# Advanced Hierarchical Models

Russ Salakhutdinov
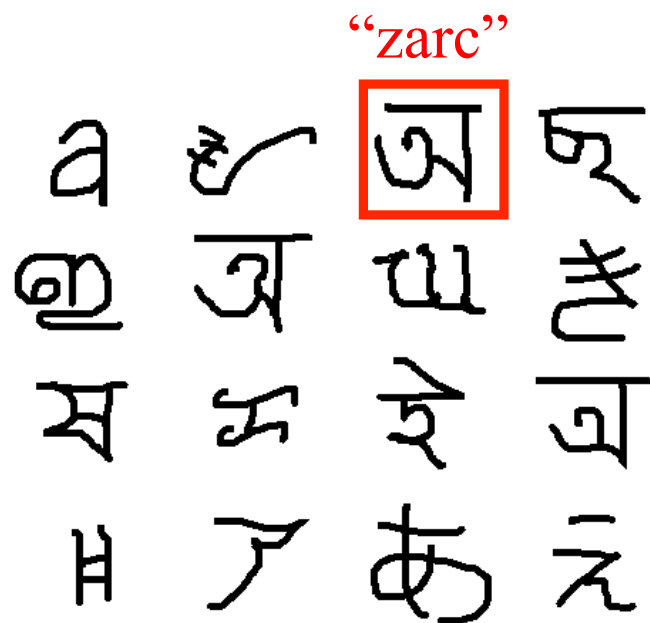
Department of Statistics and Computer Science
University of Toronto

# Motivation

- Learning abstract representations that support transfer to novel tasks, lies at the core of many problems in computer vision, speech perception, natural language processing, and machine learning.

- In many machine learning applications performance is measured using hundreds or thousands of training examples.

- For human learners, a single example of a novel category is often sufficient to make meaningful generalizations to novel instances.

Goal: Transfer higher-order knowledge abstracted from previously learned concept to infer parameters of a novel concept from few examples.

# One-shot Learning



"zarc"   "segway"

How can we learn a novel concept – a high dimensional statistical object – from few examples.

(Lake, Salakhutdinov, Gross, Tenenbaum, CogSci 2011)

# Traditional Supervised Learning



Segway

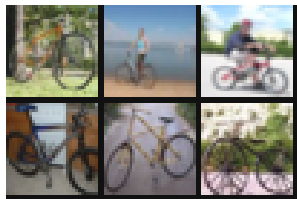Motorcycle

Test:
What is this?

# Learning to Transfer
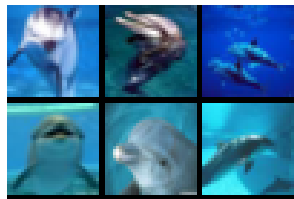
## Background Knowledge

Millions of unlabeled images



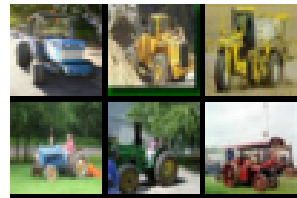Some labeled images



Bicycle

Dolphin

Elephant

Tractor

Learn to Transfer Knowledge



Learn novel concept from one example

Test:
What is this?

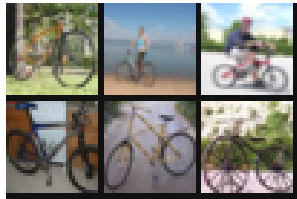# Learning to Transfer
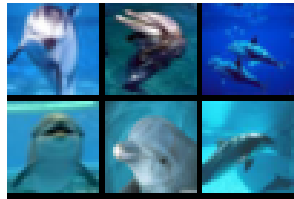
Background Knowledge

Millions of unlabeled images

Learn to Transfer Knowledge

Key problem in computer vision, speech perception, natural language processing, and many other domains.
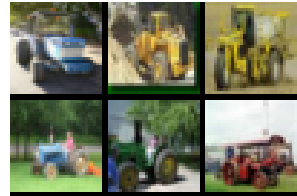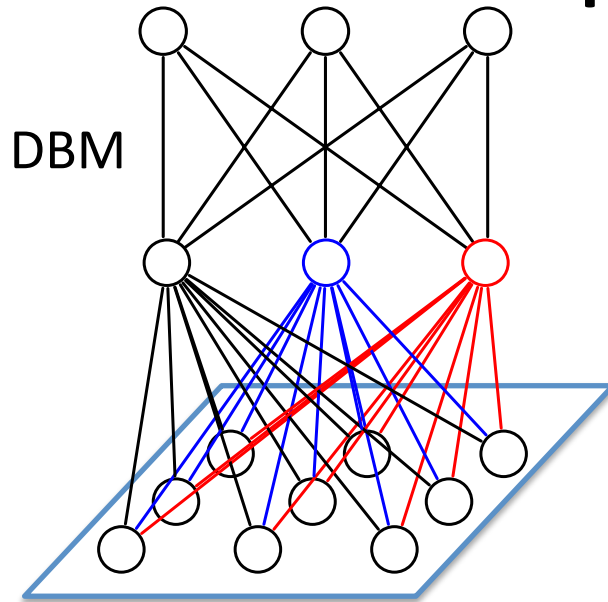
Some labeled images

Bicycle

Dolphin

Elephant

Tractor

Learn novel concept from one example
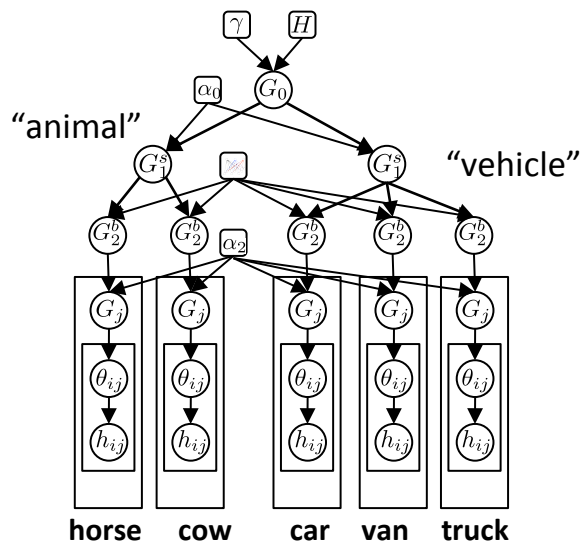
Test:
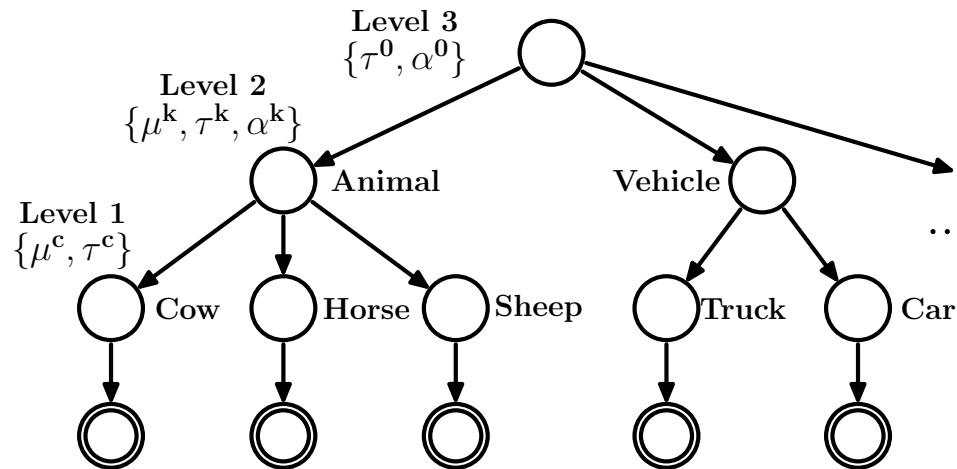What is this?

# Talk Roadmap



DBM

## Part 2: Advanced Hierarchical Models

- Introduction: Transfer Learning/ One-Shot Learning.

- **Compound Hierarchical Deep Models:**
  - Deep Boltzmann Machines.
  - Hierarchical Latent Dirichlet Allocation Model.

- Applications.

- Conclusions

# Hierarchical Bayes



Level 3
$\{\tau^0, \alpha^0\}$

Level 2
$\{\mu^k, \tau^k, \alpha^k\}$

Level 1
$\{\mu^c, \tau^c\}$

Animal

Vehicle

...

Cow   Horse   Sheep   Truck   Car

Hierarchical Bayesian
Models

**Hierarchical Prior.**

Probability of observed
data given parameters

Prior probability of
weight vector W

Posterior probability of
parameters given the
training data D.

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})}$$

- Fei-Fei, Fergus, and Perona, TPAMI 2006
- E. Bart, I. Porteous, P. Perona, and M. Welling, CVPR 2007
- Miller, Matsakis, and Viola, CVPR 2000
- Sivic, Russell, Zisserman, Freeman, and Efros, CVPR 2008

# Hierarchical-Deep Models

**HD Models:** Compose hierarchical Bayesian models with deep networks, two influential approaches from unsupervised learning
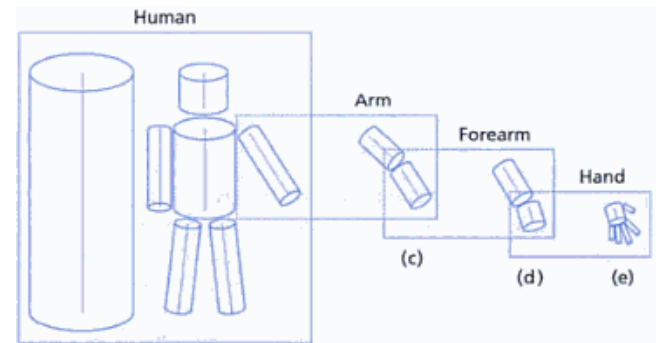
## Deep Networks:

• learn multiple **layers of nonlinearities.**
• trained in unsupervised fashion -- **unsupervised feature learning** – no need to rely on human-crafted input representations.
• **labeled data** is used to slightly adjust the model for a specific task.
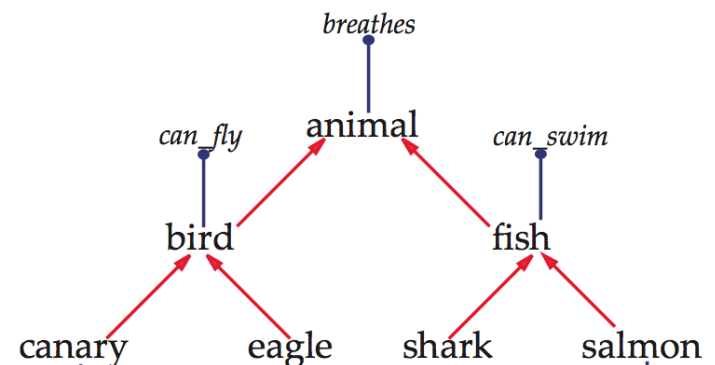
## Hierarchical Bayes:

• **explicitly represent category hierarchies** for sharing abstract knowledge.
• explicitly identify only a **small number of parameters** that are relevant to the new concept being learned.

