



Agenda

- History, continued
- Deep Learning
- TensorFlow

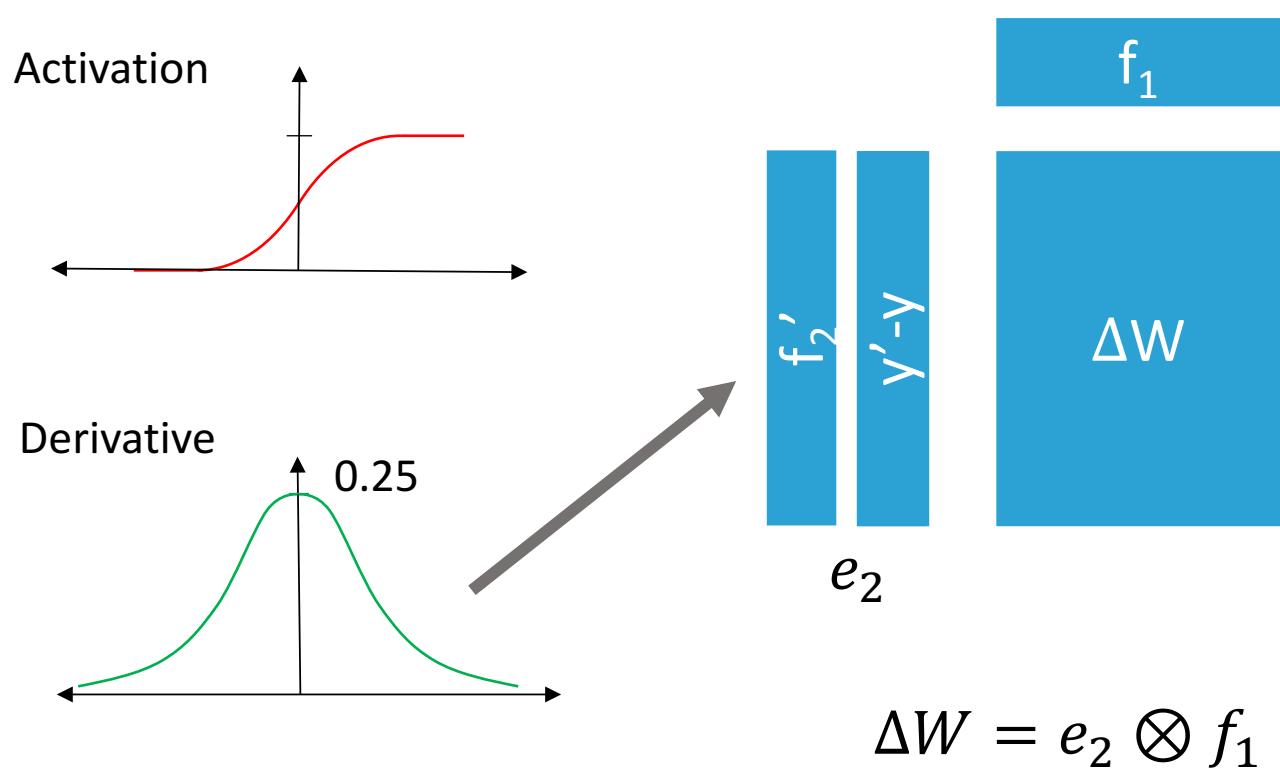
A woman with dark hair and a tattooed arm is rappelling down the side of a modern skyscraper. She is wearing a black tank top and grey pants. The background shows a dense city skyline with various buildings and flying vehicles. The overall aesthetic is futuristic and dynamic.

History of Neural Networks, Part 2

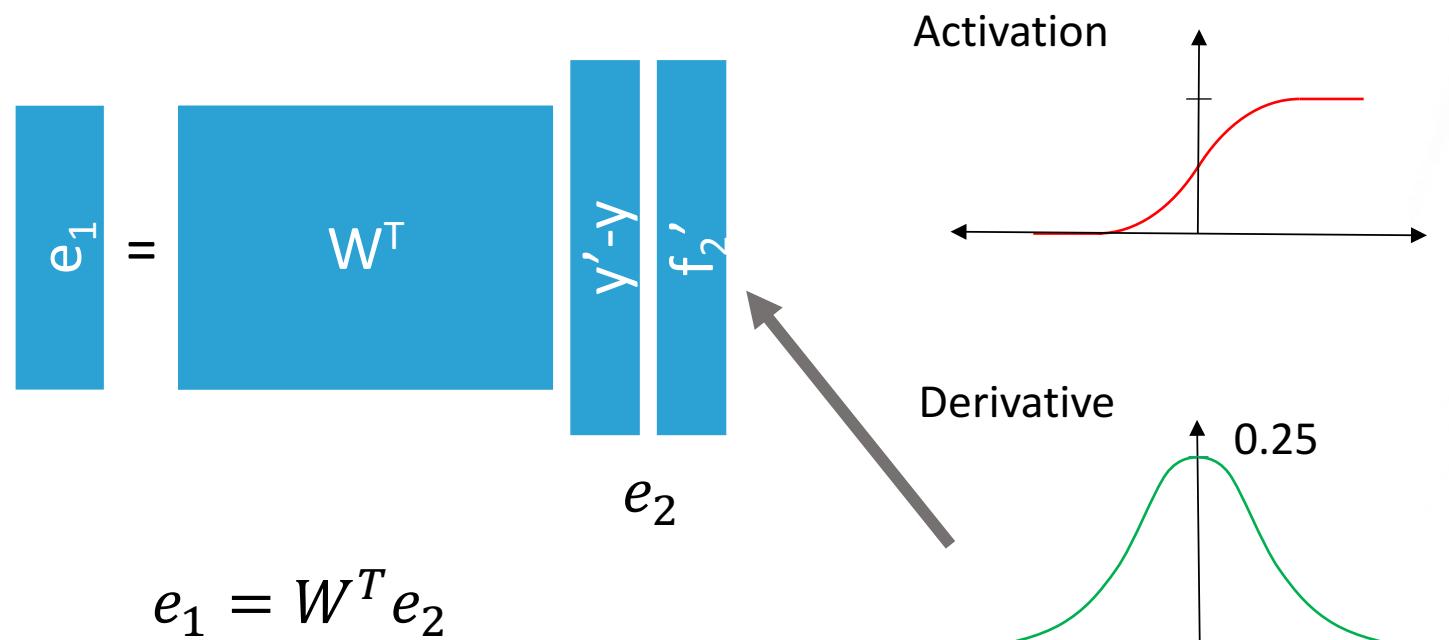
Three Problems

- Vanishing gradients in deep and recurrent networks
- High dimensional parameter spaces
- Internal representations

