

# Deep Learning II

## Unsupervised Learning

Russ Salakhutdinov

Machine Learning Department  
Carnegie Mellon University  
Canadian Institute of Advanced Research

**Carnegie  
Mellon  
University**



**CIFAR**  
CANADIAN INSTITUTE  
for ADVANCED RESEARCH

# Talk Roadmap

Part 1: Supervised Learning: Deep Networks

Part 2: Unsupervised Learning: Learning Deep Generative Models

Part 3: Open Research Questions

# Unsupervised Learning

## Non-probabilistic Models

- Sparse Coding
- Autoencoders
- Others (e.g. k-means)

## Probabilistic (Generative) Models

### Tractable Models

- Fully observed Belief Nets
- NADE
- PixelRNN

### Non-Tractable Models

- Boltzmann Machines
- Variational Autoencoders
- Helmholtz Machines
- Many others...

- Generative Adversarial Networks
- Moment Matching Networks

Explicit Density  $p(x)$

Implicit Density

# Talk Roadmap

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



# 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

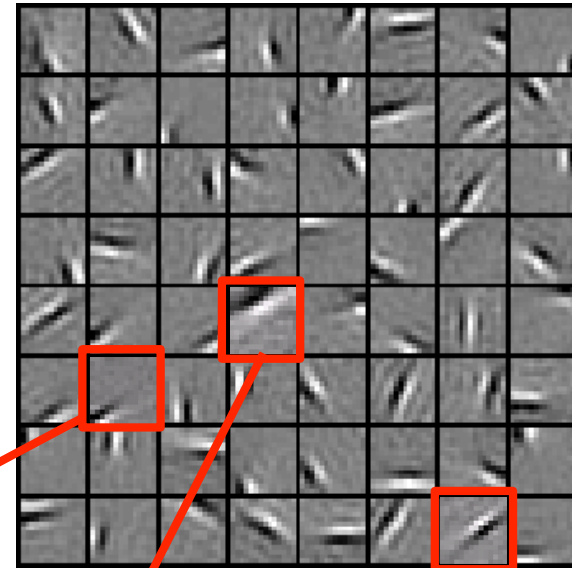
- Each data vector is represented as a sparse linear combination of bases.

# Sparse Coding

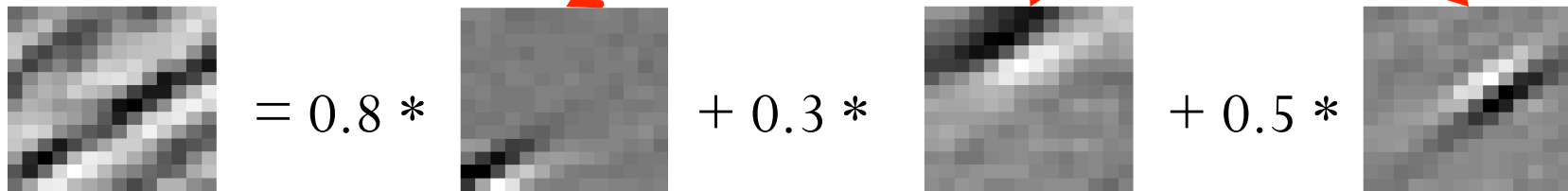
Natural Images



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:  $\phi_1, \phi_2, \dots, \phi_K \in \mathbb{R}^D$

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

- 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  $\mathbf{x}^*$ .

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

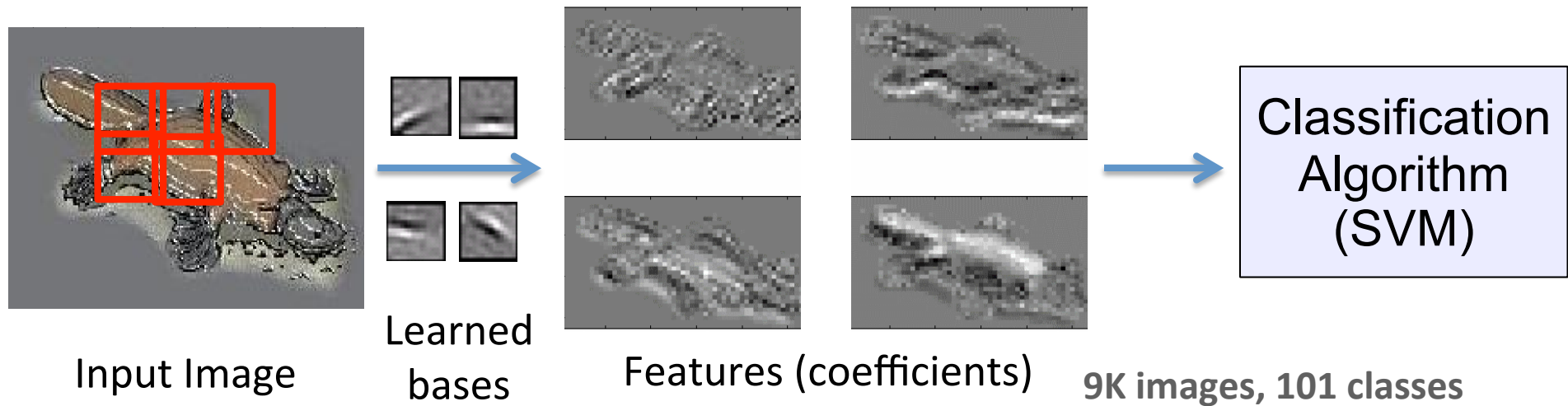


$x^* = 0.8 * \phi_{36} + 0.3 * \phi_{42} + 0.5 * \phi_{65}$

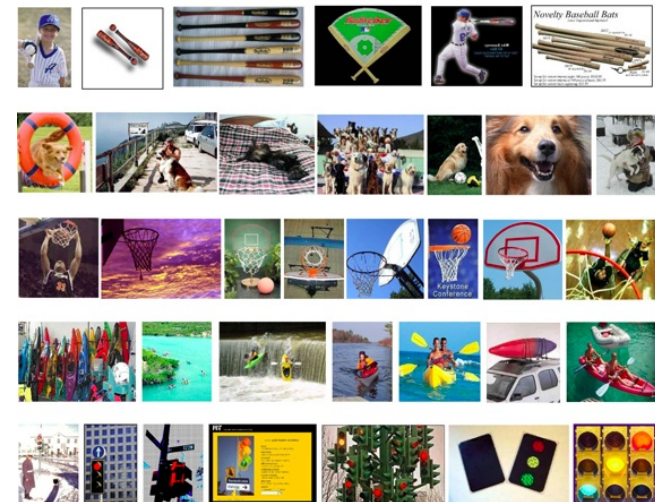
$[0, 0, \dots, \mathbf{0.8}, \dots, \mathbf{0.3}, \dots, \mathbf{0.5}, \dots]$  = coefficients (feature representation)

# Image Classification

Evaluated on Caltech101 object category dataset.



9K images, 101 classes

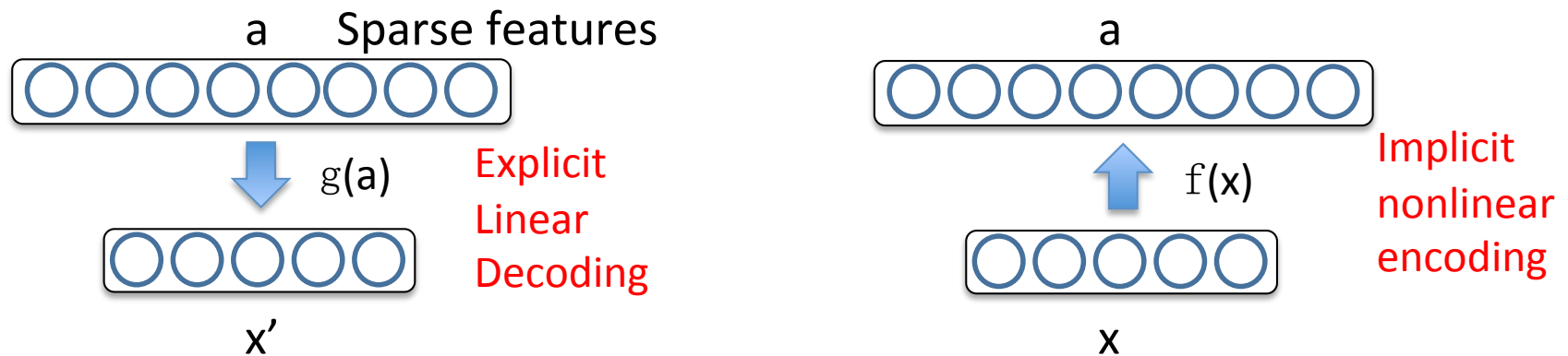


Algorithm	Accuracy
Baseline (Fei-Fei et al., 2004)	16%
PCA	37%
<b>Sparse Coding</b>	<b>47%</b>

(Lee, Battle, Raina, Ng, NIPS 2007)

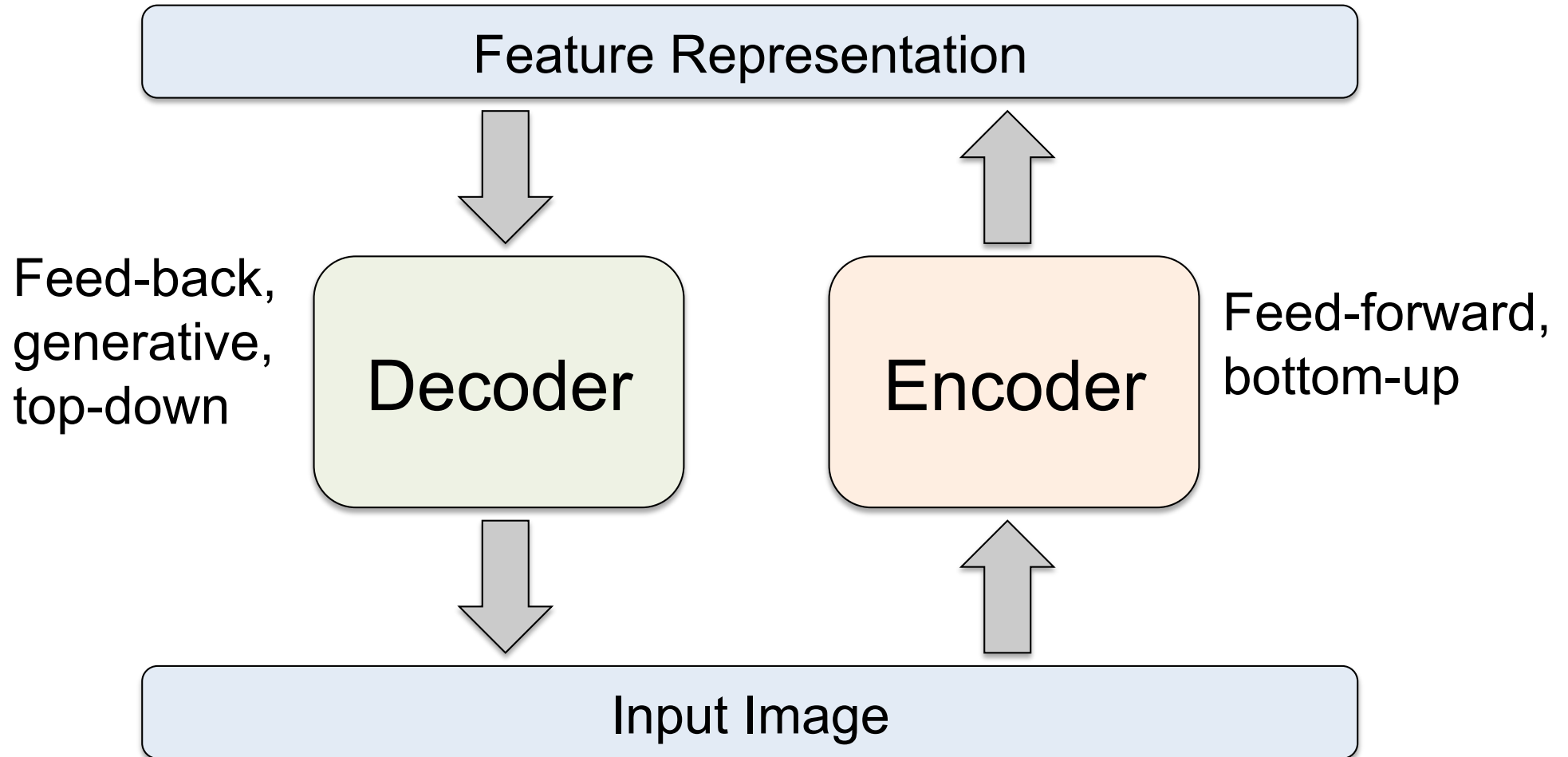
# Interpreting Sparse Coding

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



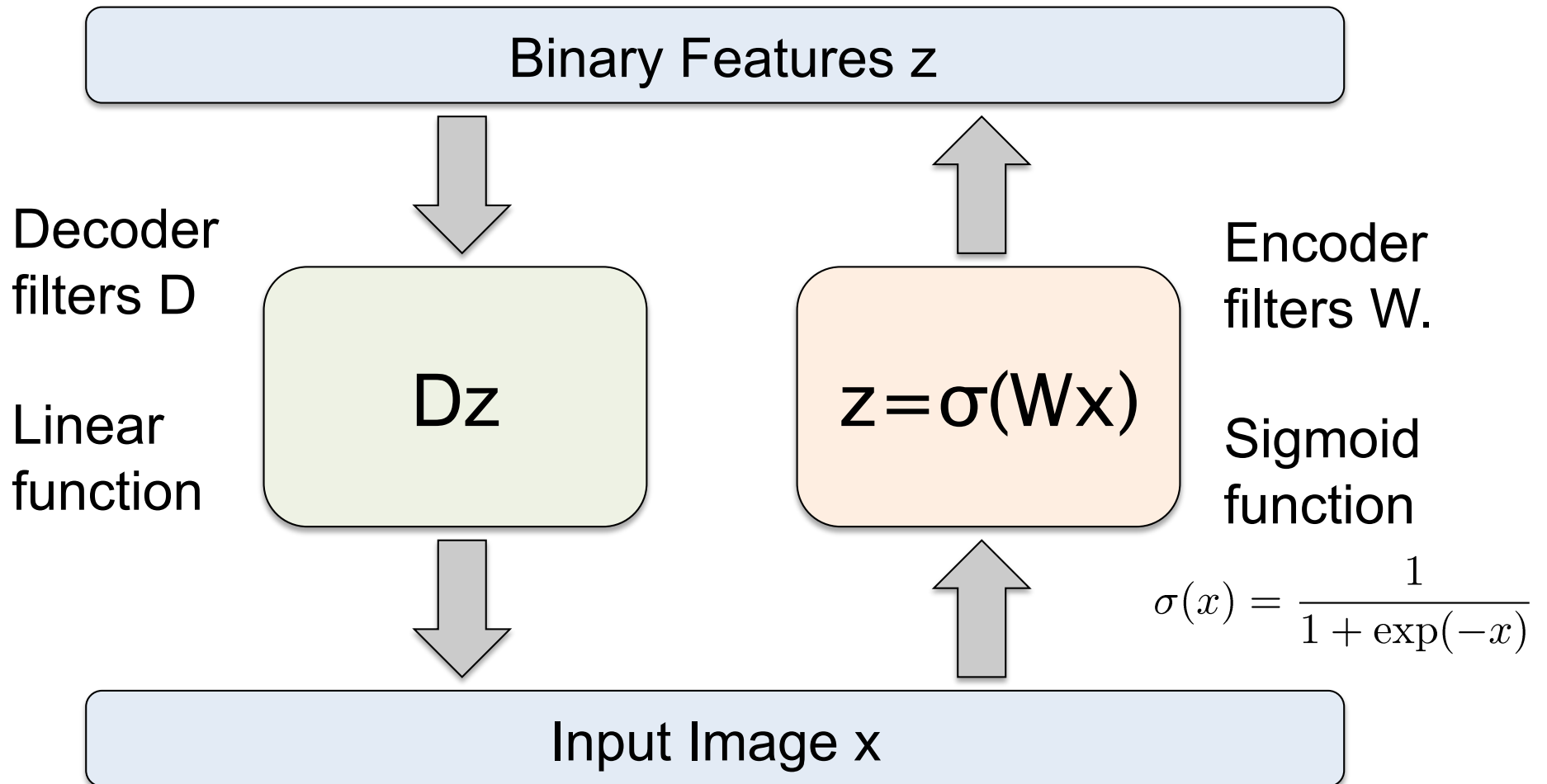
- Sparse, over-complete representation  $\mathbf{a}$ .
- **Encoding**  $\mathbf{a} = f(\mathbf{x})$  is implicit and nonlinear function of  $\mathbf{x}$ .
- **Reconstruction** (or decoding)  $\mathbf{x}' = g(\mathbf{a})$  is linear and explicit.

# Autoencoder



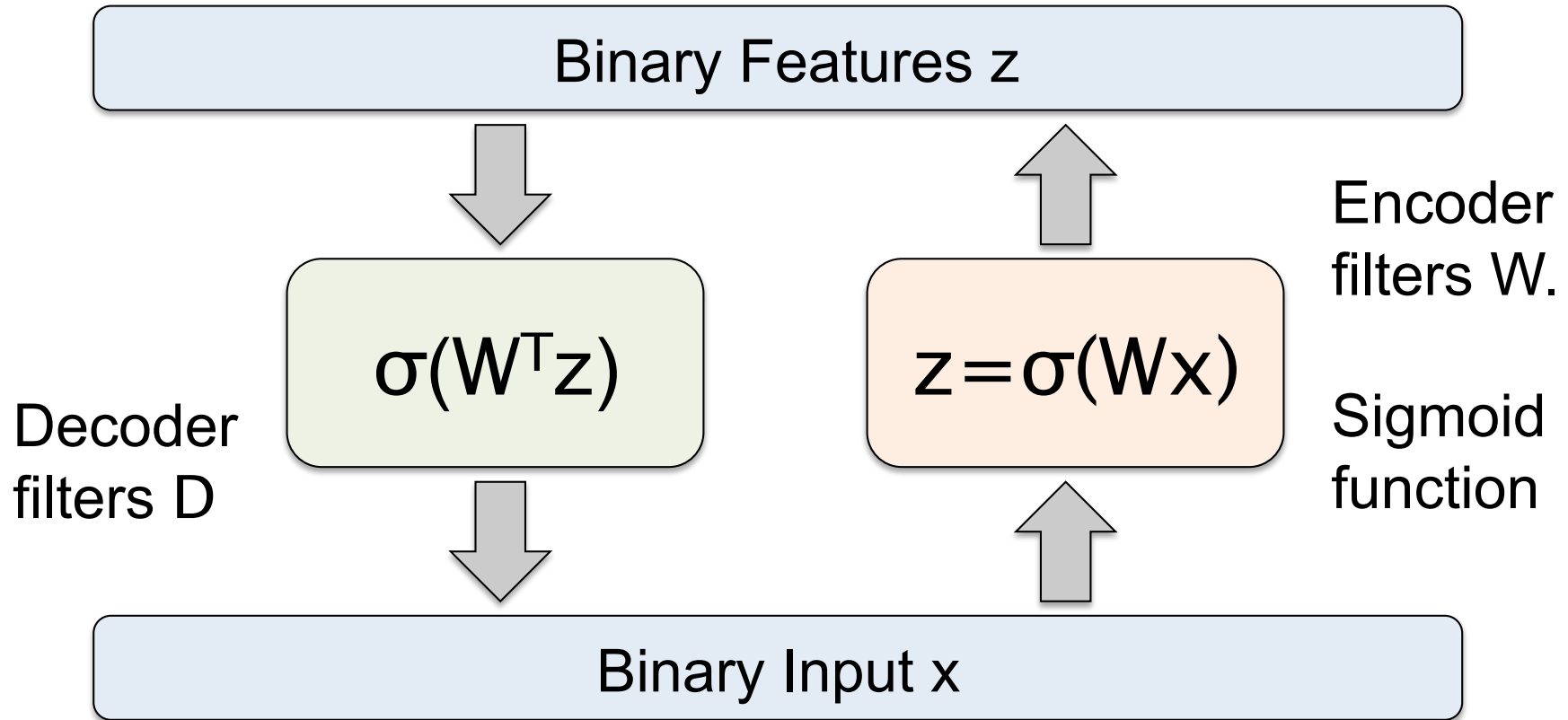
- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity.