### Deep Nets
### Part-based Hierarchy



Marr and Nishihara (1978)

### Hierarchical Bayes
### Category-based Hierarchy



Collins & Quillian (1969)

(Salakhutdinov, Tenenbaum, Torralba, NIPS 2011)

# Motivation for Our Approach

Learning to transfer knowledge:

**Hierarchical**

- Super-category: "A segway looks like a funny kind of vehicle".

- Higher-level features, or parts, shared with other classes:
  - ➢ wheel, handle, post

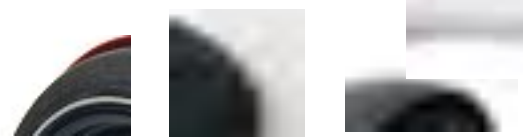- Lower-level features:
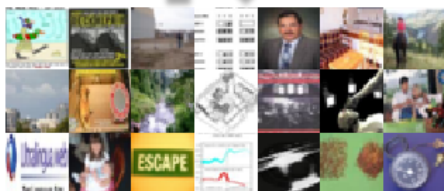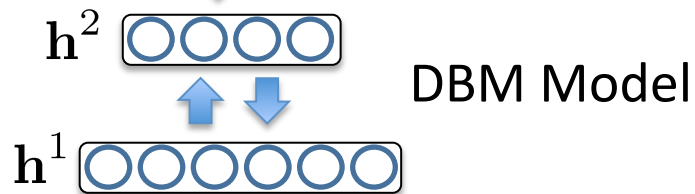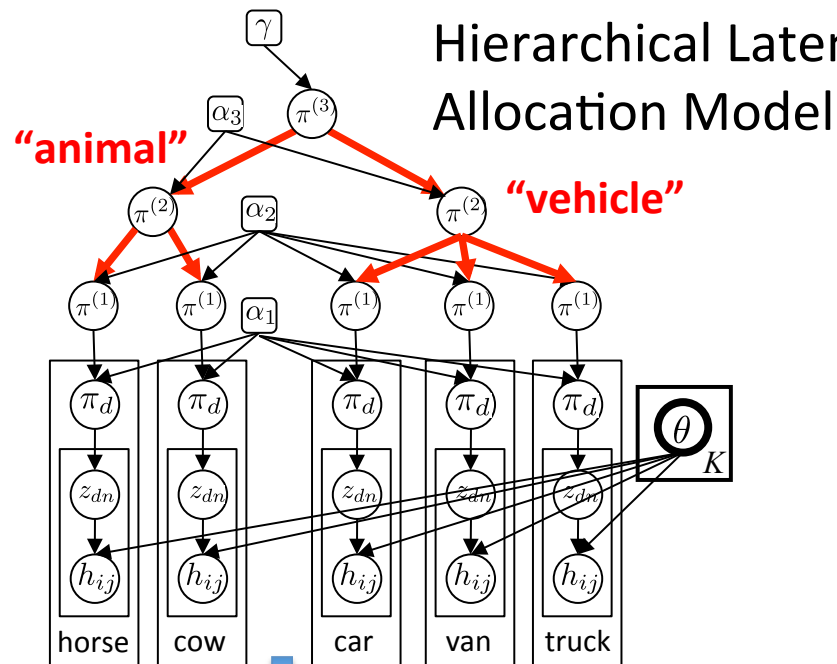  - ➢ edges, composition of edges

**Deep**

Super-class

Segway

Parts

Edges

# Hierarchical Generative Model



Hierarchical Latent Dirichlet Allocation Model

"animal"

"vehicle"

horse    cow    car    van    truck

DBM Model

$\mathbf{h}^2$
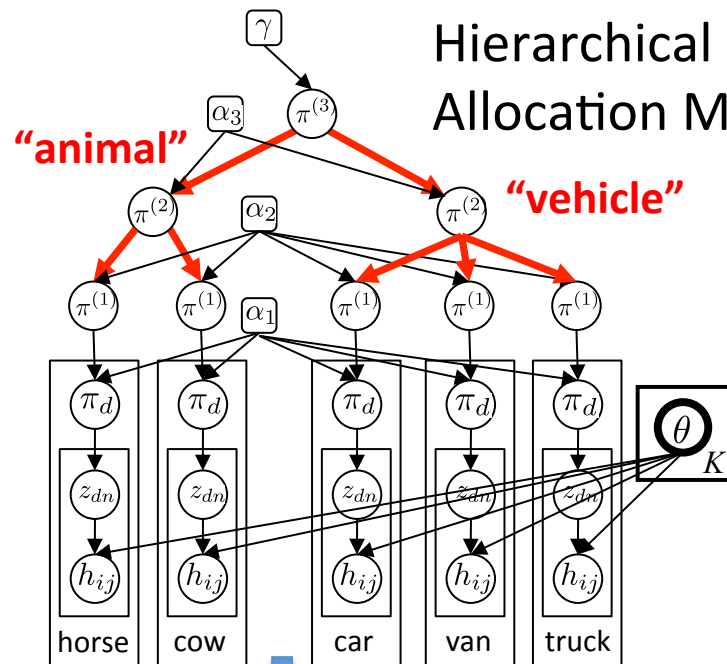
$\mathbf{h}^1$

Images

**Lower-level generic features:**
- edges, combination of edges

# Hierarchical Generative Model



Hierarchical Latent Dirichlet Allocation Model

"animal"

"vehicle"

horse    cow    car    van    truck

$\mathbf{h}^2$

DBM Model

$\mathbf{h}^1$

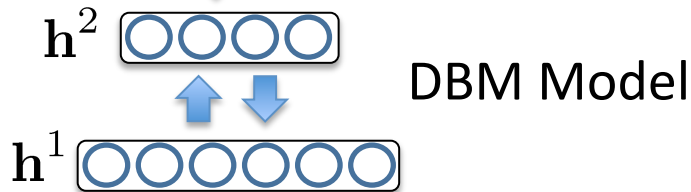Images

**Hierarchical Organization of Categories:**

- express priors on the features that are typical of different kinds of concepts
- modular data-parameter relations
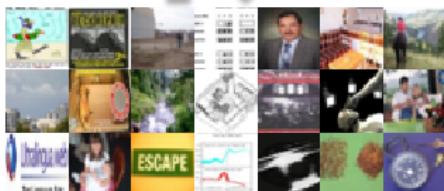
**Higher-level class-sensitive features:**

