

Deep Learning II

Russ Salakhutdinov

Department of Statistics and Computer Science
University of Toronto

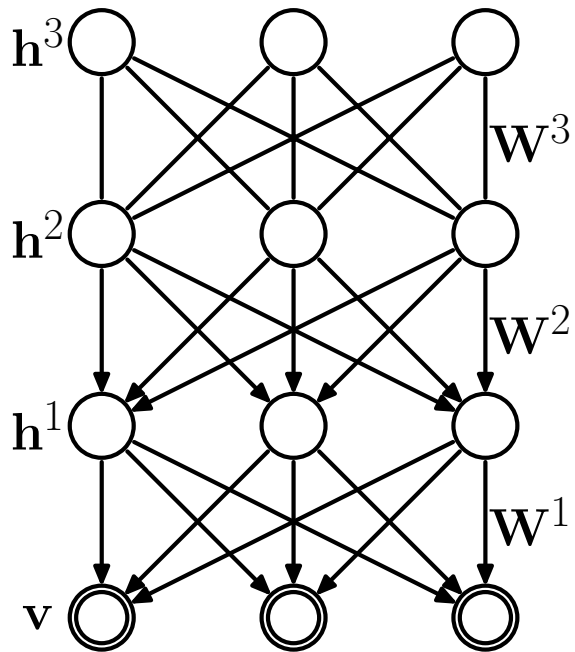
<http://www.utstat.toronto.edu/~rsalakhu/isbi.html>

Talk Roadmap

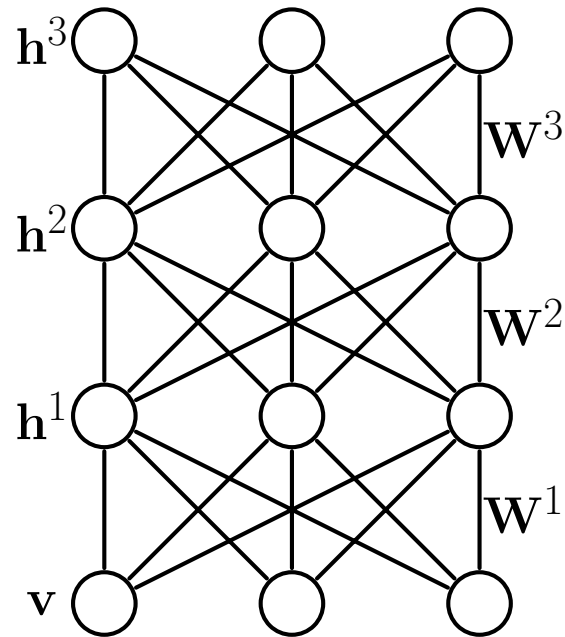
- Advanced Deep Models
 - Deep Boltzmann Machines
 - One-Shot and Transfer Learning
 - Learning Structured and Robust Deep Models
- Multimodal Learning
- Conclusions

DBNs vs. DBMs

Deep Belief Network



Deep Boltzmann Machine



DBNs are hybrid models:

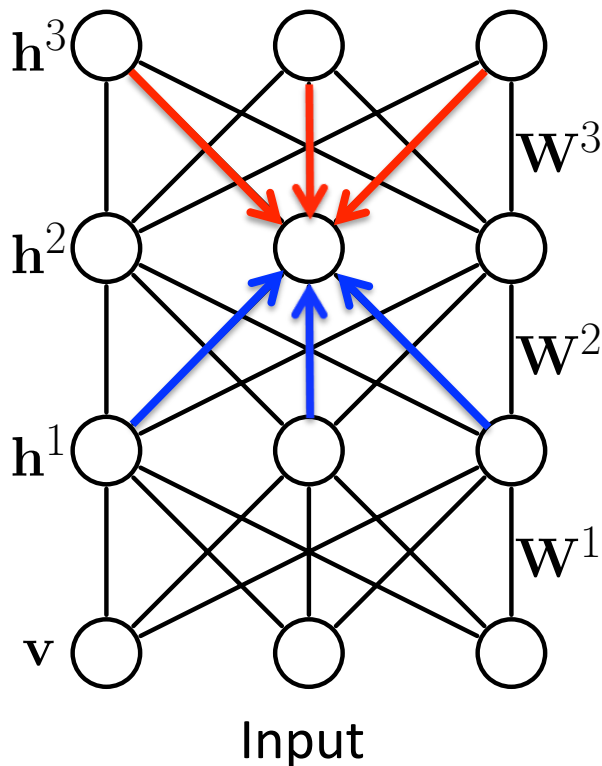
- Inference in DBNs is problematic due to **explaining away**.
- Only greedy pretraining, **no joint optimization over all layers**.
- Approximate inference is feed-forward: **no bottom-up and top-down**.

Mathematical Formulation

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

Deep Boltzmann Machine

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



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

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

Bottom-up

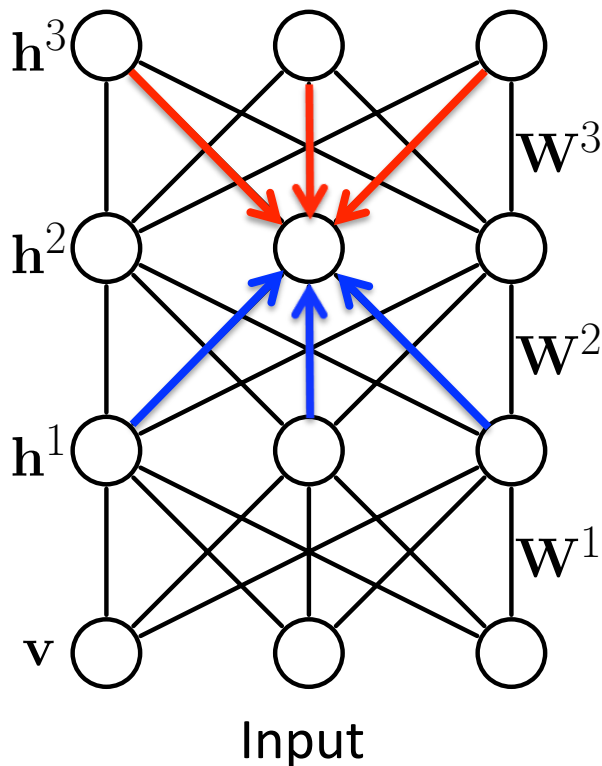
Top-Down

Unlike many existing feed-forward models: ConvNet (LeCun), HMAX (Poggio et.al.), Deep Belief Nets (Hinton et.al.)

Mathematical Formulation

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

Deep Boltzmann Machine



- Conditional Distributions:

$$P(h_j^1 = 1 | \mathbf{v}, \mathbf{h}^2) = \sigma \left(\sum_i W_{ij}^1 v_i + \sum_k W_{jk}^2 h_k^2 \right)$$

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

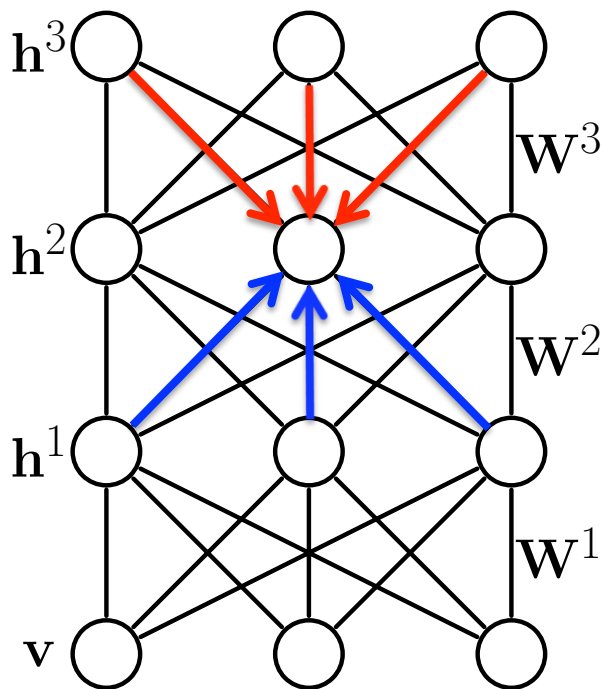
$$P(h_m^3 = 1 | \mathbf{h}^2) = \sigma \left(\sum_k W_{km}^3 h_k^2 \right)$$

- Note that exact computation of $P(\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3 | \mathbf{v})$ is intractable.

Mathematical Formulation

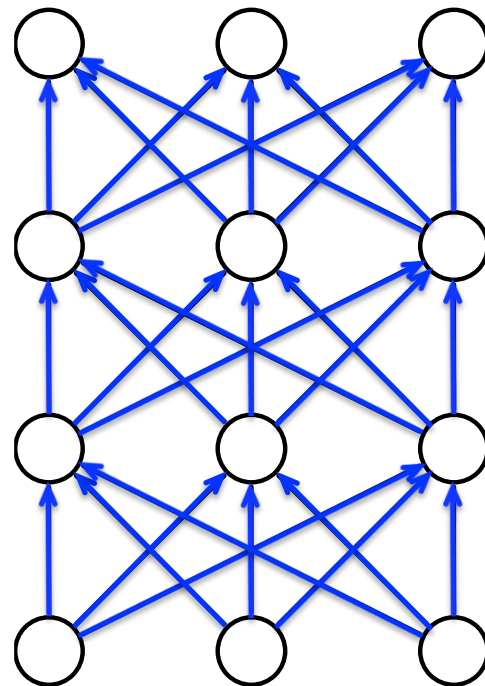
$$P_{\theta}(\mathbf{v}) = \frac{P^*(\mathbf{v})}{Z(\theta)} = \frac{1}{Z(\theta)} \sum_{\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3} \exp \left[\mathbf{v}^{\top} W^1 \mathbf{h}^1 + \mathbf{h}^1{}^{\top} W^2 \mathbf{h}^2 + \mathbf{h}^2{}^{\top} W^3 \mathbf{h}^3 \right]$$

Deep Boltzmann Machine

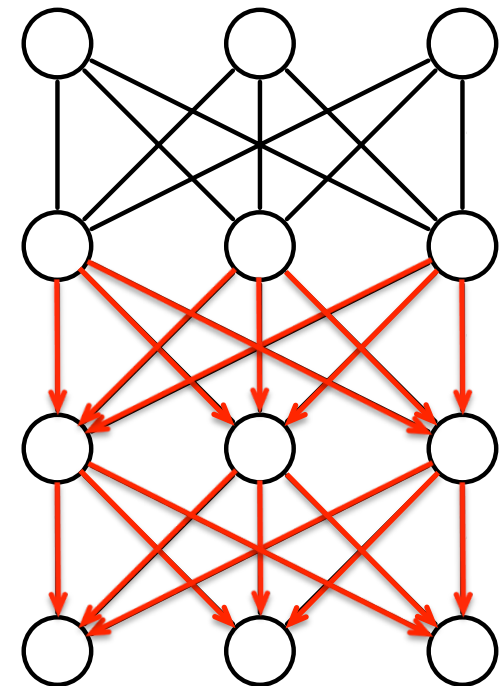


Input

Neural Network
Output



Deep Belief Network

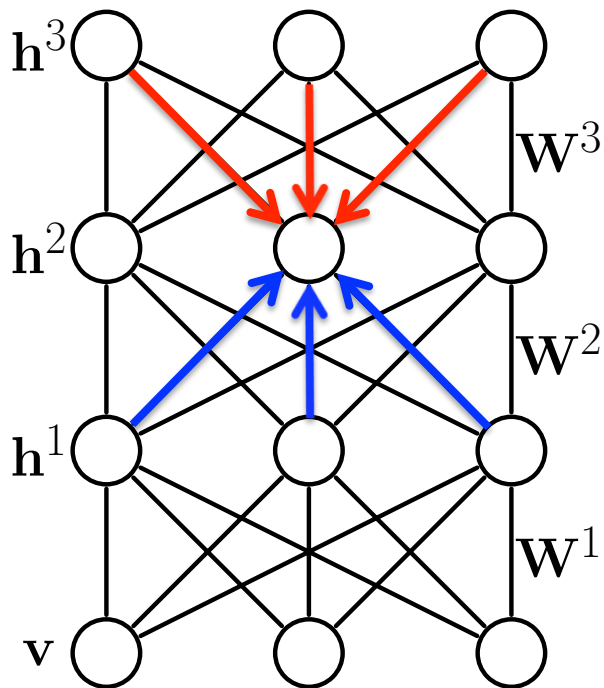


Unlike many existing feed-forward models: ConvNet (LeCun), HMAX (Poggio), Deep Belief Nets (Hinton)

Mathematical Formulation

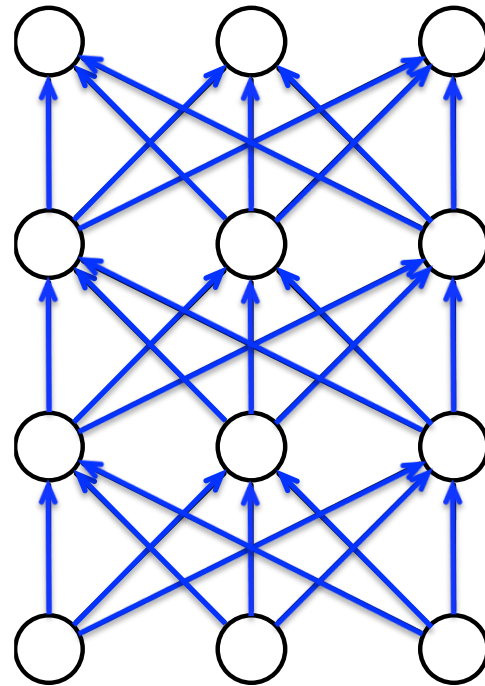
$$P_{\theta}(\mathbf{v}) = \frac{P^*(\mathbf{v})}{Z(\theta)} = \frac{1}{Z(\theta)} \sum_{\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3} \exp \left[\mathbf{v}^{\top} W^1 \mathbf{h}^1 + \mathbf{h}^1{}^{\top} W^2 \mathbf{h}^2 + \mathbf{h}^2{}^{\top} W^3 \mathbf{h}^3 \right]$$

Deep Boltzmann Machine

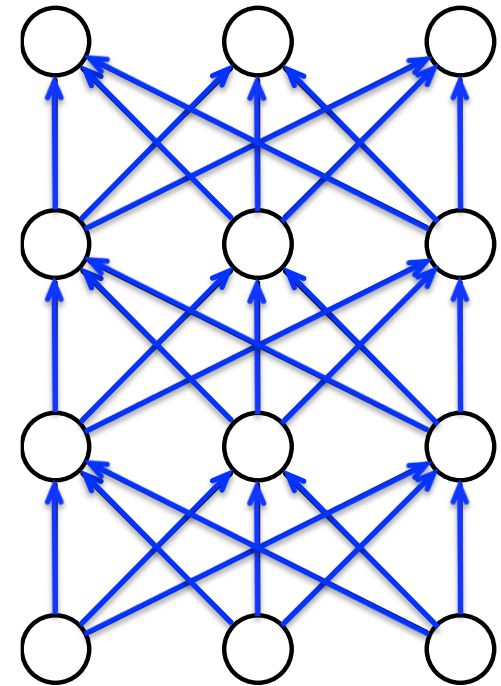


Input

Neural Network
Output



Deep Belief Network



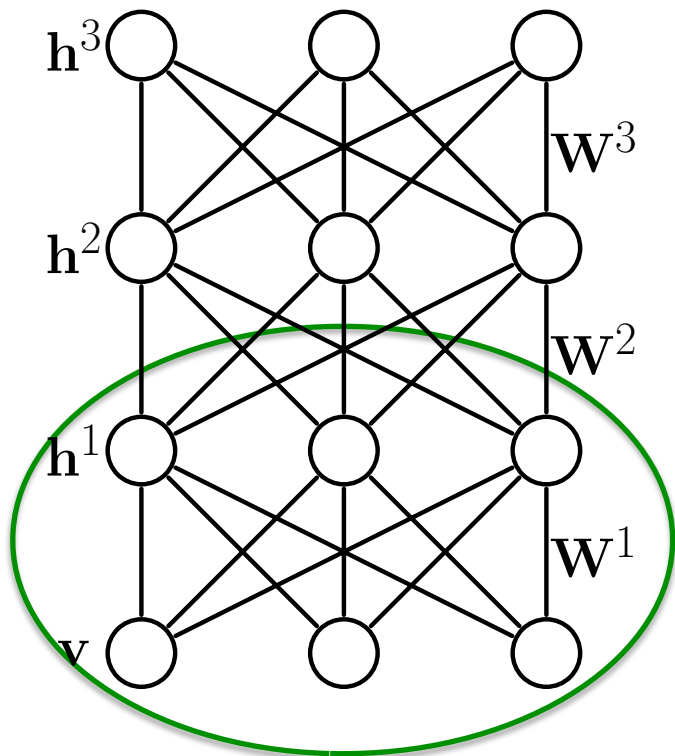
inference

Unlike many existing feed-forward models: ConvNet (LeCun), HMAX (Poggio), Deep Belief Nets (Hinton)

Mathematical Formulation

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

Deep Boltzmann Machine



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

- Dependencies between hidden variables.

Maximum likelihood learning:

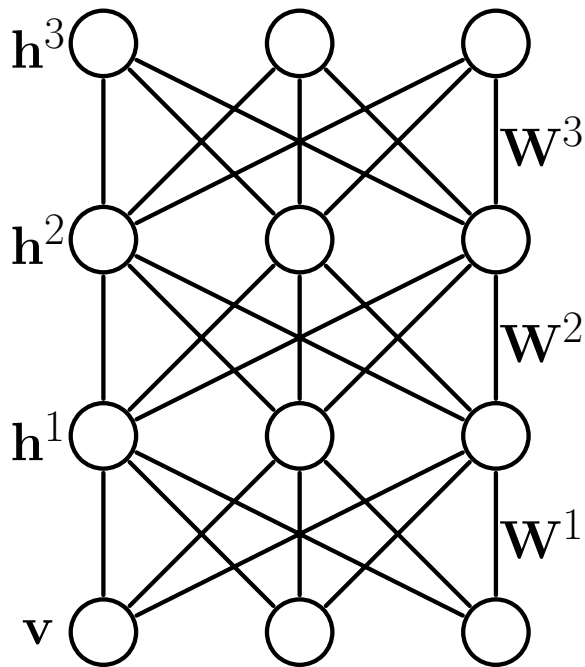
$$\frac{\partial \log P_{\theta}(\mathbf{v})}{\partial W^1} = \mathbb{E}_{P_{data}}[\mathbf{v} \mathbf{h}^1{}^{\top}] - \mathbb{E}_{P_{\theta}}[\mathbf{v} \mathbf{h}^1{}^{\top}]$$

Problem: Both expectations are intractable!

Learning rule for undirected graphical models:
MRFs, CRFs, Factor graphs.

Approximate Learning

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



(Approximate) Maximum Likelihood:

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

- Both expectations are intractable!

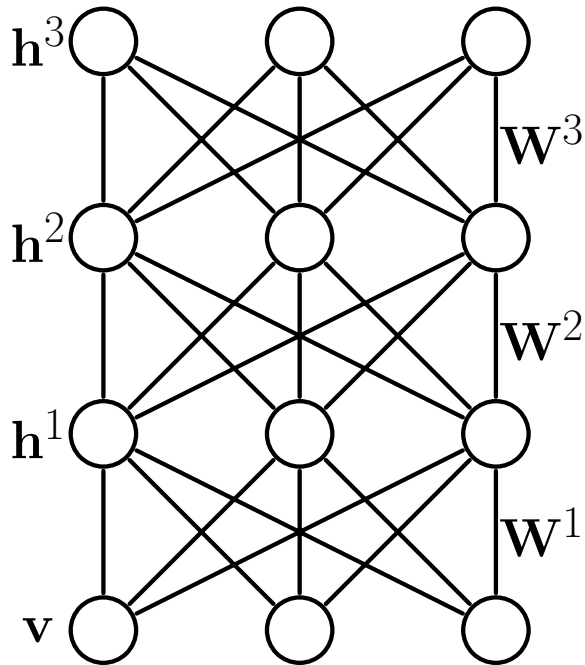
$$P_{data}(\mathbf{v}, \mathbf{h}^1) = P_{\theta}(\mathbf{h}^1 | \mathbf{v}) P_{data}(\mathbf{v})$$

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

Not factorial any more!

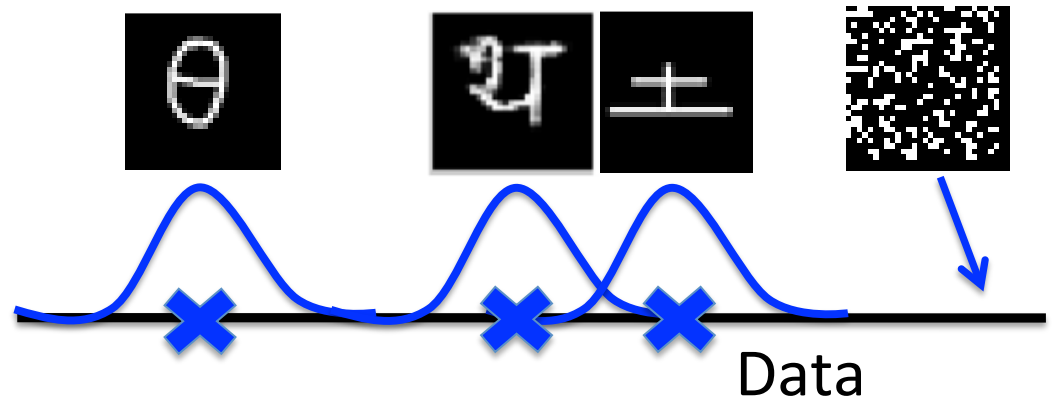
Approximate Learning

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



(Approximate) Maximum Likelihood:

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



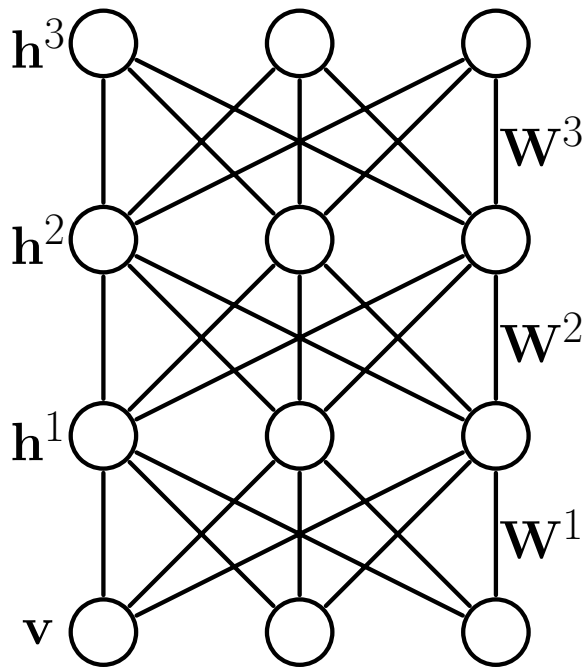
$$P_{data}(\mathbf{v}, \mathbf{h}^1) = P_{\theta}(\mathbf{h}^1 | \mathbf{v}) P_{data}(\mathbf{v})$$

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

Not factorial any more!

Approximate Learning

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



(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!

Previous Work

Many approaches for learning Boltzmann machines have been proposed over the last 20 years:

- Hinton and Sejnowski (1983),
- Peterson and Anderson (1987)
- Galland (1991)
- Kappen and Rodriguez (1998)
- Lawrence, Bishop, and Jordan (1998)
- Tanaka (1998)
- Welling and Hinton (2002)
- Zhu and Liu (2002)
- Welling and Teh (2003)
- Yasuda and Tanaka (2009)

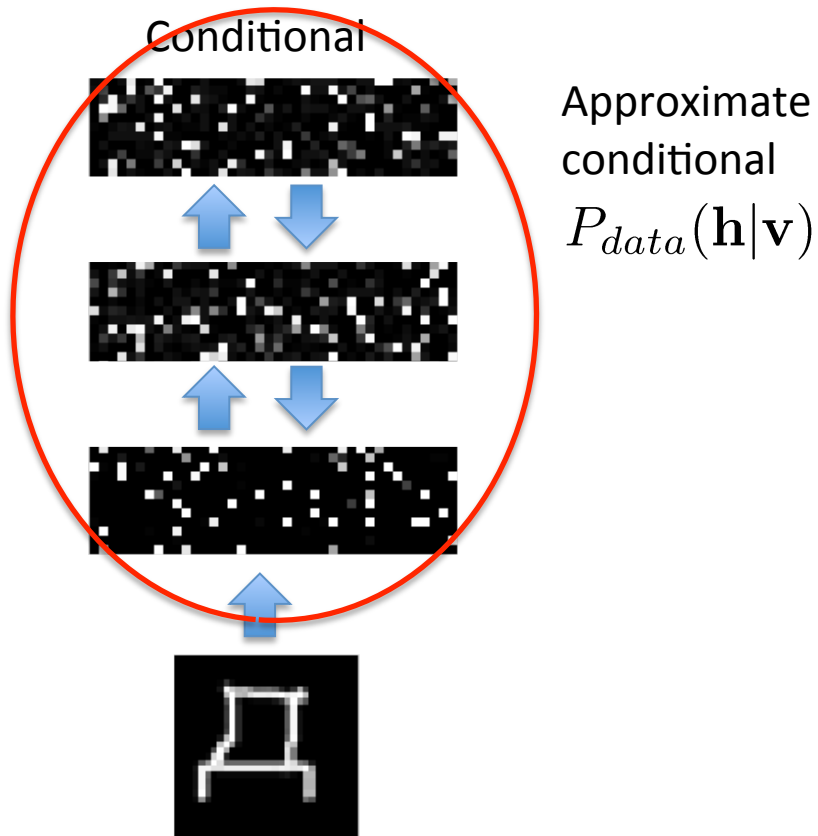
Real-world applications – thousands of hidden and observed variables with millions of parameters.

Many of the previous approaches were not successful for learning general Boltzmann machines with **hidden variables**.

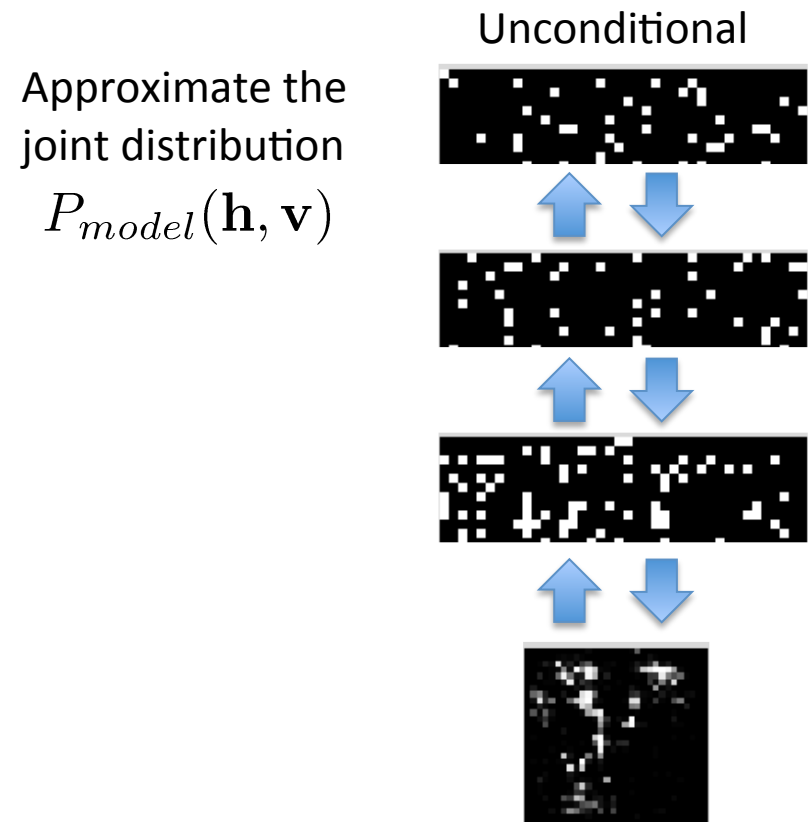
Algorithms based on Contrastive Divergence, Score Matching, Pseudo-Likelihood, Composite Likelihood, MCMC-MLE, Piecewise Learning, cannot handle multiple layers of hidden variables.

