

UCL Tutorial on:
Deep Belief Nets

(An updated and extended version of my 2007 NIPS tutorial)

Geoffrey Hinton

Canadian Institute for Advanced Research

&

Department of Computer Science

University of Toronto

Schedule for the Tutorial

- 2.00 – 3.30 Tutorial part 1
- 3.30 – 3.45 Questions
- 3.45 - 4.15 Tea Break
- 4.15 – 5.45 Tutorial part 2
- 5.45 – 6.00 Questions

Some things you will learn in this tutorial

- How to learn multi-layer generative models of unlabelled data by learning one layer of features at a time.
 - How to add Markov Random Fields in each hidden layer.
- How to use generative models to make discriminative training methods work much better for classification and regression.
 - How to extend this approach to Gaussian Processes and how to learn complex, domain-specific kernels for a Gaussian Process.
- How to perform non-linear dimensionality reduction on very large datasets
 - How to learn binary, low-dimensional codes and how to use them for very fast document retrieval.
- How to learn multilayer generative models of high-dimensional sequential data.

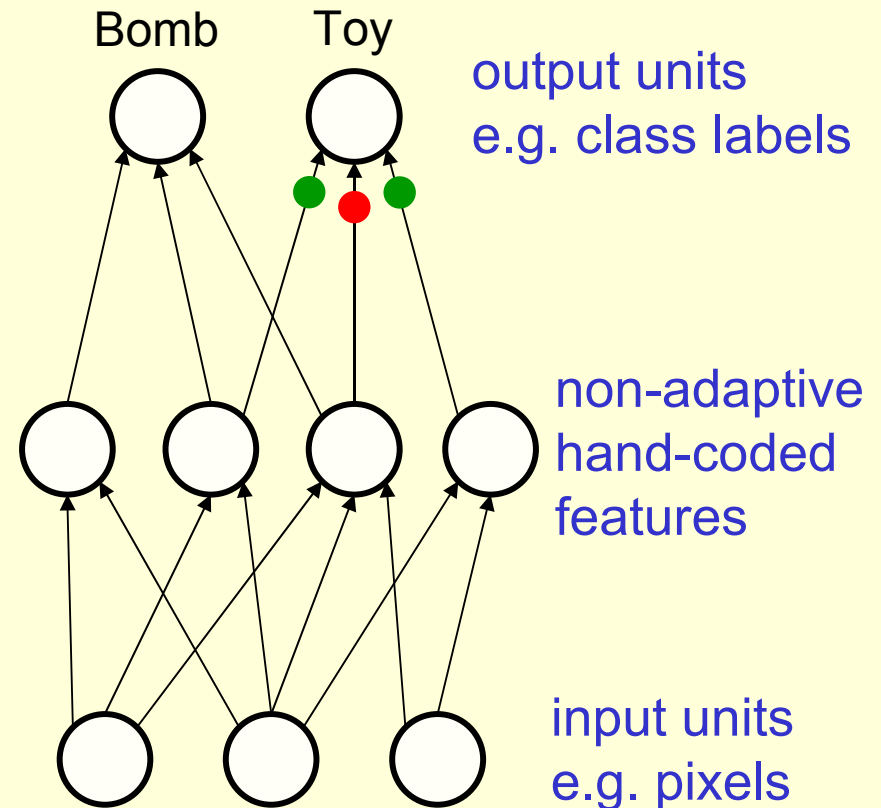
A spectrum of machine learning tasks

Typical Statistics-----Artificial Intelligence

- Low-dimensional data (e.g. less than 100 dimensions)
- Lots of noise in the data
- There is not much structure in the data, and what structure there is, can be represented by a fairly simple model.
- The main problem is distinguishing true structure from noise.
- High-dimensional data (e.g. more than 100 dimensions)
- The noise is not sufficient to obscure the structure in the data if we process it right.
- There is a huge amount of structure in the data, but the structure is too complicated to be represented by a simple model.
- The main problem is figuring out a way to represent the complicated structure so that it can be learned.

Historical background: First generation neural networks

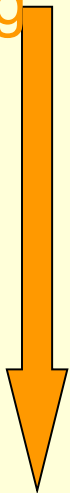
- Perceptrons (~1960) used a layer of hand-coded features and tried to recognize objects by learning how to weight these features.
 - There was a neat learning algorithm for adjusting the weights.
 - But perceptrons are fundamentally limited in what they can learn to do.



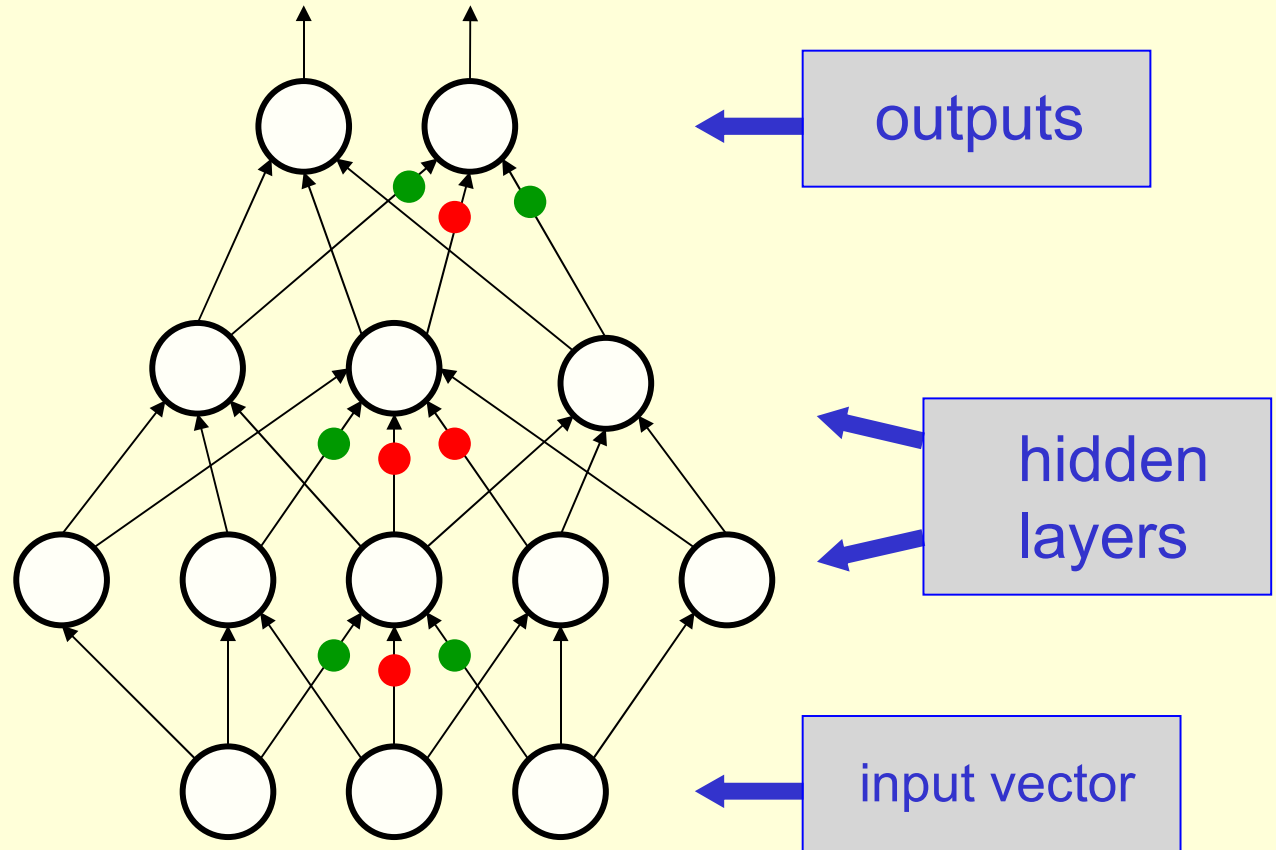
Sketch of a typical perceptron from the 1960's

Second generation neural networks (~1985)

Back-propagate
error
signal to get
derivatives for
learning



Compare outputs with
correct answer to get
error signal



A temporary digression

- Vapnik and his co-workers developed a very clever type of perceptron called a Support Vector Machine.
 - Instead of hand-coding the layer of non-adaptive features, each training example is used to create a new feature using a fixed recipe.
 - The feature computes how similar a test example is to that training example.
 - Then a clever optimization technique is used to select the best subset of the features and to decide how to weight each feature when classifying a test case.
 - But its just a perceptron and has all the same limitations.
- In the 1990's, many researchers abandoned neural networks with multiple adaptive hidden layers because Support Vector Machines worked better.

What is wrong with back-propagation?

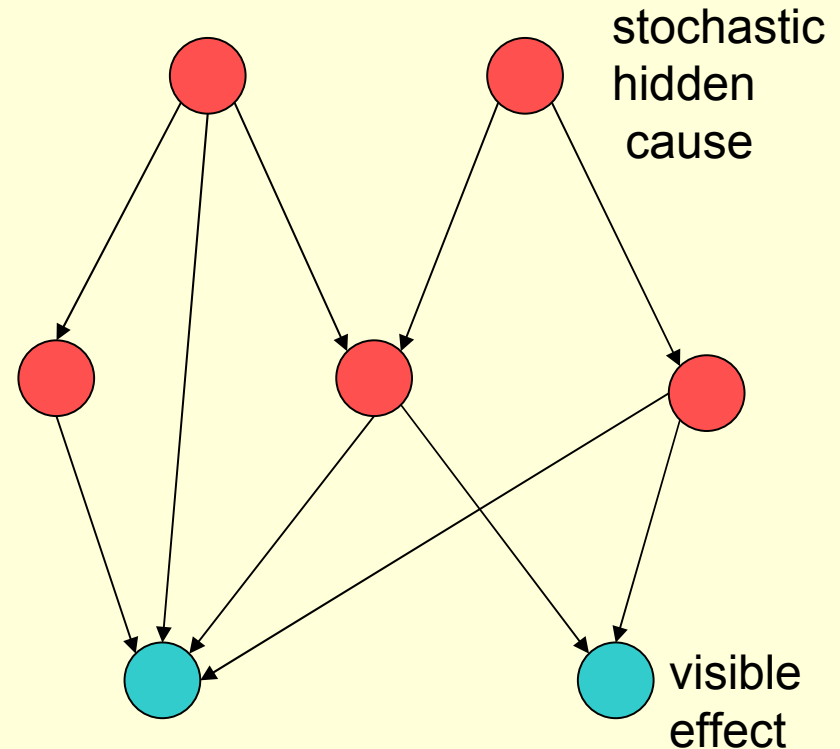
- It requires labeled training data.
 - Almost all data is unlabeled.
- The learning time does not scale well
 - It is very slow in networks with multiple hidden layers.
- It can get stuck in poor local optima.
 - These are often quite good, but for deep nets they are far from optimal.

Overcoming the limitations of back-propagation

- Keep the efficiency and simplicity of using a gradient method for adjusting the weights, but use it for modeling the structure of the sensory input.
 - Adjust the weights to maximize the probability that a generative model would have produced the sensory input.
 - Learn $p(\text{image})$ not $p(\text{label} | \text{image})$
 - If you want to do computer vision, first learn computer graphics
- What kind of generative model should we learn?

Belief Nets

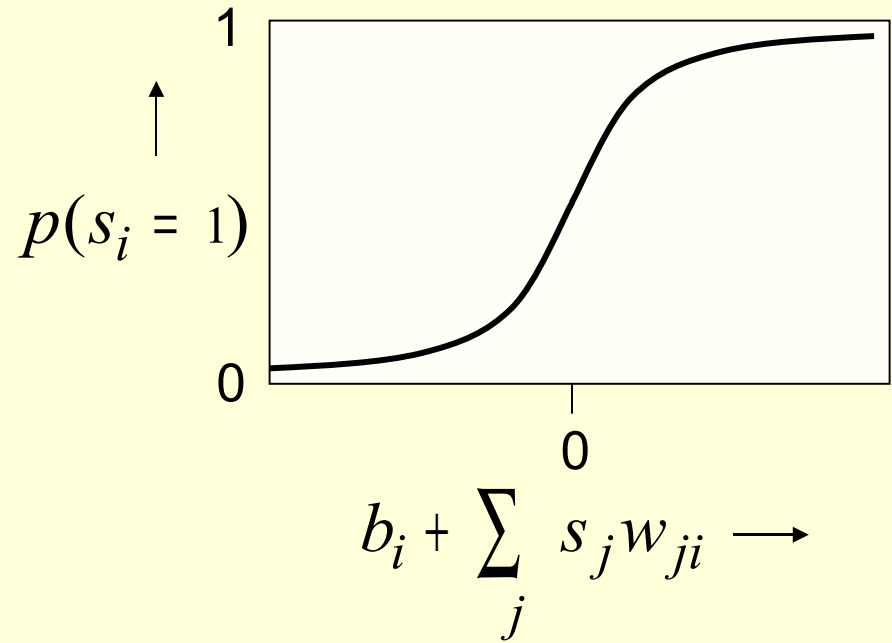
- A belief net is a directed acyclic graph composed of stochastic variables.
- We get to observe some of the variables and we would like to solve two problems:
 - **The inference problem:** Infer the states of the unobserved variables.
 - **The learning problem:** Adjust the interactions between variables to make the network more likely to generate the observed data.



We will use nets composed of layers of stochastic binary variables with weighted connections. Later, we will generalize to other types of variable.

Stochastic binary units (Bernoulli variables)

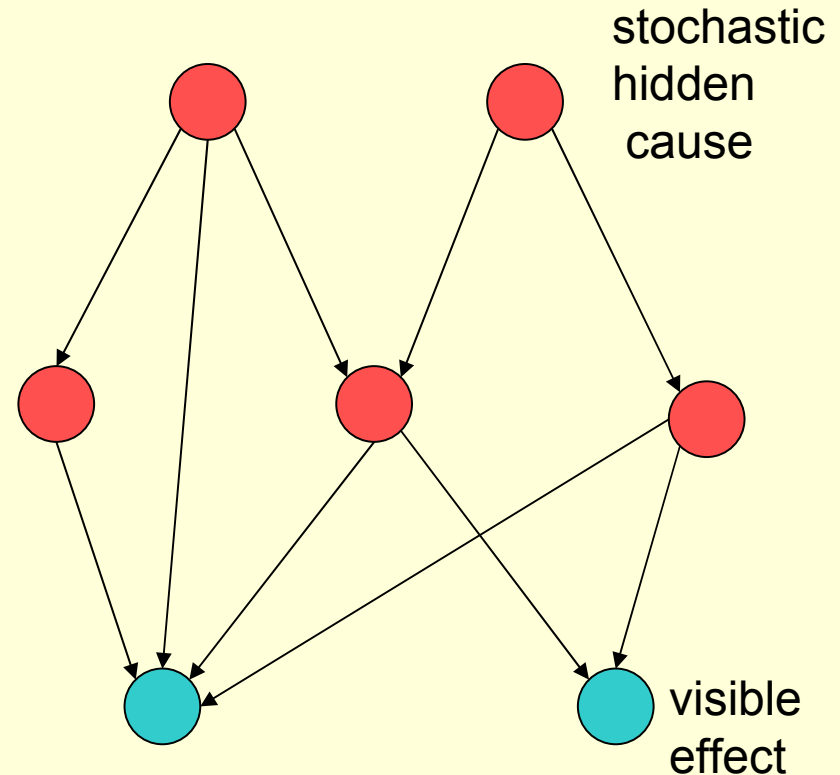
- These have a state of 1 or 0.
- The probability of turning on is determined by the weighted input from other units (plus a bias)



$$p(s_i = 1) = \frac{1}{1 + \exp(-b_i - \sum_j s_j w_{ji})}$$

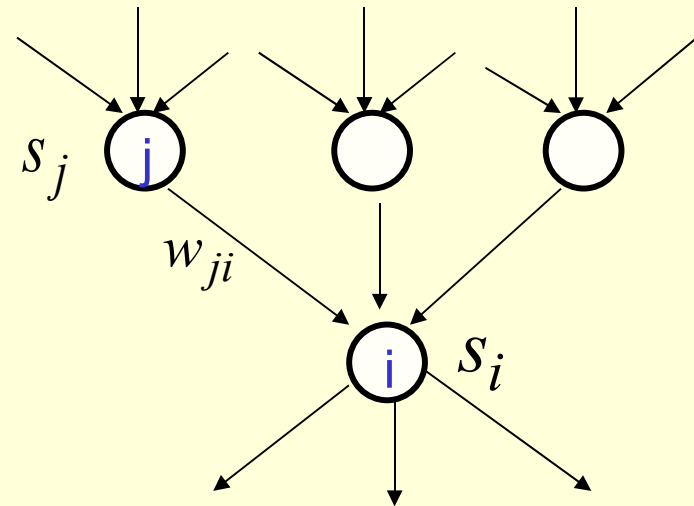
Learning Deep Belief Nets

- It is easy to generate an unbiased example at the leaf nodes, so we can see what kinds of data the network believes in.
- It is hard to infer the posterior distribution over all possible configurations of hidden causes.
- It is hard to even get a sample from the posterior.
- So how can we learn deep belief nets that have millions of parameters?



The learning rule for sigmoid belief nets

- Learning is easy if we can get an unbiased sample from the posterior distribution over hidden states given the observed data.
- For each unit, maximize the log probability that its binary state in the sample from the posterior would be generated by the sampled binary states of its parents.



$$p_i \equiv p(s_i = 1) = \frac{1}{1 + \exp(-\sum_j s_j w_{ji})}$$

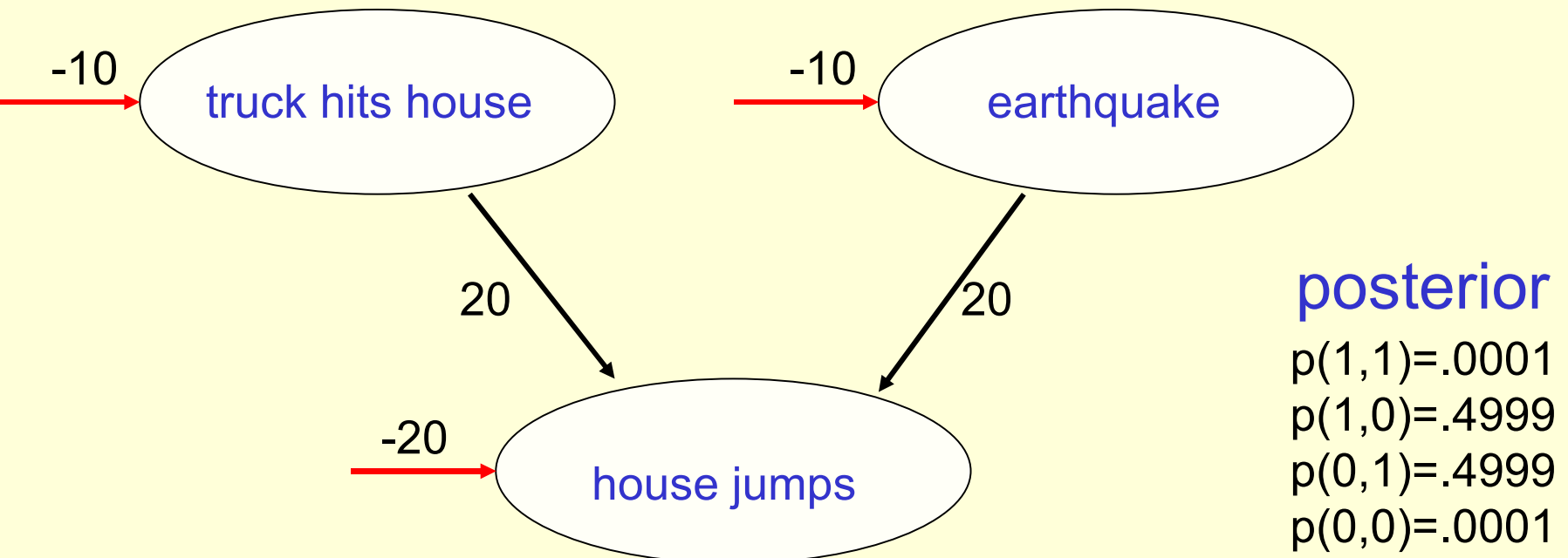
$$\Delta w_{ji} = \varepsilon s_j (s_i - p_i)$$



learning
rate

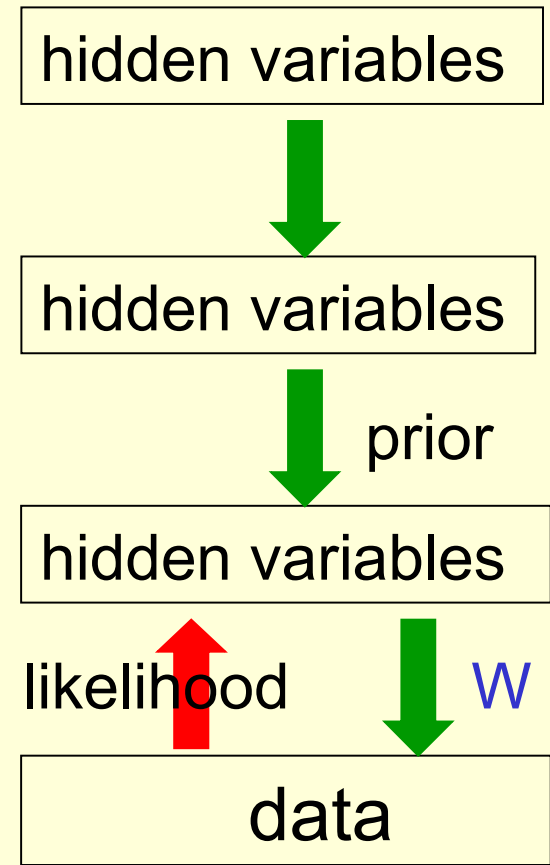
Explaining away (Judea Pearl)

- Even if two hidden causes are independent, they can become dependent when we observe an effect that they can both influence.
 - If we learn that there was an earthquake it reduces the probability that the house jumped because of a truck.



Why it is usually very hard to learn sigmoid belief nets one layer at a time

- To learn W , we need the posterior distribution in the first hidden layer.
- **Problem 1:** The posterior is typically complicated because of “explaining away”.
- **Problem 2:** The posterior depends on the prior as well as the likelihood.
 - So to learn W , we need to know the weights in higher layers, even if we are only approximating the posterior. All the weights interact.
- **Problem 3:** We need to integrate over all possible configurations of the higher variables to get the prior for first hidden layer. Yuk!



Some methods of learning deep belief nets

- Monte Carlo methods can be used to sample from the posterior.
 - But its painfully slow for large, deep models.
- In the 1990's people developed variational methods for learning deep belief nets
 - These only get approximate samples from the posterior.
 - Nevertheless, the learning is still guaranteed to improve a variational bound on the log probability of generating the observed data.

The breakthrough that makes deep learning efficient

- To learn deep nets efficiently, we need to learn one layer of features at a time. This does not work well if we assume that the latent variables are independent in the prior :
 - The latent variables are not independent in the posterior so inference is hard for non-linear models.
 - The learning tries to find independent causes using one hidden layer which is not usually possible.
- We need a way of learning one layer at a time that takes into account the fact that we will be learning more hidden layers later.
 - We solve this problem by using an undirected model.

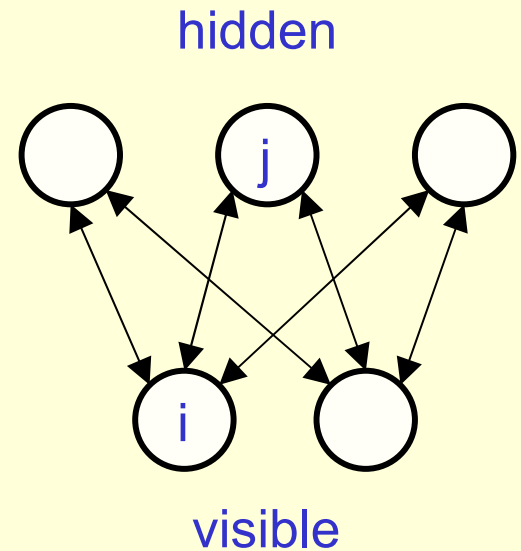
Two types of generative neural network

- If we connect binary stochastic neurons in a directed acyclic graph we get a Sigmoid Belief Net (Radford Neal 1992).
- If we connect binary stochastic neurons using symmetric connections we get a Boltzmann Machine (Hinton & Sejnowski, 1983).
 - If we restrict the connectivity in a special way, it is easy to learn a Boltzmann machine.