- capture distinctive perceptual structure of a specific concept

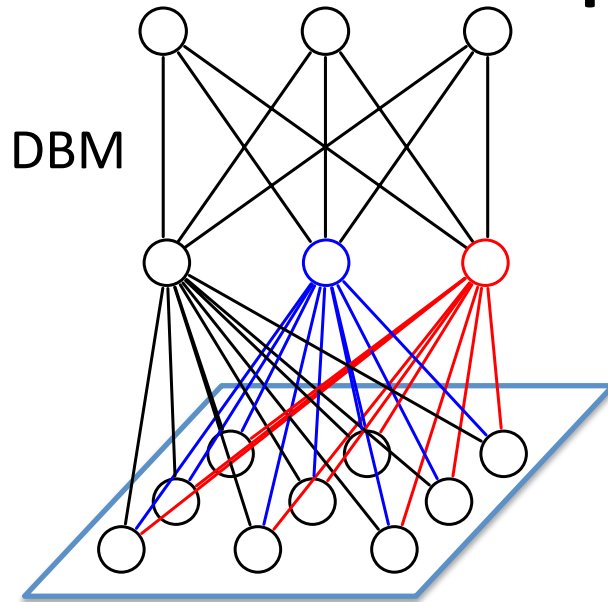**Lower-level generic features:**

- edges, combination of edges

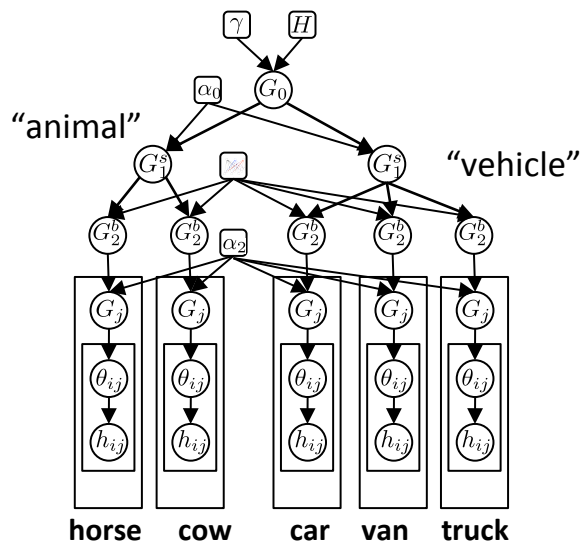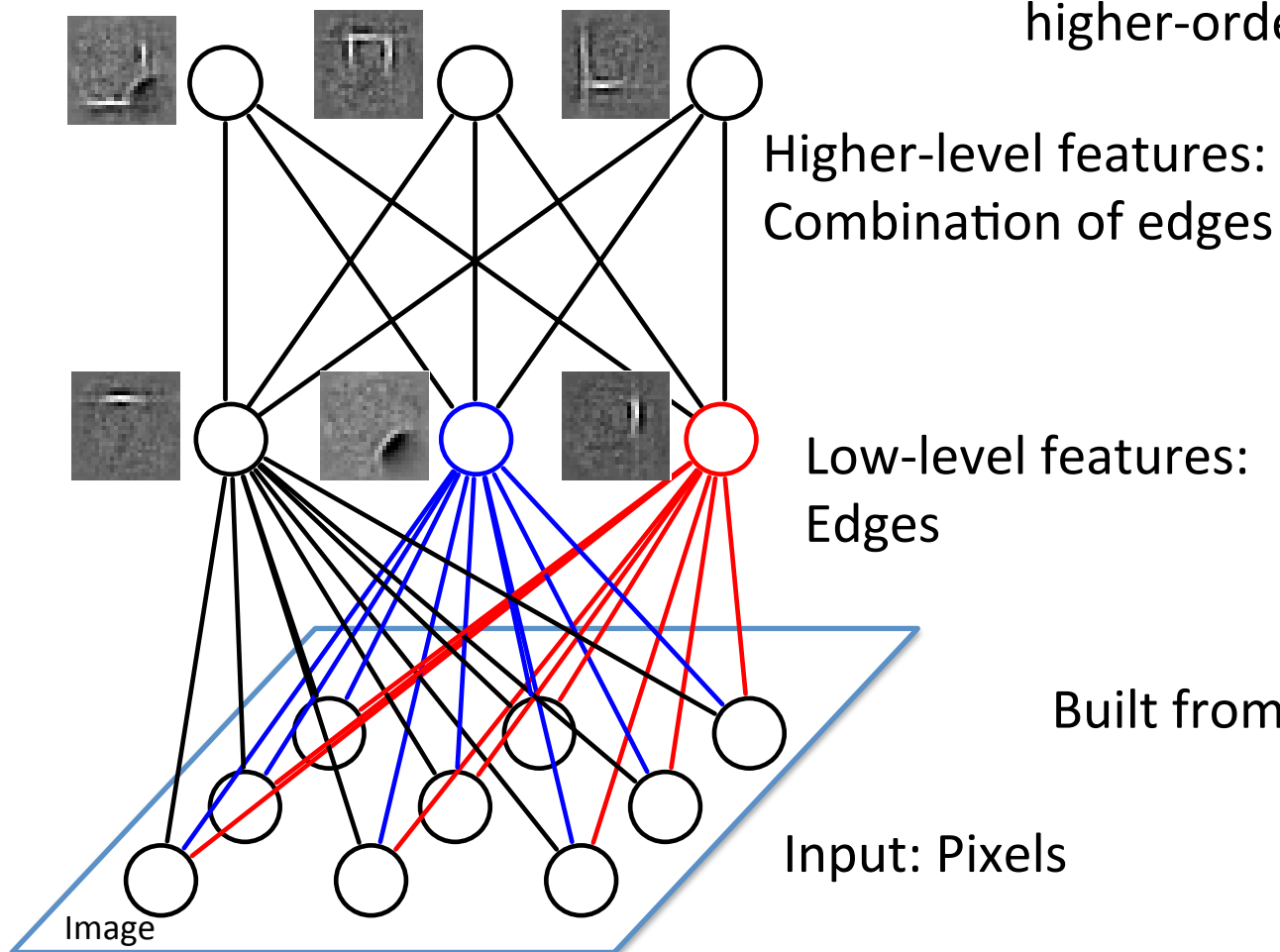# Talk Roadmap

DBM

## Part 2: Advanced Hierarchical Models

- Introduction: Transfer Learning/ One-Shot Learning.

- **Hierarchical Deep Models:**

  - **Deep Boltzmann Machines.**

  - Hierarchical Latent Dirichlet Allocation Model.

- Applications.

- Conclusions

# Deep Boltzmann Machines



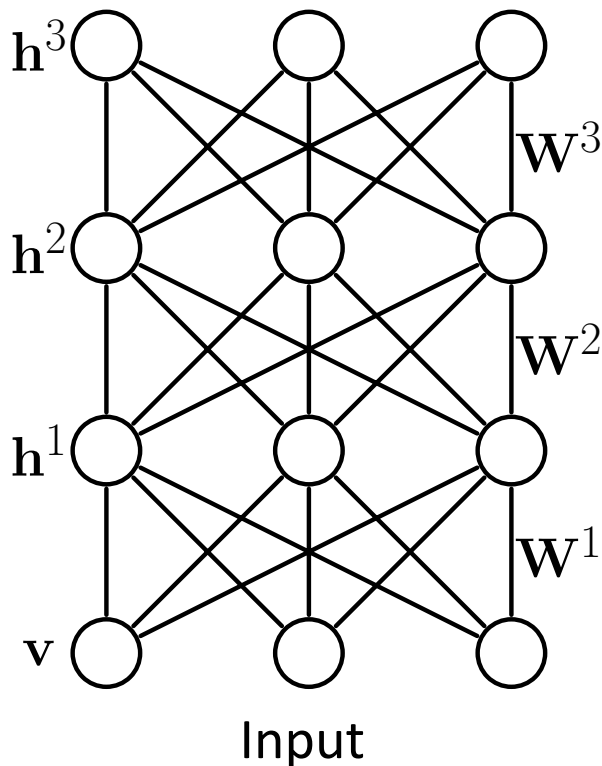Internal representations capture higher-order statistical structure

Higher-level features:
Combination of edges

Low-level features:
Edges

Built from **unlabeled** inputs.

Input: Pixels

Image

# A Brief Review

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

Deep Boltzmann Machine

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



$\mathbf{h}^3$

$\mathbf{W}^3$

$\mathbf{h}^2$

$\mathbf{W}^2$

$\mathbf{h}^1$
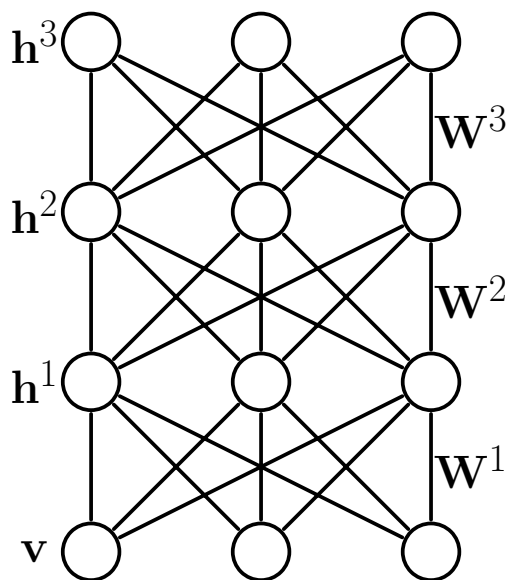
$\mathbf{W}^1$

$\mathbf{v}$

Input

- Dependencies between hidden variables.

- All connections are undirected.

- Bottom-up and Top-down:

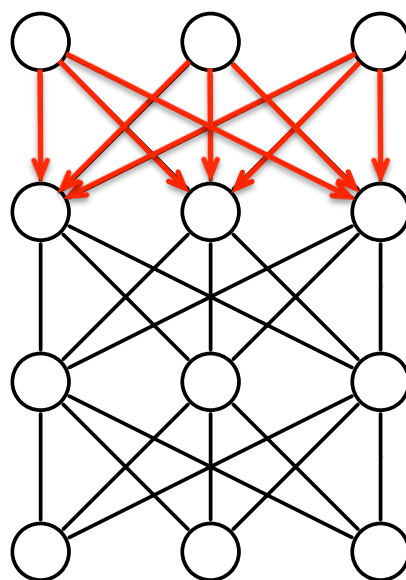# Decomposition

The joint probability can be decomposed:

$$P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3) = \underbrace{P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2 | \mathbf{h}^3)}_{\text{Conditional DBM}} \underbrace{P_\theta(\mathbf{h}^3)}_{\text{Prior term}}$$



DBM

Conditional DBM

Replace the last term with more structured hierarchical prior.

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

# Stage-wise Learning

The joint probability can be decomposed:

$$P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3) = \underbrace{P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2 | \mathbf{h}^3)}_{\text{Conditional DBM}} \underbrace{P_\theta(\mathbf{h}^3)}_{\text{Prior term}}$$

DBMs approximate intractable posterior $P_\theta(\mathbf{h}|\mathbf{v})$ with fully factorized tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$. The variational lower-bound takes form:
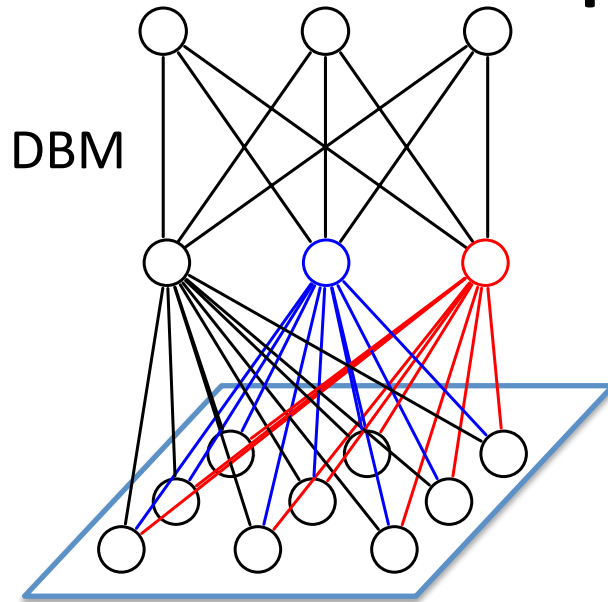
$$\log P_\theta(\mathbf{v}) \geq \sum_{\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3} \underbrace{Q_\mu(\mathbf{h}^1, \mathbf{h}^2 | \mathbf{v}) \left[ \log P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2 | \mathbf{h}^3) \right]}_{\text{Likelihood term}} + \underbrace{\mathcal{H}(Q_\mu(\mathbf{h}|\mathbf{v}))}_{\text{Entropy functional}}$$

$$+ \underbrace{\sum_{\mathbf{h}^3} Q_\mu(\mathbf{h}^3 | \mathbf{v}) \log P_\theta(\mathbf{h}^3)}_{\textcolor{blue}{\textbf{Fit Hierarchical LDA prior}}}$$

$$\mathcal{H}(Q_\mu(\mathbf{h}|\mathbf{v})) = \sum_{\mathbf{h}} Q_\mu(\mathbf{h}|\mathbf{v}) \log \frac{1}{Q_\mu(\mathbf{h}|\mathbf{v})}$$

# Stage-wise Learning

The joint probability can be decomposed:

$$P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3) = \underbrace{P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2|\mathbf{h}^3)}_{\text{Conditional DBM}}\underbrace{P_\theta(\mathbf{h}^3)}_{\text{Prior term}}$$

DBMs approximate intractable posterior $P_\theta(\mathbf{h}|\mathbf{v})$ with fully factorized tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$. The variational lower-bound takes form:

$$\log P_\theta(\mathbf{v}) \geq \sum_{\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3} Q_\mu(\mathbf{h}^1, \mathbf{h}^2|\mathbf{v})\underbrace{\left[\log P_\theta(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2|\mathbf{h}^3)\right]}_{\text{Likelihood term}} + \underbrace{\mathcal{H}(Q_\mu(\mathbf{h}|\mathbf{v}))}_{\text{Entropy functional}}$$

$Q_\mu(\mathbf{h}^3|\mathbf{v})\log P_\theta(\mathbf{h}^3)$

- Learn DBM.
- Using variational inference, infer the states of the top-level variables and fit an LDA prior.

