

**CMP717**

**Image Processing**

# **Introduction to Deep Learning**

Erkut Erdem  
Hacettepe University  
Computer Vision Lab (HUCVL)



# What is deep learning?

Y. LeCun, Y. Bengio, G. Hinton, "Deep Learning", Nature, Vol. 521, 28 May 2015

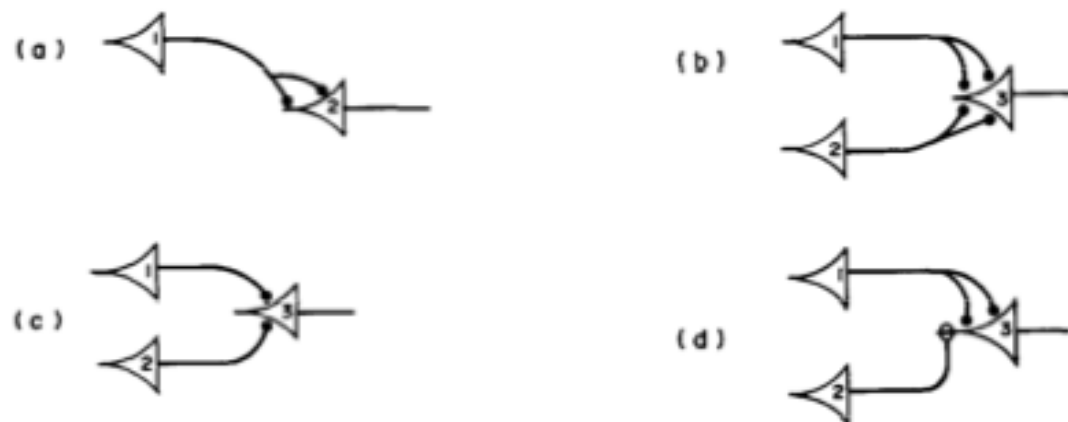
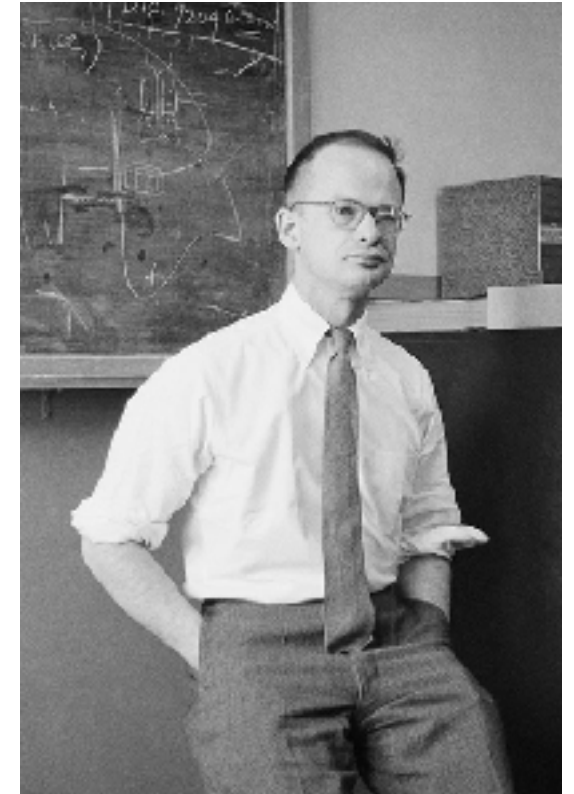
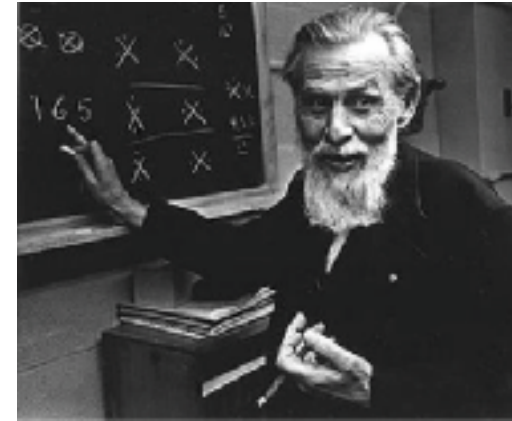


“Deep learning allows computational models that are composed of **multiple processing layers to learn representations of data with multiple levels of abstraction.**”  
– Yann LeCun, Yoshua Bengio and Geoff Hinton 2