# Autoencoder





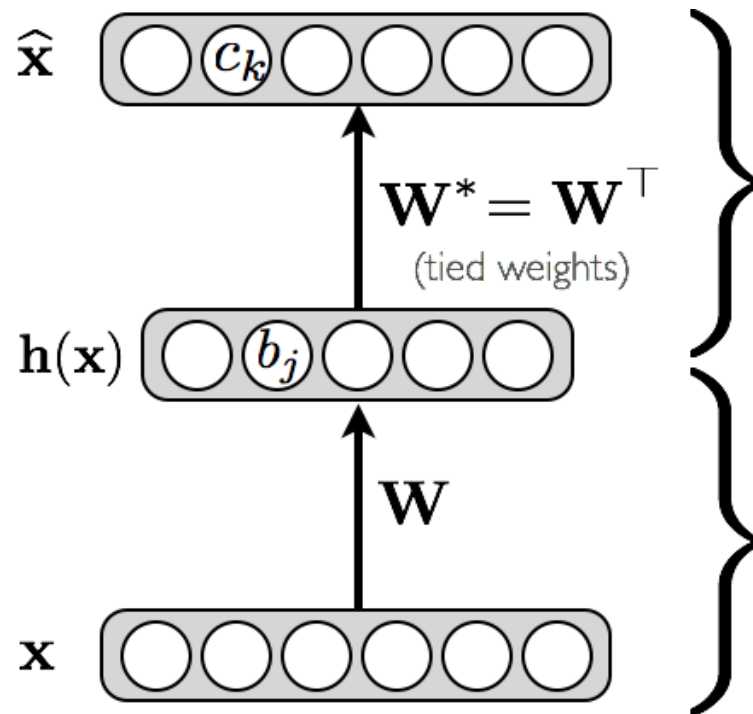
# Autoencoder



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

# Autoencoders

- Feed-forward neural network trained to reproduce its input at the output layer



## Decoder

$$\begin{aligned}\hat{\mathbf{x}} &= o(\hat{\mathbf{a}}(\mathbf{x})) \\ &= \text{sigm}(\underbrace{\mathbf{c} + \mathbf{W}^* \mathbf{h}(\mathbf{x})}_{\text{For binary units}})\end{aligned}$$

## Encoder

$$\begin{aligned}\mathbf{h}(\mathbf{x}) &= g(\mathbf{a}(\mathbf{x})) \\ &= \text{sigm}(\mathbf{b} + \mathbf{W}\mathbf{x})\end{aligned}$$

# Loss Function

- **Loss function** for binary inputs

$$l(f(\mathbf{x})) = - \sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$$

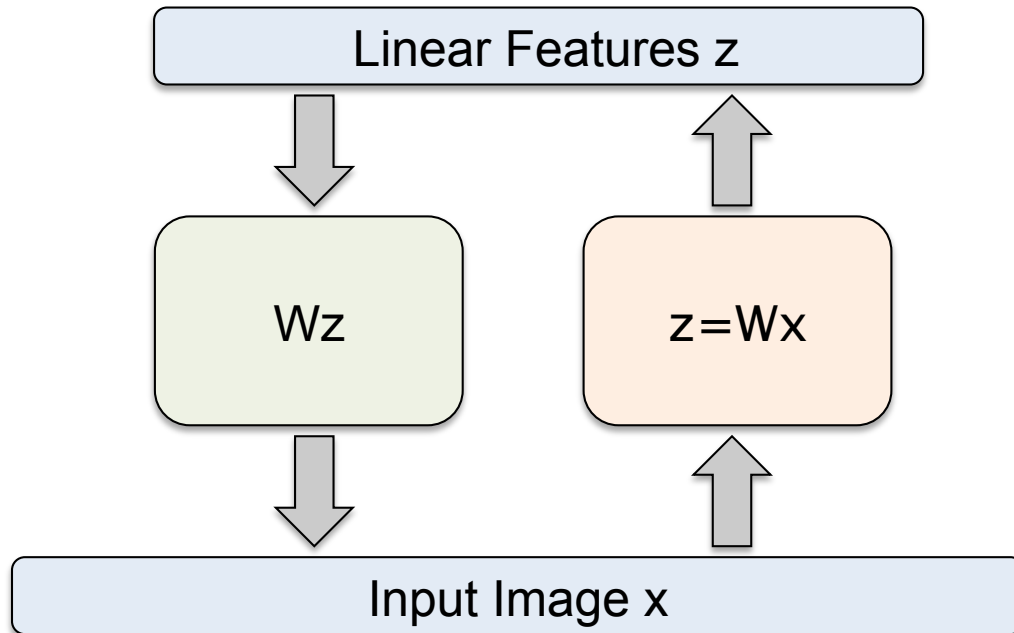
- Cross-entropy error function (reconstruction loss)  $f(\mathbf{x}) \equiv \hat{\mathbf{x}}$

- **Loss function** for real-valued inputs

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

- sum of squared differences (reconstruction loss)
- we use a linear activation function at the output

# Autoencoder



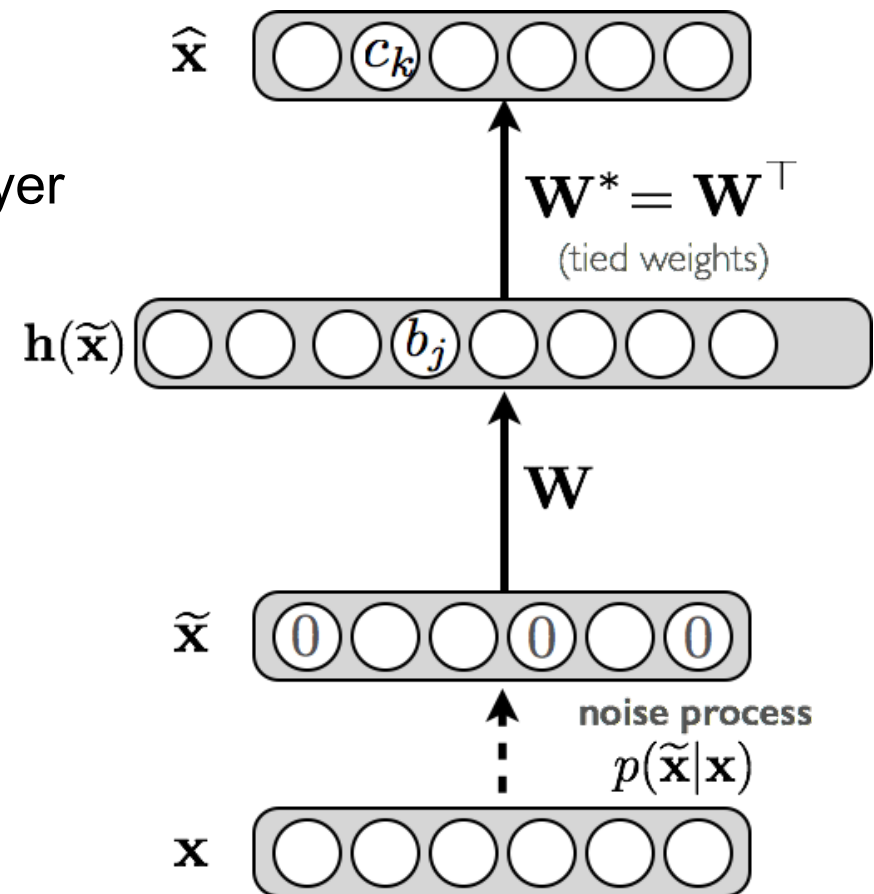
- If the **hidden and output layers are linear**, it will learn hidden units that are a linear function of the data and minimize the squared error.
- The  $K$  hidden units will span the same space as the first  $k$  principal components. The weight vectors may not be orthogonal.

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

# Denoising Autoencoder

- **Idea**: representation should be robust to introduction of noise:
  - random assignment of subset of inputs to 0, with probability  $\nu$
  - Similar to dropouts on the input layer
  - Gaussian additive noise

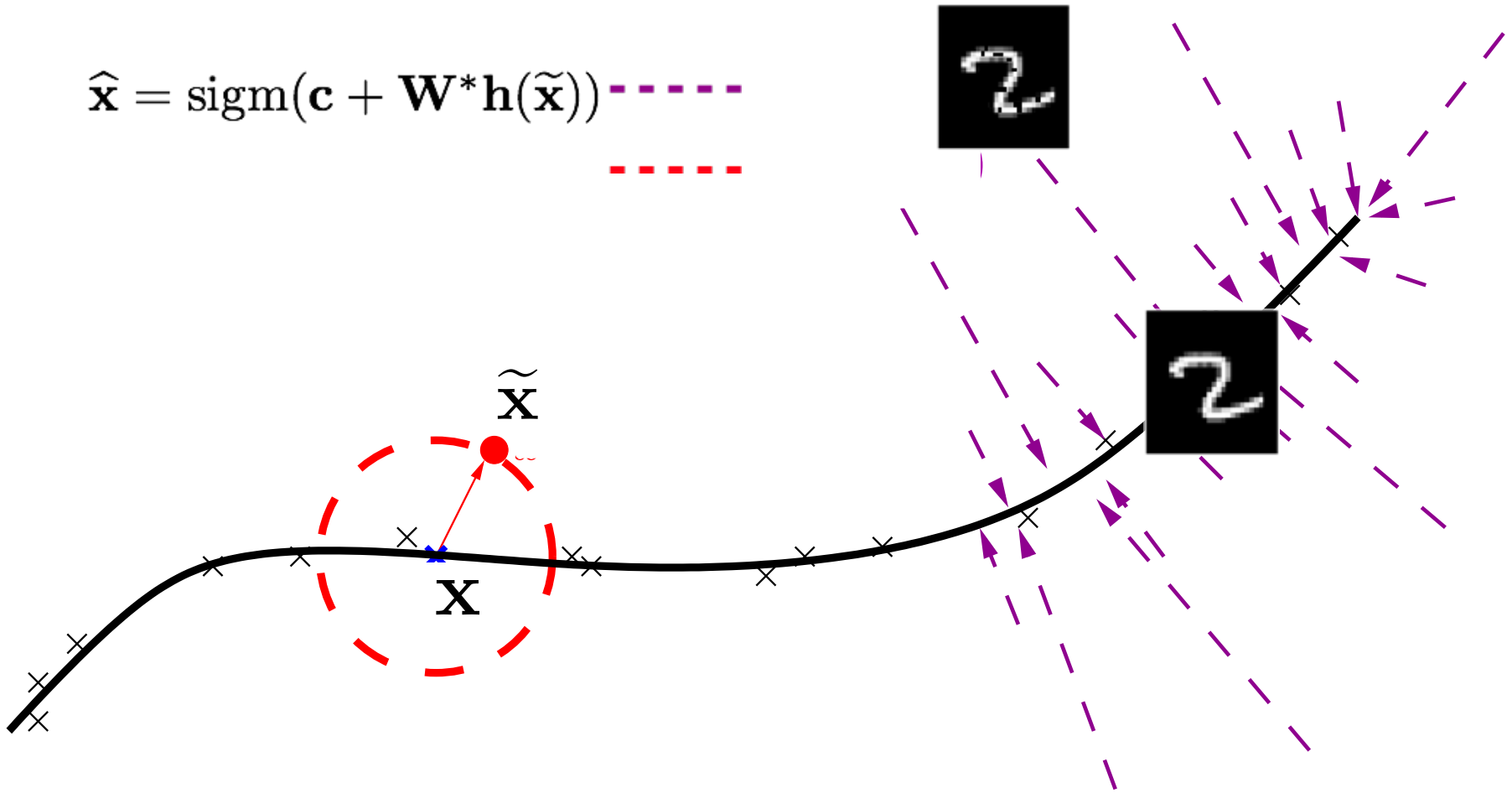
- **Reconstruction**  $\hat{\mathbf{x}}$  computed from the corrupted input  $\tilde{\mathbf{x}}$
- **Loss function** compares  $\hat{\mathbf{x}}$  reconstruction with the noiseless input  $\mathbf{x}$



(Vincent et al., ICML 2008)

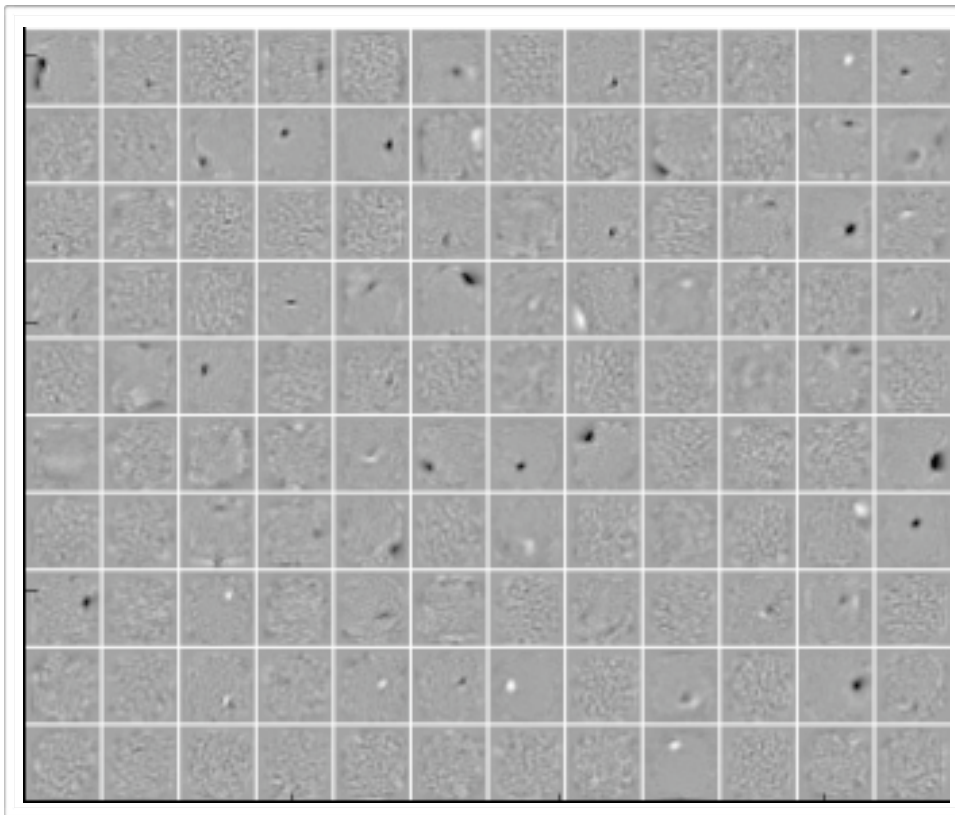
# Denoising Autoencoder

$$\hat{\mathbf{x}} = \text{sigm}(\mathbf{c} + \mathbf{W}^* \mathbf{h}(\tilde{\mathbf{x}}))$$