Restricted Boltzmann Machines

(Smolensky ,1986, called them “harmoniums”)

- We restrict the connectivity to make learning easier.
 - Only one layer of hidden units.
 - We will deal with more layers later
 - No connections between hidden units.
- In an RBM, the hidden units are conditionally independent given the visible states.
 - So we can quickly get an unbiased sample from the posterior distribution when given a data-vector.
 - This is a big advantage over directed belief nets



The Energy of a joint configuration

(ignoring terms to do with biases)

$$E(v, h) = - \sum_{i, j} v_i h_j w_{ij}$$

binary state of visible unit i binary state of hidden unit j

Energy with configuration v on the visible units and h on the hidden units

weight between units i and j

$$- \frac{\partial E(v, h)}{\partial w_{ij}} = v_i h_j$$

Weights \rightarrow Energies \rightarrow Probabilities

- Each possible joint configuration of the visible and hidden units has an energy
 - The energy is determined by the weights and biases (as in a Hopfield net).
- The energy of a joint configuration of the visible and hidden units determines its probability:

$$p(v, h) \propto e^{-E(v, h)}$$

- The probability of a configuration over the visible units is found by summing the probabilities of all the joint configurations that contain it.

Using energies to define probabilities

- The probability of a joint configuration over both visible and hidden units depends on the energy of that joint configuration compared with the energy of all other joint configurations.

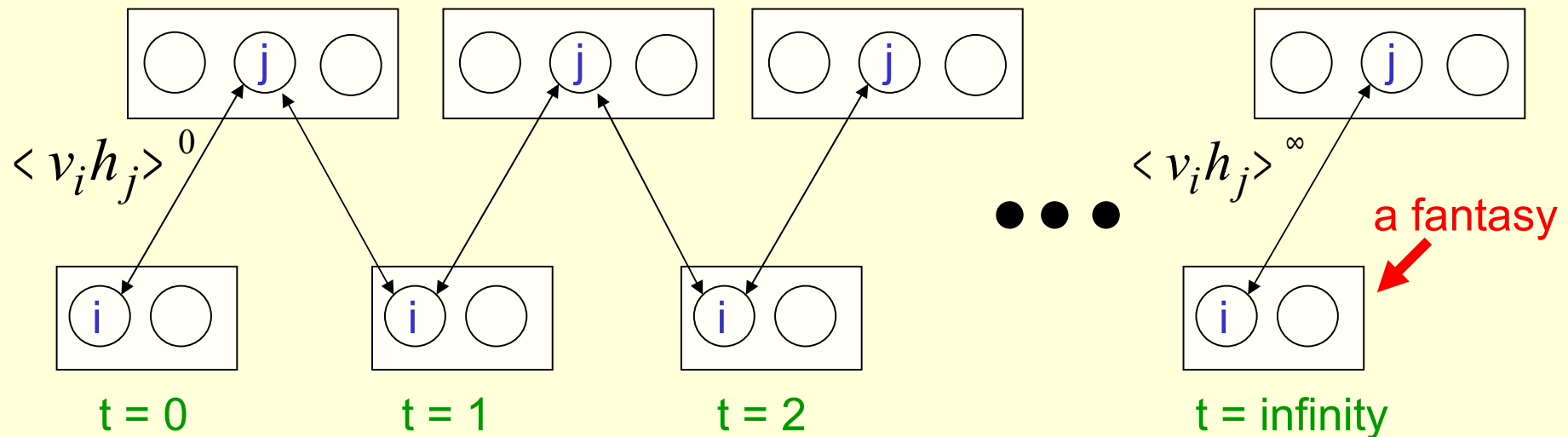
$$p(v, h) = \frac{e^{-E(v, h)}}{\sum_{u, g} e^{-E(u, g)}}$$

partition function

- The probability of a configuration of the visible units is the sum of the probabilities of all the joint configurations that contain it.

$$p(v) = \frac{\sum_h e^{-E(v, h)}}{\sum_{u, g} e^{-E(u, g)}}$$

A picture of the maximum likelihood learning algorithm for an RBM

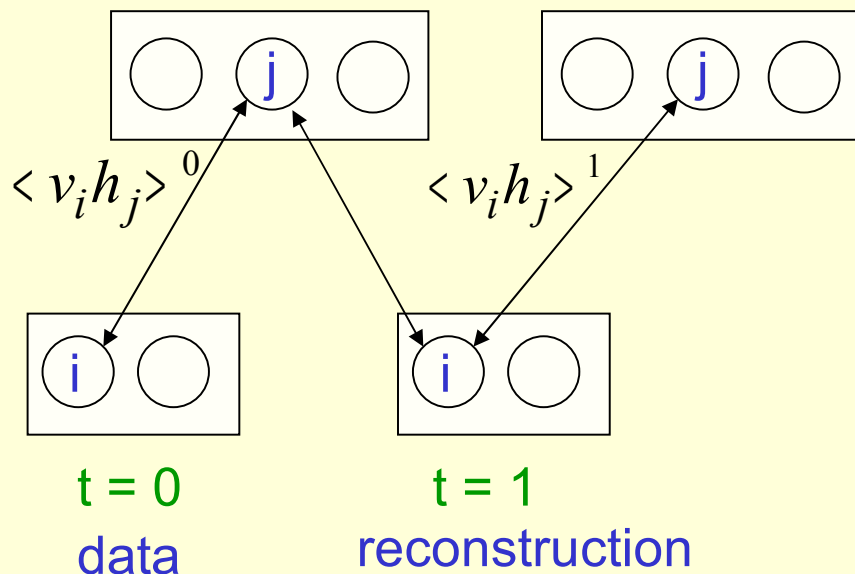


Start with a training vector on the visible units.

Then alternate between updating all the hidden units in parallel and updating all the visible units in parallel.

$$\frac{\partial \log p(v)}{\partial w_{ij}} = \langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^\infty$$

A quick way to learn an RBM



Start with a training vector on the visible units.

Update all the hidden units in parallel

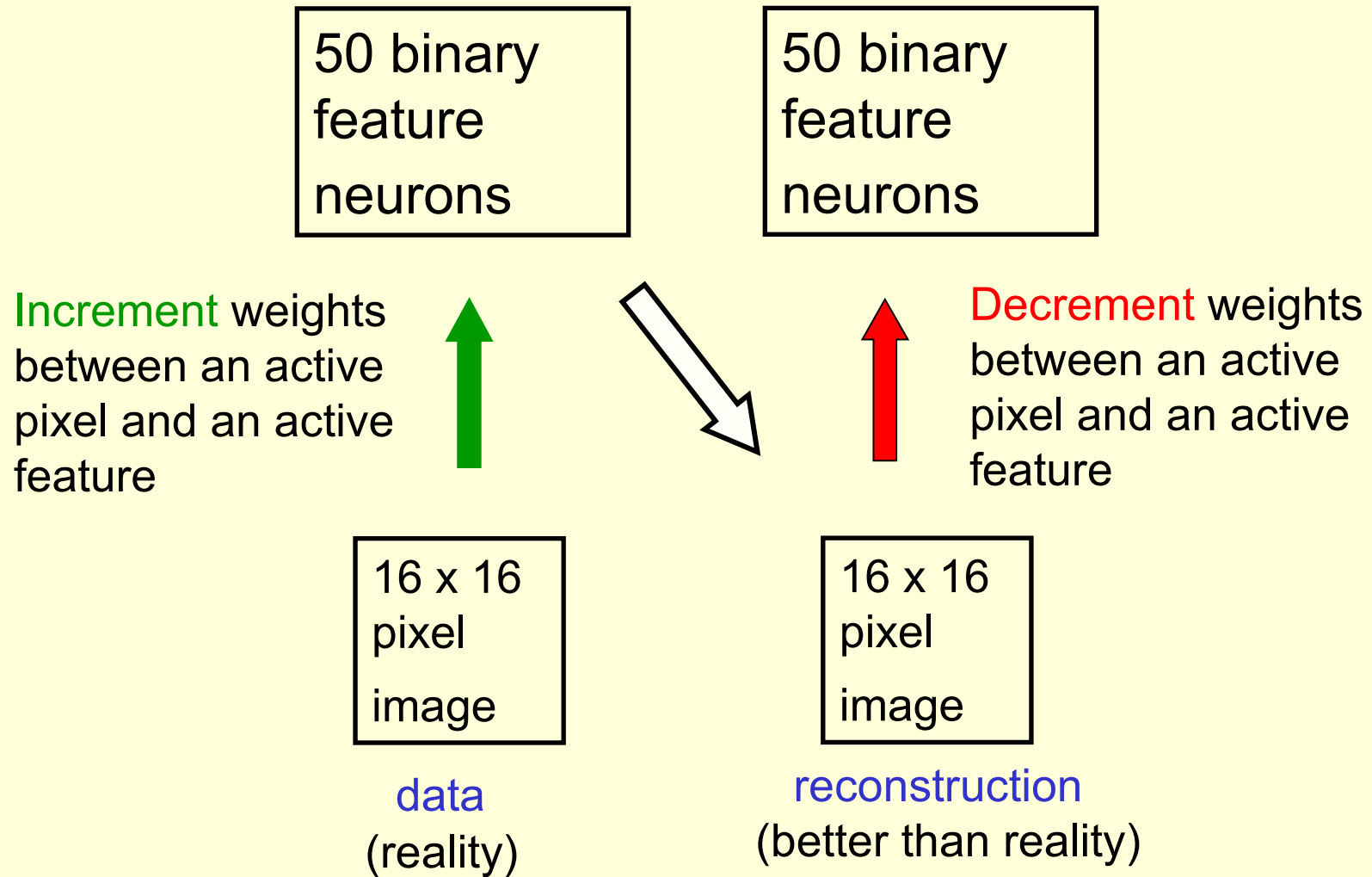
Update the all the visible units in parallel to get a “reconstruction”.

Update the hidden units again.

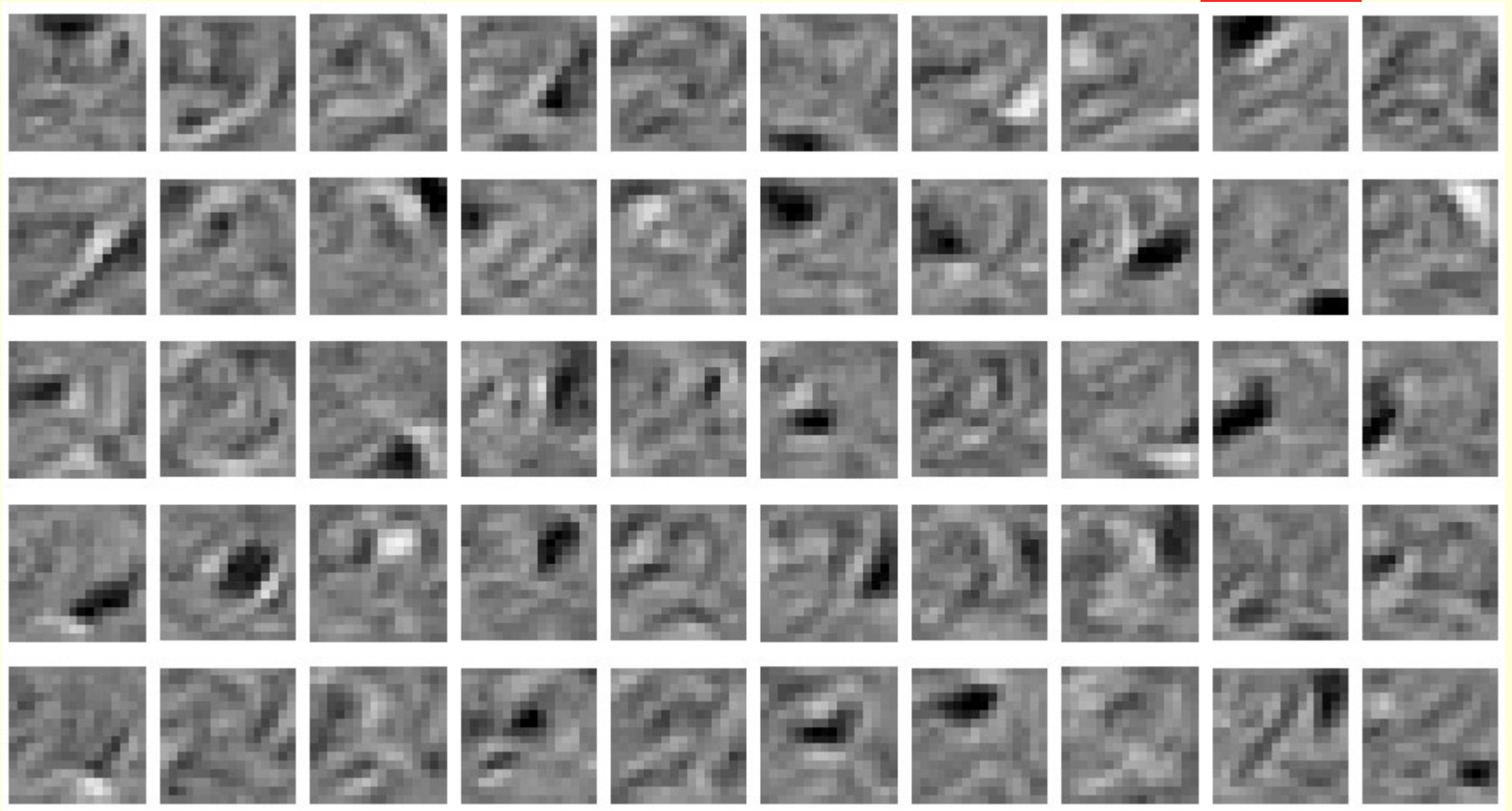
$$\Delta w_{ij} = \varepsilon (\langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^1)$$

This is not following the gradient of the log likelihood. But it works well. It is approximately following the gradient of another objective function (Carreira-Perpinan & Hinton, 2005).

How to learn a set of features that are good for reconstructing images of the digit 2

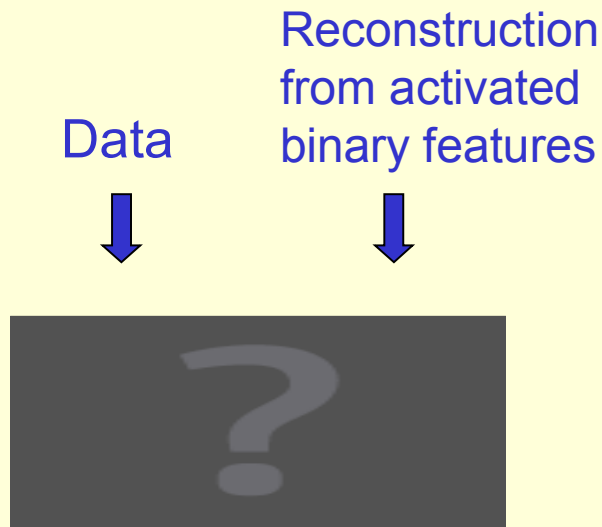


The final 50 x 256 weights

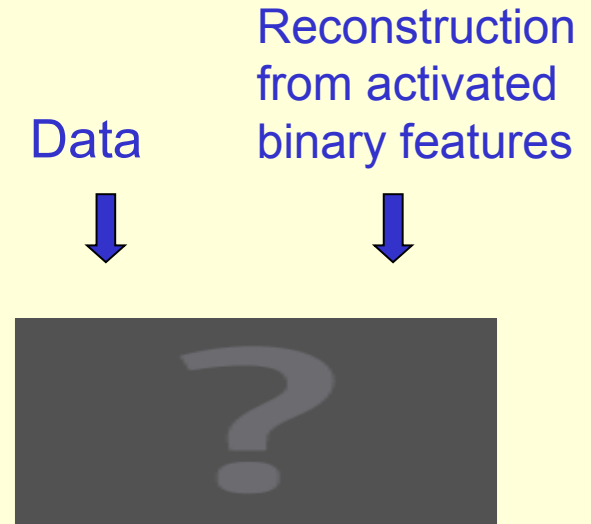


Each neuron grabs a different feature.

How well can we reconstruct the digit images from the binary feature activations?



New test images from the digit class that the model was trained on



Images from an unfamiliar digit class (the network tries to see every image as a 2)

Three ways to combine probability density models (an underlying theme of the tutorial)

- **Mixture:** Take a weighted average of the distributions.
 - It can never be sharper than the individual distributions. It's a very weak way to combine models.
- **Product:** Multiply the distributions at each point and then renormalize (this is how an RBM combines the distributions defined by each hidden unit)
 - Exponentially more powerful than a mixture. The normalization makes maximum likelihood learning difficult, but approximations allow us to learn anyway.
- **Composition:** Use the values of the latent variables of one model as the data for the next model.
 - Works well for learning multiple layers of representation, but only if the individual models are undirected.

Training a deep network

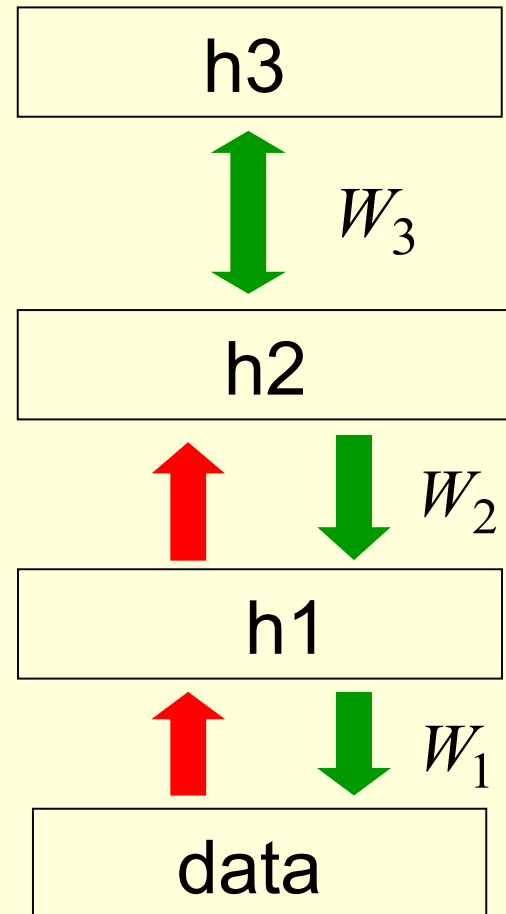
(the main reason RBM's are interesting)

- First train a layer of features that receive input directly from the pixels.
- Then treat the activations of the trained features as if they were pixels and learn features of features in a second hidden layer.
- It can be proved that each time we add another layer of features we improve a variational lower bound on the log probability of the training data.
 - The proof is slightly complicated.
 - But it is based on a neat equivalence between an RBM and a deep directed model (described later)

The generative model after learning 3 layers

- To generate data:
 1. Get an equilibrium sample from the top-level RBM by performing alternating Gibbs sampling for a long time.
 2. Perform a top-down pass to get states for all the other layers.

So the lower level bottom-up connections are **not** part of the generative model. They are just used for inference.



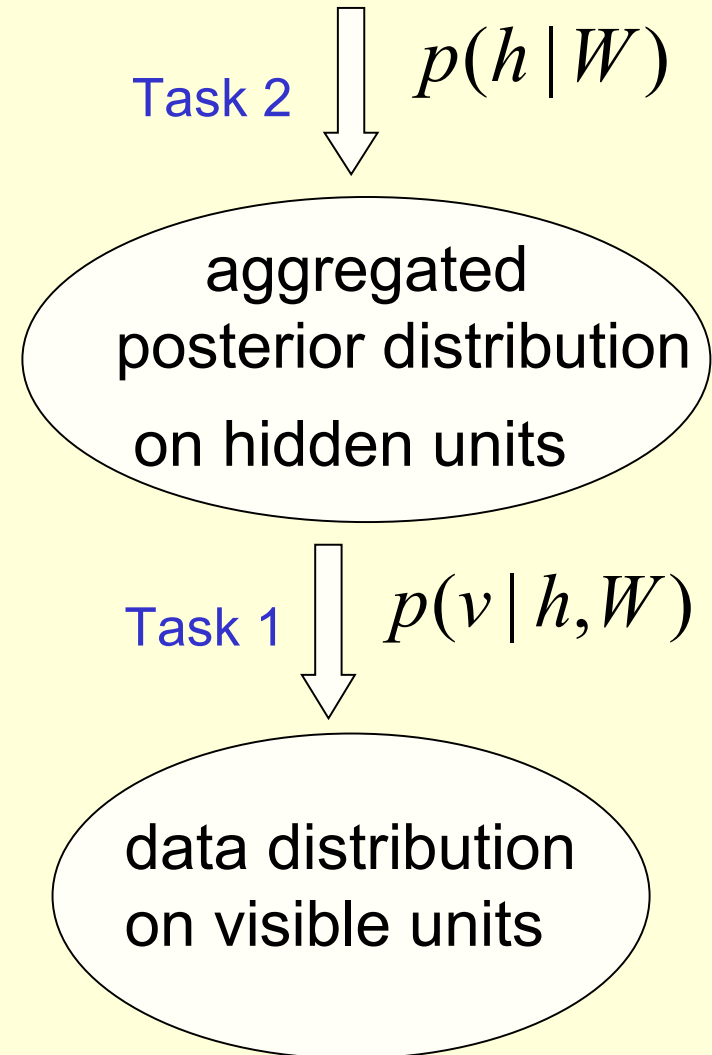
Why does greedy learning work?

An aside: Averaging factorial distributions

- If you average some factorial distributions, you do NOT get a factorial distribution.
 - In an RBM, the posterior over the hidden units is factorial for each visible vector.
 - But the aggregated posterior over all training cases is not factorial (even if the data was generated by the RBM itself).

Why does greedy learning work?

- Each RBM converts its data distribution into an aggregated posterior distribution over its hidden units.
- This divides the task of modeling its data into two tasks:
 - **Task 1:** Learn generative weights that can convert the aggregated posterior distribution over the hidden units back into the data distribution.
 - **Task 2:** Learn to model the aggregated posterior distribution over the hidden units.
 - The RBM does a good job of task 1 and a moderately good job of task 2.
- Task 2 is easier (for the next RBM) than modeling the original data because the aggregated posterior distribution is closer to a distribution that an RBM can model perfectly.



Why does greedy learning work?

The weights, W , in the bottom level RBM define $p(v|h)$ and they also, indirectly, define $p(h)$.

So we can express the RBM model as

$$p(v) = \sum_h p(h) p(v | h)$$

If we leave $p(v|h)$ alone and improve $p(h)$, we will improve $p(v)$.

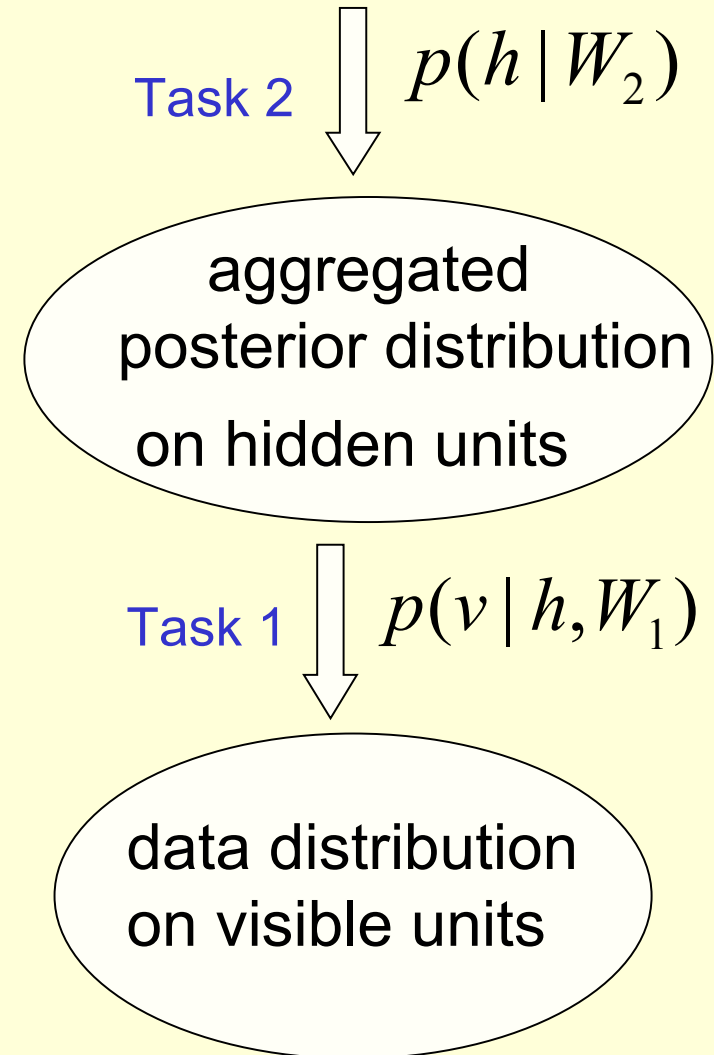
To improve $p(h)$, we need it to be a better model of the **aggregated posterior** distribution over hidden vectors produced by applying W to the data.

