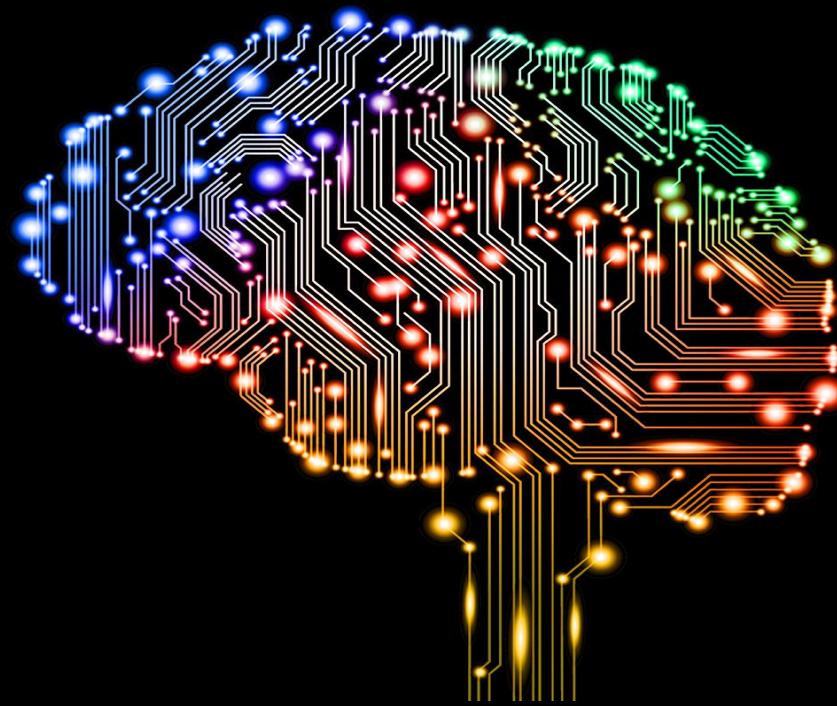


数据驱动的人工智能（6b）基于能量的神经网络

Data Driven Artificial Intelligence

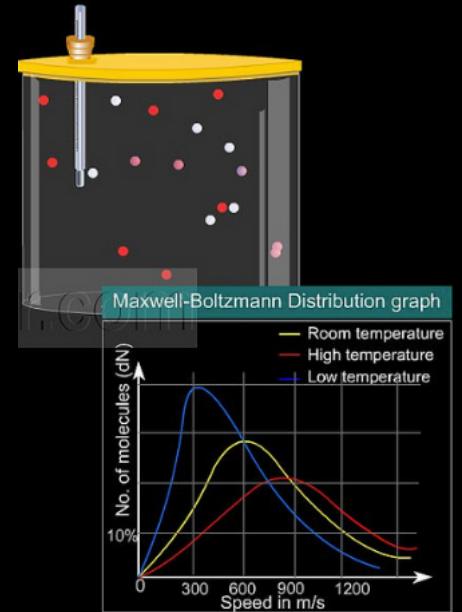
邬学宁 SAP硅谷创新中心

2017 / 03



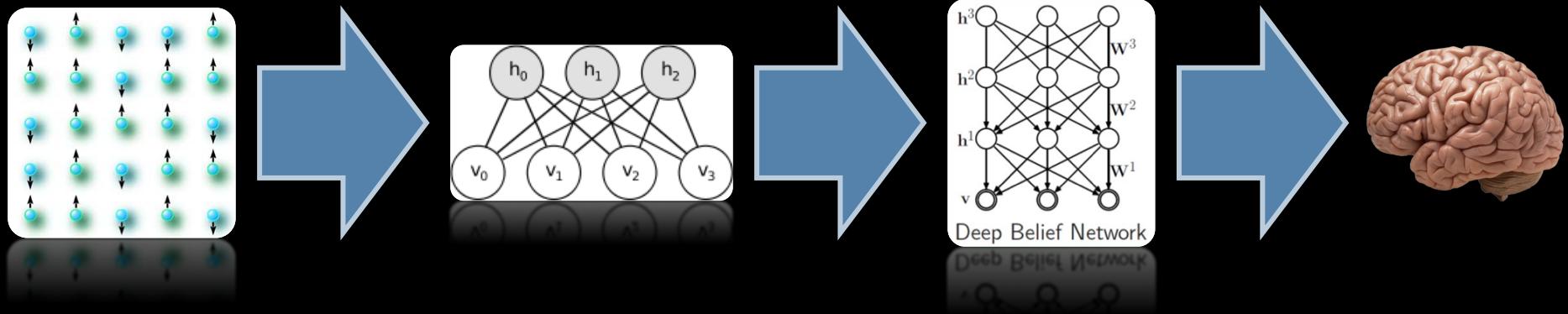
日程: EBN

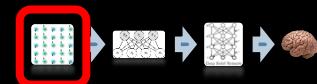
- Bayesian Belief Network
- Auto-encoder / PCA
- Hopfield Network
- Boltzmann Machine
- Restricted Boltzmann Machine
- Deep Boltzmann Machine
- Deep Belief Network
- Sparse Coding
- Game Theory & Generative Adversarial Network



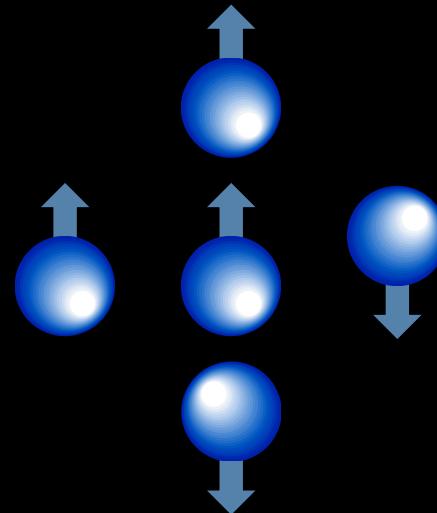
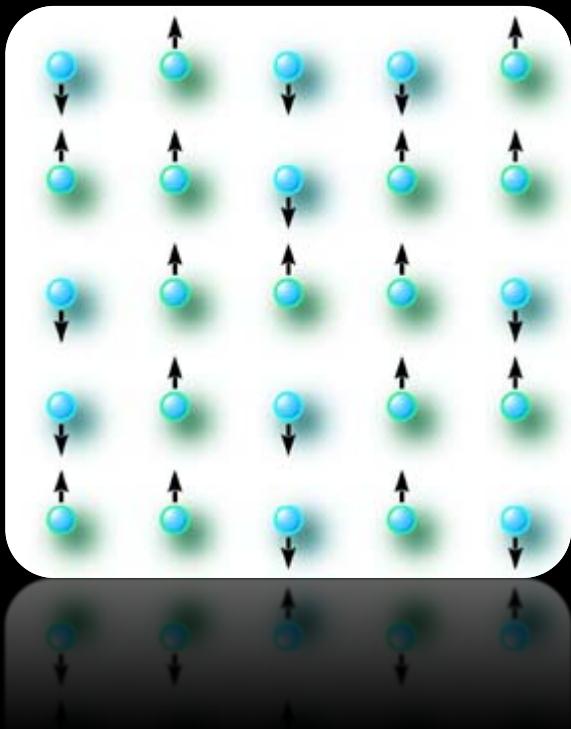


From Ising Model to brain...



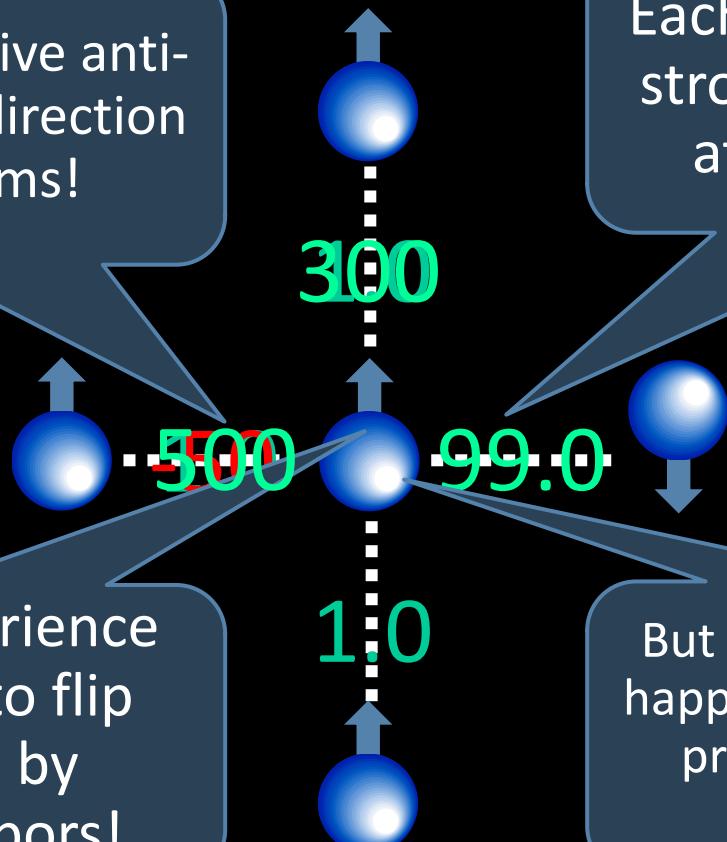


The Ising Model (Ernst Ising, 1922)



The Ising Model

Negative weights drive anti-correlations in the direction of adjacent atoms!



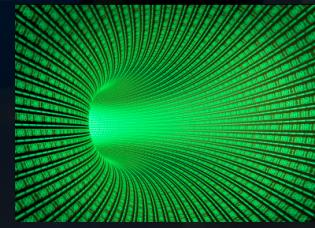
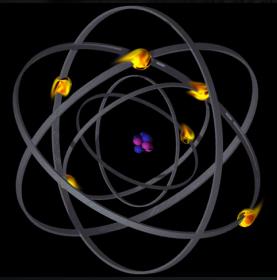
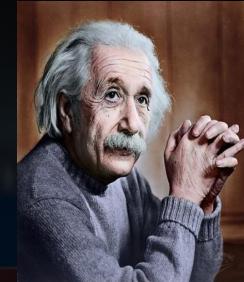
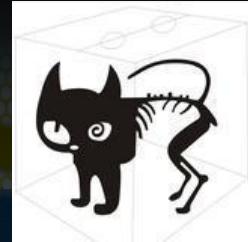
Each weight indicates how strongly the two adjacent atoms are correlated.

And atoms experience peer-pressure to flip when coaxed by multiple neighbors!