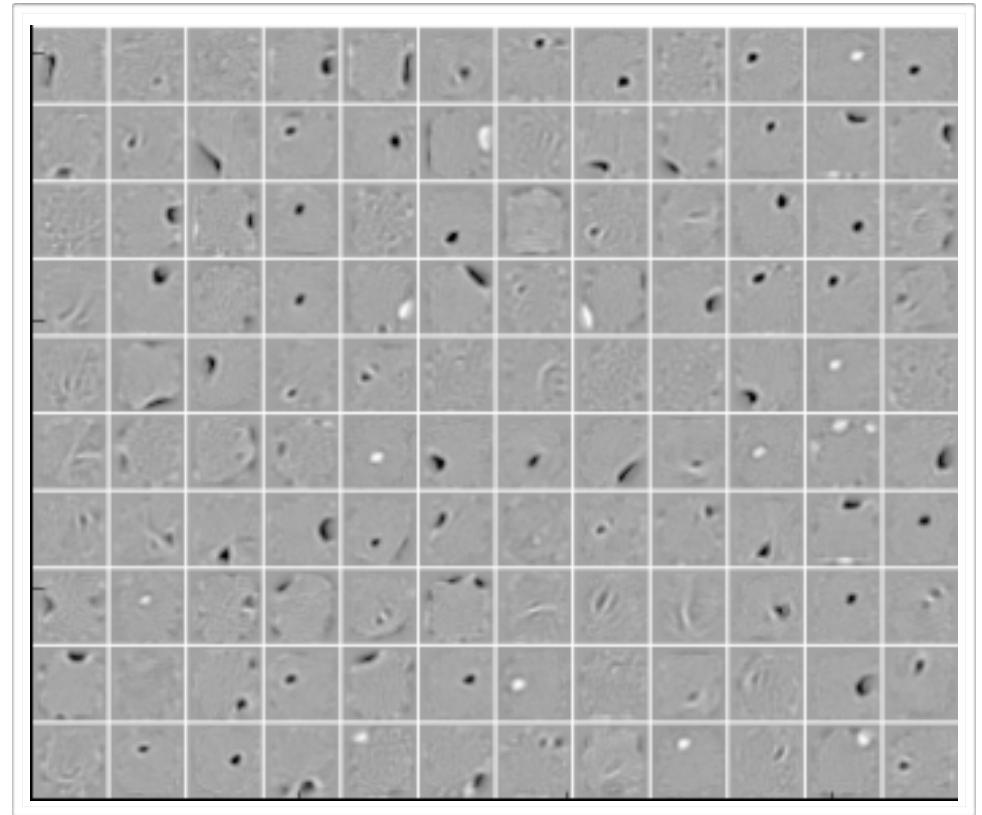


# Learned Filters

Non-corrupted

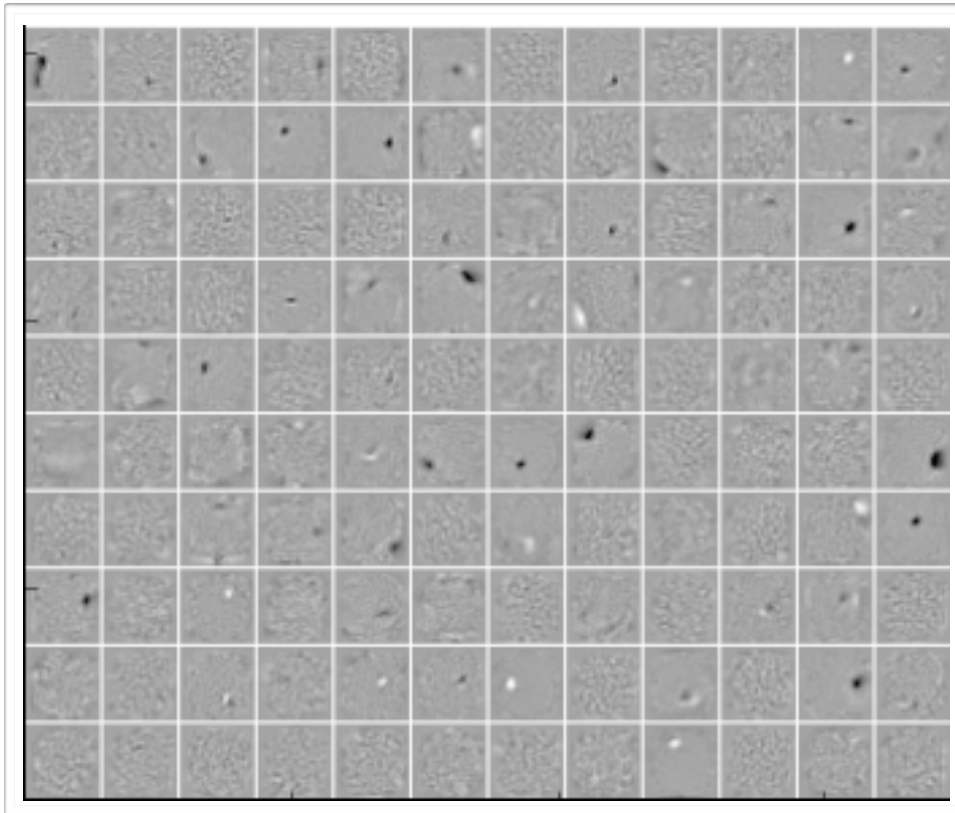


25% corrupted input

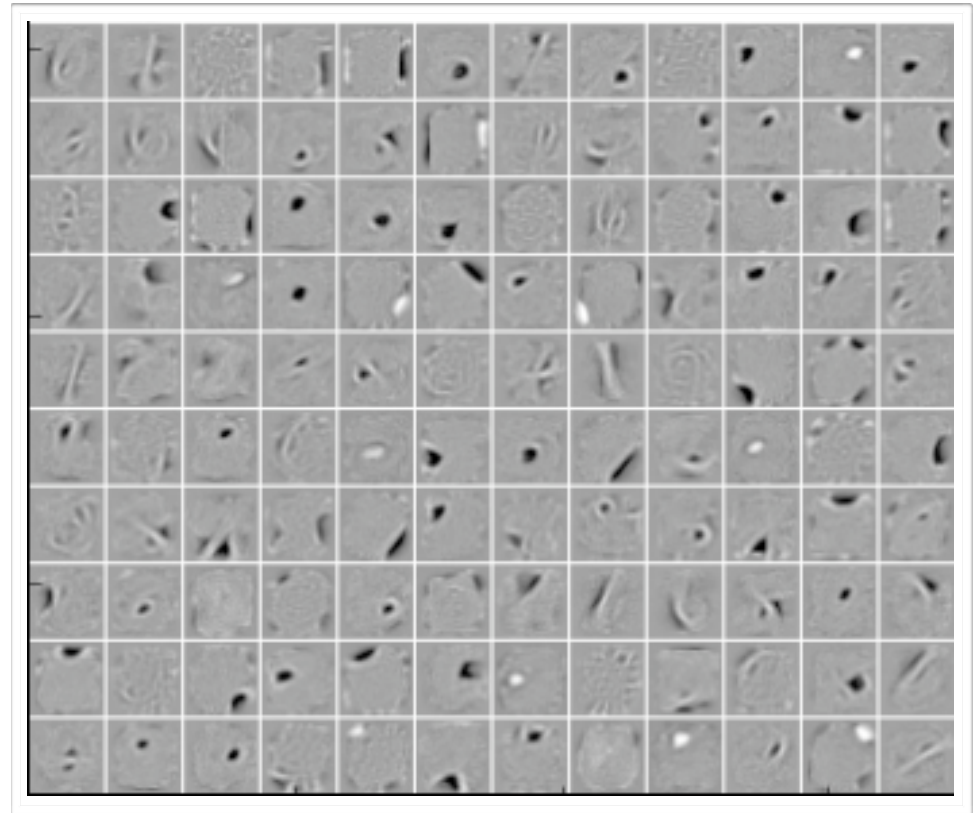


# Learned Filters

Non-corrupted



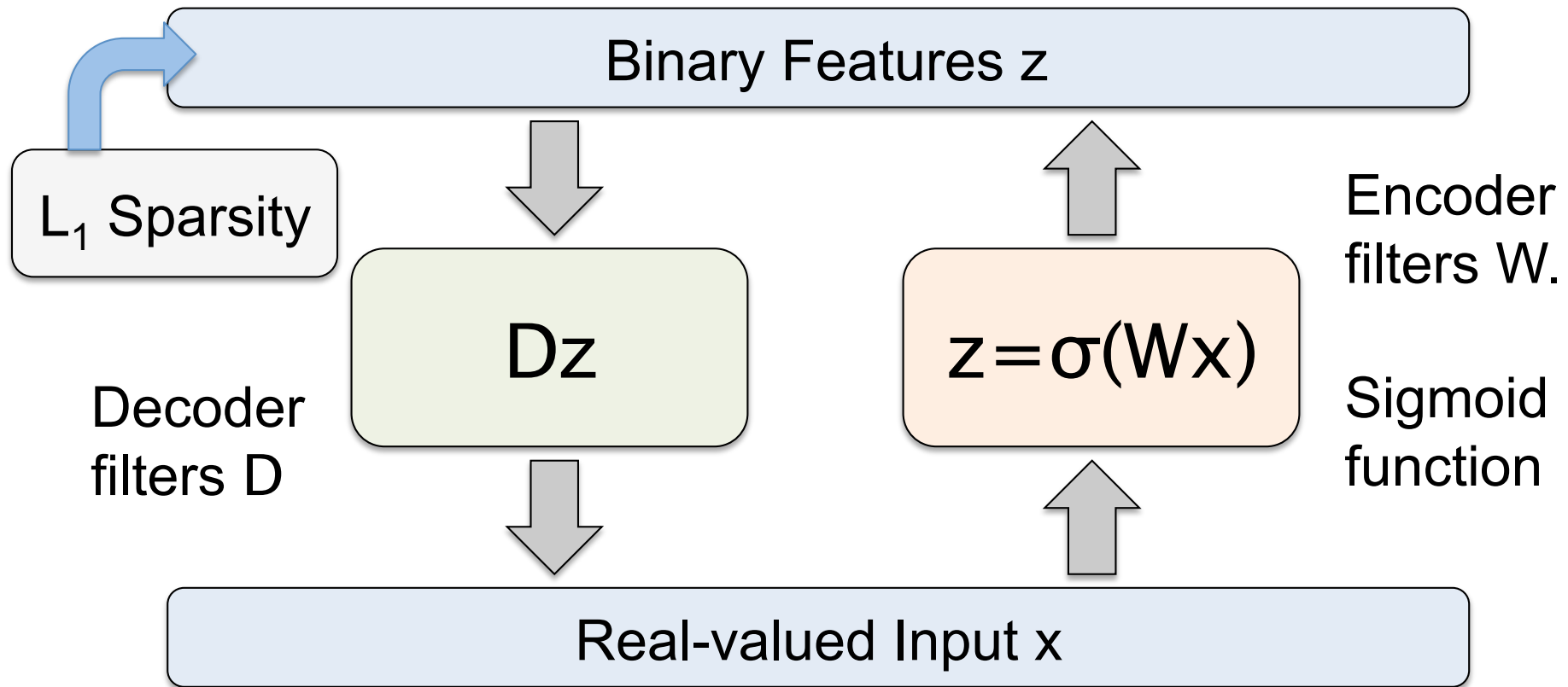
50% corrupted input



(Vincent et al., ICML 2008)



# Predictive Sparse Decomposition



At training time

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

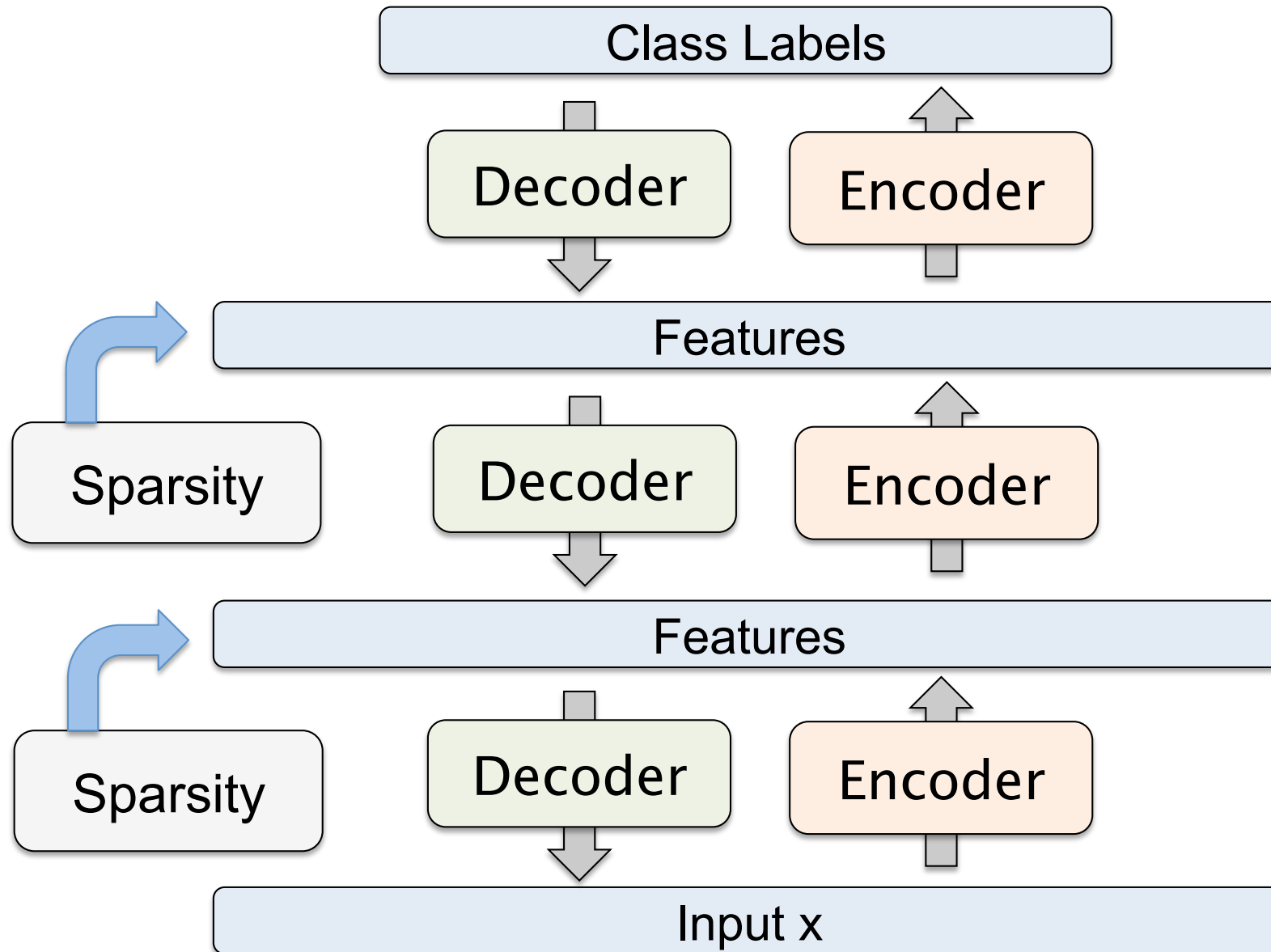
Decoder

Encoder

(Kavukcuoglu, Ranzato, Fergus, LeCun, 2009)