New Learning Algorithm

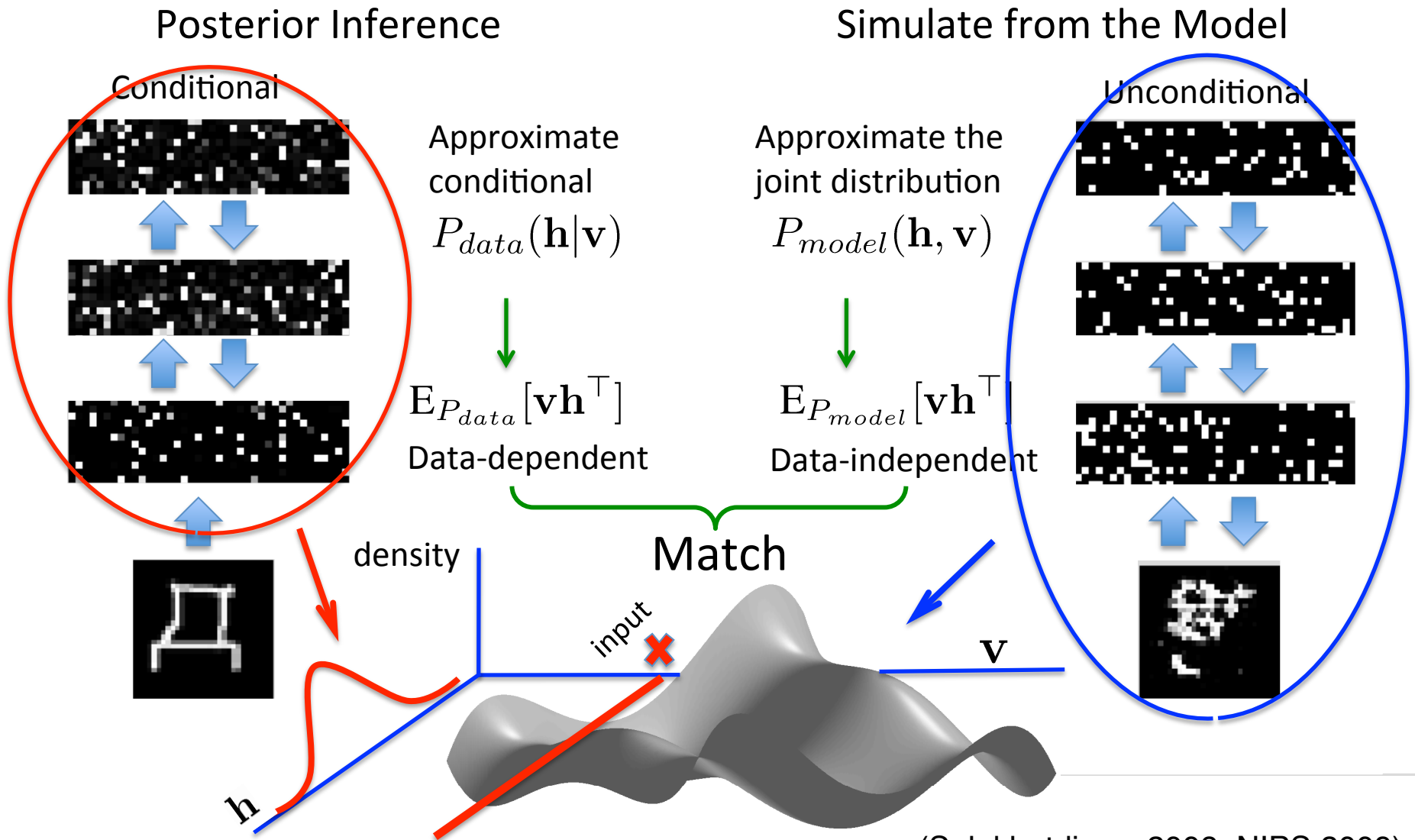
Posterior Inference



Simulate from the Model

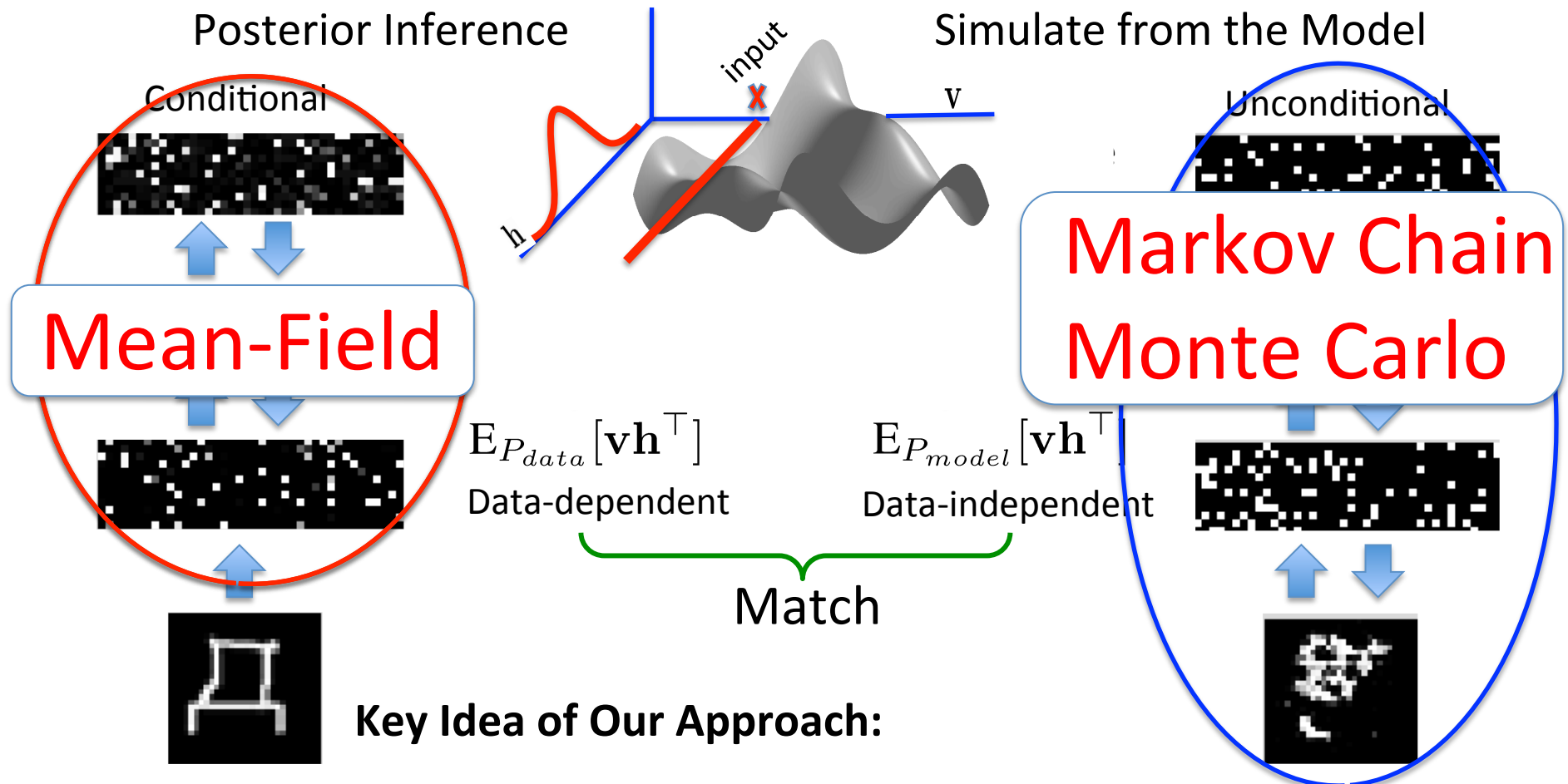


New Learning Algorithm



(Salakhutdinov, 2008; NIPS 2009)

New Learning Algorithm



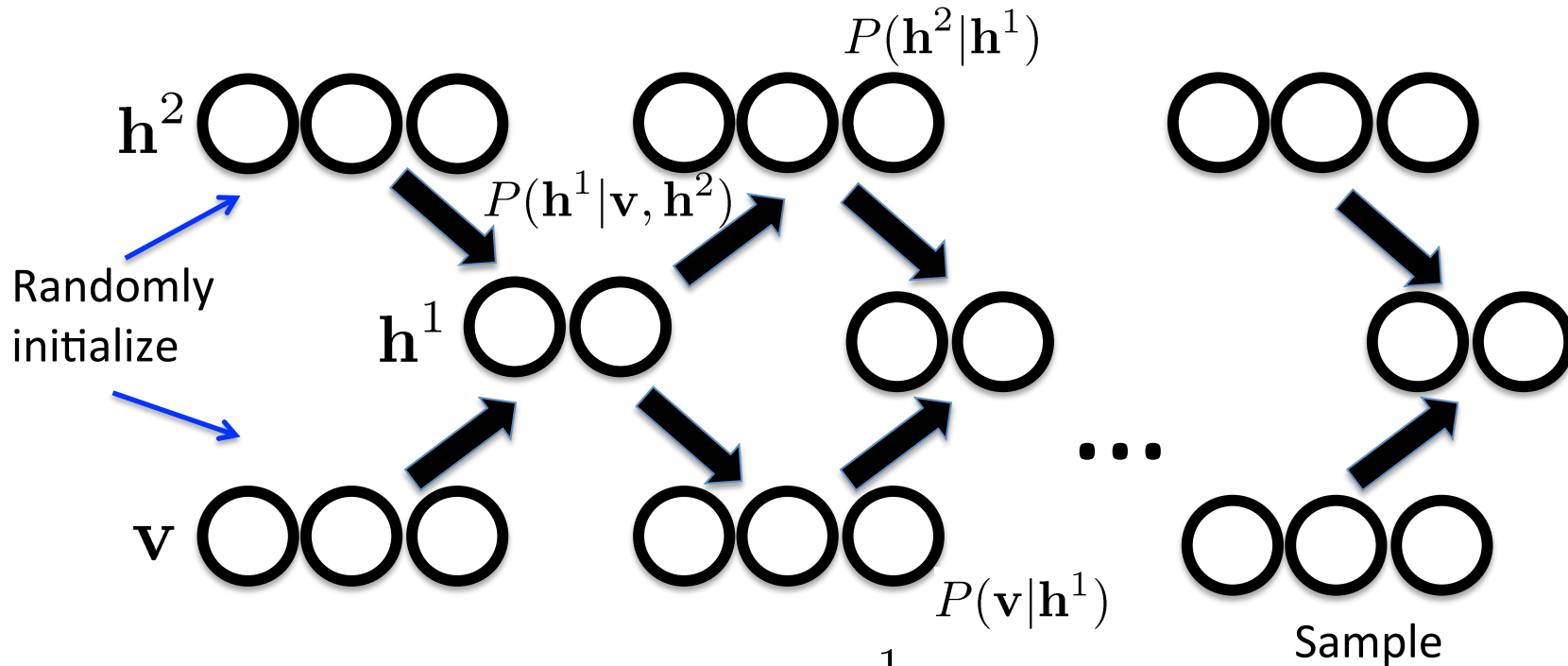
Key Idea of Our Approach:

Data-dependent: **Variational Inference**, mean-field theory

Data-independent: **Stochastic Approximation**, MCMC based

Sampling from DBMs

Sampling from two-hidden layer DBM by running a Markov chain:



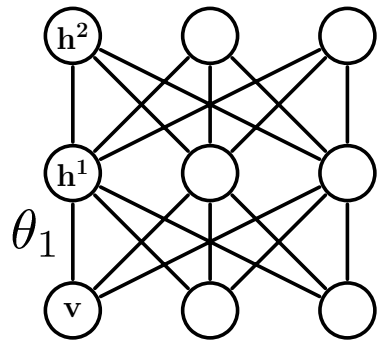
$$P(h_m^1 = 1|\mathbf{v}, \mathbf{h}^2) = \frac{1}{1 + \exp(-\sum_i W_{im}^1 v_i - \sum_j W_{mj}^2 h_j^2)}$$

$$P(h_j^2 = 1|\mathbf{h}^1) = \frac{1}{1 + \exp(-\sum_m W_{mj}^2 h_m^1)}$$

$$P(v_i = 1|\mathbf{h}^1) = \frac{1}{1 + \exp(-\sum_m W_{im}^1 h_m^1)}$$

Stochastic Approximation

Time t=1

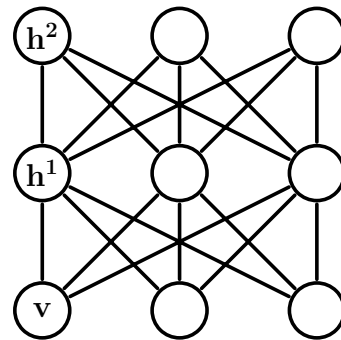


$$\mathbf{x}_1 \sim T_{\theta_1}(\mathbf{x}_1 \leftarrow \mathbf{x}_0)$$

Update θ_1



t=2

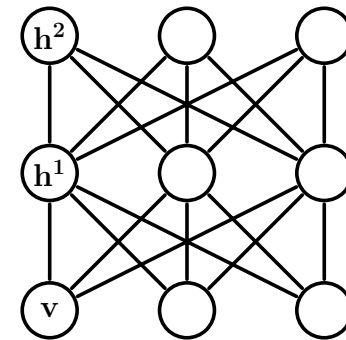


$$\mathbf{x}_2 \sim T_{\theta_2}(\mathbf{x}_2 \leftarrow \mathbf{x}_1)$$

Update θ_2



t=3



$$\mathbf{x}_3 \sim T_{\theta_3}(\mathbf{x}_3 \leftarrow \mathbf{x}_2)$$

Update θ_t and \mathbf{x}_t sequentially, where $\mathbf{x} = \{\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2\}$

- Generate $\mathbf{x}_t \sim T_{\theta_t}(\mathbf{x}_t \leftarrow \mathbf{x}_{t-1})$ by simulating from a Markov chain that leaves P_{θ_t} invariant (e.g. Gibbs or M-H sampler)
- Update θ_t by replacing intractable $E_{P_{\theta_t}}[\mathbf{v}\mathbf{h}^\top]$ with a point estimate $[\mathbf{v}_t\mathbf{h}_t^\top]$

In practice we simulate several Markov chains in parallel.

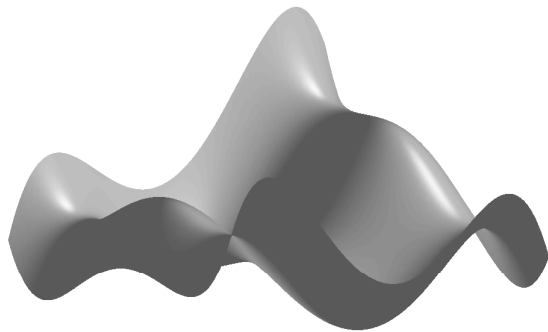
Robbins and Monro, Ann. Math. Stats, 1957
L. Younes, Probability Theory 1989

Stochastic Approximation

Update rule decomposes:

$$\theta_{t+1} = \theta_t + \underbrace{\alpha_t \left(\mathbb{E}_{P_{data}}[\mathbf{v}\mathbf{h}^\top] - \mathbb{E}_{P_{\theta_t}}[\mathbf{v}\mathbf{h}^\top] \right)}_{\text{True gradient}} + \underbrace{\alpha_t \left(\mathbb{E}_{P_{\theta_t}}[\mathbf{v}\mathbf{h}^\top] - \sum_{m=1}^M \mathbf{v}_t^{(m)} \mathbf{h}_t^{(m)\top} \right)}_{\text{Noise term } \epsilon_t}$$

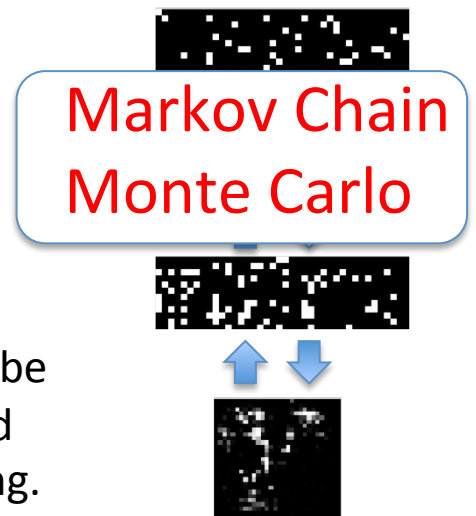
Almost sure convergence guarantees as learning rate $\alpha_t \rightarrow 0$



Salakhutdinov,
ICML 2010

Problem: High-dimensional data:
the energy landscape is highly
multimodal

Key insight: The transition operator can be
any valid transition operator – Tempered
Transitions, Parallel/Simulated Tempering.



Connections to the theory of stochastic approximation and adaptive MCMC.

Variational Inference

Approximate intractable distribution $P_\theta(\mathbf{h}|\mathbf{v})$ with simpler, tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$:

$$\log P_\theta(\mathbf{v}) = \log \sum_{\mathbf{h}} P_\theta(\mathbf{h}, \mathbf{v}) = \log \sum_{\mathbf{h}} Q_\mu(\mathbf{h}|\mathbf{v}) \frac{P_\theta(\mathbf{h}, \mathbf{v})}{Q_\mu(\mathbf{h}|\mathbf{v})}$$

$$\geq \sum_{\mathbf{h}} Q_\mu(\mathbf{h}|\mathbf{v}) \log \frac{P_\theta(\mathbf{h}, \mathbf{v})}{Q_\mu(\mathbf{h}|\mathbf{v})}$$

$$= \sum_{\mathbf{h}} Q_\mu(\mathbf{h}|\mathbf{v}) \log P_\theta^*(\mathbf{h}, \mathbf{v}) - \log \mathcal{Z}(\theta) + \sum_{\mathbf{h}} Q_\mu(\mathbf{h}|\mathbf{v}) \log \frac{1}{Q_\mu(\mathbf{h}|\mathbf{v})}$$

$$\underbrace{\mathbf{v}^\top W^1 \mathbf{h}^1 + \mathbf{h}^1^\top W^2 \mathbf{h}^2 + \mathbf{h}^2^\top W^3 \mathbf{h}^3}_{\text{Variational Lower Bound}}$$

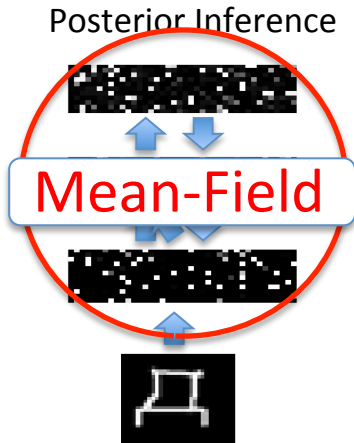
Variational Lower Bound

$$= \log P_\theta(\mathbf{v}) - \text{KL}(Q_\mu(\mathbf{h}|\mathbf{v}) || P_\theta(\mathbf{h}|\mathbf{v}))$$

$$\text{KL}(Q||P) = \int Q(x) \log \frac{Q(x)}{P(x)} dx$$

Minimize KL between approximating and true distributions with respect to variational parameters μ .

(Salakhutdinov, 2008; Salakhutdinov & Larochelle, AI & Statistics 2010)



Variational Inference

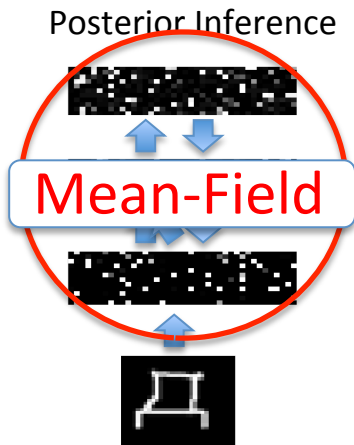
Approximate intractable distribution $P_\theta(\mathbf{h}|\mathbf{v})$ with simpler, tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$:

$$\text{KL}(Q||P) = \int Q(x) \log \frac{Q(x)}{P(x)} dx$$

$$\log P_\theta(\mathbf{v}) \geq \log P_\theta(\mathbf{v}) - \text{KL}(Q_\mu(\mathbf{h}|\mathbf{v})||P_\theta(\mathbf{h}|\mathbf{v}))$$



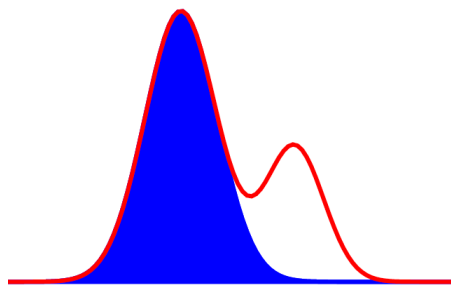
Variational Lower Bound



Mean-Field: Choose a fully factorized distribution:

$$Q_\mu(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^F q(h_j|\mathbf{v}) \text{ with } q(h_j = 1|\mathbf{v}) = \mu_j$$

Variational Inference: Maximize the lower bound w.r.t. Variational parameters μ .



Nonlinear fixed-point equations:

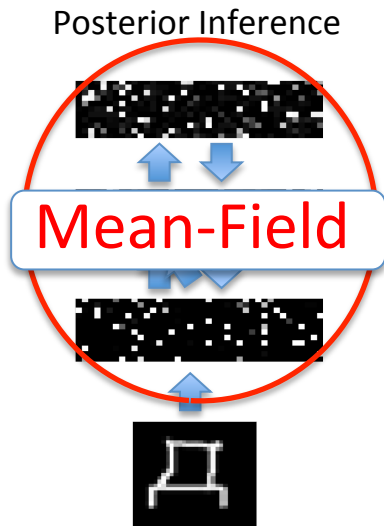
$$\begin{aligned} \mu_j^{(1)} &= \sigma \left(\sum_i W_{ij}^1 v_i + \sum_k W_{jk}^2 \mu_k^{(2)} \right) \\ \mu_k^{(2)} &= \sigma \left(\sum_j W_{jk}^2 \mu_j^{(1)} + \sum_m W_{km}^3 \mu_m^{(3)} \right) \\ \mu_m^{(3)} &= \sigma \left(\sum_k W_{km}^3 \mu_k^{(2)} \right) \end{aligned}$$

Variational Inference

Approximate intractable distribution $P_\theta(\mathbf{h}|\mathbf{v})$ with simpler, tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$:

$$\text{KL}(Q||P) = \int Q(x) \log \frac{Q(x)}{P(x)} dx$$

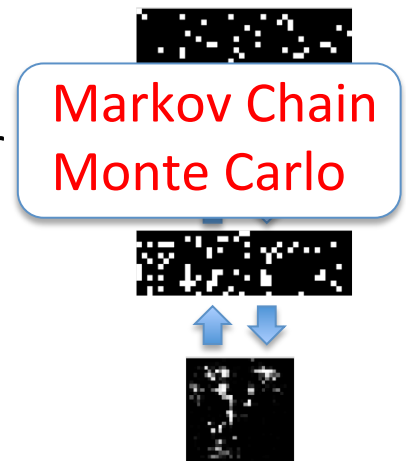
$$\log P_\theta(\mathbf{v}) \geq \log P_\theta(\mathbf{v}) - \underbrace{\text{KL}(Q_\mu(\mathbf{h}|\mathbf{v})||P_\theta(\mathbf{h}|\mathbf{v}))}_{\text{Variational Lower Bound}}$$



Variational Lower Bound

1. **Variational Inference:** Maximize the lower bound w.r.t. variational parameters
2. **MCMC:** Apply stochastic approximation to update model parameters

Unconditional Simulation



Almost sure convergence guarantees to an asymptotically stable point.

Variational Inference

Approximate intractable distribution $P_\theta(\mathbf{h}|\mathbf{v})$ with simpler, tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$:

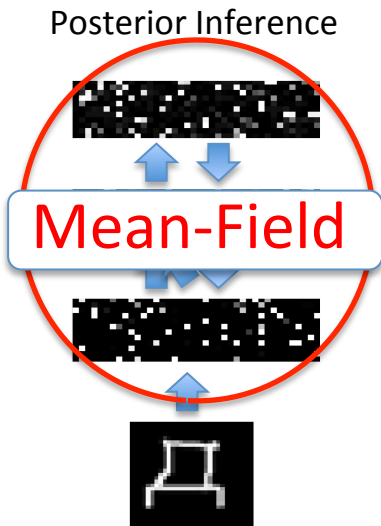
$$\text{KL}(Q||P) = \int Q(x) \log \frac{Q(x)}{P(x)} dx$$

$$\log P_\theta(\mathbf{v}) \geq \log P_\theta(\mathbf{v}) - \text{KL}(Q_\mu(\mathbf{h}|\mathbf{v})||P_\theta(\mathbf{h}|\mathbf{v}))$$



Variational Lower Bound

Unconditional Simulation



1. v
bou

Fast Inference

wer

Markov Chain
Monte Carlo

2. M
to u

Learning can scale to
millions of examples



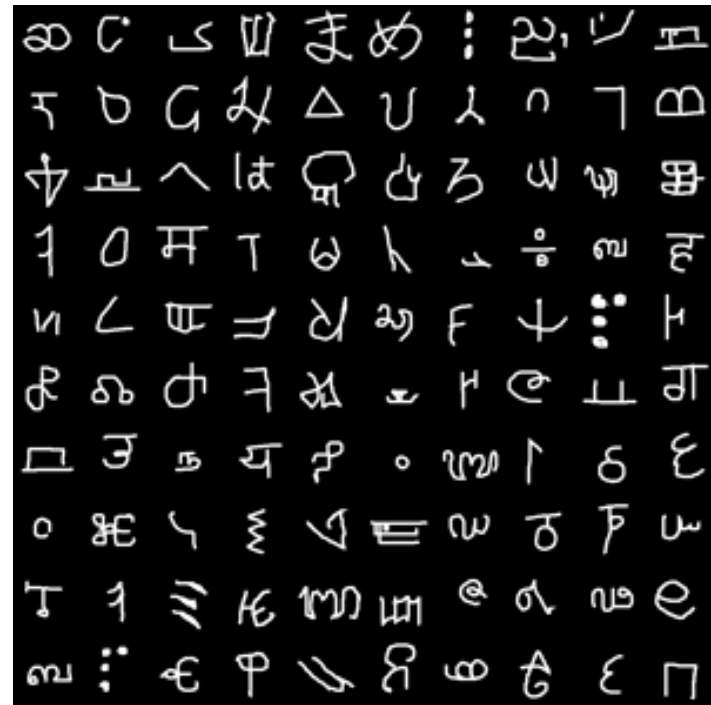
Almost sure convergence guarantees to an asymptotically stable point.

Good Generative Model?

Handwritten Characters

Good Generative Model?

Handwritten Characters



Good Generative Model?

Handwritten Characters

Simulated

Real Data

Good Generative Model?

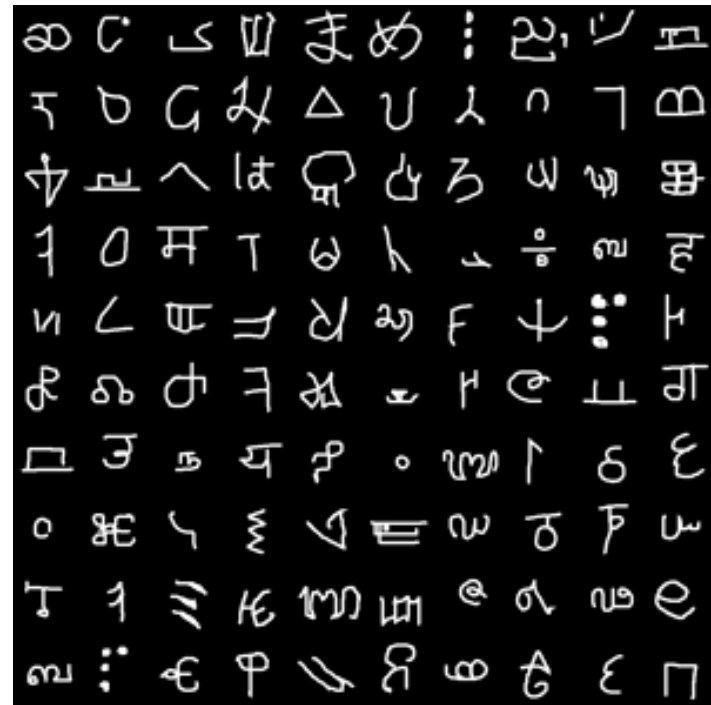
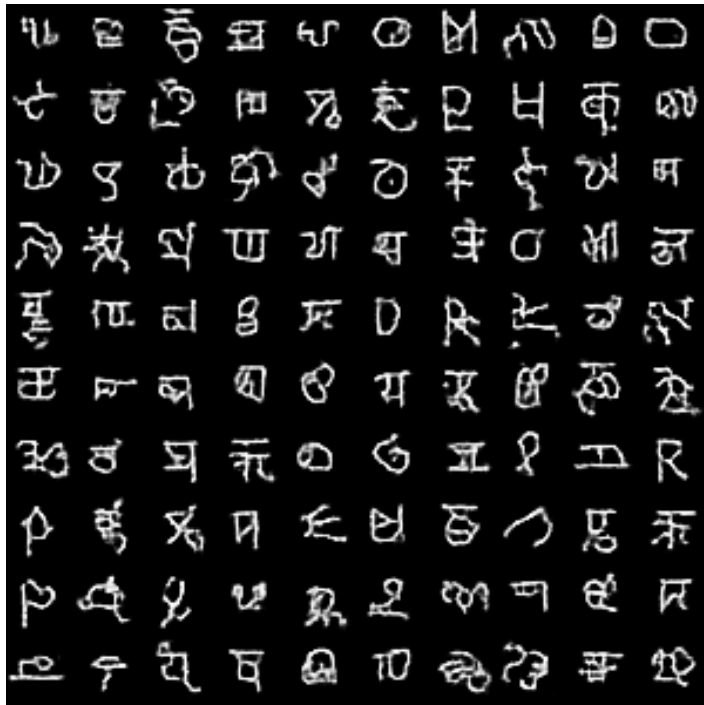
Handwritten Characters

Real Data

Simulated

Good Generative Model?

Handwritten Characters



Good Generative Model?