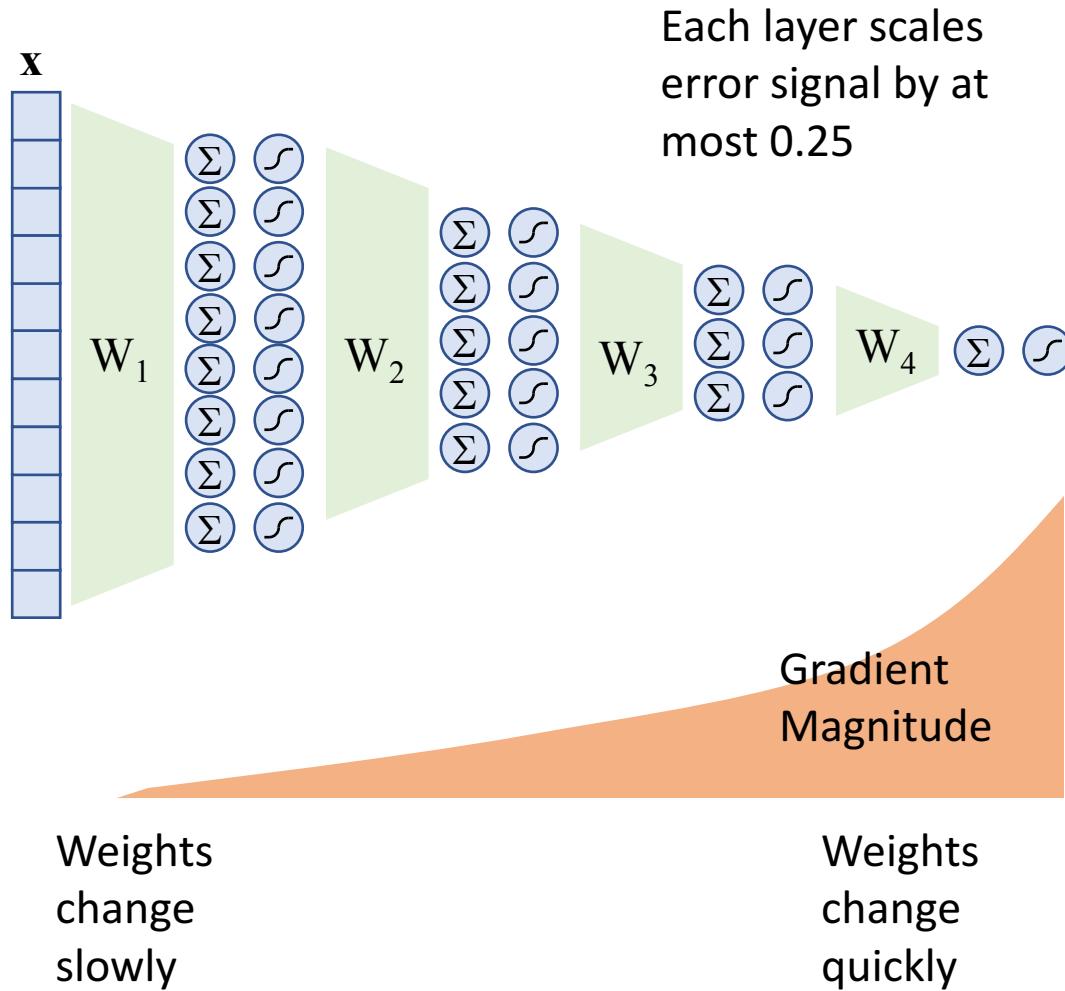
Vanishing Gradient



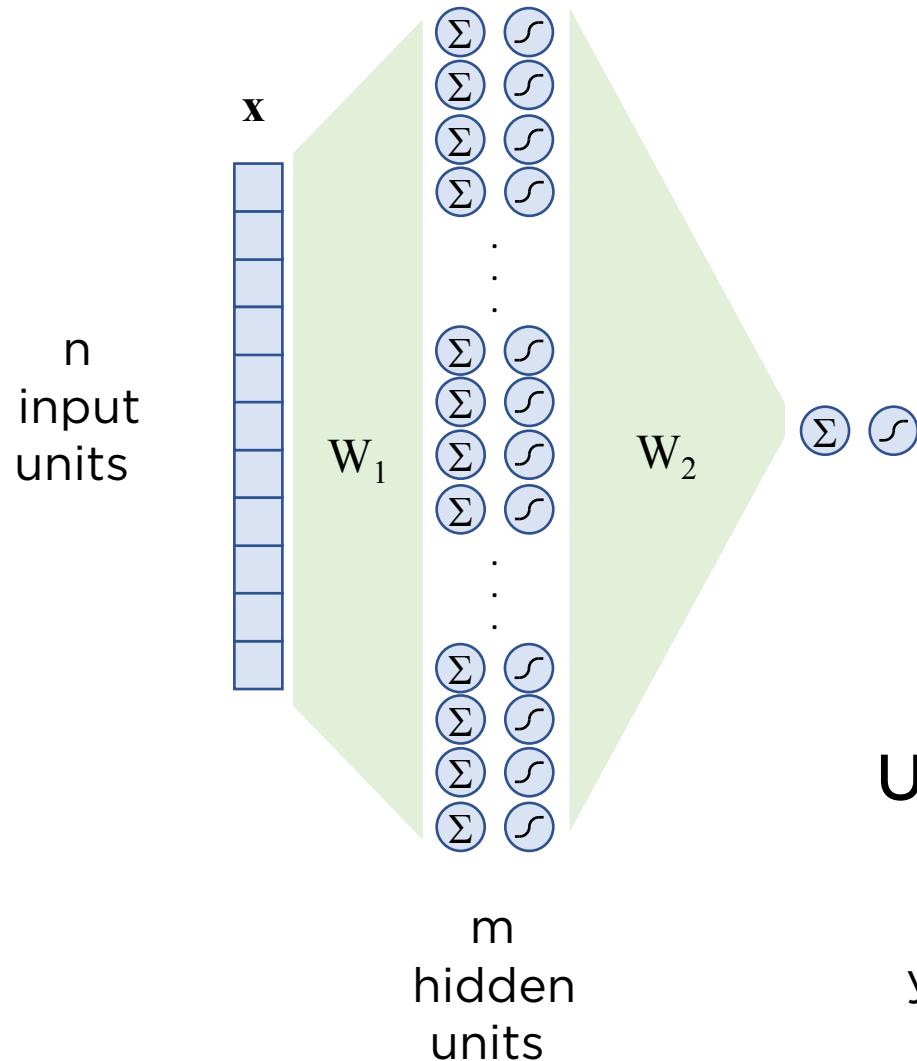
Vanishing Gradient



Vanishing Gradient



High Dimensional Search Spaces

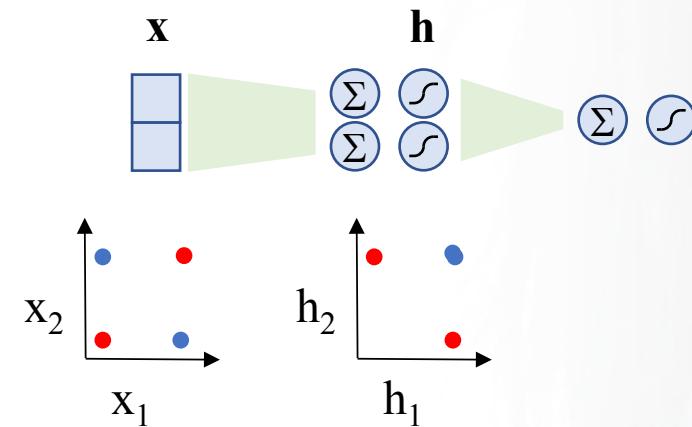


Universal Approximation Theorem

With enough hidden units
you can approximate anything

Internal Representations

- Activations of hidden units
- Totally dependent upon
- Training Set
- Architecture Selection
(units in hidden layer)
- Random Initial Weights
- Hard to change poor/unlucky initial representational commitments



Ten Year Intermission

- Neural networks are black boxes
 - No explanations
- Hyperparameter search is hard
 - Complex interactions between learning rates, network sizes, randomization parameters
- Support Vector Machines
 - Return to linear transformations
 - Kernel Trick

A close-up, low-angle shot of a character from the video game Cyberpunk 2077. The character has short, light-colored hair and is wearing a detailed, purple-toned leather jacket over a dark top. They are looking upwards and to the right. The background is dark and atmospheric, featuring glowing blue spheres, a large satellite-like structure, and neon lights reflecting off surfaces.