**Fit Hierarchical LDA prior**
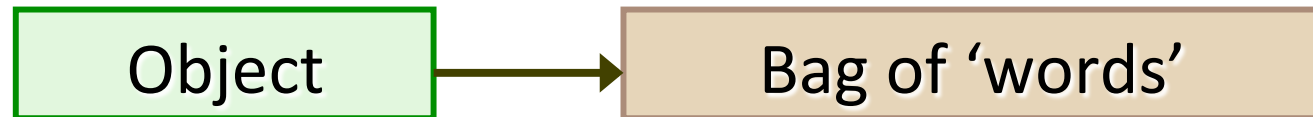
# Talk Roadmap



DBM

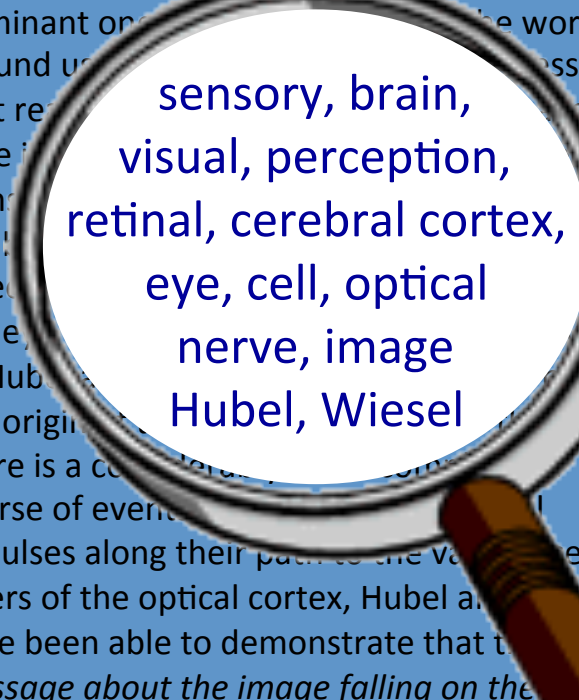## Part 2: Advanced Hierarchical Models

- Introduction: Transfer Learning/ One-Shot Learning.

- **Compound Hierarchical Deep Models:**

  - Deep Boltzmann Machines.

  - Hierarchical Latent Dirichlet Allocation Model.

- Applications.

- Conclusions.



"animal"          "vehicle"

horse    cow    car    van    truck

# Bag of Words Representation

Object ⟶ Bag of 'words'

# Analogy to Documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant one. Our perception of the world around us is based essentially on the messages that reach the brain from our eyes. For a long time it was thought that the retinal image was transmitted point by point to visual centers in the brain; the cerebral cortex was a movie screen, so to speak, upon which the image in the eye was projected. Through the discoveries of Hubel and Wiesel we now know that behind the origin of the visual perception in the brain there is a considerably more complicated course of events. By following the visual impulses along their path to the various cell layers of the optical cortex, Hubel and Wiesel have been able to demonstrate that the *message about the image falling on the retina undergoes a step-wise analysis in a system of nerve cells stored in columns*

**sensory, brain, visual, perception, retinal, cerebral cortex, eye, cell, optical nerve, image Hubel, Wiesel**

China is forecasting a trade surplus of $90bn (£51bn) to $100bn this year, a threefold increase on 2004's $32bn. The Commerce Ministry said the surplus would be created by a predicted 30% jump in exports to $750bn, compared with a 18% rise in imports to $660bn. The figures are likely to further annoy the US, which has long argued that China's exports are unfairly helped by a deliberately undervalued yuan. Beijing agrees the surplus is too high, but says the yuan is only one factor. Bank of China governor Zhou Xiaochuan said the country also needed to do more to boost domestic demand so more goods stayed within the country. China increased the value of the yuan against the dollar by 2.1% in July and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**

**Intuition**: Documents contain multiple topics.

# Latent Dirichlet Allocation

**Text document**

The William Randolph Hearst Foundation will give $1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be $200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive $400,000 each. The Juilliard School, where music and the performing arts are taught, will get $250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual $100,000 donation, too.

**Discovered topics**

| "Arts" | "Budgets" | "Children" | "Education" |
|--------|-----------|------------|-------------|
| NEW | MILLION | CHILDREN | SCHOOL |
| FILM | TAX | WOMEN | STUDENTS |
| SHOW | PROGRAM | PEOPLE | SCHOOLS |
| MUSIC | BUDGET | CHILD | EDUCATION |
| MOVIE | BILLION | YEARS | TEACHERS |
| PLAY | FEDERAL | FAMILIES | HIGH |
| MUSICAL | YEAR | WORK | PUBLIC |
| BEST | SPENDING | PARENTS | TEACHER |
| ACTOR | NEW | SAYS | BENNETT |
| FIRST | STATE | FAMILY | MANIGAT |
| YORK | PLAN | WELFARE | NAMPHY |
| OPERA | MONEY | MEN | STATE |
| THEATER | PROGRAMS | PERCENT | PRESIDENT |
| ACTRESS | GOVERNMENT | CARE | ELEMENTARY |
| LOVE | CONGRESS | LIFE | HAITI |

Blei, et al. 2003

# Latent Dirichlet Allocation

Generative Process: $\mathbf{w} \sim \mathrm{LDA}$

Draw each topic $\theta_k \sim \mathrm{Dir}(\eta)$ for $k = 1..., K$

For each document d:

- Draw topic proportions $\pi_d \sim \mathrm{Dir}(\alpha)$
- For each word:
  - Draw topic indicator $z_{d,n} \sim \mathrm{Mult}(\pi_d)$
  - Draw word $\quad\quad w_{d,n} \sim \mathrm{Mult}(\theta_{z_{d,n}})$



Pr(topic | doc)

Pr(word | topic)

# Latent Dirichlet Allocation

Generative Process: $\mathbf{w} \sim \mathrm{LDA}$

Draw each topic $\theta_k \sim \mathrm{Dir}(\eta)$ for $k = 1 \ldots, K$

For each document:

- Draw topic proportions $\pi_d \sim \mathrm{Dir}(\alpha)$

- For each word:

  - Draw topic indicator $z_{d,n} \sim \mathrm{Mult}(\pi_d)$

  - Draw word $\quad w_{d,n} \sim \mathrm{Mult}(\theta_{z_{d,n}})$



Pr(topic | doc)

Pr(word | topic)



The William Randolph Hearst Foundation will give $1.25
tan Opera Co., New York Philharmonic and Juilliard So
real opportunity to make a mark on the future of the per
every bit as important as our traditional areas of support i
and the social services," Hearst Foundation President 1
announcing the grants. Lincoln Center's share will be $
will house young artists and provide new public faciliti
New York Philharmonic will receive $400,000 each. Th
the performing arts are taught, will get $250,000. The H₁
of the Lincoln Center Consolidated Corporate Fund, 
donation, too.

| "Arts" | "Budgets" | "Children" | "Education" |
|---|---|---|---|
| NEW | MILLION | CHILDREN | SCHOOL |
| FILM | TAX | WOMEN | STUDENTS |
| SHOW | PROGRAM | PEOPLE | SCHOOLS |
| MUSIC | BUDGET | CHILD | EDUCATION |
| MOVIE | BILLION | YEARS | TEACHERS |
| PLAY | FEDERAL | FAMILIES | HIGH |
| MUSICAL | YEAR | WORK | PUBLIC |
| BEST | SPENDING | PARENTS | TEACHER |
| ACTOR | NEW | SAYS | BENNETT |
| FIRST | STATE | FAMILY | MANIGAT |
| YORK | PLAN | WELFARE | NAMPHY |
| OPERA | MONEY | MEN | STATE |
| THEATER | PROGRAMS | PERCENT | PRESIDENT |
| ACTRESS | GOVERNMENT | CARE | ELEMENTARY |
| LOVE | CONGRESS | LIFE | HAITI |

# Latent Dirichlet Allocation

Generative Process: $\mathbf{w} \sim \mathrm{LDA}$

Draw each topic $\theta_k \sim \mathrm{Dir}(\eta)$ for $k = 1..., K$

For each document:

- Draw topic proportions $\pi_d \sim \mathrm{Dir}(\alpha)$
- For each word:
  - Draw topic indicator $z_{d,n} \sim \mathrm{Mult}(\pi_d)$
  - Draw word $\quad w_{d,n} \sim \mathrm{Mult}(\theta_{z_{d,n}})$

**Remember**: compound HD model:

$\mathbf{h}^3 \sim \mathrm{LDA}$ prior

Words ⇔ activations of DBM's top-level units.
Topics ⇔ distributions over top-level units, or higher-level parts.



Pr(topic | doc)

Pr(word | topic)

# Intuition

$$\mathbf{h}^3 \sim \text{LDA prior}$$

Words ⇔ activations of DBM's top-level units.
Topics ⇔ distributions over top-level units, or higher-level parts.