MNIST Handwritten Digit Dataset



Handwriting Recognition

MNIST Dataset
60,000 examples of 10 digits

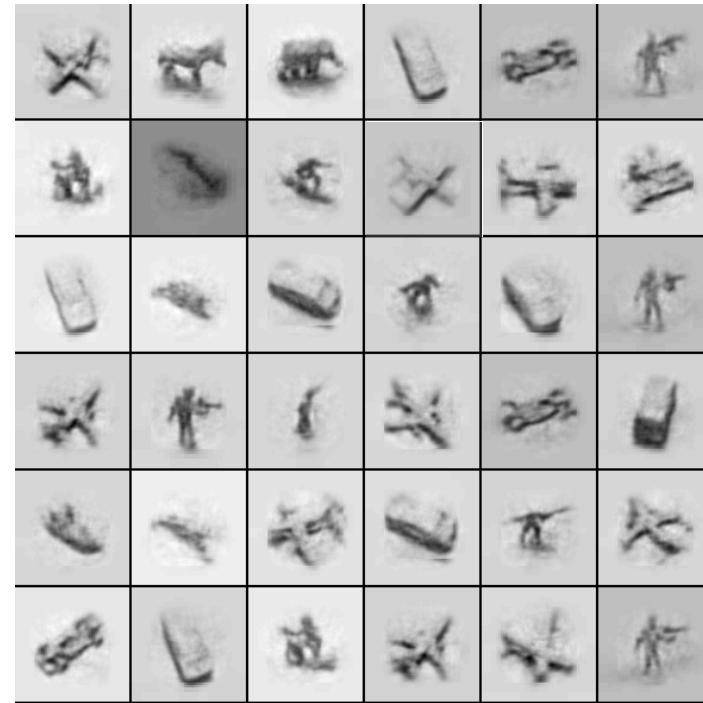
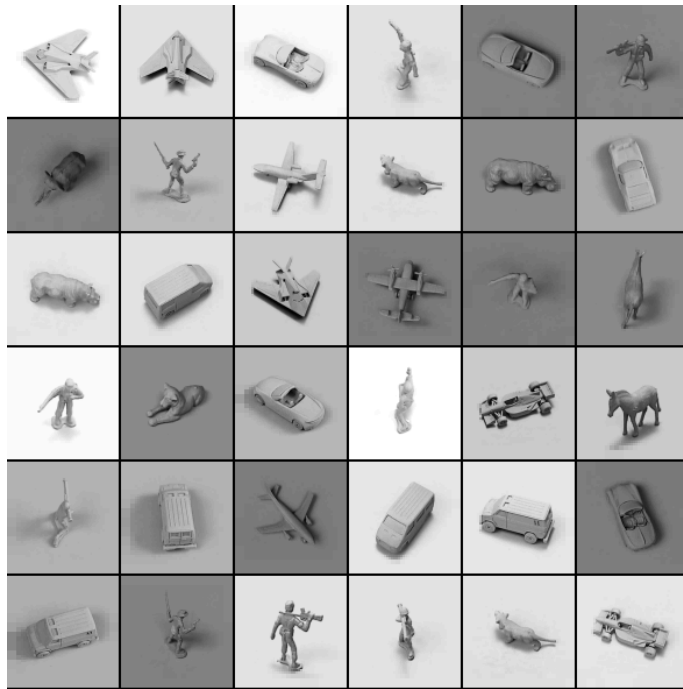
Learning Algorithm	Error
Logistic regression	12.0%
K-NN	3.09%
Neural Net (Platt 2005)	1.53%
SVM (Decoste et.al. 2002)	1.40%
Deep Autoencoder (Bengio et. al. 2007)	1.40%
Deep Belief Net (Hinton et. al. 2006)	1.20%
DBM	0.95%

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%
DBM	8.40%

Permutation-invariant version.

Generative Model of 3-D Objects

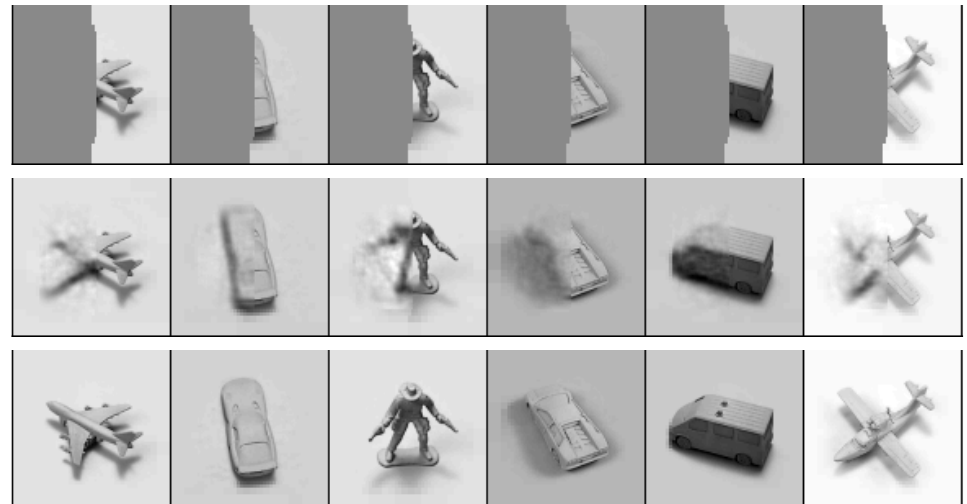


24,000 examples, 5 object categories, 5 different objects within each category, 6 lightning conditions, 9 elevations, 18 azimuths.

3-D Object Recognition

Pattern Completion

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%
DBM	7.2%

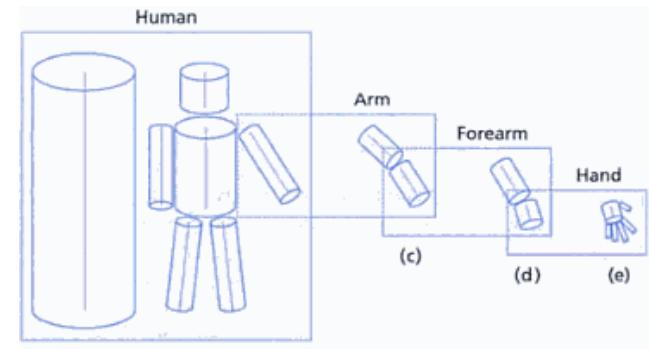


Permutation-invariant version.

Learning Hierarchical Representations

Deep Boltzmann Machines:

Learning Hierarchical Structure
in Features: edges, combination
of edges.



- Performs well in many application domains
- Fast Inference: fraction of a second
- Learning scales to millions of examples

Learning Hierarchical Representations

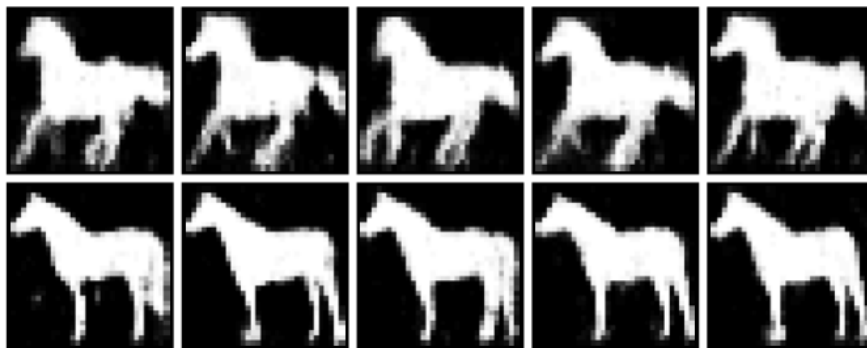
Deep Boltzmann Machines:

Learning H
in Features
of edges.

Need more structured
and robust models

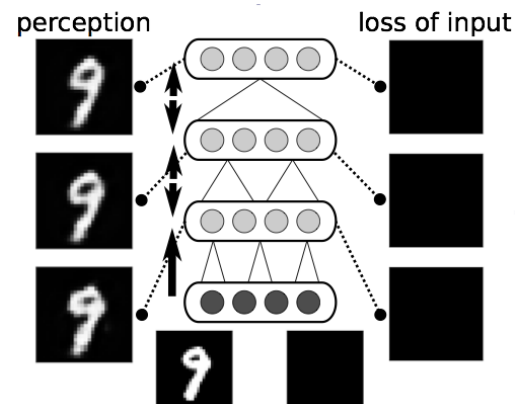


The Shape Boltzmann Machine: a Strong Model of Object Shape
(Eslami, Heess, Winn, CVPR 2012).



Demo DBM

Hallucinations in Charles Bonnet Syndrome Induced by Homeostasis: a Deep Boltzmann Machine Model
(Reichert, Series, Storkey, NIPS 2012)



Talk Roadmap

- Advanced Deep Models
 - Deep Boltzmann Machines
 - **One-Shot and Transfer Learning**
 - Learning Structured and Robust Deep Models
- Multimodal Learning
- Conclusions

One-shot Learning

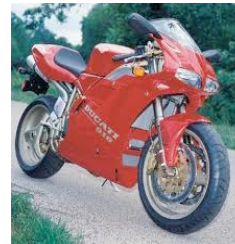


How can we learn a novel concept – a high dimensional statistical object – from few examples.

Supervised Learning



Segway



Motorcycle

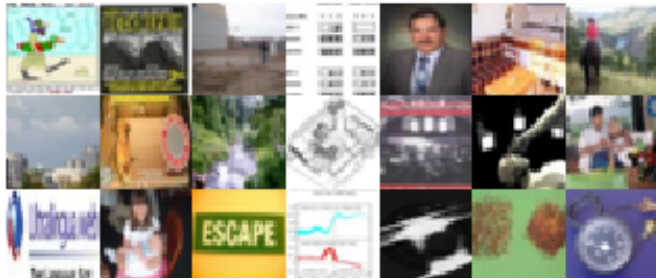
Test:



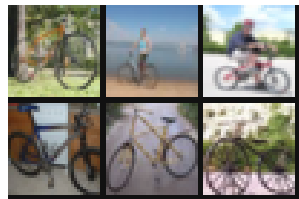
Learning to Learn

Background Knowledge

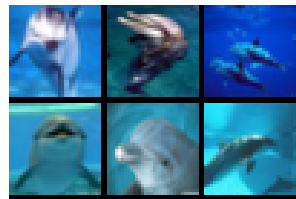
Millions of unlabeled images



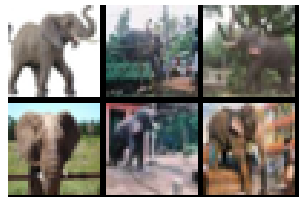
Some labeled images



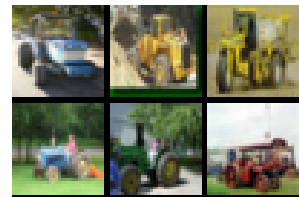
Bicycle



Dolphin



Elephant



Tractor

Learn to Transfer Knowledge



Learn novel concept from one example

Test:



Learning to Learn

Background Knowledge

Millions of unlabeled images

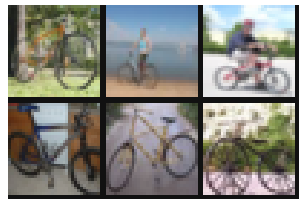


Learn to Transfer Knowledge

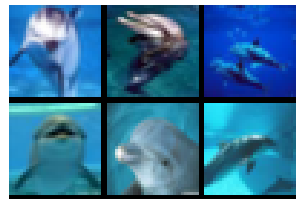
Key problem in computer vision, speech perception, natural language processing, and many other domains.



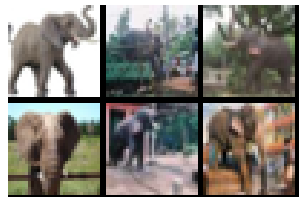
Some labeled images



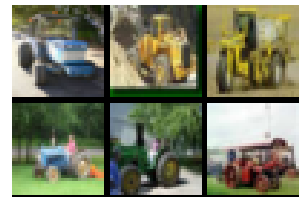
Bicycle



Dolphin



Elephant



Tractor

Learn novel concept from one example

Test:



Hierarchical-Deep Model

HD Models: Integrate hierarchical Bayesian models with deep models.

Hierarchical Bayes:

- Learn **hierarchies of categories** for sharing abstract knowledge.

Deep Models:

- Learn **hierarchies of features**.
- **Unsupervised feature learning** – no need to rely on human-crafted input features.

One-Shot Learning



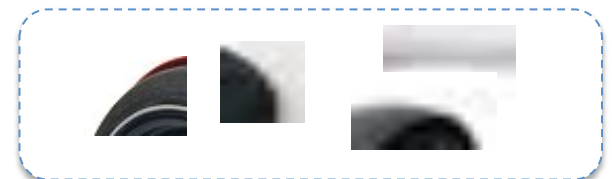
Super-category



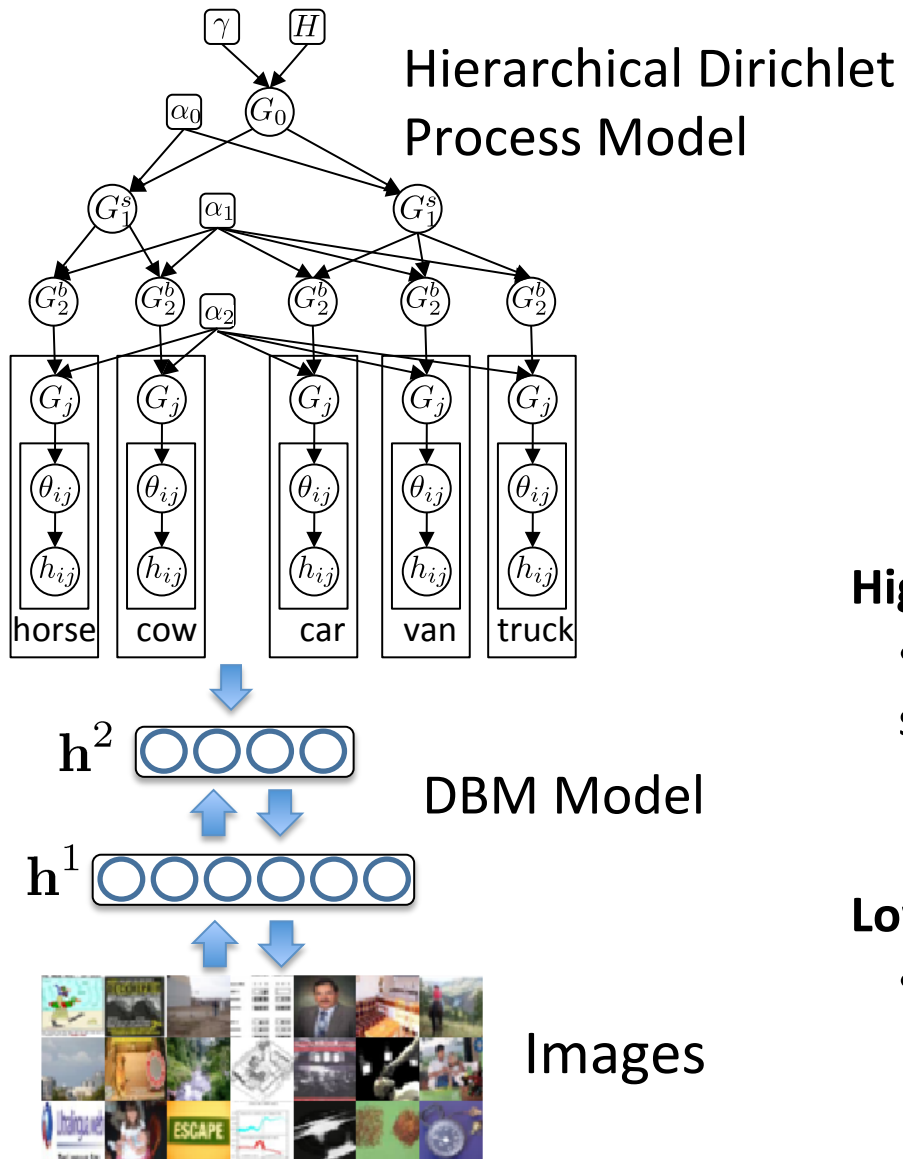
Shared higher-level features



Shared low-level features



Hierarchical-Deep Model



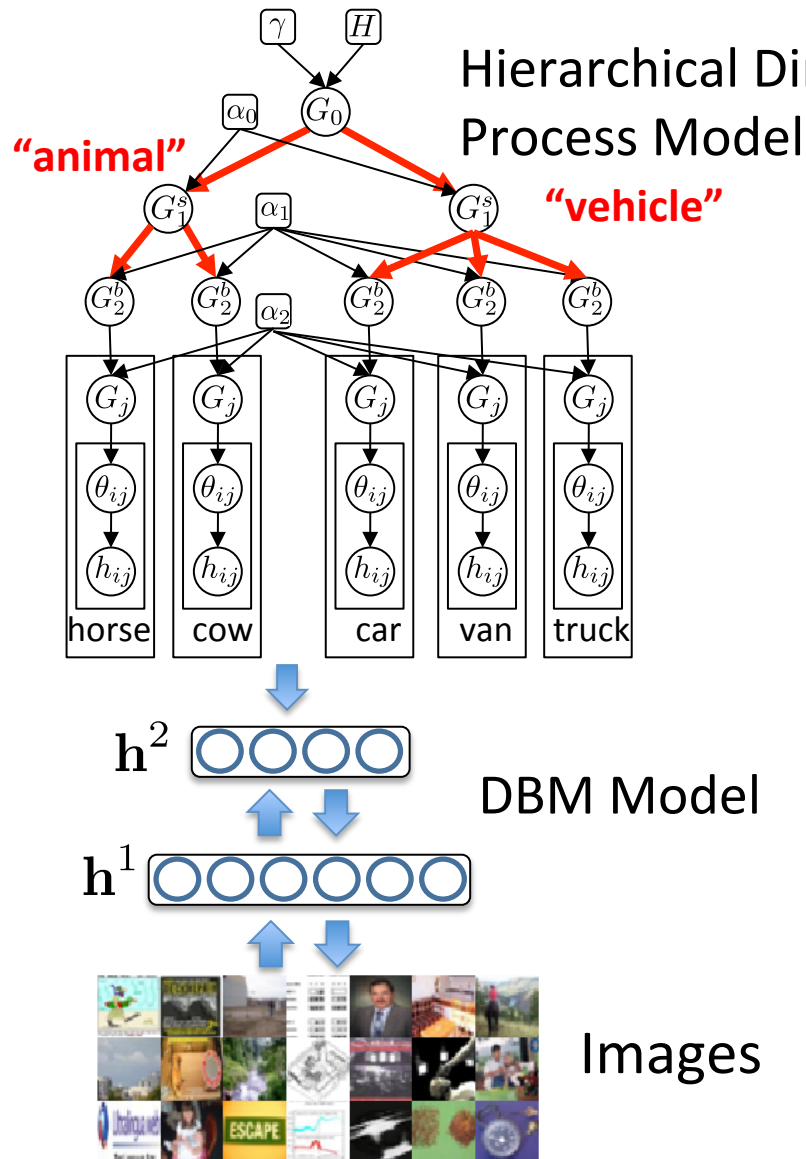
Higher-level class-sensitive features:

- capture distinctive perceptual structure of a specific concept

Lower-level generic features:

- edges, combination of edges

Hierarchical-Deep Model



Hierarchical Organization of Categories:

- express priors on the features that are typical of different kinds of concepts
- modular data-parameter relations

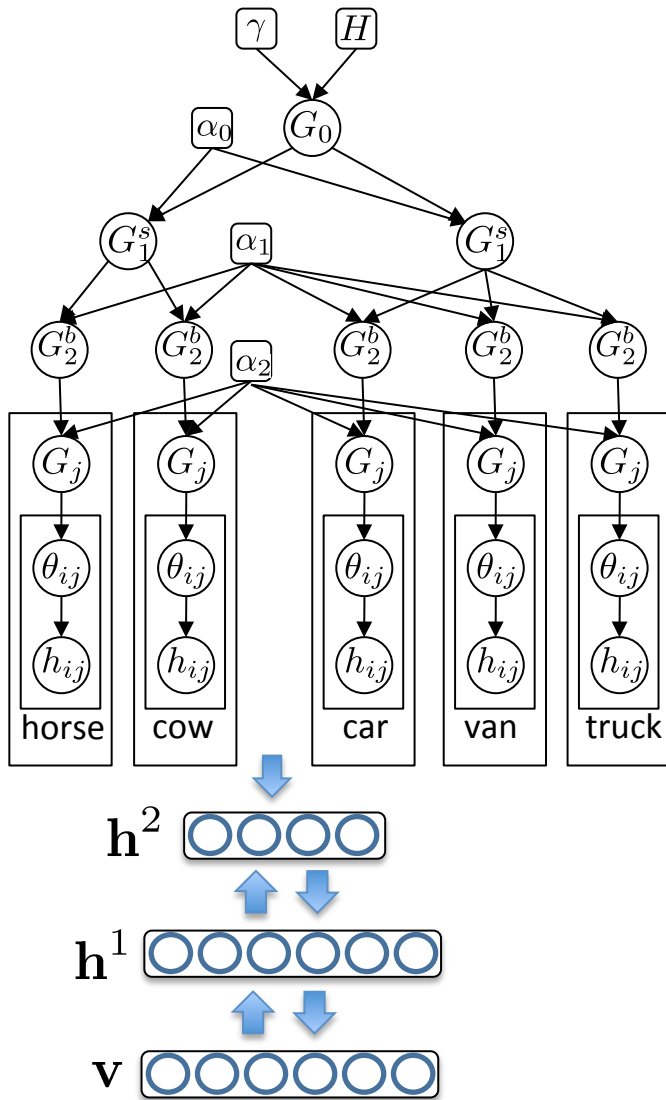
Higher-level class-sensitive features:

- capture distinctive perceptual structure of a specific concept

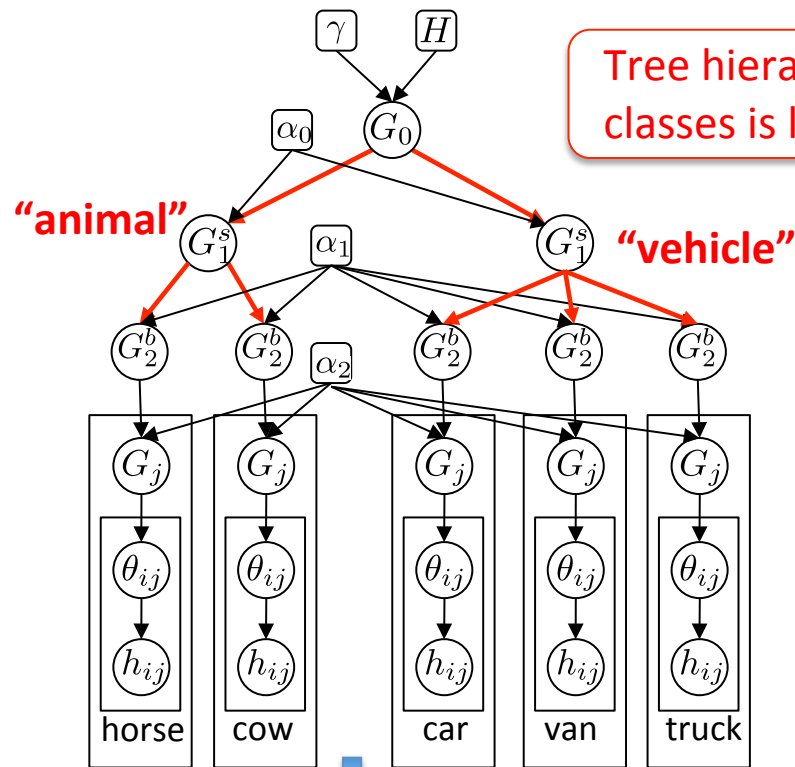
Lower-level generic features:

- edges, combination of edges

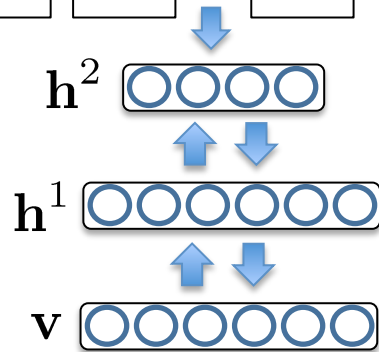
Hierarchical-Deep Model



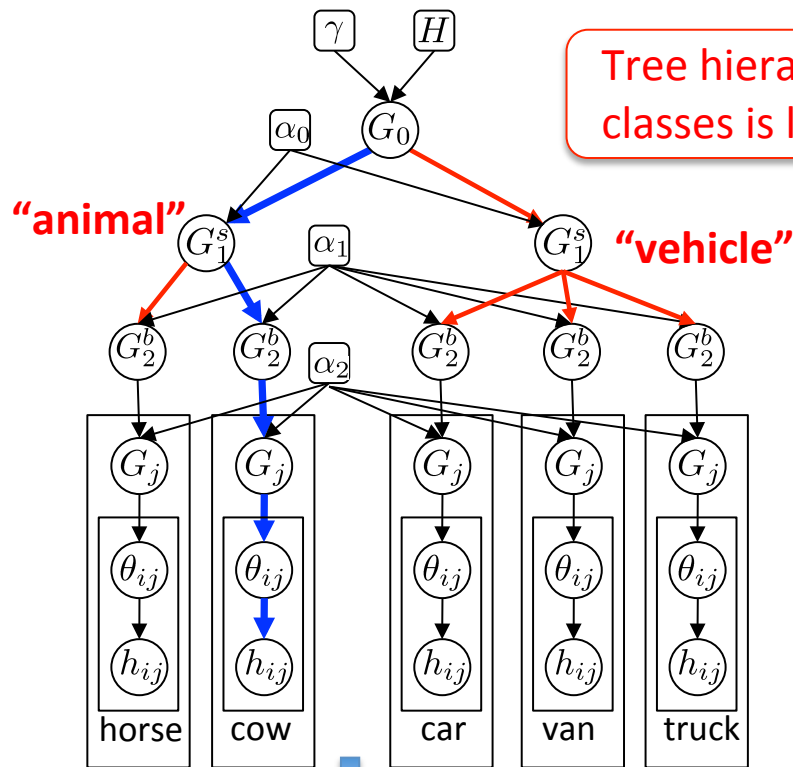
Hierarchical-Deep Model



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

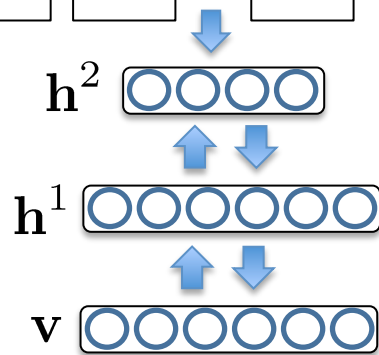


Hierarchical-Deep Model

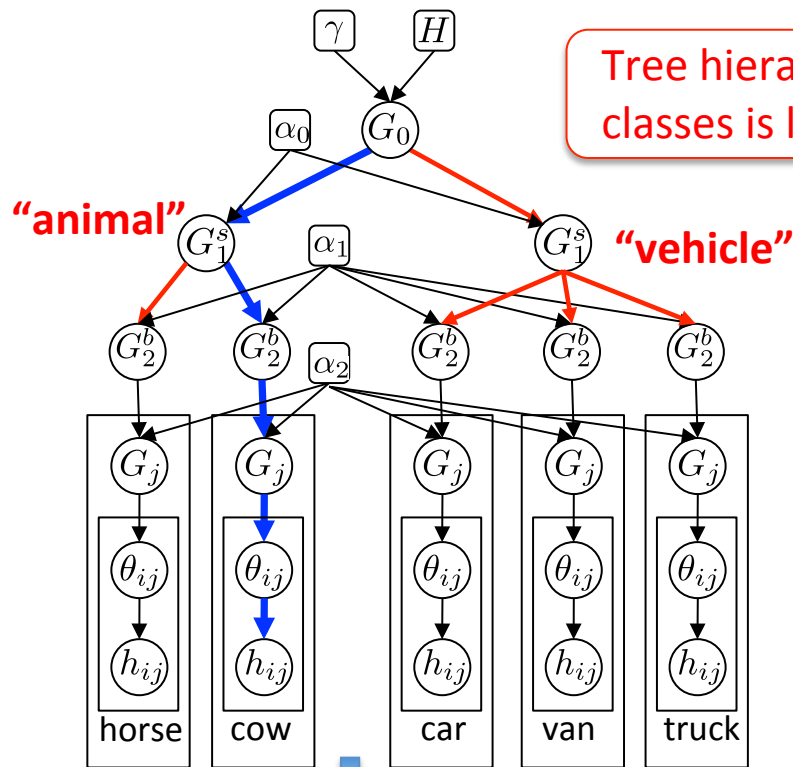


$\mathbf{z} \sim \text{nCRP}$ (**Nested Chinese Restaurant Process**)
 prior: a nonparametric prior over tree structures

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



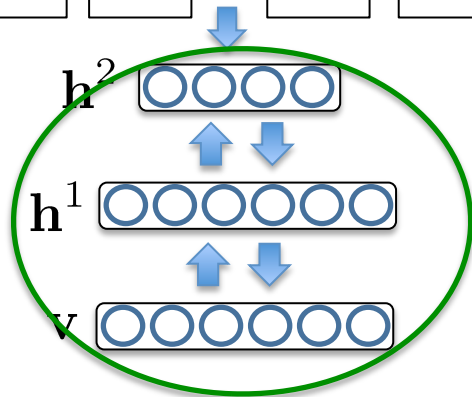
Hierarchical-Deep Model



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

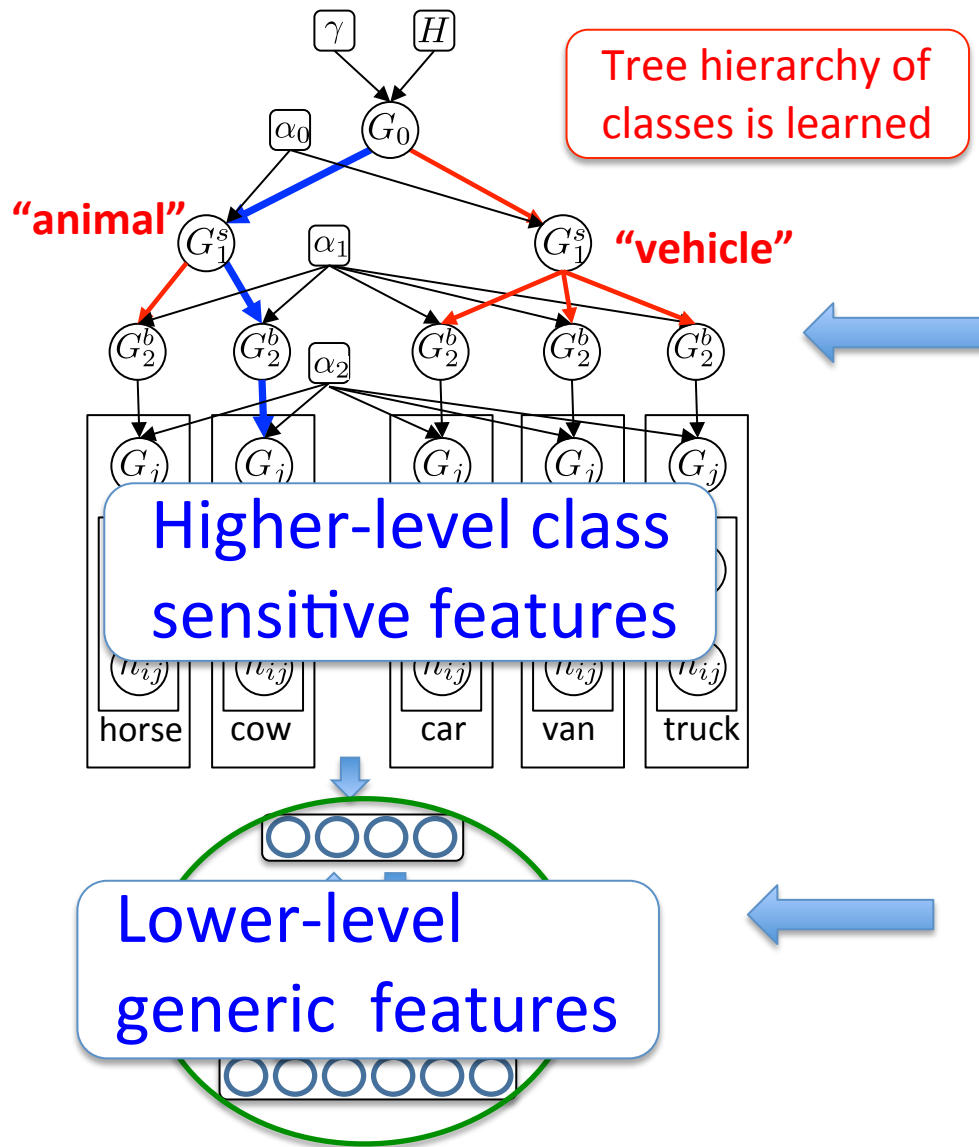
$\mathbf{h}^3 | \mathbf{z} \sim \text{HDP}$ (**Hierarchical Dirichlet Process**) prior:
 a nonparametric prior allowing categories to share higher-level features, or parts.