Deep Learning

Question

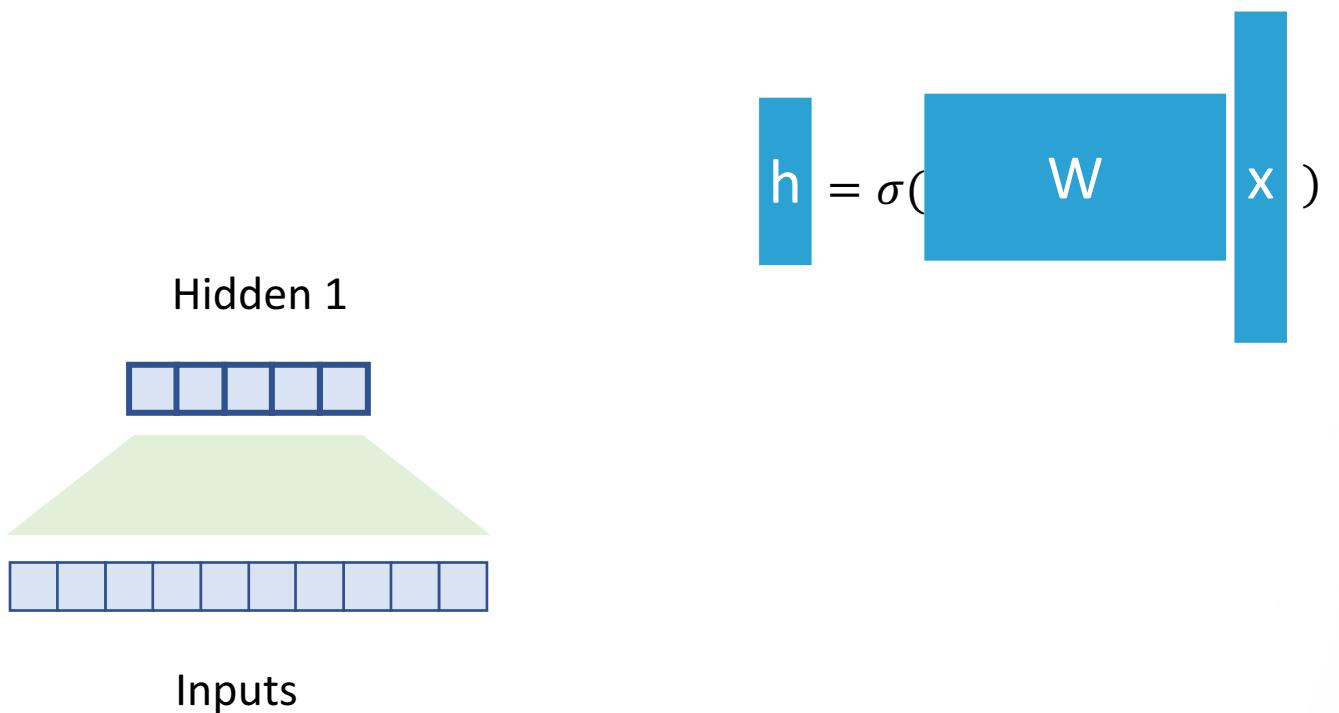
How do you train deep networks
with back propagation?

You don't.



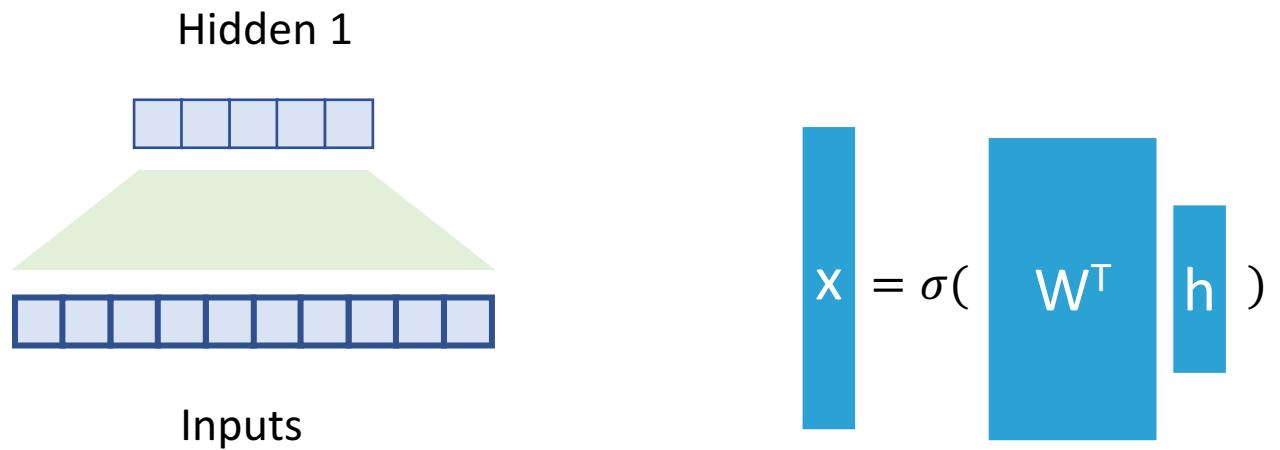
Deep Learning (2006)

Incrementally stack unsupervised networks
Restricted Boltzman Machines



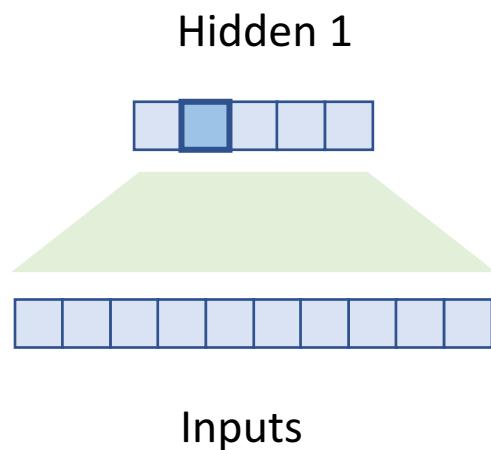
Deep Learning (2006)

Incrementally stack unsupervised networks
Restricted Boltzman Machines



Deep Learning (2006)