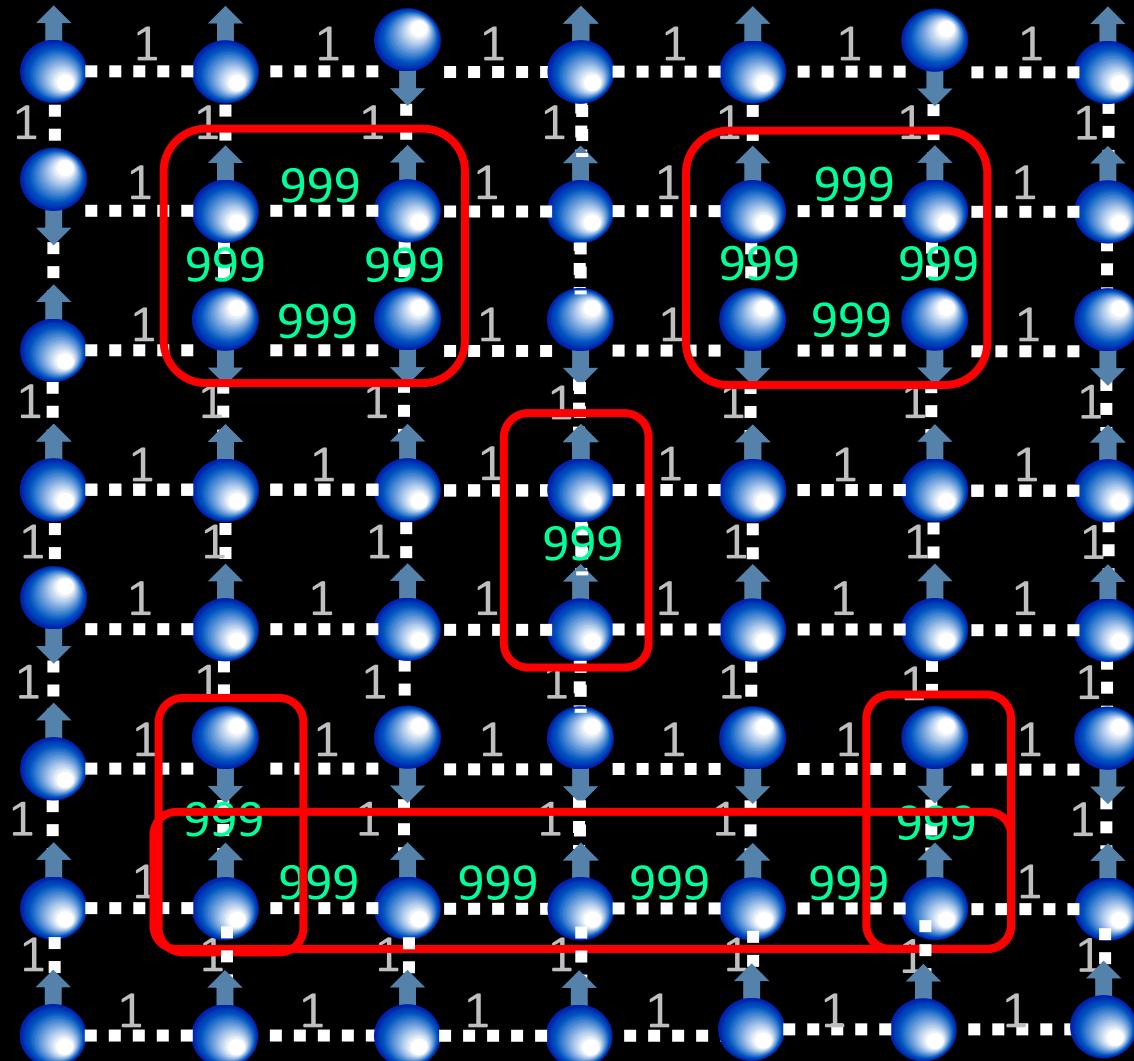
But ultimately, anything can happen – the interactions are probabilistic rather than deterministic!

偏题：人类社会范式转移：从决定论到概率论 / 从原子到比特

宇宙间，一个技术文明等级的重要标志，是它能够控制和使用的微观维度



Let's See a Simulation!



Ising Model Energy

Every configuration of atoms and weights has a specific “energy” associated with it.

give each up-facing “atom” a value of $x = +1$

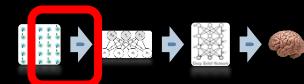
give each down-facing “atom” a value of $x = -1$

And between atoms x_i and x_j we'll call the weight w_{ij} .

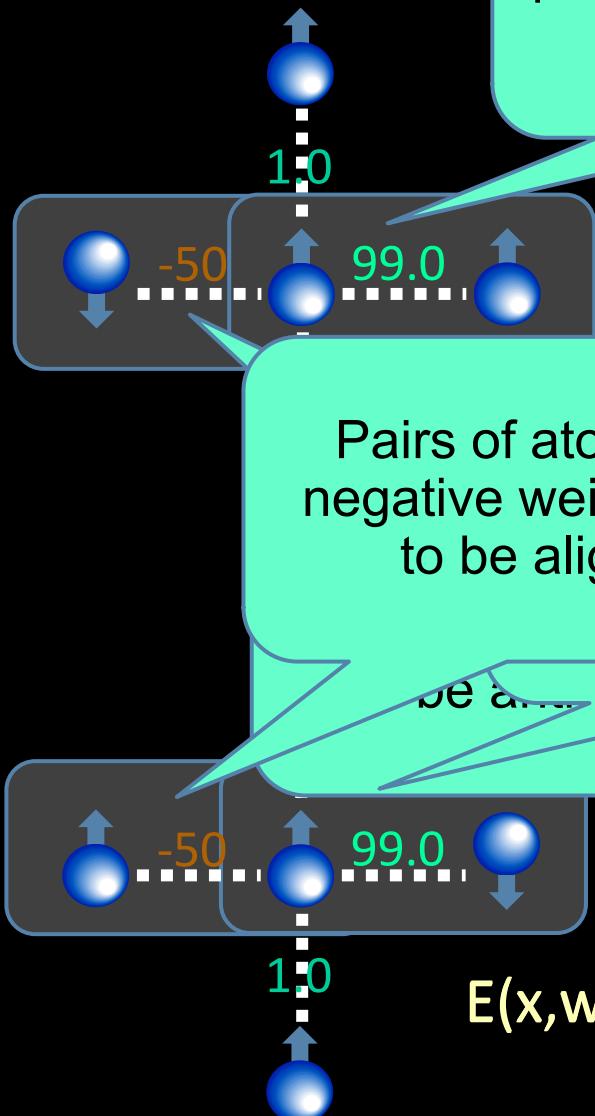
We can then define the “energy” of a configuration of atoms and weights...

$$E(x,w) = - \sum_{\text{all adjacent pairs } i,j \text{ of atoms}} x_i * x_j * w_{ij}$$

$$E(x,w) = -((1*1*1.0) + (1*(-1)*-50) + (1*(-1)*1.0) + (1*1*99))$$



Ising Model Energy



Pairs of atoms with positive weights tend to be aligned.

$$E(x, w) = - \sum_{\text{all adjacent pairs } i,j \text{ of atoms}} x_i * x_j * w_{ij}$$

Low-energy configurations are those where atoms connected by positive weights are aligned, and atoms connected by negative weights are anti-aligned.

Pairs of atoms with negative weights tend to be aligned.

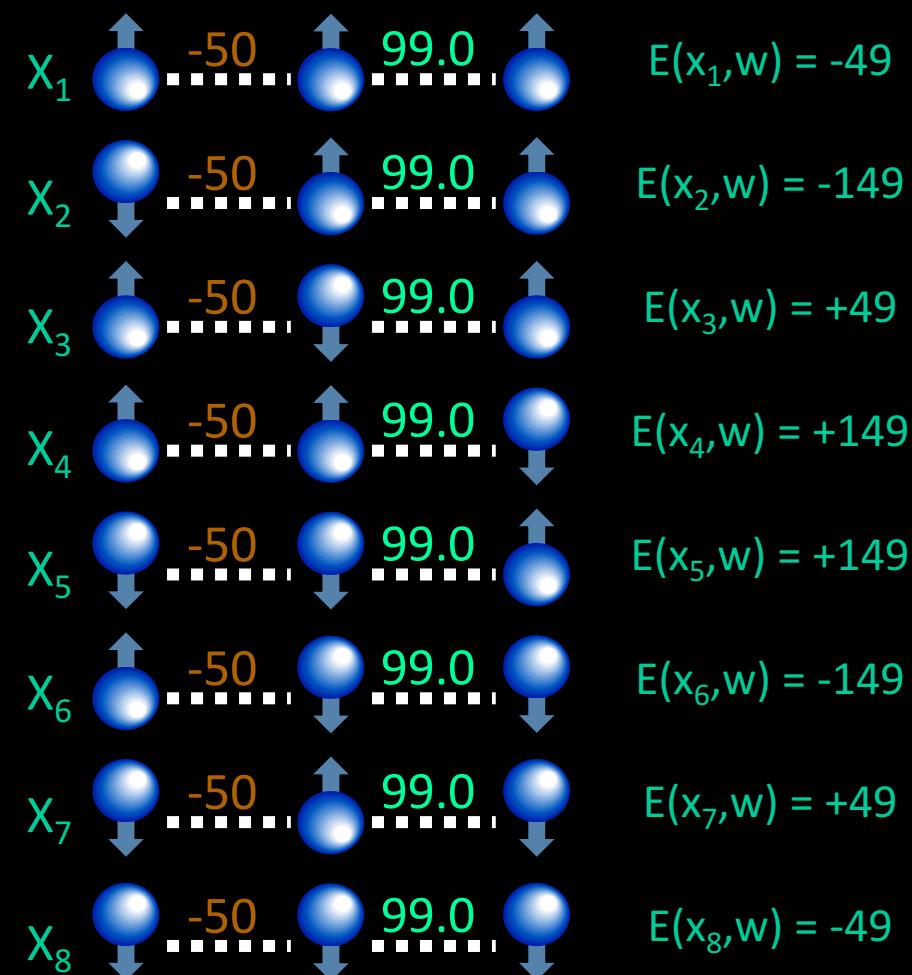
Pairs of atoms with positive weights tend to be aligned.

These configurations tend to be more “stable.”

High-energy configurations are the opposite – their alignment goes against the weights.

These configurations tend to be more “unstable.”

Ising Model Probabilities

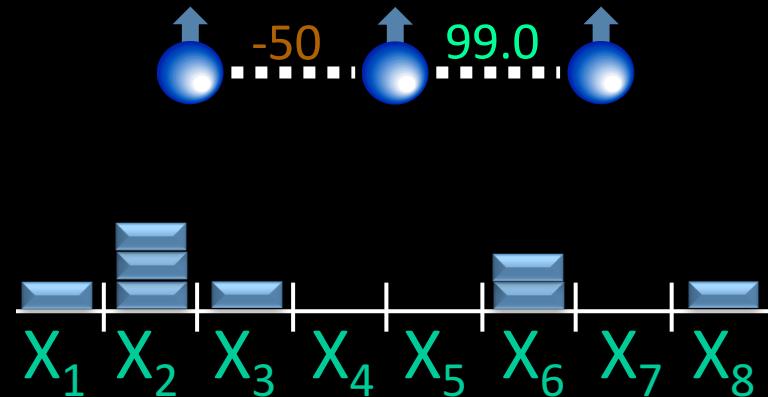


Consider an Ising model with 3 atoms...

There are 2^3 or 8 possible configurations of the atoms.

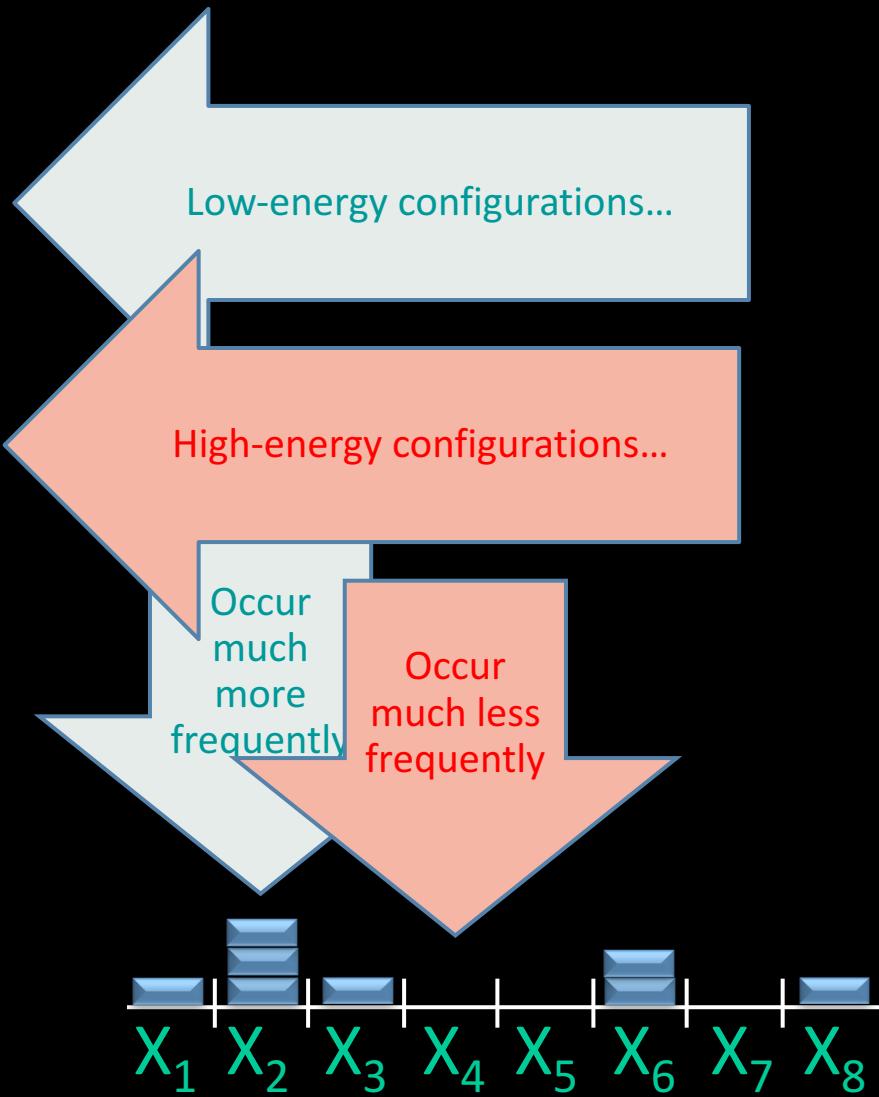
If we let our atoms vibrate (flip up and down) naturally, they will oscillate through these various configurations...