1943 – 2006:  
A Prehistory of Deep Learning

# 1943: Warren McCulloch and Walter Pitts

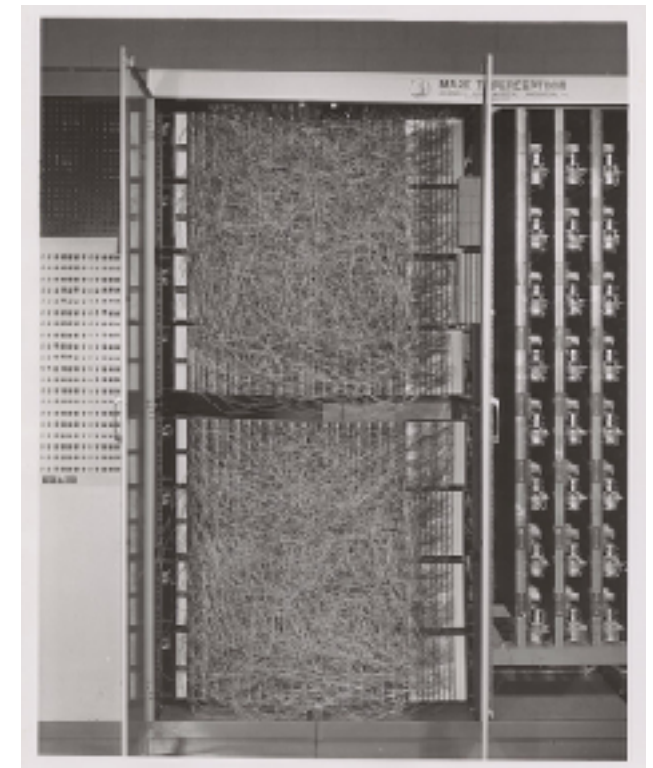
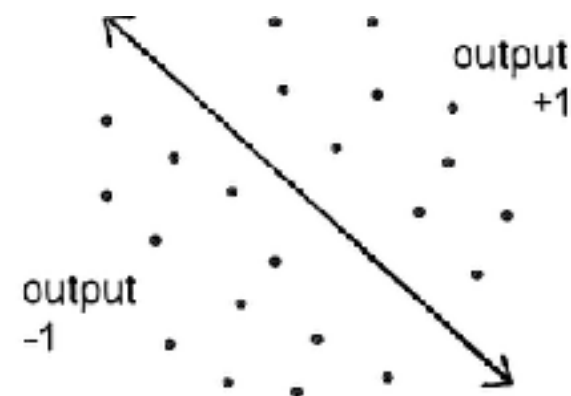
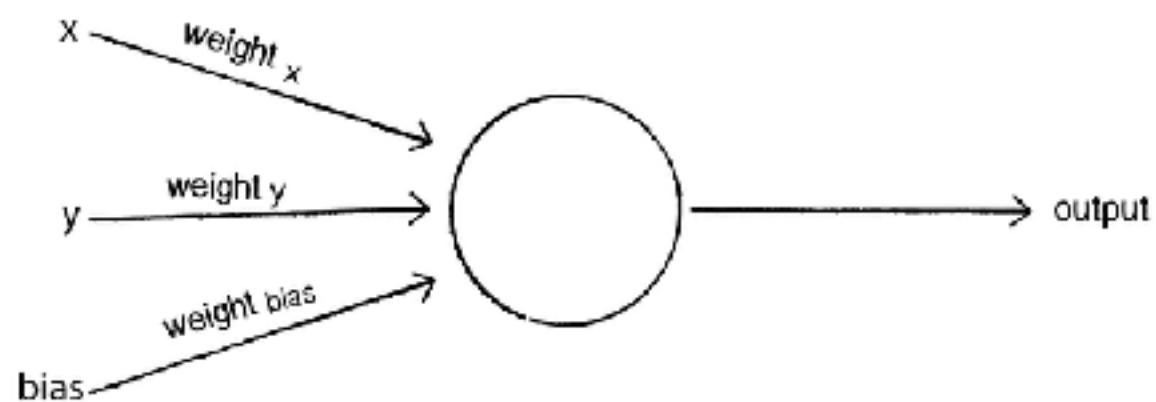
- First computational model
- Neurons as logic gates (AND, OR, NOT)
- A neuron model that sums binary inputs and outputs 1 if the sum exceeds a certain threshold value, and otherwise outputs 0





# 1958: Frank Rosenblatt's Perceptron

- A computational model of a **single neuron**
- Solves a **binary classification problem**
- Simple training algorithm
- Built using specialized hardware

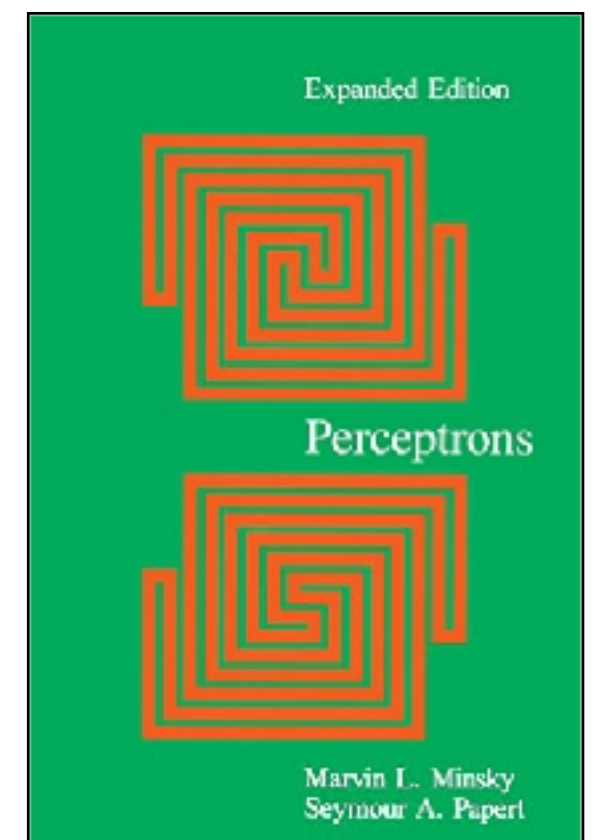
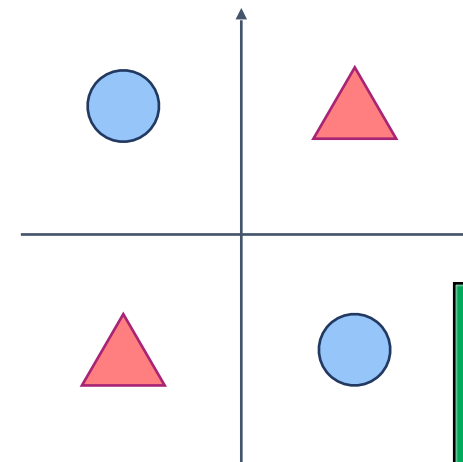


# 1969: Marvin Minsky and Seymour Papert

*“No machine can learn to recognize  $X$  unless it possesses, at least potentially, some scheme for representing  $X$ .” (p. xiii)*