Which distributions are factorial in a directed belief net?

- In a directed belief net with one hidden layer, the posterior over the hidden units $p(h|v)$ is non-factorial (due to explaining away).
 - The aggregated posterior is factorial if the data was generated by the directed model.
 - It's the opposite way round from an undirected model which has factorial posteriors and a non-factorial prior $p(h)$ over the hiddens.
 - The intuitions that people have from using directed models are very misleading for undirected models.

Why does greedy learning fail in a directed module?

- A directed module also converts its data distribution into an aggregated posterior
 - Task 1 The learning is now harder because the posterior for each training case is non-factorial.
 - Task 2 is performed using an independent prior. This is a very bad approximation unless the aggregated posterior is close to factorial.
- A directed module attempts to make the aggregated posterior factorial in one step.
 - This is too difficult and leads to a bad compromise. There is also no guarantee that the aggregated posterior is easier to model than the data distribution.



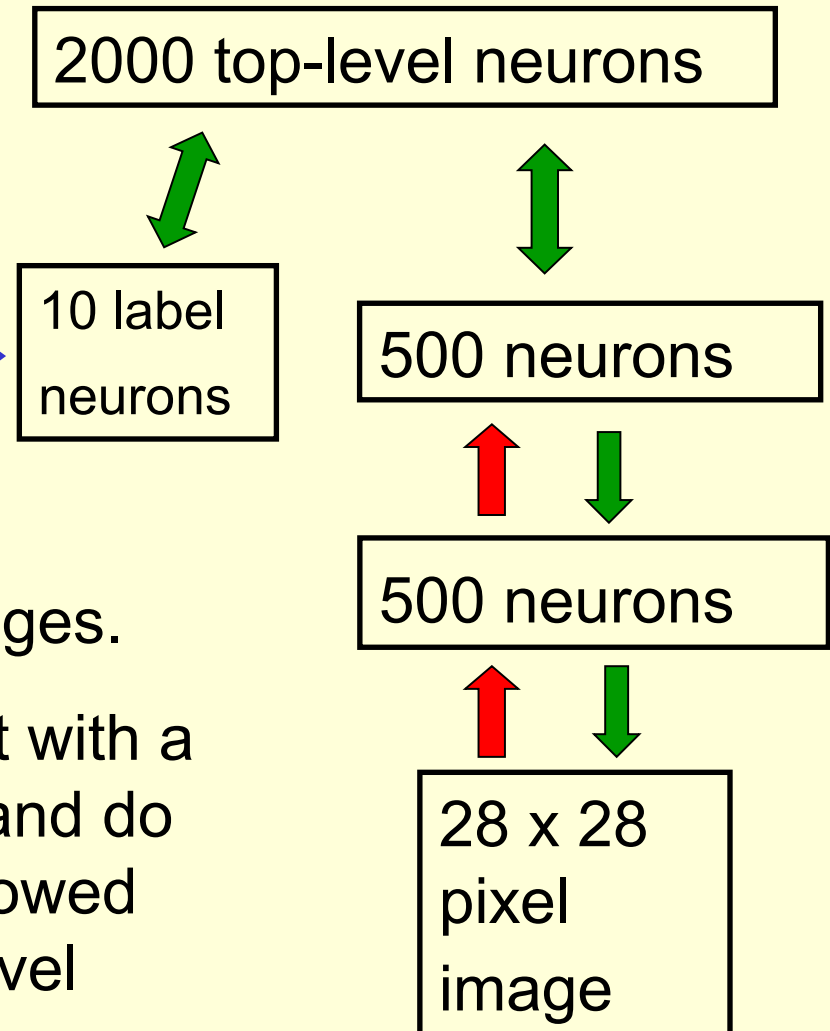
A model of digit recognition

The top two layers form an associative memory whose energy landscape models the low dimensional manifolds of the digits.

The energy valleys have names →

The model learns to generate combinations of labels and images.

To perform recognition we start with a neutral state of the label units and do an up-pass from the image followed by a few iterations of the top-level associative memory.



Fine-tuning with a contrastive version of the “wake-sleep” algorithm

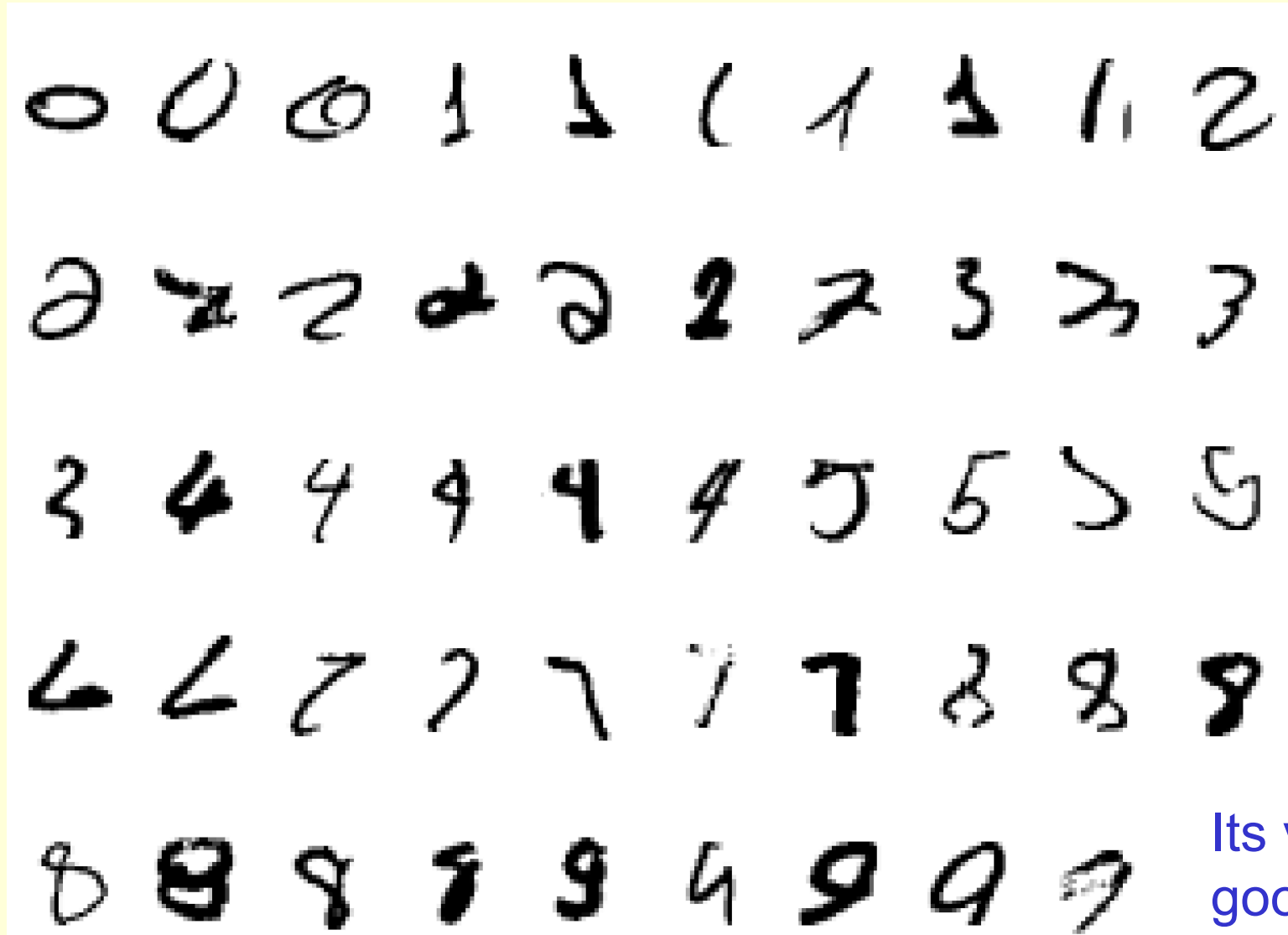
After learning many layers of features, we can fine-tune the features to improve generation.

1. Do a stochastic bottom-up pass
 - Adjust the top-down weights to be good at reconstructing the feature activities in the layer below.
3. Do a few iterations of sampling in the top level RBM
 - Adjust the weights in the top-level RBM.
4. Do a stochastic top-down pass
 - Adjust the bottom-up weights to be good at reconstructing the feature activities in the layer above.

Show the movie of the network
generating digits

(available at www.cs.toronto/~hinton)

Examples of correctly recognized handwritten digits that the neural network had never seen before



Its very good

How well does it discriminate on MNIST test set with no extra information about geometric distortions?

- Generative model based on RBM's 1.25%
- Support Vector Machine (Decoste et. al.) 1.4%
- Backprop with 1000 hiddens (Platt) ~1.6%
- Backprop with 500 -->300 hiddens ~1.6%
- K-Nearest Neighbor ~ 3.3%
- See Le Cun et. al. 1998 for more results
- Its better than backprop and much more neurally plausible because the neurons only need to send one kind of signal, and the teacher can be another sensory input.

Unsupervised “pre-training” also helps for models that have more data and better priors

- Ranzato et. al. (NIPS 2006) used an additional 600,000 distorted digits.
- They also used convolutional multilayer neural networks that have some built-in, local translational invariance.

Back-propagation alone: 0.49%

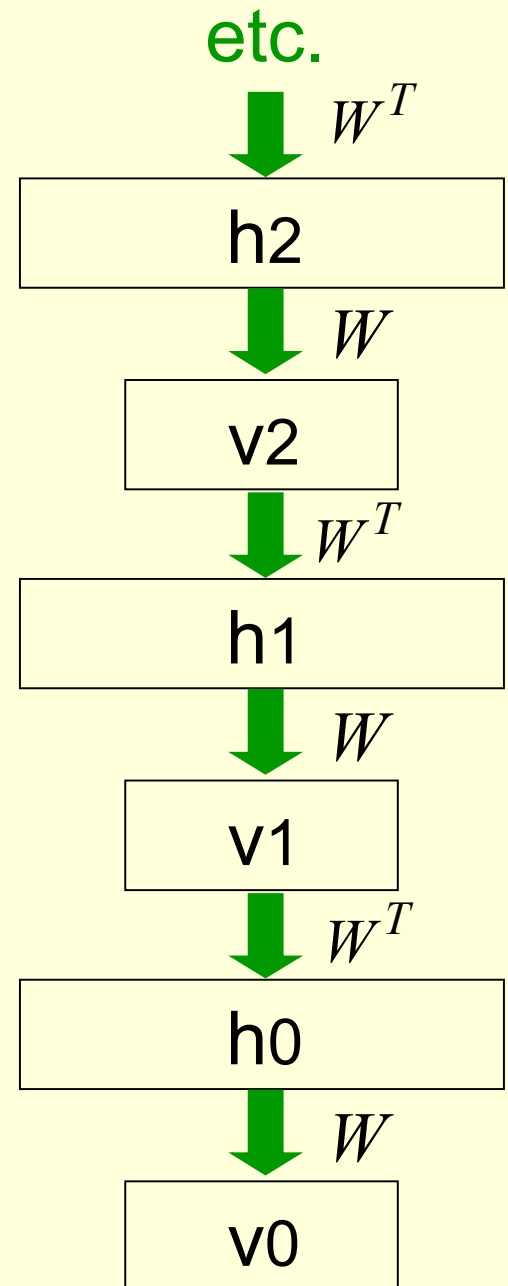
Unsupervised layer-by-layer
pre-training followed by backprop: 0.39% (record)

Another view of why layer-by-layer learning works (Hinton, Osindero & Teh 2006)

- There is an unexpected equivalence between RBM's and directed networks with many layers that all use the same weights.
 - This equivalence also gives insight into why contrastive divergence learning works.

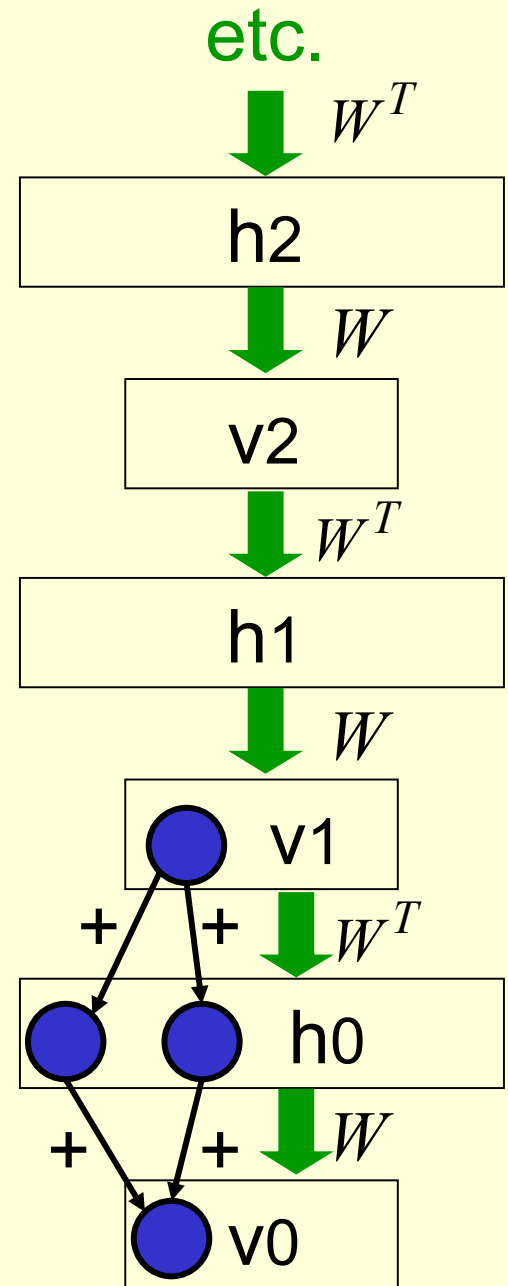
An infinite sigmoid belief net that is equivalent to an RBM

- The distribution generated by this infinite directed net with replicated weights is the equilibrium distribution for a compatible pair of conditional distributions: $p(v|h)$ and $p(h|v)$ that are both defined by W
 - A top-down pass of the directed net is exactly equivalent to letting a Restricted Boltzmann Machine settle to equilibrium.
 - So this infinite directed net defines the same distribution as an RBM.



Inference in a directed net with replicated weights

- The variables in h_0 are conditionally independent given v_0 .
 - Inference is trivial. We just multiply v_0 by W transpose.
 - The model above h_0 implements a complementary prior.
 - Multiplying v_0 by W transpose gives the product of the likelihood term and the prior term.
- Inference in the directed net is exactly equivalent to letting a Restricted Boltzmann Machine settle to equilibrium starting at the data.



- The learning rule for a sigmoid belief net is:

$$\Delta w_{ij} \propto s_j (s_i - \hat{s}_i)$$

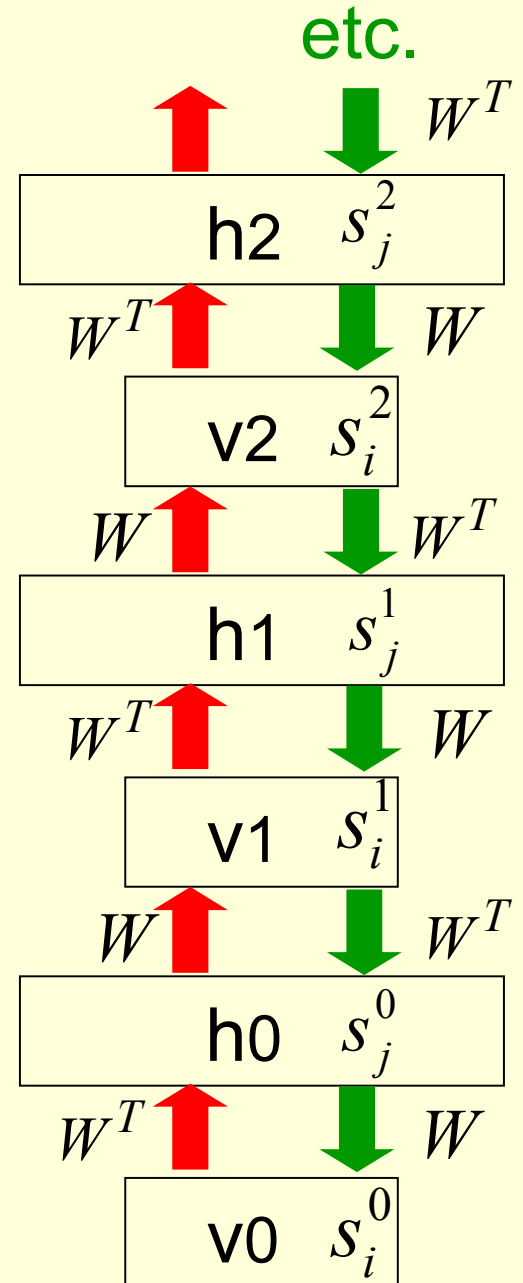
- With replicated weights this becomes:

$$s_j^0 (s_i^0 - s_i^1) +$$

$$s_i^1 (s_j^0 - s_j^1) +$$

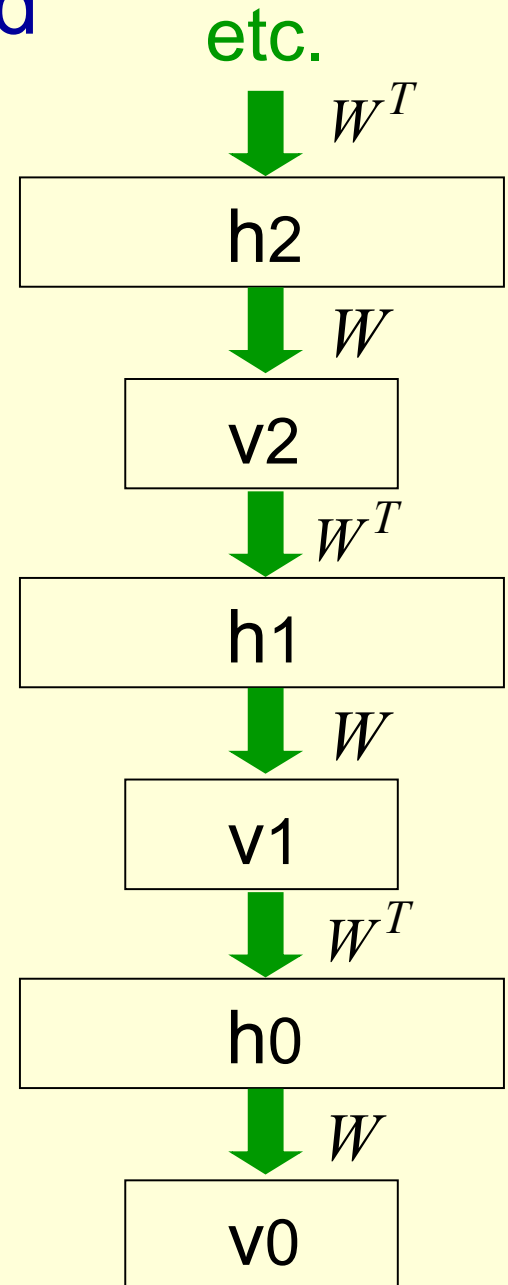
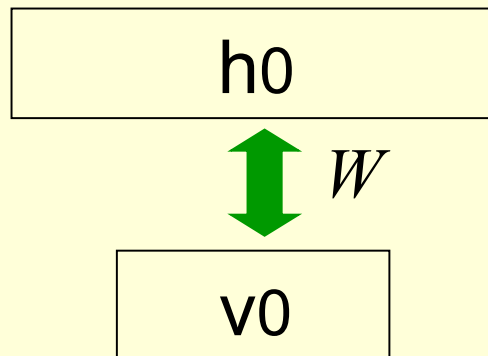
$$s_j^1 (s_i^1 - s_i^2) + \dots$$

$$s_j^\infty s_i^\infty$$

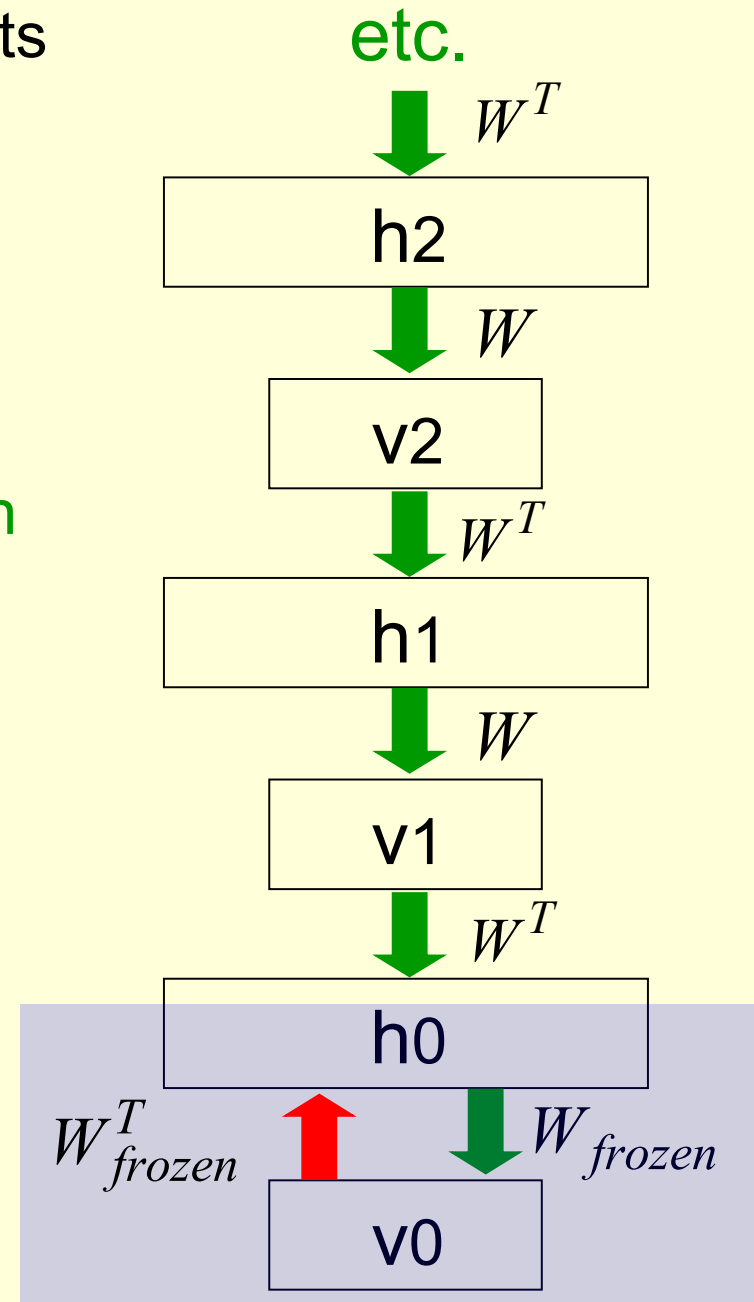
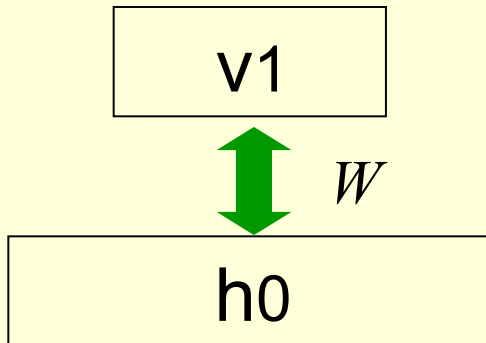


Learning a deep directed network

- First learn with all the weights tied
 - This is exactly equivalent to learning an RBM
 - Contrastive divergence learning is equivalent to ignoring the small derivatives contributed by the tied weights between deeper layers.