And the probability of a particular configuration occurring depends entirely on the energy of that configuration vs. the energy of all of the other configurations.

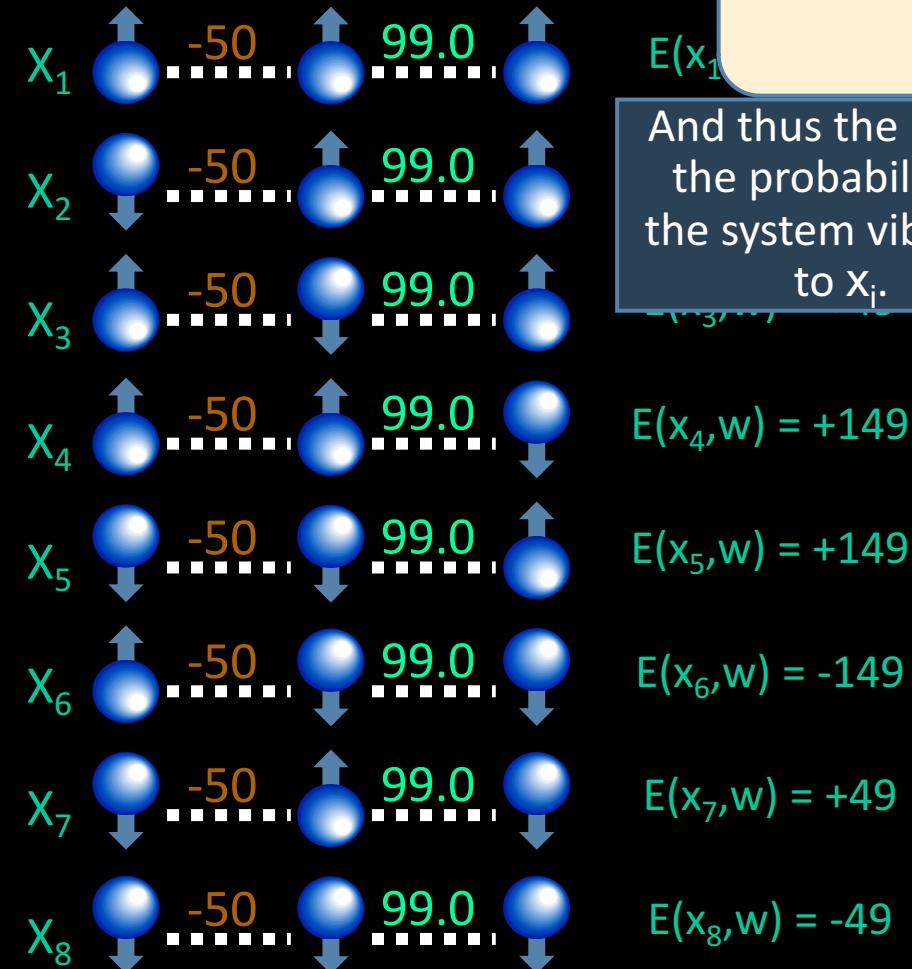


Ising Model Probabilities

x_1		-50		99.0		$E(x_1, w) = -49$
x_2		-50		99.0		$E(x_2, w) = -149$
x_3		-50		99.0		$E(x_3, w) = +49$
x_4		-50		99.0		$E(x_4, w) = +149$
x_5		-50		99.0		$E(x_5, w) = +149$
x_6		-50		99.0		$E(x_6, w) = -149$
x_7		-50		99.0		$E(x_7, w) = +49$
x_8		-50		99.0		$E(x_8, w) = -49$



Ising Model Probability



And thus the higher
the probability of
the system vibrating
to x_i .

$$E(x_4, w) = +149$$

$$E(x_5, w) = +149$$

$$E(x_6, w) = -149$$

$$E(x_7, w) = +49$$

$$E(x_8, w) = -49$$

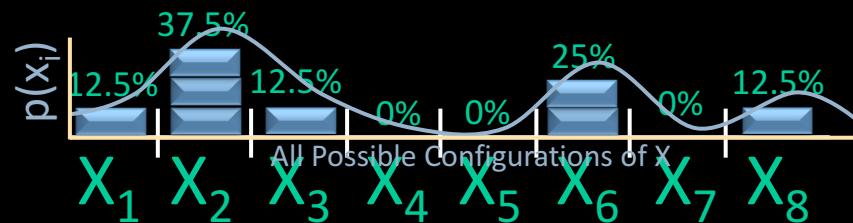
$$E(x_1, w)$$

$$p(x_i | w) = \frac{e^{-E(x_i, w)}}{\sum_{k=1}^n e^{-E(x_k, w)}}$$

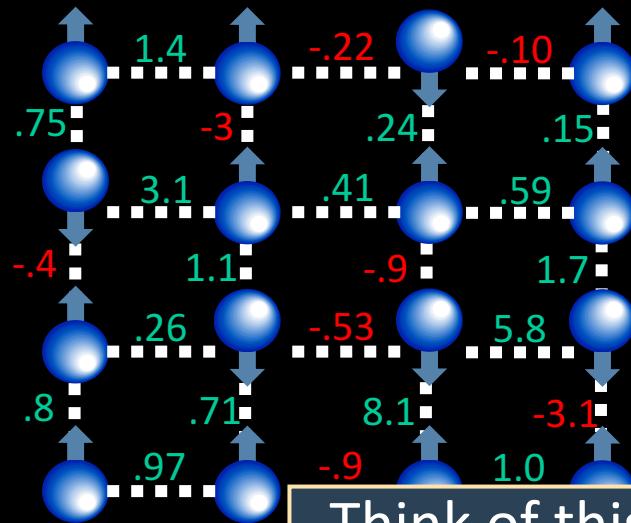
$p(\text{the system vibrates to config. } x_i) =$

$$\frac{\text{eNRG of config. } x_i}{\text{sum of eNRG across all configs}}$$

The probability distribution described by
this formula, and shown below, is called a
“Boltzmann Distribution”.



A 3-D Ising Model



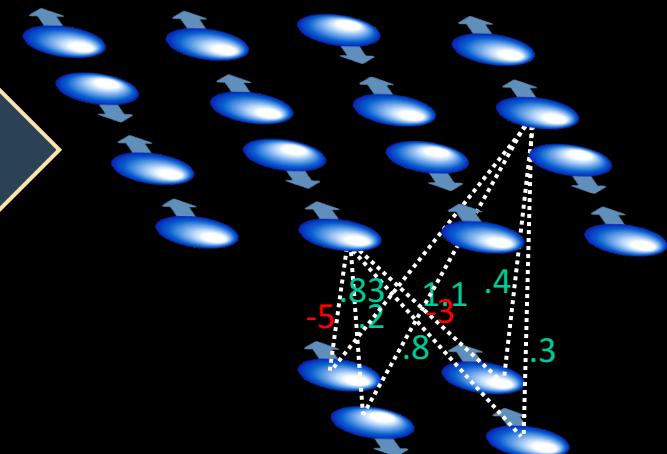
What if we swap
a "visible" plane

Think of this plane
like the pixels on a
touch-screen

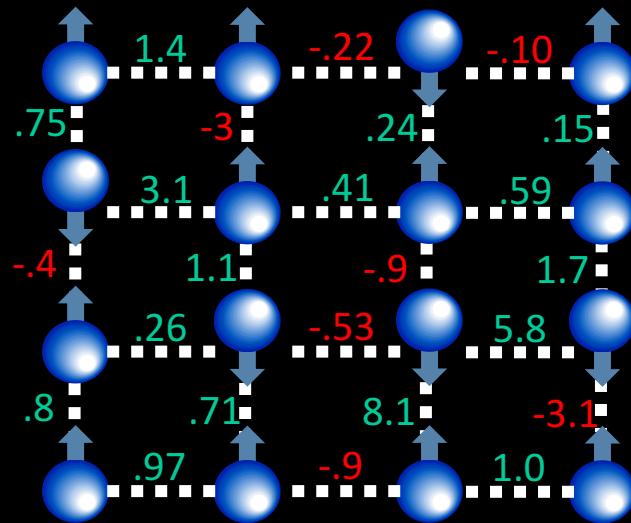
In this case, we'll have weights between
each atom on the visible plane and every
atom on the hidden plane.

Atoms on each plane are now only
influenced by the atoms on the other plane.

Now what if instead of having the
weights between adjacent atoms in a
plane, we changed things up a bit...



A 3-D Ising Model

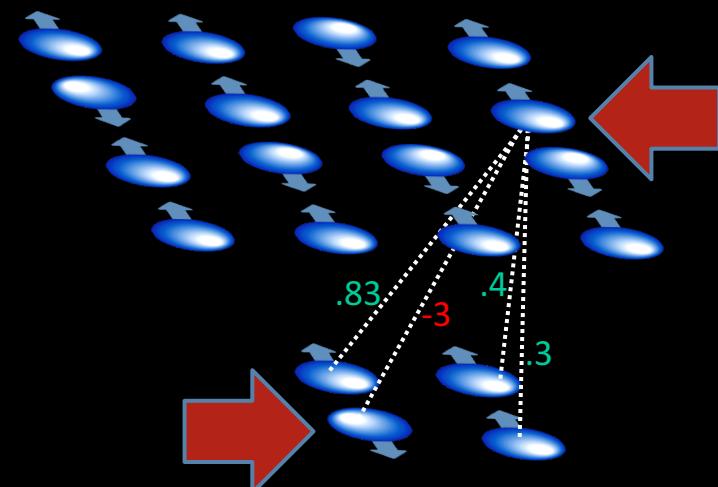


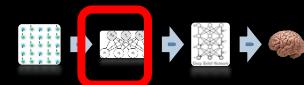
Now what if instead of having the weights between adjacent atoms in a plane, we changed things up a bit...

What if we switched to 2 planes...
a “visible” plane and a “hidden” plane

In this case, we’ll have weights between
each atom on the visible plane and every
atom on the hidden plane.

Atoms on each plane are now only
influenced by the atoms on the other plane.



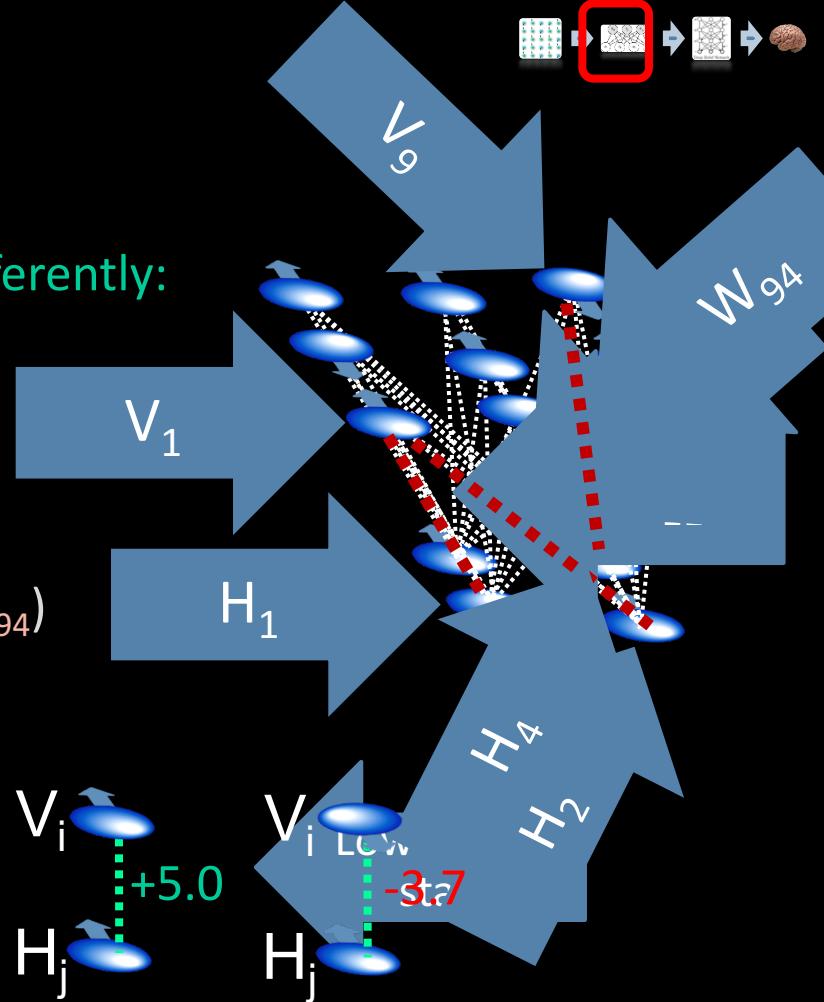


A 3-D Ising Model

In such a model, energy is computed a bit differently:

$$E(V, H, W) = - \sum_{i=1}^V \sum_{j=1}^H V_i * H_j * W_{ij}$$
$$= - (V_1 * H_1 * W_{11} + V_1 * H_2 * W_{12} + \dots + V_9 * H_4 * W_{94})$$

Just like before, if we have a visible at odd index facing the positive direction, it is oriented by a negative weight. This also contributes to contributing to low energy...



