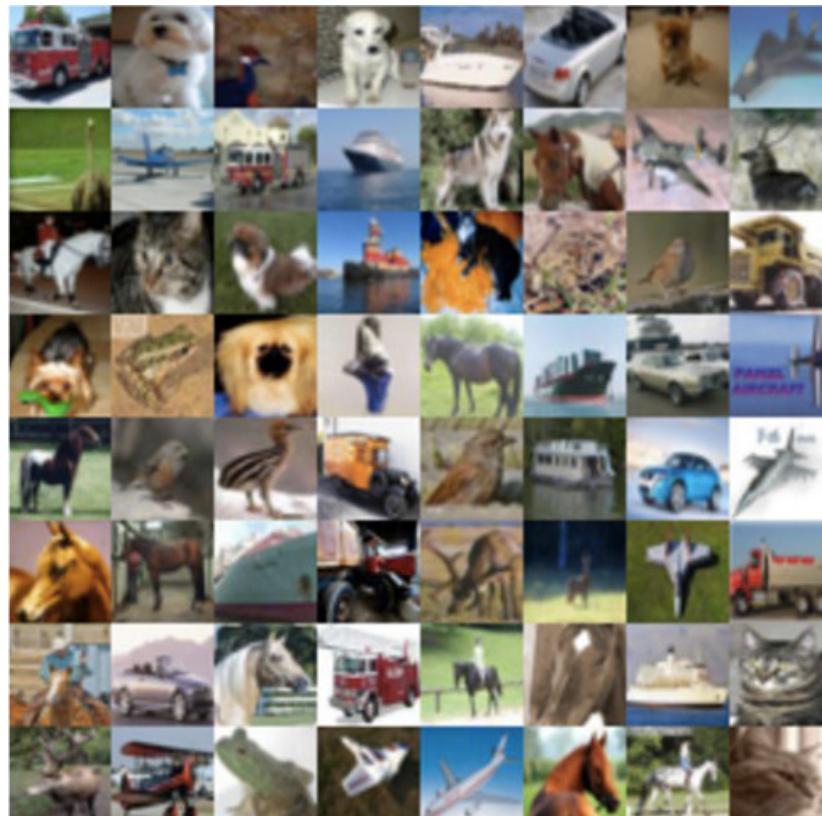
Training



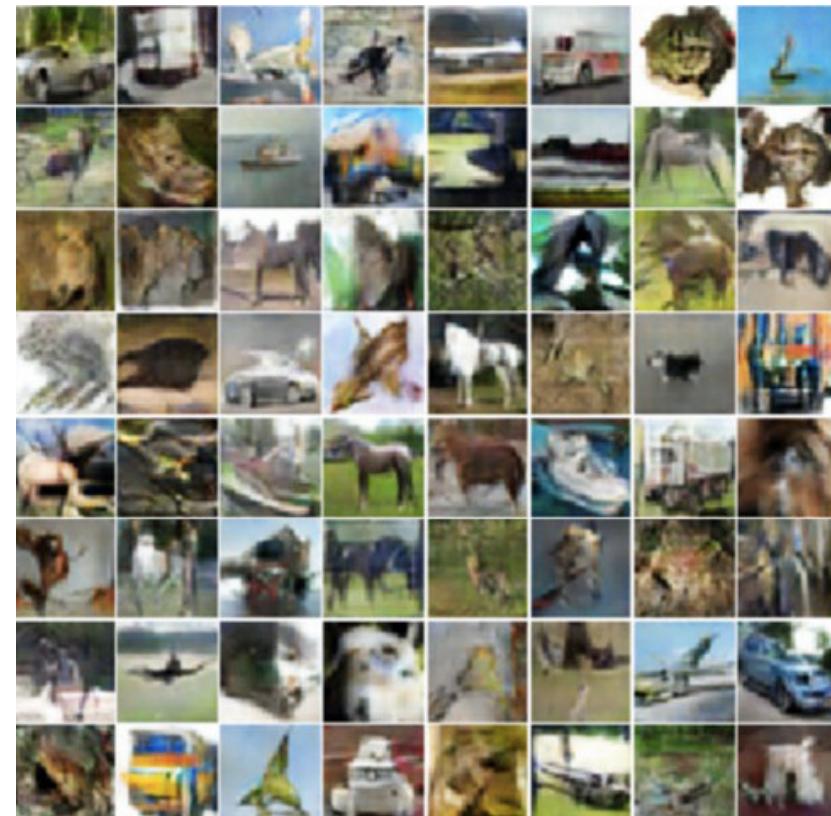
Samples

(Salimans et. al., 2016)

# IMAGENET



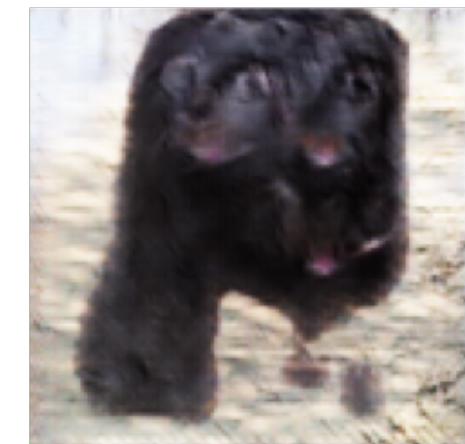