Pr(topic | doc)

Pr(word | topic)

DBM generic features: **Words**

LDA high-level features: **Topics**

Images **Documents**

**Each topic is made up of words.**

**Each document is made up of topics.**

# Hierarchical LDA
# Modeling Super-Category Structure



- Draw **global** topic proportions: $\pi^{(3)} \sim \mathrm{Dir}(\gamma)$

- Draw **super-class specific** topic proportions:
$$\pi^{(2)} | \pi^{(3)} \sim \mathrm{Dir}(\alpha^{(3)} \pi^{(3)})$$

- Draw **class-class specific** topic proportions:
$$\pi^{(1)} | \pi^{(2)} \sim \mathrm{Dir}(\alpha^{(2)} \pi^{(2)})$$

- Draw **document specific** topic proportions:
$$\pi_d | \pi^{(1)} \sim \mathrm{Dir}(\alpha^{(1)} \pi^{(1)})$$

Nonparametric extension:
**Hierarchical Dirichlet Process (HDP).**

# Hierarchical LDA: Example

**Global** topic proportions:



**Super-class specific** topic proportions:

**Fruits:** apples, oranges, pears

**Aquatic animals:** dolphins, sharks.



**Class specific** topic proportions:

**Apples:**

**Oranges:**



**Image specific** topic proportions:

# Hierarchical LDA: Example

**Global** topic proportions:



**S**uper-class specific topic proportions:

**Fruits:** apples, oranges, pears

**Aquatic animals:** dolphins, sharks.



**Class specific** topic proportions:

**Apples:**

**Oranges:**



**I**n...
p...

> So far we have assumed
> a **fixed** hierarchy

# Modeling the Number of Super-Categories

Place Chinese Restaurant Process (CRP) Prior over the number of super-classes.

CRP defines a distribution on partition of integers.

Generating from $\mathrm{CRP}(\alpha)$:

Customers enter a restaurant with an unbounded number of tables, where the n$^{th}$ customer occupies a table k drawn from:

$$P(z_n = k | z_1, ..., z_{n-1}) = \begin{cases} \frac{n^k}{n-1+\alpha} & n^k > 0 \\ \frac{\alpha}{n-1+\alpha} & k \text{ is new} \end{cases}$$

where $n^k$ is the number of previous customers at table k and $\alpha$ is the concentration parameter.

Customers ⇔ integers, tables ⇔ clusters.

# Modeling the Hierarchy

# Modeling the Hierarchy

# Modeling the Hierarchy

# Modeling the Hierarchy



Expected number of clusters: $O(\alpha \log n)$

The nested CRP, nCRP, extends CRP to nested sequence of partitions, one for each level of the tree (Blei et.al. NIPS 2003).

# Hierarchical Deep Model

# Hierarchical Deep Model



Tree hierarchy of classes is learned

$\mathbf{z} \sim \text{nCRP}$ (**Nested Chinese Restaurant Process**)
prior: a nonparametric prior over tree structures.

# Hierarchical Deep Model

Tree hierarchy of classes is learned



$\mathbf{z} \sim \mathrm{nCRP}$ (**Nested Chinese Restaurant Process**) prior: a nonparametric prior over tree structures.

$\mathbf{h}^3 | \mathbf{z} \sim \mathrm{HDP}$(**Hierarchical Dirichlet Process**) prior: a nonparametric prior allowing categories to share higher-level features, or parts.

# Hierarchical Deep Model



$\mathbf{z} \sim \mathrm{nCRP}$ (**Nested Chinese Restaurant Process**) prior: a nonparametric prior over tree structures

$\mathbf{h}^3 | \mathbf{z} \sim \mathrm{HDP}$(**Hierarchical Dirichlet Process**) prior: a nonparametric prior allowing categories to share higher-level features, or parts.

$\mathbf{v} | \mathbf{h}^3 \sim \mathrm{DBM}$ **Conditional Deep Boltzmann Machine.**

Enforce (approximate) global consistency through many local constraints.

# Hierarchical Deep Model



Tree hierarchy of classes is learned

"animal"   "vehicle"

Unlike standard statistical models, in addition to inferring parameters, we also infer the hierarchy for sharing those parameters.

Topics

share higher-level features, or parts.

$$\mathbf{v}|\mathbf{h}^3 \sim \mathrm{DBM}$$ **Conditional Deep Boltzmann Machine.**

Enforce (approximate) global consistency through many local constraints.

# CIFAR Object Recognition



Tree hierarchy of classes is learned

"animal"

"vehicle"

Higher-level class sensitive features

Lower-level generic features

horse   cow   car   van   truck

50,000 images of 100 classes

**Inference: Markov chain Monte Carlo – Later!**

4 million unlabeled images

32 x 32 pixels x 3  RGB

# Learning to Learn

The model learns how to share the knowledge across many visual categories.



"global" — **Learned super-class hierarchy**

"aquatic animal"

"fruit"

"human"

dolphin     turtle   shark   ray      apple    orange   sunflower      girl    baby    man

**Basic level class**

woman

**Learned higher-level class-sensitive features**

...

**Learned low-level generic features**

...

# Learning to Learn

The model learns how to share the knowledge across many visual

# Sharing Features



**Learning to Learn:** Learning a hierarchy for sharing parameters – rapid learning of a novel concept.

# Object Recognition



Area under ROC curve for same/different
(1 new class vs. 99 distractor classes)

**Our model outperforms standard computer vision features (e.g. GIST).**

# Handwritten Character Recognition



"alphabet 1"

"alphabet 2"

Learned lower-level features

25,000 characters

**Edges**

# Handwritten Character Recognition

Area under ROC curve for same/different
(1 new class vs. 1000 distractor classes)



[Averaged over 40 test classes]

# Simulating New Characters



Real data within super class

Simulated new characters

# Simulating New Characters



Real data within super class

Simulated new characters

# Simulating New Characters



Real data within super class

Simulated new characters

# Simulating New Characters



Real data within super class

Simulated new characters

# Simulating New Characters

Real data within super class





Simulated new characters

# Simulating New Characters



Global

Super class 1

Super class 2

Class 1    Class 2    New class

Real data within super class

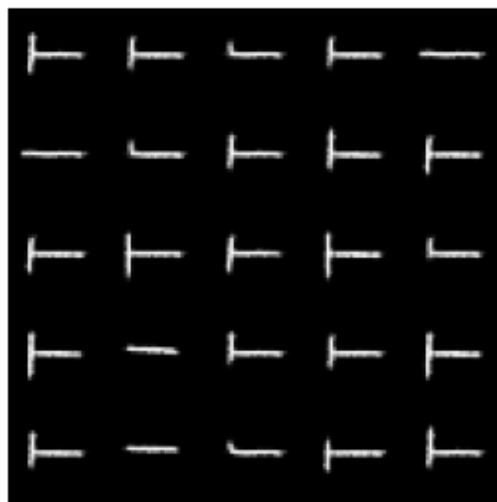Simulated new characters
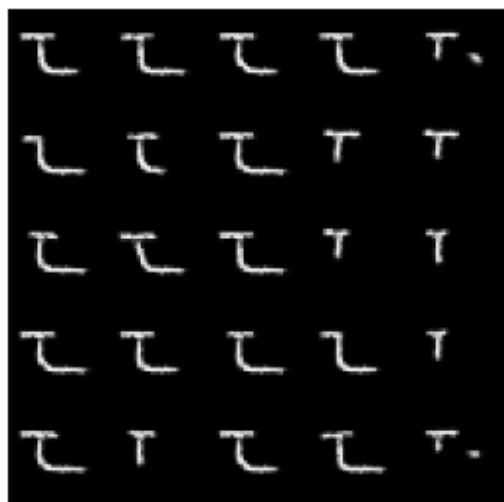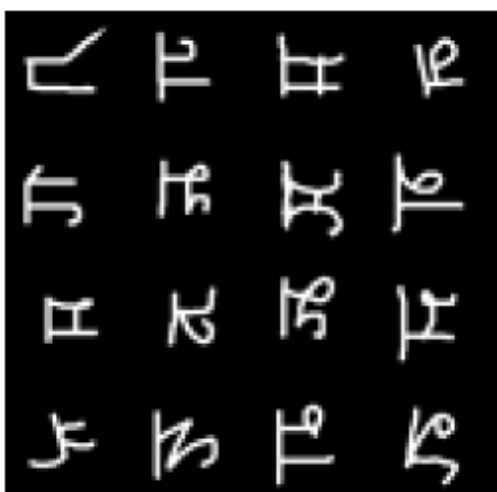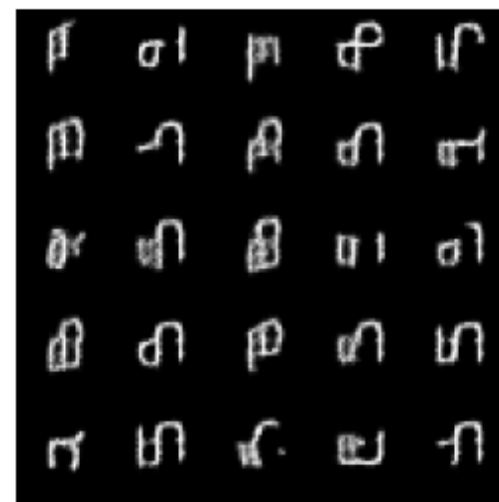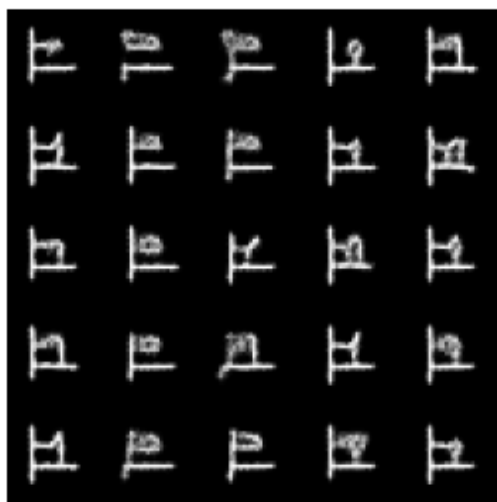
# Simulating New Characters