Incrementally stack unsupervised networks
Restricted Boltzman Machines

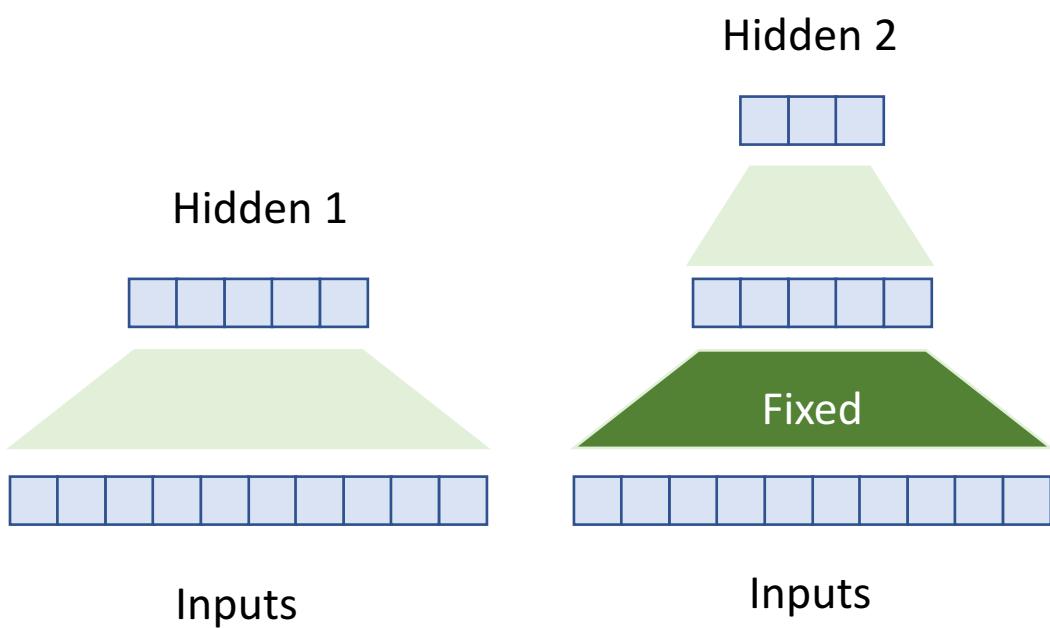


$$h = \sigma(Wx)$$
$$x = \sigma(W^T h)$$

Update units randomly to avoid oscillations

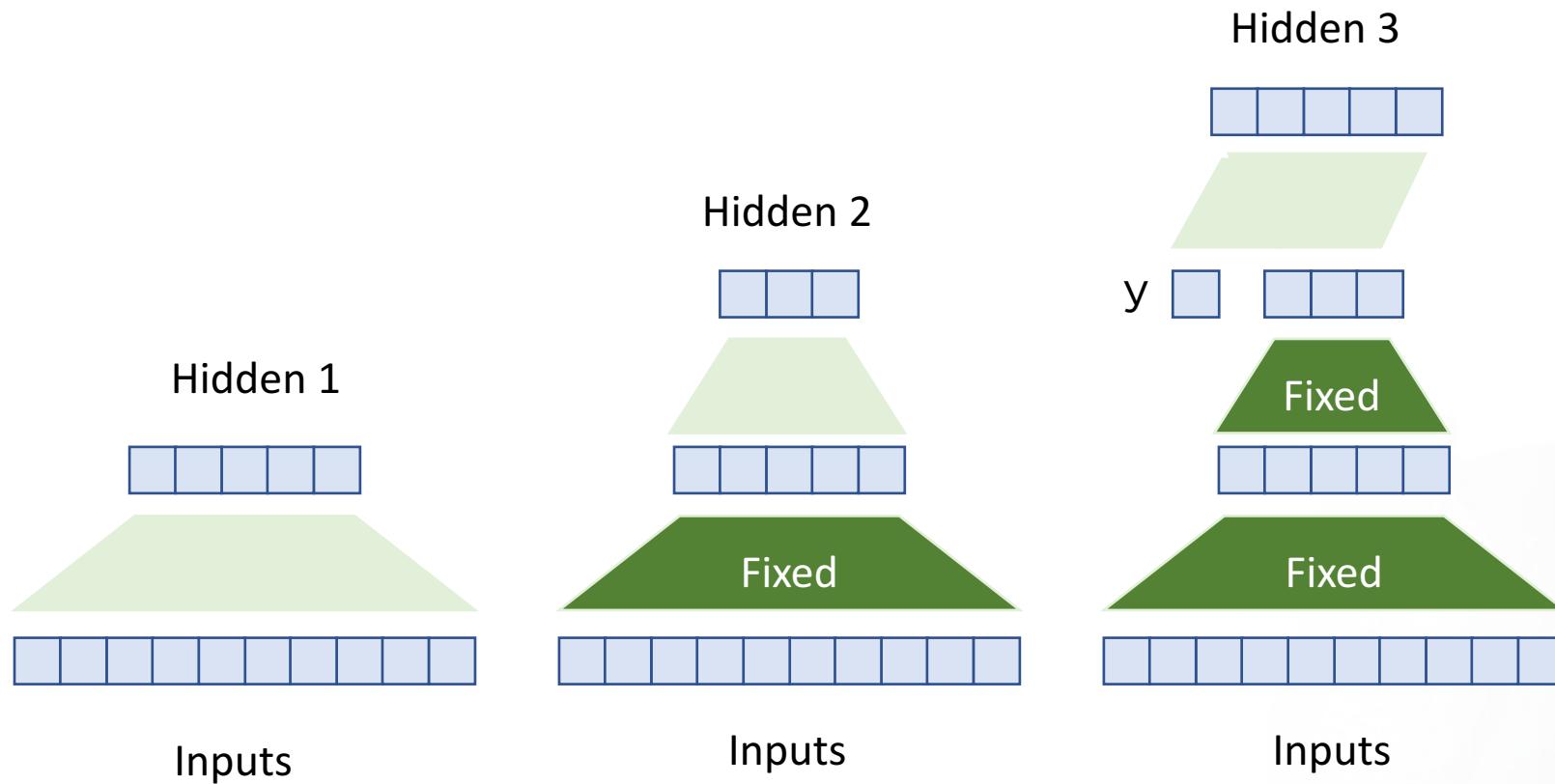
Deep Learning (2006)

Incrementally stack unsupervised networks
Restricted Boltzman Machines



Deep Learning (2006)

Incrementally stack unsupervised networks
Restricted Boltzman Machines

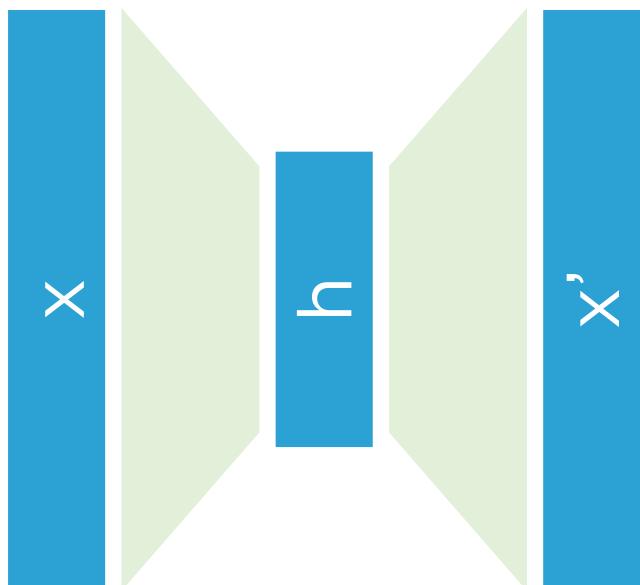


Autoencoders

Input = Output

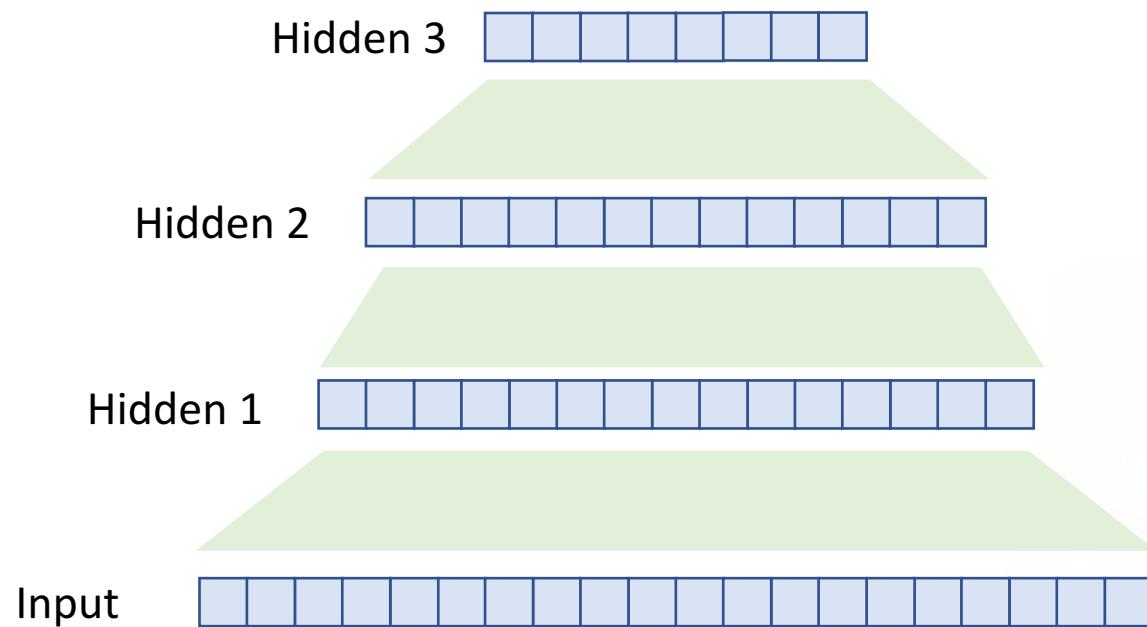
Capture distribution of data set

Bottleneck forces dimensionality reduction



Deep Belief Networks

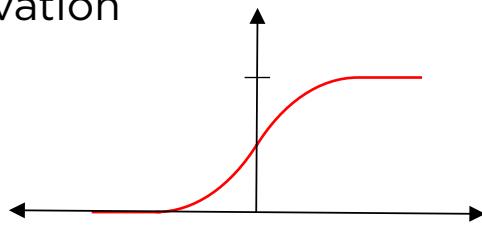
Stacking autoencoders creates layers of abstractions of *input space*