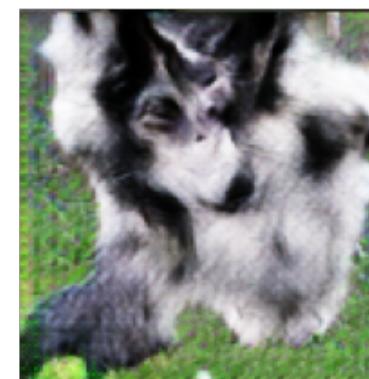
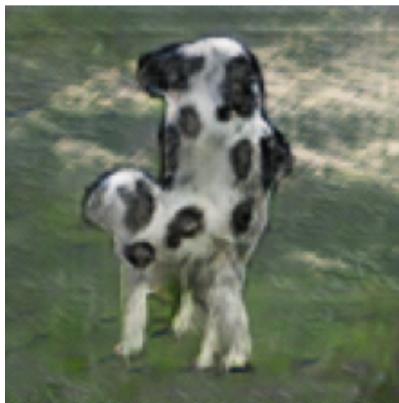
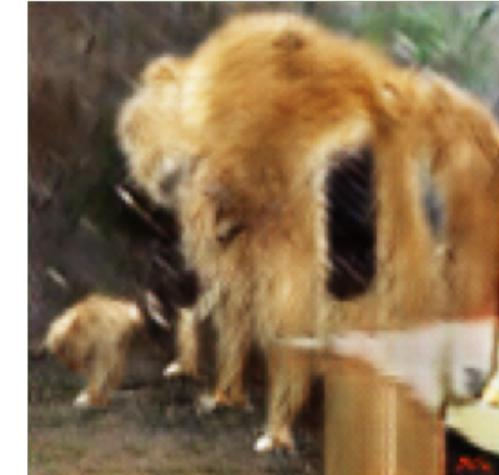
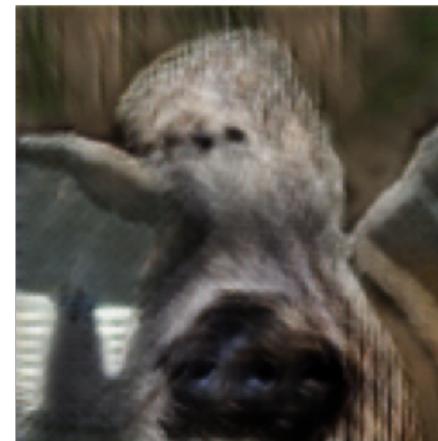
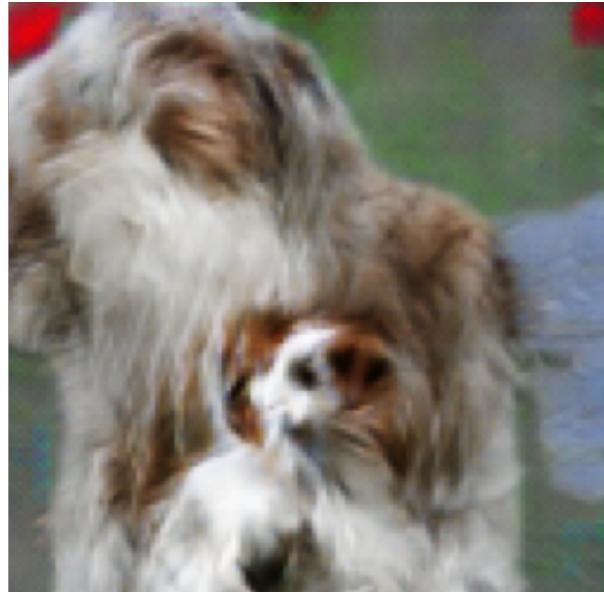
Training