- Then freeze the first layer of weights in both directions and learn the remaining weights (still tied together).
 - This is equivalent to learning another RBM, using the aggregated posterior distribution of h_0 as the data.



How many layers should we use and how wide should they be?

- There is no simple answer.
 - Extensive experiments by Yoshua Bengio's group (described later) suggest that several hidden layers is better than one.
 - Results are fairly robust against changes in the size of a layer, but the top layer should be big.
- Deep belief nets give their creator a lot of freedom.
 - The best way to use that freedom depends on the task.
 - With enough narrow layers we can model any distribution over binary vectors (Sutskever & Hinton, 2007)

What happens when the weights in higher layers become different from the weights in the first layer?

- The higher layers no longer implement a complementary prior.
 - So performing inference using the frozen weights in the first layer is no longer correct. But its still pretty good.
 - Using this incorrect inference procedure gives a variational lower bound on the log probability of the data.
- The higher layers learn a prior that is closer to the aggregated posterior distribution of the first hidden layer.
 - This improves the network's model of the data.
 - Hinton, Osindero and Teh (2006) prove that this improvement is always bigger than the loss in the variational bound caused by using less accurate inference.

An improved version of Contrastive Divergence learning (if time permits)

- The main worry with CD is that there will be deep minima of the energy function far away from the data.
 - To find these we need to run the Markov chain for a long time (maybe thousands of steps).
 - But we cannot afford to run the chain for too long for each update of the weights.
- Maybe we can run the same Markov chain over many weight updates? (Neal, 1992)
 - If the learning rate is very small, this should be equivalent to running the chain for many steps and then doing a bigger weight update.

Persistent CD

(Tijmen Teieleman, ICML 2008 & 2009)

- Use minibatches of 100 cases to estimate the first term in the gradient. Use a single batch of 100 fantasies to estimate the second term in the gradient.
- After each weight update, generate the new fantasies from the previous fantasies by using one alternating Gibbs update.
 - So the fantasies can get far from the data.

Contrastive divergence as an adversarial game

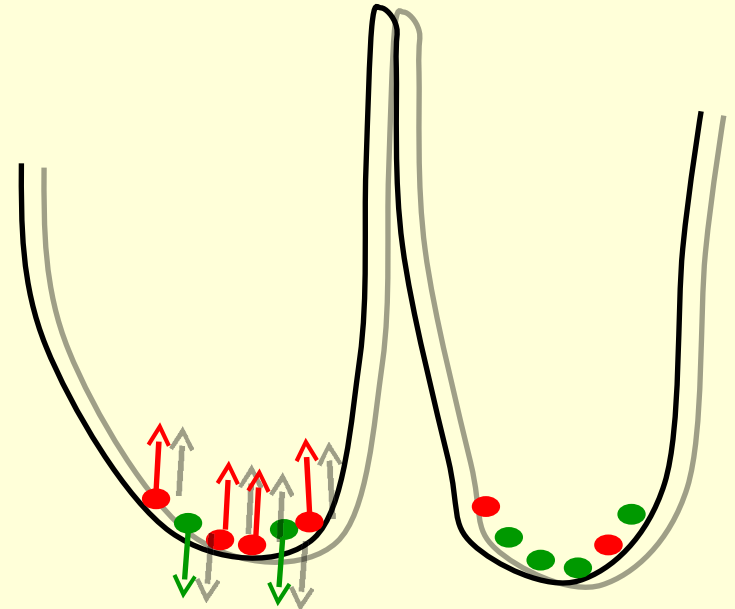
- Why does persistent CD work so well with only 100 negative examples to characterize the whole partition function?
 - For all interesting problems the partition function is highly multi-modal.
 - How does it manage to find all the modes without starting at the data?

The learning causes very fast mixing

- The learning interacts with the Markov chain.
- Persistent Contrastive Divergence cannot be analysed by viewing the learning as an outer loop.
 - Wherever the fantasies outnumber the positive data, the free-energy surface is raised. This makes the fantasies rush around hyperactively.

How persistent CD moves between the modes of the model's distribution

- If a mode has more fantasy particles than data, the free-energy surface is raised until the fantasy particles escape.
 - This can overcome free-energy barriers that would be too high for the Markov Chain to jump.
- The free-energy surface is being changed to help **mixing** in addition to defining the model.



Summary so far

- Restricted Boltzmann Machines provide a simple way to learn a layer of features without any supervision.
 - Maximum likelihood learning is computationally expensive because of the normalization term, but contrastive divergence learning is fast and usually works well.
- Many layers of representation can be learned by treating the hidden states of one RBM as the visible data for training the next RBM (a composition of experts).
- This creates good generative models that can then be fine-tuned.
 - Contrastive wake-sleep can fine-tune generation.

BREAK

Overview of the rest of the tutorial

- How to fine-tune a greedily trained generative model to be better at discrimination.
- How to learn a kernel for a Gaussian process.
- How to use deep belief nets for non-linear dimensionality reduction and document retrieval.
- How to learn a generative hierarchy of conditional random fields.
- A more advanced learning module for deep belief nets that contains multiplicative interactions.
- How to learn deep models of sequential data.

Fine-tuning for discrimination

- First learn one layer at a time greedily.
- Then treat this as “pre-training” that finds a good initial set of weights which can be fine-tuned by a local search procedure.
 - Contrastive wake-sleep is one way of fine-tuning the model to be better at generation.
- Backpropagation can be used to fine-tune the model for better discrimination.
 - This overcomes many of the limitations of standard backpropagation.

Why backpropagation works better with greedy pre-training: The optimization view

- Greedily learning one layer at a time scales well to really big networks, especially if we have locality in each layer.
- We do not start backpropagation until we already have sensible feature detectors that should already be very helpful for the discrimination task.
 - So the initial gradients are sensible and backprop only needs to perform a local search from a sensible starting point.

Why backpropagation works better with greedy pre-training: The overfitting view

- Most of the information in the final weights comes from modeling the distribution of input vectors.
 - The input vectors generally contain a lot more information than the labels.
 - The precious information in the labels is only used for the final fine-tuning.
 - The fine-tuning only modifies the features slightly to get the category boundaries right. It does not need to discover features.
- This type of backpropagation works well even if most of the training data is unlabeled.
 - The unlabeled data is still very useful for discovering good features.

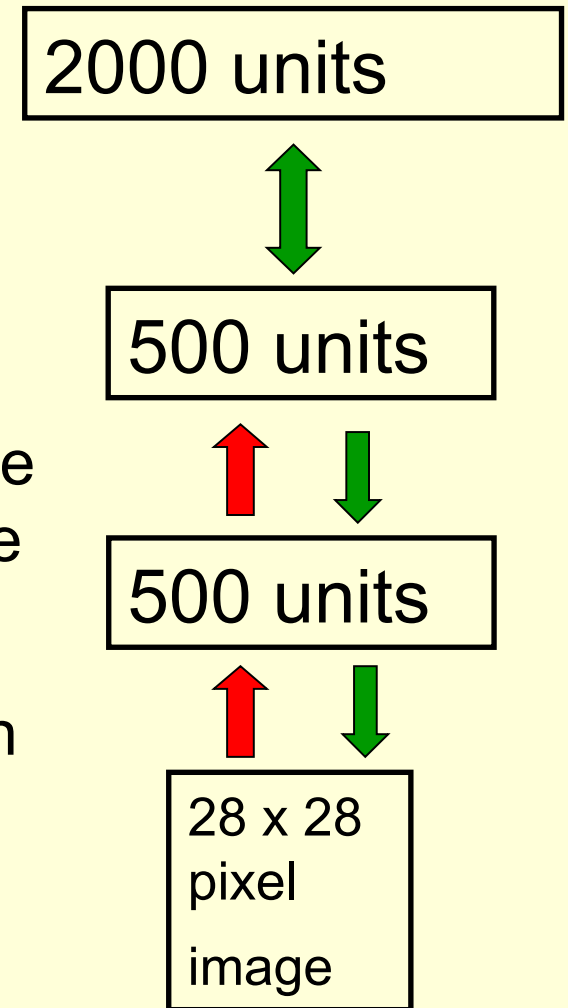
First, model the distribution of digit images

The top two layers form a restricted Boltzmann machine whose free energy landscape should model the low dimensional manifolds of the digits.

The network learns a density model for unlabeled digit images. When we generate from the model we get things that look like real digits of all classes.

But do the hidden features really help with digit discrimination?

Add 10 softmaxed units to the top and do backpropagation.

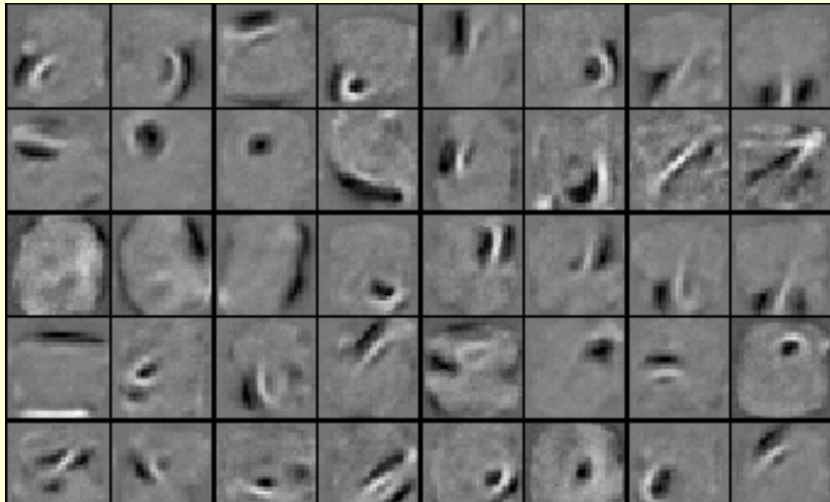


Results on permutation-invariant MNIST task

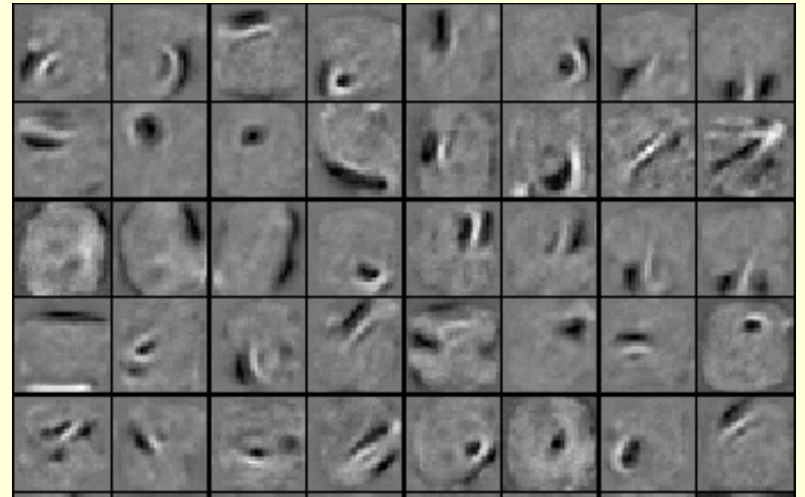
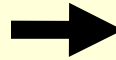
- Very carefully trained backprop net with one or two hidden layers (Platt; Hinton) 1.6%
- SVM (Decoste & Schoelkopf, 2002) 1.4%
- Generative model of joint density of images and labels (+ generative fine-tuning) 1.25%
- Generative model of unlabelled digits followed by gentle backpropagation (Hinton & Salakhutdinov, Science 2006) 1.15%

Learning Dynamics of Deep Nets

the next 4 slides describe work by Yoshua Bengio's group



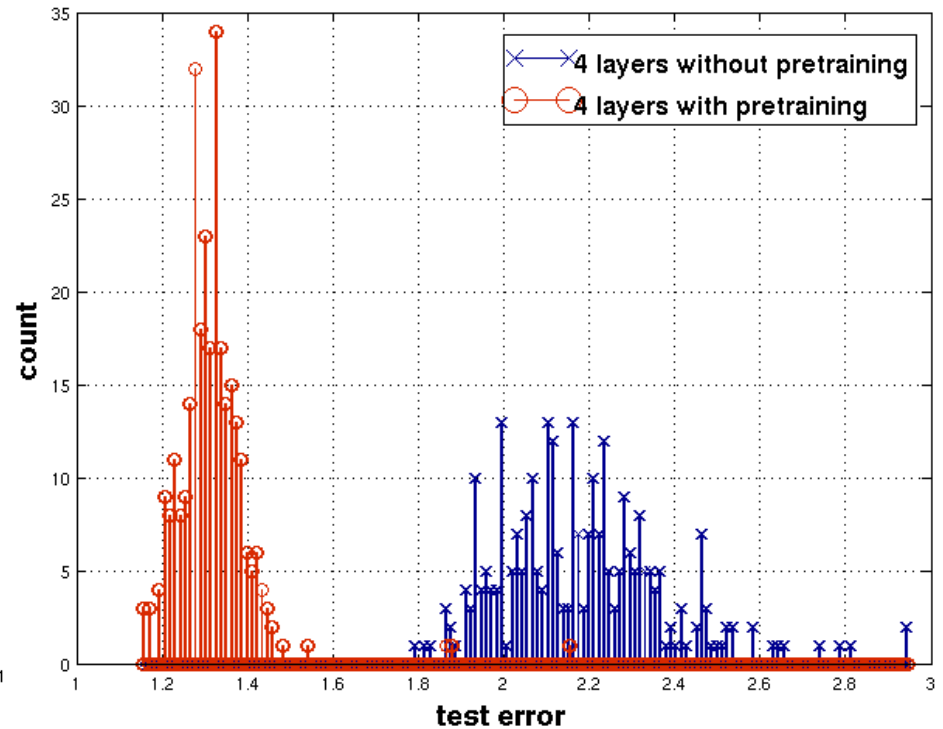
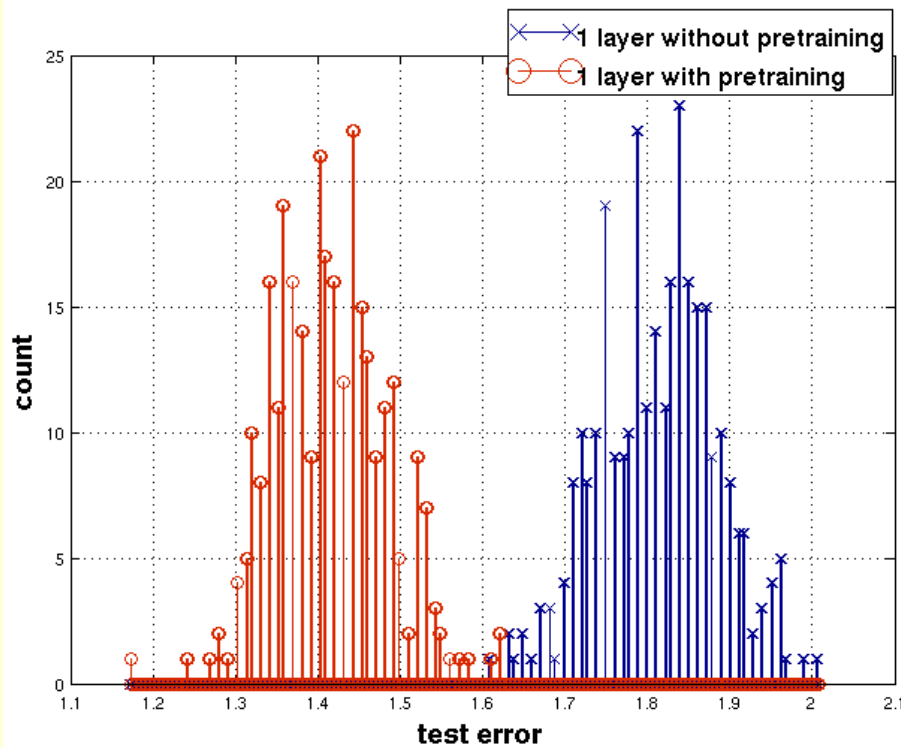
Before fine-tuning



After fine-tuning

Effect of Unsupervised Pre-training

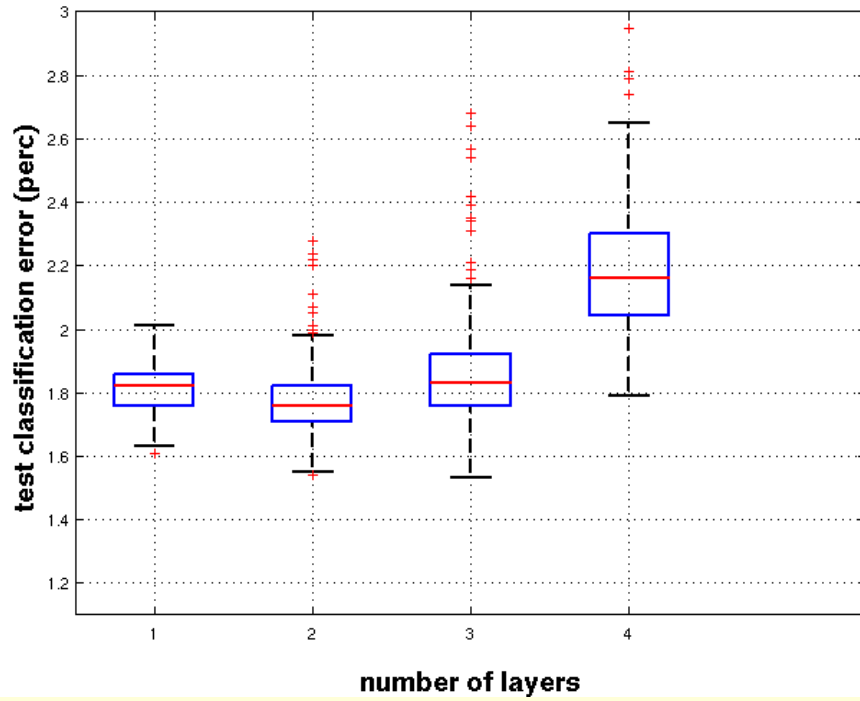
Erhan et. al. AISTATS'2009



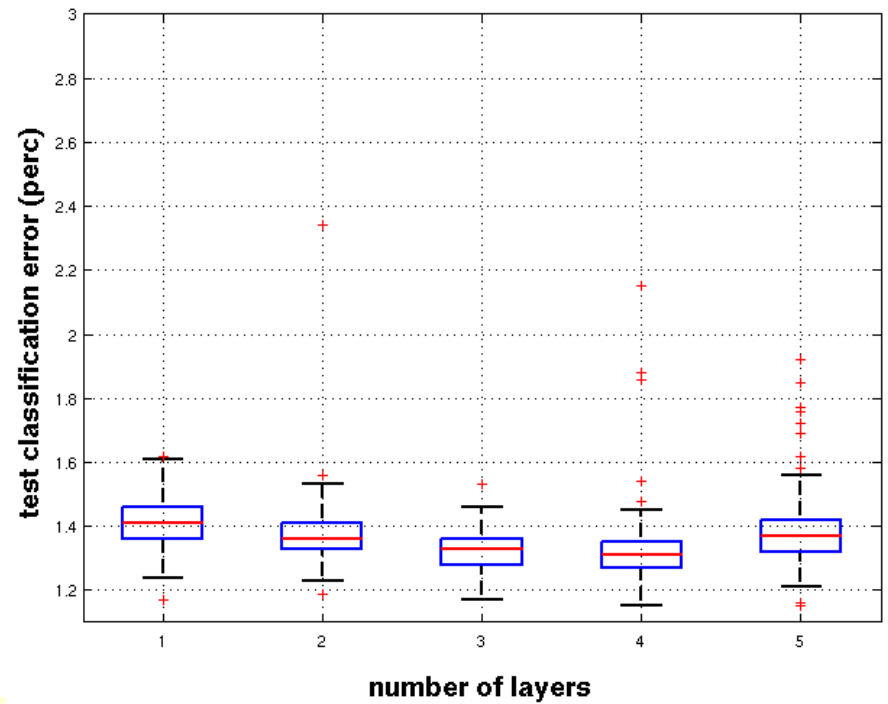
Effect of Depth

without pre-training

w/o pre-training



with pre-training

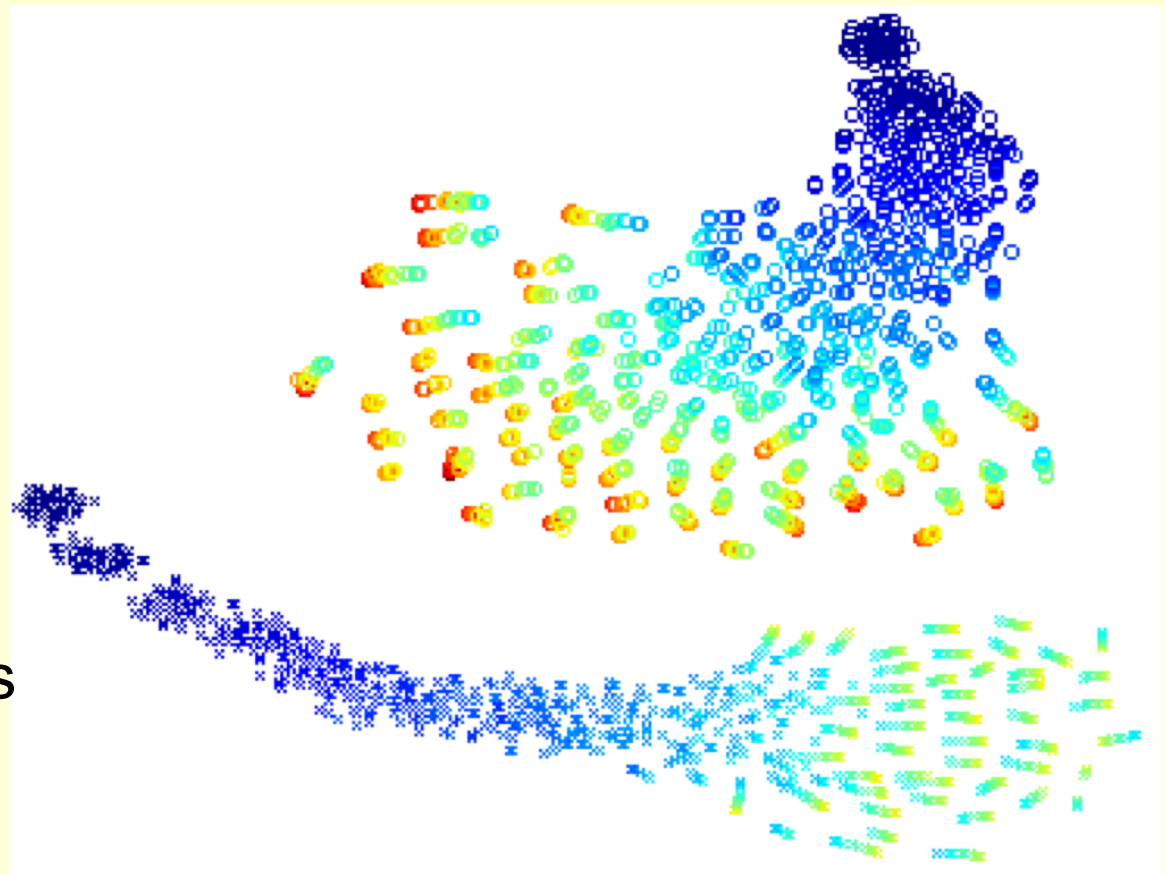


Learning Trajectories in Function Space

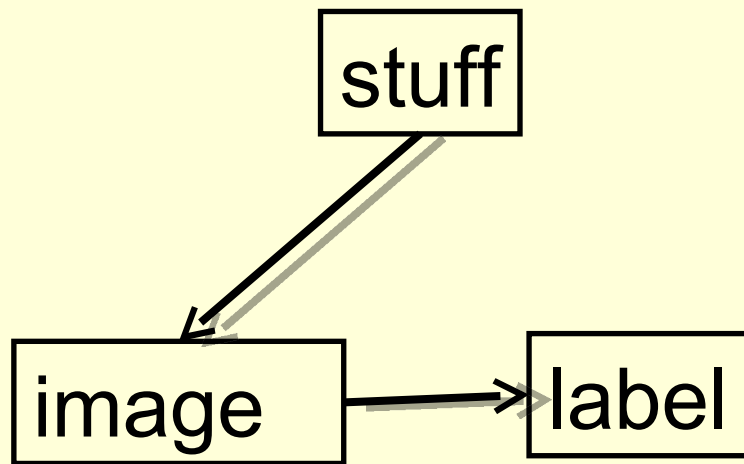
(a 2-D visualization produced with t-SNE)

Erhan et. al. AISTATS'2009

- Each point is a model in function space
- Color = epoch
- Top: trajectories without pre-training. Each trajectory converges to a different local min.
- Bottom: Trajectories with pre-training.
- No overlap!