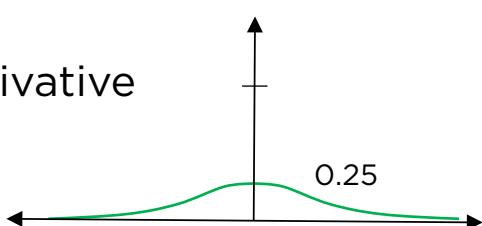
Deep Learning (2010)

Why deep networks can't learn

Activation



Derivative

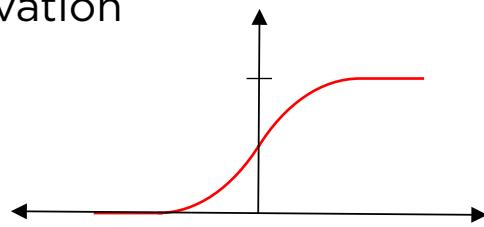


Sigmoid
Vanishing Gradient

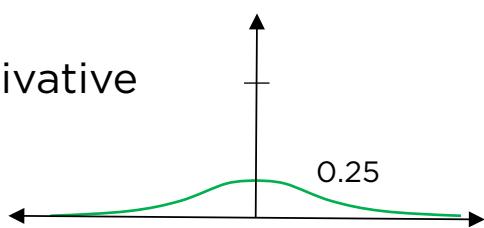
Deep Learning (2010)

Now deep networks can learn

Activation

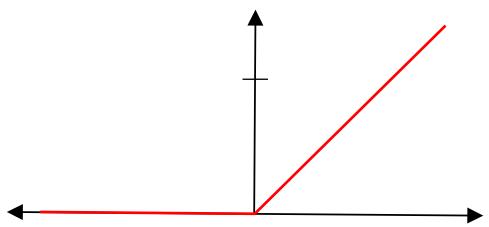


Derivative

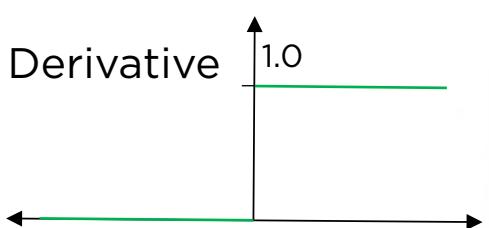


Sigmoid
Vanishing Gradient

Activation



Derivative



Rectified Linear Unit
Gated Gradient

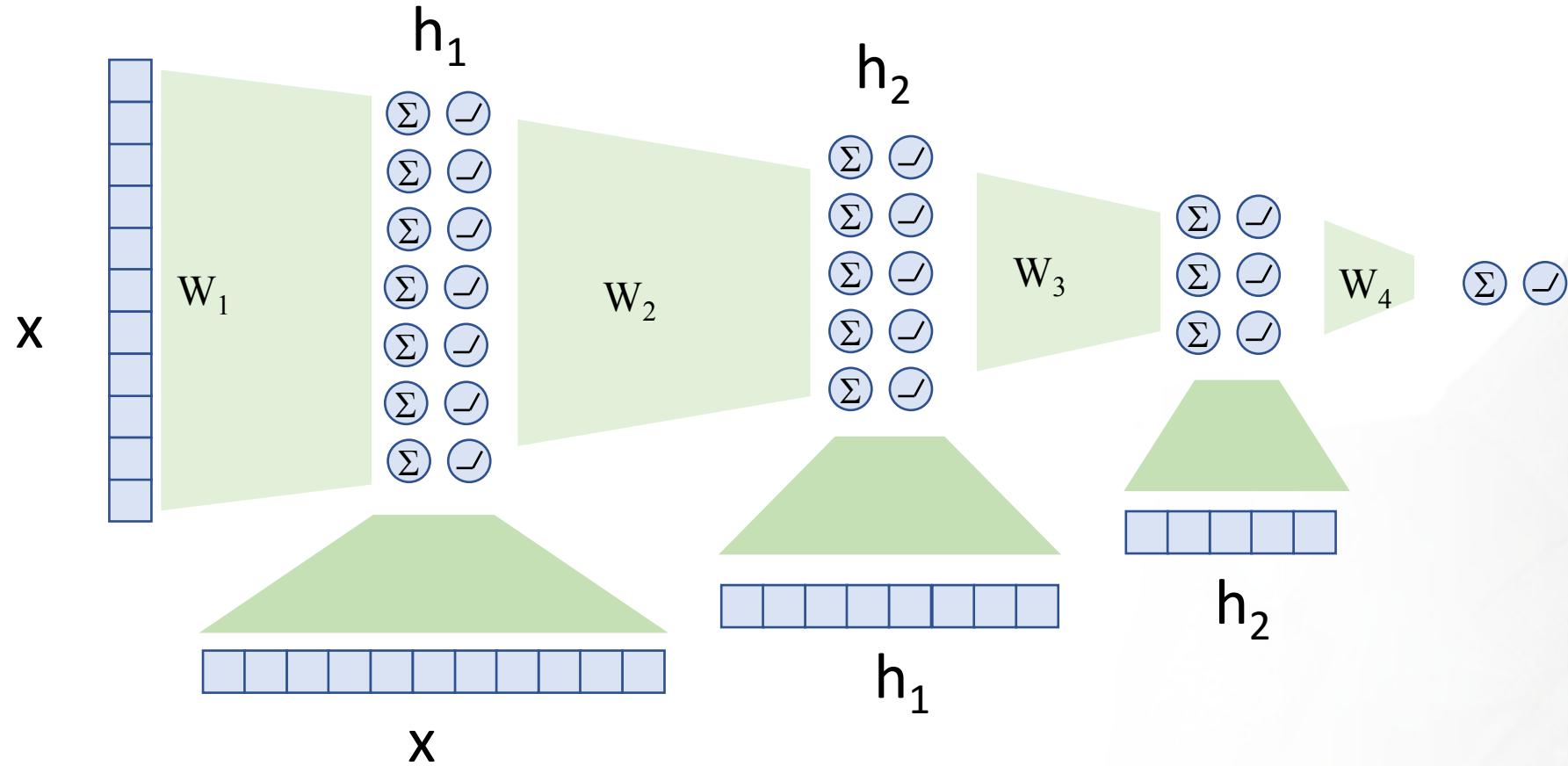
Solved Problems

- Internal Representations
 - Autoencoder constraints
 - Input space knowledge
- Vanishing Gradient
 - Rectified Linear Unit (ReLU)
 - Functional knowledge