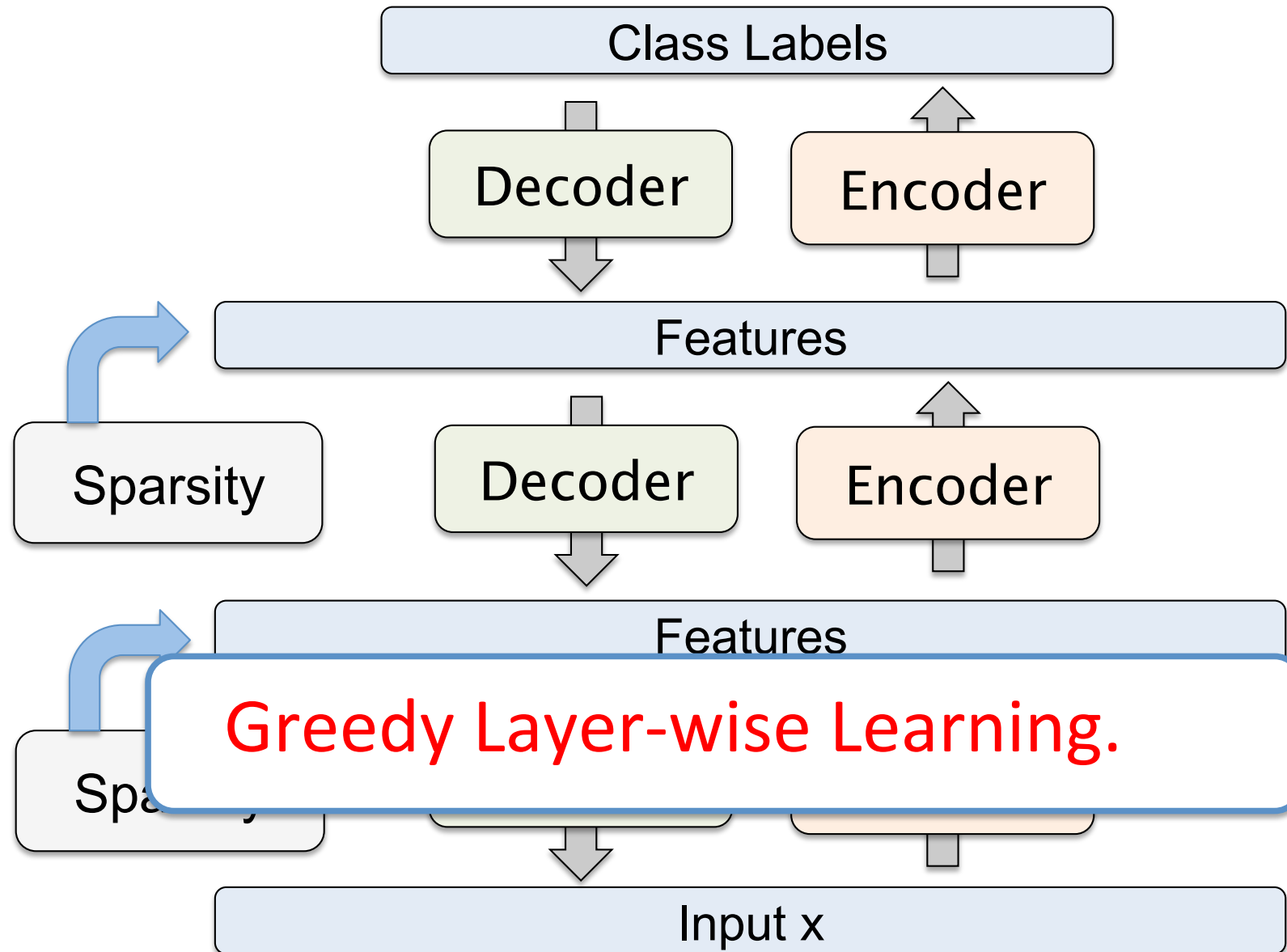
# Stacked Autoencoders

---



# Stacked Autoencoders

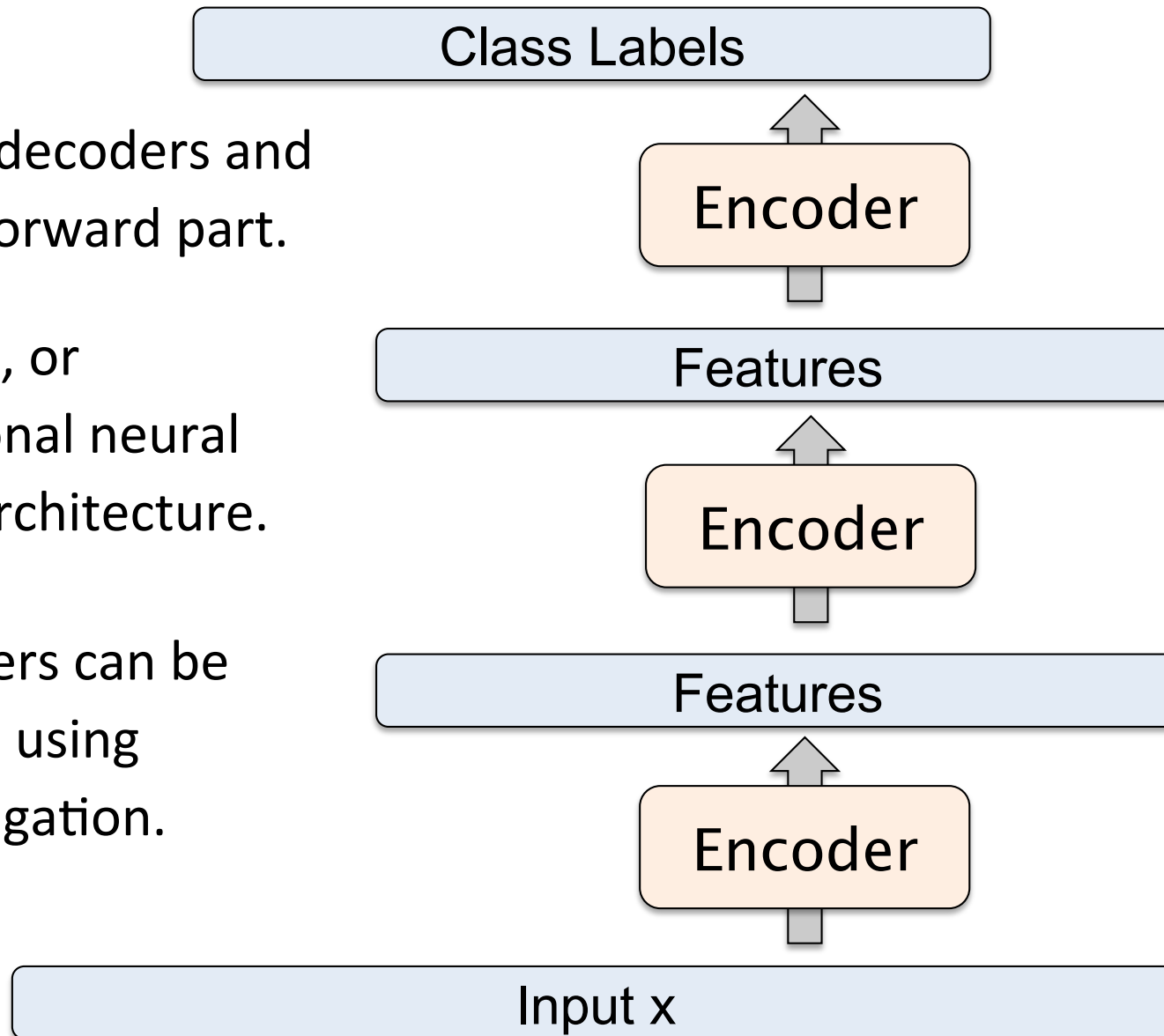
---



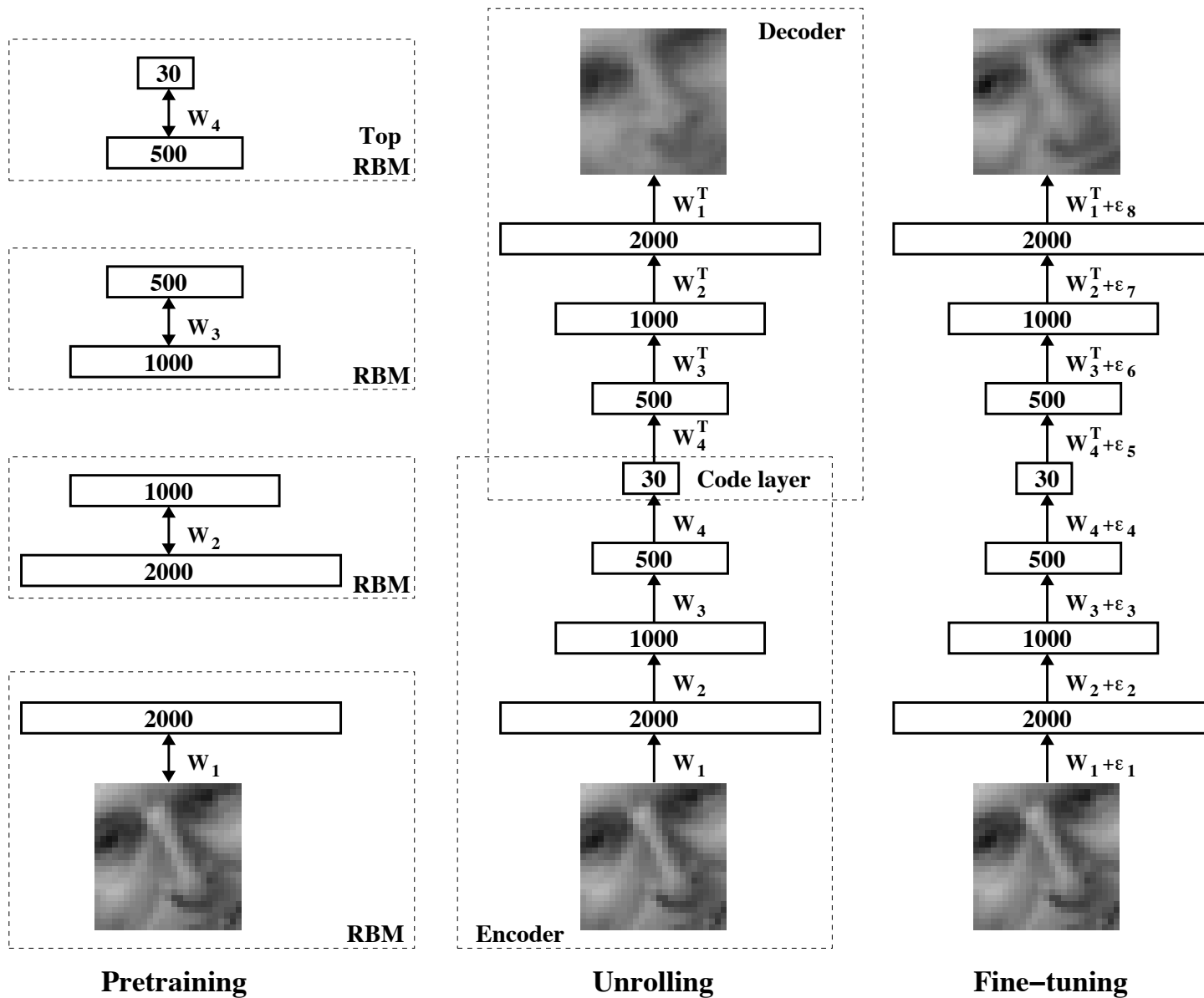
# Stacked Autoencoders

---

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



# Deep Autoencoders



# Deep Autoencoders

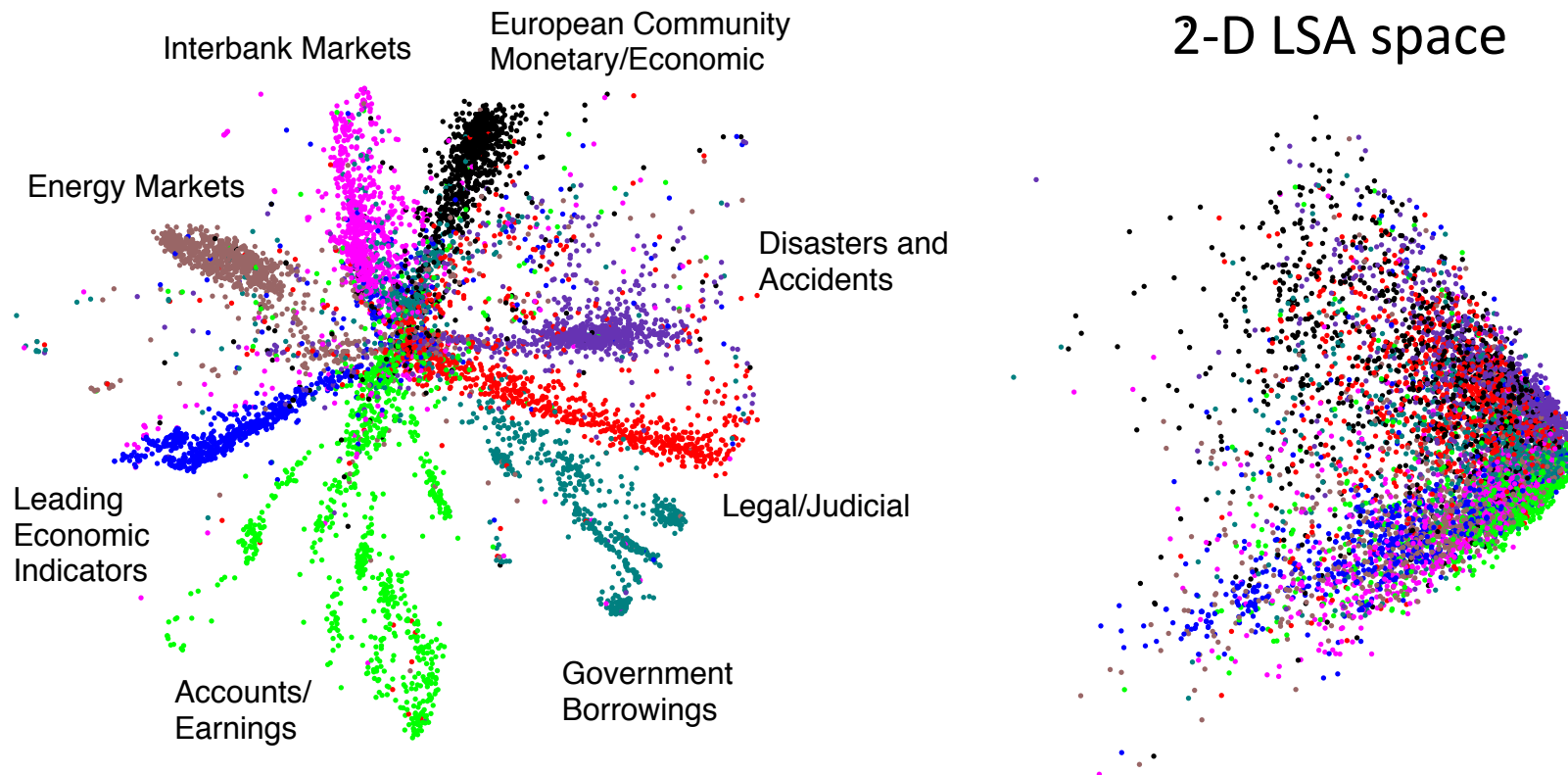
- 25x25 – 2000 – 1000 – 500 – 30 autoencoder to extract 30-D real-valued codes for Olivetti face patches.



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

(Hinton and Salakhutdinov, Science 2006)

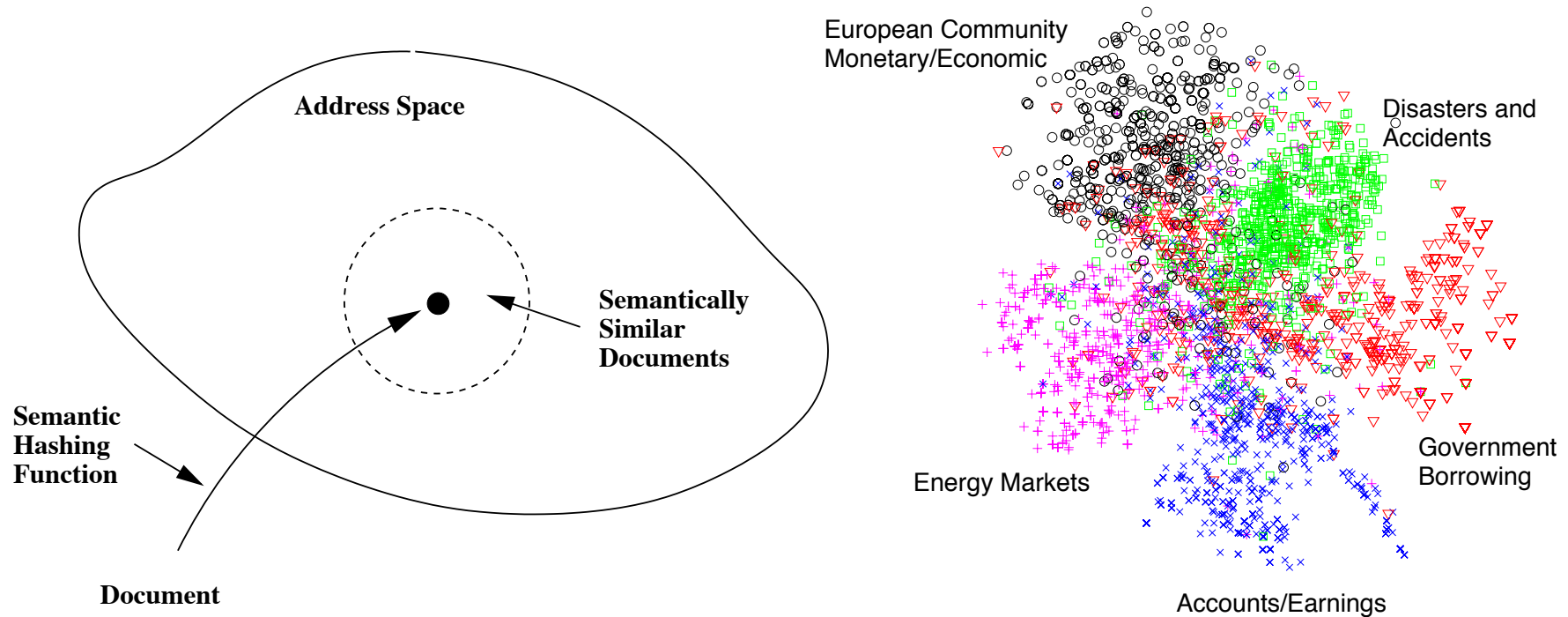
# Information Retrieval



- 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,

# Unsupervised Learning

## Non-probabilistic Models

- Sparse Coding
- Autoencoders
- Others (e.g. k-means)

## Probabilistic (Generative) Models

### Tractable Models

- Fully observed Belief Nets
- NADE
- PixelRNN

### Non-Tractable Models

- Boltzmann Machines
- Variational Autoencoders
- Helmholtz Machines
- Many others...

- Generative Adversarial Networks
- Moment Matching Networks

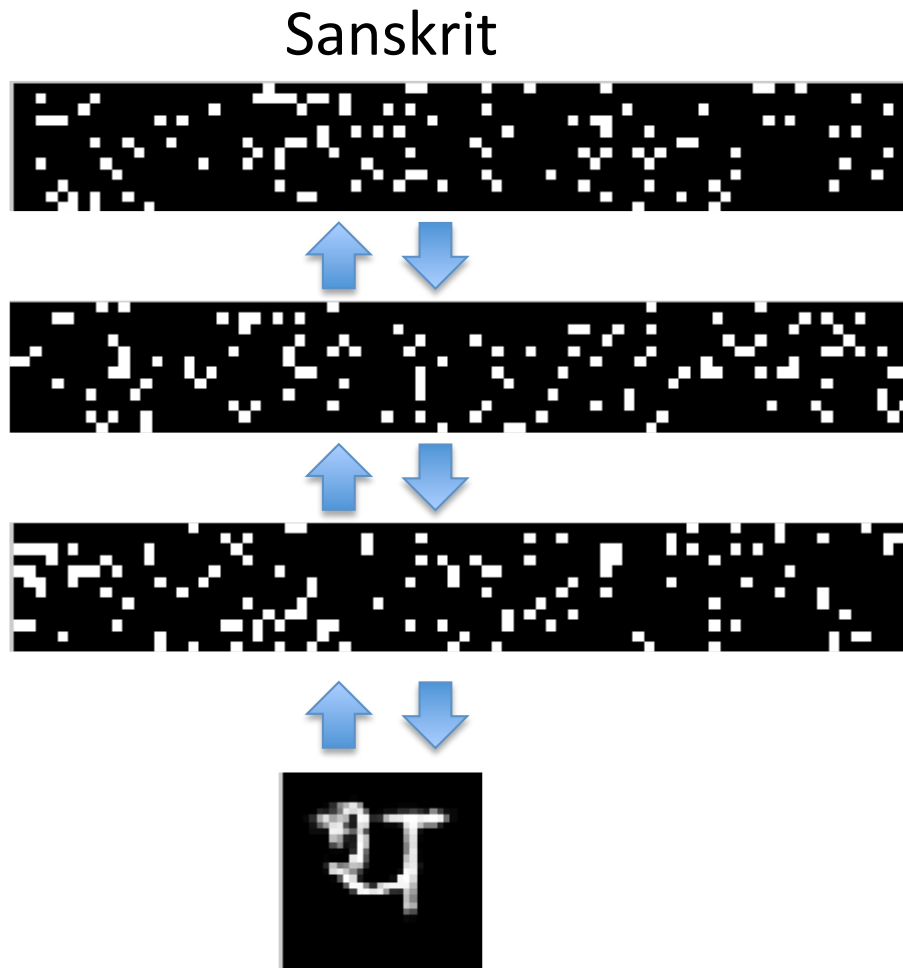
Explicit Density  $p(x)$

Implicit Density

# Talk Roadmap

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

# Deep Generative Model



Model  $P(\text{image})$

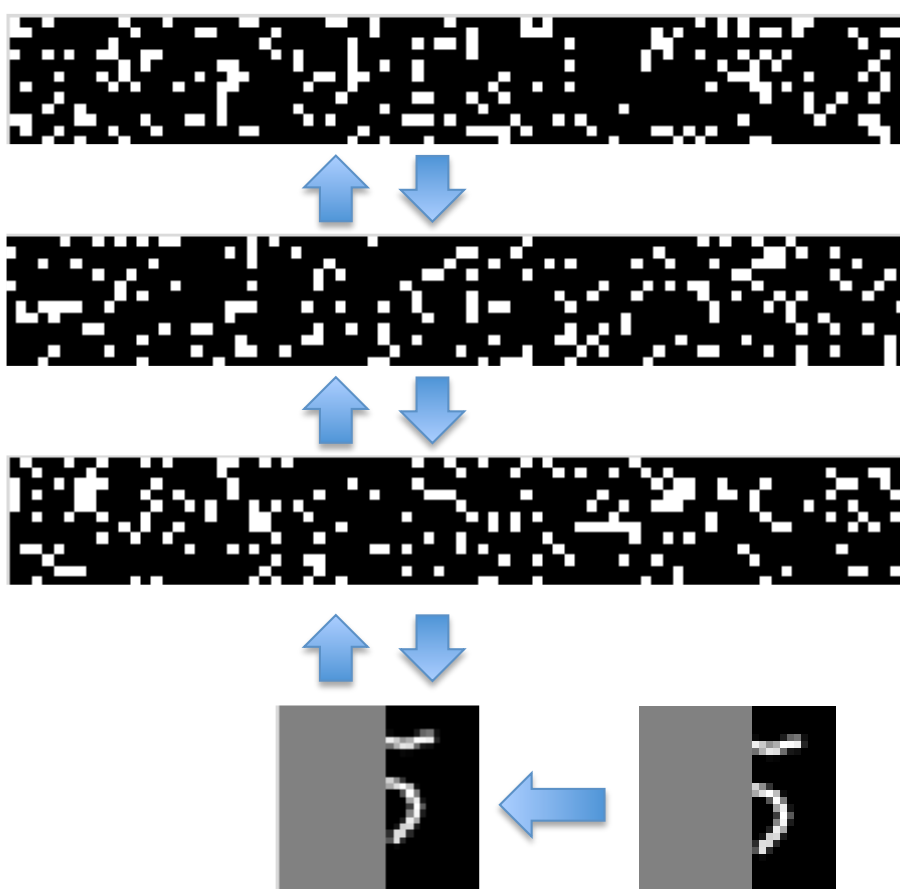
ल च थ श म छ ण ञ  
ट ढ ब आ ल ओ ट र  
ऋ इ ल ब ष अ उ आ  
ए ष श य ऋ ष इ त्र

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

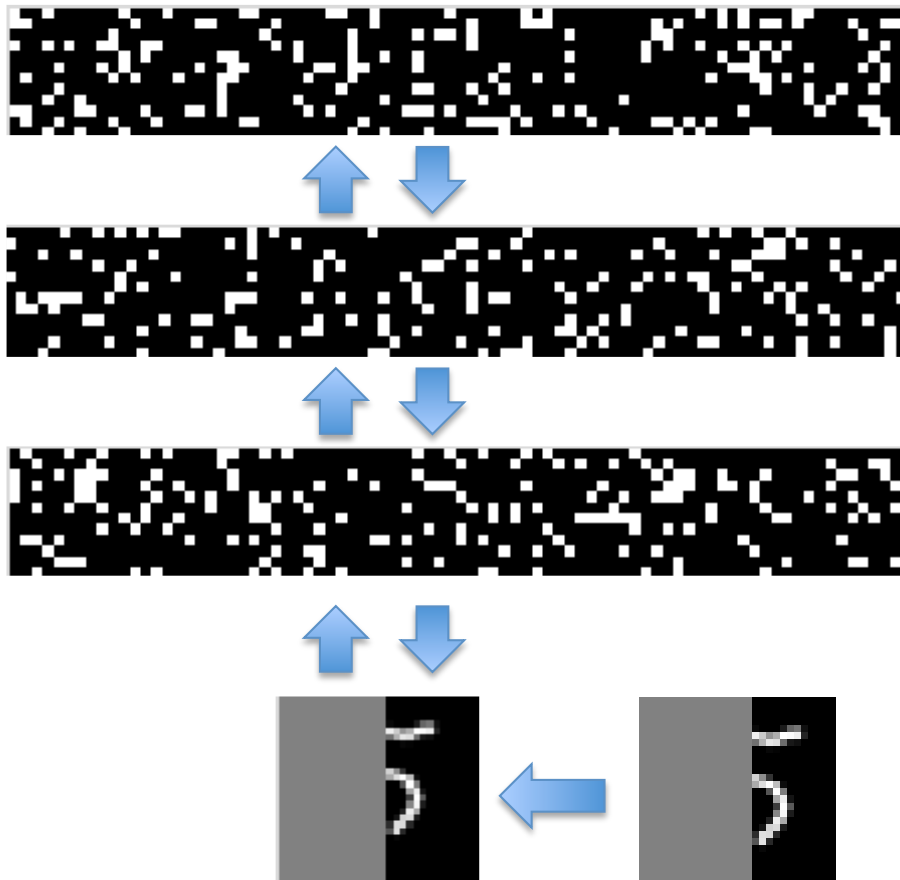


Conditional  
Simulation

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

Bernoulli Markov Random Field

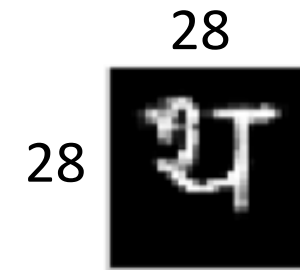
# Deep Generative Model



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

Conditional Simulation

Why so difficult?