And as before, configurations with lower energy will occur with higher probability than those with higher energy.

A 3-D Ising Model

As before, once we can find the probability of vibrating to this specific configuration, we can also find the probability of vibrating to any other configuration.

“What are the odds that these 9 visible atoms will vibrate to this specific configuration?”

$p(\text{the system vibrates to a given 3-D config}) =$

$$\frac{\text{eNRG of that configuration}}{\text{Total eNRG across all possible configs}}$$

And we can compute the probability of the system vibrating to a specific visual configuration irrespective of what the hidden side looks like...

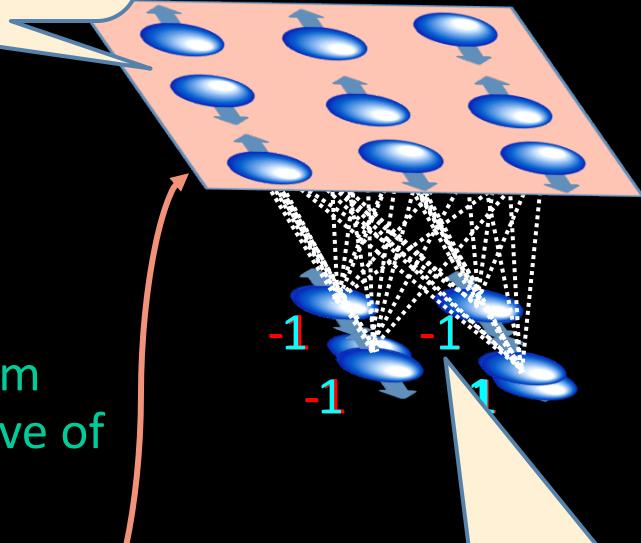
To do so, we simply sum up the probability of config v showing up for *all* possible hidden configurations:

$p(\text{the system displays visible config } v) =$

$p(\text{the system displays } v \text{ when } h \text{ is } \{-1, -1, -1, -1\}) +$
 $p(\text{the system displays } v \text{ when } h \text{ is } \{-1, -1, -1, 1\}) +$

...

$p(\text{the system displays } v \text{ when } h \text{ is } \{1, 1, 1, 1\})$

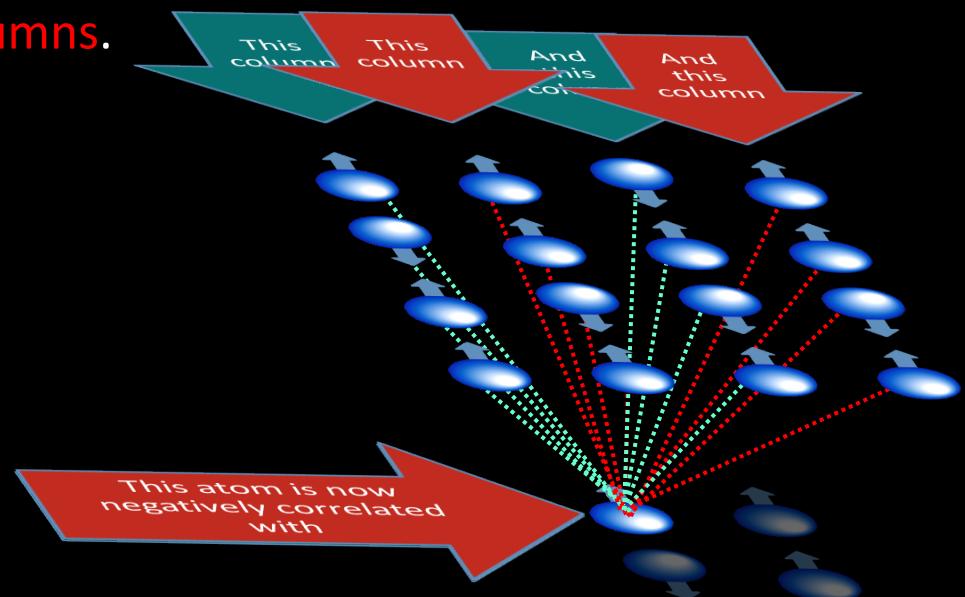


“Irrespective of what these atoms look like...”

A 3-D Ising Model

Just as with our 2-D Ising model,
we can pick weights to bias certain
configurations of atoms.

For instance, we can **positively connect** the
upper-left hidden atom to **odd columns**...
and **negatively tie** it to **even columns**.



A 3-D Ising Model

Just as with our 2-D Ising model, we can pick weights to bias certain configurations

For instance, we can **positively connect** the upper-left hidden atom and **negatively tie** it to even columns.

We'd now expect our visible layer to display stripes (on, off, on, off) more often than the weights were

But why?

That'll pressure the positively-connected hidden atom to flip too!

If a hidden atom just happens to randomly flip up...

And similarly, if a bunch of visible atoms just happen to randomly flip in a particular direction...

Further reinforcing alignment on the

This creates a self-reinforcing set of stable configurations!

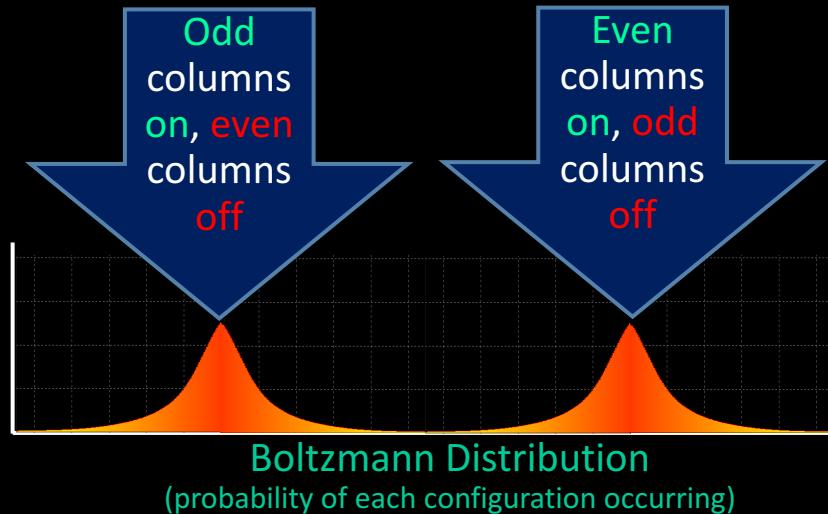
A 3-D Ising Model: A Generative Model

In fact, we say that our 3-D Ising model is “generative”.

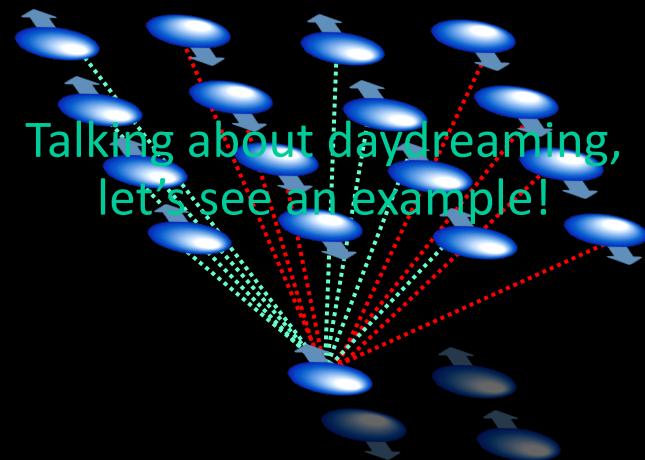
We call it generative because if we allow the system to vibrate...

it will naturally tend to generate patterns on its visible plane...

based on the specific weights connecting the atoms in the hidden and visible planes.



When left to vibrate, it's almost as if the system is daydreaming!





A 3-D Ising Model

Here's a demo with 5 atoms in the hidden layer:

Hidden Atom Has **positive** weights to ...



two sets of “eye” atoms near the top of the visible plane.



a set of “nose” atoms in the middle of the visible plane.



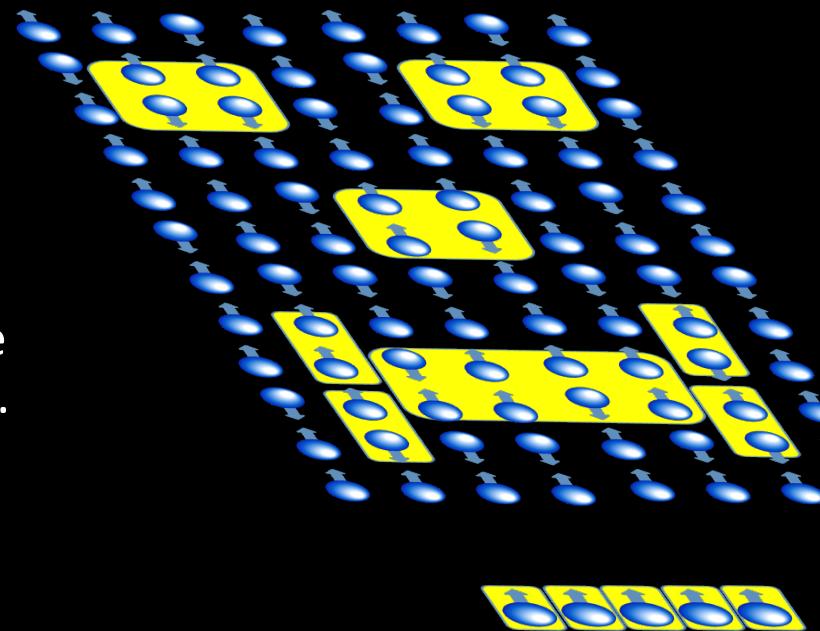
a line of “mouth” atoms in the lower part of the visible plane.



a set of “happy” atoms just above the mouth.



a set of “frown” atoms just below the mouth.



So by manually adjusting the weights between the two planes, we can bias the distribution of what “images” are generated on the visible plane.



An Aha Moment: The 3-D Ising Model becomes the Restricted Boltzmann Machine (RBM)



Geoffrey Hinton



Terry Sejnowski

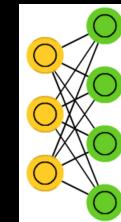
“Let’s use a simulation of an Ising network to encode memories...”

- Approximate quote fabricated by Carey



The Restricted Boltzmann Machine

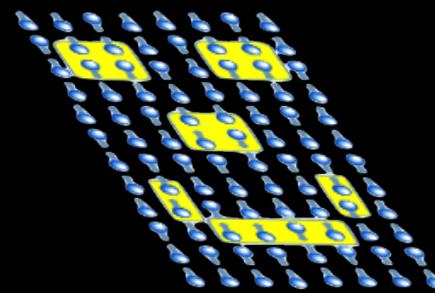
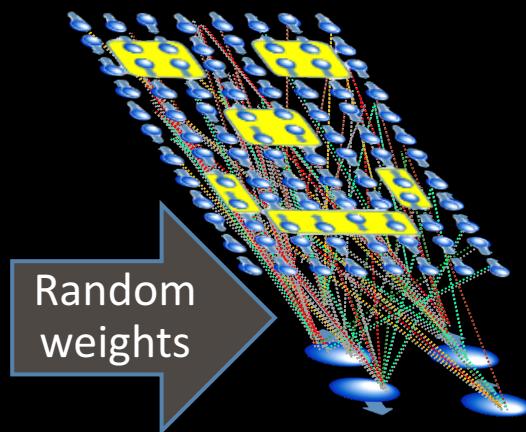
Learning an Image and Daydreaming



In the early 1980's, Hinton and Sejnowski started experimenting with 3-D Ising Models, also called Restricted Boltzmann Machines (RBMs).