Now we can build deep functional networks

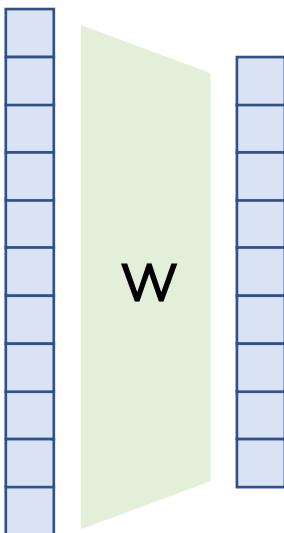
Deep Functional Networks

Can be trained with back propagation and autoencoder constraints

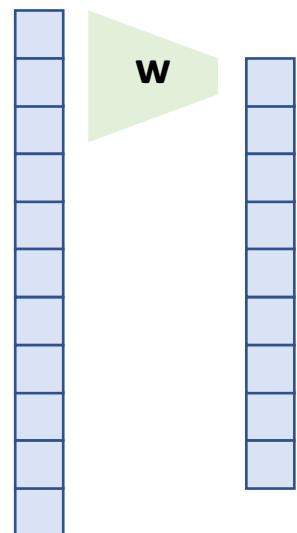


Convolutions

Reducing the number of weights
Spatial distribution of a simple function



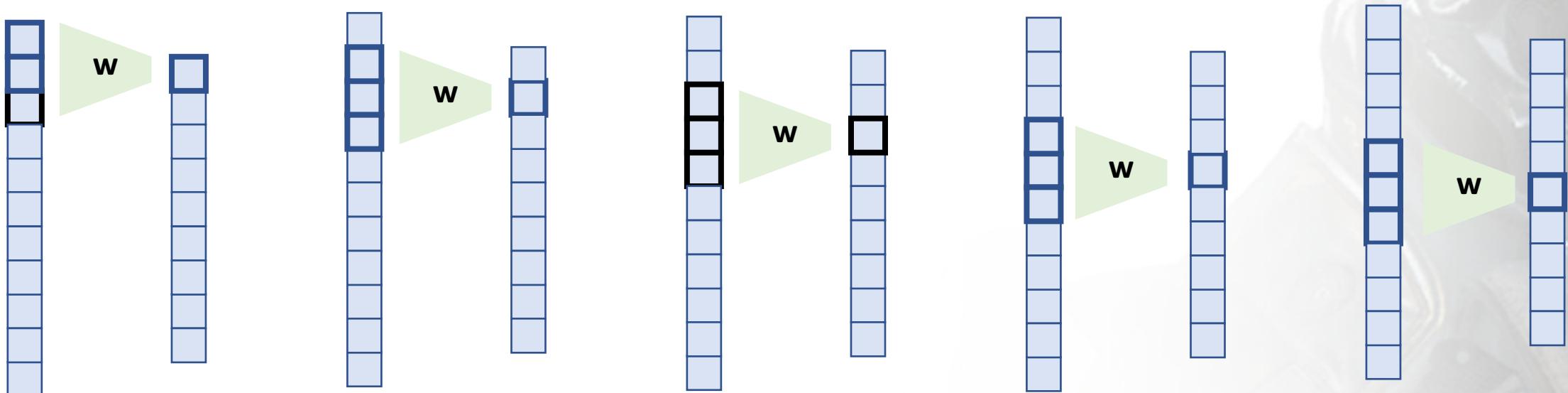
$O(nm)$ weights



$O(k)$ weights

Convolutions

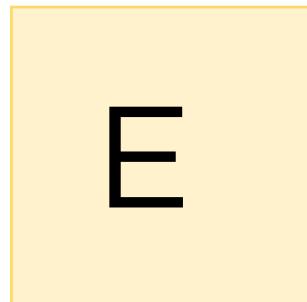
Reducing the number of weights
Spatial distribution of a simple function



O(k) weights
Applied m times

Convolutions

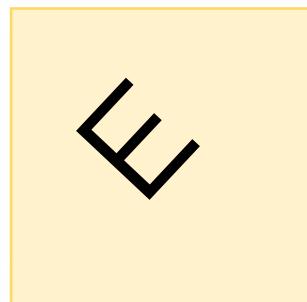
Capture invariances



Translation



Scale



Rotation

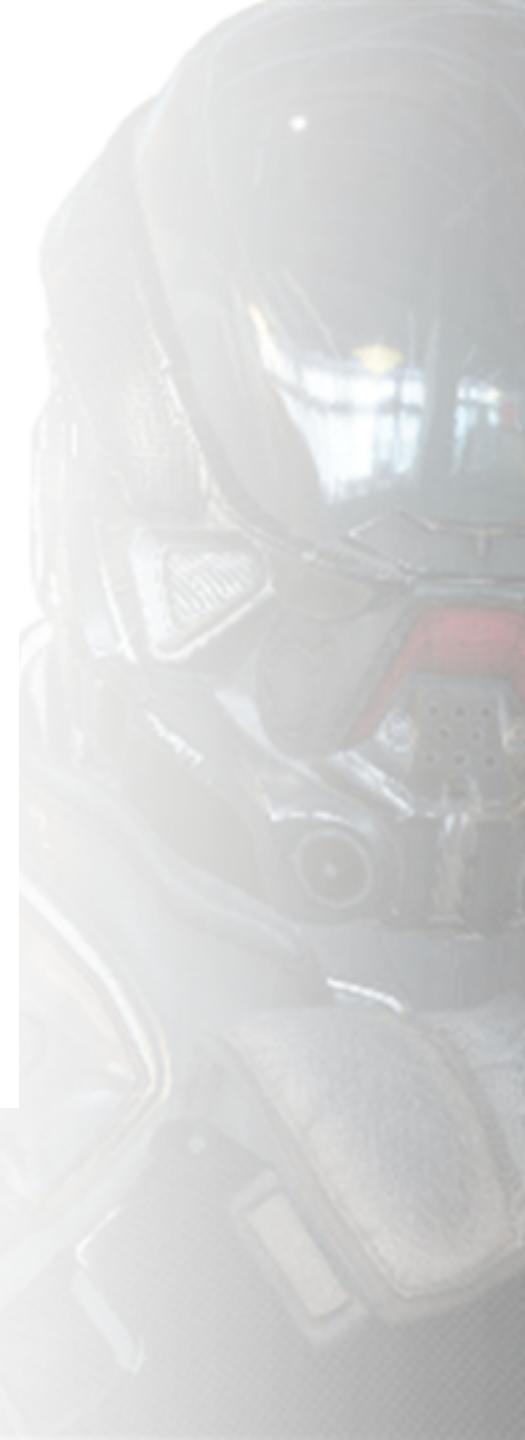
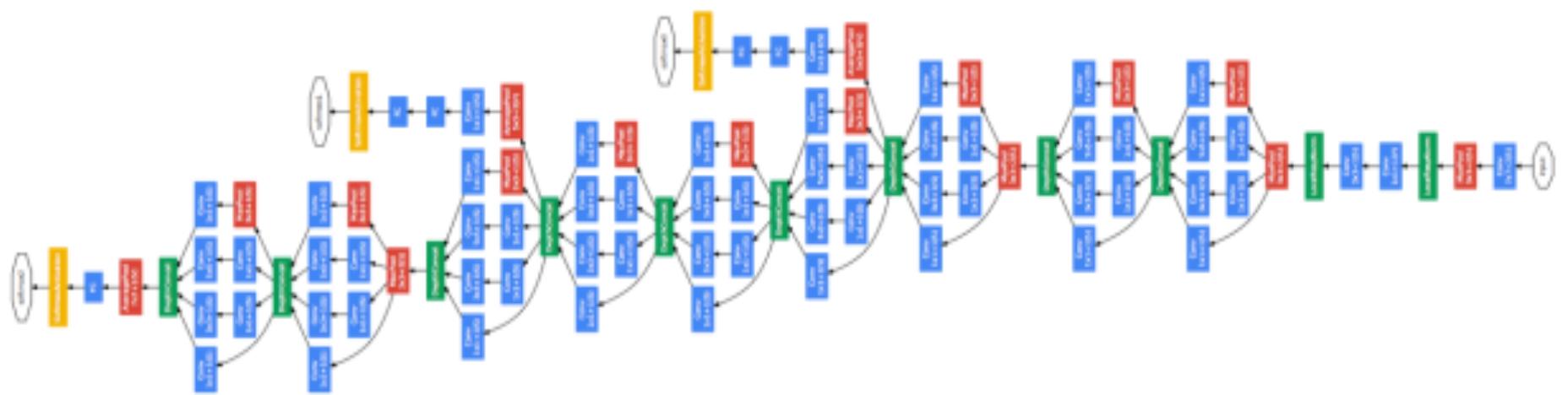
Solved Problems

- Internal Representations
 - Autoencoder constraints
 - Input space knowledge
- Vanishing Gradient
 - Rectified Linear Unit (ReLU)
 - Functional knowledge
- High dimensional parameter spaces
 - Convolutions
 - Invariances