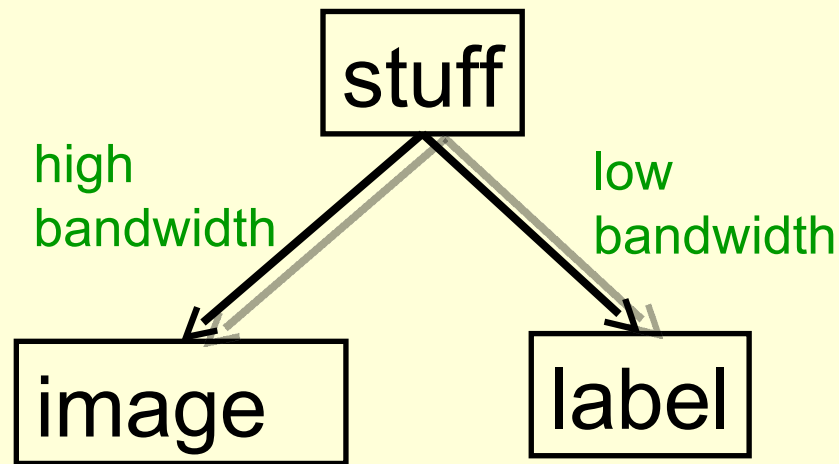


Why unsupervised pre-training makes sense



If image-label pairs were generated this way, it would make sense to try to go straight from images to labels.

For example, do the pixels have even parity?



If image-label pairs are generated this way, it makes sense to first learn to recover the stuff that caused the image by inverting the high bandwidth pathway.

Modeling real-valued data

- For images of digits it is possible to represent intermediate intensities as if they were probabilities by using “mean-field” logistic units.
 - We can treat intermediate values as the probability that the pixel is inked.
- This will not work for real images.
 - In a real image, the intensity of a pixel is almost always almost exactly the average of the neighboring pixels.
 - Mean-field logistic units cannot represent precise intermediate values.

Replacing binary variables by integer-valued variables

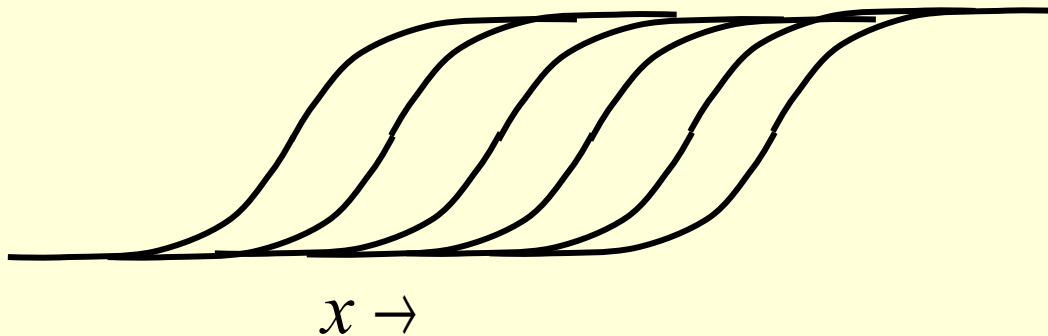
(Teh and Hinton, 2001)

- One way to model an integer-valued variable is to make N identical copies of a binary unit.
- All copies have the same probability, of being “on” : $p = \text{logistic}(x)$
 - The total number of “on” copies is like the firing rate of a neuron.
 - It has a binomial distribution with mean $N p$ and variance $N p(1-p)$

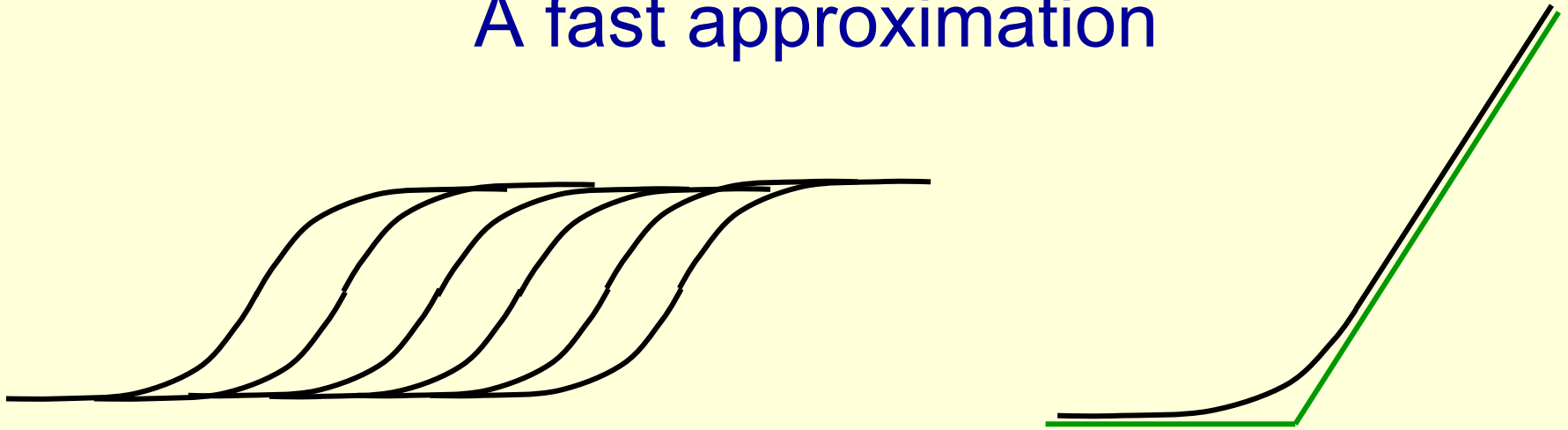
A better way to implement integer values

- Make many copies of a binary unit.
- All copies have the same weights and the same adaptive bias, b , but they have different fixed offsets to the bias:

$$b - 0.5, b - 1.5, b - 2.5, b - 3.5, \dots$$



A fast approximation



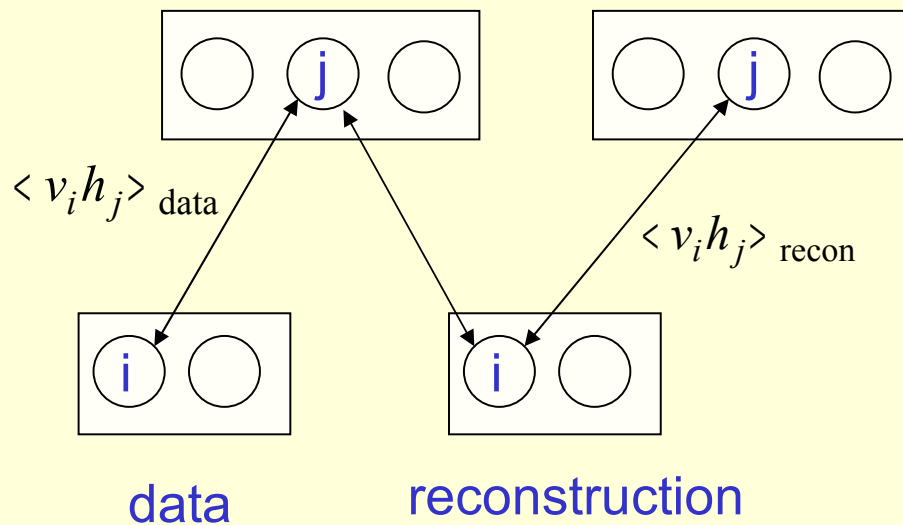
$$\sum_{n=1}^{n=\infty} \text{logistic}(x + 0.5 - n) \approx \log(1 + e^x)$$

- Contrastive divergence learning works well for the sum of binary units with offset biases.
- It also works for rectified linear units. These are much faster to compute than the sum of many logistic units.

$$\text{output} = \max(0, x + \text{randn} * \text{sqrt}(\text{logistic}(x)))$$

How to train a bipartite network of rectified linear units

- Just use contrastive divergence to lower the energy of data and raise the energy of nearby configurations that the model prefers to the data.



Start with a training vector on the visible units.

Update all hidden units in parallel **with sampling noise**

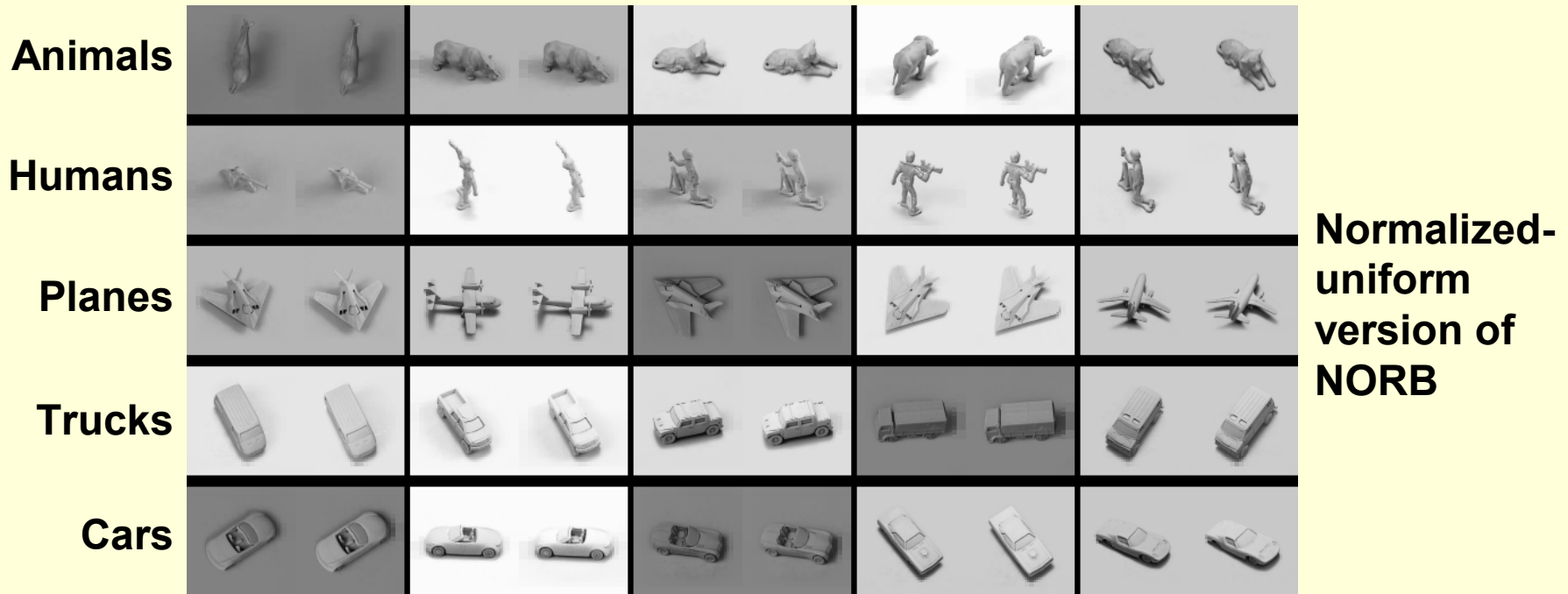
Update the visible units in parallel to get a “reconstruction”.

Update the hidden units again

$$\Delta w_{ij} = \varepsilon (\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{recon}})$$

3D Object Recognition: The NORB dataset

Stereo-pairs of grayscale images of toy objects.



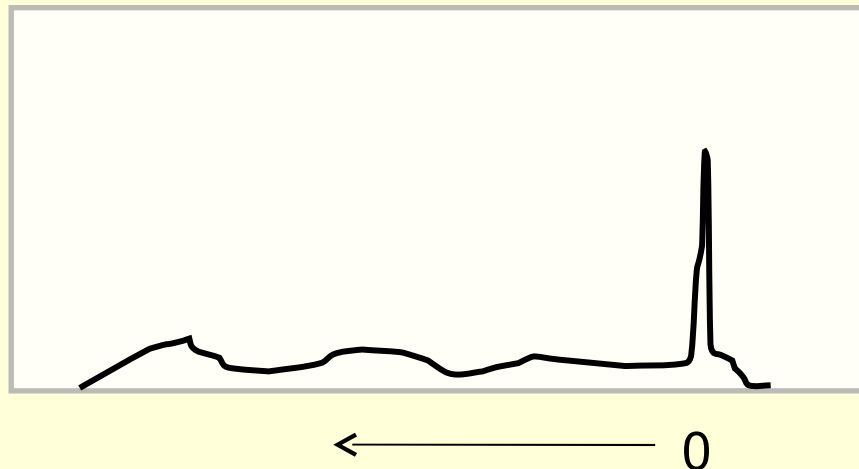
- 6 lighting conditions, 162 viewpoints
- Five object instances per class in the training set
- A *different* set of five instances per class in the test set
- 24,300 training cases, 24,300 test cases

Simplifying the data

- Each training case is a stereo-pair of 96x96 images.
 - The object is centered.
 - The edges of the image are mainly blank.
 - The background is uniform and bright.
- To make learning faster I used simplified the data:
 - Throw away one image.
 - Only use the middle 64x64 pixels of the other image.
 - Downsample to 32x32 by averaging 4 pixels.

Simplifying the data even more so that it can be modeled by rectified linear units

- The intensity histogram for each 32x32 image has a sharp peak for the bright background.
- Find this peak and call it zero.
- Call all intensities brighter than the background zero.
- Measure intensities downwards from the background intensity.



Test set error rates on NORB after greedy learning of one or two hidden layers using **rectified linear units**

Full NORB (2 images of 96x96)

- Logistic regression on the raw pixels 20.5%
 - Gaussian SVM (trained by Leon Bottou) 11.6%
 - Convolutional neural net (Le Cun's group) 6.0%
- (convolutional nets have knowledge of translations built in)

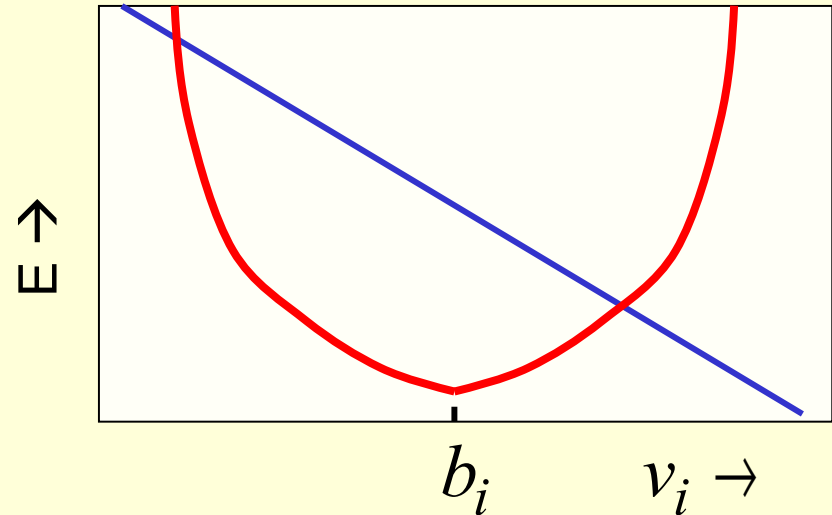
Reduced NORB (1 image 32x32)

- Logistic regression on the raw pixels 30.2%
- Logistic regression on first hidden layer 14.9%
- Logistic regression on second hidden layer 10.2%

The
receptive
fields of
some
rectified
linear
hidden
units.

A standard type of real-valued visible unit

- We can model pixels as Gaussian variables. Alternating Gibbs sampling is still easy, though learning needs to be much slower.

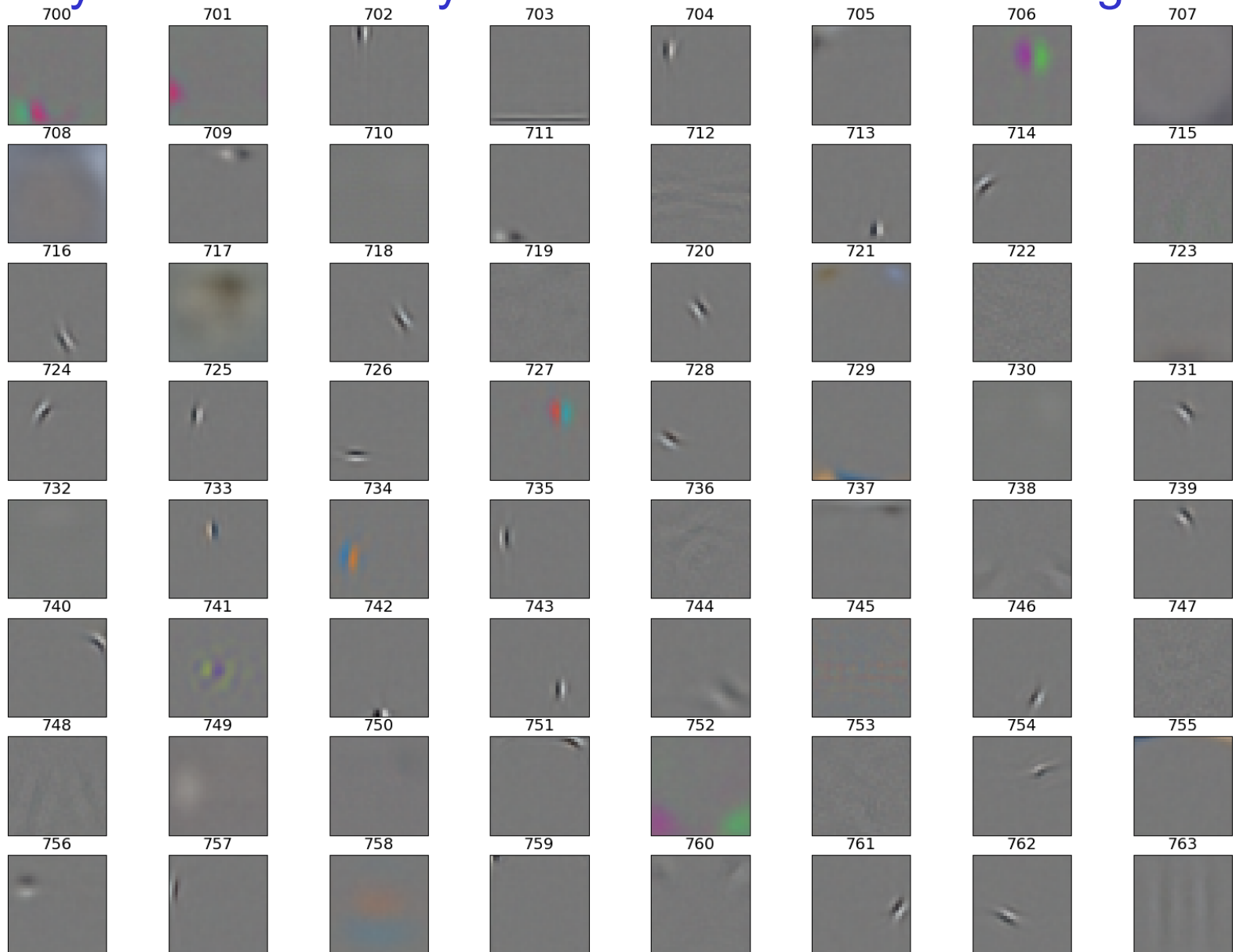


$$E(\mathbf{v}, \mathbf{h}) = \sum_{i \in \text{vis}} \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_{j \in \text{hid}} b_j h_j - \sum_{i,j} \frac{v_i}{\sigma_i} h_j w_{ij}$$

A red arrow points from the text 'parabolic containment function' to the first term of the equation. A blue arrow points from the text 'energy-gradient produced by the total input to a visible unit' to the third term of the equation.

Welling et. al. (2005) show how to extend RBM's to the exponential family. See also Bengio et. al. (2007)

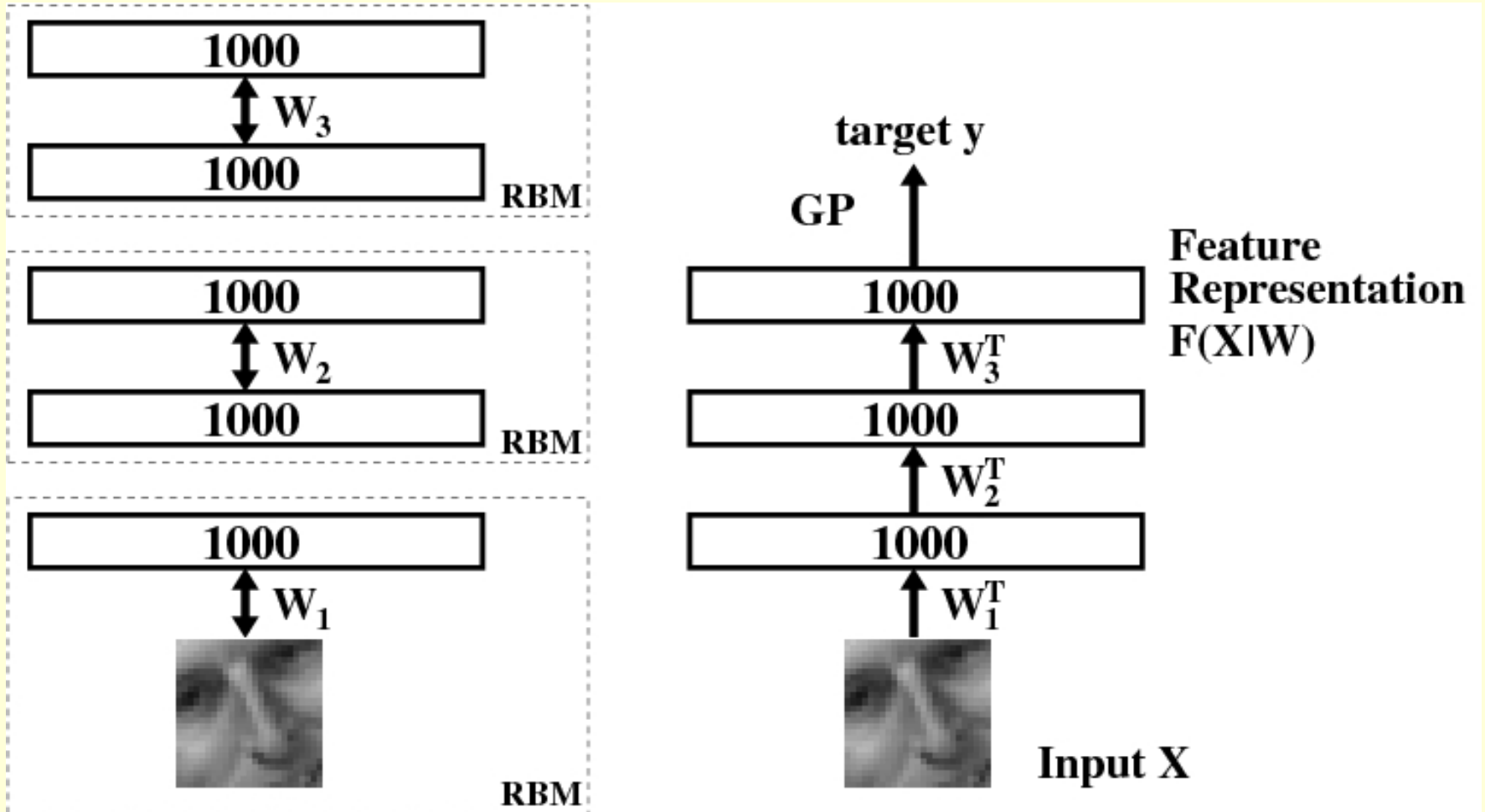
A random sample of 10,000 binary filters learned by Alex Krizhevsky on a million 32x32 color images.