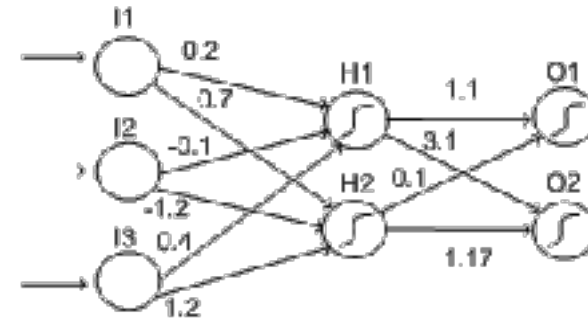


- Perceptrons can only represent linearly separable functions.
  - such as **XOR** Problem
- Wrongly attributed as the reason behind the **AI winter**, a period of reduced funding and interest in AI research



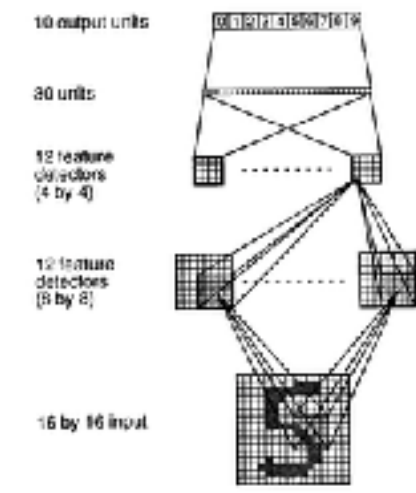
# 1990s

- **Multi-layer perceptrons** can theoretically learn any function (Cybenko, 1989; Hornik, 1991)



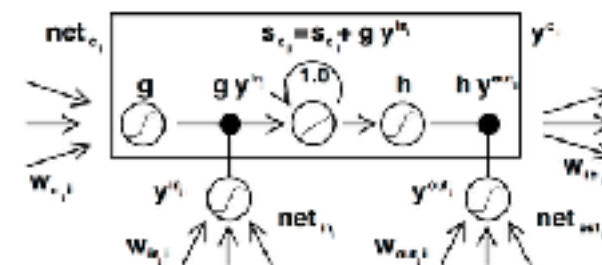
- Training multi-layer perceptrons

- **Back-propagation** (Rumelhart, Hinton, Williams, 1986)
- **Back-propagation through time (BPTT)** (Werbos, 1988)



- New neural architectures

- **Convolutional neural nets** (LeCun et al., 1989)
- **Long-short term memory networks (LSTM)** (Schmidhuber, 1997)



# Why it failed then

- Too many parameters to learn from few labeled examples.
- “I know my features are better for this task”.
- Non-convex optimization? No, thanks.
- Black-box model, no interpretability.
- Very slow and inefficient
- Overshadowed by the success of SVMs (Cortes and Vapnik, 1995)



A major breakthrough in 2006

# 2006 Breakthrough: Hinton and Salakhutdinov

## Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton\* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. **Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights** that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

- The first solution to the **vanishing gradient problem**.
- Build the model in a layer-by-layer fashion using unsupervised learning
  - The features in early layers are already initialized or “pretrained” with some suitable features (weights).
  - Pretrained features in early layers only need to be adjusted slightly during supervised learning to achieve good results.

G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks”, Science, Vol. 313, 28 July 2006.



# The 2012 revolution

# ImageNet Challenge

- **IMAGENET** Large Scale Visual Recognition Challenge (ILSVRC)

- **1.2M** training images with **1K** categories
- Measure top-5 classification error