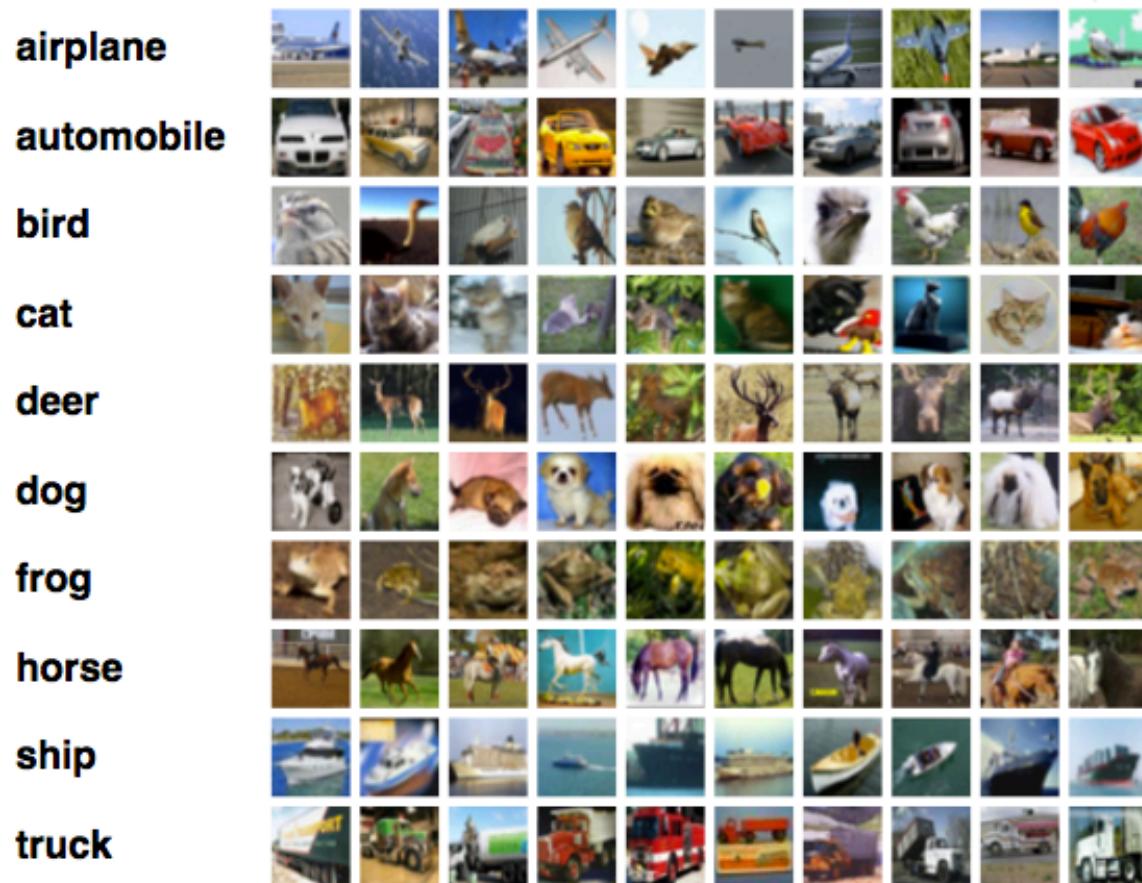
Now we can build really deep functional networks



Deep Learning



Neural Networks Are Really Great!



Deep Learning on GPUs

Why GPUs?

Very fast matrix multiplies

- Hundreds of simple cores
- Thousands of hardware threads
- Maximized floating point throughput
- Increased memory bandwidth

NVIDIA's Kepler and Maxwell microarchitecture based GPUs

Other Hardware Solutions

- Google – Tensor Processing Unit (TPU)
 - Reduced precision (8 bit arithmetic)
 - GPU architecture, with graphics specializations removed
- Cerebras – mesh architecture
- Deep Vision – convolution multiprocessor
- Graphcore – intelligent processing unit (IPU)
- Wave Computing – data flow architecture (no CPU)

Don't Forget Big Data

- Lots of data to work with
- Millions of images and documents
- Capturing distributions easier





Are We There Yet?

Deep Learning Applications

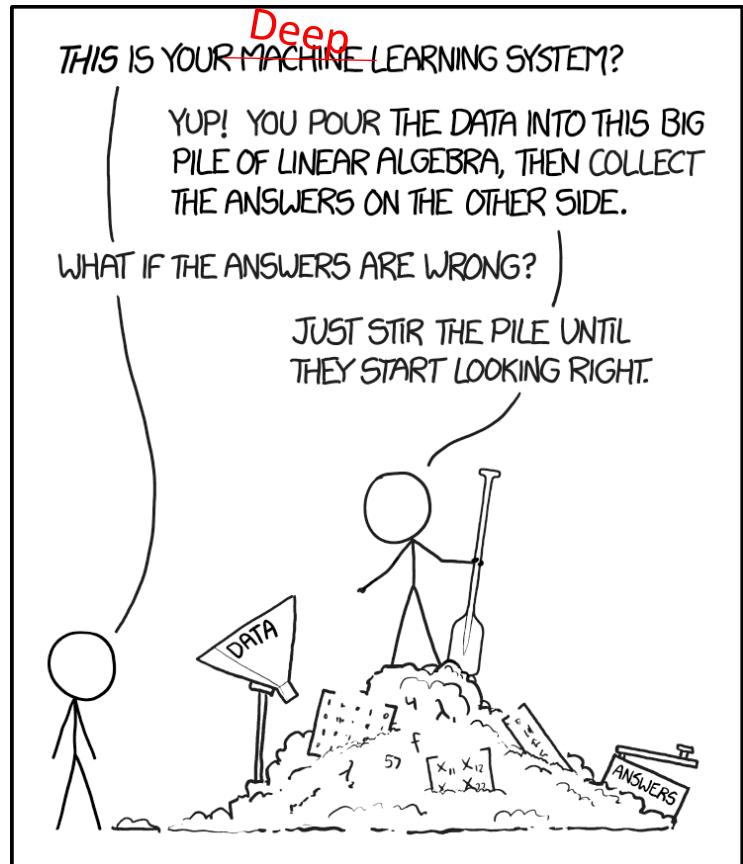
- Speech Recognition
- Alpha Go
- Self Driving Vehicles
- Text Analysis
- Image Analysis

Large input and/or output representations (e.g. images)

*Transformation invariances on inputs (e.g.
translation/scale/rotation)*



Nothing But Matrix Multiplication



Easy

Hard

Stirring The Pile

The deep learning algorithms find a set of weights for the network, but

The model selection problem has not gone away

Developers still have to

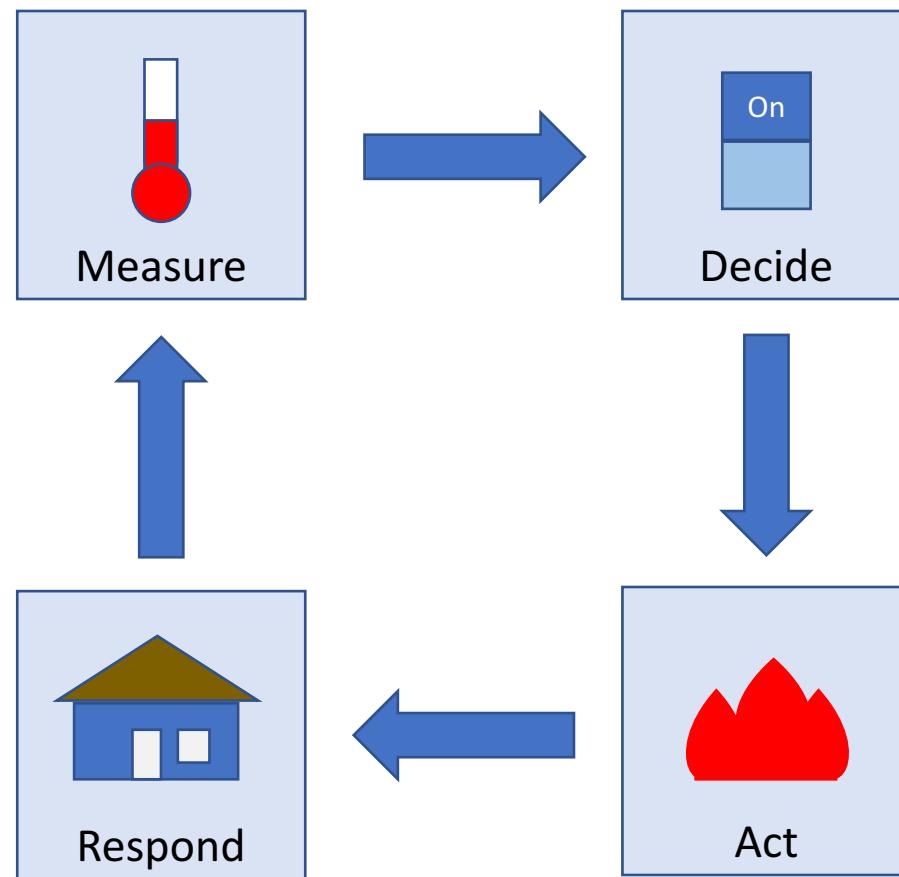
- Specify the data flow graph
- Identify type of processing occurring at each node and layer
- Supply more hyper-parameters
 - convolution size, hidden units, activation functions, etc.