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

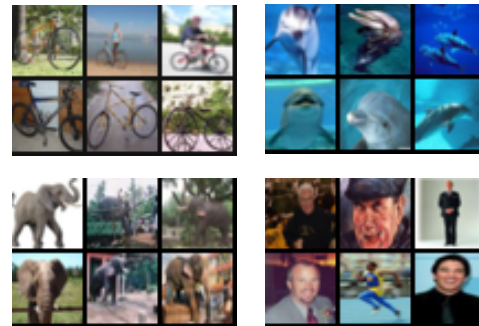


Enforce (approximate) global consistency through many local constraints.

CIFAR Object Recognition

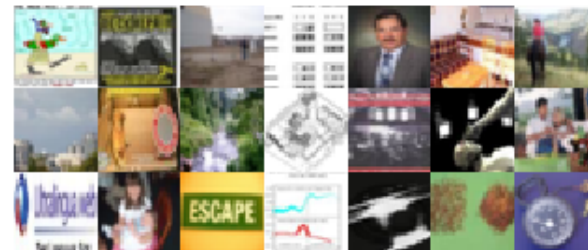


50,000 images of 100 classes



Inference: Markov chain Monte Carlo

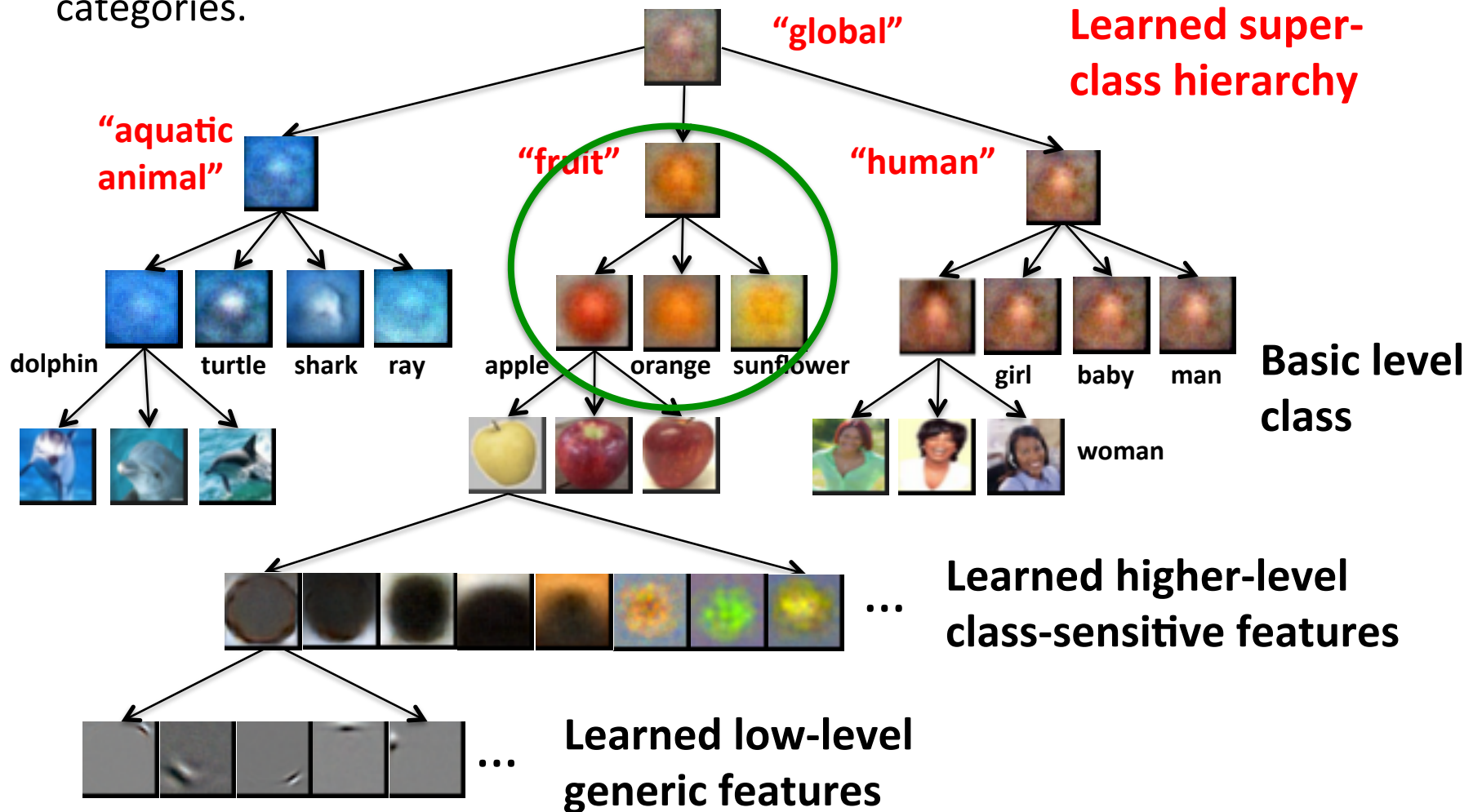
4 million unlabeled images



32 x 32 pixels x 3 RGB

Learning the Hierarchy

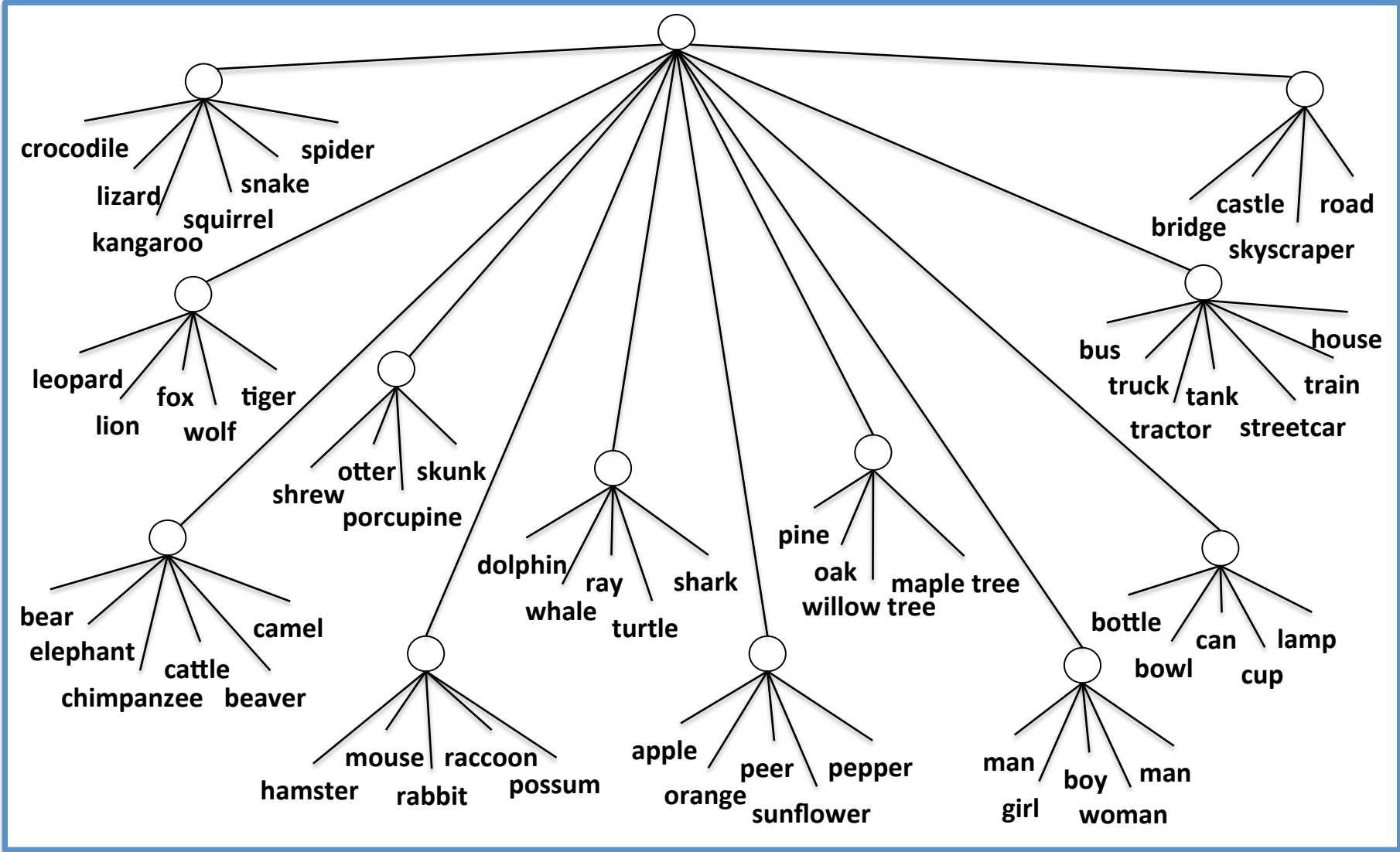
The model learns how to share the knowledge across many visual categories.



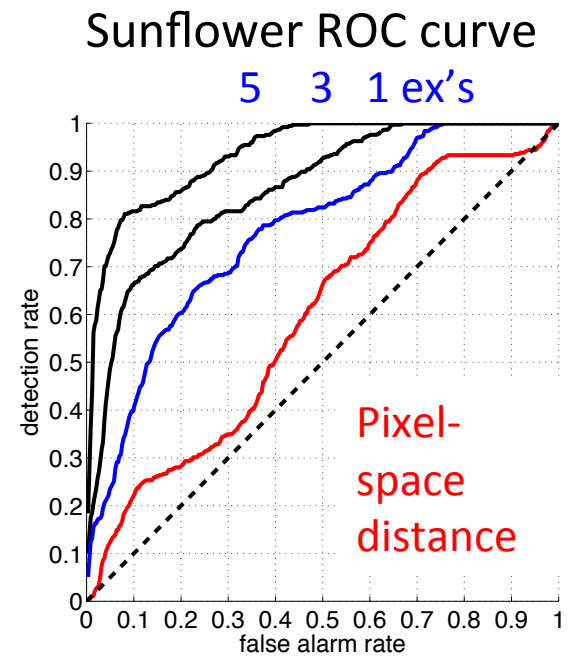
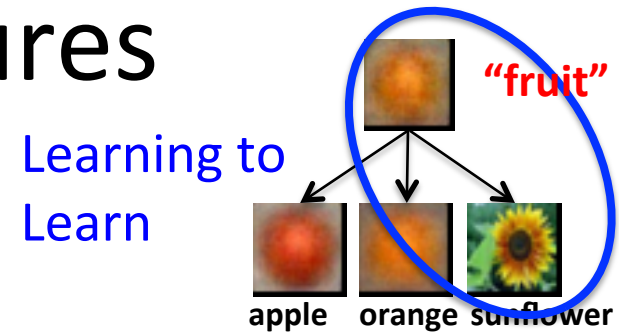
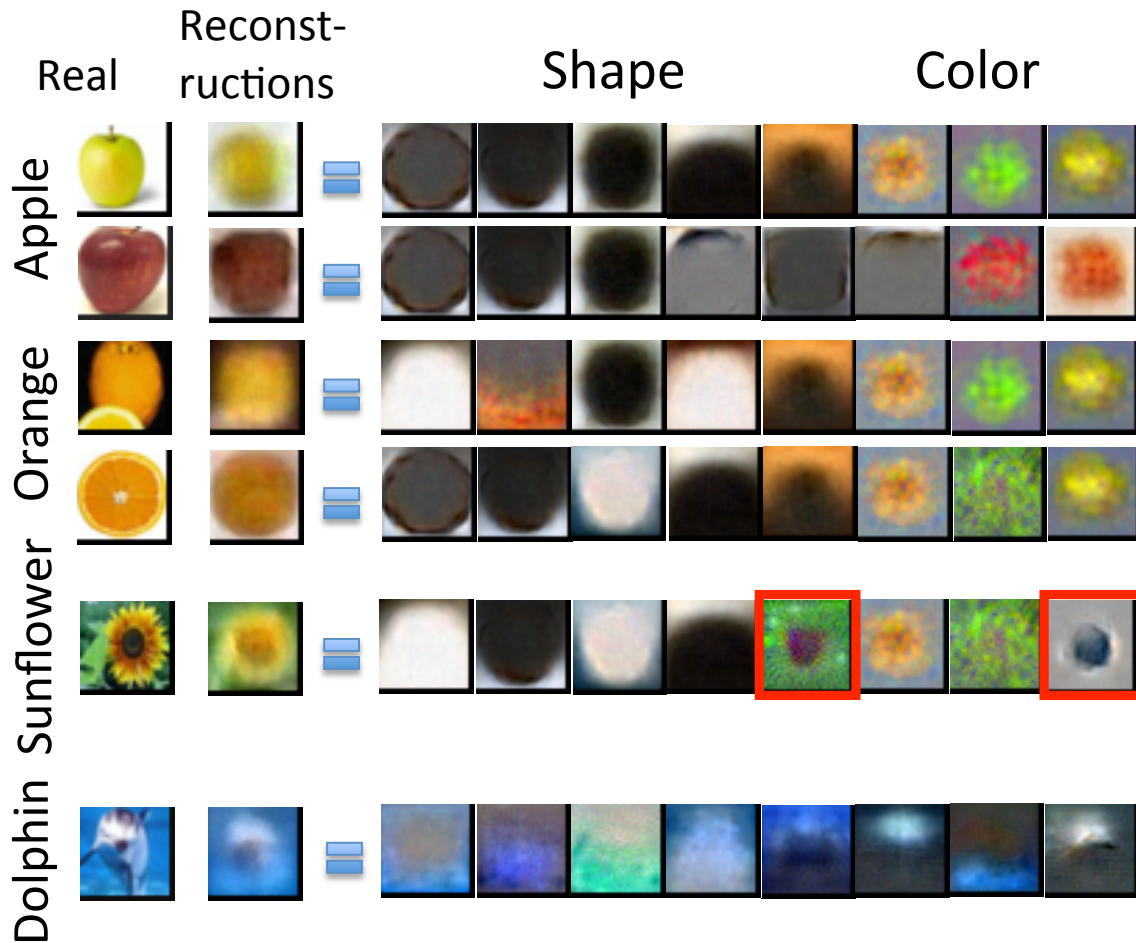


Learning the Hierarchy

The model learns how to share the knowledge across many visual



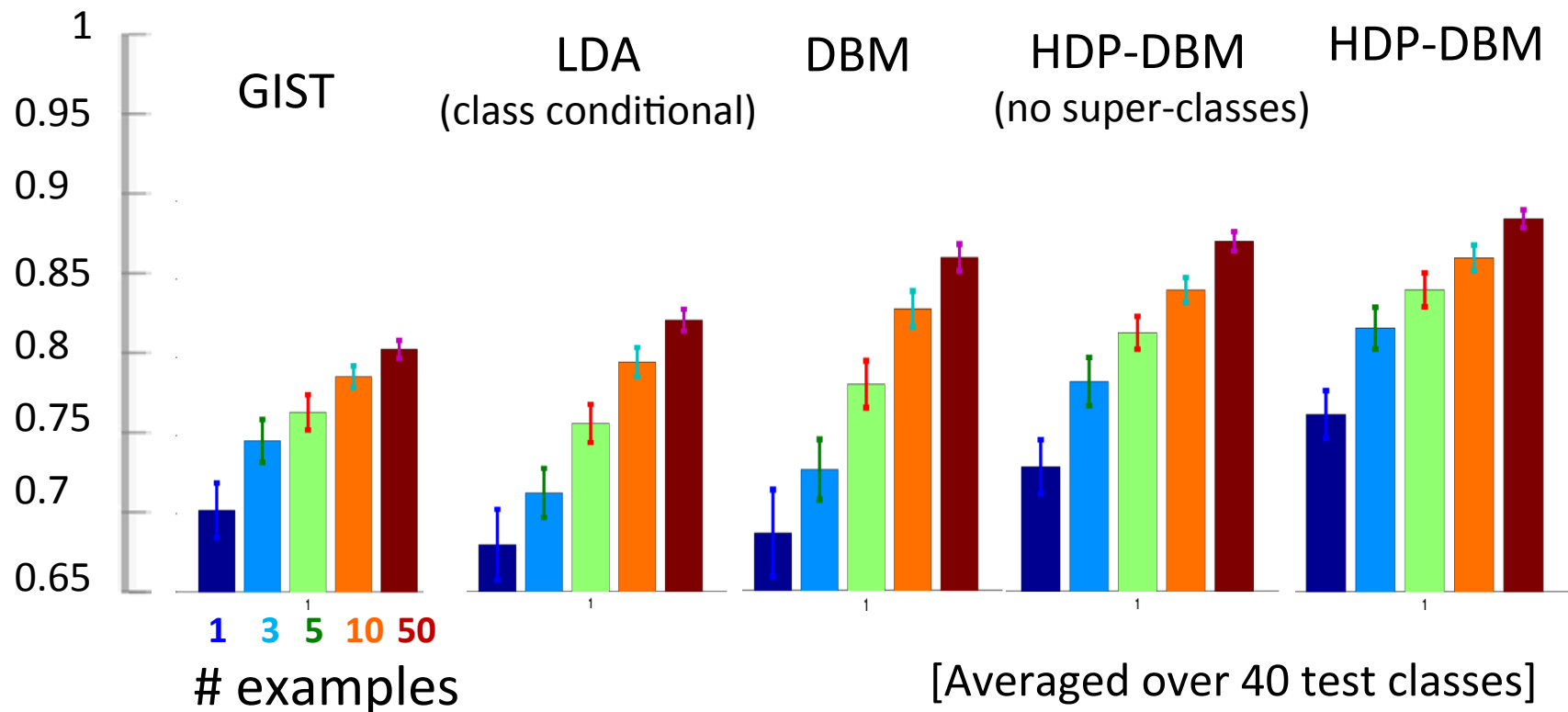
Sharing Features



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)



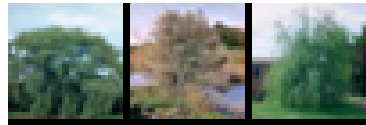
Our model outperforms standard computer vision features (e.g. GIST).

Learning from 3 Examples

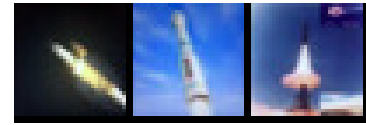
Given only 3 Examples



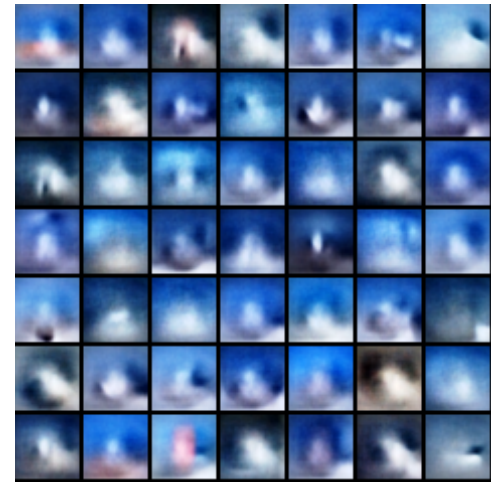
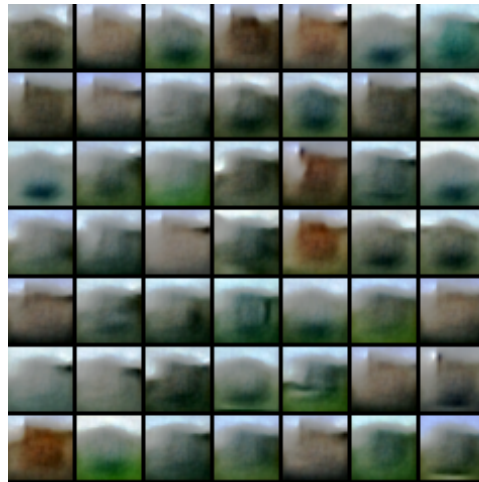
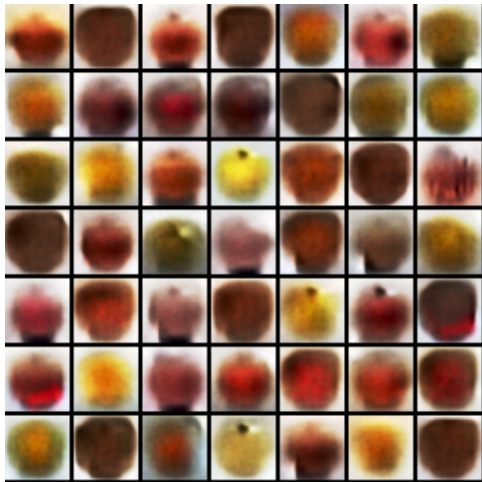
Willow Tree