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

Bernoulli Markov Random Field

# Fully Observed Models

- Explicitly model conditional probabilities:

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

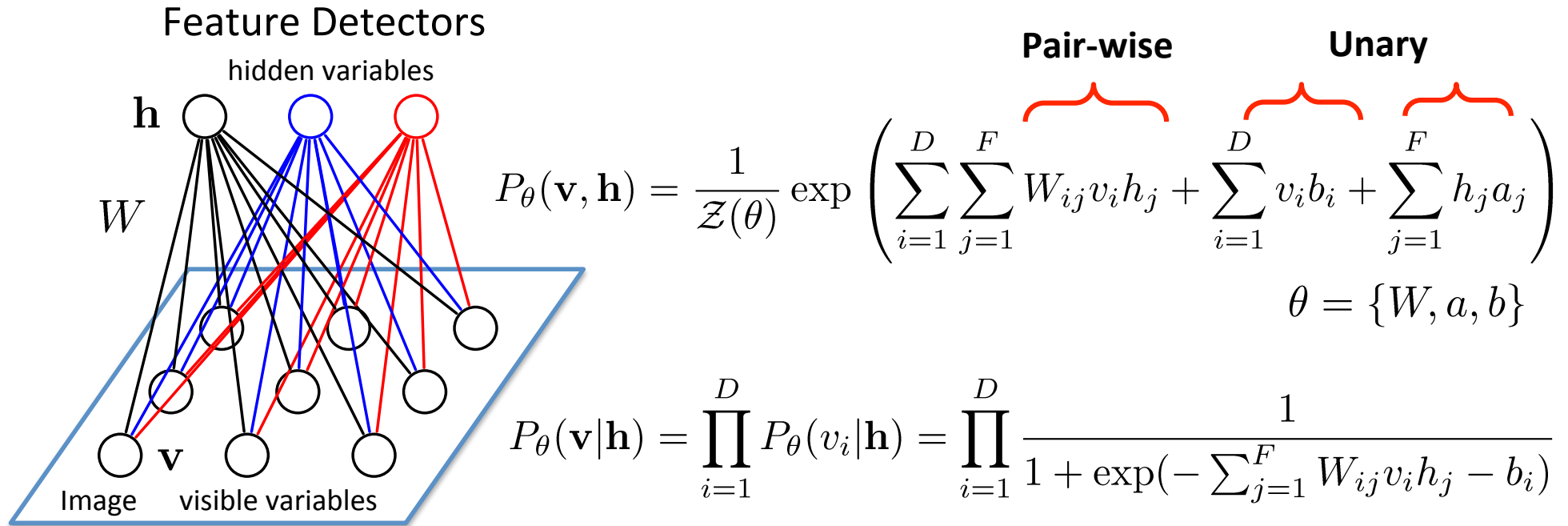
Each conditional can be a complicated neural network

- A number of successful models, including
  - NADE, RNADE (Larochelle, et.al. 20011)
  - Pixel CNN (van den Ord et. al. 2016)
  - Pixel RNN (van den Ord et. al. 2016)



Pixel CNN

# Restricted Boltzmann Machines



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.

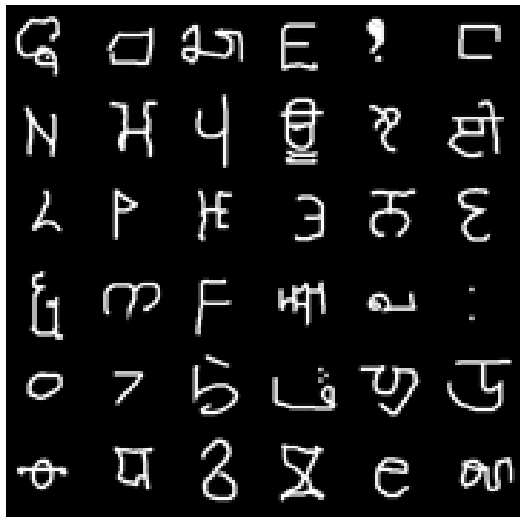
Markov random fields, Boltzmann machines, log-linear models.



# Learning Features

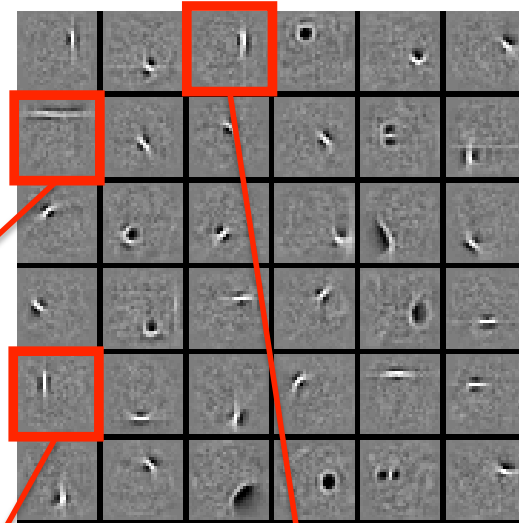
Observed Data

Subset of 25,000 characters

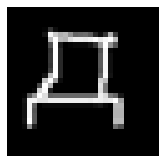


Learned W: "edges"

Subset of 1000 features



New Image:



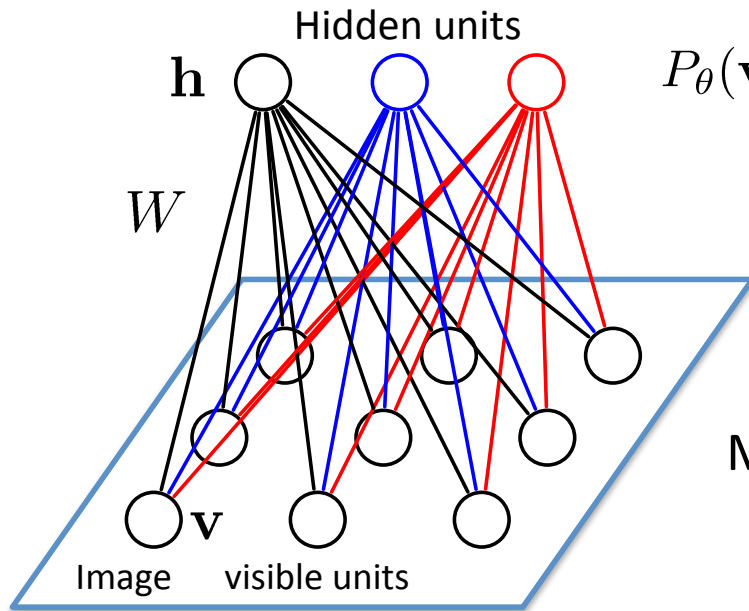
$$p(h_7 = 1|v) \quad p(h_{29} = 1|v)$$

$$= \sigma \left( 0.99 \times \text{[edge detector 1]} + 0.97 \times \text{[edge detector 2]} + 0.82 \times \text{[edge detector 3]} + \dots \right)$$

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

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:

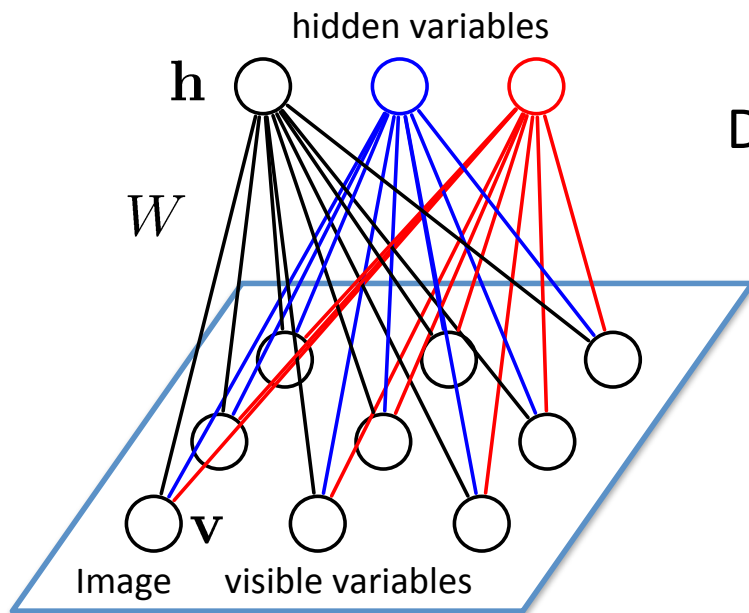
$$\begin{aligned} \frac{\partial L(\theta)}{\partial W_{ij}} &= \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial W_{ij}} \log \left( \sum_{\mathbf{h}} \exp \left[ \mathbf{v}^{(n)\top} W \mathbf{h} + \mathbf{a}^{\top} \mathbf{h} + \mathbf{b}^{\top} \mathbf{v}^{(n)} \right] \right) - \frac{\partial}{\partial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

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

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

Difficult to compute: exponentially many configurations

# Model Learning



Derivative of the log-likelihood:

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

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

Easy to compute exactly

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:

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

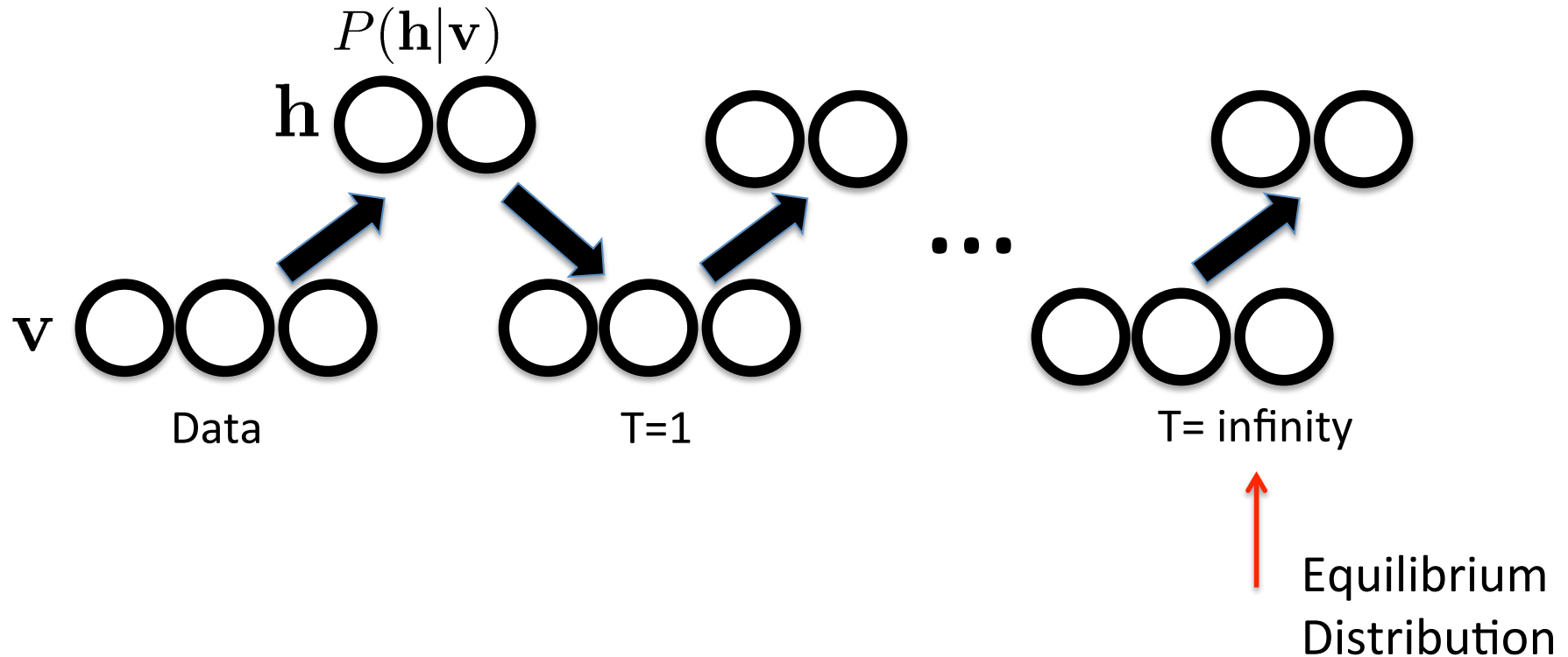
$$\sum_{\mathbf{v}, \mathbf{h}} v_i h_j P_\theta(\mathbf{v}, \mathbf{h})$$

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

# Approximate ML Learning for RBMs

Run Markov chain (alternating Gibbs Sampling):

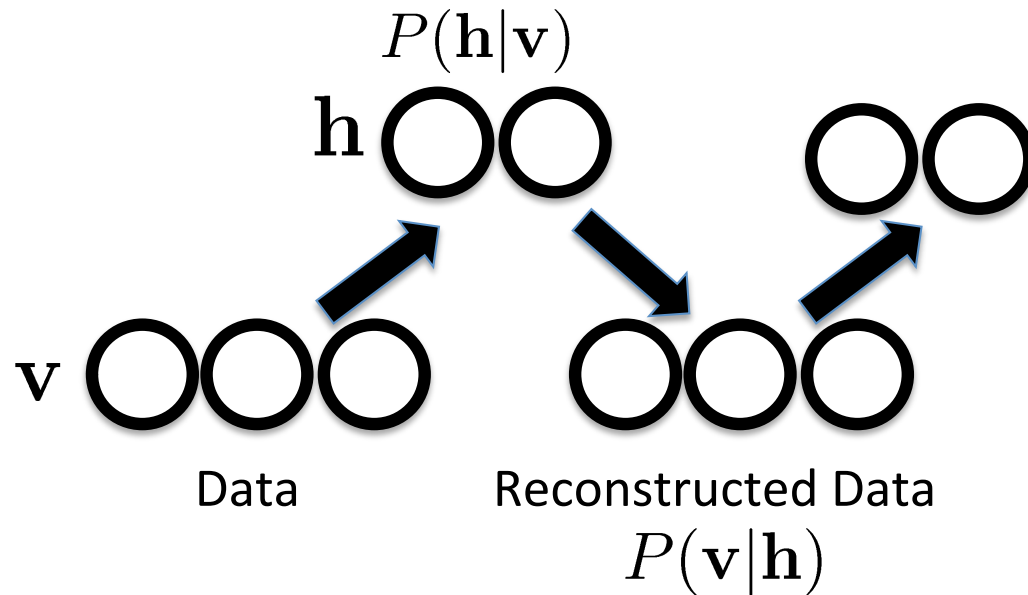


$$P(\mathbf{h}|\mathbf{v}) = \prod_j P(h_j|\mathbf{v}) \quad P(h_j = 1|\mathbf{v}) = \frac{1}{1 + \exp(-\sum_i W_{ij}v_i - a_j)}$$

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

# 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 the all the visible units in parallel to get a “reconstruction”.
- Update the hidden units again.

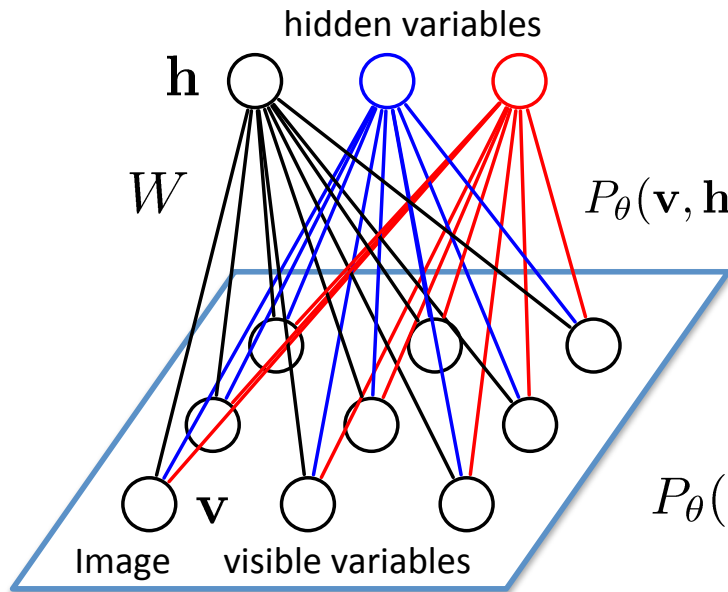
Update model parameters:

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

Implementation: ~10 lines of Matlab code.