Back To The Cursed Loop

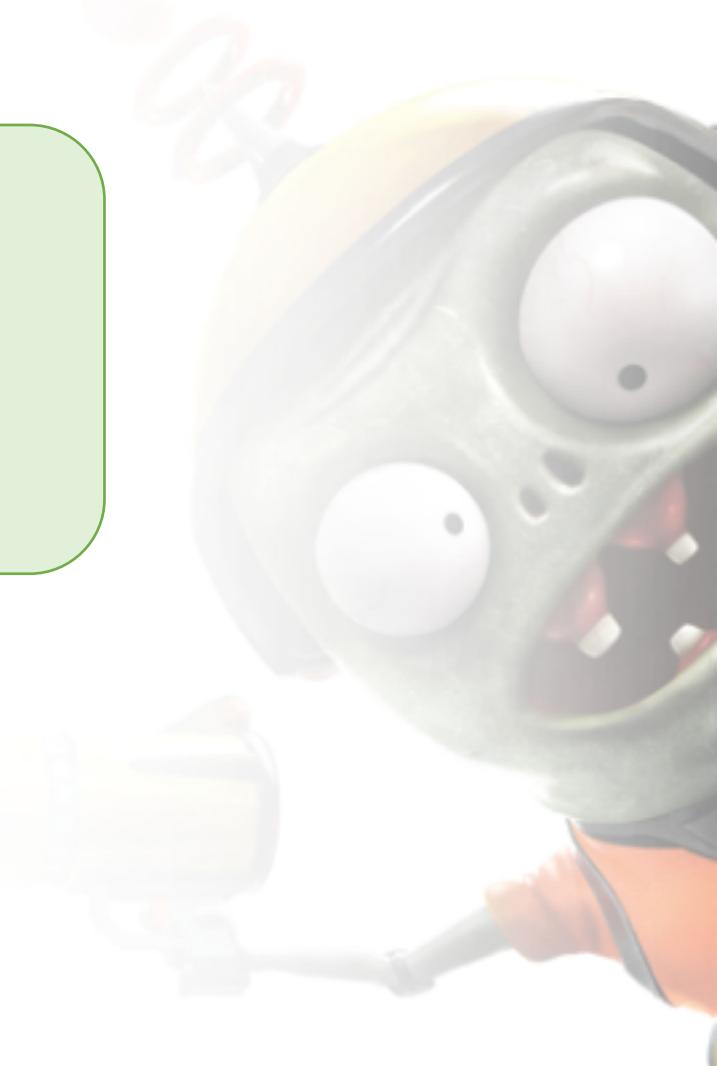
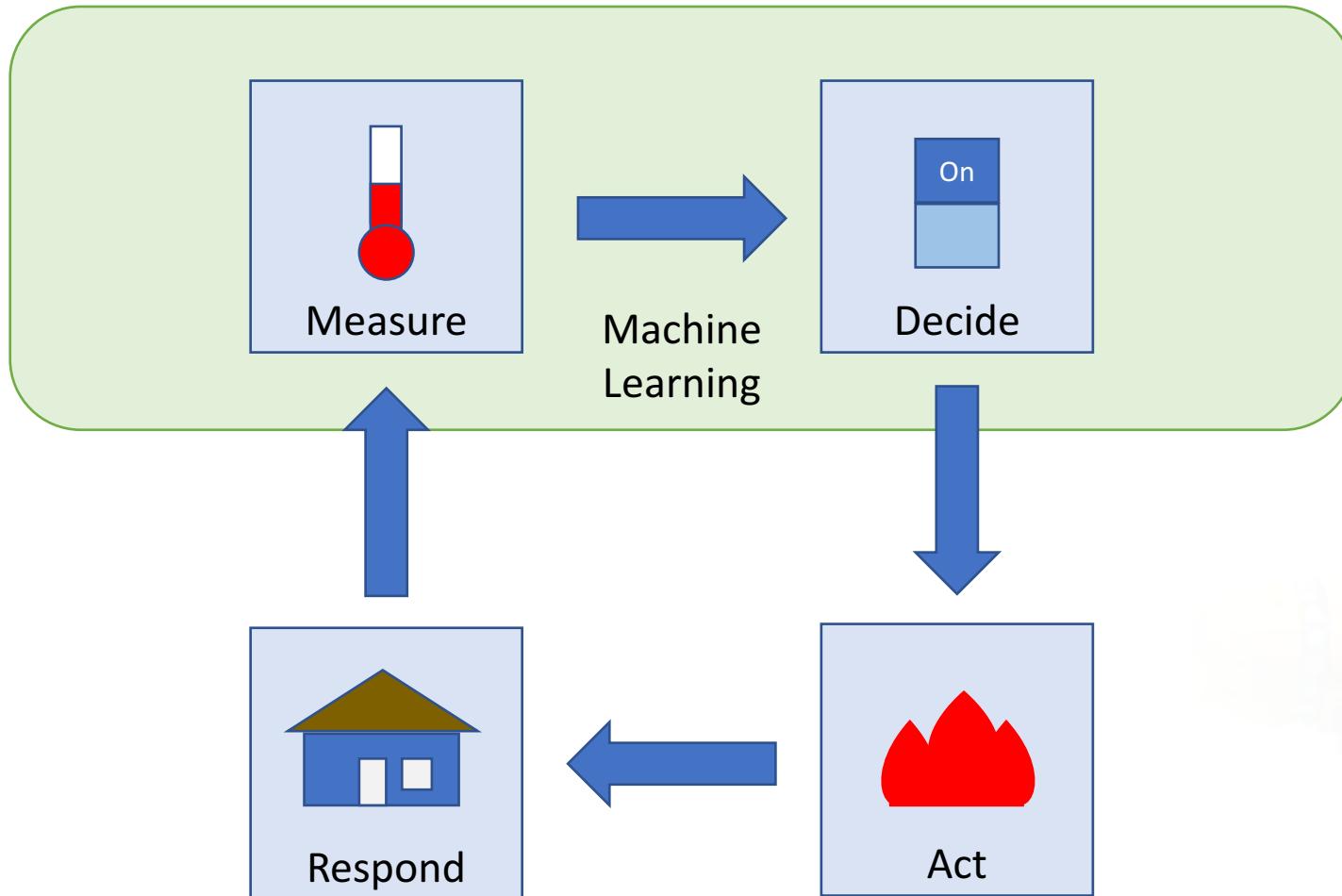
```
do  
    select p from possible hyperparameters  
    model = train(data, p)  
until model is acceptable
```



Systems and Models



Systems and Models



TensorFlow Demo

MNIST Data



00000000000000000000
/11111111111111111111
22222222222222222222
33333333333333333333
44444444444444444444
55555555555555555555
66666666666666666666
77777777777777777777
88888888888888888888
99999999999999999999



Next Steps