Real data within super class

Simulated new characters

# Learning from very few examples



3 examples of
a new class

Conditional samples
in the same class

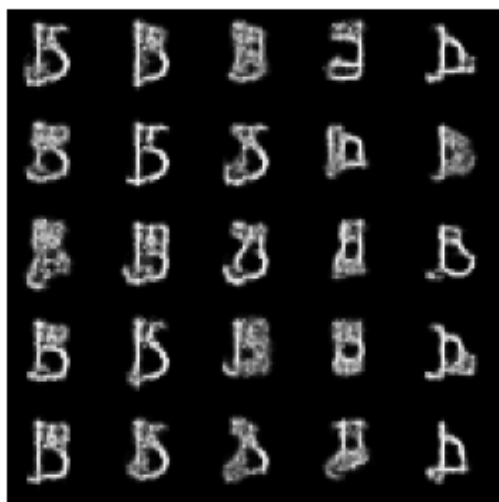Inferred super-class

# Learning from very few examples

# Learning from very few examples

# Learning from very few examples

# Learning from very few examples

# Learning from very few examples

# Learning from very few examples

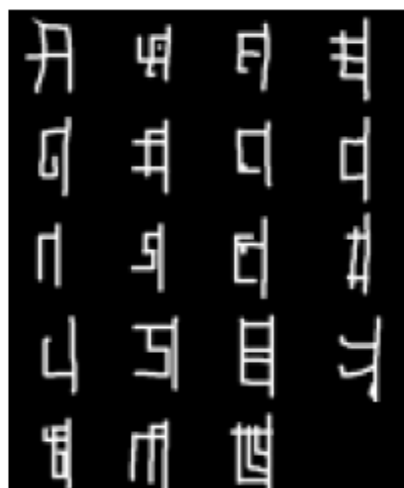# Learning from very few examples
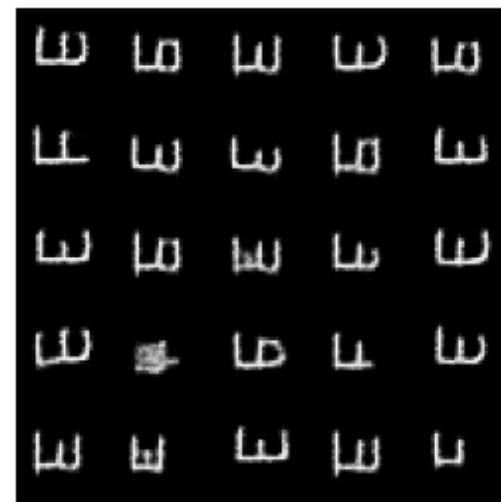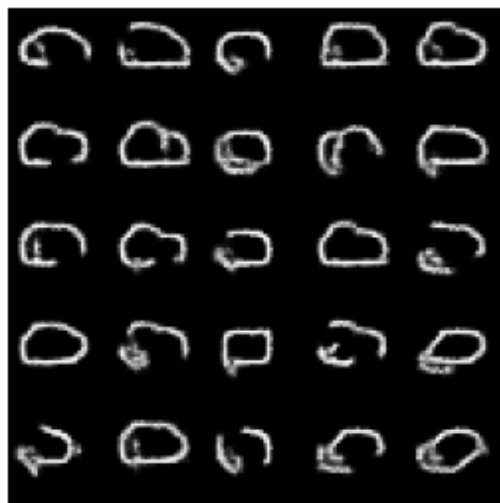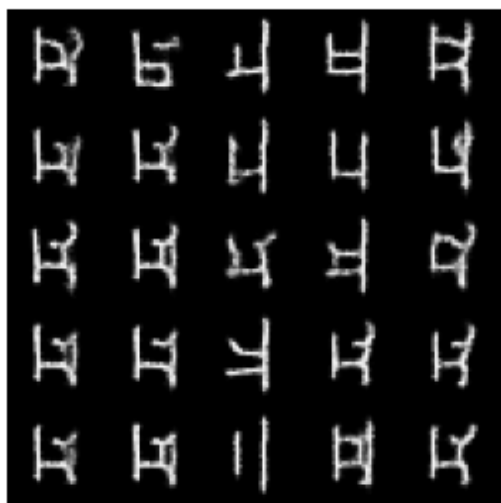
# Learning from very few examples

# Motion Capture

# Motion Capture



Walk

Drunken Walk

$\gamma$ $H$ $\alpha_0$ $G_0$ $G_1^s$ $\alpha_1$ $G_1^s$ $G_2^b$ $G_2^b$ $\alpha_2$ $G_2^b$ $G_2^b$ $G_2^b$ $G_j$ $G_j$ $G_j$ $G_j$ $G_j$ $\theta_{ij}$ $\theta_{ij}$ $\theta_{ij}$ $\theta_{ij}$ $\theta_{ij}$ $h_{ij}$ $h_{ij}$ $h_{ij}$ $h_{ij}$ $h_{ij}$

Time →

Sexy walk ROC curve

HDP-DBM

Input space distance
(no hierarchy)

detection rate

false alarm rate

# Motion Capture



Walk

Drunken Walk

The same model can be applied to speech, text, video, or any other high-dimensional data.

HDP-DBM

Input space distance (no hierarchy)

detection rate

false alarm rate

Time

# Talk Roadmap

DBM

## Part 2: Advanced Hierarchical Models

- Introduction: Transfer Learning/ One-Shot Learning.

- Compound Hierarchical Deep Models:

  – Deep Boltzmann Machines.

  – Hierarchical Latent Dirichlet Allocation Model.

- Applications.

- Conclusions

# Other Hierarchical Models

At a minimum, object categorization requires information about
• category mean (prototype)
• variances along each dimension (similarity metric)



Color features vary strongly, whereas shape features vary weakly.



A single example provides some information about the prototype, but not about the variances.

# Learning Class-Specific Similarity Metrics



**Dog**

**Sheep**

**Horse**

**Novel Category: Cow**

**Car** **Van** **Truck**

(Salakhutdinov, Tenenbaum, & Torralba, JMLR WC&P 2012)

# Learning Class-Specific Similarity Metrics



Dog
Sheep
Horse
Cow
Car  Van  Truck

(Salakhutdinov, Tenenbaum, & Torralba, JMLR WC&P 2012)

# Learning Class-Specific Similarity Metrics



In order to transfer appropriate similarity metric, the model needs to discover how to group related categories into super-categories.

# Hierarchical Bayes



- Probabilistic linear model with Gaussian observation noise:

$$P(x|z = c) = N(\mu^c, 1/\tau^c)$$

- Place a conjugate Normal-Gamma prior over the means and precision parameters:

$$P(\mu^c, \tau^c) = \mathcal{N}(\mu^k 1/(\nu\tau^c))\Gamma(\alpha^k, \tau^k)$$

**Hierarchical Prior.**

As before, infer the hierarchy.

# Image Retrieval

MSR Cambridge
Dataset



aeroplanes
benches and chairs
bicycles/single
cars/front
cars/rear
cars/side
signs

buildings
chimneys
doors'
scenes/office
scenes/urban
windows

trees
birds
flowers
leaves
scenes/countryside

forks
knives
spoons

animals/cows
animals/sheep

clouds

Retrieved images with our model



Query image



Given only one
examples of a cpw

Nearest neighbor

# Unsupervised Category Discovery

Can we discover when the model has encountered novel categories, and how can we break up new instances into novel categories?

The test set consists of **many unlabeled examples from an unknown number of basic-level classes.**

**Existing Categories**

| | | |
|---|---|---|
| Novel: 0.01<br>Car: 0.99 | Novel: 0.02<br>Plane: 0.97 | Novel: 0.02<br>Bench: 0.92 |



**Novel Categories**

| | | |
|---|---|---|
| Novel: 0.28<br>Countryside: 0.53 | Novel: 0.42<br>Building: 0.49 | Novel: 0.87<br>Bird: 0.11 |



**Existing Categories**



**Novel Categories**



With 18 unlabeled test images the model correctly places nine familiar images in nine different basic-level categories, while also correctly forming three novel categories with 3 examples each.

# Object Detection Challenge

Consider challenging object detection task.



By looking at the output of a detector, can you guess which object is it trying to detect?

# Learning from Few Examples

SUN database



Classes sorted by frequency

Rare objects are similar to frequent objects

(Salakhutdinov, Torralba, & Tenenbaum, CVPR 2011)

# Learning from Few Examples



chair

armchair

Swivel chair

Deck chair

Number of training examples

300
250
200
150
100
50

Classes sorted by frequency

# Generative Model of Classifier Parameters

Many state-of-the-art object detection systems use sophisticated models, based on multiple parts with separate appearance and shape components.

$$y = \beta^\top \Phi(\mathbf{x})$$



Detect objects by testing sub-windows and scoring corresponding test patches with a linear function.

**We can define hierarchical prior over parameters of discriminative model and learn the hierarchy.**

**Image Specific:** concatenation of the HOG feature pyramid at multiple scales.

Felzenszwalb, McAllester & Ramanan, 2008

# Generative Model of Classifier Parameters

By learning hierarchical structure, we can improve the current state-of-the-art.

Sun Dataset: 32,855 examples of 200 categories



Hierarchical Bayes

Level 1

$\theta_1^{(1)}$ Animal

Global

Vehicle

Level 2

$\theta_1^{(2)}$ Horse

$\theta_2^{(2)}$ Cow

$\theta_3^{(2)}$ Car

$\theta_4^{(2)}$ Van

Truck

185 ex    27 ex    12 ex

Hierarchical Model



Single Class

# Truck

**Single classifier**



**Hierarchical Model**

**Dome**

Single classifier

Hierarchical Model

Number of training examples

person
column
personsitting
sculpture
pole
bottle
bottles
people
tombstone
bicycle
statue

painting
picture
mirror
text
screen
poster
television
monitor
microwave
oven
speaker

car
truck
airplane
hat
van
bus
...cars

# Generative Model of Matrix Factorizations

Image bases

Relational data

Gene expression data



Karklin and Lewicki (2009)

Kemp et.al. (2006)

Meeds et al. (2007)

How can we automatically choose the right structure from raw data?

Context free grammar:

US Senate votes:

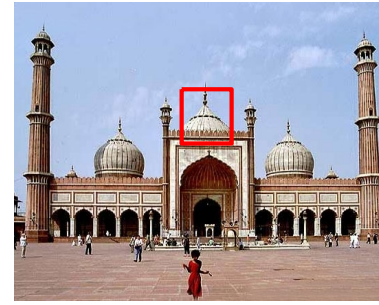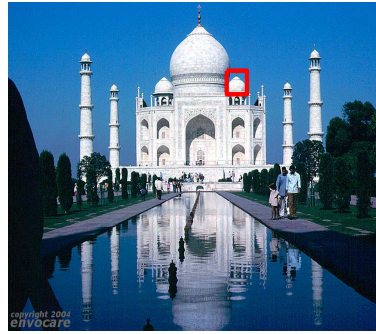| | | |
|---|---|---|
| low-rank | $G \to GG + G$ | (1) |
| clustering | $G \to MG + G \mid GM^T + G$ | (2) |
| | $M \to MG + G$ | (3) |
| linear dynamics | $G \to CG + G \mid GC^T + G$ | (4) |
| | $C \to CG + G$ | (5) |
| sparsity | $G \to \exp(G) \circ G$ | (6) |
| binary factors | $G \to BG + G \mid GB^T + G$ | (7) |
| | $M \to B$ | (8) |



(a) Level 1: $GG + G$    (b) Level 2: $(MG + G)G + G$    (c) Level 3: $(MG + G)(GM^T + G) + G$

Evolution of structure discovery

Grosse, Salakhutdinov, Freeman, and Tenenbaum, UAI 2012

# Talk Roadmap

DBM

## Part 2: Advanced Hierarchical Models

- Introduction: Transfer Learning/ One-Shot Learning.