Output  
Scale  
T-shirt  
**Steel drum**  
Drumstick  
Mud turtle



Output  
Scale  
T-shirt  
Giant panda  
Drumstick  
Mud turtle



Easiest classes



Hardest classes



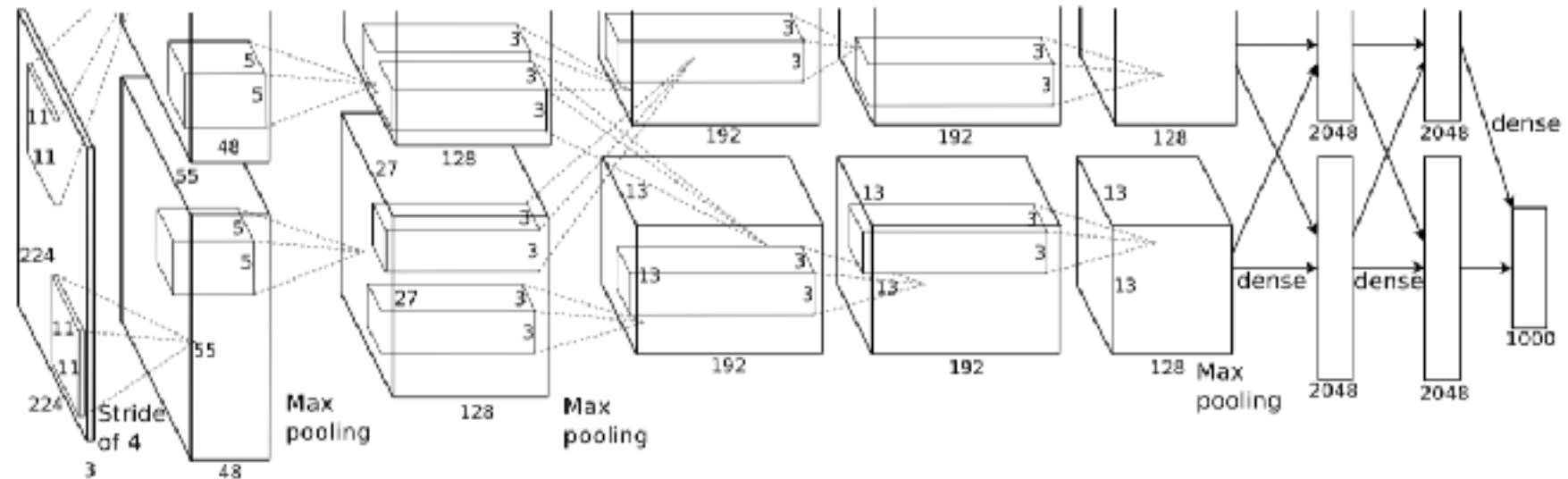
Image classification



# ILSVRC 2012 Competition

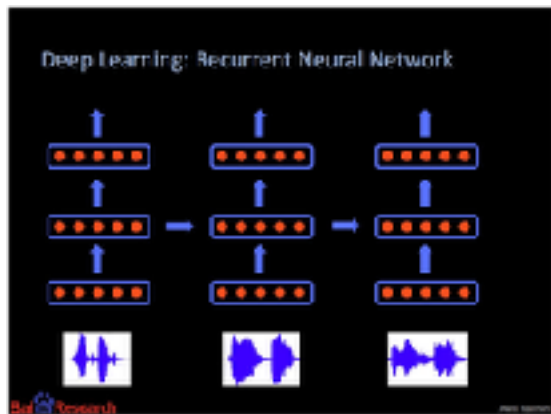
2012 Teams	%Error
Supervision (Toronto)	15.3
ISI (Tokyo)	26.1
VGG (Oxford)	26.9
XRCE/INRIA	27.0
UvA (Amsterdam)	29.6
INRIA/LEAR	33.4

CNN based,  
non-CNN based

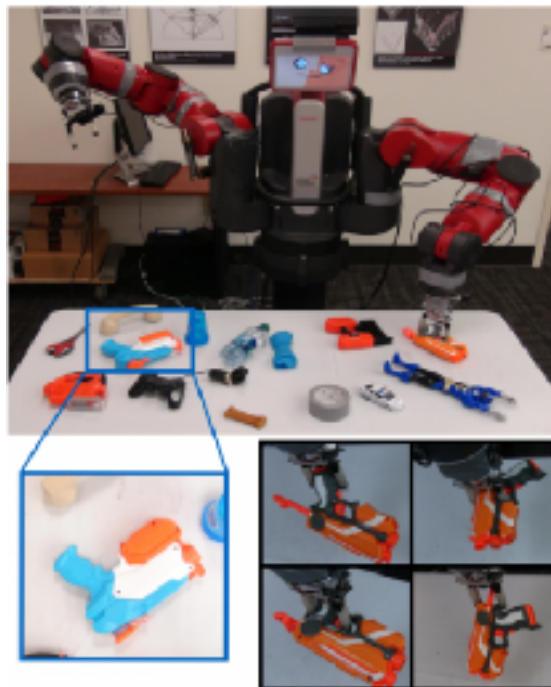


- The success of AlexNet, a deep convolutional network
  - 7 hidden layers (not counting some max pooling layers)
  - 60M parameters
- Combined several tricks
  - ReLU activation function, data augmentation, dropout

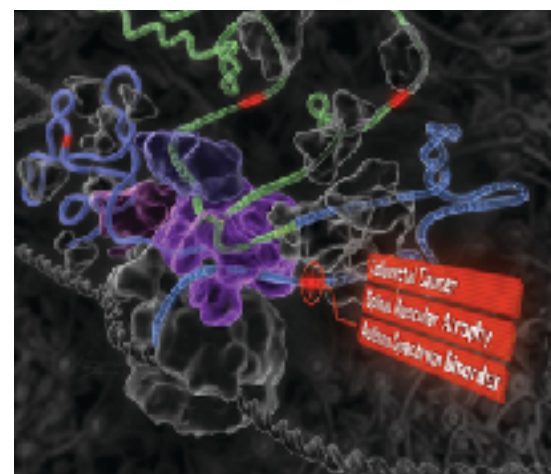
2012 – now  
A Cambrian explosion in  
deep learning



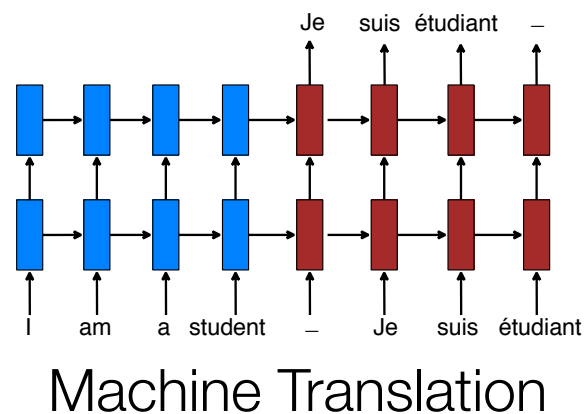
Speech recognition



Robotics



Genomics



Machine Translation



Game Playing



Audio Generation



Self-Driving Cars

Amodei et al., "Deep Speech 2: End-to-End Speech Recognition in English and Mandarin", In CoRR 2015

M.-T. Luong et al., "Effective Approaches to Attention-based Neural Machine Translation", EMNLP 2015

M. Bojarski et al., "End to End Learning for Self-Driving Cars", In CoRR 2016

D. Silver et al., "Mastering the game of Go with deep neural networks and tree search", Nature 529, 2016

L. Pinto and A. Gupta, "Supersizing Self-supervision: Learning to Grasp from 50K Tries and 700 Robot Hours" ICRA 2015

H. Y. Xiong et al., "The human splicing code reveals new insights into the genetic determinants of disease", Science 347, 2015