Their goal was to come up with an algorithm that could take a "blank" RBM network...

And teach that network (by modifying its weights) how to "encode" an image...



So if allowed to vibrate, the RBM would vibrate up the encoded image on its visible plane with higher probability than other images...

"I p

"Right! Then the Ising model will naturally vibrate to our training image more often."

"Well, as we know, an Ising model naturally vibrates to low-energy configurations more frequently than to high-energy ones..."

"Remember a picture?"

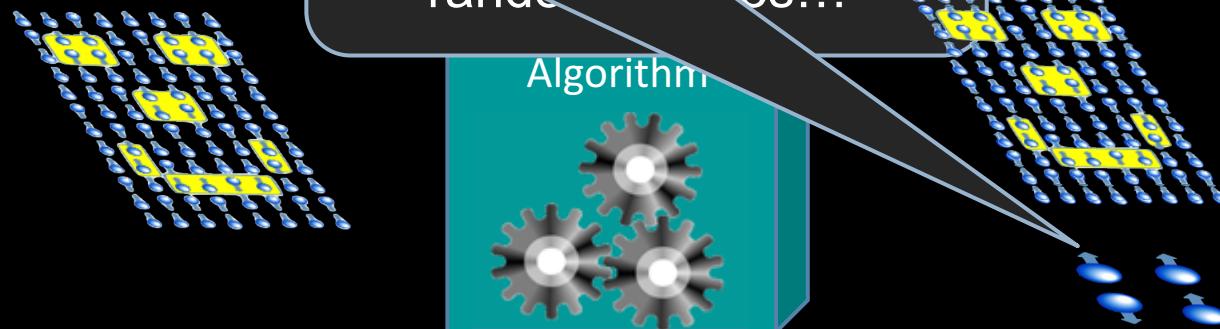
And then find some automated way,

The resulting Ising model is a "generative model!"

Why? Because once it's been trained, if we allow it to vibrate through its natural states, it'll tend to generate the image that it's learned!

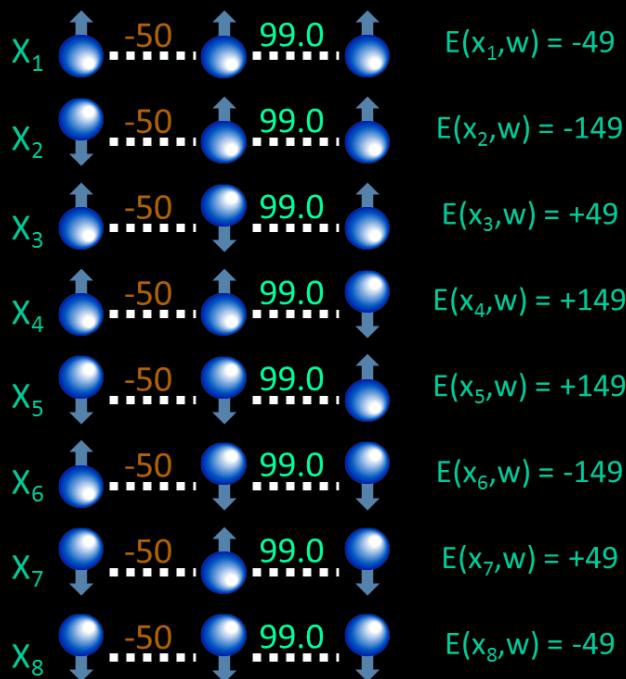
on Ising model

• pick a set of images our image as low energy as possible."



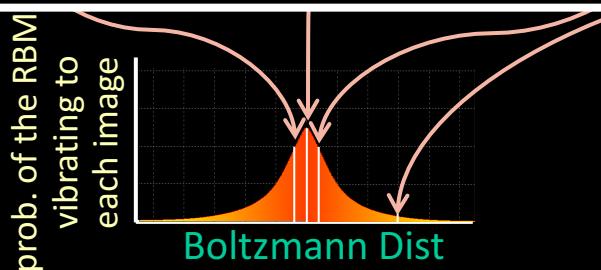
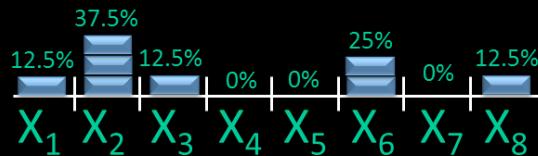


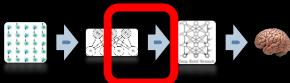
Ising Model Probabilities



The probability of vibrating to each configuration x_i (given fixed weights w) can also be predicted by this formula:

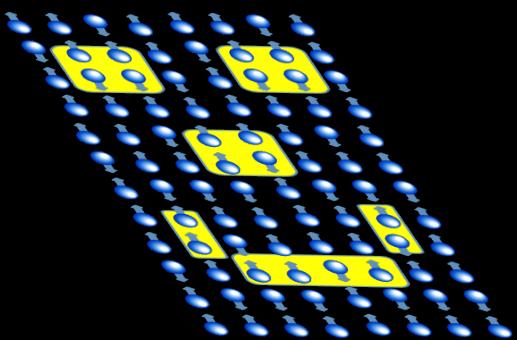
$$p(x_i | w) = \frac{e^{-E(x_i, w)}}{\sum_{k=1}^n e^{-E(x_k, w)}}$$





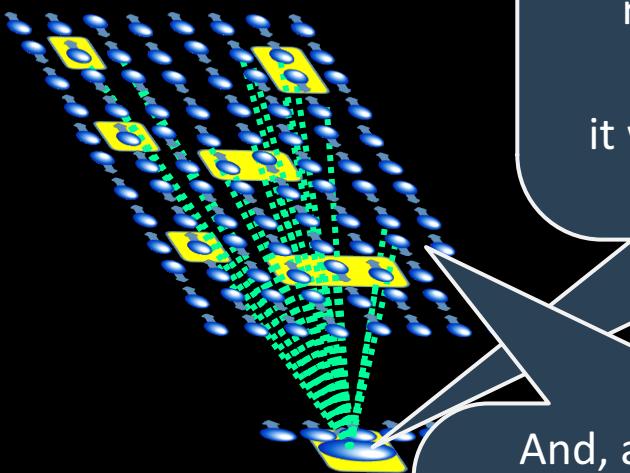
Learning an Image - Intuition

So what is actually happening to the weights
Gradient Descent algorithm?



Training Image

Let's imagine for a moment that we only have one hidden atom...
...our Gradient Descent algorithm would likely move the weights to positive atoms that need to be activated and negative weights (not depicted).



This will bias “vibration” of the system to configs resembling the trained image.

Why? Because in those random cases where our hidden atom activates... it will generate our image on the visible plane.

And, any time the visible plane vibrates into a config that (remotely) resembles our training image... this will likely result in the activation of our hidden atom... which will then increase the odds of the trained image showing up, etc...

"OK, so we know how to train an RBM on an image... But can it learn (to date) multiple images?"

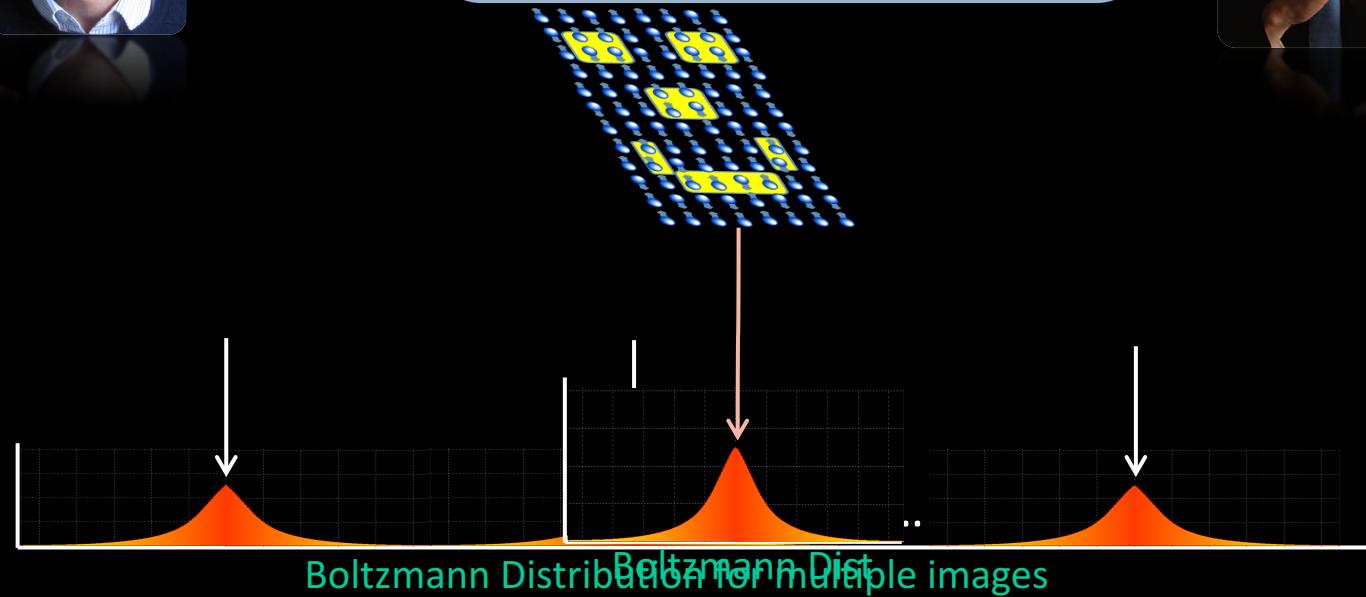
"You're saying it's possible to find a single set of weights that will create multiple such low-energy points?"

"Why not? Instead of using our Gradient Descent algorithm to learn weights to create a single low-energy state..."



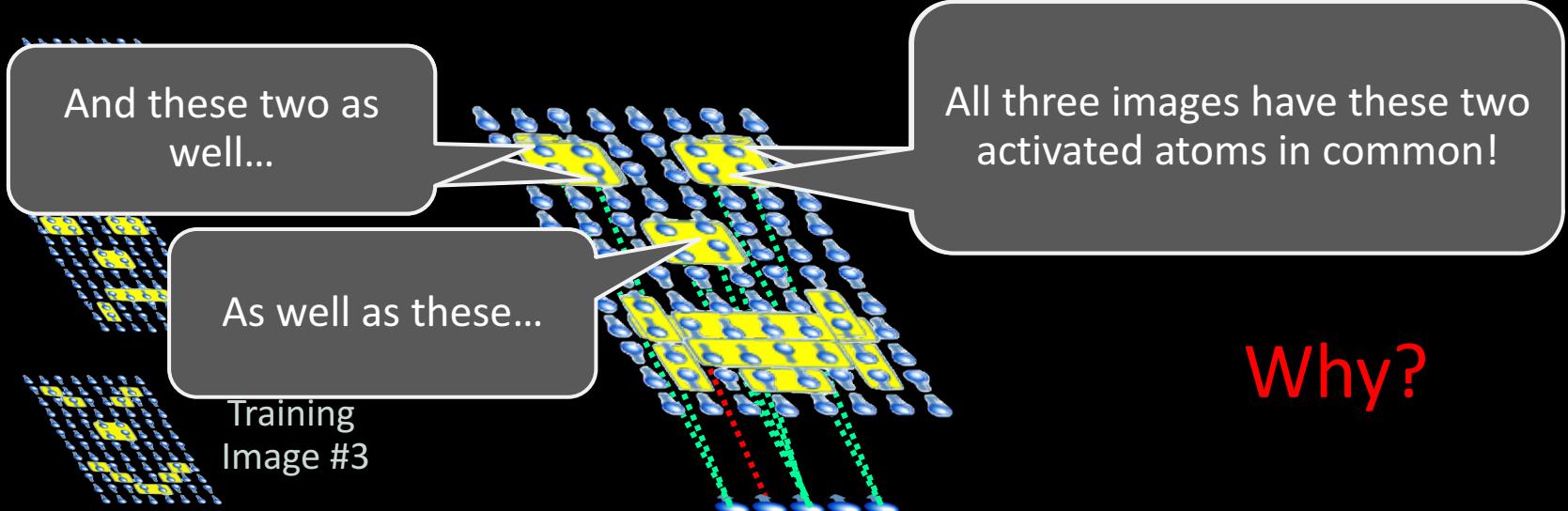
w
m

"That's correct – if you have enough atoms/weights, the system can represent many low-energy states."