Combining deep belief nets with Gaussian processes

- Deep belief nets can benefit a lot from unlabeled data when labeled data is scarce.
 - They just use the labeled data for fine-tuning.
- Kernel methods, like Gaussian processes, work well on small labeled training sets but are slow for large training sets.
- So when there is a lot of unlabeled data and only a little labeled data, combine the two approaches:
 - First learn a deep belief net without using the labels.
 - Then apply a Gaussian process model to the deepest layer of features. This works better than using the raw data.
 - Then use GP's to get the derivatives that are back-propagated through the deep belief net. This is a further win. It allows GP's to fine-tune complicated domain-specific kernels.

Learning to extract the orientation of a face patch (Salakhutdinov & Hinton, NIPS 2007)



The training and test sets for predicting face orientation

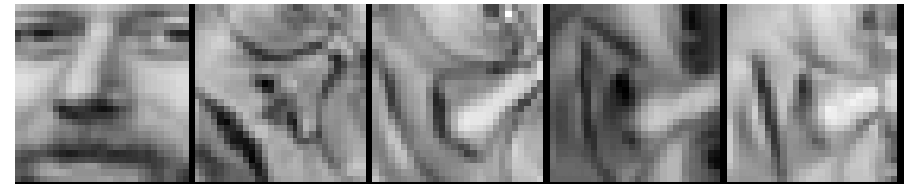
Training Data

-22.07 32.99 -41.15 66.38 27.49



100, 500, or 1000 labeled cases

Unlabeled



11,000 unlabeled cases

Test Data



face patches from new people

The root mean squared error in the orientation when combining GP's with deep belief nets

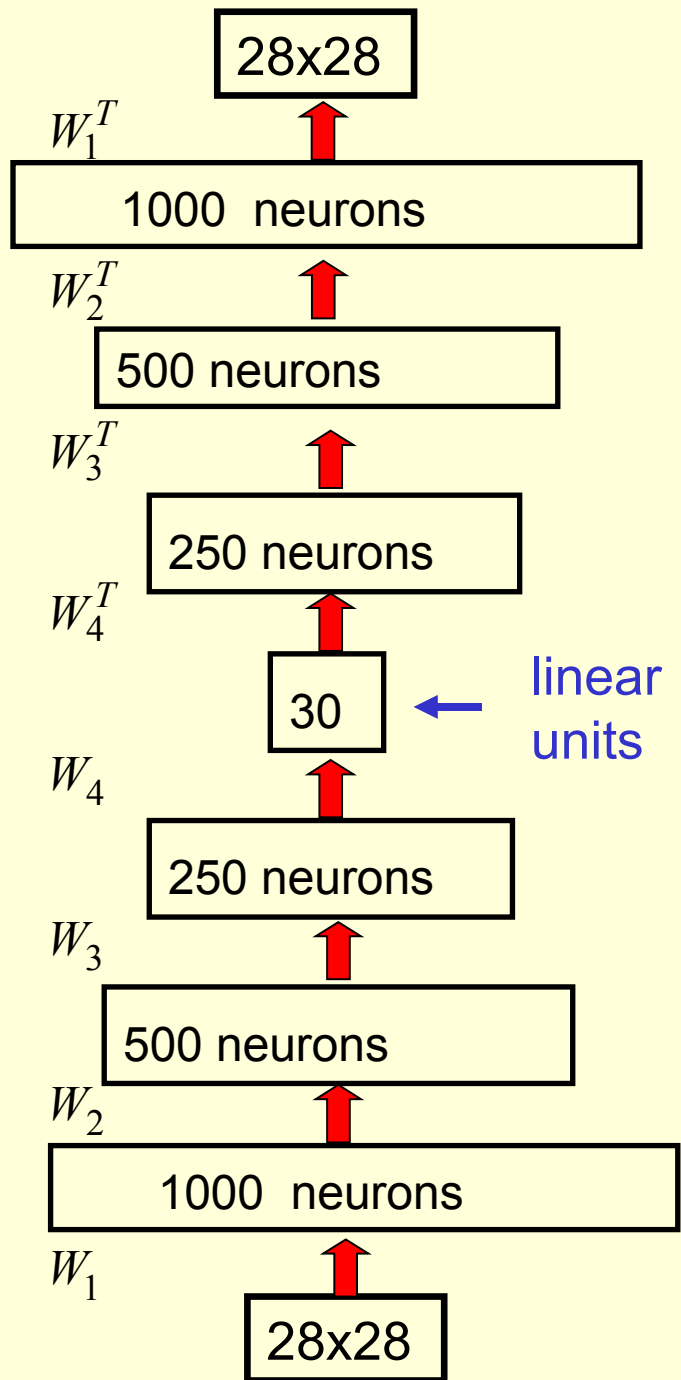
	GP on the pixels	GP on top-level features	GP on top-level features with fine-tuning
100 labels	22.2	17.9	15.2
500 labels	17.2	12.7	7.2
1000 labels	16.3	11.2	6.4

Conclusion: The deep features are much better than the pixels. Fine-tuning helps a lot.

Deep Autoencoders

(Hinton & Salakhutdinov, 2006)

- They always looked like a really nice way to do non-linear dimensionality reduction:
 - But it is **very** difficult to optimize deep autoencoders using backpropagation.
- We now have a much better way to optimize them:
 - First train a stack of 4 RBM's
 - Then “unroll” them.
 - Then fine-tune with backprop.



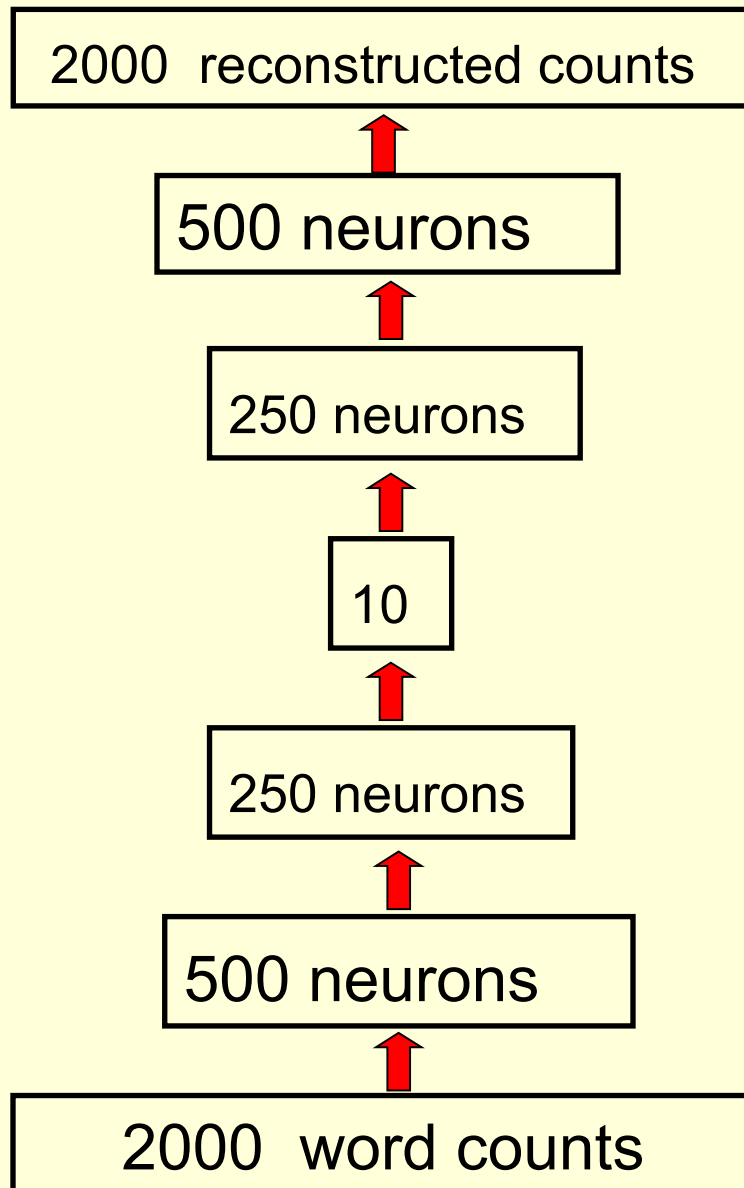
A comparison of methods for compressing digit images to 30 real numbers.



Retrieving documents that are similar to a query document

- We can use an autoencoder to find low-dimensional codes for documents that allow fast and accurate retrieval of similar documents from a large set.
- We start by converting each document into a “bag of words”. This a 2000 dimensional vector that contains the counts for each of the 2000 commonest words.

How to compress the count vector



output
vector

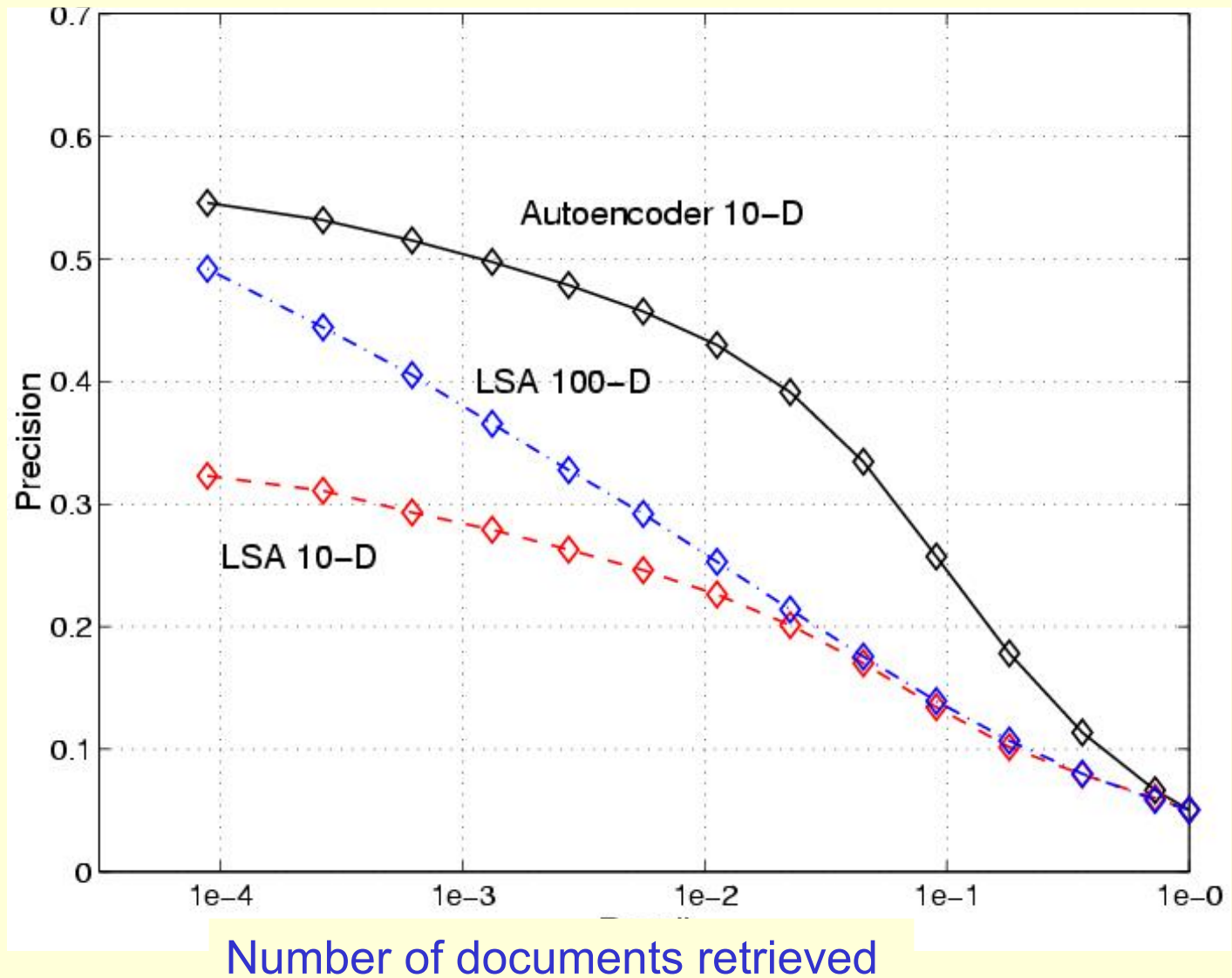
- We train the neural network to reproduce its input vector as its output
- This forces it to compress as much information as possible into the 10 numbers in the central bottleneck.
- These 10 numbers are then a good way to compare documents.

input
vector

Performance of the autoencoder at document retrieval

- Train on bags of 2000 words for 400,000 training cases of business documents.
 - First train a stack of RBM's. Then fine-tune with backprop.
- Test on a separate 400,000 documents.
 - Pick one test document as a query. Rank order all the other test documents by using the cosine of the angle between codes.
 - Repeat this using each of the 400,000 test documents as the query (requires 0.16 trillion comparisons).
- Plot the number of retrieved documents against the proportion that are in the same hand-labeled class as the query document.

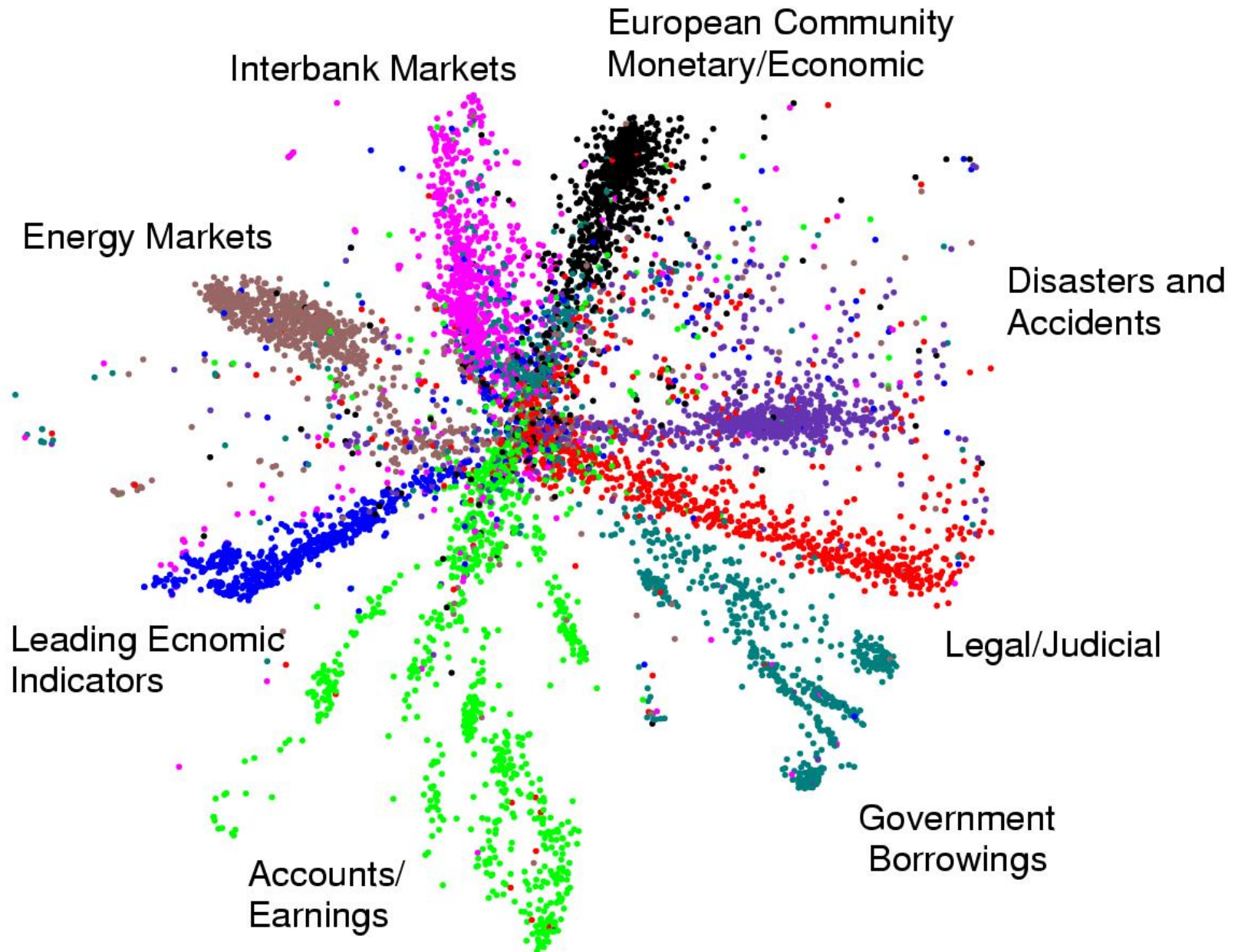
Proportion of retrieved documents in same class as query



First compress all documents to 2 numbers using a type of PCA
Then use different colors for different document categories

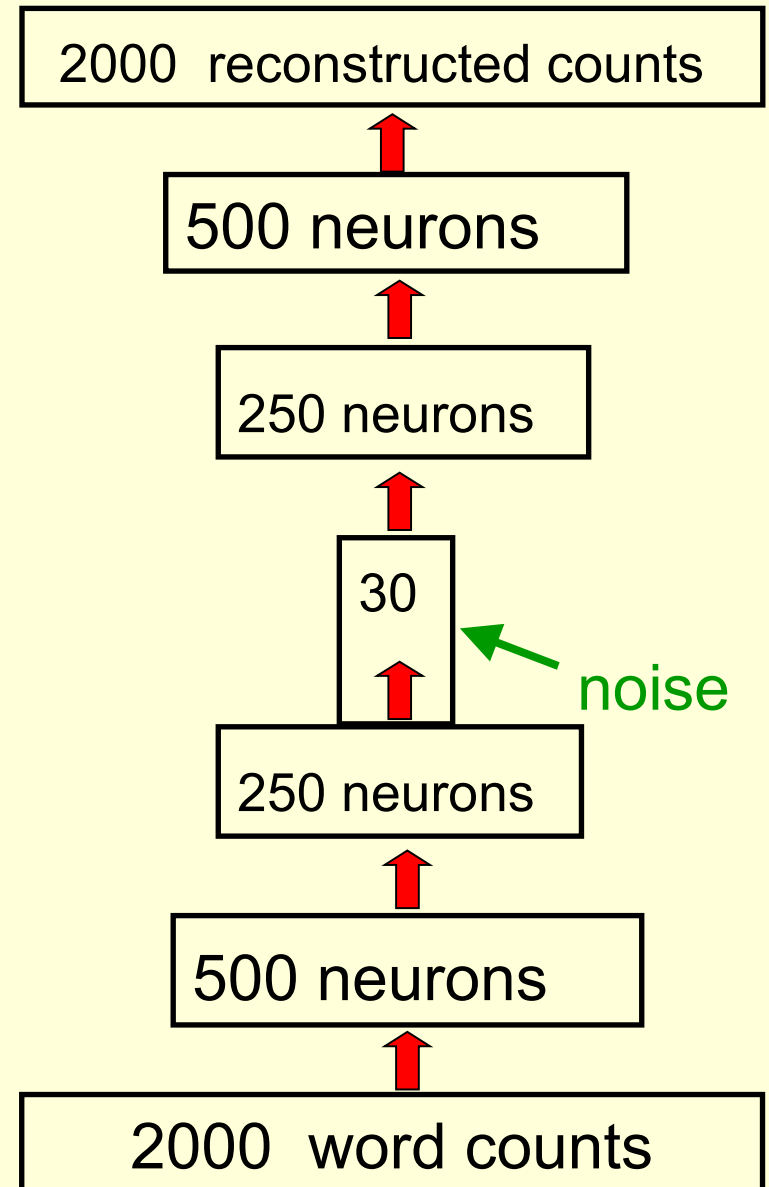


First compress all documents to 2 numbers.
Then use different colors for different document categories

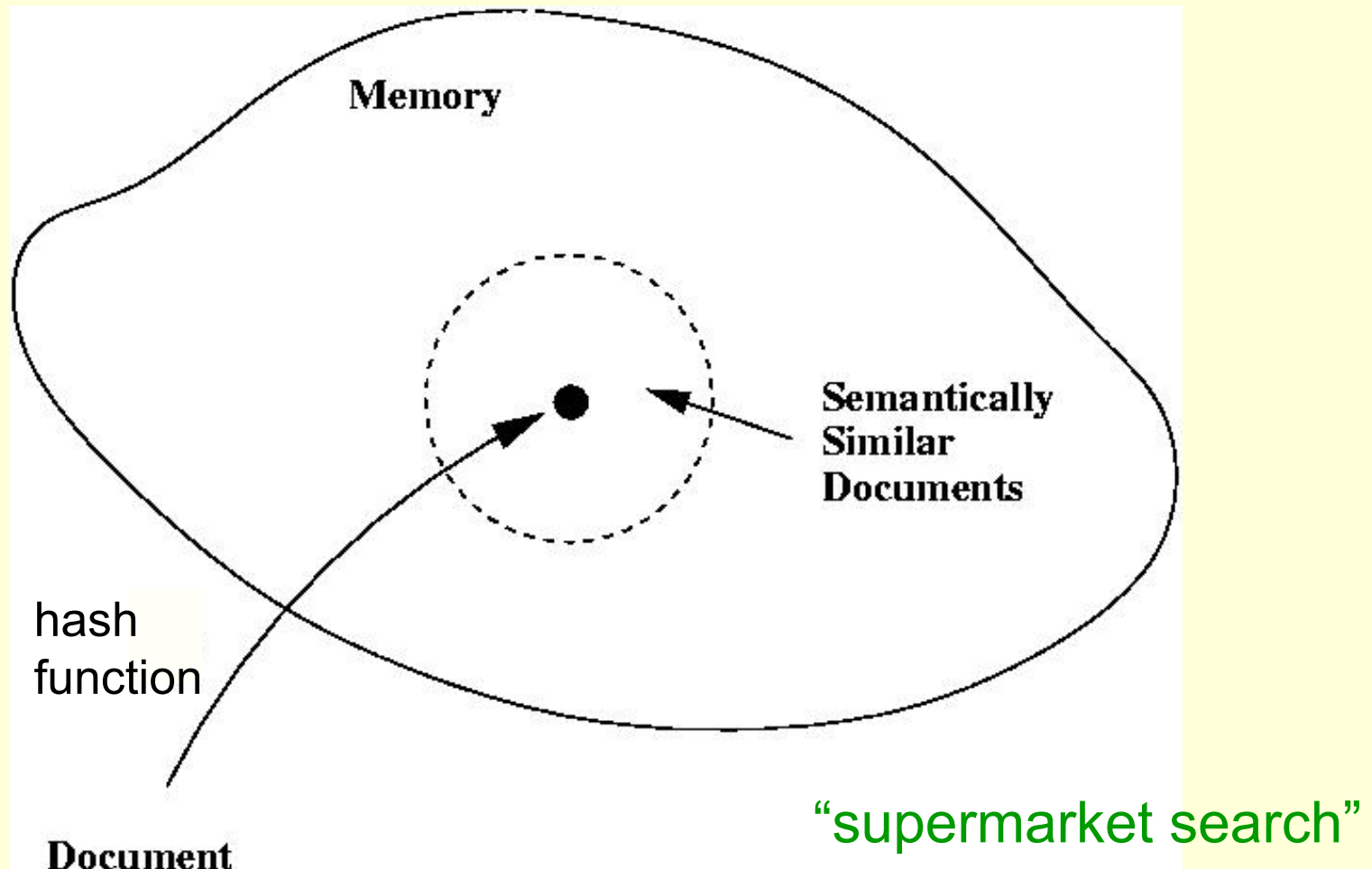


Finding binary codes for documents

- Train an auto-encoder using 30 logistic units for the code layer.
- During the fine-tuning stage, add noise to the inputs to the code units.
 - The “noise” vector for each training case is fixed. So we still get a deterministic gradient.
 - The noise forces their activities to become bimodal in order to resist the effects of the noise.
 - Then we simply round the activities of the 30 code units to 1 or 0.