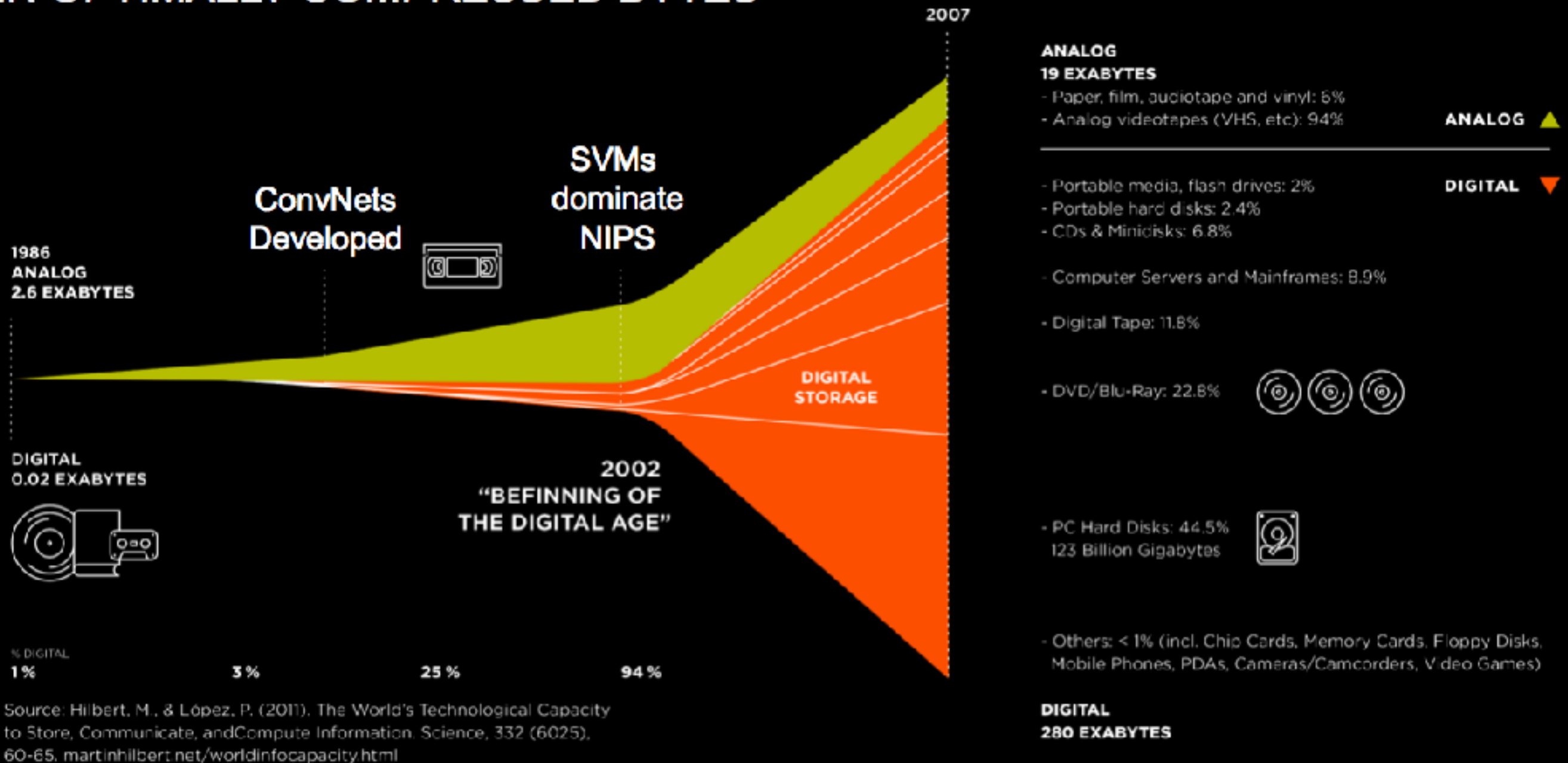
M. Ramona et al., "Capturing a Musician's Groove: Generation of Realistic Accompaniments from Single Song Recordings", In IJCAI 2015

And many more... 15

Why now?



# GLOBAL INFORMATION STORAGE CAPACITY IN OPTIMALLY COMPRESSED BYTES



# Datasets vs. Algorithms

Year	Breakthroughs in AI	Datasets (First Available)	Algorithms (First Proposed)
1994	Human-level spontaneous speech recognition	Spoken Wall Street Journal articles and other texts (1991)	Hidden Markov Model (1984)
1997	IBM Deep Blue defeated Garry Kasparov	700,000 Grandmaster chess games, aka “The Extended Book” (1991)	Negascout planning algorithm (1983)
2005	Google’s Arabic-and Chinese-to-English translation	1.8 trillion tokens from Google Web and News pages (collected in 2005)	Statistical machine translation algorithm (1988)
2011	IBM Watson became the world Jeopardy! champion	8.6 million documents from Wikipedia, Wiktionary, and Project Gutenberg (updated in 2010)	Mixture-of-Experts (1991)
2014	Google’s GoogLeNet object classification at near-human performance	ImageNet corpus of 1.5 million labeled images and 1,000 object categories (2010)	Convolutional Neural Networks (1989)
2015	Google’s DeepMind achieved human parity in playing 29 Atari games by learning general control from video	Arcade Learning Environment dataset of over 50 Atari games (2013)	Q-learning (1992)
<b>Average No. of Years to Breakthrough:</b>		<b>3 years</b>	<b>18 years</b>

# Powerful Hardware

## GOOGLE DATACENTER



1,000 CPU Servers  
2,000 CPUs • 16,000 cores

600 kWatts  
\$5,000,000

## STANFORD AI LAB

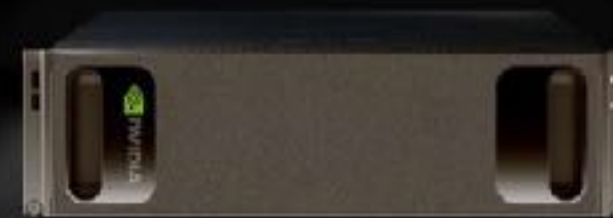


3 GPU Accelerated Servers  
12 GPUs • 18,432 cores

4 kWatts  
\$33,000

## NVIDIA DGX-1

### WORLD'S FIRST DEEP LEARNING SUPERCOMPUTER

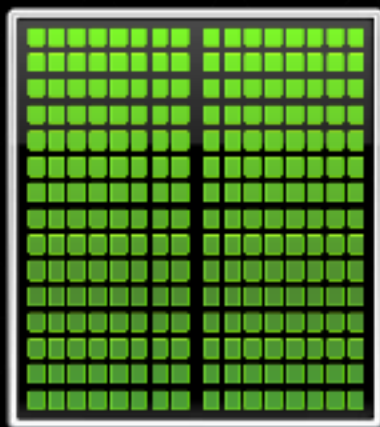


170 TFLOPS FP16  
8x Tesla P100 16GB  
NVLink Hybrid Cube Mesh  
Accelerates Major AI Frameworks  
Dual Xeon  
7 TB SSD Deep Learning Cache  
Dual 10GbE, Quad IB 100Gb  
3RU - 3200W