Learning Multiple Images - Intuition

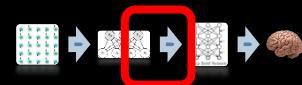
So intuitively, how does the Gradient Descent algorithm discover a combination of weights that creates low energy states for multiple images?



Notice that some atoms are actually consistently on across multiple images...

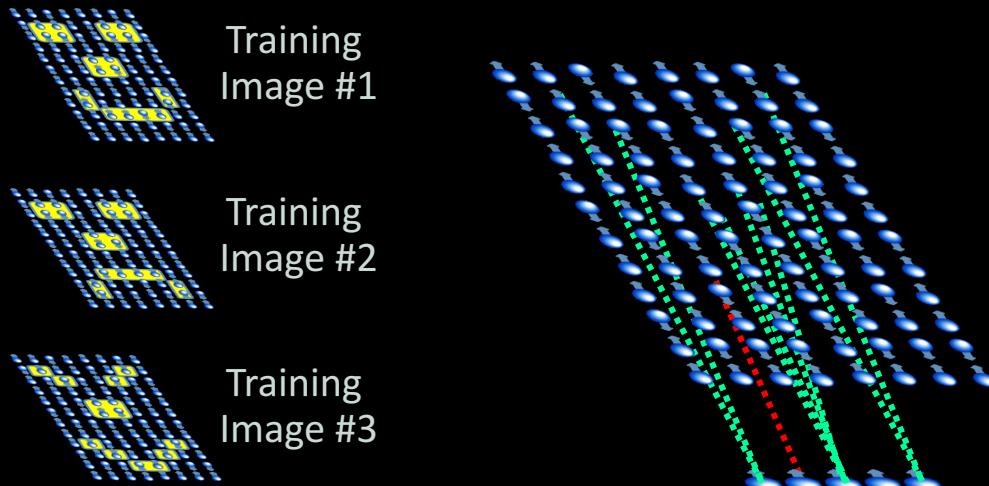
Our Gradient Descent algorithm identifies these “regularities” and ties them to hidden atoms with strong positive weights!

It’ll similarly tie visible atoms that are consistently off across multiple training images to hidden atoms with strong negative weights.



Learning Multiple Images - Intuition

So intuitively, how does the Gradient Descent algorithm discover a combination of weights that creates low energy states for multiple images?

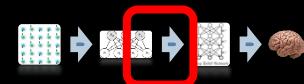


Why?

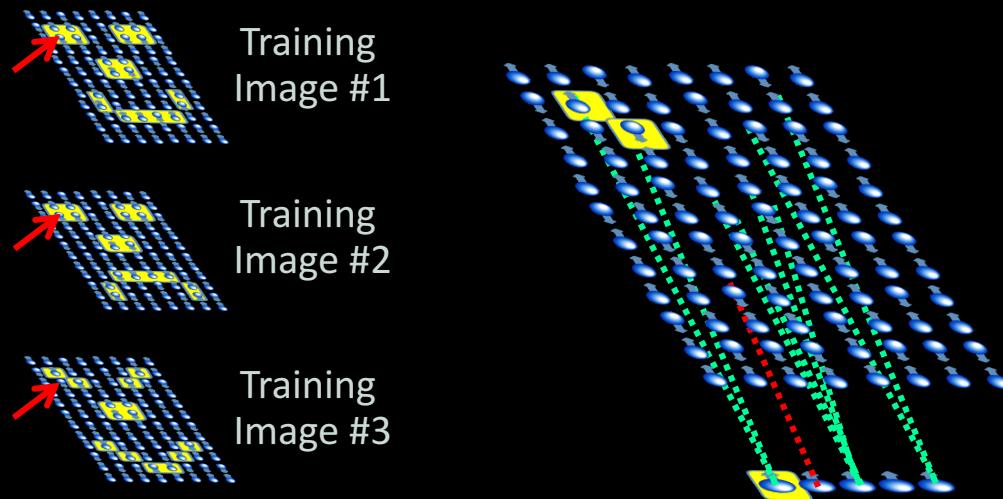
Notice that some atoms are actually consistently on across multiple images...

Our Gradient Descent algorithm identifies these “regularities” and ties them to hidden atoms with strong positive weights!

It’ll similarly tie visible atoms that are consistently off across multiple training images to hidden atoms with strong negative weights.



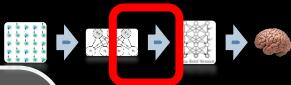
Learning Multiple Images - Intuition



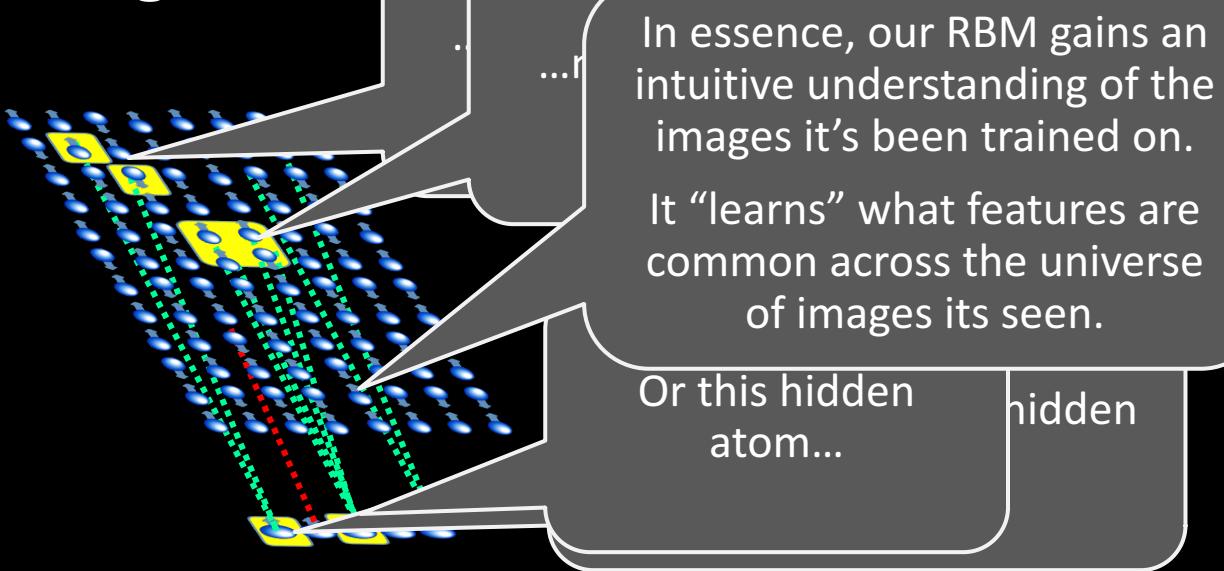
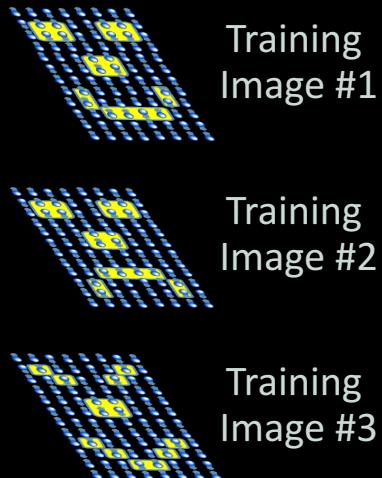
Associating a positive weight between a hidden atom
and one or more visible atom(s) that are frequently on
across multiple training images...

means that any time that hidden atom activates,
there's a good chance it'll also activate visible atom(s) that are useful in
reconstructing many of our training images ...

this increases the odds that one of those multiple trained images
will vibrate up on the RBM's visible surface!



Learning Multiple Images - Int

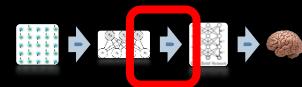


So the gradient descent process identifies “regularities” – common clusters of activated (or deactivated) visible atoms – across multiple training images...

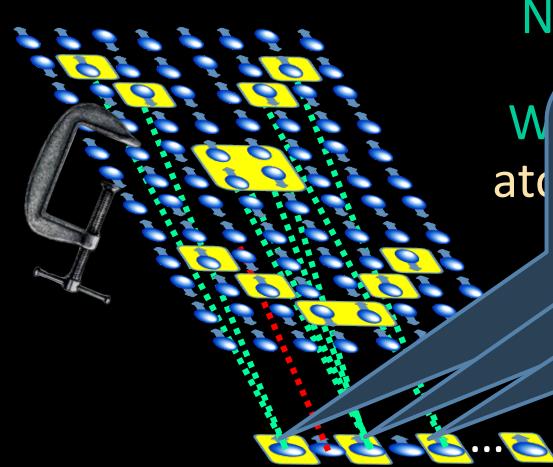
And it associates these via stronger weights with various hidden atoms.

In practice, the identified regularities will be very simple – not complete eyes but rather simple edges and curves.

It's this exploitation of regularities that causes our RBM to display its learned images more frequently than arbitrary/random images.



An Interesting Intuition



Now once we've trained our RBM on a bunch of images...

We can...
atc...
to

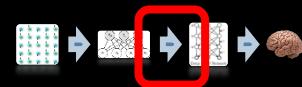
And it appears to have
part of a right eye...

As expected, this would cause our hidden atoms to
adjust based on our learned weights...

ally flip our visible
hold them steady...
s they like?

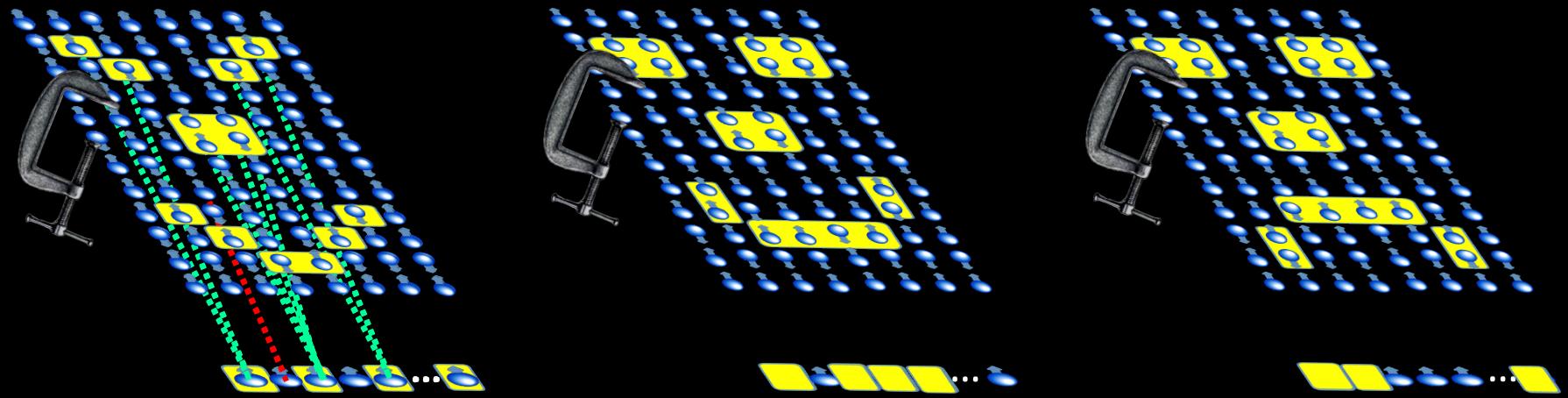
In essence, the resulting hidden atom configuration is a “code” that describes the set of regularities that were found within the visible image.

The hidden atom configuration thus represents a “high-level description” of the original image.



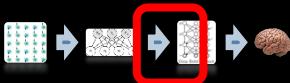
An Interesting Intuition

OK, let's take our trained RBM and use it to compute codes for each training image.



So now we have high-level descriptions for each of our original images:

Let's do something wacky with them!



An Interesting Intuition

... to these specific configurations of atoms.

create a
Is
Recall that's
associate

And, as before, to do so, we'll need to find regularities across our inputs...

mann Machine

This RB
ut train

So we'll learn the

... and associate them with various hidden atoms.