Semantic hashing: Using a deep autoencoder as a hash-function for finding **approximate** matches
(Salakhutdinov & Hinton, 2007)

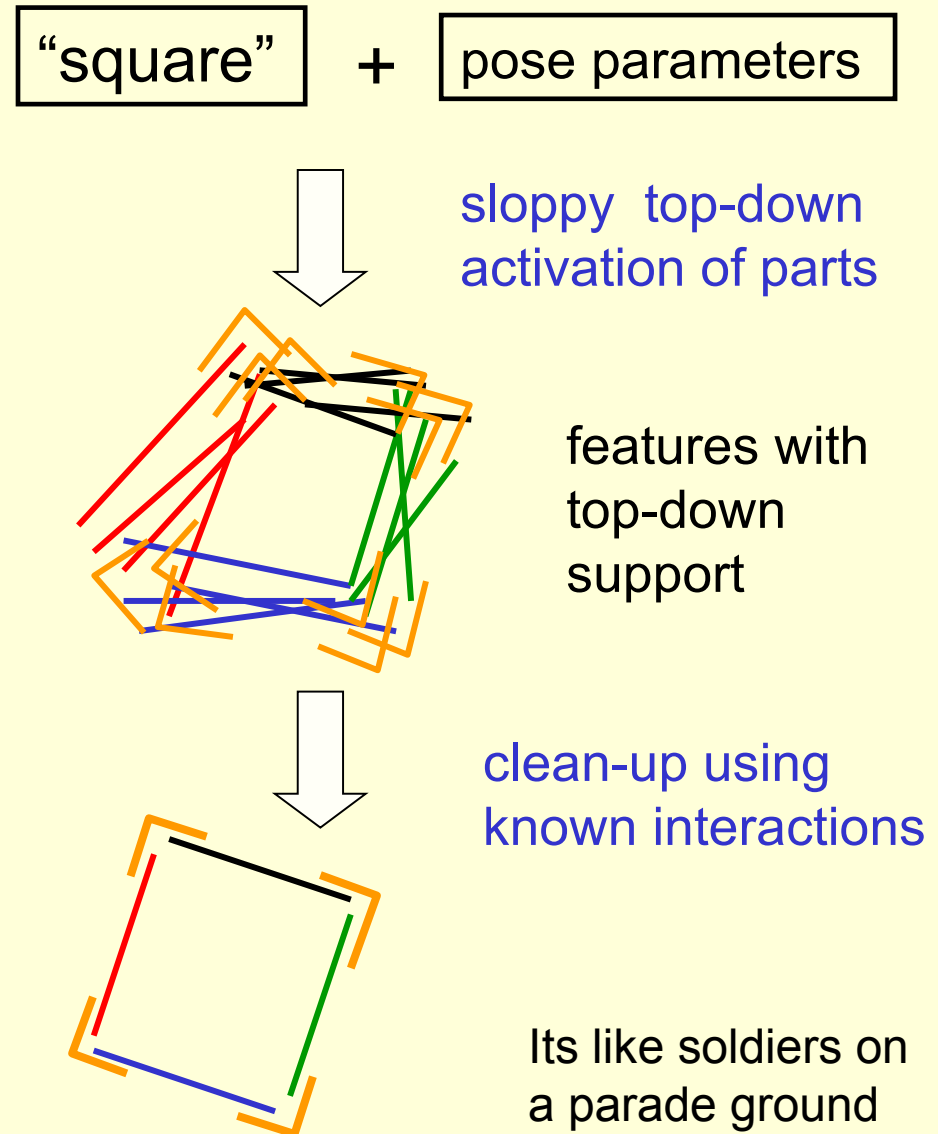


How good is a shortlist found this way?

- We have only implemented it for a million documents with 20-bit codes --- but what could possibly go wrong?
 - A 20-D hypercube allows us to capture enough of the similarity structure of our document set.
- The shortlist found using binary codes actually improves the precision-recall curves of TF-IDF.
 - Locality sensitive hashing (the fastest other method) is 50 times slower and has worse precision-recall curves.

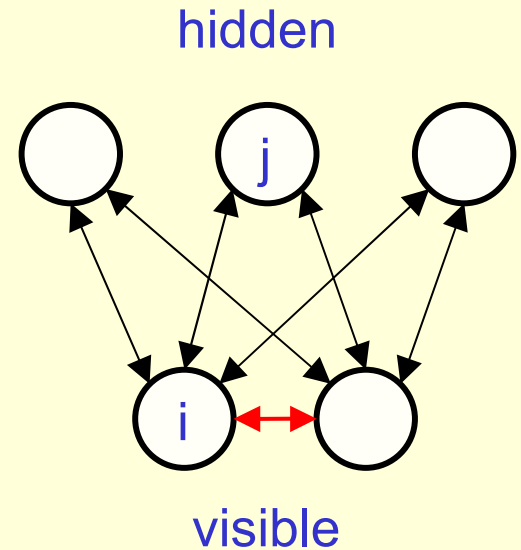
Generating the parts of an object

- One way to maintain the constraints between the parts is to generate each part very accurately
 - But this would require a lot of communication bandwidth.
- Sloppy top-down specification of the parts is less demanding
 - but it messes up relationships between features
 - so use redundant features and use lateral interactions to clean up the mess.
- Each transformed feature helps to locate the others
 - This allows a noisy channel

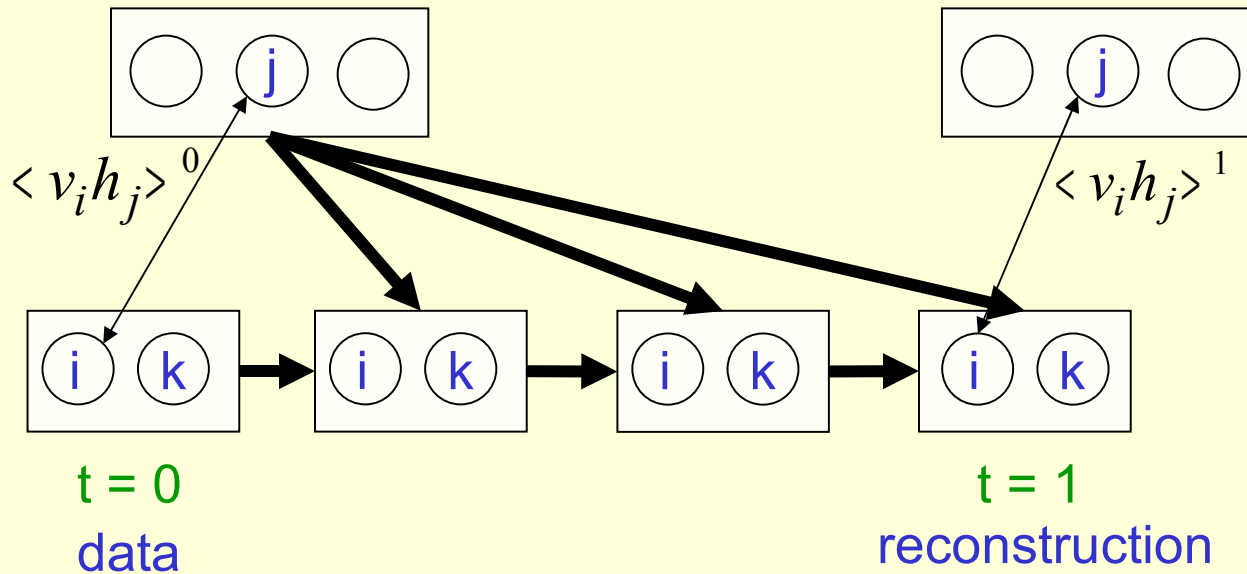


Semi-restricted Boltzmann Machines

- We restrict the connectivity to make learning easier.
- Contrastive divergence learning requires the hidden units to be in conditional equilibrium with the visibles.
 - But it does not require the visible units to be in conditional equilibrium with the hiddens.
 - All we require is that the visible units are closer to equilibrium in the reconstructions than in the data.
- So we can allow connections between the visibles.



Learning a semi-restricted Boltzmann Machine



1. Start with a training vector on the visible units.

2. Update all of the hidden units in parallel

3. Repeatedly update all of the visible units in parallel using mean-field updates (with the hiddens fixed) to get a “reconstruction”.

4. Update all of the hidden units again.

$$\Delta w_{ij} = \varepsilon (\langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^1)$$

$$\Delta l_{ik} = \varepsilon (\langle v_i v_k \rangle^0 - \langle v_i v_k \rangle^1)$$

↑
update for a
lateral weight

Learning in Semi-restricted Boltzmann Machines

- **Method 1:** To form a reconstruction, cycle through the visible units updating each in turn using the top-down input from the hidden units plus the lateral input from the other visible units.
- **Method 2:** Use “mean field” visible units that have real values. Update them all in parallel.
 - Use damping to prevent oscillations

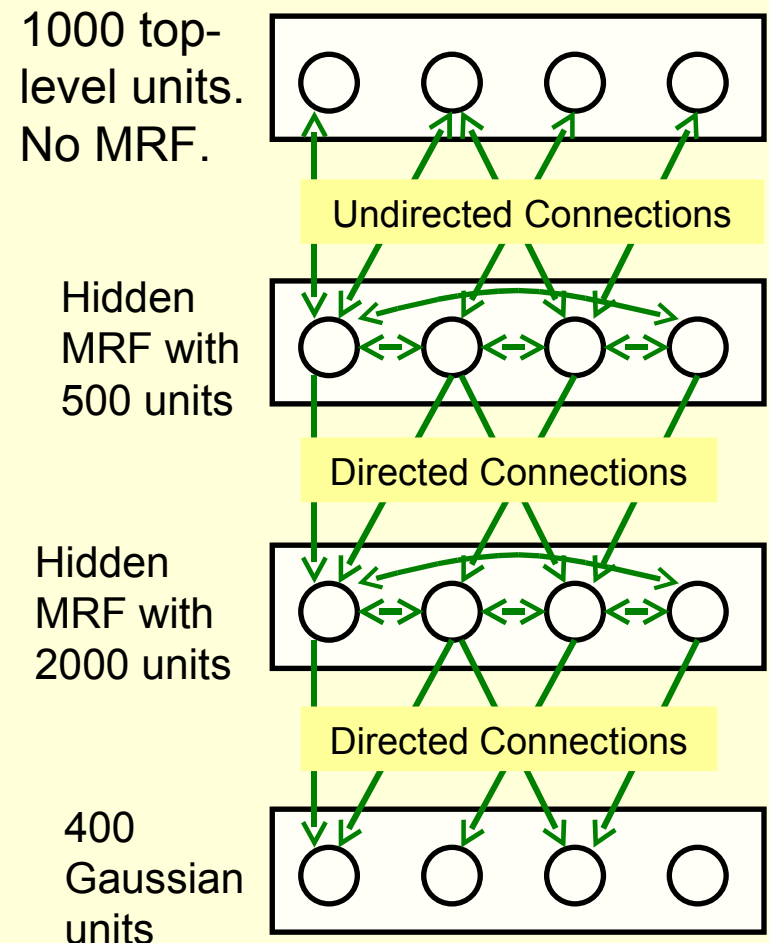
$$p_i^{t+1} = \lambda p_i^t + (1 - \lambda) \sigma(x_i)$$

↑
damping

↑
total input to i

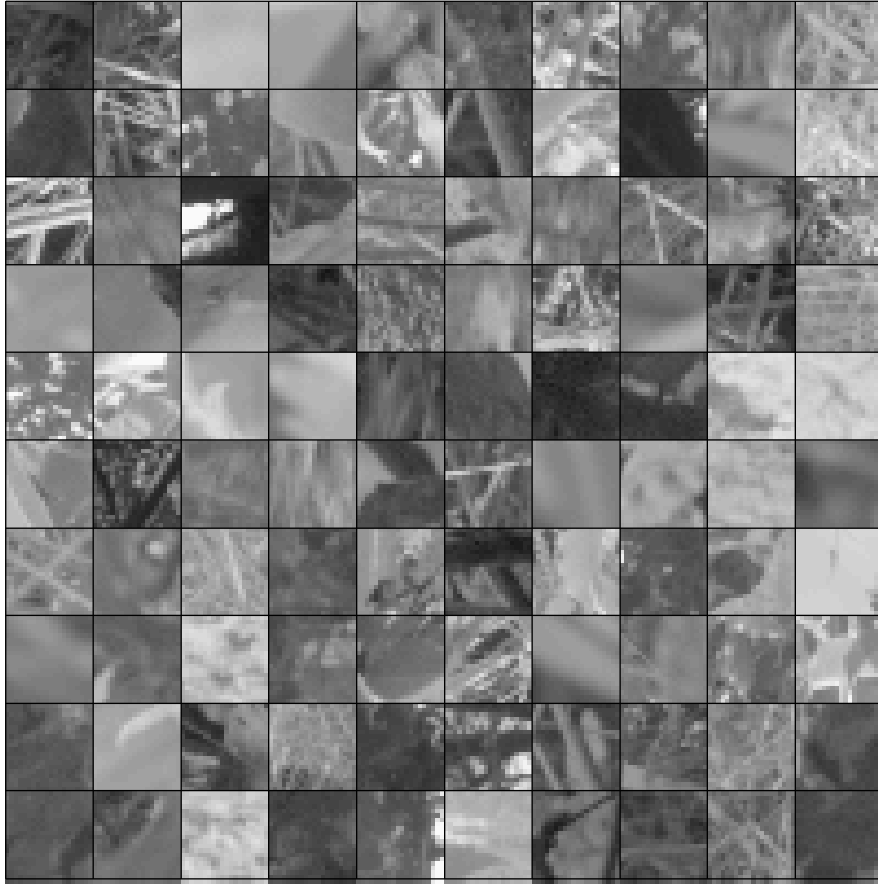
Results on modeling natural image patches using a stack of RBM's (Osindero and Hinton)

- Stack of RBM's learned one at a time.
- 400 Gaussian visible units that see whitened image patches
 - Derived from 100,000 Van Hateren image patches, each 20x20
- The hidden units are all binary.
 - The lateral connections are learned when they are the visible units of their RBM.
- Reconstruction involves letting the visible units of each RBM settle using mean-field dynamics.
 - The already decided states in the level above determine the effective biases during mean-field settling.

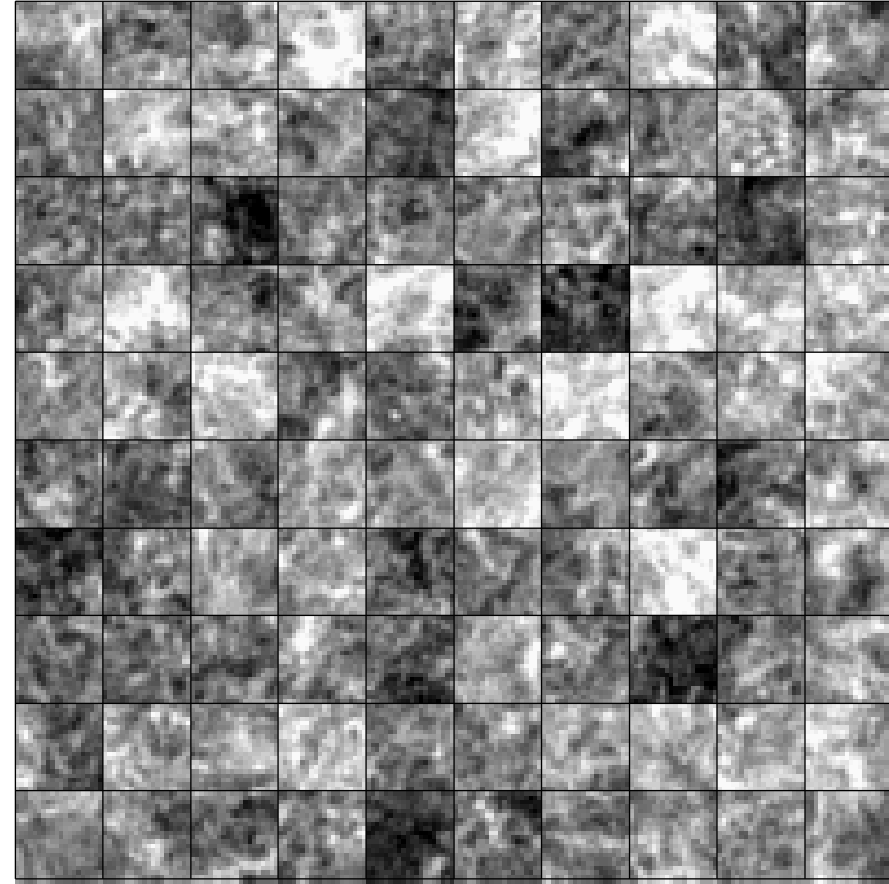


Without lateral connections

real data

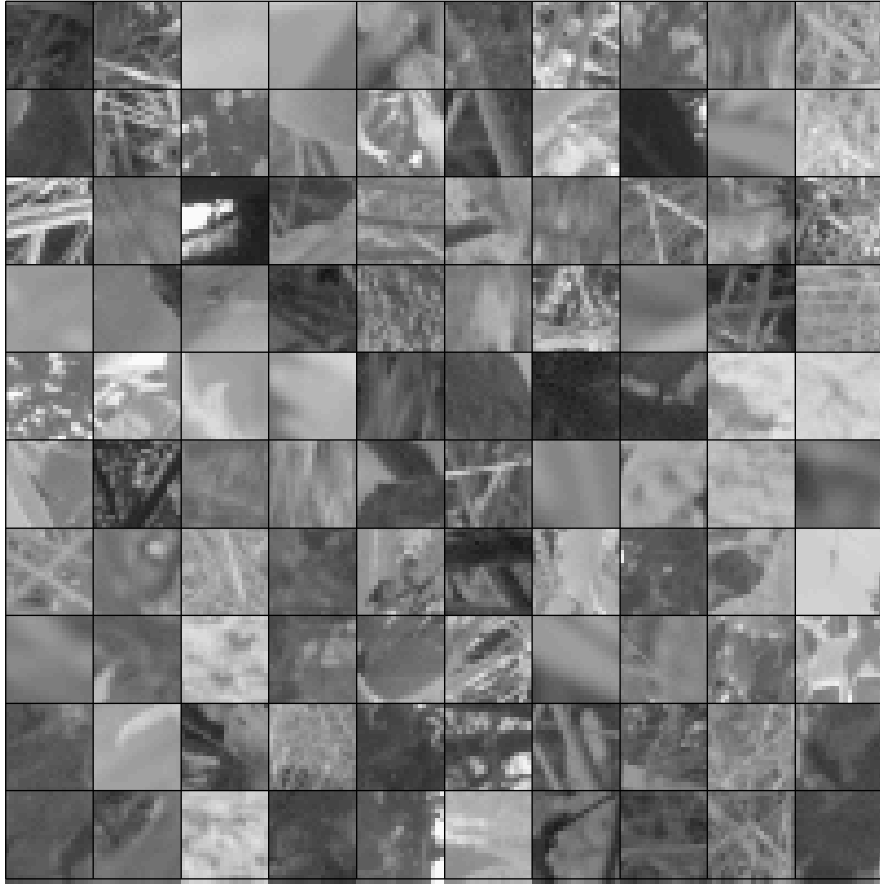


samples from model

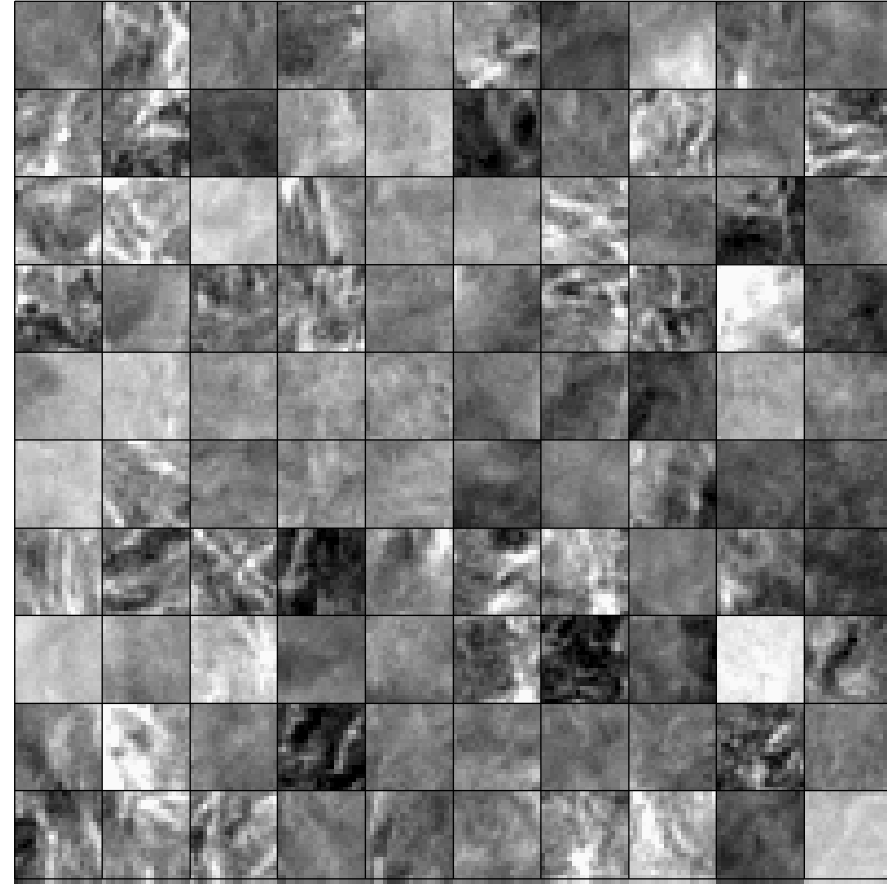


With lateral connections

real data



samples from model



A funny way to use an MRF

- The lateral connections form an MRF.
- The MRF is used during learning and generation.
- The MRF is **not** used for inference.
 - This is a novel idea so vision researchers don't like it.
- The MRF enforces constraints. During inference, constraints do not need to be enforced because the data obeys them.
 - The constraints only need to be enforced during generation.
- Unobserved hidden units cannot enforce constraints.
 - To enforce constraints requires lateral connections or observed descendants.

Why do we whiten data?

- Images typically have strong pair-wise correlations.
- Learning higher order statistics is difficult when there are strong pair-wise correlations.
 - Small changes in parameter values that improve the modeling of higher-order statistics may be rejected because they form a slightly worse model of the much stronger pair-wise statistics.
- So we often remove the second-order statistics before trying to learn the higher-order statistics.

Whitening the learning signal instead of the data

- Contrastive divergence learning can remove the effects of the second-order statistics **on the learning** without actually changing the data.
 - The lateral connections model the second order statistics
 - If a pixel can be reconstructed correctly using second order statistics, its will be the same in the reconstruction as in the data.
 - The hidden units can then focus on modeling high-order structure that cannot be predicted by the lateral connections.
 - For example, a pixel close to an edge, where interpolation from nearby pixels causes incorrect smoothing.

Towards a more powerful, multi-linear stackable learning module

- So far, the states of the units in one layer have only been used to determine the effective biases of the units in the layer below.
- It would be much more powerful to modulate the pair-wise interactions in the layer below.
 - A good way to design a hierarchical system is to allow each level to determine the objective function of the level below.
- To modulate pair-wise interactions we need higher-order Boltzmann machines.