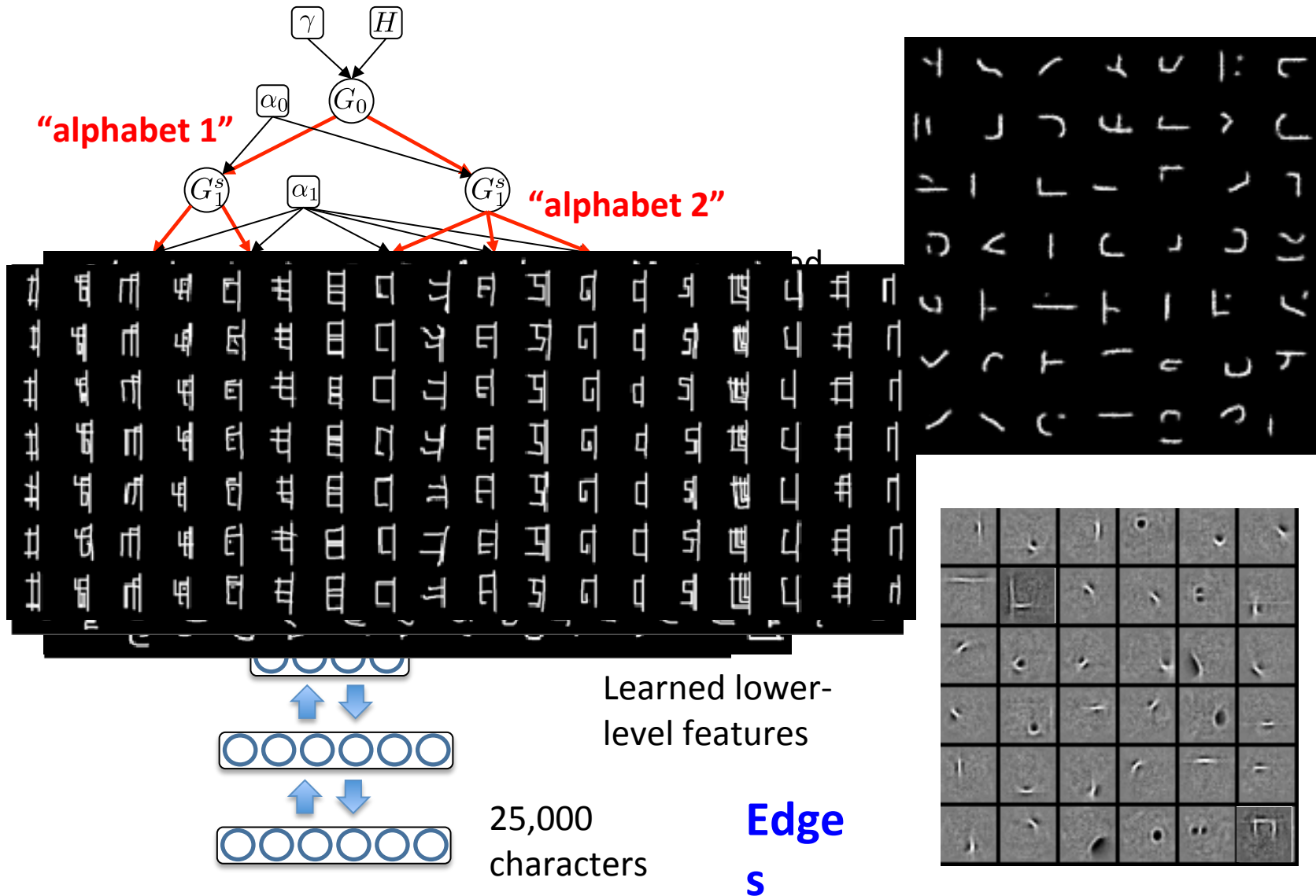
Rocket



Generated Samples

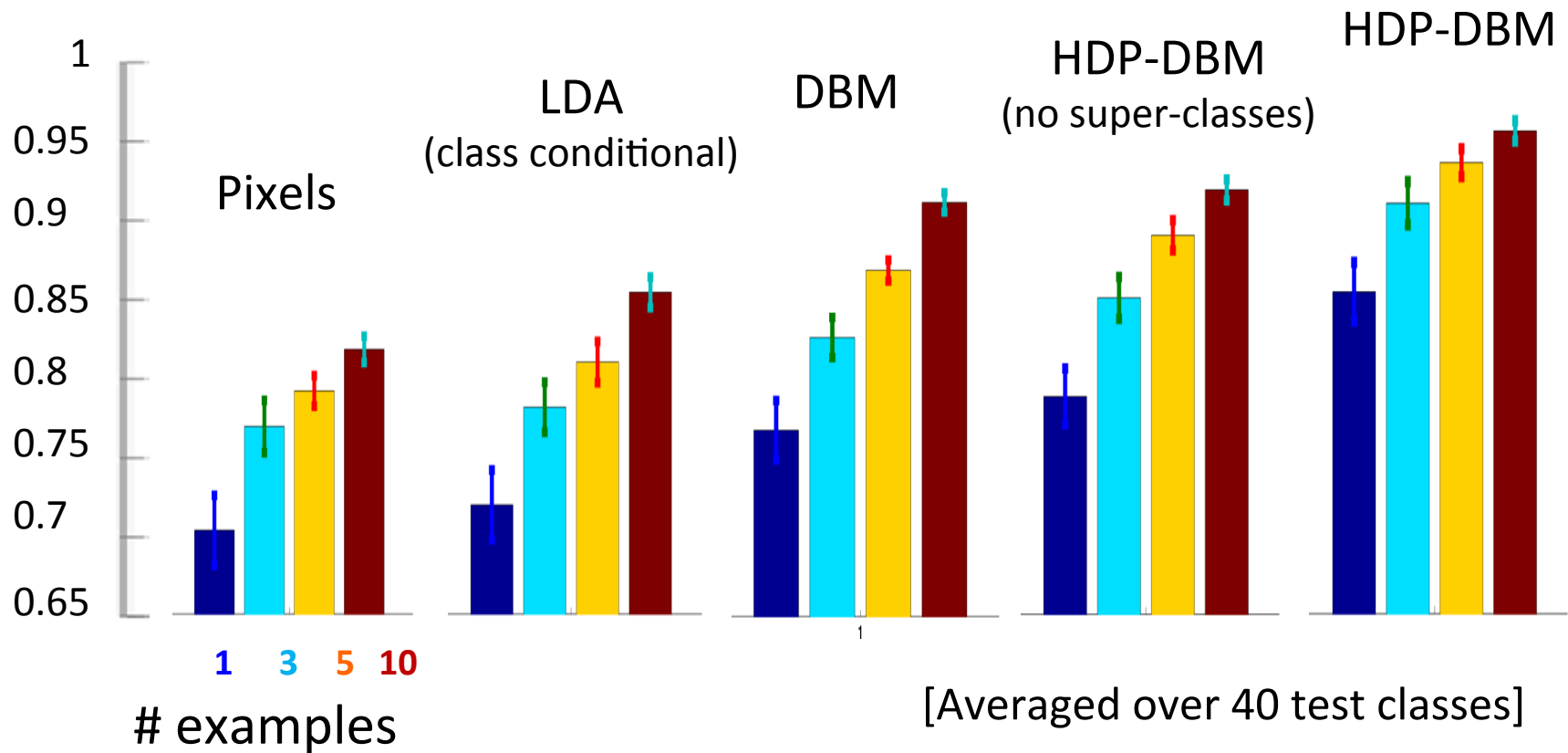


Handwritten Character Recognition

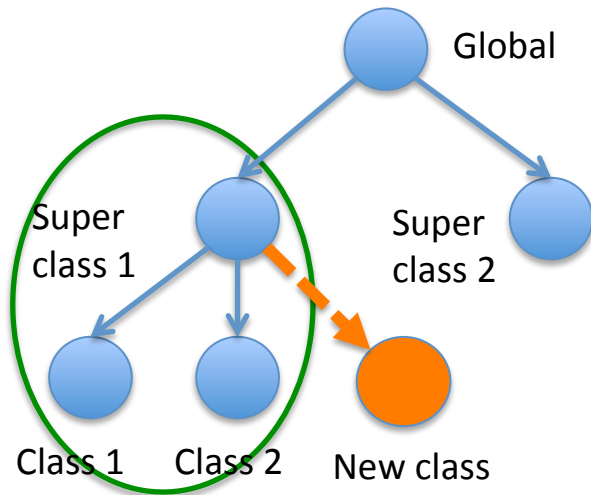


Handwritten Character Recognition

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



Simulating New Characters



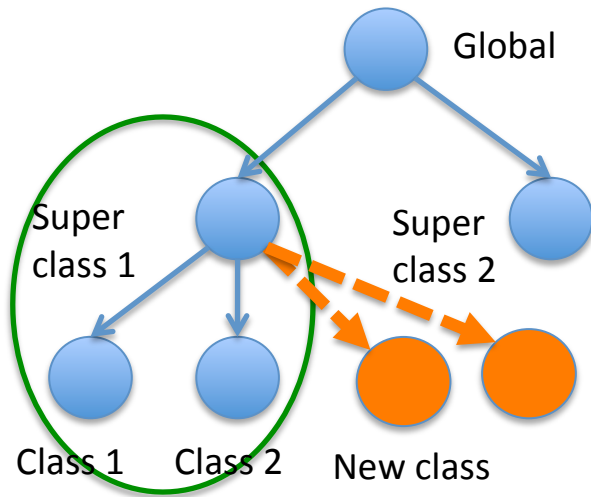
Real data within super class



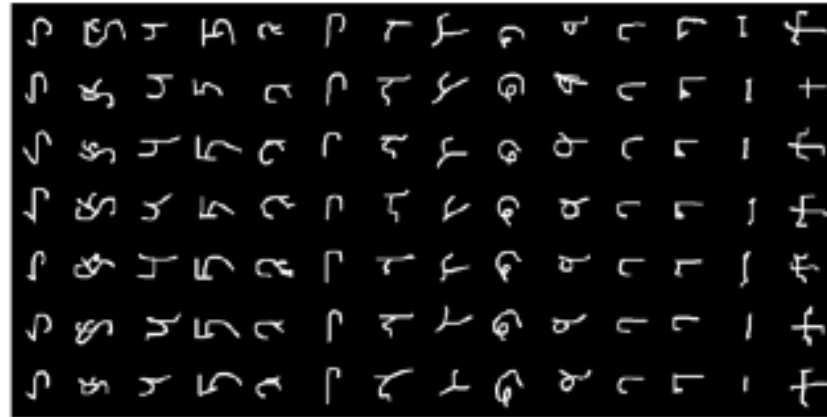
Simulated new characters



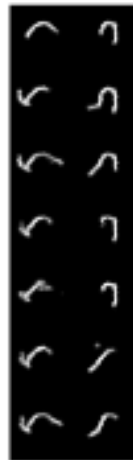
Simulating New Characters



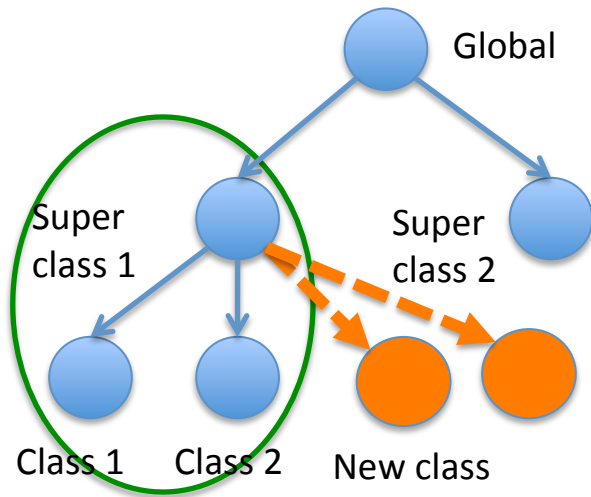
Real data within super class



Simulated new characters



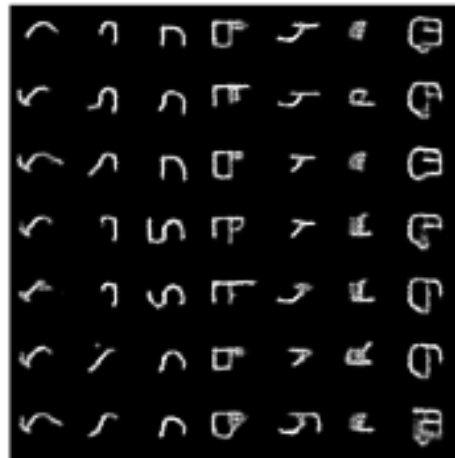
Simulating New Characters



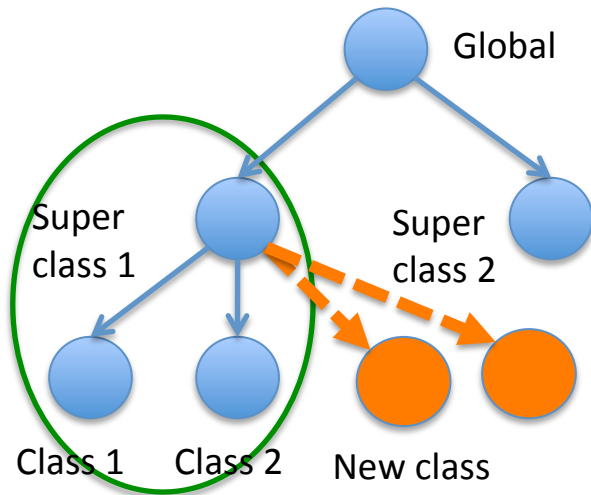
Real data within super class



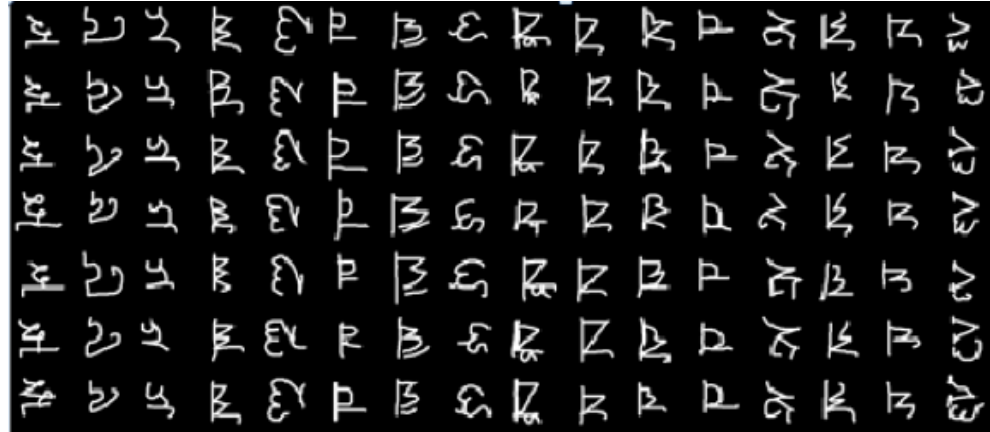
Simulated new characters



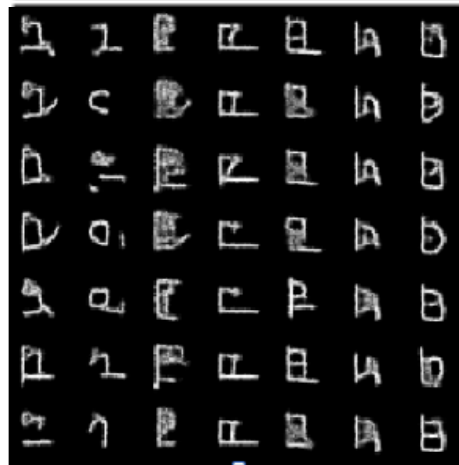
Simulating New Characters



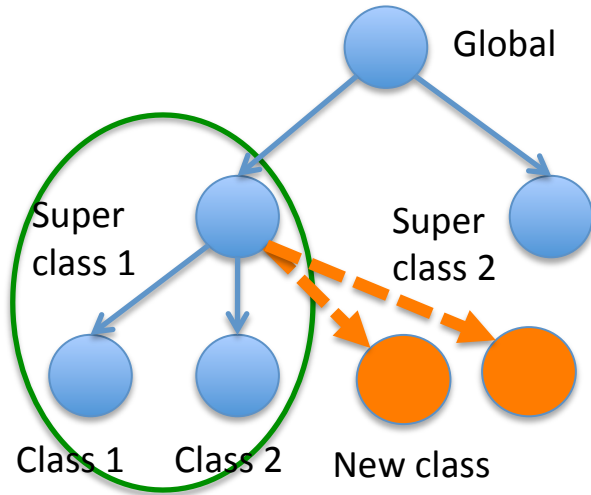
Real data within super class



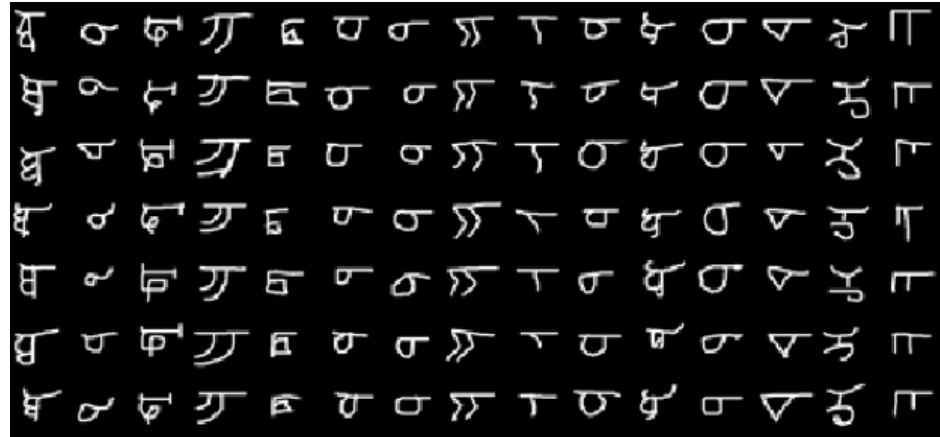
Simulated new characters



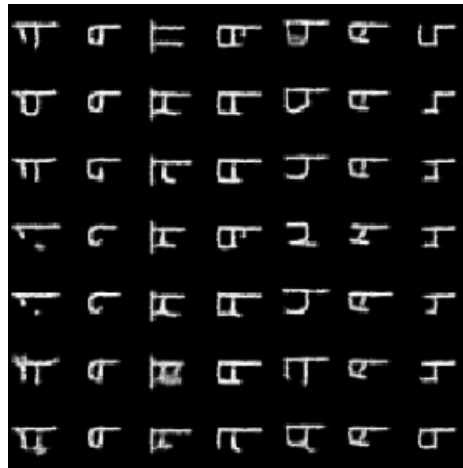
Simulating New Characters



Real data within super class

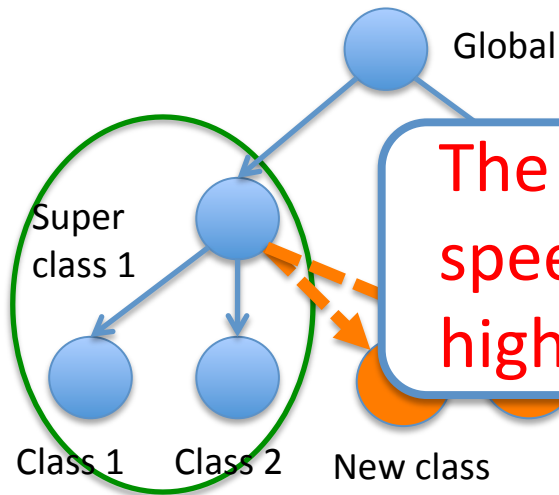


Simulated new characters

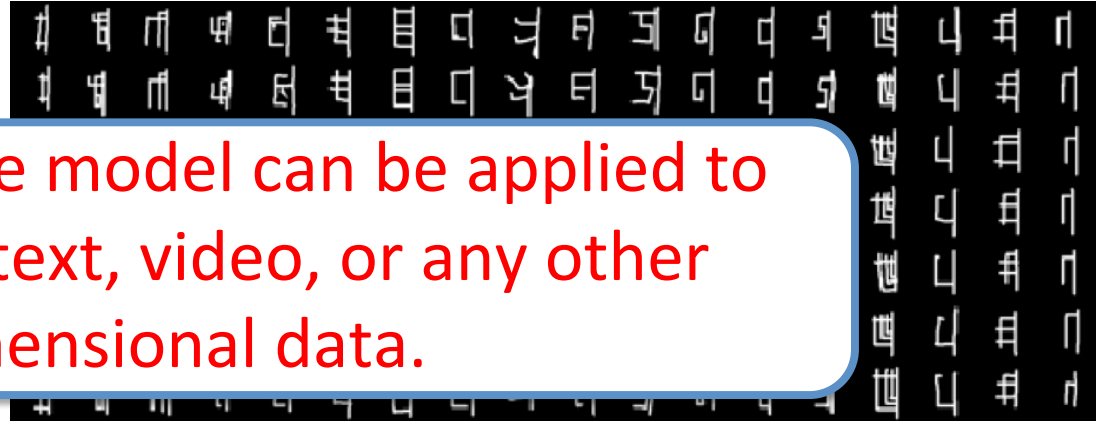


Simulating New Characters

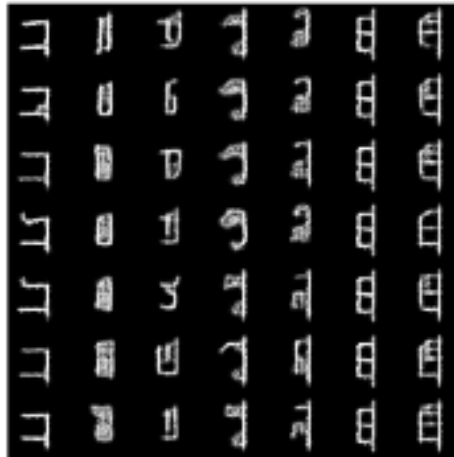
Real data within super class



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

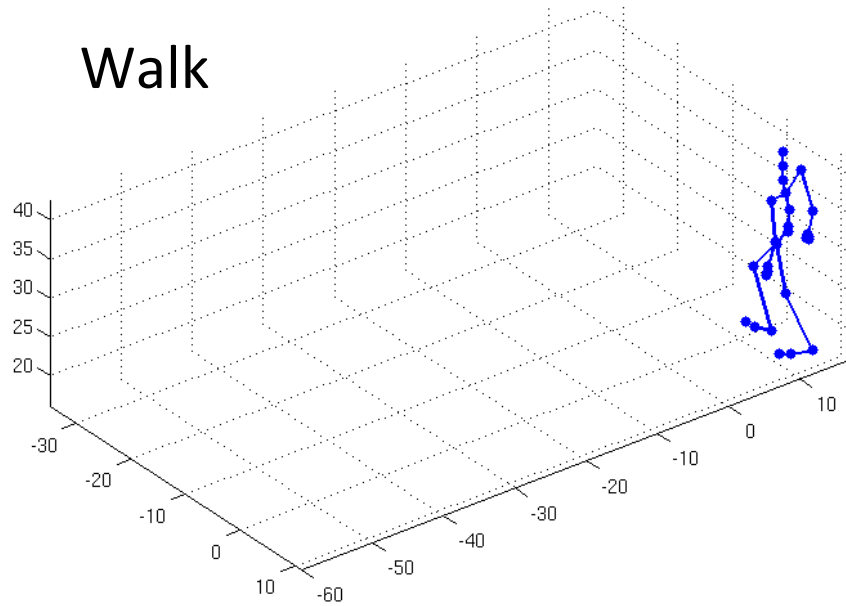


Simulated new characters

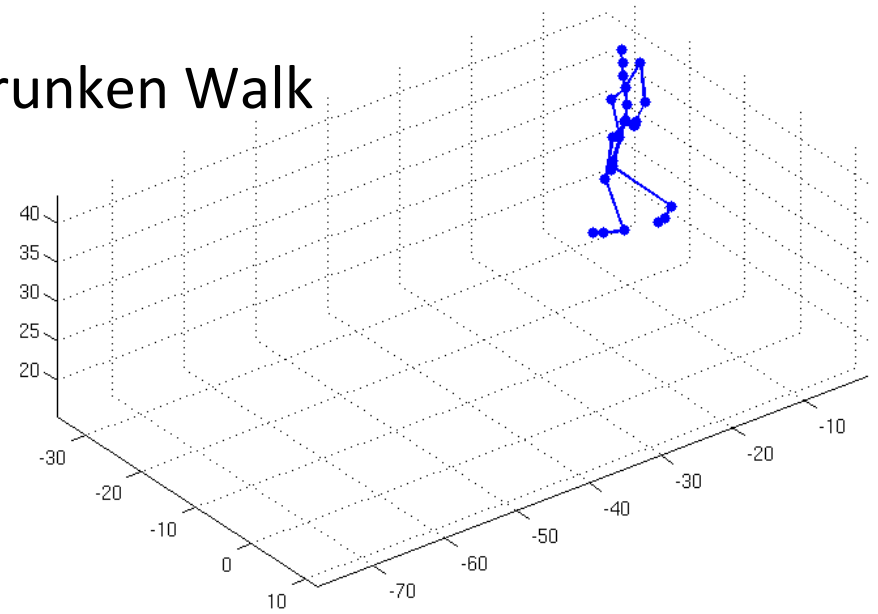


Motion Capture

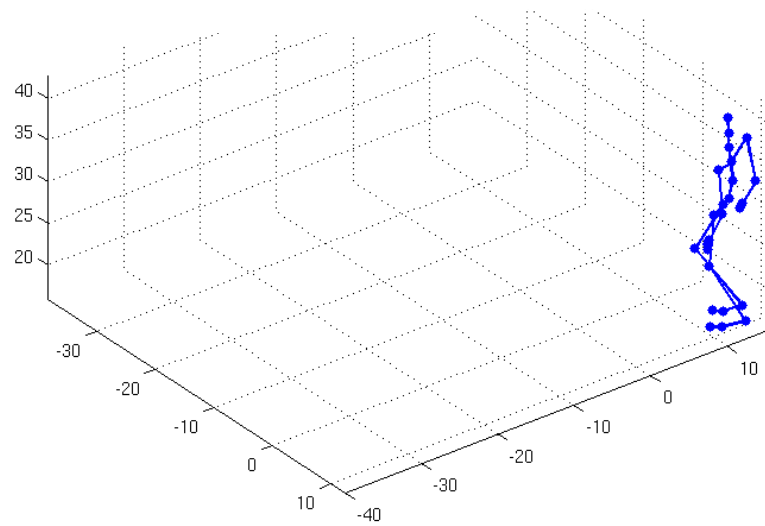
Walk



Drunken Walk



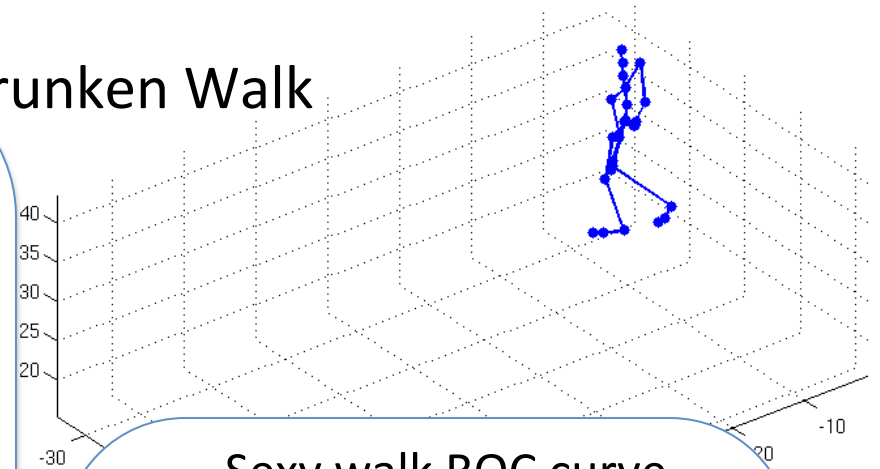
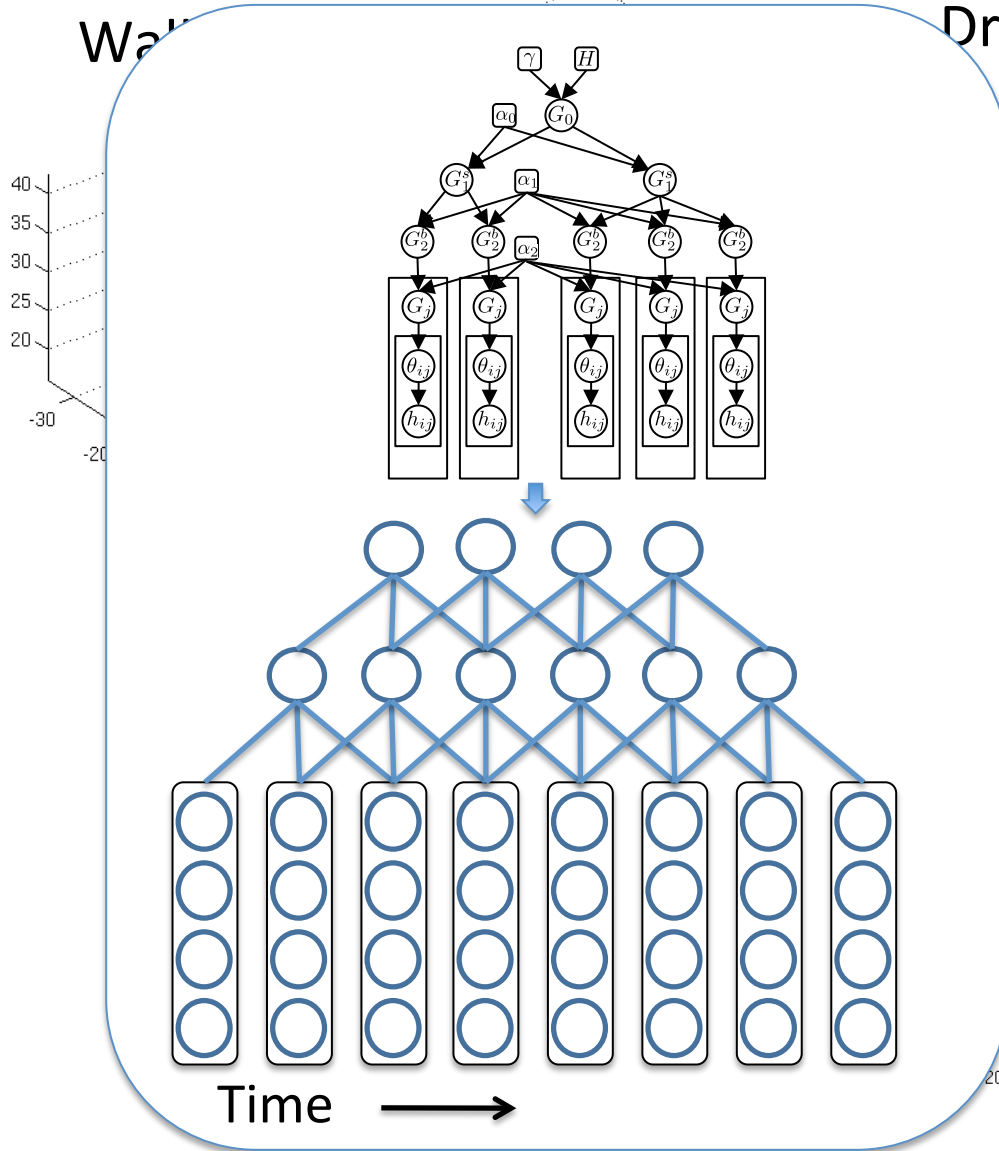
Sexy Walk



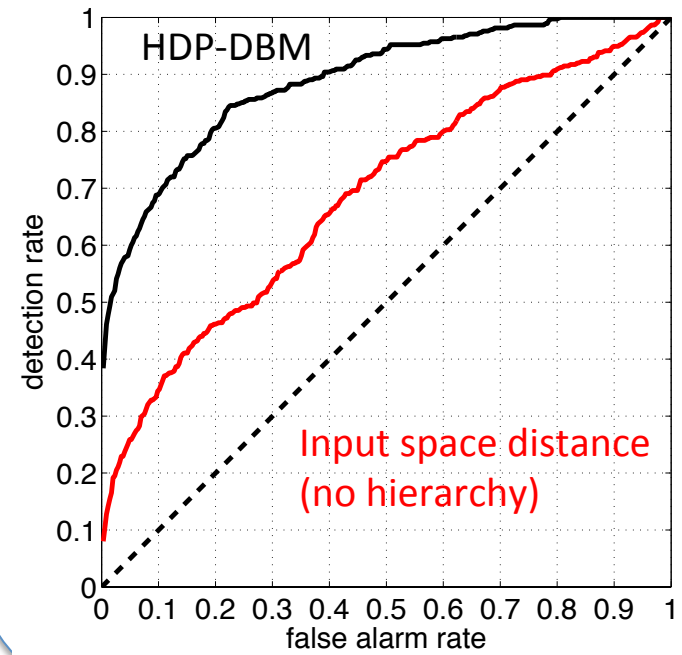
Motion Capture

Walk

Drunken Walk



Sexy walk ROC curve



Talk Roadmap

- Advanced Deep Models
 - Deep Boltzmann Machines
 - One-Shot and Transfer Learning
 - Learning Structured and Robust Deep Models
- Multimodal Learning
- Conclusions

Face Recognition

Yale B Extended Face Dataset

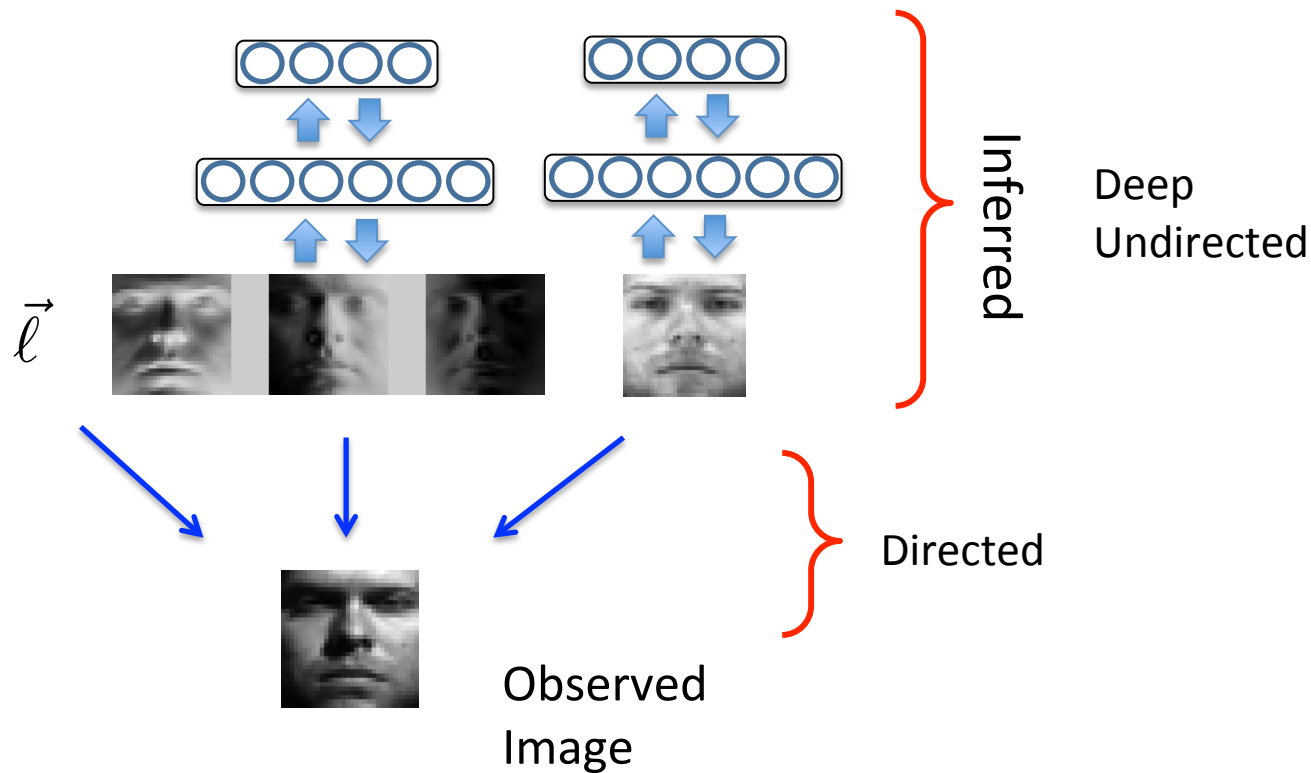
4 subsets of increasing illumination variations



Due to extreme illumination variations, deep models perform quite poorly on this dataset.