(Hinton, Neural Computation 2002)

# RBM for Real-valued Data



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

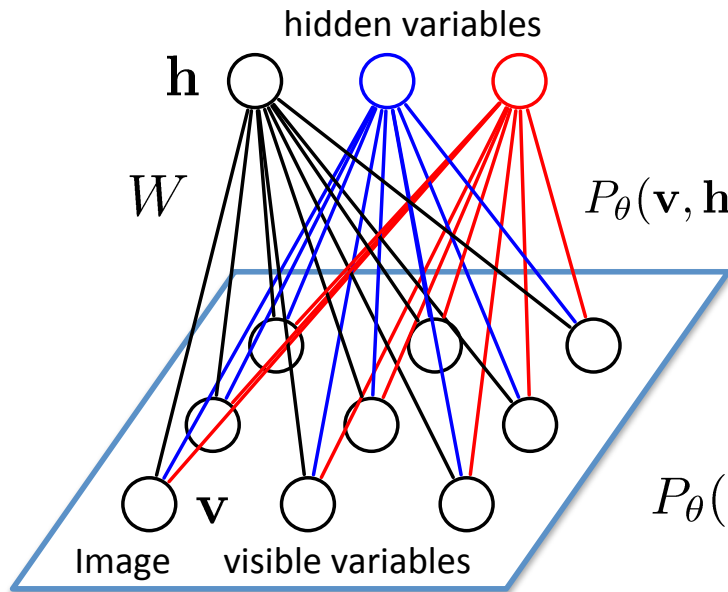
$$\theta = \{W, a, b\}$$

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

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 for Real-valued Data



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

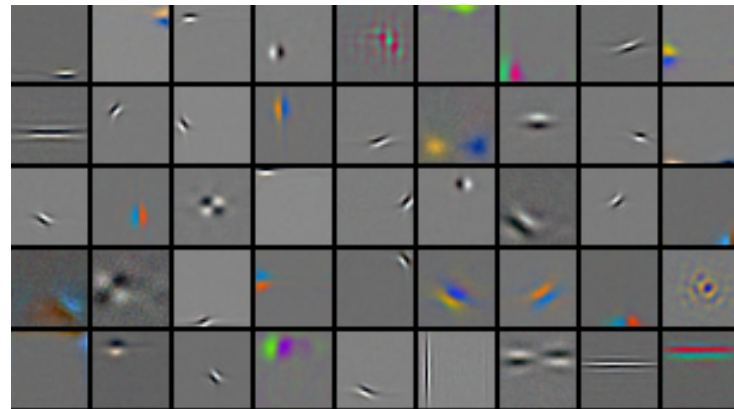
$$\theta = \{W, a, b\}$$

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

4 million **unlabelled** images



Learned features (out of 10,000)



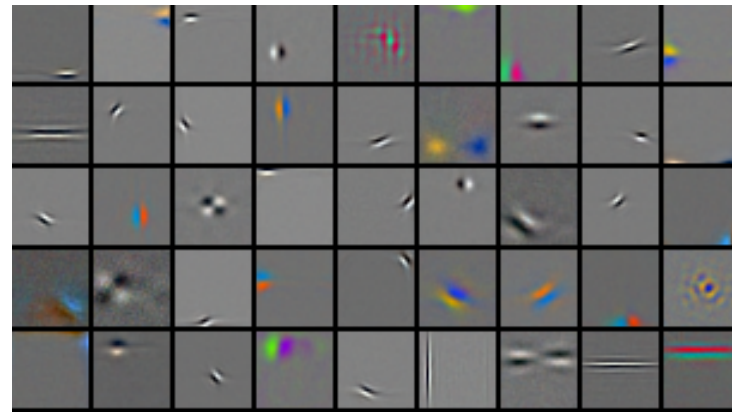



# RBM for Real-valued Data

4 million **unlabelled** images



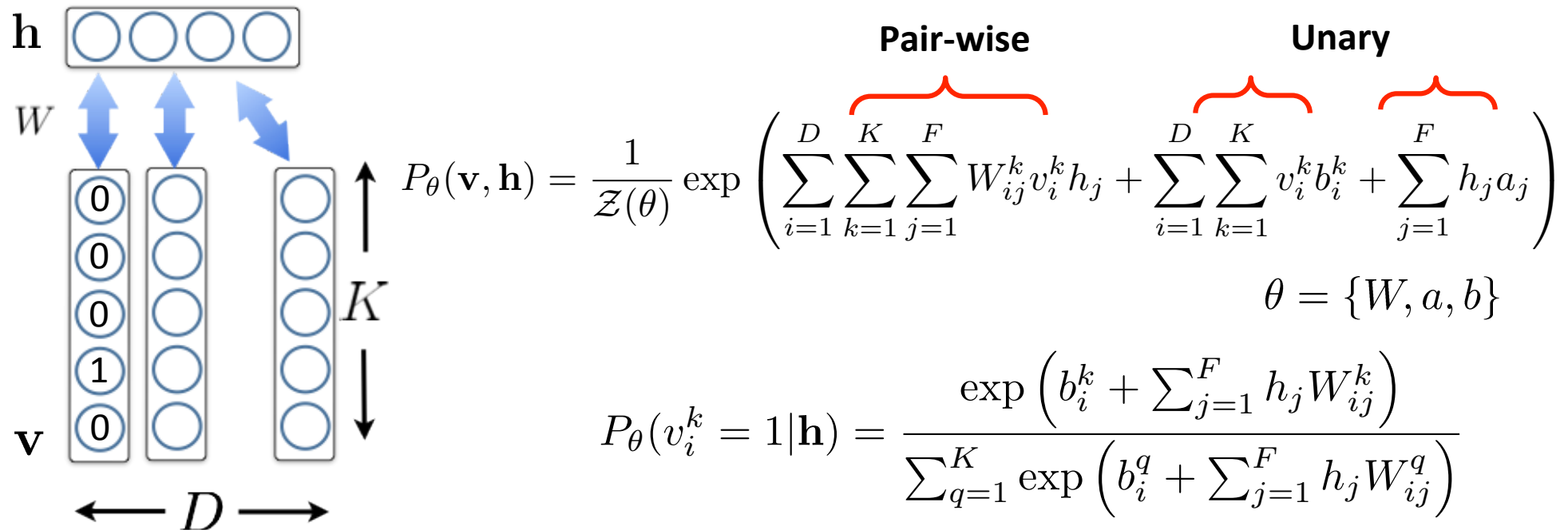
Learned features (out of 10,000)



  
New Image

$$= 0.9 * \begin{matrix} p(h_7 = 1|v) \\ \downarrow \\ \text{feature 7} \end{matrix} + 0.8 * \begin{matrix} p(h_{29} = 1|v) \\ \downarrow \\ \text{feature 29} \end{matrix} + 0.6 * \begin{matrix} \text{feature 3} \\ \text{feature 1} \end{matrix} \dots$$

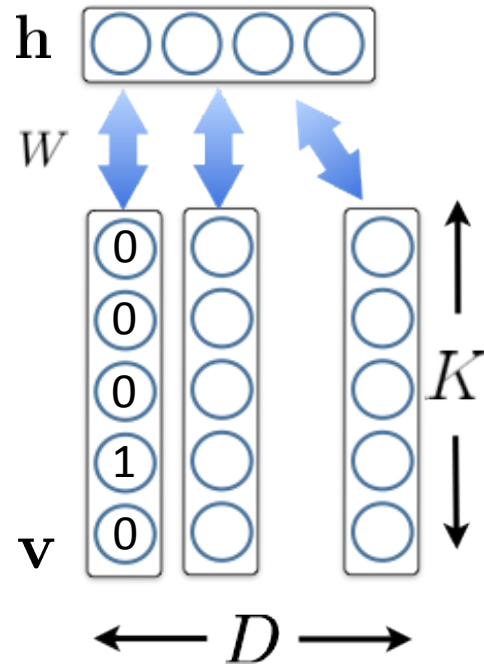
# RBM for Word Counts



Replicated Softmax Model: undirected topic model:

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

# RBMMs for Word Counts



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

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

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



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

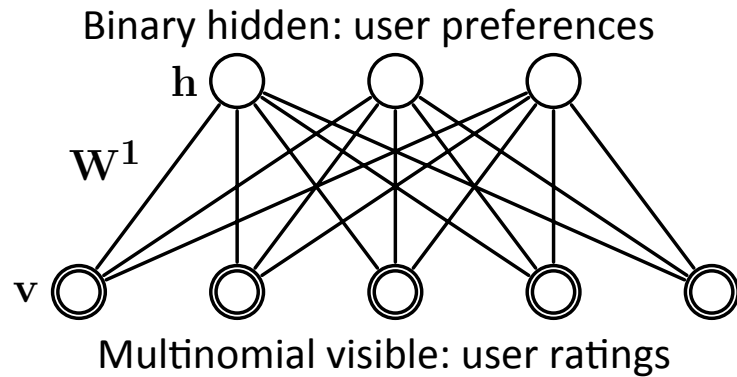
# RBM for Word Counts

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

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



Learned features: ``genre``

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

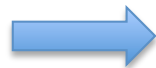
Independence Day  
The Day After Tomorrow  
Con Air  
Men in Black II  
Men in Black

Netflix dataset:

480,189 users

17,770 movies

Over 100 million ratings



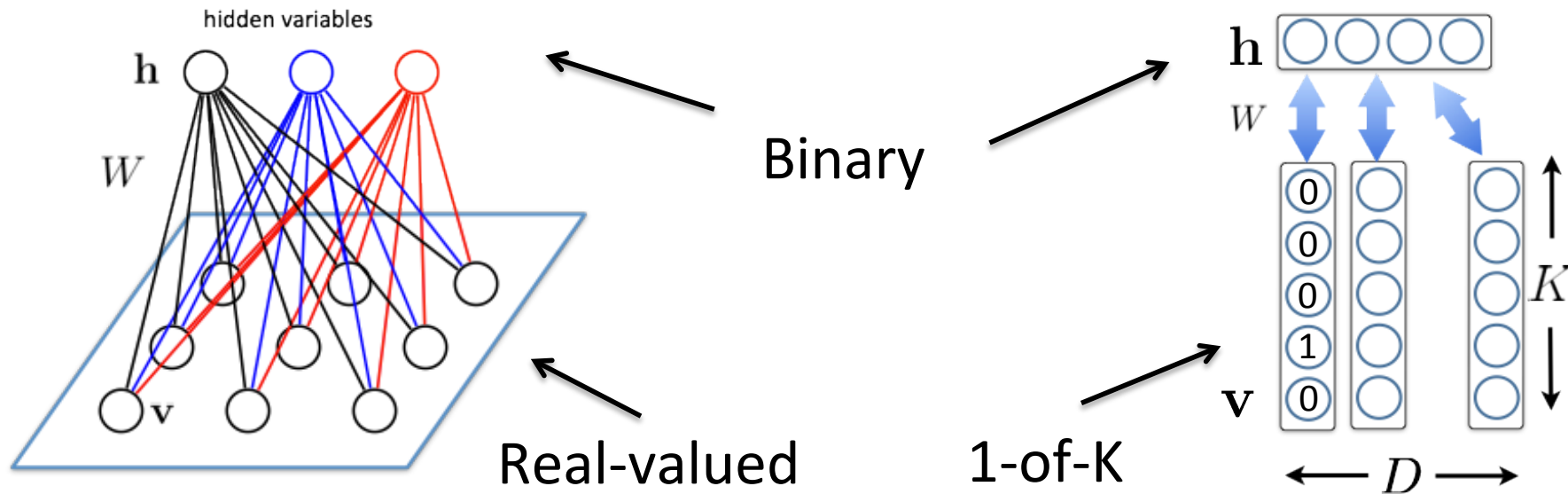
Friday the 13th  
The Texas Chainsaw Massacre  
Children of the Corn  
Child's Play  
The Return of Michael Myers

Scary Movie  
Naked Gun  
Hot Shots!  
American Pie  
Police Academy



# 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:

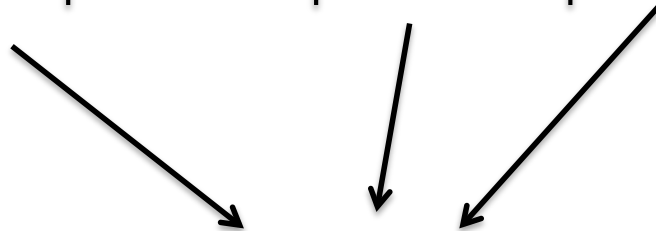
$$P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \exp \left( \sum_{ij} W_{ij} v_i h_j + \sum_i b_i v_i + \sum_j a_j h_j \right)$$

Marginalizing over hidden variables:

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

**Product of Experts**

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



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:

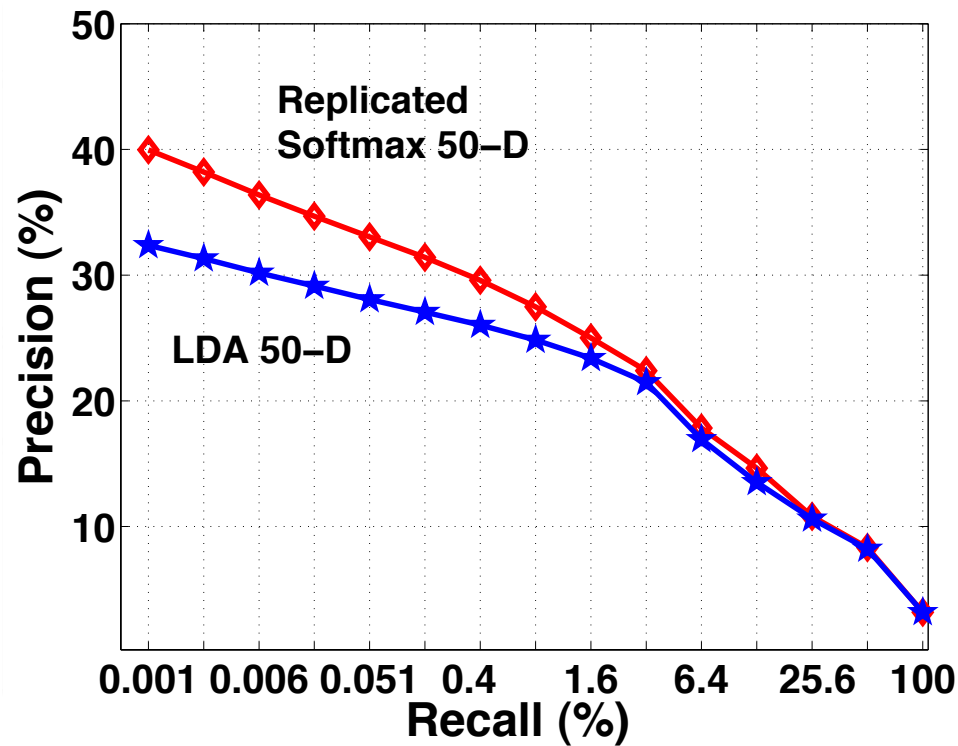
$$P_{\theta}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} \exp \left( \sum_{ij} W_{ij} v_i h_j + \sum_i b_i v_i + \sum_j a_j h_j \right)$$

Marginalizing over  $\mathbf{h}$

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

government  
authority  
power  
empire  
federation

clint  
hou  
pres  
bill  
cong



Product of Experts

$$\left( \prod_i W_{ij} v_i \right)$$

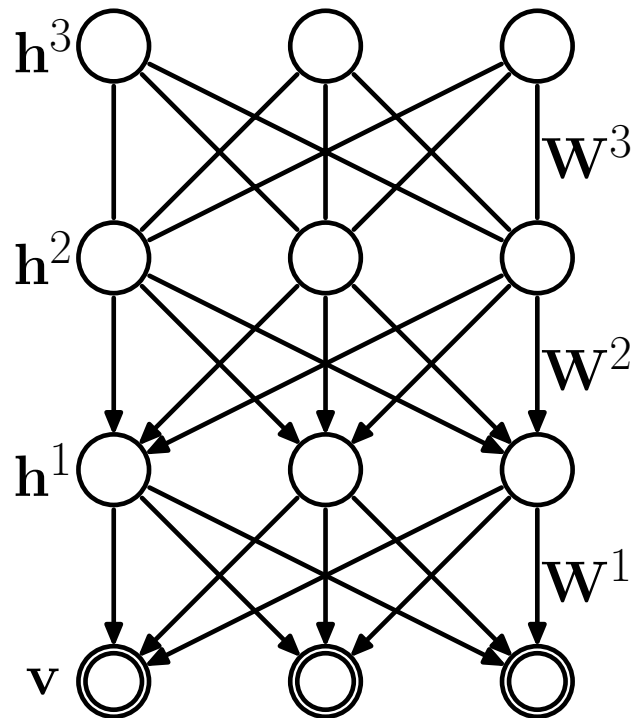
,"corruption"  
bine to give very  
word "Silvio



# Talk Roadmap

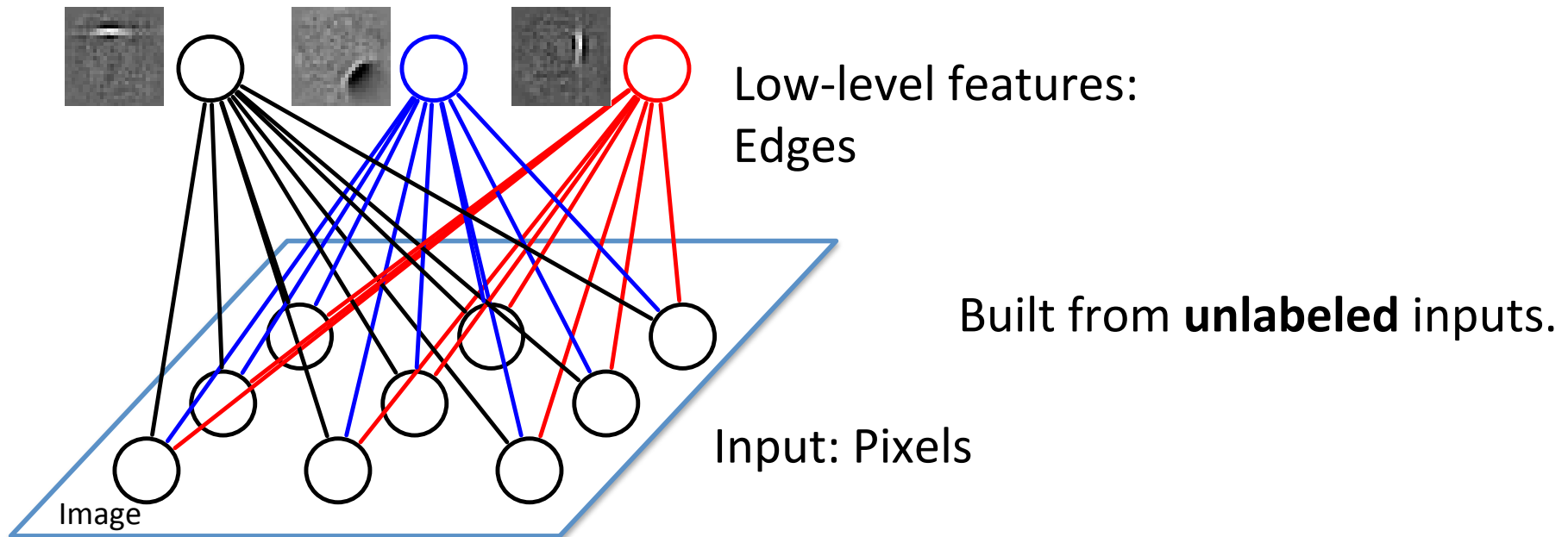
- Basic Building Blocks (non-probabilistic models):
  - Sparse Coding
  - Autoencoders
- Deep Generative Models
  - Restricted Boltzmann Machines
  - Deep Boltzmann Machines
  - Helmholtz Machines / Variational Autoencoders
- Generative Adversarial Networks

# Deep Belief Network



- Probabilistic Generative model.
  - Contains multiple layers of nonlinear representation.
  - Fast, greedy layer-wise pretraining algorithm.
  - Inferring the states of the latent variables in highest layers is easy.
- 
- Inferring the states of the latent variables in highest layers is easy.