But wait – when we train our RBM, its visible atoms don't represent raw dots; they represent raw dots plus visible atoms don't activate when an image has two eyes!

Instead, as we learned earlier, these atoms represent more abstract concepts like parts of faces, curves of the mouth, etc.

As such, when our second RBM discovers regularities in its inputs, these will be even higher-order regularities - for instance that every face has two eyes.

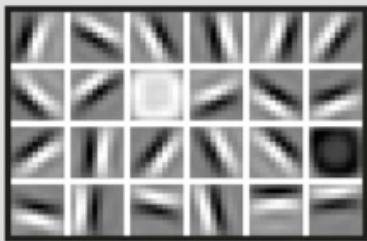
FACIAL RECOGNITION



Deep-learning neural networks use layers of increasingly complex rules to categorize complicated shapes such as faces.



Layer 1: The computer identifies pixels of light and dark.



Layer 2: The computer learns to identify edges and simple shapes.



Layer 3: The computer learns to identify more complex shapes and objects.



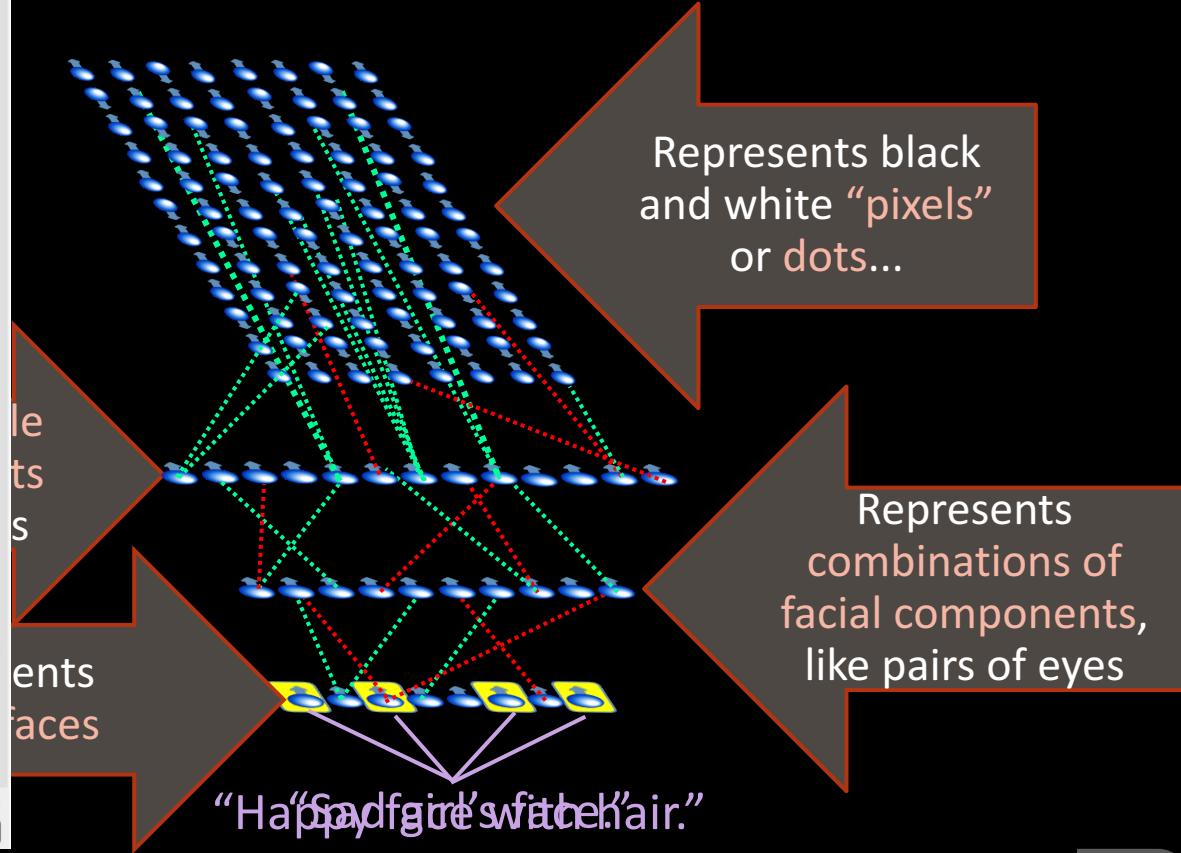
Layer 4: The computer learns which shapes and objects can be used to define a human face.

IMAGES: ANDREW NG

Neural Networks

You can repeat this process over and over!

In an RBM layer, its hidden atoms “learn” successive abstractions about the original images!



The Deep Belief Network

This network is called a “Deep Belief Network.”

Once we have such a network, we can present an image to it at the top visible layer and fix it in place...

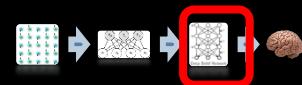
And then see how the atoms in the lower levels activate!

If trained with enough input images...

the DBN can actually learn to recognize common elements in those images (e.g., cats)...

without ever being told what they are!

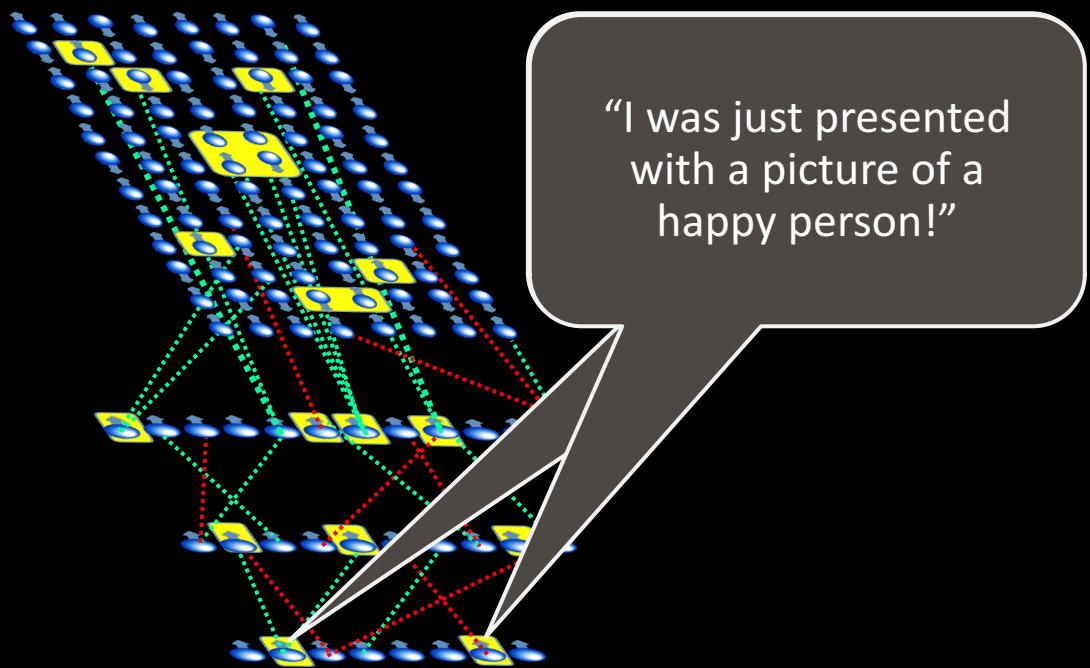




The Deep Belief Network: Recognition of Noisy Inputs

And the network is resilient to noise!

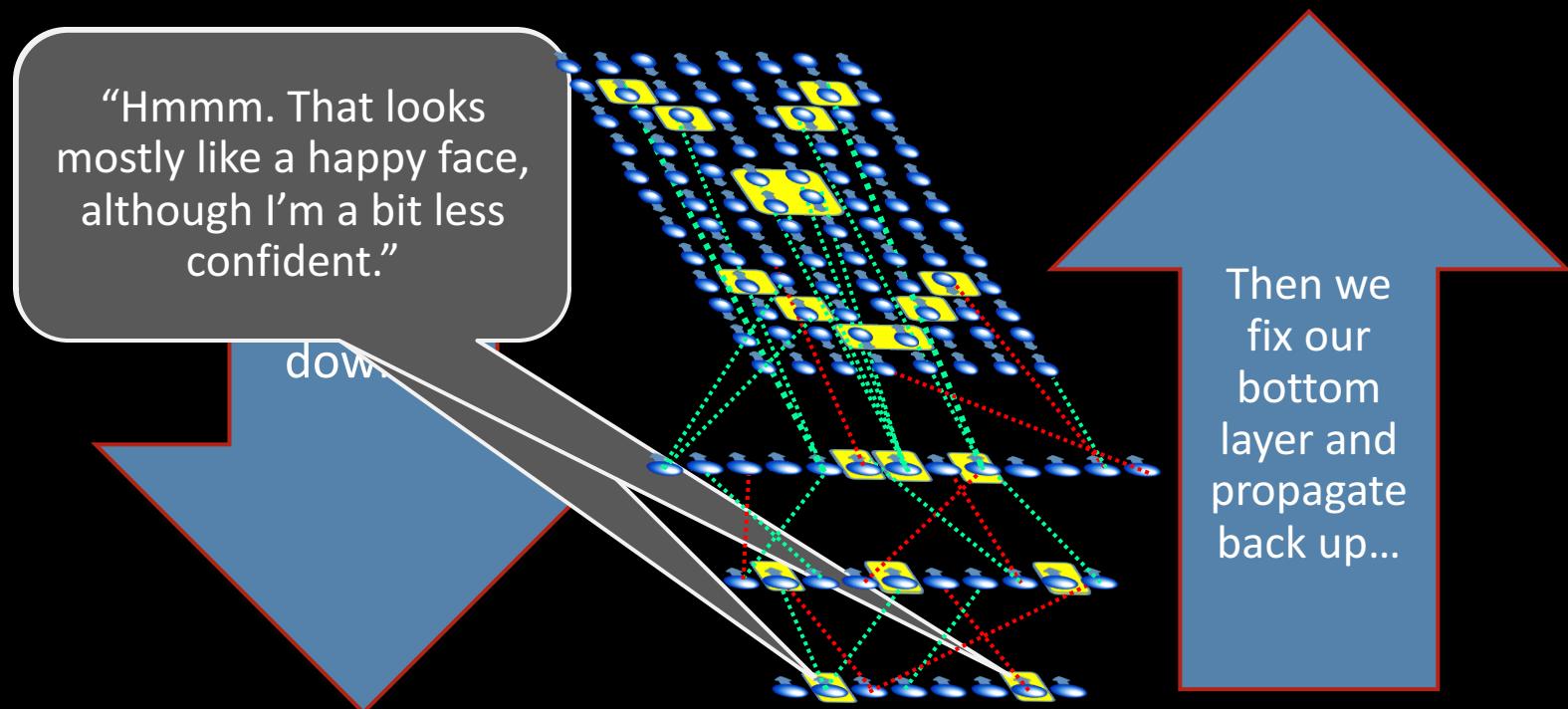
Even when presented with a partial image or a similar (but not identical) image, the network generally activates as expected.



The Deep Belief Network: Reconstruction

Even more exciting, the network can reproduce the original image from an incomplete or partial input!

All you have to do is propagate down and back up!