Deep Lambertian Model

Consider More Structured Models: undirected + directed models.



Combines the elegant properties of the Lambertian model with the Gaussian DBM model.

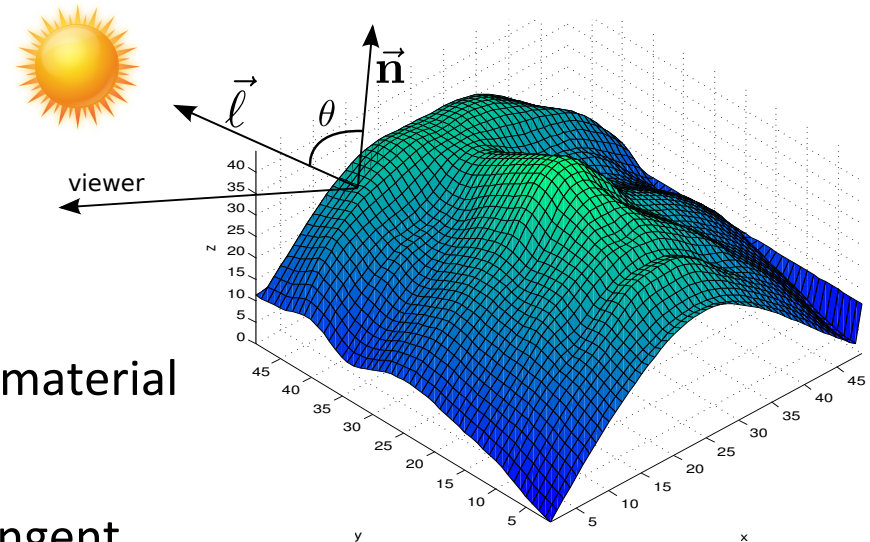
(Tang et. Al., ICML 2012, Tang et. al. CVPR 2012)

Lambertian Reflectance Model

- A simple model of the image formation process.

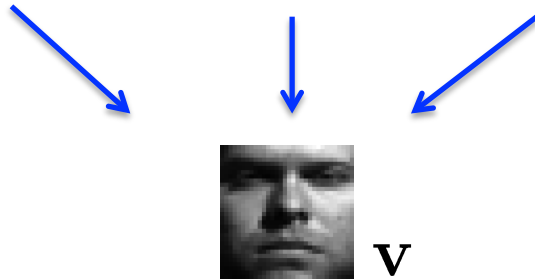
$$I = a \times |\vec{\ell}| |\vec{n}| \cos(\theta)$$

Image albedo Light source Surface normal



- Albedo -- diffuse reflectivity of a surface, material dependent, illumination independent.
- Surface normal -- perpendicular to the tangent plane at a point on the surface.
- Images with different illumination can be generated by varying light directions

Deep Lambertian Model



Observed Image

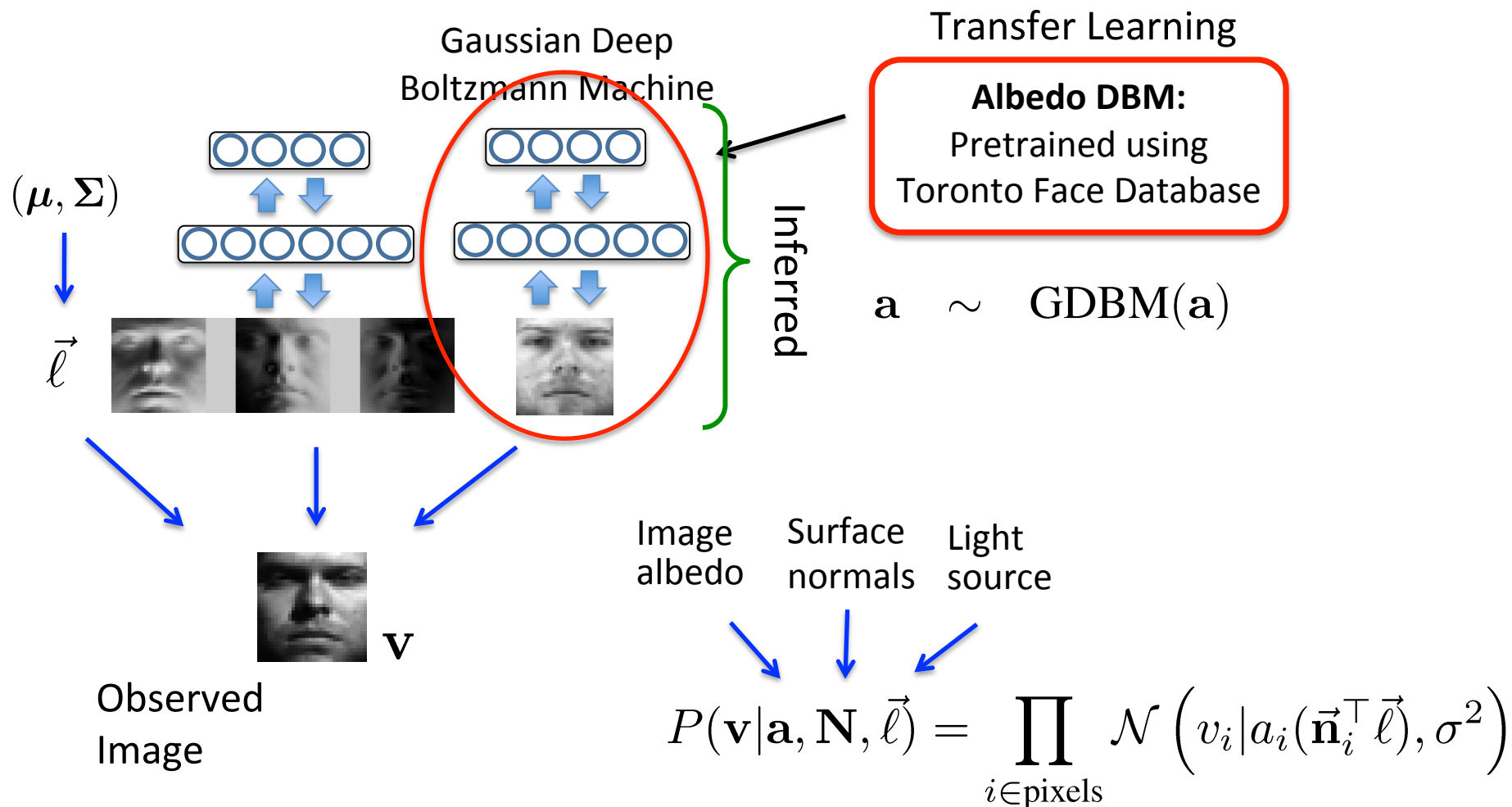
Image albedo Surface normals Light source

Three blue arrows point from the labels 'Image albedo', 'Surface normals', and 'Light source' to the corresponding variables \mathbf{a} , \mathbf{N} , and $\vec{\ell}$ in the probability function below.

$$P(\mathbf{v} | \mathbf{a}, \mathbf{N}, \vec{\ell}) = \prod_{i \in \text{pixels}} \mathcal{N}(v_i | a_i (\vec{\mathbf{n}}_i^\top \vec{\ell}), \sigma^2)$$

$$\mathbf{a} \in \mathbb{R}^D, \quad \mathbf{N} \in \mathbb{R}^{D \times 3}, \quad \ell \in \mathbb{R}^3$$

Deep Lambertian Model



Inference: Variational Inference.
Learning: Stochastic Approximation

$$\mathbf{a} \in \mathbb{R}^D, \quad \mathbf{N} \in \mathbb{R}^{D \times 3}, \quad \ell \in \mathbb{R}^3$$

Yale B Extended Face Dataset



- 38 subjects, ~ 45 images of varying illuminations per subject, divided into 4 subsets of increasing illumination variations.
- 28 subjects for training, and 10 for testing.

Face Relighting

One Test Image

Observed Inferred
 albedo

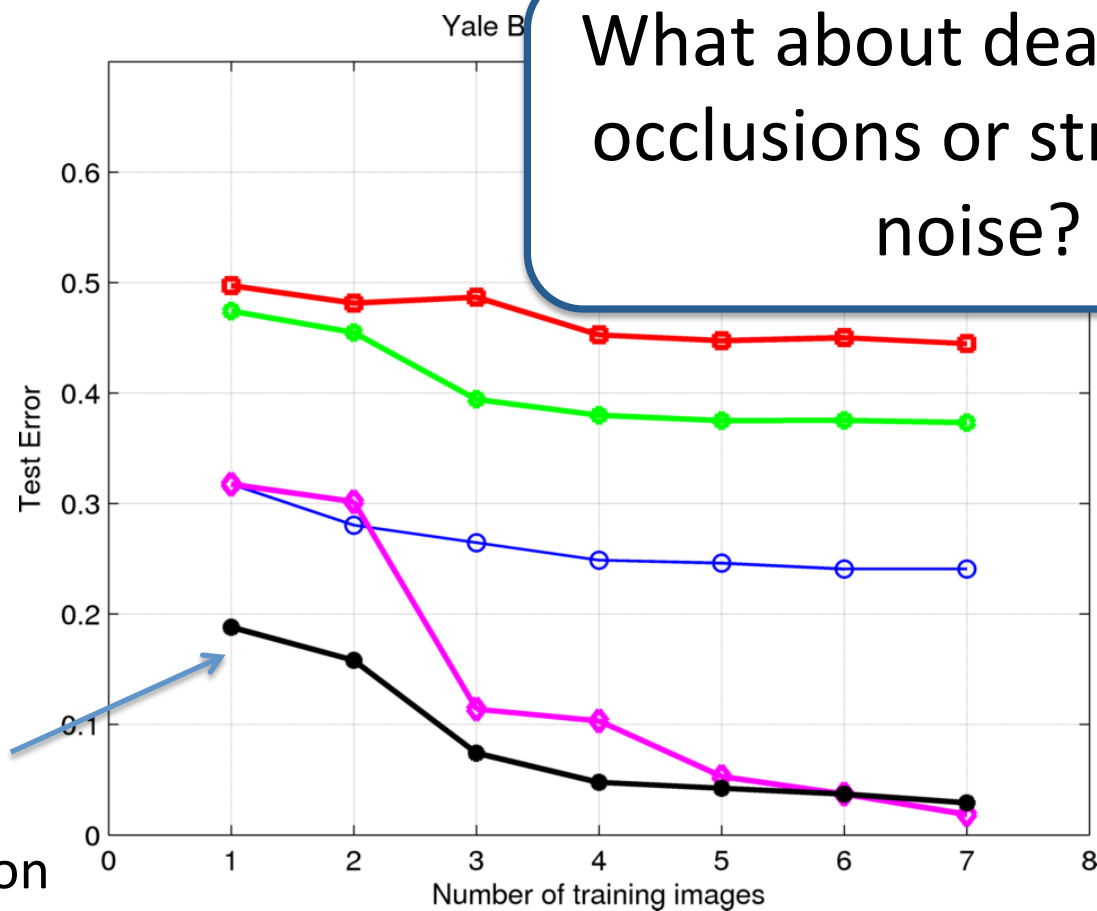


Face Relighting



Recognition Results

Recognition as function of the number of training images for 10 test subjects.

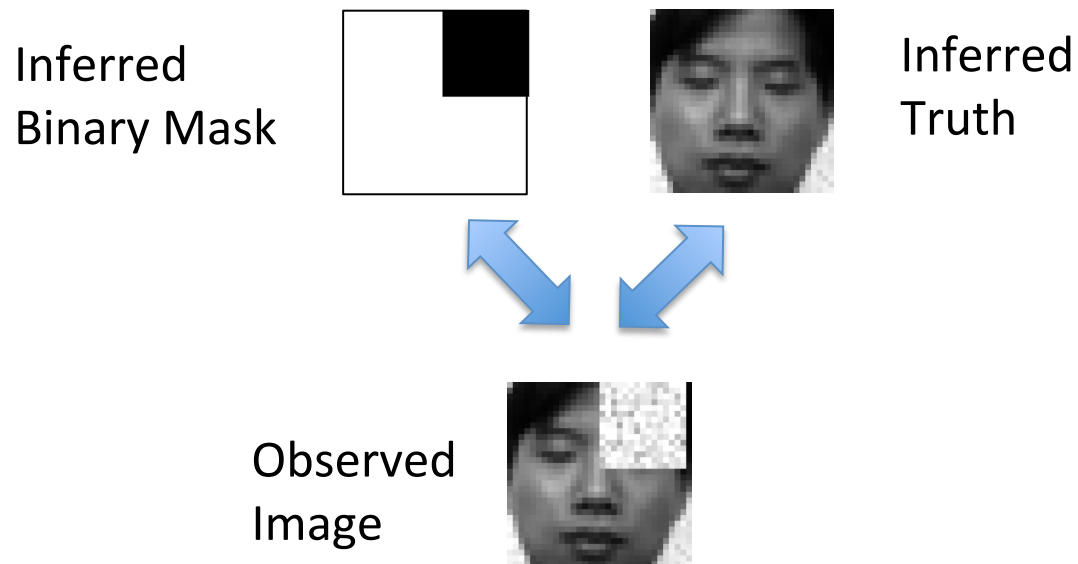


What about dealing with occlusions or structured noise?

Robust Boltzmann Machines

- Build more structured models that can deal with occlusions or structured noise.

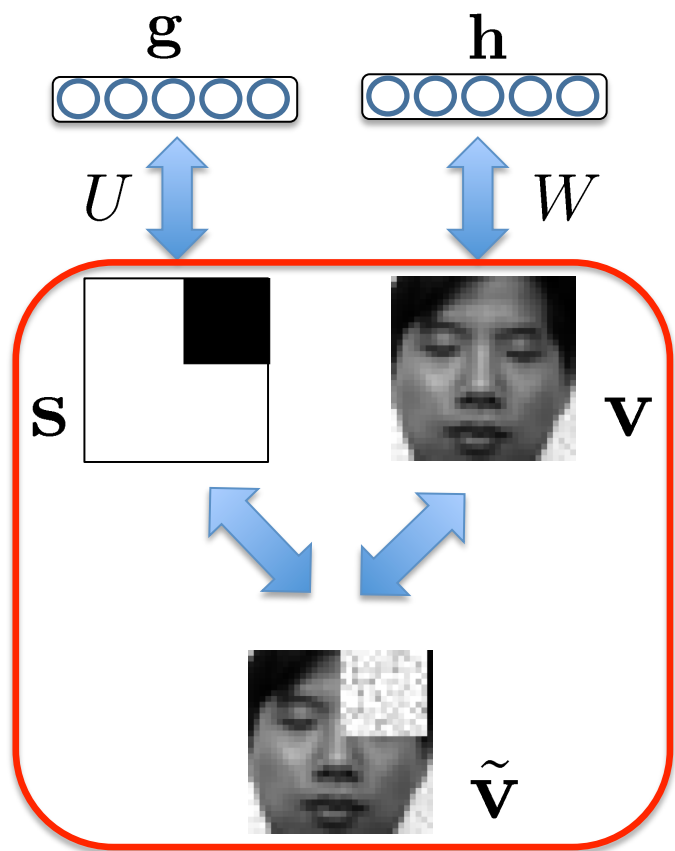
$$\log P(\tilde{\mathbf{v}}, \mathbf{v}, \mathbf{s}, \mathbf{h}, \mathbf{g}) \sim$$



(Tang et. Al., ICML 2012, Tang et. al. CVPR 2012)

Robust Boltzmann Machines

- Build more structured models that can deal with occlusions or structured noise.



$$\log P(\tilde{\mathbf{v}}, \mathbf{v}, \mathbf{s}, \mathbf{h}, \mathbf{g}) \sim$$

$$-\frac{1}{2} \sum_{i \in \text{pixels}} \frac{(v_i - b_i)^2}{\sigma_i^2} + \mathbf{v}^\top \mathbf{W} \mathbf{h} + \mathbf{s}^\top \mathbf{U} \mathbf{g}$$

Gaussian RBM, modeling clean faces

Binary RBM modeling occlusions

$$-\frac{1}{2} \sum_{i \in \text{pixels}} \gamma_i s_i (v_i - \tilde{v}_i)^2 - \frac{1}{2} \sum_{i \in \text{pixels}} \frac{(\tilde{v}_i - \tilde{b}_i)^2}{\tilde{\sigma}_i^2}$$

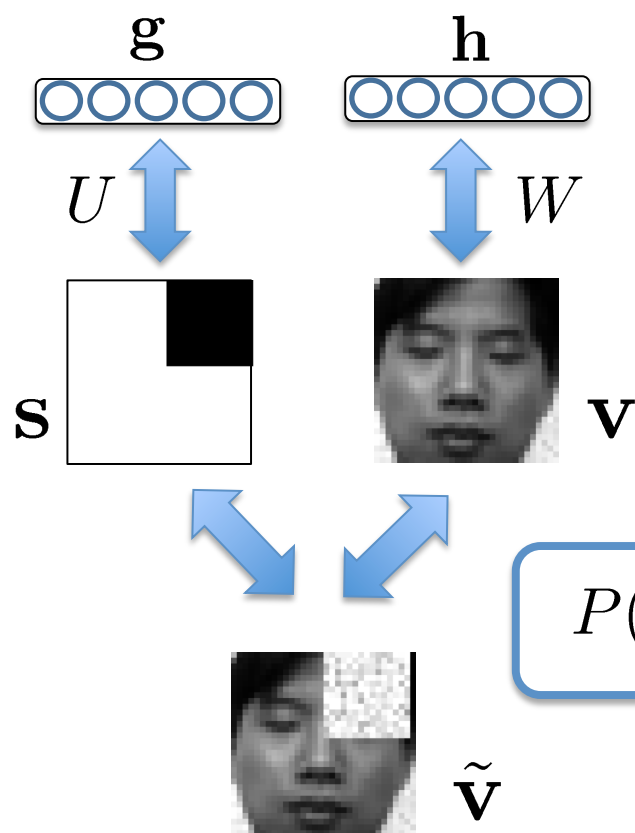
Binary pixel-wise Mask

Gaussian noise model

Observed Image

Robust Boltzmann Machines

- Build more structured models that can deal with occlusions or structured noise.



$$\log P(\tilde{\mathbf{v}}, \mathbf{v}, \mathbf{s}, \mathbf{h}, \mathbf{g}) \sim$$

$$-\frac{1}{2} \sum_{i \in \text{pixels}} \frac{(v_i - b_i)^2}{\sigma_i^2} + \mathbf{v}^\top W \mathbf{h} + \mathbf{s}^\top U \mathbf{g}$$

Gaussian RBM, modeling
clean faces

Binary RBM
modeling occlusions

$P(\tilde{\mathbf{v}} | \mathbf{h}, \mathbf{g})$ is a heavy-tailed distribution

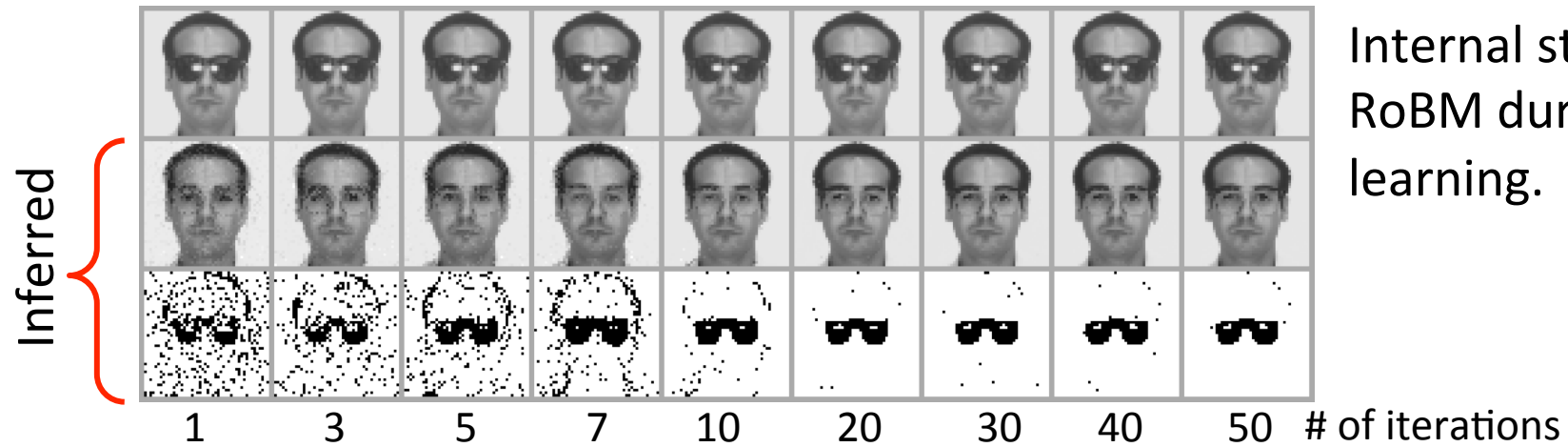
Binary pixel-wise
Mask

Gaussian noise model

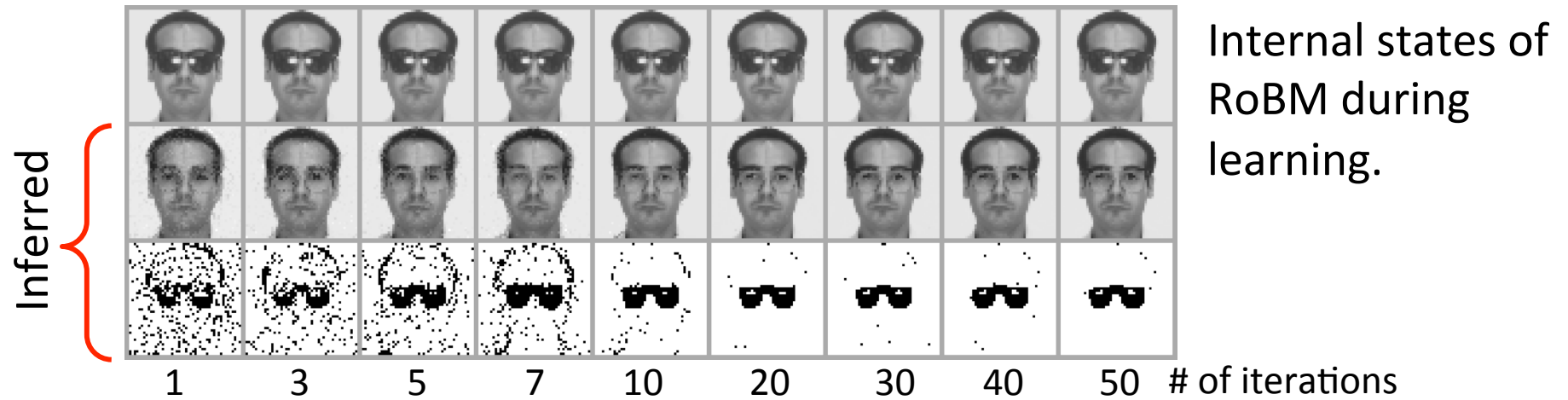
Inference: Variational Inference.

Learning: Stochastic Approximation

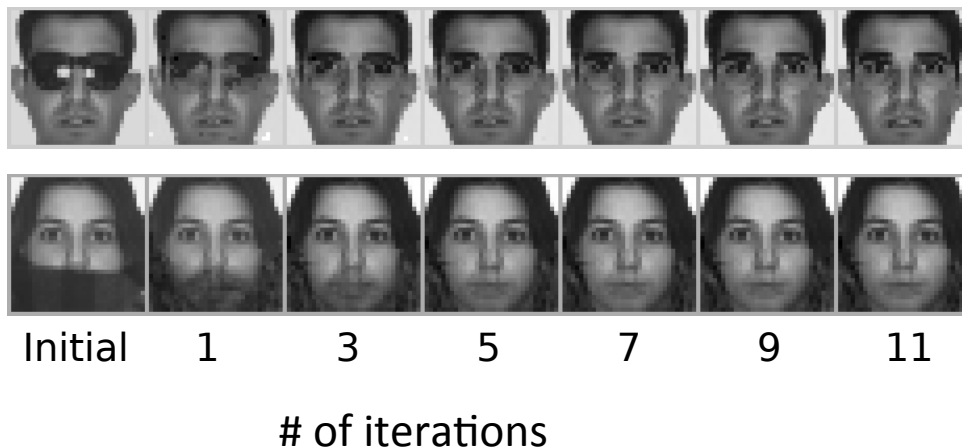
Recognition Results on AR Face Database



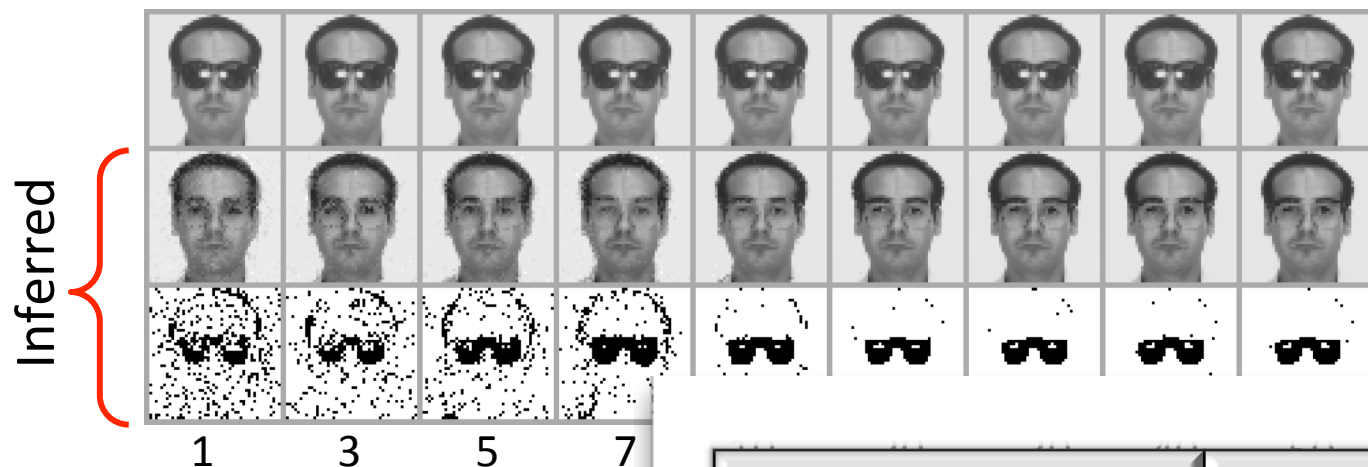
Recognition Results on AR Face Database



Inference on the test subjects



Recognition Results on AR Face Database



Internal states of RoBM during learning.

Inference on the



Initial 1 3 5

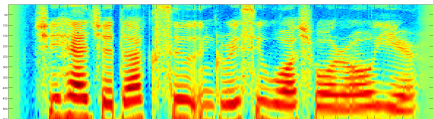
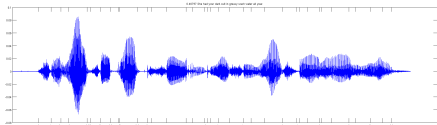
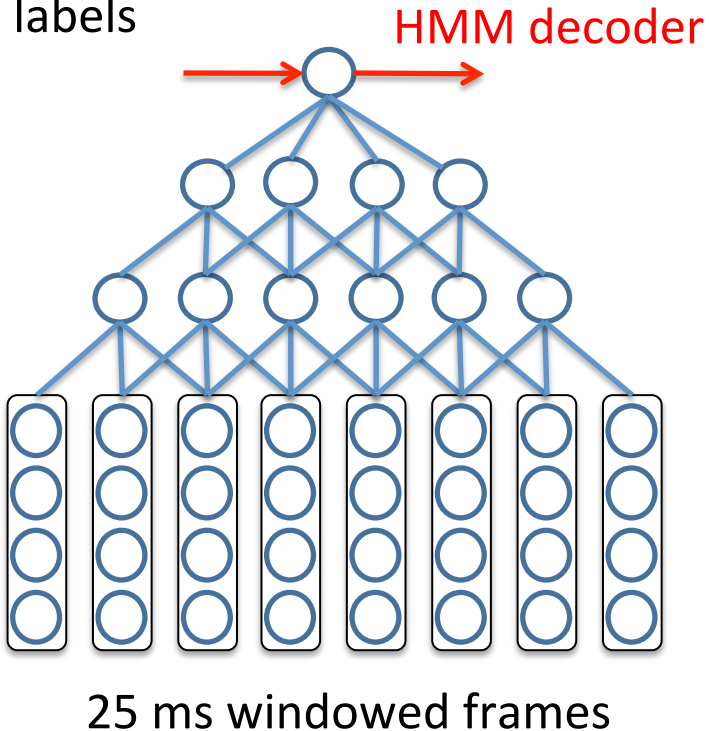
of iteration

Learning Algorithm	Sunglasses	Scarf
Robust BM	84.5%	80.7%
RBM	61.7%	32.9%
Eigenfaces	66.9%	38.6%
LDA	56.1%	27.0%
Pixel	51.3%	17.5%

Speech Recognition

(Zhang, Salakhutdinov, Chang, Glass, ICASSP 2012)

61 phonetic labels



- 630 speaker TIMIT corpus: 3,696 training and 944 test utterances.
- **Spoken Query Detection:**
For each keyword, estimate utterance's probability of containing that keyword.
- Performance: Average equal error rate (EER).