**CPU**  
Optimized for  
Serial Tasks



**GPU Accelerator**  
Optimized for  
Parallel Tasks



## TITAN X

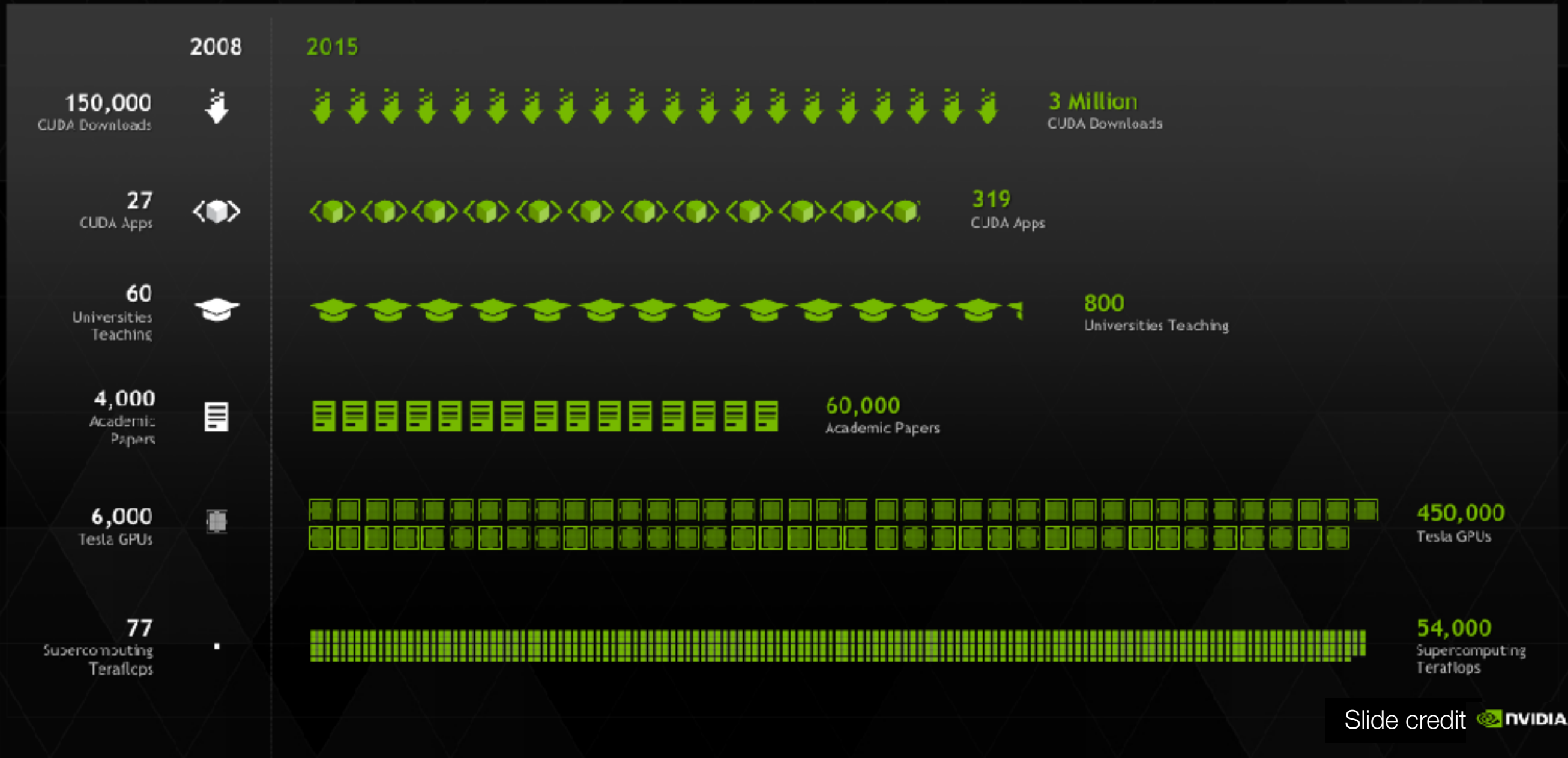
### THE WORLD'S FASTEST GPU

8 Billion Transistors  
3,072 CUDA Cores  
7 TFLOPS SP / 0.2 TFLOPS DP  
12GB Memory



Slide credit  NVIDIA

# 10X GROWTH IN GPU COMPUTING



Slide credit  NVIDIA



# Working ideas on how to train deep architectures

## Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Nitish Srivastava  
Geoffrey Hinton  
Alex Krizhevsky  
Ilya Sutskever  
Ruslan Salakhutdinov

NITISH@CS.TORONTO.EDU  
HINTON@CS.TORONTO.EDU  
KRIZ@CS.TORONTO.EDU  
ILYA@CS.TORONTO.EDU  
RSALAKHU@CS.TORONTO.EDU

### Abstract

Deep neural nets with a large number of parameters are very powerful machine learning systems. However, overfitting is a serious problem in such networks. Large networks are also slow to use, making it difficult to deal with overfitting by combining the predictions of many different large neural nets at test time. Dropout is a technique for addressing this problem. The key idea is to randomly drop units (along with their connections) from the neural network during training. This prevents units from co-adapting too much. During training, dropout samples from an exponential number of different “thinned” networks. At test time,



- Better Learning Regularization (e.g. **Dropout**)

N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, “Dropout: A Simple Way to Prevent Neural Networks from Overfitting”, JMLR Vol. 15, No. 1,

# Working ideas on how to train deep architectures

## Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe  
Google Inc., [sioffe@google.com](mailto:sioffe@google.com)