- Compound Hierarchical Deep Models:

  - Deep Boltzmann Machines.

  - Hierarchical Latent Dirichlet Allocation Model.

- Applications.

- MCMC techniques.

"animal"

"vehicle"

horse   cow   car   van   truck

# Inference



**Gibbs Sampler**



**Problem:** When dealing with complex high-dimensional data: the probability landscape is highly multimodal.

Inability to efficiently explore a distribution with many isolated modes.

Problem for both directed and undirected graphical models.

- Posterior distribution: $P(\theta|\mathcal{D}) = \frac{1}{P(\mathcal{D})}P(\mathcal{D}|\theta)P(\theta)$

- Boltzmann machine: $P(z) = \frac{1}{Z}\exp(-E(z))$

# Tempered Transitions

(Radford Neal, 1994)

Define a sequence of intermediate probability distributions $p_0, \ldots, p_S$ where:

- $p_S = p(\mathbf{x}; \theta)$ is the original complicated distribution.
- $p_0$ is more spread out and easier to sample from.

One way is to define:

$$p_s(\mathbf{x}) \; \propto \; p^*(\mathbf{x}; \theta)^{\beta_s},$$

where "inverse temperatures" $\beta_0 < \beta_1 < \ldots < \beta_S = 1$ are chosen by the user.

$\beta = 0 \qquad \beta = 0.01 \qquad \beta = 0.1 \qquad \beta = 0.25 \qquad \beta = 0.5 \qquad \beta = 1$

For each $s = 1, .., S - 1$ we define a transition operator $T_s(\mathbf{x}' \leftarrow \mathbf{x})$ that leaves $p_s$ invariant.

# Tempered Transitions

Define reverse transition operator $p_s(\mathbf{x})T_s(\mathbf{x}' \leftarrow \mathbf{x}) = \tilde{T}_s(\mathbf{x} \leftarrow \mathbf{x}')p_s(\mathbf{x}')$.



- Given a current state, apply a sequence of transition operators:

$$T_{S-1} \ldots T_0 \tilde{T}_0 \ldots \tilde{T}_{S-1}.$$

- Systematically "move" the sample from the complicated distribution to the easily sampled distribution and back.

- Accept a new state $\tilde{\mathbf{x}}^S$ with probability:

$$\min \left[ 1, \prod_{s=1}^{S} p^*(\mathbf{x}_s)^{\beta_{s-1} - \beta_s} p^*(\tilde{\mathbf{x}}_s)^{\beta_s - \beta_{s-1}} \right].$$

# Learning MRFs using Tempered Transitions

Training data

Samples with
Tempered Transitions

Samples without
Tempered Transitions



Plain stochastic approximation using simple Gibbs works badly.

A large fraction of the model's probability mass is placed on images of humans.

# Simulated Tempering (ST)

- Simulated tempering: Sample from the joint distribution:

$$p(\mathbf{x}, k) \propto w_k \exp(-\beta_k E(\mathbf{x})),$$

where $w_k$ are pre-specified constants, and $0 < \beta_K < \beta_{K-1} < ... < \beta_1 = 1$ represent the K "inverse temperatures".

# Simulated Tempering (ST)

- Simulated tempering: Sample from the joint distribution:

$$p(\mathbf{x}, k) \propto w_k \exp(-\beta_k E(\mathbf{x})),$$

where w$_k$ are pre-specified constants, and $0 < \beta_K < \beta_{K-1} < ... < \beta_1 = 1$ represent the K "inverse temperatures".

- The main problem of ST:

$$p(k) \propto \sum_{\mathbf{x}} w_k \exp(-\beta_k E(\mathbf{x})) = w_k \mathcal{Z}_k$$

- To be efficient, it is important for the Markov chain to spend roughly equal amount of time at each temperature level.

- Hence w$_k$ needs to be proportional to $1/\mathcal{Z}_k$.

# Adaptive Simulated Tempering (AST)

- Partitioning the state space into K sets $\{k\} \cup \mathcal{X},$ each corresponding to a different temperature value.

- If the move into a different partition (temperature) is rejected:

  - The adaptive weight $g_k$ for the current partition k will increase.

  - This will (exponentially) increase the probability of accepting the next move into a different temperature level.



Atchade and Liu, 2004,
Famong Liang, 2005

# Adaptive Simulated Tempering (AST)

- Given k$^t$, sample k$^{t+1}$ from proposal distribution: $q(k^{t+1} \leftarrow k^t)$

  Accept with probability:

  $$\min\left(1, \underbrace{\frac{p(\mathbf{x}^t, k^{t+1})q(k^t \leftarrow k^{t+1})}{p(\mathbf{x}^t, k^t)q(k^{t+1} \leftarrow k^t)}}_{\text{Standard M-H update}} \times \underbrace{\frac{g_{k^t}}{g_{k^{t+1}}}}_{\text{Adaptive factor}}\right)$$

- Update adaptive weights:

  $$g_i^{t+1} = g_i^t(1 + \gamma_t \mathbb{I}(k^{t+1} \in \{i\})), \; i = 1, ..., K.$$

- It can be verified: $g_i^t / g_j^t \longrightarrow \mathcal{Z}_i / \mathcal{Z}_j$ as $\gamma_t \longrightarrow 0$.



Atchade and Liu, 2004,
Famong Liang, 2005

# Fast-Slow AST

- When using AST for learning, it is hard to balance between:

  - Exploration: waiting until adaptive ST escapes from the local mode.
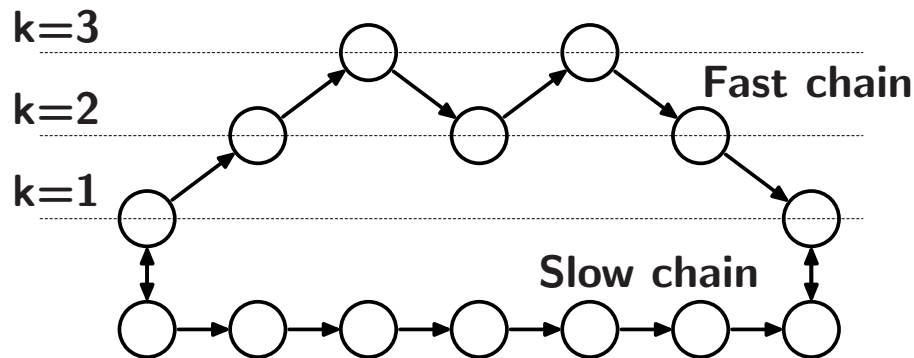  - Exploitation: learning model parameters.

# Fast-Slow AST

- When using AST for learning, it is hard to balance between:

  - Exploration: waiting until adaptive ST escapes from the local mode.
  - Exploitation: learning model parameters.

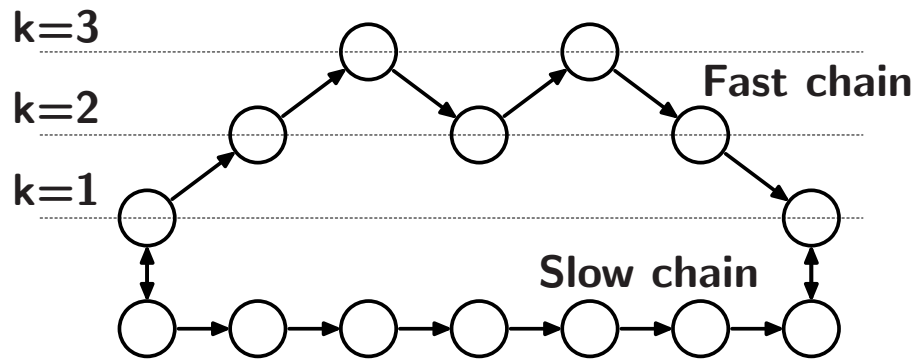- Consider two chains, sampling from the same target distribution.



Slow chain evolves according to the standard Gibbs updates.

Fast chain uses adaptive ST.

- Parameters are updated based on the slow chain. The role of the fast chain is to explore different modes.

# Fast-Slow AST

k=3

k=2

k=1

**Fast chain**

**Slow chain**

Slow chain evolves according to the standard Gibbs updates.

Fast chain uses adaptive ST.

• The algorithm is only twice as expensive compared to the standard stochastic approximation algorithm.

• Parameters are updated after every Gibbs update, while the fast chain runs in parallel, adaptively mixing between different modes of the energy landscape.