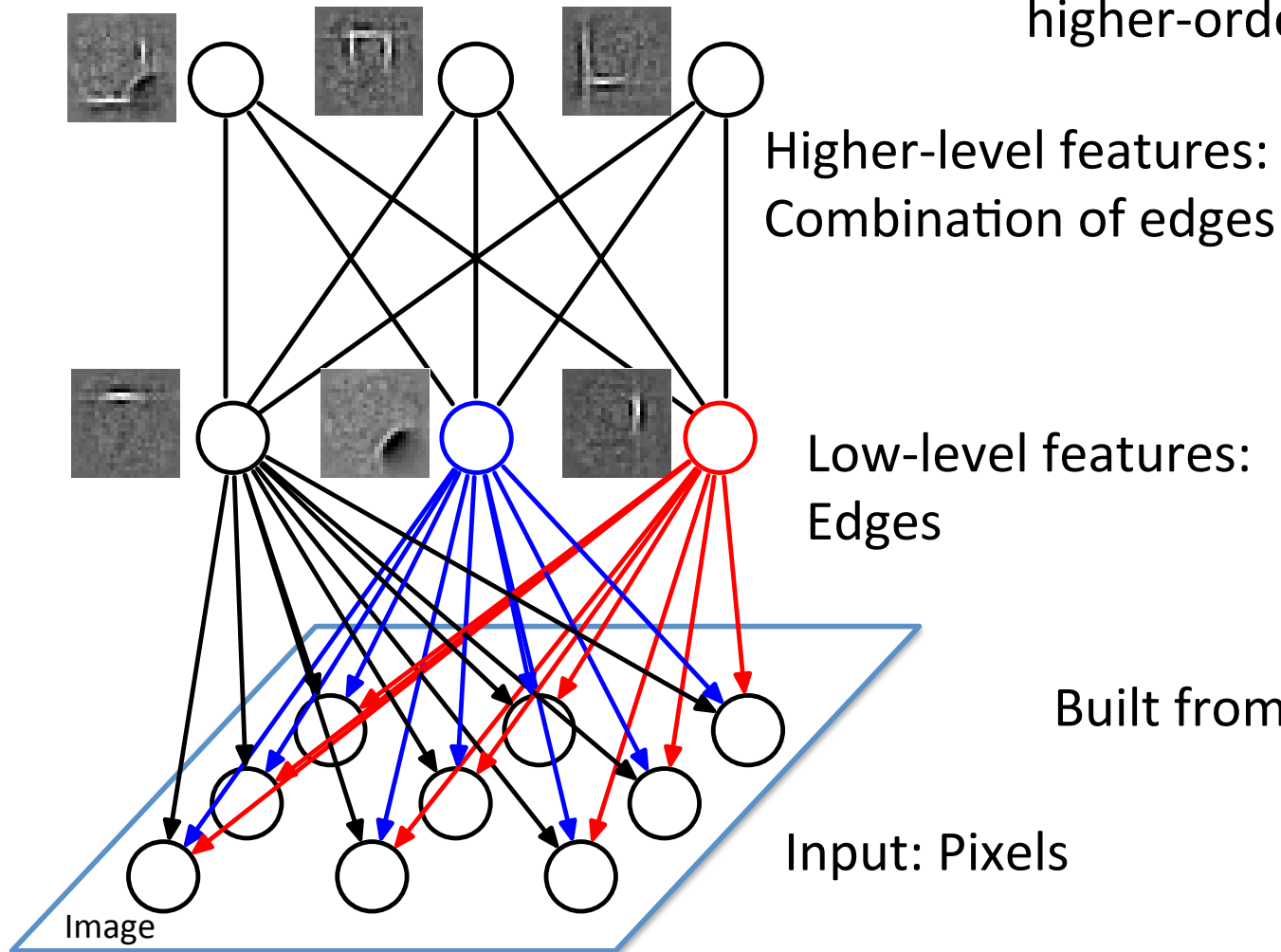
# Deep Belief Network



(Hinton et al. Neural Computation 2006)

# Deep Belief Network

Internal representations capture higher-order statistical structure



Higher-level features:  
Combination of edges

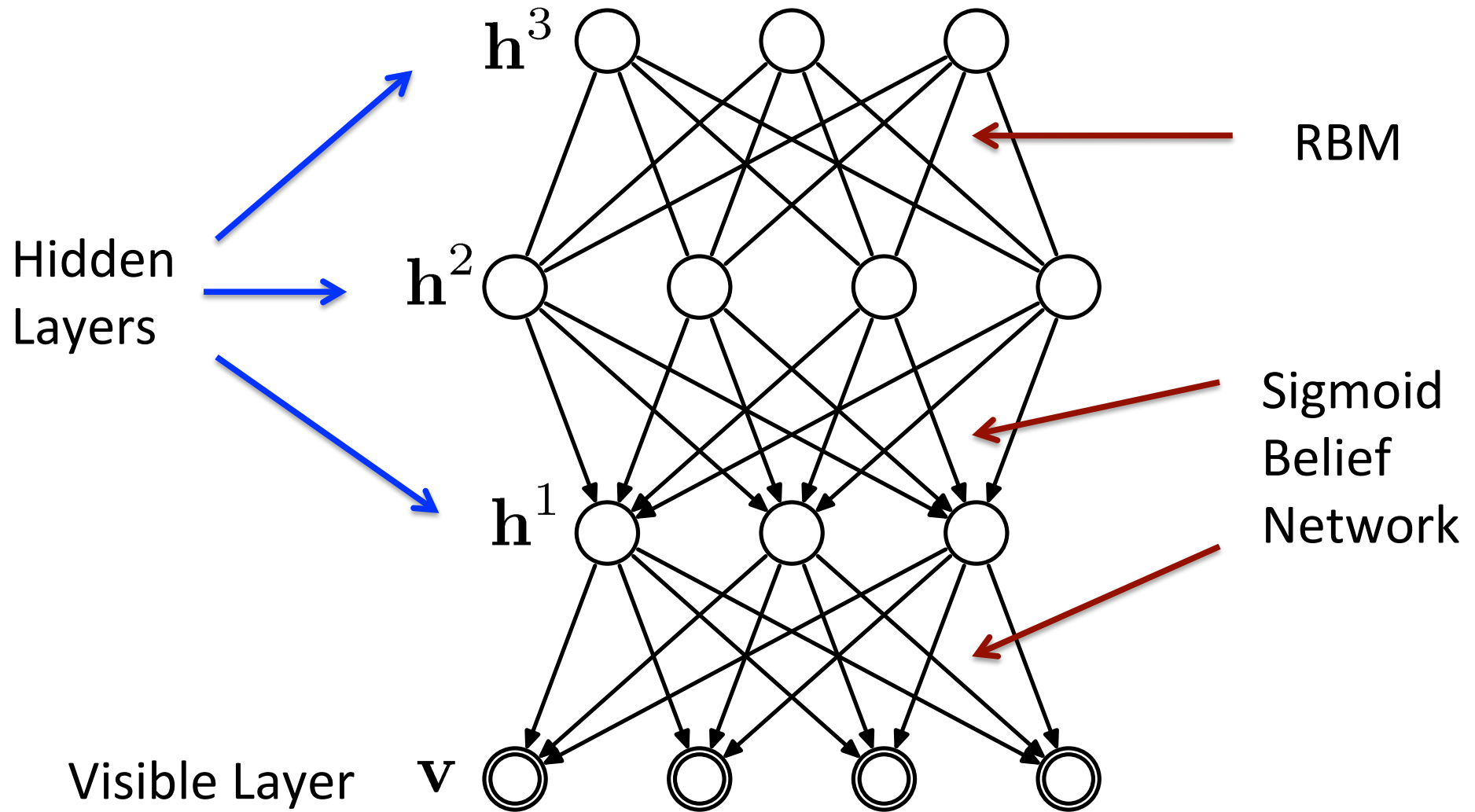
Low-level features:  
Edges

Built from **unlabeled** inputs.

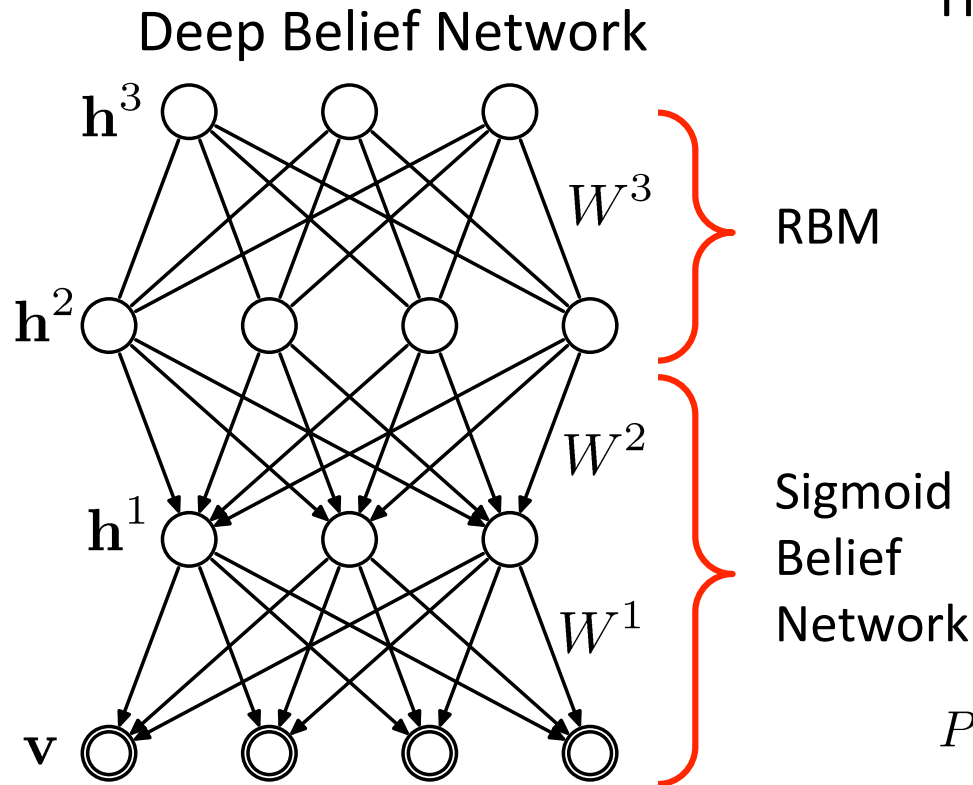
Input: Pixels

(Hinton et al. Neural Computation 2006)

# Deep Belief Network



# Deep Belief Network



The joint probability distribution factorizes:

$$P(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3) = P(\mathbf{v}|\mathbf{h}^1)P(\mathbf{h}^1|\mathbf{h}^2)P(\mathbf{h}^2, \mathbf{h}^3)$$

Sigmoid Belief Network

RBM

$$P(\mathbf{h}^2, \mathbf{h}^3) = \frac{1}{Z(W^3)} \exp[\mathbf{h}^{2\top} W^3 \mathbf{h}^3]$$

$$P(\mathbf{h}^1|\mathbf{h}^2) = \prod_j P(h_j^1|\mathbf{h}^2)$$

$$P(h_j^1 = 1|\mathbf{h}^2) = \frac{1}{1 + \exp\left(-\sum_k W_{jk}^2 h_k^2\right)}$$

$$P(\mathbf{v}|\mathbf{h}^1) = \prod_i P(v_i|\mathbf{h}^1)$$

$$P(v_i = 1|\mathbf{h}^1) = \frac{1}{1 + \exp\left(-\sum_j W_{ij}^1 h_j^1\right)}$$

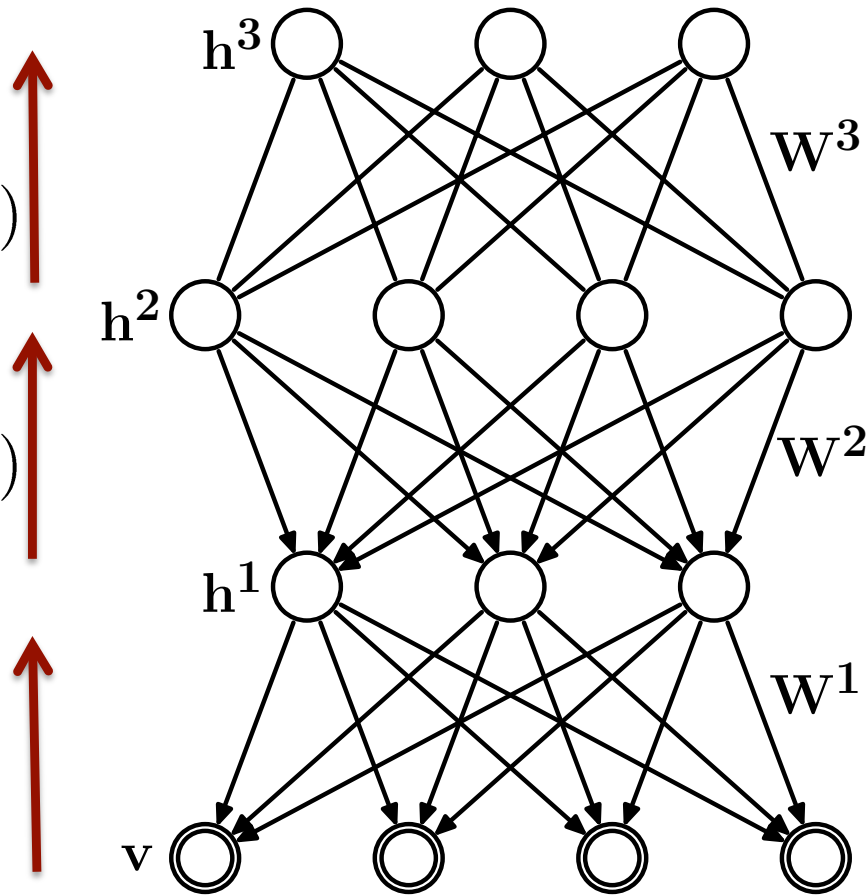
# Deep Belief Network

Approximate  
Inference

$$Q(\mathbf{h}^3 | \mathbf{h}^2)$$

$$Q(\mathbf{h}^2 | \mathbf{h}^1)$$

$$Q(\mathbf{h}^1 | \mathbf{v})$$



Generative  
Process

$$P(\mathbf{h}^2, \mathbf{h}^3)$$

$$P(\mathbf{h}^1 | \mathbf{h}^2)$$

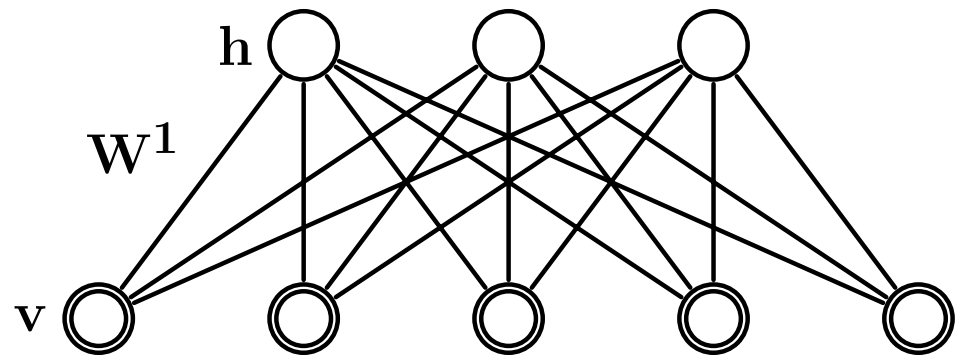
$$P(\mathbf{v} | \mathbf{h}^1)$$

$$Q(\mathbf{h}^t | \mathbf{h}^{t-1}) = \prod_j \sigma \left( \sum_i W^t h_i^{t-1} \right)$$

$$P(\mathbf{h}^{t-1} | \mathbf{h}^t) = \prod_j \sigma \left( \sum_i W^t h_i^t \right)$$

# DBN Layer-wise Training

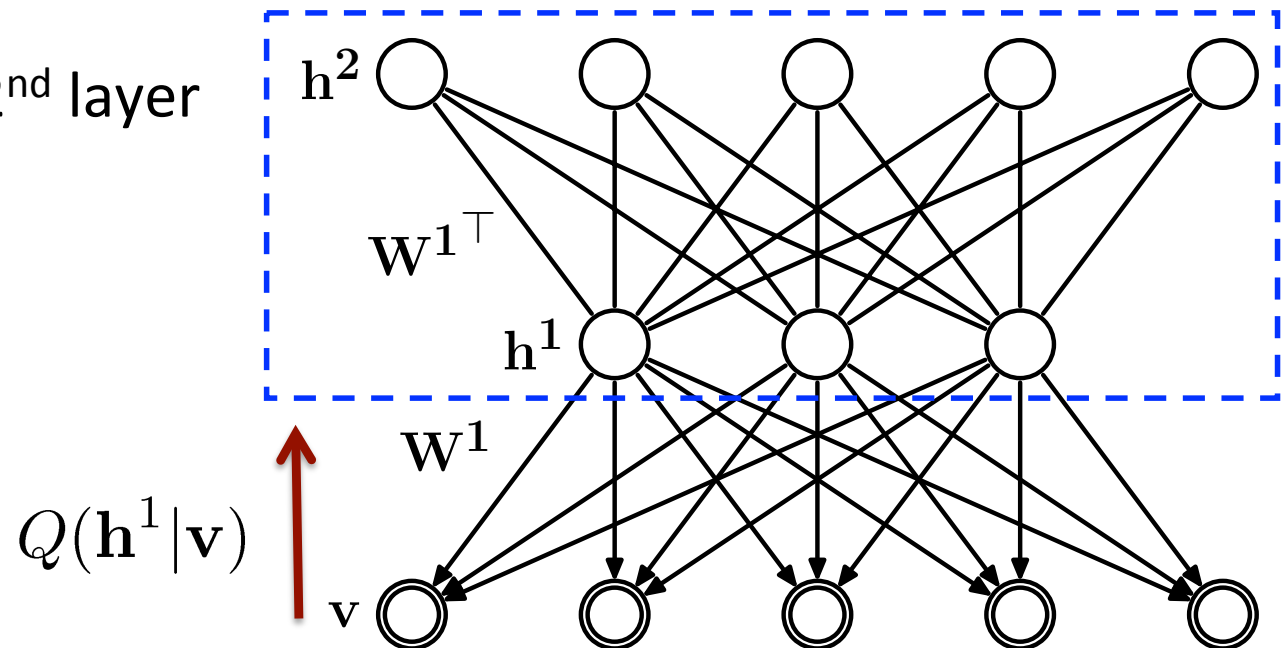
- Learn an RBM with an input layer  $v$  and a hidden layer  $h$ .