Learning Algorithm	AVG EER
GMM Unsupervised	16.4%
DBM Unsupervised	14.7%
DBM (1% labels)	13.3%
DBM (30% labels)	10.5%
DBM (100% labels)	9.7%

Talk Roadmap

- Advanced Deep Models
 - Deep Boltzmann Machines
 - One-Shot and Transfer Learning
 - Learning Structured and Robust Deep Models
- Multimodal Learning
- Conclusions

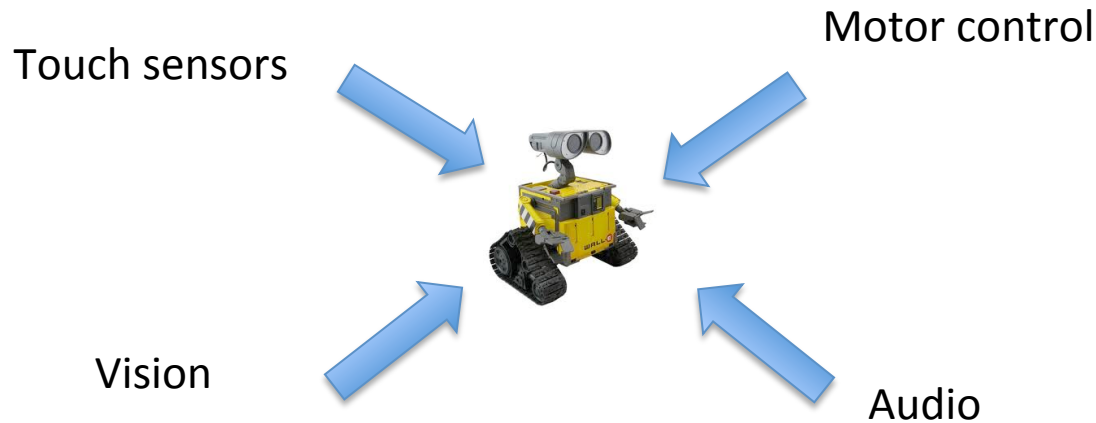
Data – Collection of Modalities

- Multimedia content on the web - image + text + audio.
- Product recommendation systems.
- Robotics applications.

You Tube Google flickr



amazon



Shared Concept

“Modality-free” representation

“Concept”



sunset, pacific ocean,
baker beach, seashore,
ocean

“Modality-full” representation

Multi-Modal Input

- Improve Classification



pentax, k10d, kangarooisland
southaustralia, sa australia
australiansealion 300mm



SEA / NOT SEA

- Fill in Missing Modalities



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

- Retrieve data from one modality when queried using data from another modality

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



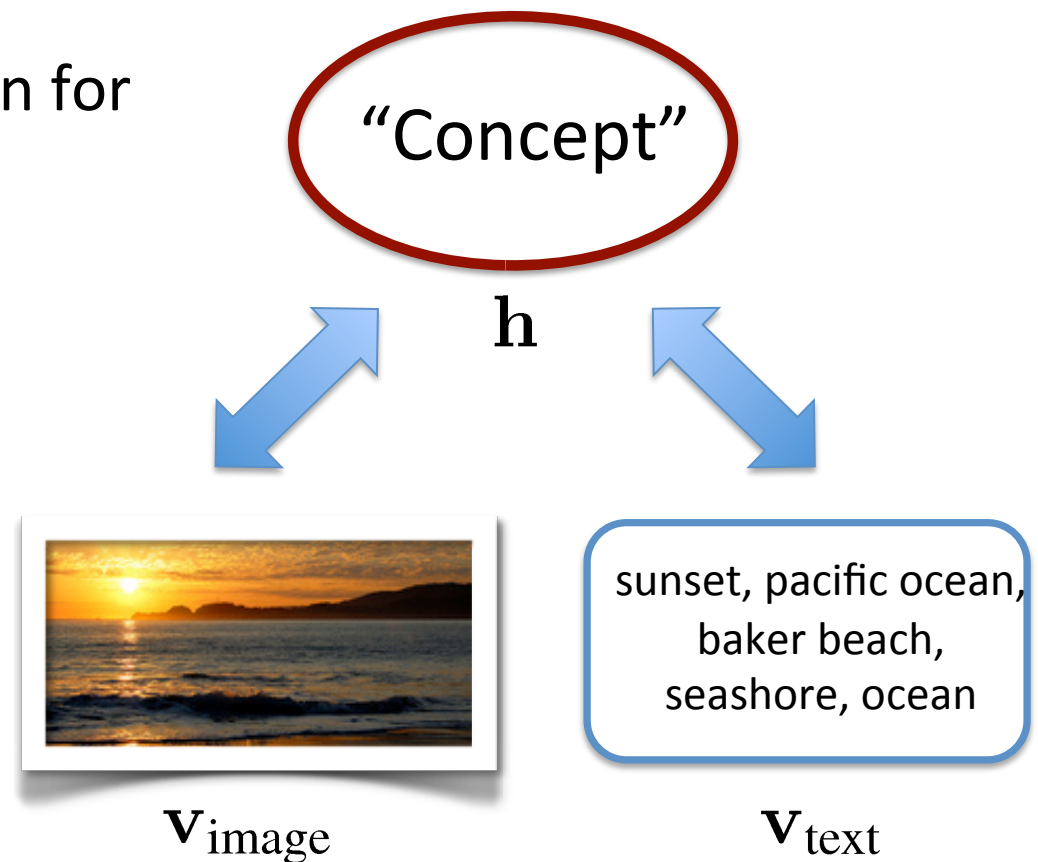
Building a Probabilistic Model

- Learn a joint density model:

$$P(\mathbf{h}, \mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}}).$$

$$P(\mathbf{h} | \mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}})$$

- \mathbf{h} : “fused” representation for classification, retrieval.



Building a Probabilistic Model

- Learn a joint density model:

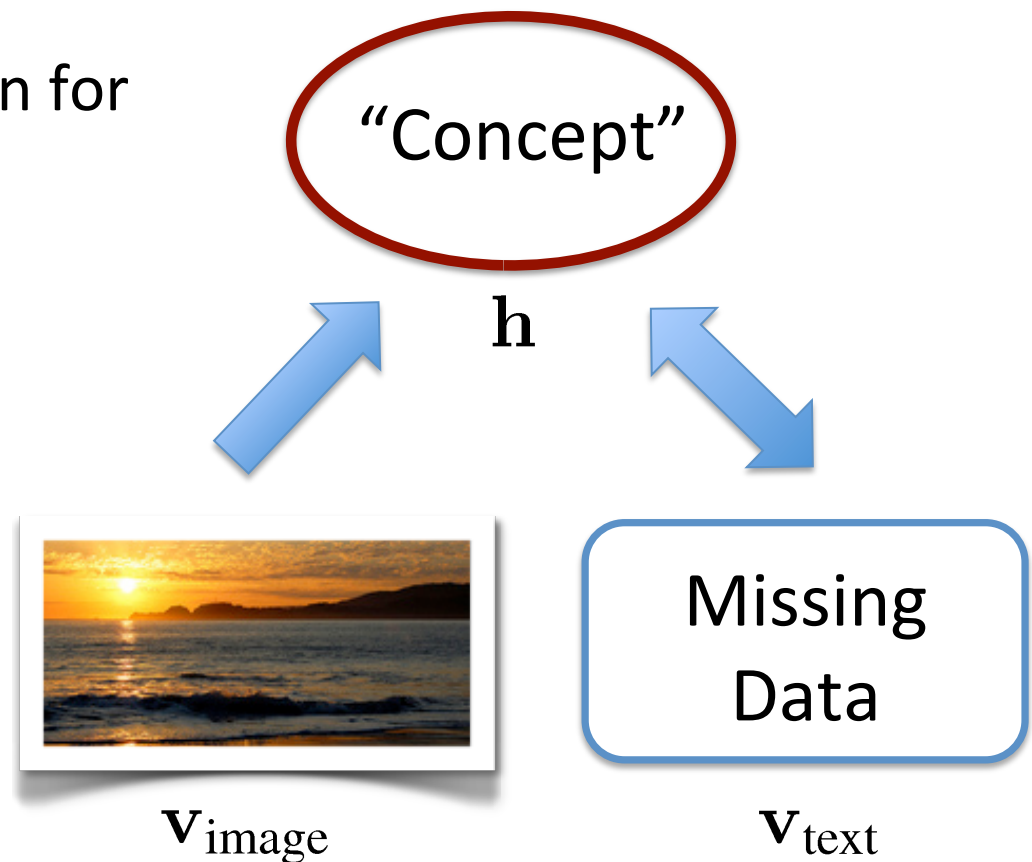
$$P(\mathbf{h}, \mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}}).$$

$$P(\mathbf{h}, \mathbf{v}_{\text{text}} | \mathbf{v}_{\text{image}})$$

- \mathbf{h} : “fused” representation for classification, retrieval.

- Generate data from conditional distributions for

- Image Annotation



Building a Probabilistic Model

- Learn a joint density model:

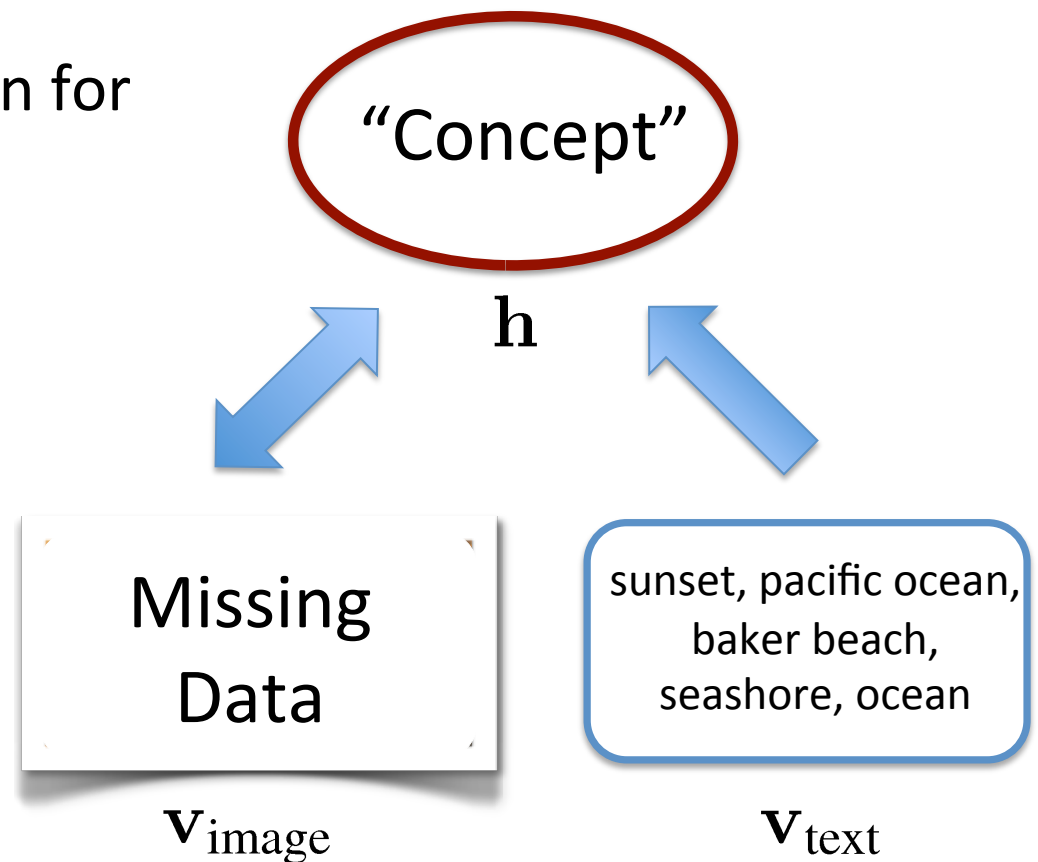
$$P(\mathbf{h}, \mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}}).$$

$$P(\mathbf{h}, \mathbf{v}_{\text{image}} | \mathbf{v}_{\text{text}})$$

- \mathbf{h} : “fused” representation for classification, retrieval.

- Generate data from conditional distributions for

- Image Annotation
- Image Retrieval

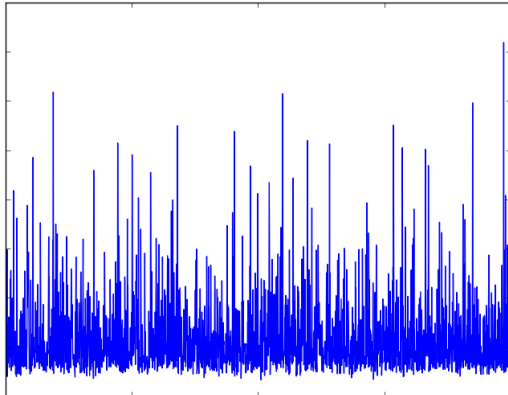


Challenges - I

Image



Dense

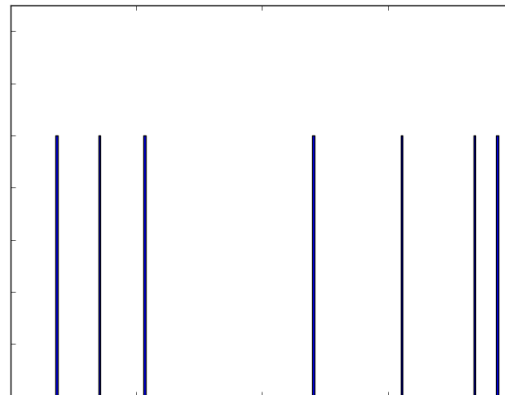


Text

sunset, pacific ocean,
baker beach, seashore,
ocean



Sparse



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

mickikrimmel,
mickipedia,
headshot

< no text >

unseulpixel,
naturey, crap

Noisy and missing data

Challenges - II

Image

Text

Text generated by the model



pentax, k10d,
pentaxda50200,
kangarooisland, sa,
australiansealion

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

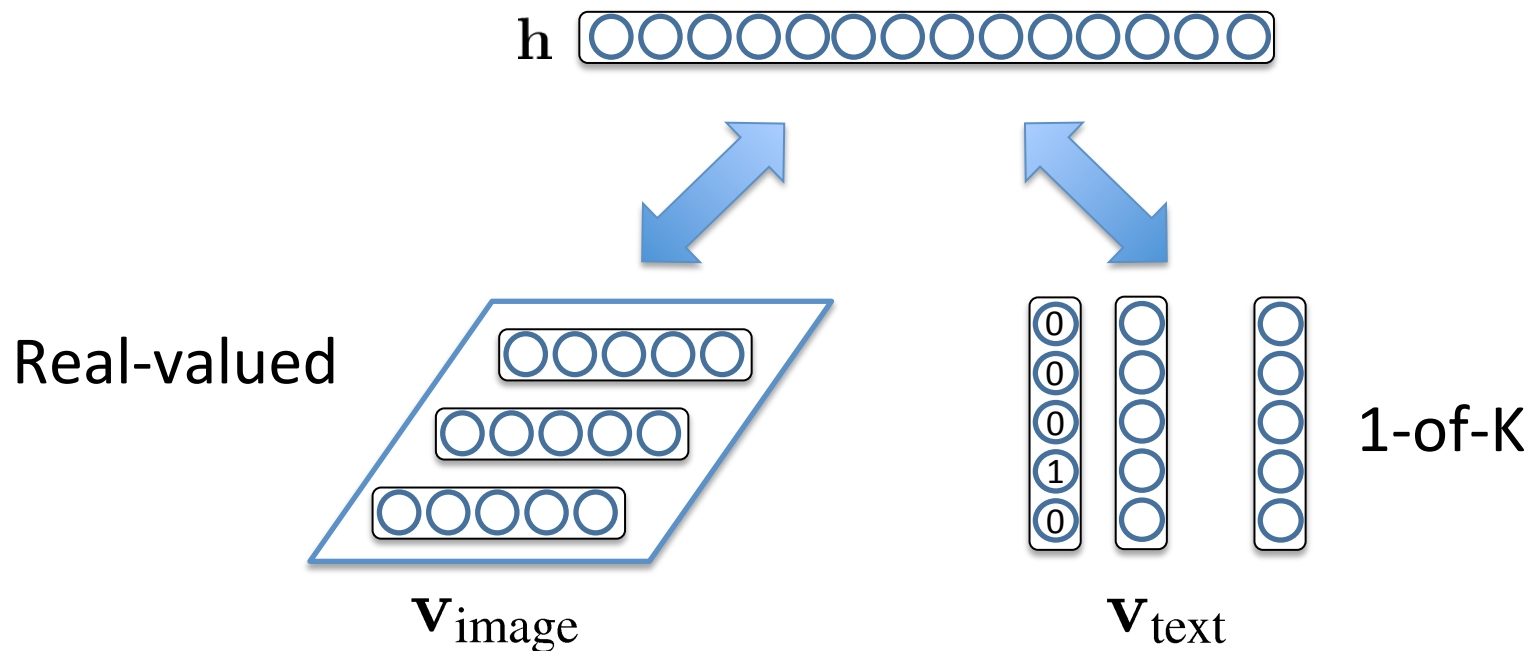


unseulpixel,
naturey, crap

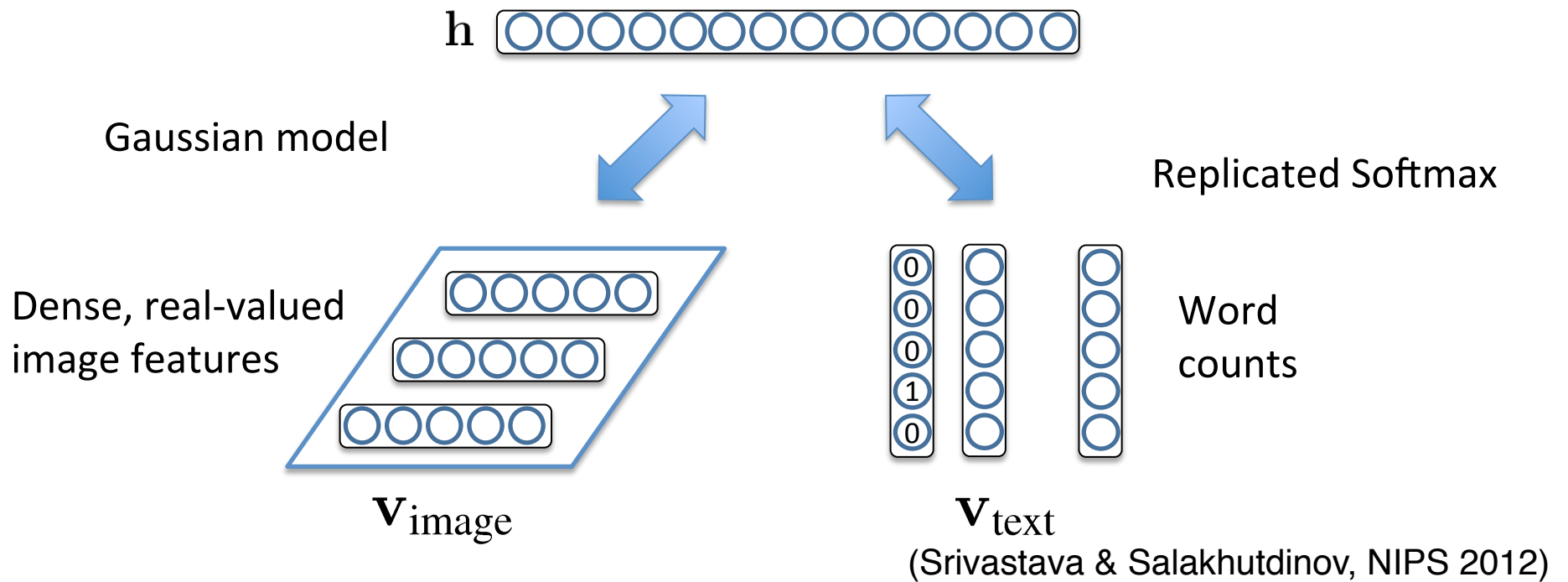
fall, autumn, trees, leaves,
foliage, forest, woods,
branches, path

A Simple Multimodal Model

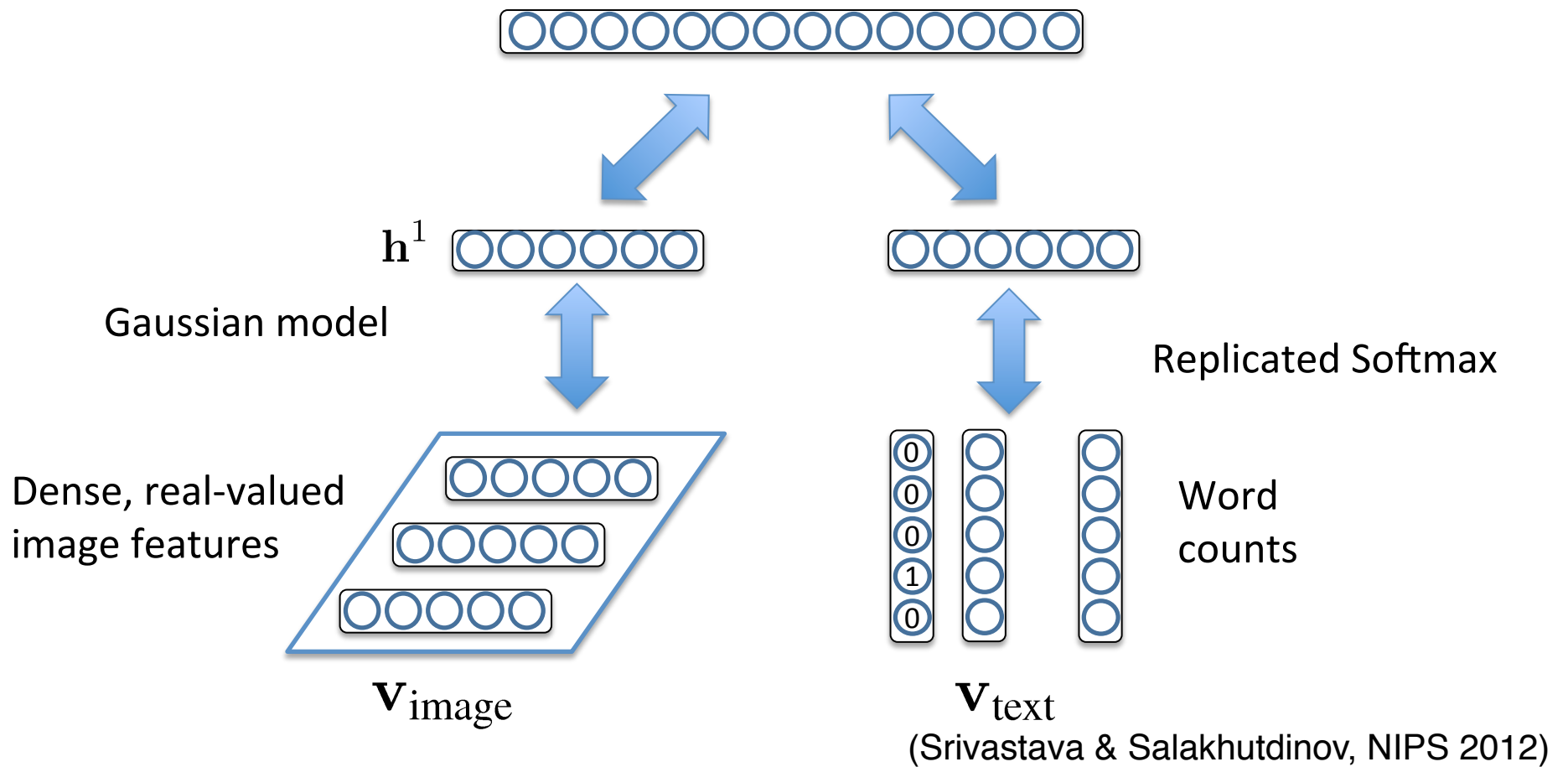
- Use a joint binary hidden layer.
- **Problem:** Inputs have very different statistical properties.
- Difficult to learn cross-modal features.