• Unlike fast Persistent Contrastive Divergence (PCD), the fast chain is likely to visit spurious modes that may reside far away from the data.
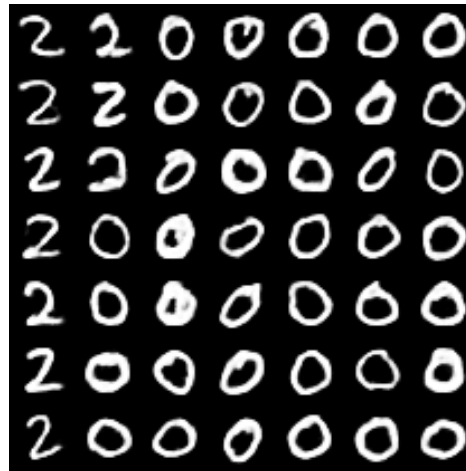
# MNIST Dataset

1000 latents

500 latents

28 x 28 pixel image

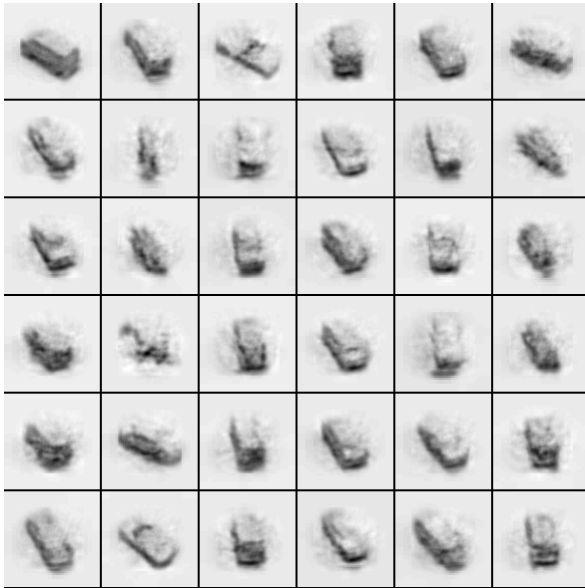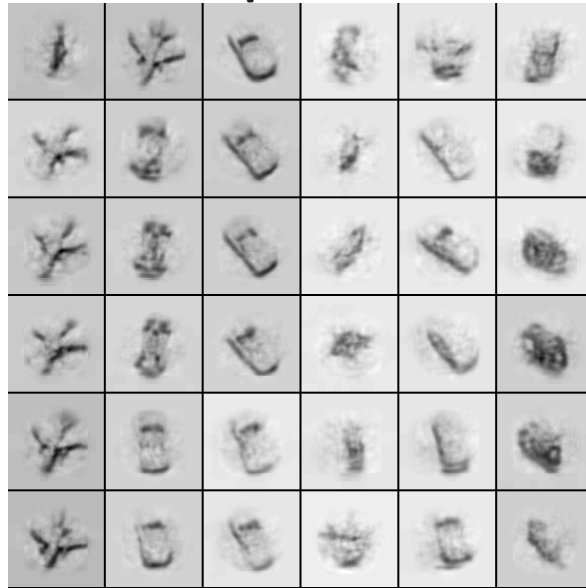About 890,000 parameters

Gibbs

Adaptive ST

- Samples from the two-hidden-layer DBM (1000-500-784) produced by the Gibbs and adaptive ST with 300 Gibbs steps between consecutive images (by column).

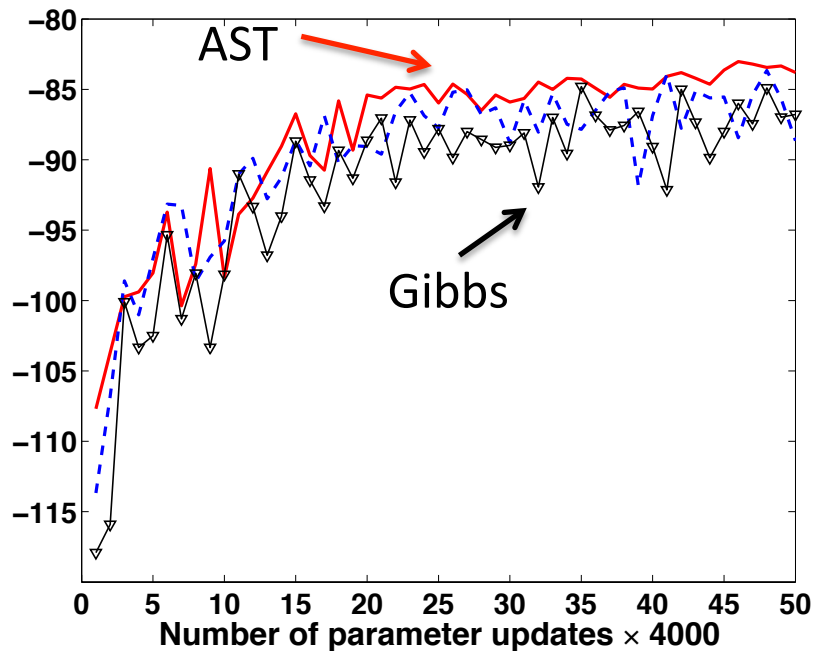# NORB Dataset

Gibbs                          Adaptive ST



• Samples from two-hidden-layer DBM: 4000-4000-(96x96), produced by the Gibbs and fast-slow adaptive ST with 500 Gibbs steps between consecutive images (by column). About 3 million parameters.

# Learning DBMs

The estimates of the average test log-probabilities per image (in nats) for different learning algorithms.



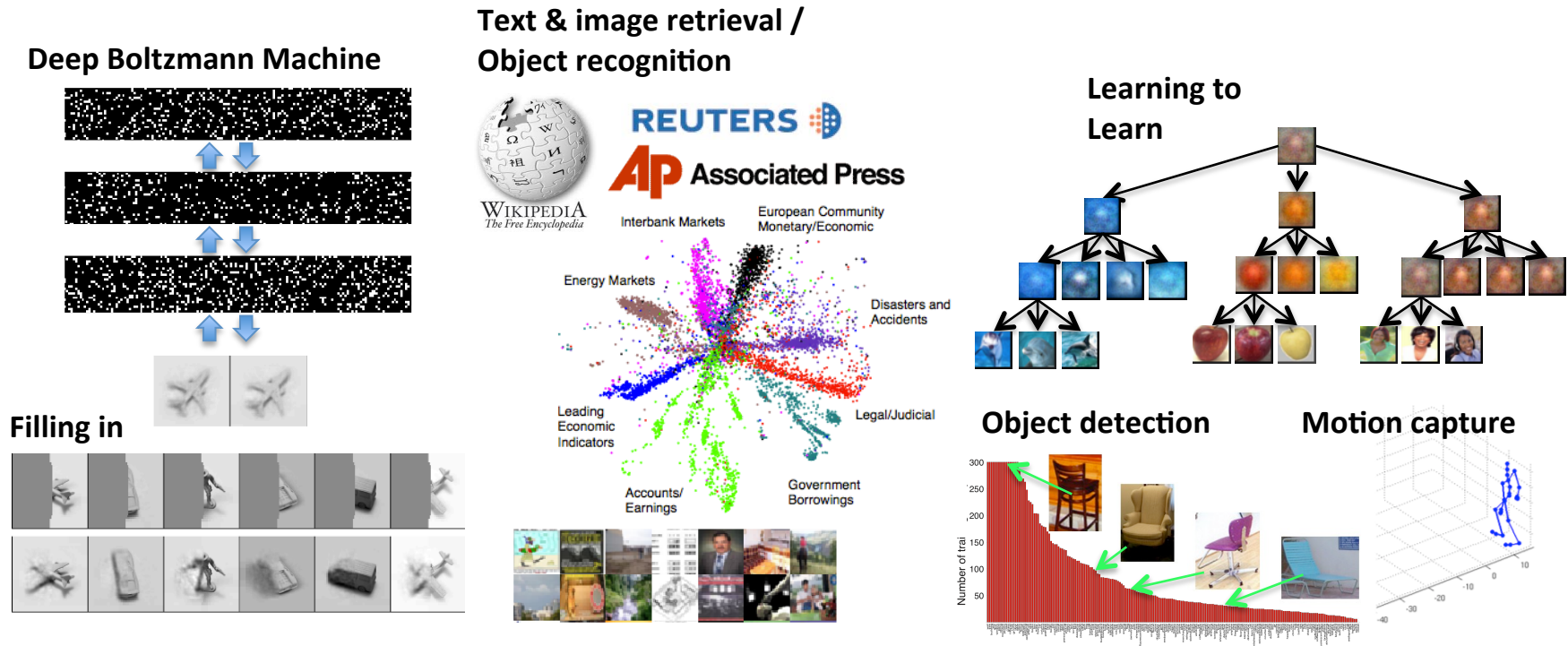| Algorithm | MNIST (+/- 0.5) | NORB (+/- 1.1) |
|---|---|---|
| Gibbs | -87.23 | -596.92 |
| Fast PCD | -86.72 | -597.12 |
| Tempered Transitions | -85.41 | -595.54 |
| Fast-Slow AST | -84.12 | -591.18 |

• Fast-Slow AST  tends to exhibit a more stable behavior during learning.

# Recap

- Efficient learning algorithms for Hierarchical Generative Models.

**Deep Boltzmann Machine**

**Text & image retrieval / Object recognition**

**Learning to Learn**

**Filling in**

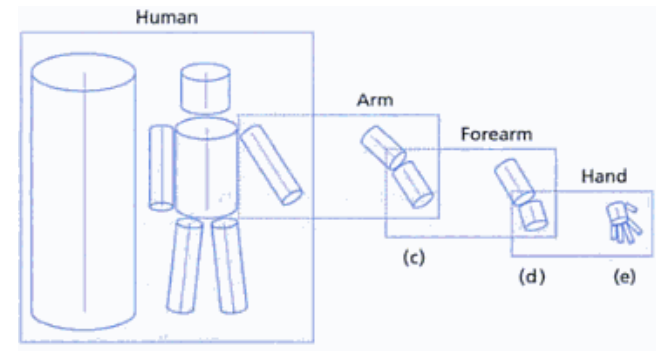**Object detection**

**Motion capture**



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

# Summary

Compose hierarchical Bayesian models with deep networks for transfer learning / one-shot learning.
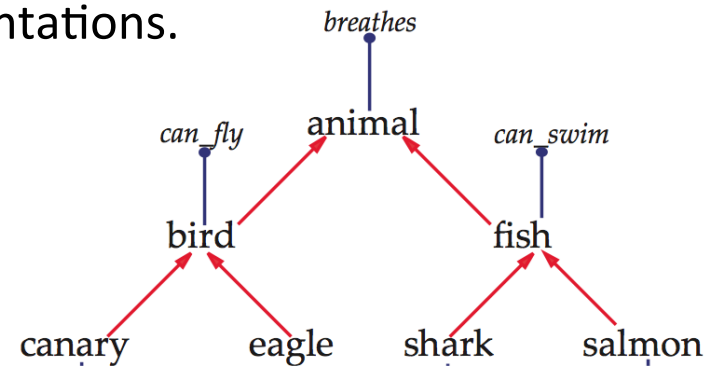
## Deep Networks: Learning Part-based Hierarchy:



- multiple **layers of nonlinearities.**
- **distributed representations.**
- **unsupervised learning of generic features** -- no need to rely on human-crafted input representations.

## Hierarchical Bayes: Learning Category Hierarchy:



- **explicitly learn category hierarchies** for sharing abstract knowledge.
- **modular data-parameter relations**.
- higher-level **class sensitive features.**

# Thank you