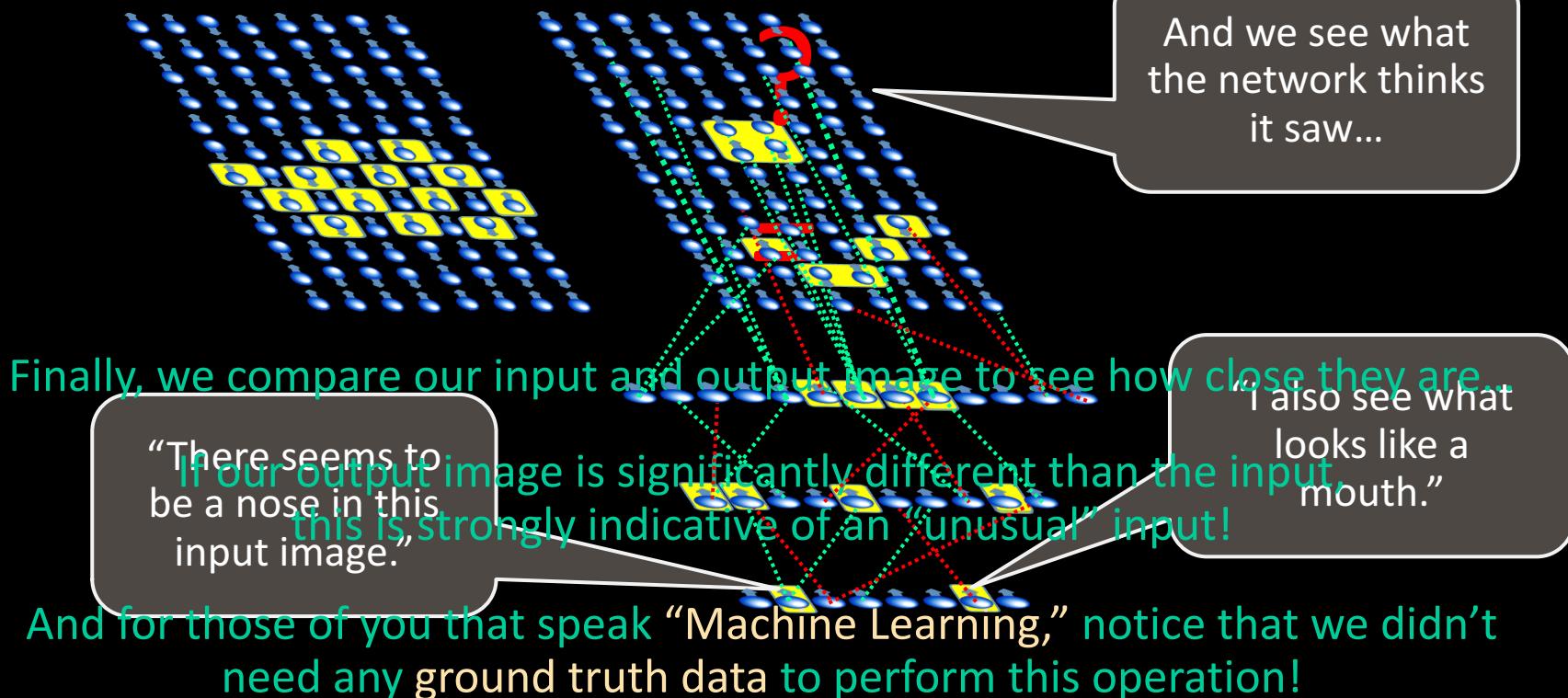


The Deep Belief Network: Anomaly Detection

Finally, we can use our DBN to detect anomalous inputs that don't match what our network has been trained on!

We take a new input image and propagate down...

Then we propagate back up.

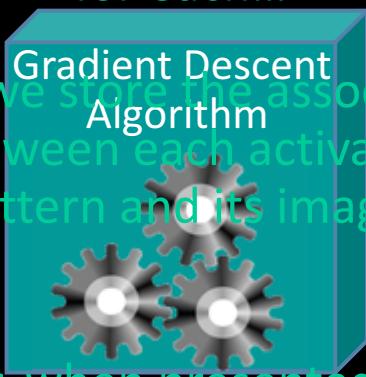


The Deep Belief Network: Fuzzy Searching

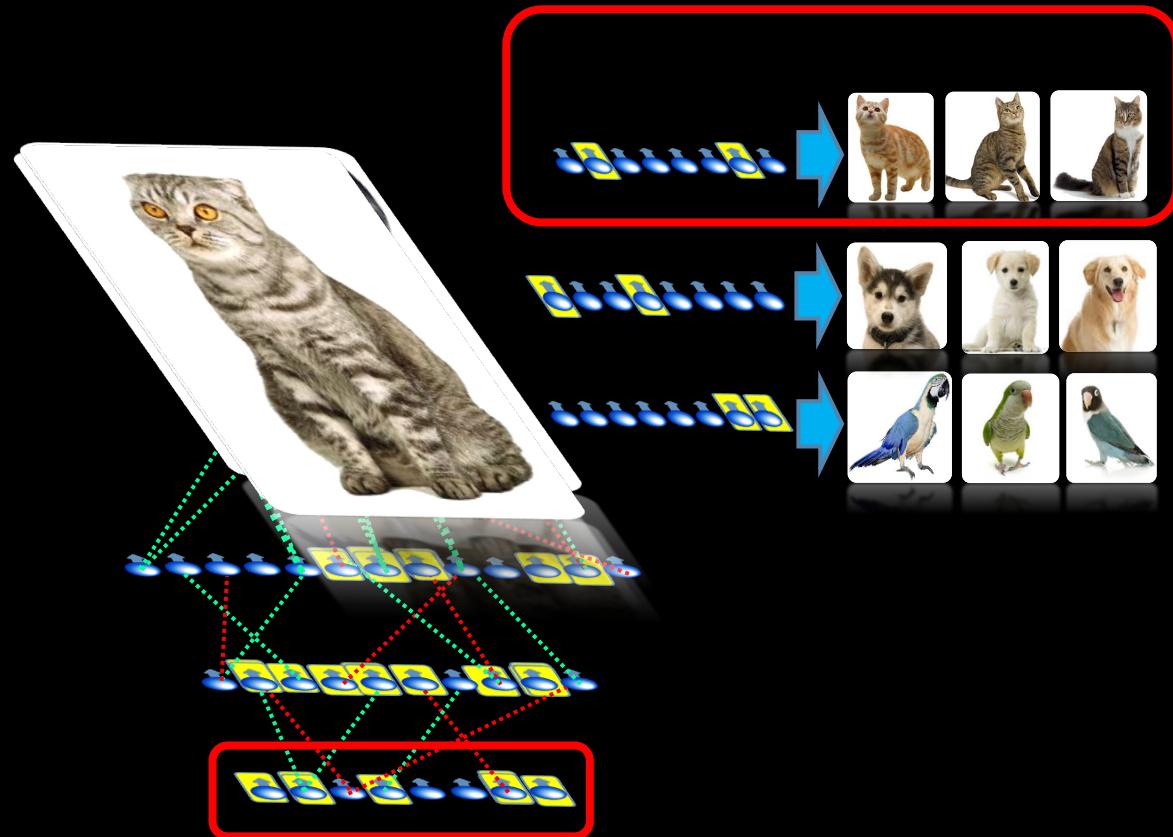
We can also use the network to search for similar items!

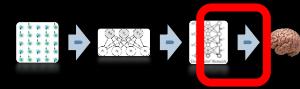
First we train our belief network on our millions of images to determine its weights for each...

And we store the association between each activation pattern and its image...



Later, when presented with a new image, we can then identify similar matches!



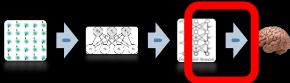


And Here's the Kicker

Notice that we literally “pumped in” raw input data into our Restricted Boltzmann Machines!

They figured out the rest!

Most traditional Machine Learning approaches require the engineer to do extensive manual feature extraction...
prior to applying Machine Learning on the input.



Deep Learning: Recent Wins



Baidu's Deep Speech project uses Recurrent Neural Networks and massive amounts of speech data (100x) to achieve a 29% improvement in speech recognition in 2 weeks!

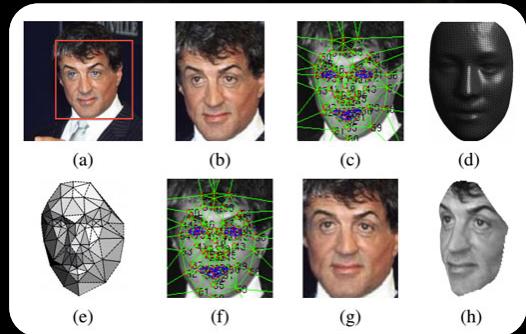
System	Clean (94)	Noisy (82)	Combined (176)
Apple Dictation	14.24	43.76	26.73
Bing Speech	11.73	36.12	22.05
Google API	6.64	30.47	16.72
wit.ai	7.94	35.06	19.41
Deep Speech	6.56	19.06	11.85

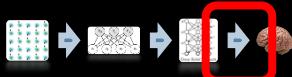


Microsoft Research recently published image recognition results of 4.94% (top-5 error rate) on the 2012 ImageNet corpus... besting the established human error rate of 5.1%!



Facebook's DeepFace is capable of facial recognition rates of 97.25% - roughly .28% less than human-level accuracy!





Deep Learning: The Catalysts

2006 – Geoffrey Hinton discovers
Contrastive Divergence and
greedy layer-wise pre-training.



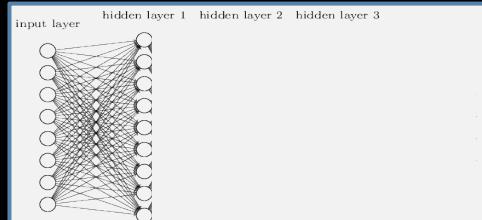
LETTER Communicated by Yann Le Cun

A Fast Learning Algorithm for Deep Belief Nets

Geoffrey E. Hinton
hinton@cs.toronto.edu
Simon Osindero
osindero@cs.toronto.edu
Department of Computer Science, University of Toronto, Toronto, Canada M5S 3G4

Yee-Whye Teh
teluya@comp.nus.edu.sg
Department of Computer Science, National University of Singapore,
Singapore 117543

We show how to use “complementary priors” to eliminate the explaining-away effects that make inference difficult in densely connected belief nets that have multiple layers of latent variables. In particular, we derive a fast greedy algorithm that can learn deep, directed belief nets layer by layer at a time, where the top two layers form an undirected associative memory. The fast, greedy algorithm is used to initialize a slower learning procedure that fine-tunes the weights using a contrastive version of wake-sleep. The final model is a hierarchical tree of three hidden layers of digit images and their labels. This general model distribution gives better digit classification than the best discriminative learning algorithms, and it is learned by long ravines in the free-energy landscape of top-level connections. By using the directed connections to display what the associative memory has in mind.

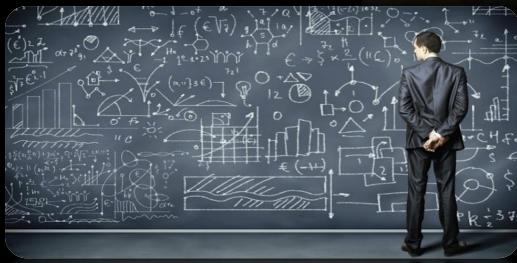


Huge increases in computing power from GPUs and cloud computing.

The Internet **yields** orders of magnitude more data **to train on**.



A series of advances in machine learning algorithms: weight initialization, regularization, new activation functions

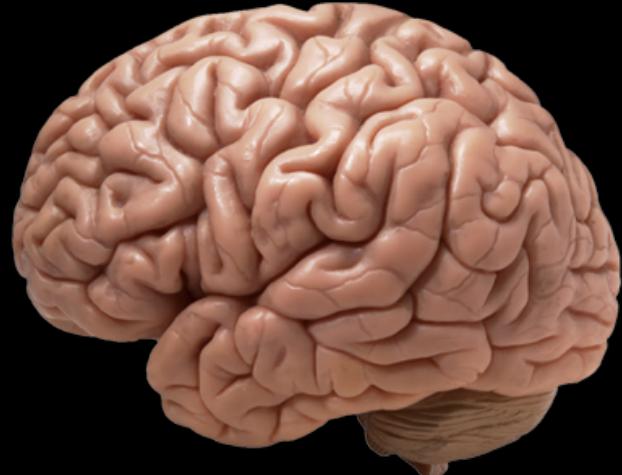


One Last Thought...

Have you ever thought about how amazing your visual memory is?

Do you think your brain literally stores away high-resolution snapshots of everything it sees?

Well maybe not!



This could explain why...

Given a memory cue you can suddenly remember a whole thought.

When viewing a partial image, your brain can visualize the rest...

When you see something unusual your brain knows it immediately...

After seeing hundreds of cats as a kid, but not knowing what they were called...

the moment your mom told you “that’s called a ‘cat’” you could instantly classify every cat you’d ever seen.

Conclusion

Deep Learning is a Machine Learning technique that works by learning multiple levels of abstraction about its inputs.

Once a Deep Learning system learns the inherent structure of its inputs,
you can use it to...

- classify items,
- reconstruct partial inputs,
- detect anomalies,
- search for related items,
- just daydream!

And finally, the engineer doesn't have to do lots of pre-processing to apply Deep Learning to many real-world problems!



Thank you !

Contact information:

Wu Xuening (邬学宁)

SAP硅谷创新中心 首席科学家

A:上海市浦东新区晨辉路1001号

E: x.wu@sap.com

T:021-61085287

As we know, there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns—the ones we don't know we don't know.

Rumsfeld