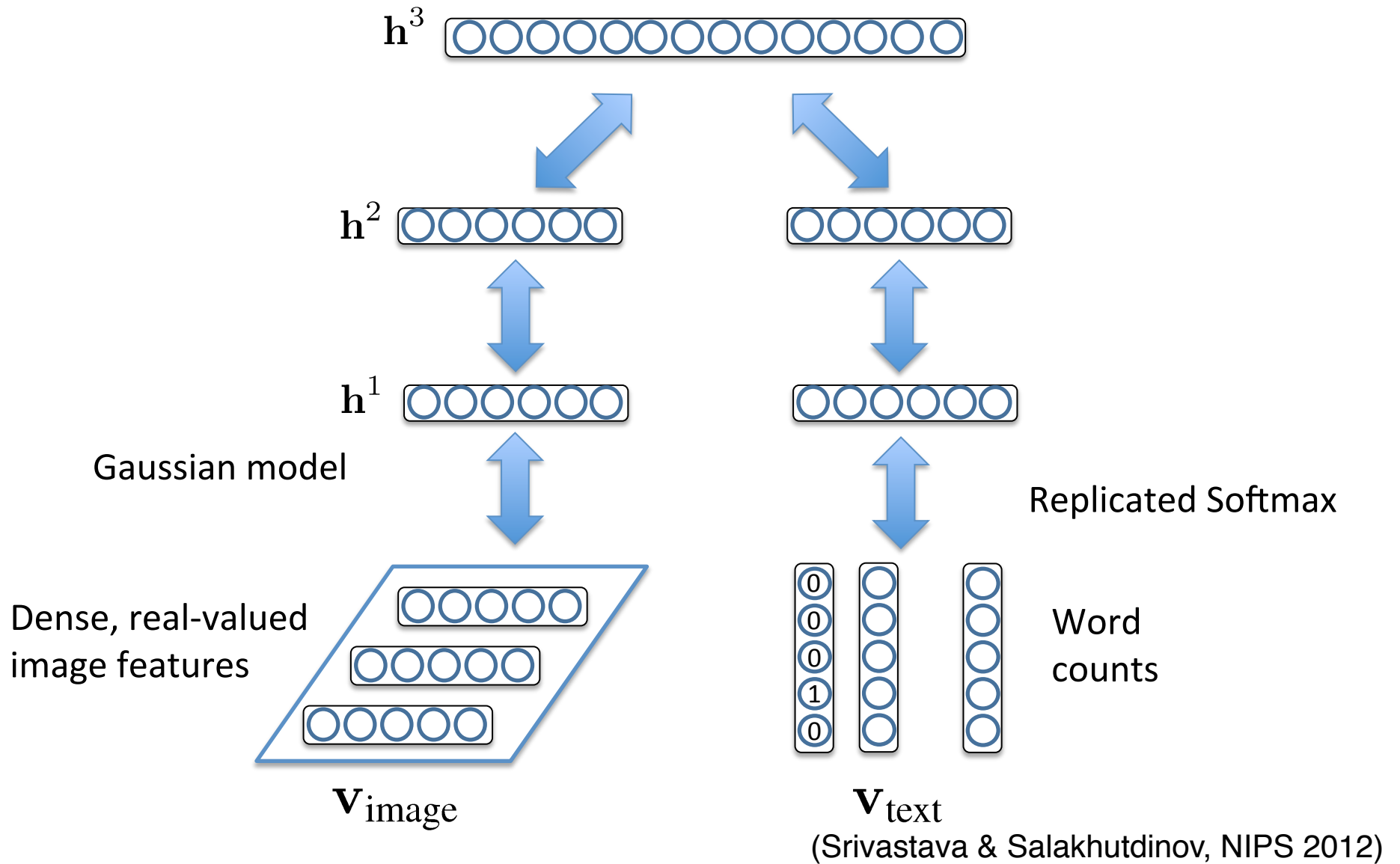
Multimodal DBM



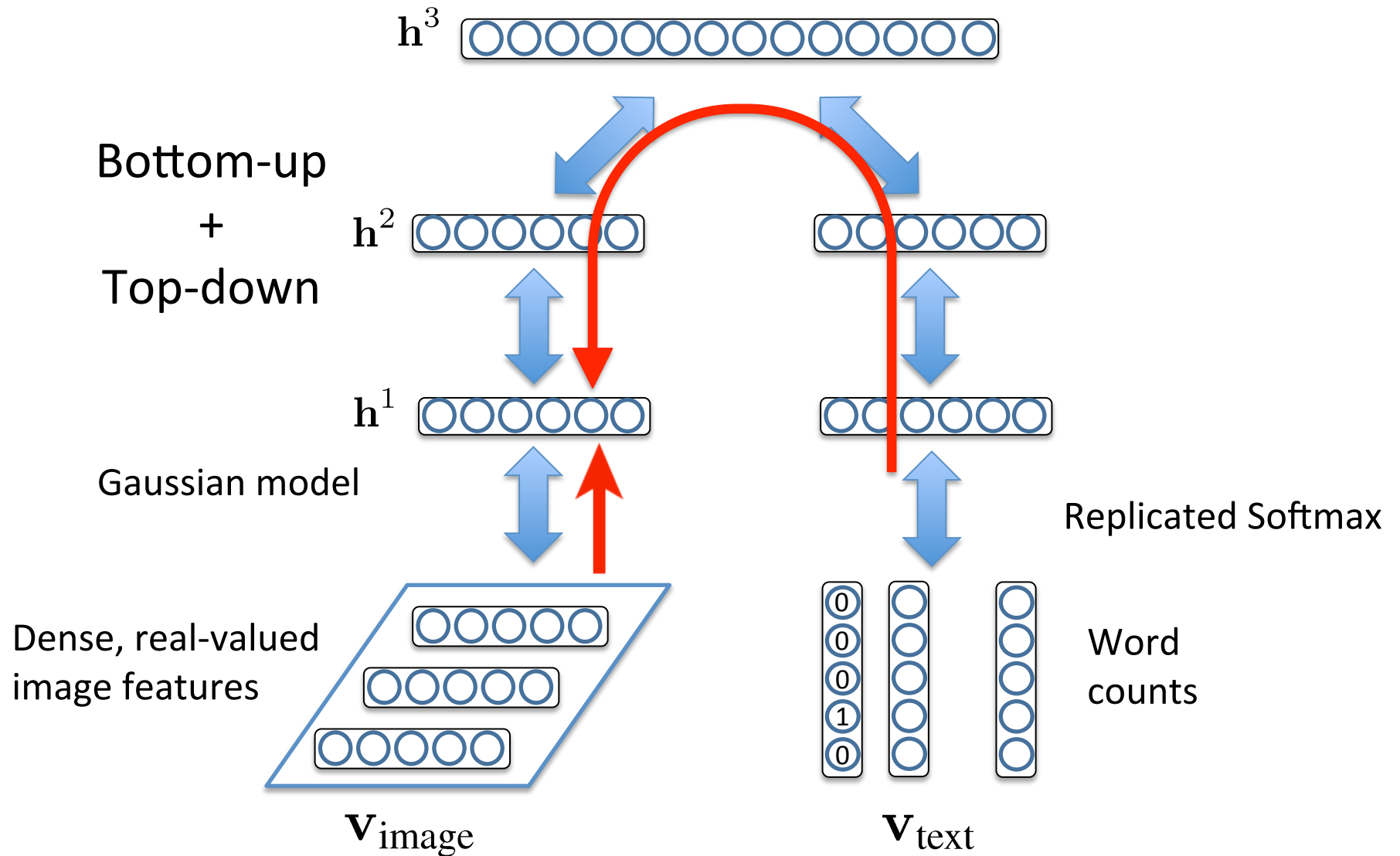
Multimodal DBM



Multimodal DBM



Multimodal DBM



Multimodal DBM

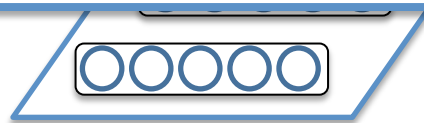


$$P(\mathbf{v}^m, \mathbf{v}^t; \theta) = \sum_{\mathbf{h}^{(2m)}, \mathbf{h}^{(2t)}, \mathbf{h}^{(3)}} P(\mathbf{h}^{(2m)}, \mathbf{h}^{(2t)}, \mathbf{h}^{(3)}) \left(\sum_{\mathbf{h}^{(1m)}} P(\mathbf{v}^m, \mathbf{h}^{(1m)} | \mathbf{h}^{(2m)}) \right) \left(\sum_{\mathbf{h}^{(1t)}} P(\mathbf{v}^t, \mathbf{h}^{(1t)} | \mathbf{h}^{(2t)}) \right)$$

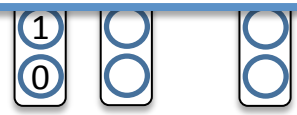
$$\frac{1}{Z(\theta, M)} \sum_{\mathbf{h}} \exp \left(\underbrace{- \sum_i \frac{(v_i^m)^2}{2\sigma_i^2} + \sum_{ij} \frac{v_i^m}{\sigma_i} W_{ij}^{(1m)} h_j^{(1m)} + \sum_{jl} W_{jl}^{(2m)} h_j^{(1m)} h_l^{(2m)}}_{\text{Gaussian Image Pathway}} \right)$$

$$\left(\underbrace{+ \sum_{jk} W_{kj}^{(1t)} h_j v_k^t + \sum_{jl} W_{jl}^{(2t)} h_j^{(1t)} h_l^{(2t)}}_{\text{Replicated Softmax Text Pathway}} + \underbrace{\sum_{lp} W^{(3t)} h_l^{(2t)} h_p^{(3)} + \sum_{lp} W^{(3m)} h_l^{(2m)} h_p^{(3)}}_{\text{Joint 3rd Layer}} \right)$$

image



$\mathbf{V}_{\text{image}}$



\mathbf{V}_{text}

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



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

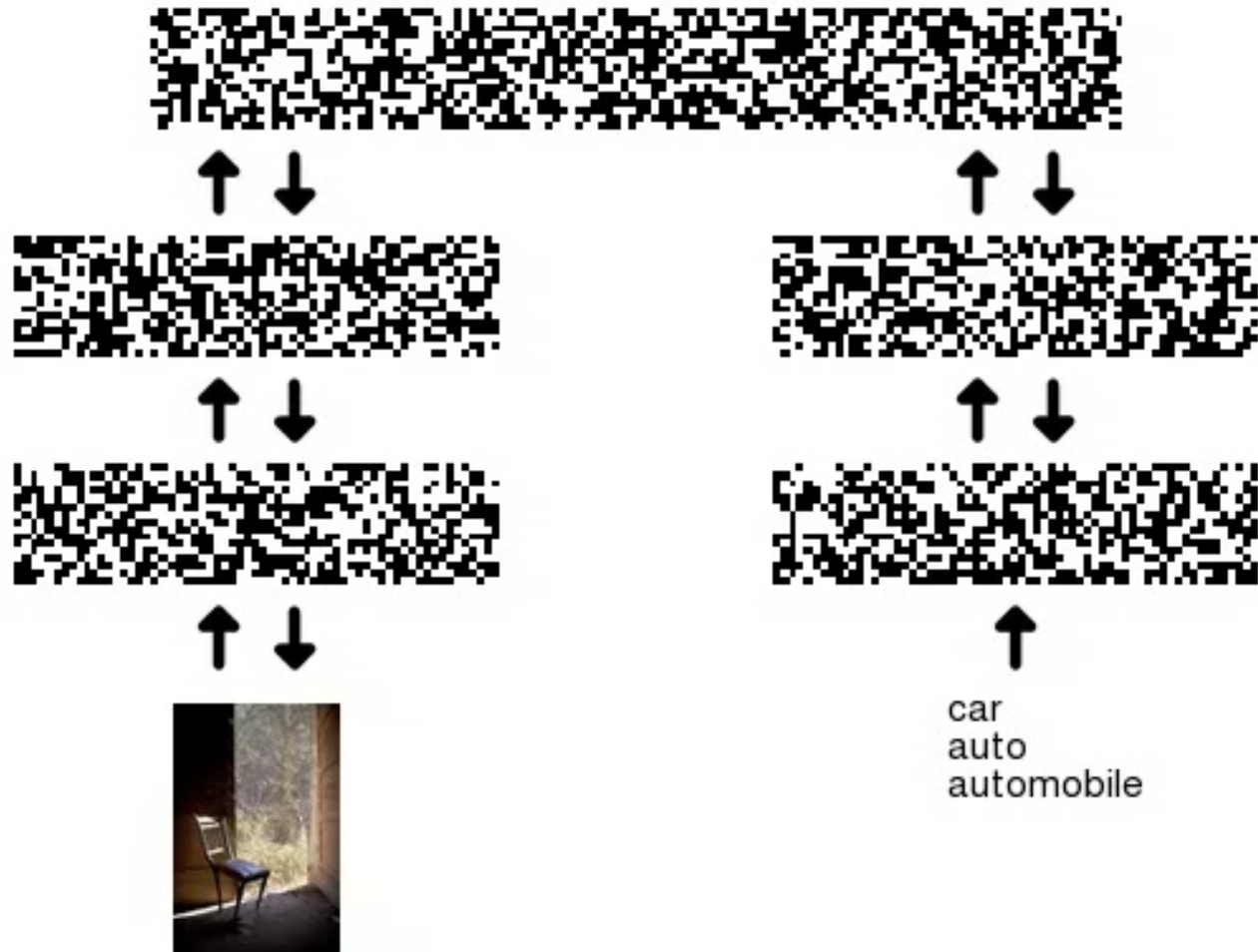
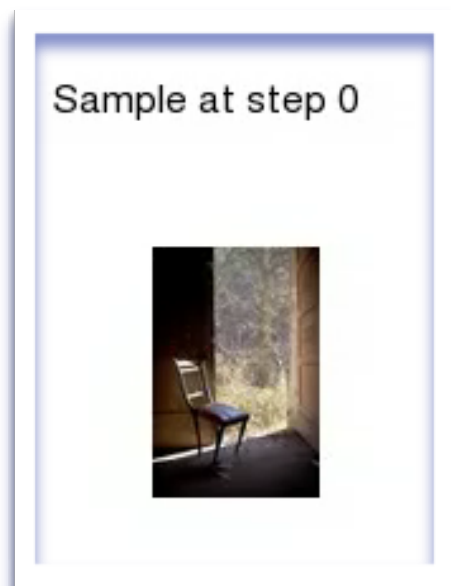


water, glass, beer, bottle,
drink, wine, bubbles, splash,
drops, drop

Images from Text

Step 0

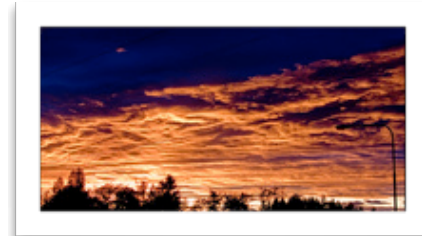
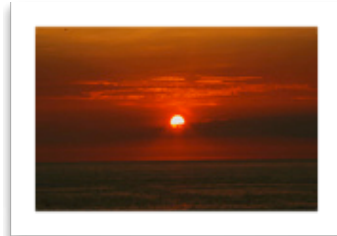
Sample drawn after every 50 steps of Gibbs sampling



Images from Text

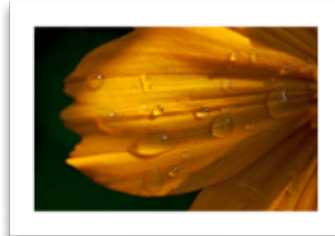
Given

water, red,
sunset

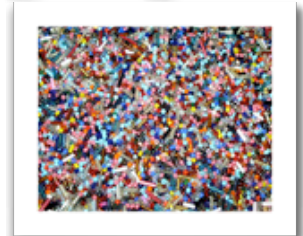


Retrieved

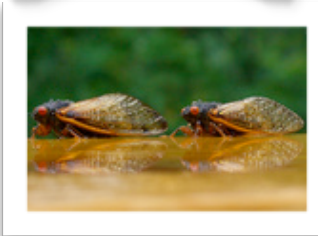
nature, flower,
red, green



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,
thepfectphotographer,
flash, damniwishidtakenthat,
spiritofphotography



nikon, green, light,
photoshop, apple, d70



white, yellow,
abstract, lines, bus,
graphic

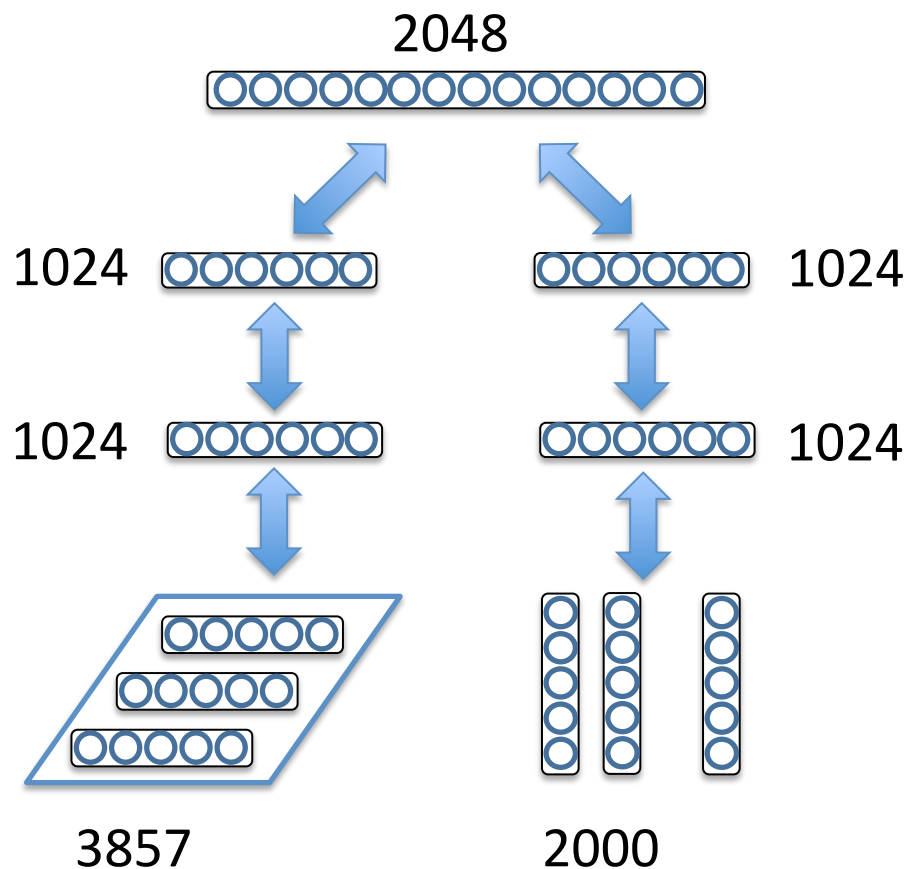


sky, geotagged,
reflection, cielo,
bilbao, reflejo

Huiskes et. al.

Data and Architecture

≈ 12 Million parameters



- 200 most frequent tags.
- 25K labeled subset (15K training, 10K testing)
- Additional 1 million unlabeled data
- 38 classes - *sky, tree, baby, car, cloud ...*

Results

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

Mean Average Precision



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

} Similar Features, 25K

Results

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

Mean Average Precision

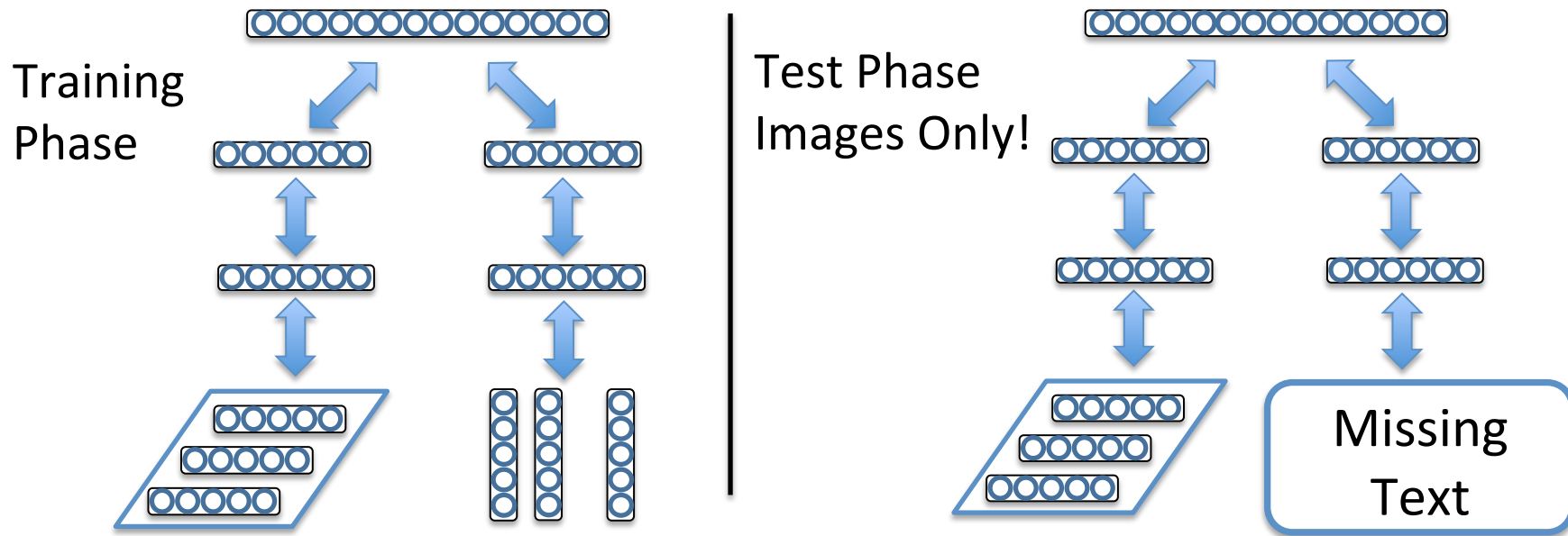


Learning Algorithm	MAP	Precision@50
Random	0.124	0.124
LDA [Huiskes et. al.]	0.492	0.754
SVM [Huiskes et. al.]	0.475	0.758
DBM-Labelled	0.526	0.791
DBM	0.609	0.863
Deep Belief Net	0.599	0.867
Autoencoder	0.600	0.875

} Similar Features, 25K

} + 1 Million Unlabelled

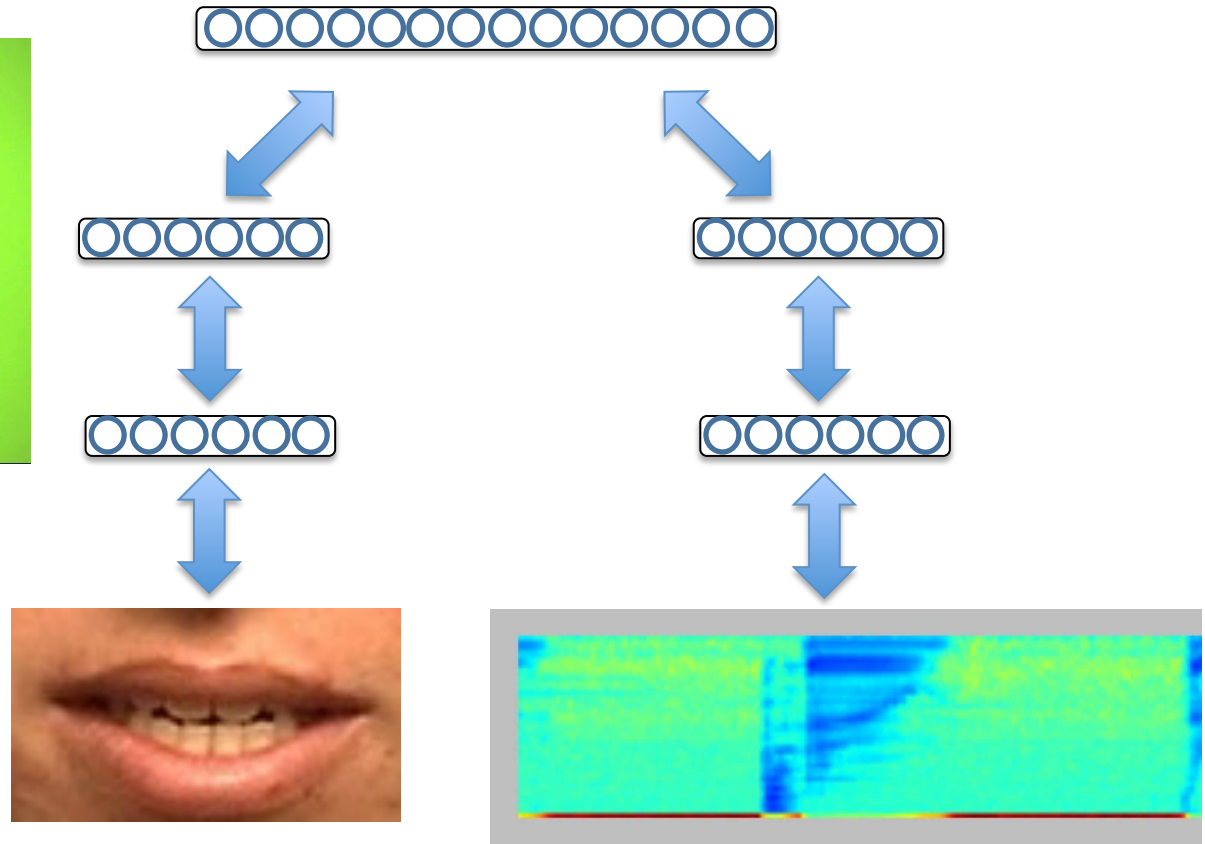
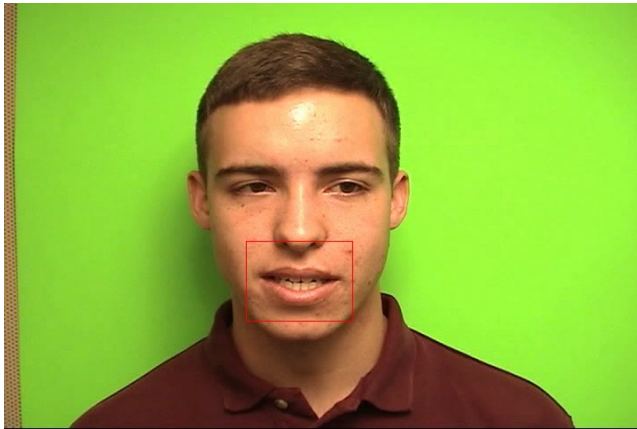
Benefits of using Multimodal Data



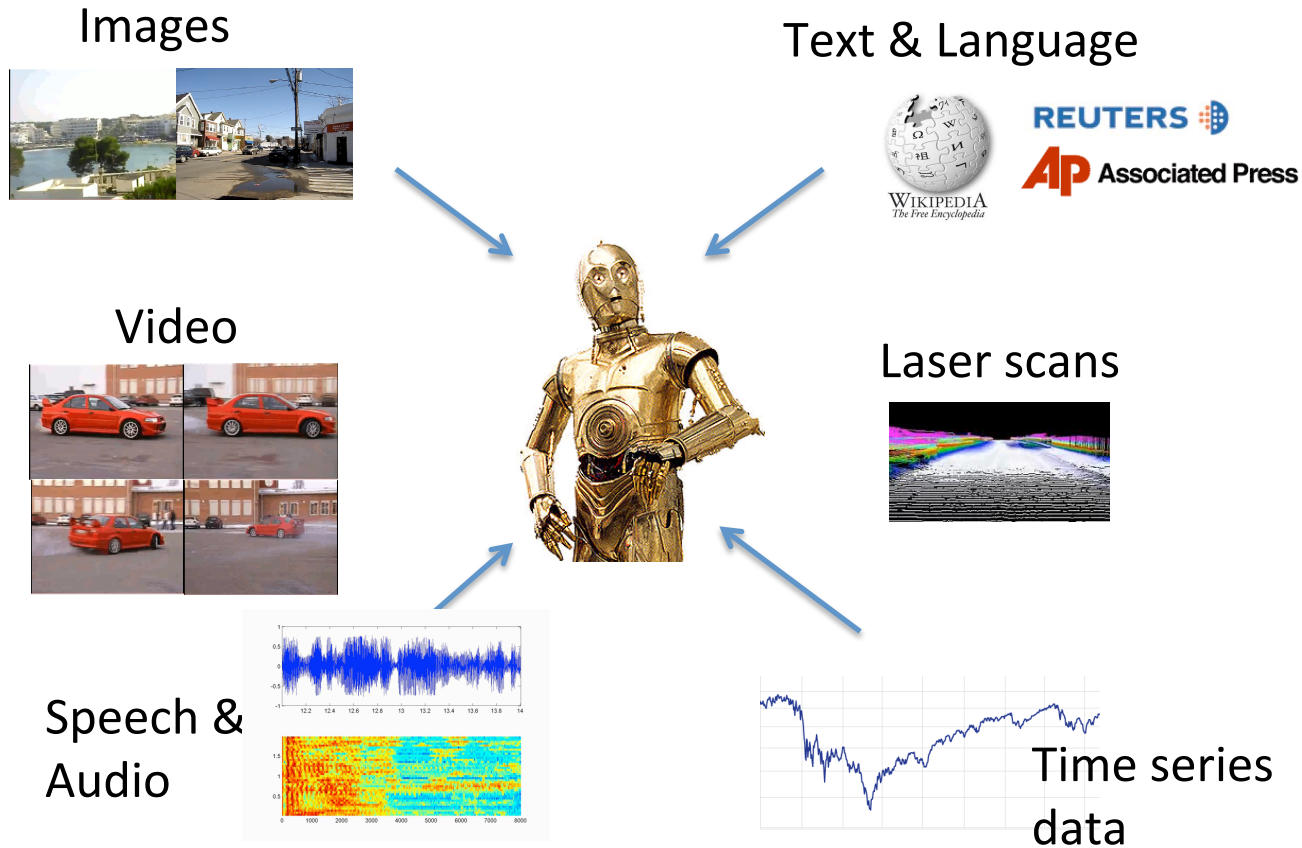
Learning Algorithm	MAP	Precision@50
Image-LDA [Huiskes et. al.]	0.315	-
Image-SVM [Huiskes et. al.]	0.375	-
Image-DBM	0.469	0.803
Multimodal-DBM (missing text)	0.531	0.832

Video and Audio

Cuave Dataset



Multi-Modal Models



Develop learning systems that come closer to displaying human like intelligence

One of Key Challenges:
Inference

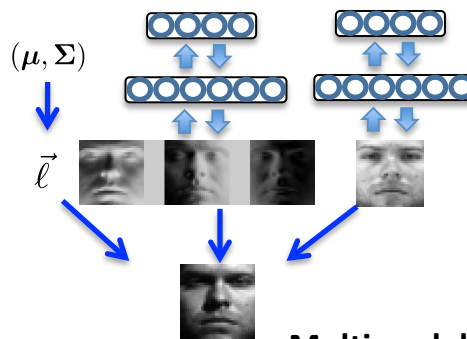
Summary

- Efficient learning algorithms for Hierarchical Generative Models. Learning more adaptive, robust, and structured representations.

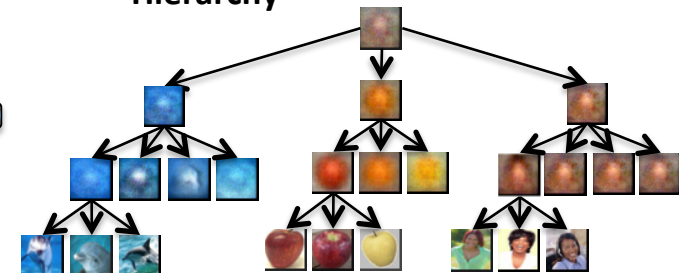
Text & image retrieval /
Object recognition



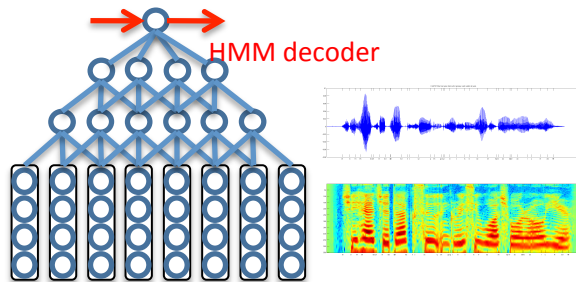
Dealing with missing/
occluded data



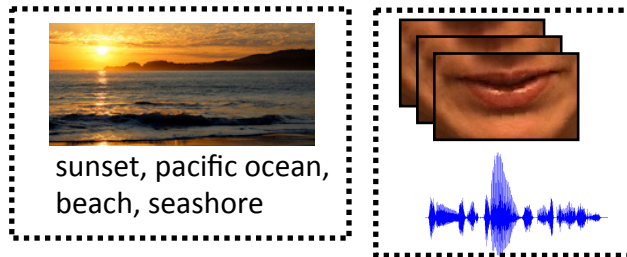
Learning a Category
Hierarchy



Speech Recognition



Multimodal Data



Object Detection



- Deep models can improve 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

Thanks to my collaborators:

Nitish Srivastava	University of Toronto
Charlie Tang	University of Toronto
Josh Tenenbaum	MIT
Geoffrey Hinton	University of Toronto
Nathan Srebro	TTI, University of Chicago
Roger Grosse	MIT
Ilya Sutskever	Google
Iain Murray	University of Edinburgh
Andriy Mnih	Gatsby Computational Neuroscience Unit, UCL
Hugo Larochelle	University of Toronto
Antonio Torralba	MIT
Bill Freeman	MIT
John Langford	Yahoo Research
Tong Zhang	Rutgers
Sham Kakade	University of Pennsylvania
Brenden Lake	MIT

Code for learning RBMs, DBNs, and DBMs is available at:

<http://www.utstat.toronto.edu/~rsalakhu/code.html>