# DBN Layer-wise Training

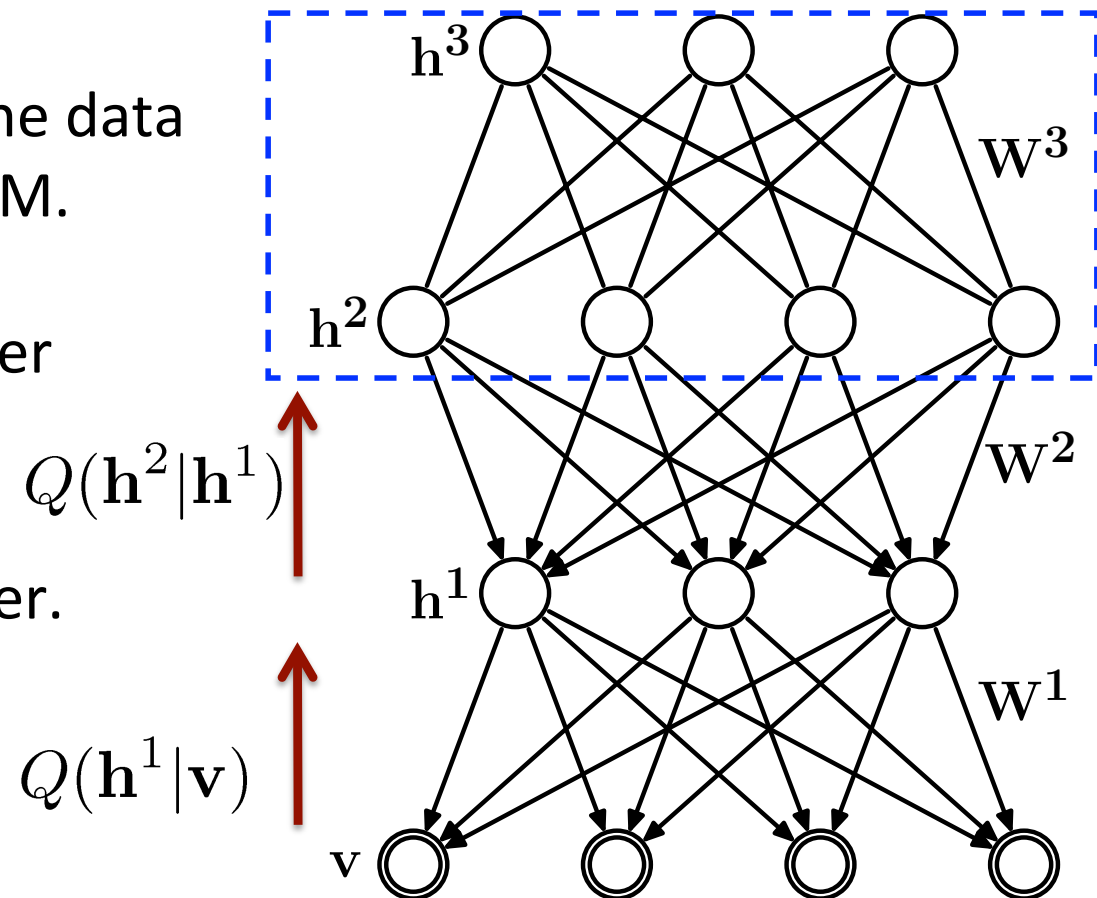
- Learn an RBM with an input layer  $v$  and a hidden layer  $h$ .
- Treat inferred values  $Q(\mathbf{h}^1 | \mathbf{v}) = P(\mathbf{h}^1 | \mathbf{v})$  as the data for training 2<sup>nd</sup>-layer RBM.
- Learn and freeze 2<sup>nd</sup> layer RBM.



# DBN Layer-wise Training

- Learn an RBM with an input layer  $v$  and a hidden layer  $h$ .
- Treat inferred values  $Q(\mathbf{h}^1 | \mathbf{v}) = P(\mathbf{h}^1 | \mathbf{v})$  as the data for training 2<sup>nd</sup>-layer RBM.
- Learn and freeze 2<sup>nd</sup> layer RBM.
- Proceed to the next layer.

Unsupervised Feature Learning.

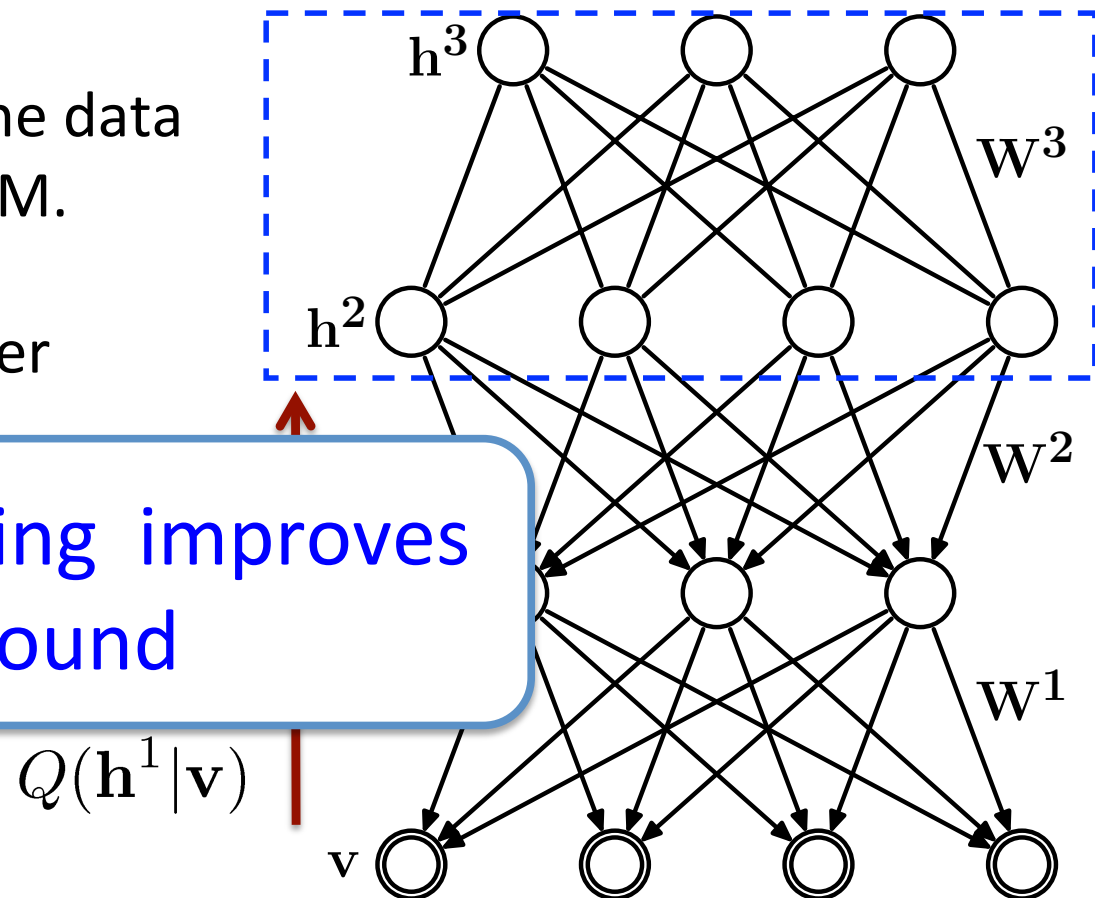


# DBN Layer-wise Training

- Learn an RBM with an input layer  $v$  and a hidden layer  $h$ .
- Treat inferred values  $Q(\mathbf{h}^1 | \mathbf{v}) = P(\mathbf{h}^1 | \mathbf{v})$  as the data for training 2<sup>nd</sup>-layer RBM.
- Learn and freeze 2<sup>nd</sup> layer RBM

Layerwise pretraining improves variational lower bound

Unsupervised Feature Learning.



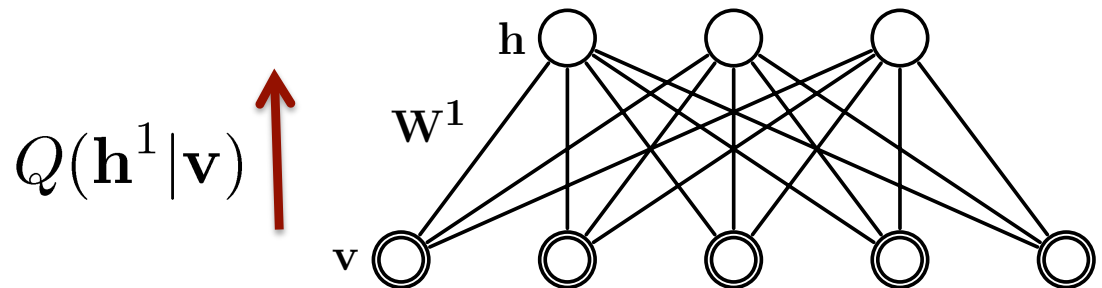
# Why this Pre-training Works?

- Greedy training improves variational lower bound!

- For any approximating distribution  $Q(\mathbf{h}^1 | \mathbf{v})$

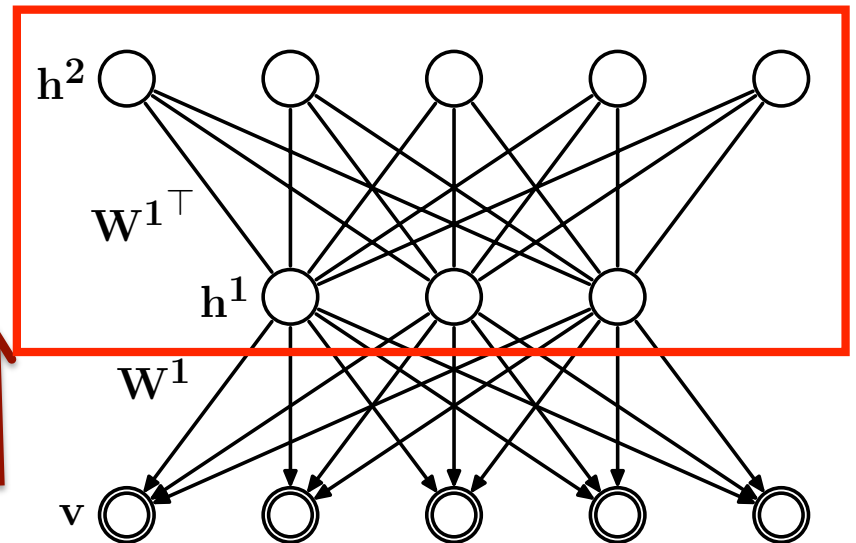
$$\log P_{\theta}(\mathbf{v}) = \sum_{\mathbf{h}^1} P_{\theta}(\mathbf{v}, \mathbf{h}^1)$$

$$\geq \sum_{\mathbf{h}^1} Q(\mathbf{h}^1 | \mathbf{v}) \left[ \log P(\mathbf{h}^1) + \log P(\mathbf{v} | \mathbf{h}^1) \right] + \mathcal{H}(Q(\mathbf{h}^1 | \mathbf{v}))$$



# Why this Pre-training Works?

- Greedy training improves variational lower bound.
- RBM and 2-layer DBN are equivalent when  $W^2 = W^{1\top}$ .
- The lower bound is tight and the log-likelihood improves by greedy training.



- For any approximating distribution  $Q(\mathbf{h}^1 | \mathbf{v})$

$$\log P_{\theta}(\mathbf{v}) = \sum_{\mathbf{h}^1} P_{\theta}(\mathbf{v}, \mathbf{h}^1)$$

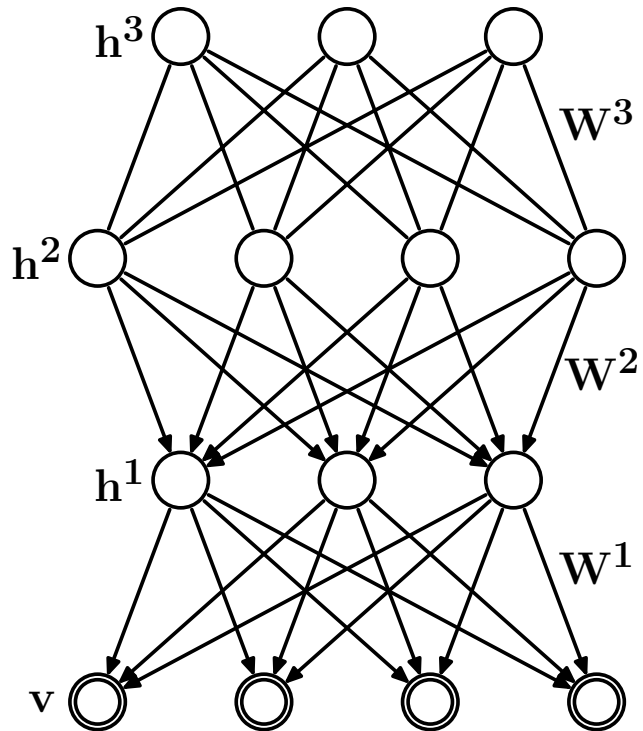
$$\geq \sum_{\mathbf{h}^1} Q(\mathbf{h}^1 | \mathbf{v}) \left[ \log P(\mathbf{h}^1) + \log P(\mathbf{v} | \mathbf{h}^1) \right] + \mathcal{H}(Q(\mathbf{h}^1 | \mathbf{v}))$$

Train 2<sup>nd</sup>-layer RBM

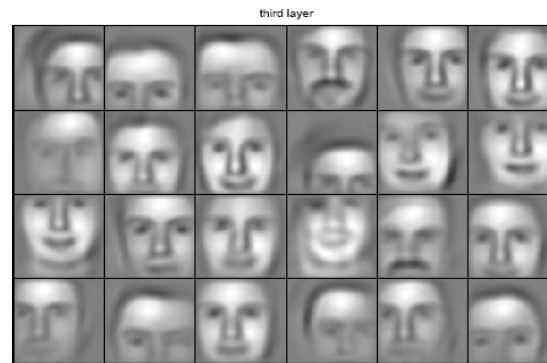
$Q(\mathbf{h}^1 | \mathbf{v})$

# Learning Part-based Representation

Convolutional DBN



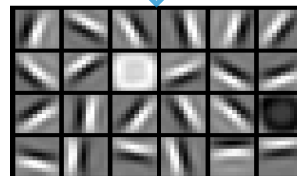
Faces



Groups of parts.



Object Parts

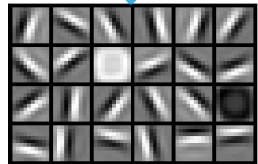


Trained on face images.

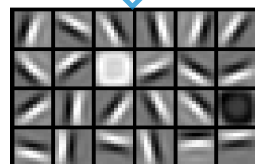
(Lee, Grosse, Ranganath, Ng, ICML 2009)

# Learning Part-based Representation

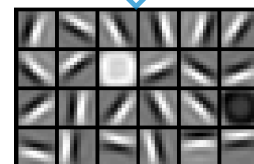
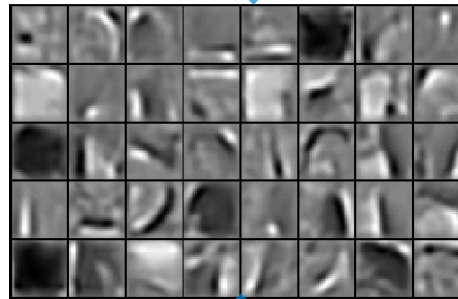
Faces



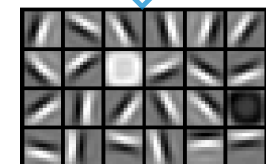
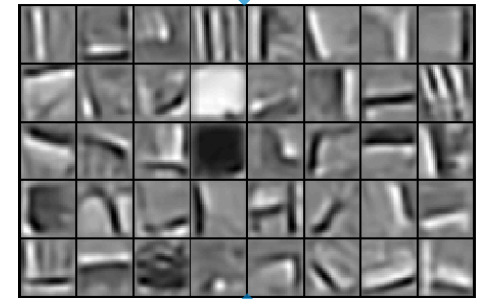
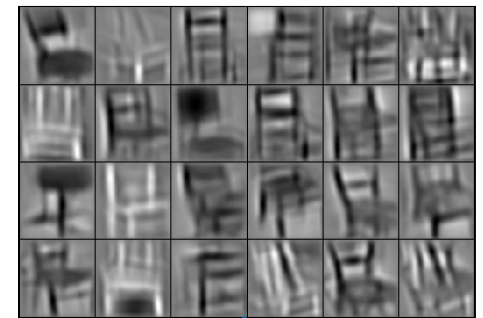
Cars



Elephants

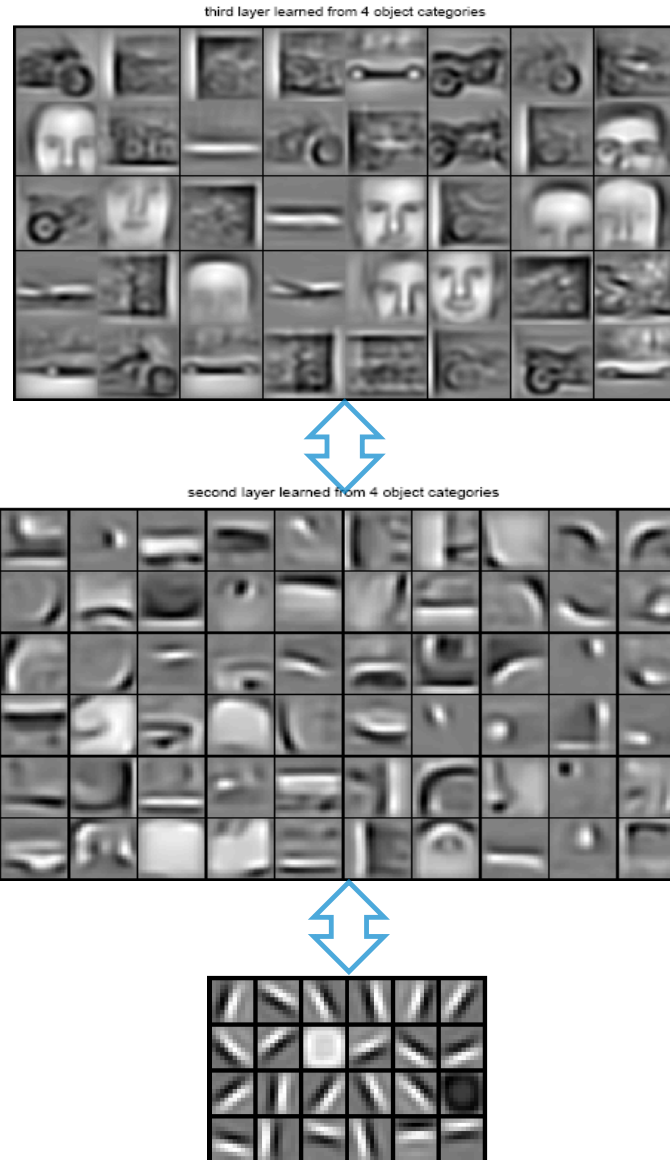


Chairs



(Lee, Grosse, Ranganath, Ng, ICML 2009)

# Learning Part-based Representation



Groups of parts.

Class-specific object parts

Trained from multiple classes (cars, faces, motorbikes, airplanes).

(Lee, Grosse, Ranganath, Ng, ICML 2009)