Higher order Boltzmann machines (Sejnowski, ~1986)

- The usual energy function is quadratic in the states:

$$E = \textit{bias terms} - \sum_{i < j} s_i s_j w_{ij}$$

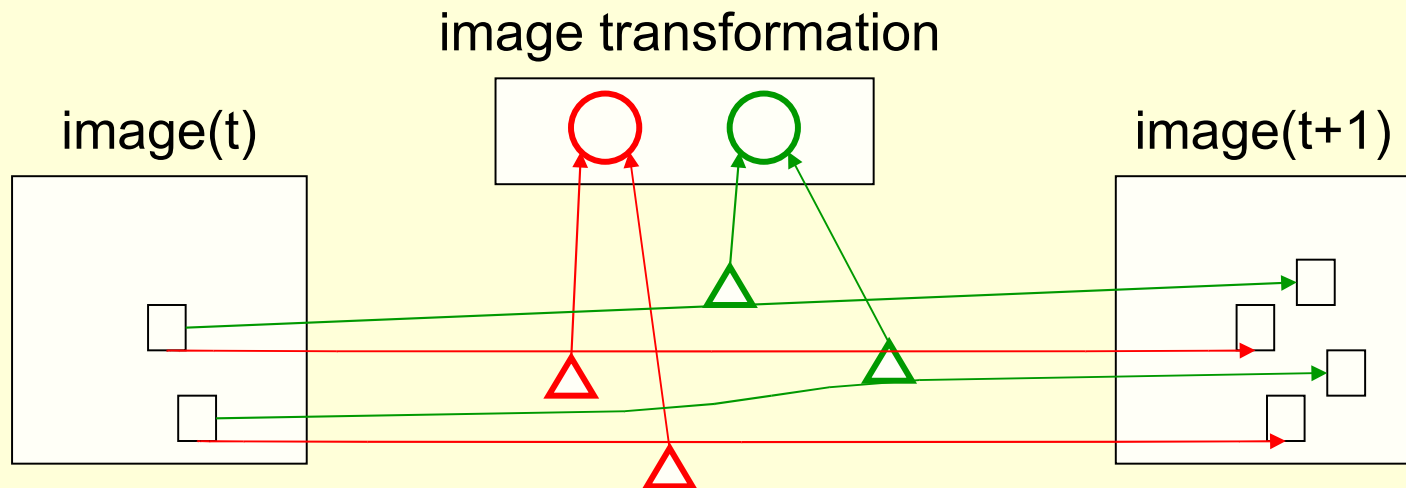
- But we could use higher order interactions:

$$E = \textit{bias terms} - \sum_{i < j < k} s_i s_j s_k w_{ijk}$$

- Unit k acts as a switch. When unit k is on, it switches in the pairwise interaction between unit i and unit j.
 - Units i and j can also be viewed as switches that control the pairwise interactions between j and k or between i and k.

Using higher-order Boltzmann machines to model image transformations (the unfactored version)

- A global transformation specifies which pixel goes to which other pixel.
- Conversely, each pair of similar intensity pixels, one in each image, votes for a particular global transformation.



Factoring three-way multiplicative interactions

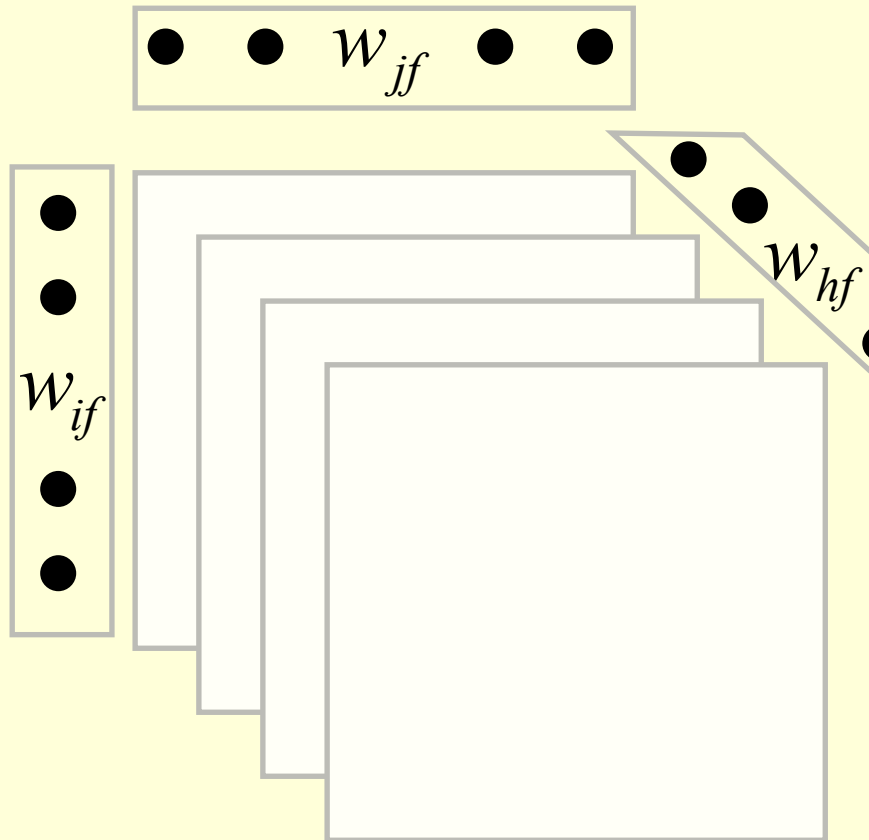
$$- E = \sum_{i,j,h} S_i S_j S_h W_{ijh}$$

← unfactored
with cubically
many parameters

$$- E = \sum_f \sum_{i,j,h} S_i S_j S_h W_{if} W_{jf} W_{hf}$$

← factored
with linearly
many parameters
per factor.

A picture of the low-rank tensor contributed by factor f



Each layer is a scaled version of the same matrix.

The basis matrix is specified as an outer product with typical term $W_{if} W_{jf}$

So each active hidden unit contributes a scalar, W_{hf} times the matrix specified by factor f .

Inference with factored three-way multiplicative interactions

$$- E_f = \sum_{i,j,h} s_i s_j s_h w_{if} w_{jf} w_{hf}$$

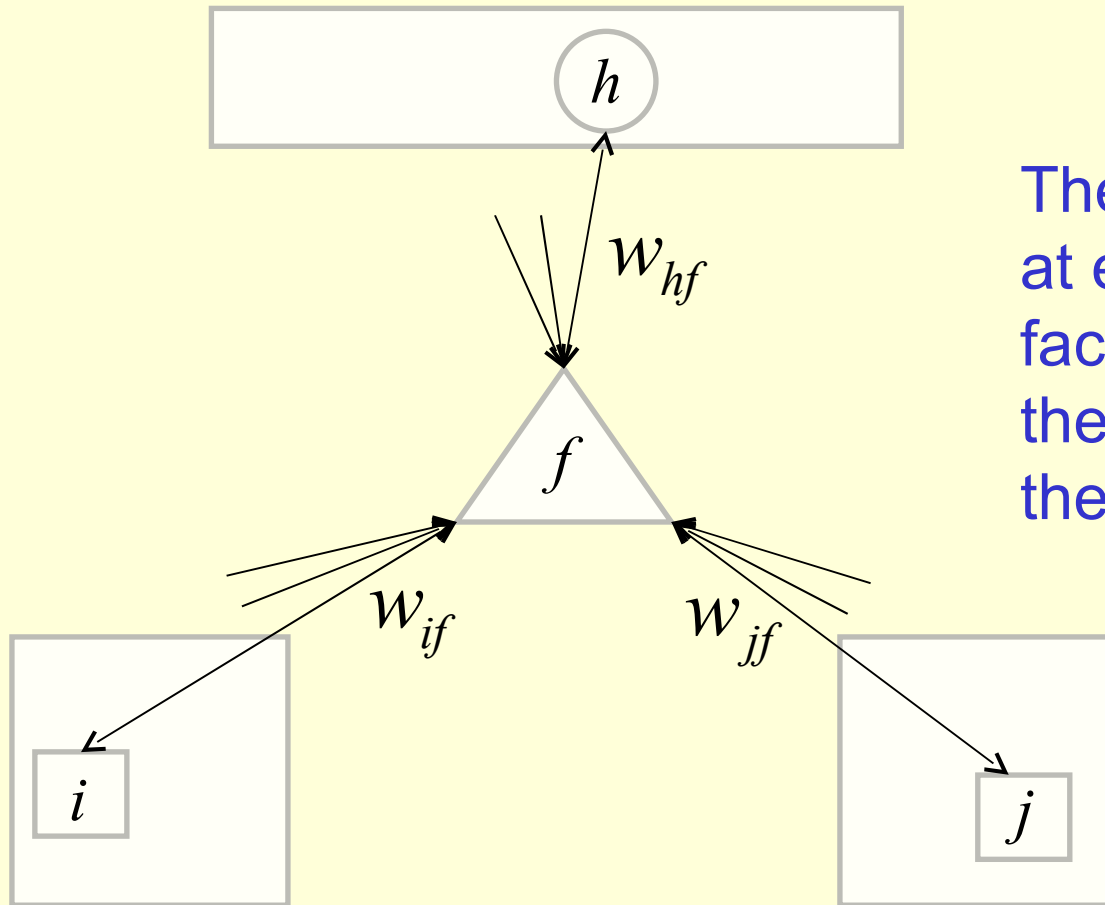
← The energy contributed by factor f.

$$\left[E_f(s_h=0) - E_f(s_h=1) \right] = w_{hf} \left(\sum_i s_i w_{if} \right) \left(\sum_j s_j w_{jf} \right)$$

↑
How changing the binary state of unit h changes the energy contributed by factor f.

↑
What unit h needs to know in order to do Gibbs sampling

Belief propagation



The outgoing message at each vertex of the factor is the product of the weighted sums at the other two vertices.

Learning with factored three-way multiplicative interactions

$$m_f^h = \left(\sum_i s_i w_{if} \right) \left(\sum_j s_j w_{jf} \right)$$

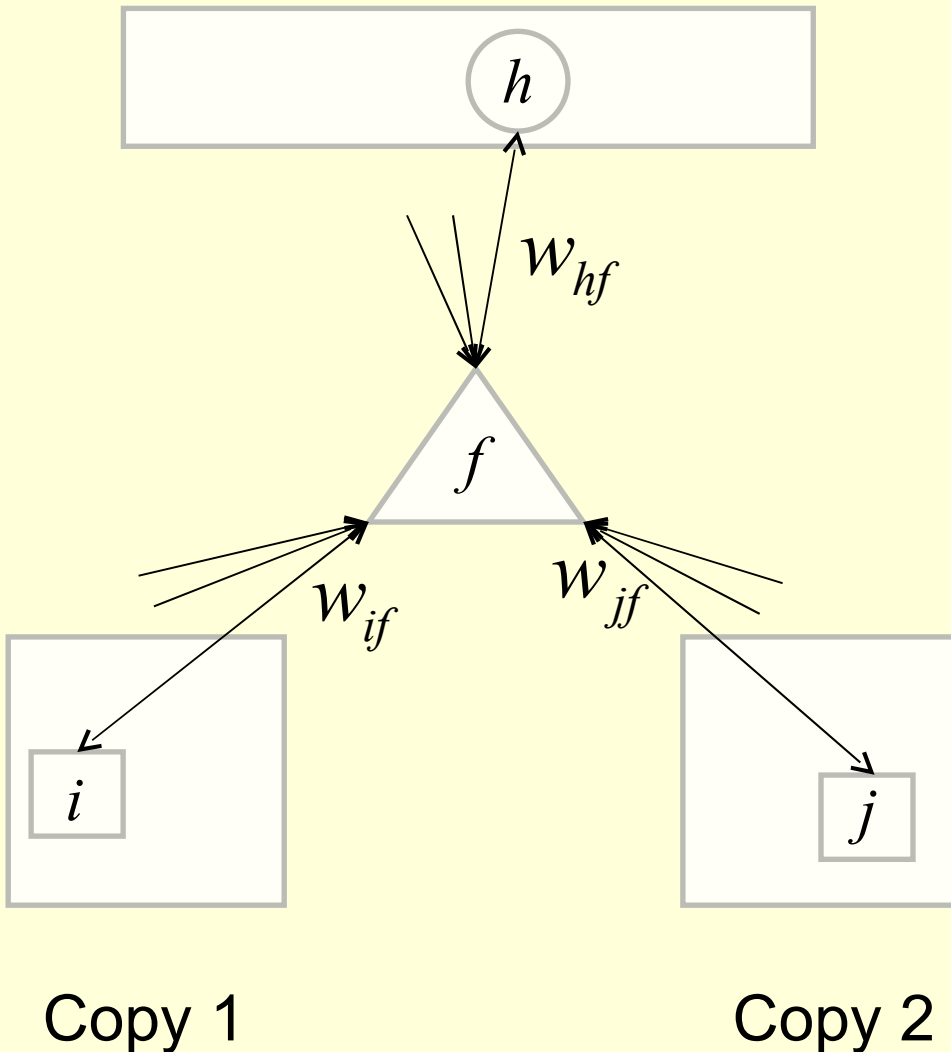
message
from factor f
to unit h

$$\Delta w_{hf} \propto \left\langle -\frac{\partial E_f}{\partial w_{hf}} \right\rangle_{\text{data}} - \left\langle -\frac{\partial E_f}{\partial w_{hf}} \right\rangle_{\text{model}}$$

$$= \left\langle s_h m_f^h \right\rangle_{\text{data}} - \left\langle s_h m_f^h \right\rangle_{\text{model}}$$

Roland data

Modeling the correlational structure of a static image by using two copies of the image



Each factor sends the squared output of a linear filter to the hidden units.

It is exactly the standard model of simple and complex cells. It allows complex cells to extract oriented energy.

The standard model drops out of doing belief propagation for a factored third-order energy function.

An advantage of modeling correlations between pixels rather than pixels

- During generation, a “vertical edge” unit can turn off the horizontal interpolation in a region without worrying about exactly where the intensity discontinuity will be.
 - This gives some translational invariance
 - It also gives a lot of invariance to brightness and contrast.
 - So the “vertical edge” unit is like a complex cell.
- By modulating the correlations between pixels rather than the pixel intensities, the generative model can still allow interpolation parallel to the edge.

A principle of hierarchical systems

- Each level in the hierarchy should not try to micro-manage the level below.
- Instead, it should create an objective function for the level below and leave the level below to optimize it.
 - This allows the fine details of the solution to be decided locally where the detailed information is available.
- Objective functions are a good way to do abstraction.

Time series models

- Inference is difficult in directed models of time series if we use non-linear distributed representations in the hidden units.
 - It is hard to fit Dynamic Bayes Nets to high-dimensional sequences (e.g motion capture data).
- So people tend to avoid distributed representations and use much weaker methods (e.g. HMM's).

Time series models

- If we really need distributed representations (which we nearly always do), we can make inference much simpler by using three tricks:
 - Use an RBM for the interactions between hidden and visible variables. This ensures that the main source of information wants the posterior to be factorial.
 - Model short-range temporal information by allowing several previous frames to provide input to the hidden units and to the visible units.
- This leads to a temporal module that can be stacked
 - So we can use greedy learning to learn deep models of temporal structure.

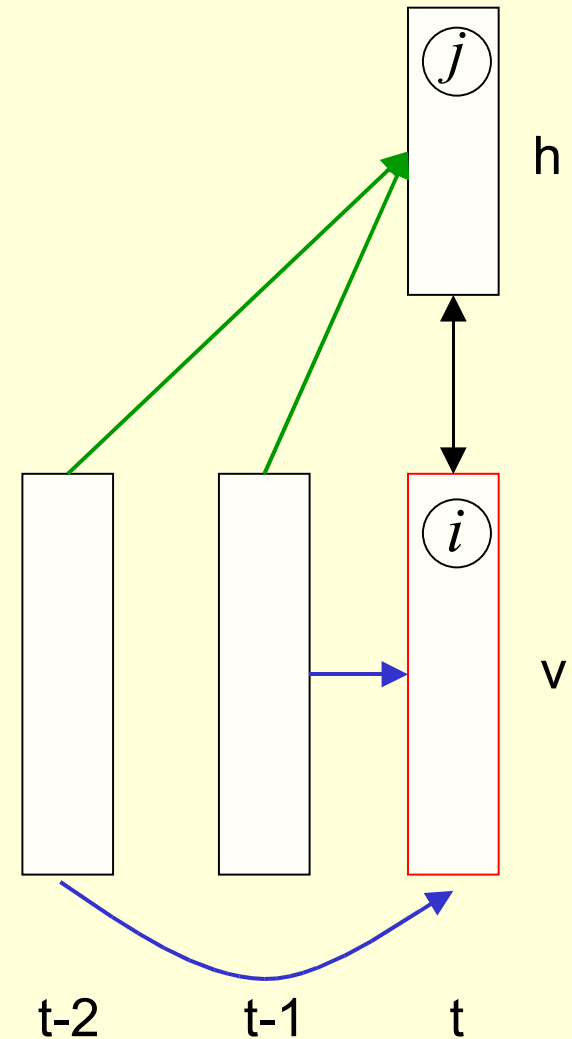
An application to modeling motion capture data

(Taylor, Roweis & Hinton, 2007)

- Human motion can be captured by placing reflective markers on the joints and then using lots of infrared cameras to track the 3-D positions of the markers.
- Given a skeletal model, the 3-D positions of the markers can be converted into the joint angles plus 6 parameters that describe the 3-D position and the roll, pitch and yaw of the pelvis.
 - We only represent **changes** in yaw because physics doesn't care about its value and we want to avoid circular variables.

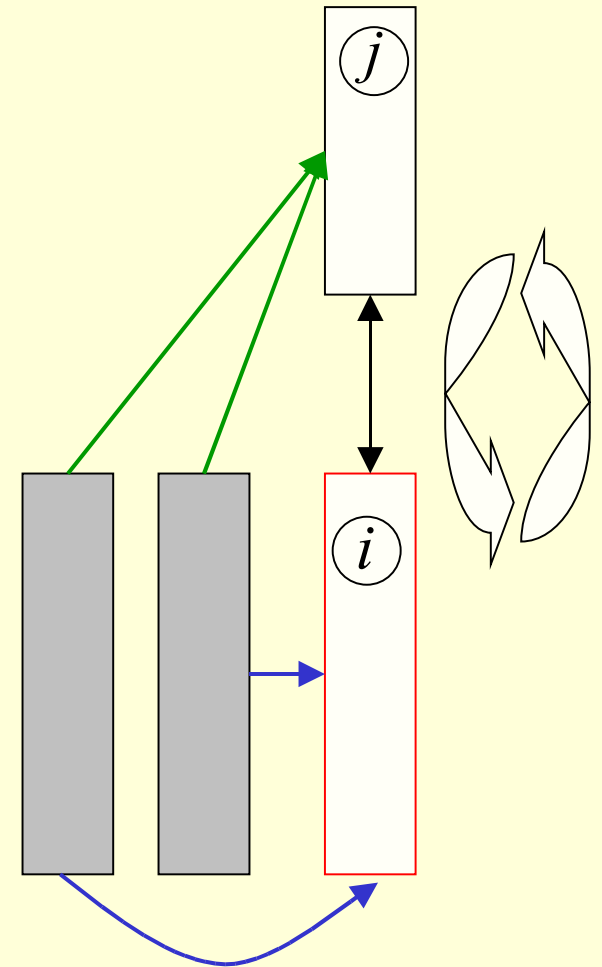
The conditional RBM model (a partially observed CRF)

- Start with a generic RBM.
- Add two types of conditioning connections.
- Given the data, the hidden units at time t are conditionally independent.
- The autoregressive weights can model most short-term temporal structure very well, leaving the hidden units to model nonlinear irregularities (such as when the foot hits the ground).



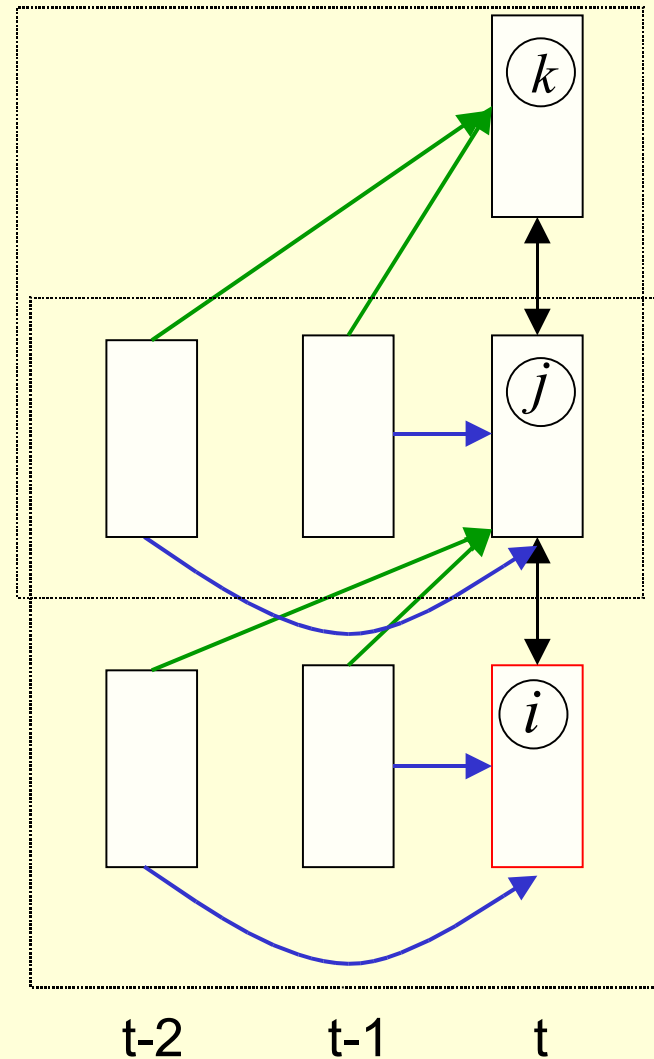
Causal generation from a learned model

- Keep the previous visible states fixed.
 - They provide a time-dependent bias for the hidden units.
- Perform alternating Gibbs sampling for a few iterations between the hidden units and the most recent visible units.
 - This picks new hidden and visible states that are compatible with each other and with the recent history.



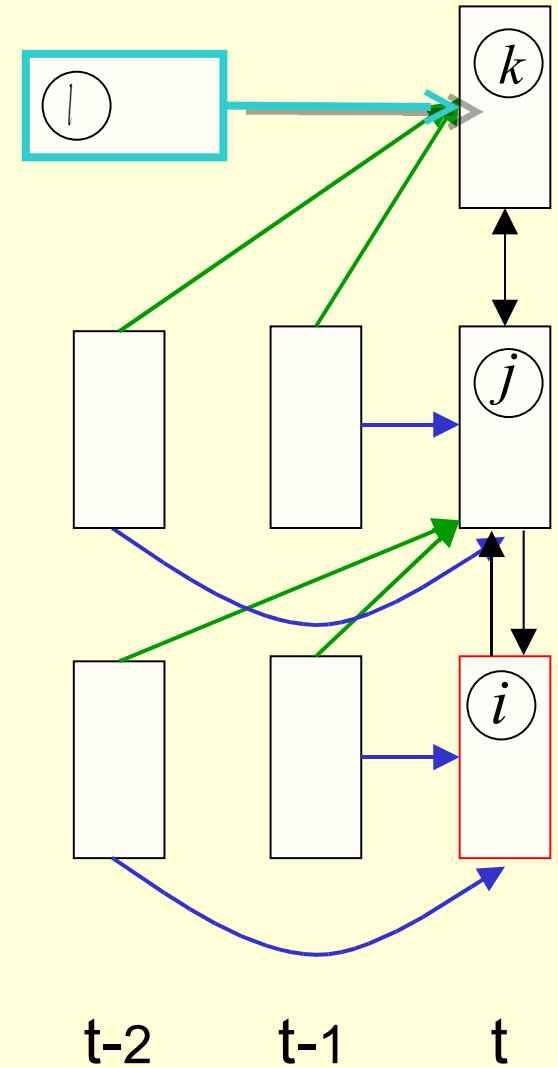
Higher level models

- Once we have trained the model, we can add layers like in a Deep Belief Network.
- The previous layer CRBM is kept, and its output, while driven by the data is treated as a new kind of “fully observed” data.
- The next level CRBM has the same architecture as the first (though we can alter the number of units it uses) and is trained the same way.
- Upper levels of the network model more “abstract” concepts.
- This greedy learning procedure can be justified using a variational bound.



Learning with “style” labels

- As in the generative model of handwritten digits (Hinton et al. 2006), style labels can be provided as part of the input to the top layer.
- The labels are represented by turning on one unit in a group of units, but they can also be blended.



Show demo's of multiple styles of walking

These can be found at
www.cs.toronto.edu/~gwtaylor/

Readings on deep belief nets

A reading list (that is still being updated) can be found at

www.cs.toronto.edu/~hinton/deeprefs.html