Samples

(Salimans et. al., 2016)

# ImageNet: Cherry-Picked Results



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

Slide Credit: Ian Goodfellow

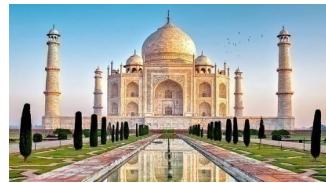
# Summary

- Efficient learning algorithms for Deep Unsupervised Models

Text & image retrieval /  
Object recognition

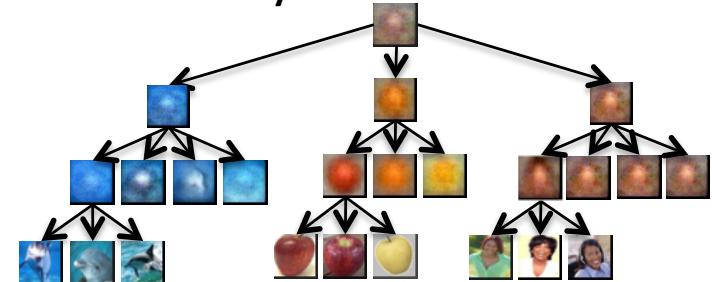


Image Tagging



mosque, tower,  
building, cathedral,  
dome, castle

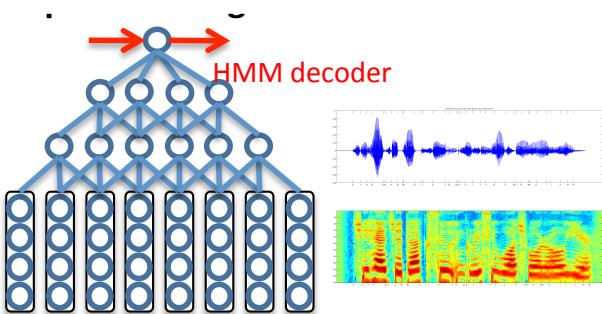
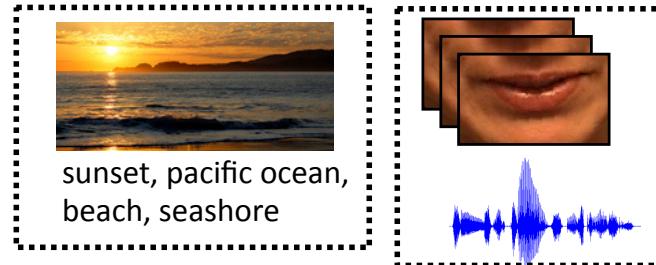
Learning a Category  
Hierarchy



Object Detection



Multimodal Data



- Deep models improve the current state-of-the art in many application domains:
  - Object recognition and detection, text and image retrieval, handwritten character and speech recognition, and others.

Thank you