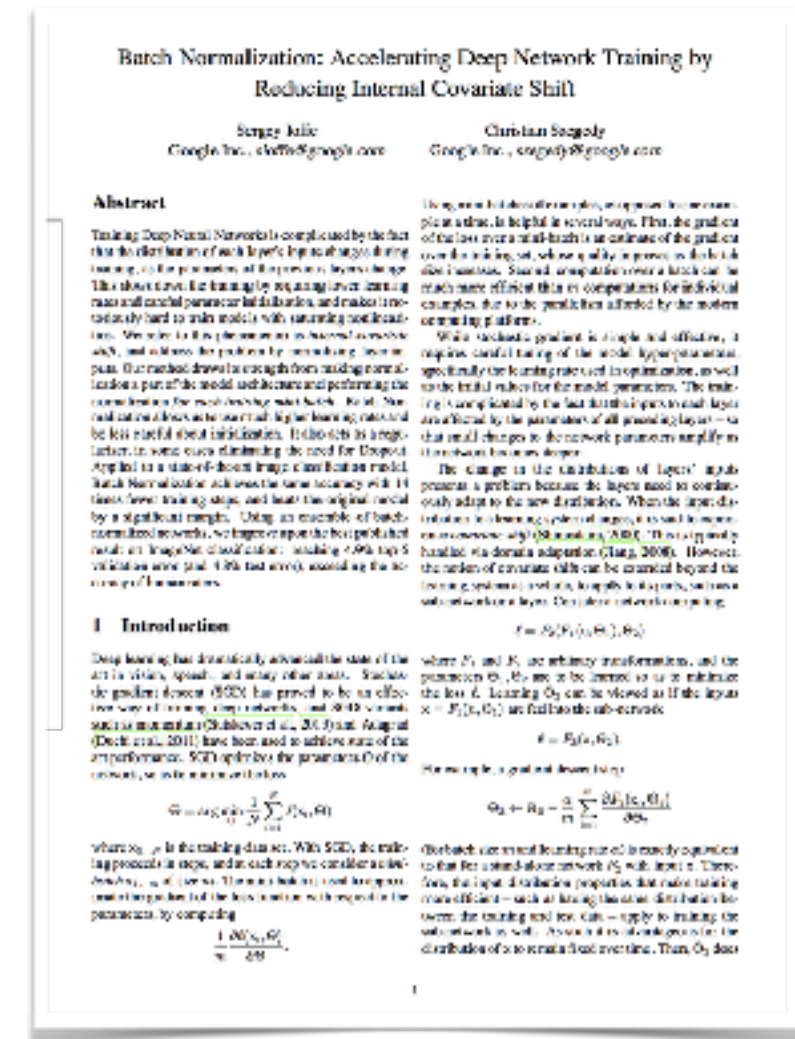
Christian Szegedy  
Google Inc., [szegedy@google.com](mailto:szegedy@google.com)

### Abstract

Training Deep Neural Networks is complicated by the fact that the distribution of each layer's inputs changes during training, as the parameters of the previous layers change. This slows down the training by requiring lower learning rates and careful parameter initialization, and makes it notoriously hard to train models with saturating nonlinearities. We refer to this phenomenon as *internal covariate shift*, and address the problem by normalizing layer inputs. Our method draws its strength from making normalization a part of the model architecture and performing the normalization for each training mini-batch. Batch Normalization

Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways. First, the gradient of the loss over a mini-batch is an estimate of the gradient over the training set, whose quality improves as the batch size increases. Second, computation over a batch can be much more efficient than  $m$  computations for individual examples, due to the parallelism afforded by the modern computing platforms.

While stochastic gradient is simple and effective, it requires careful tuning of the model hyper-parameters, specifically the learning rate used in optimization, as well as the initial values for the model parameters. The training is complicated by the fact that the inputs to each layer



- Better Optimization Conditioning (e.g. **Batch Normalization**)

S. Ioffe, C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", In ICML 2015



# Working ideas on how to train deep architectures

## Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun  
Microsoft Research  
[kahe, v-xiangz, v-shren, jiansun}@microsoft.com

### Abstract

Deeper neural networks are more difficult to train. We present a residual learning framework to ease the training of networks that are substantially deeper than those used previously. We explicitly reformulate the layers as learning residual functions with reference to the layer inputs, instead of learning unreferenced functions. We provide comprehensive empirical evidence showing that these residual networks are easier to optimize, and can gain accuracy from considerably increased depth. On the ImageNet dataset we evaluate residual nets with a depth of up to 152 layers—8× deeper than VGG nets [41] but still having lower complexity. An ensemble of these residual nets achieves 3.57% error

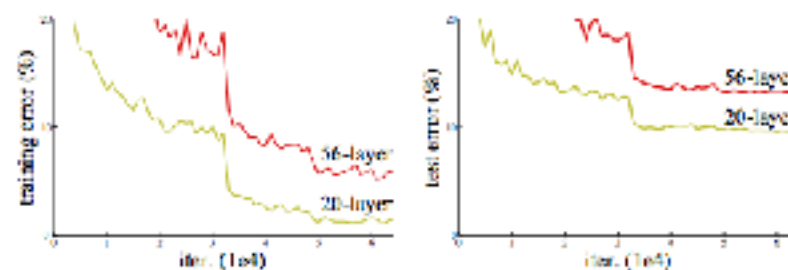
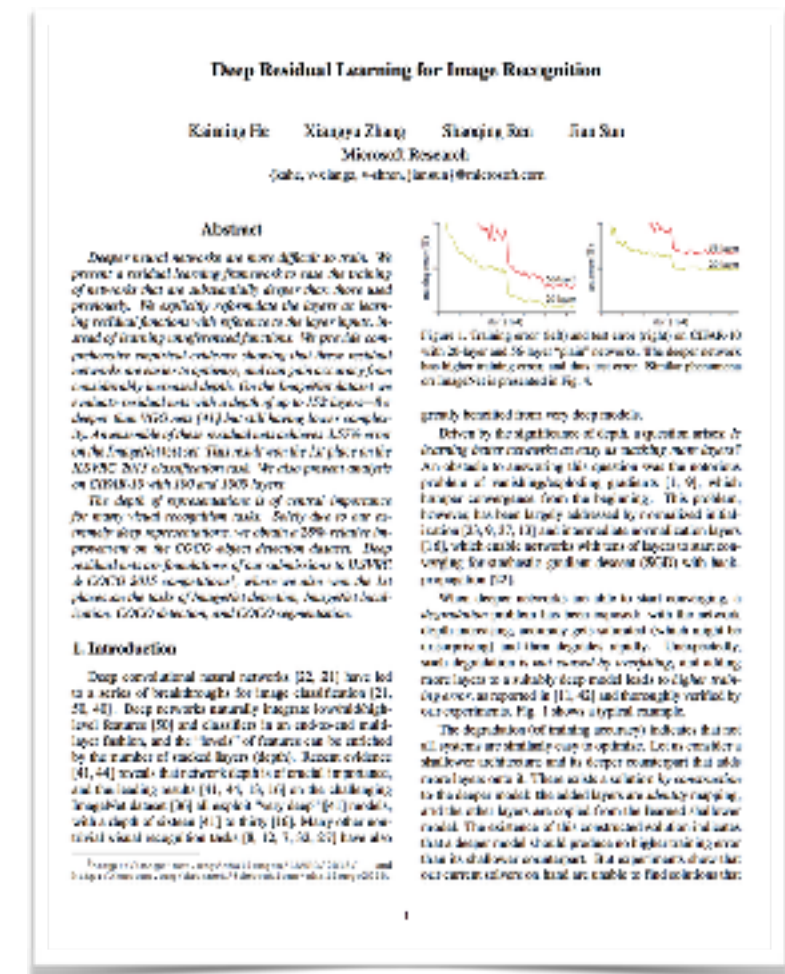


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

greatly benefited from very deep models.

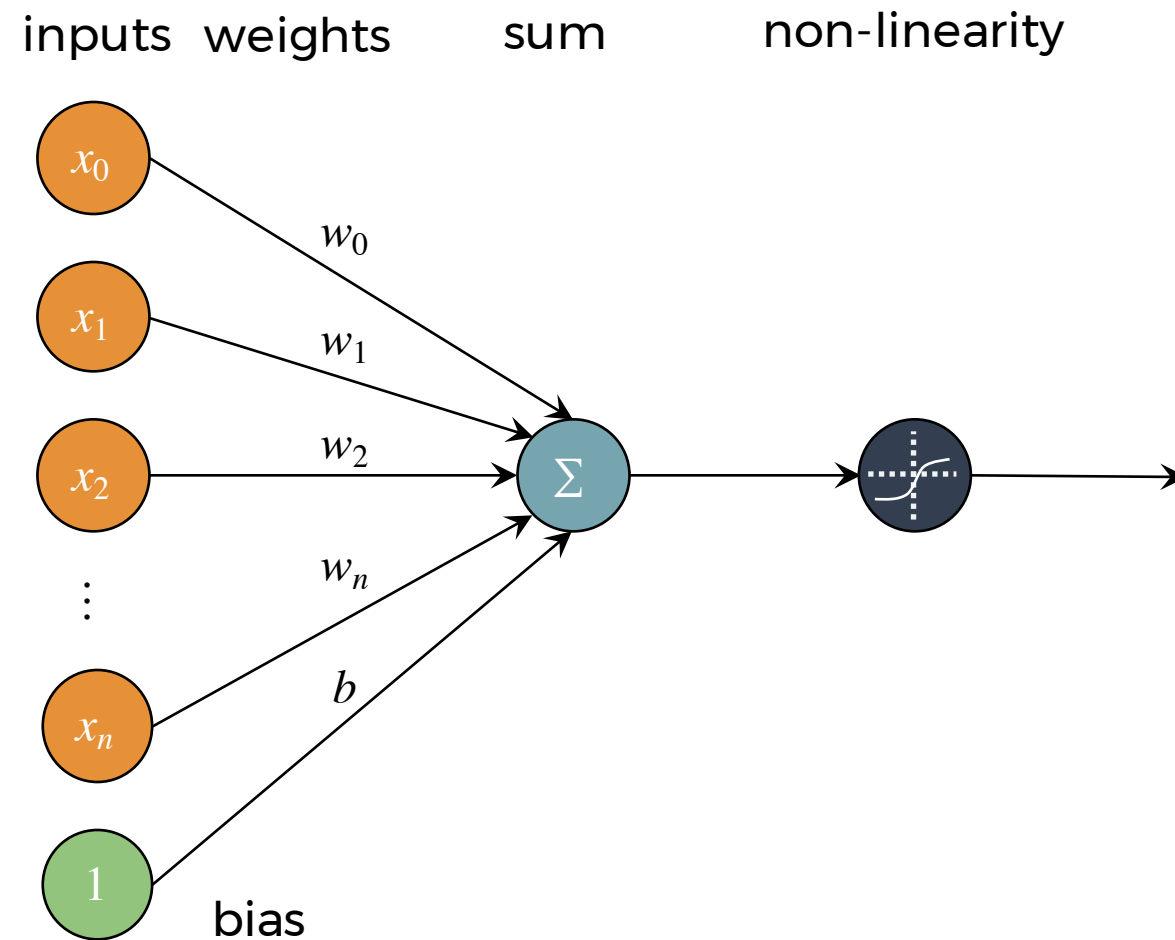
Driven by the significance of depth, a question arises: *Is*



- Better neural architectures (e.g. **Residual Nets**)

Let's make a review  
of neural networks

# The Perceptron





# Perceptron Forward Pass

- Neuron pre-activation  
(or input activation)

$$a(\mathbf{x}) = b + \sum_i w_i x_i = b + \mathbf{w}^\top \mathbf{x}$$

- Neuron output activation:

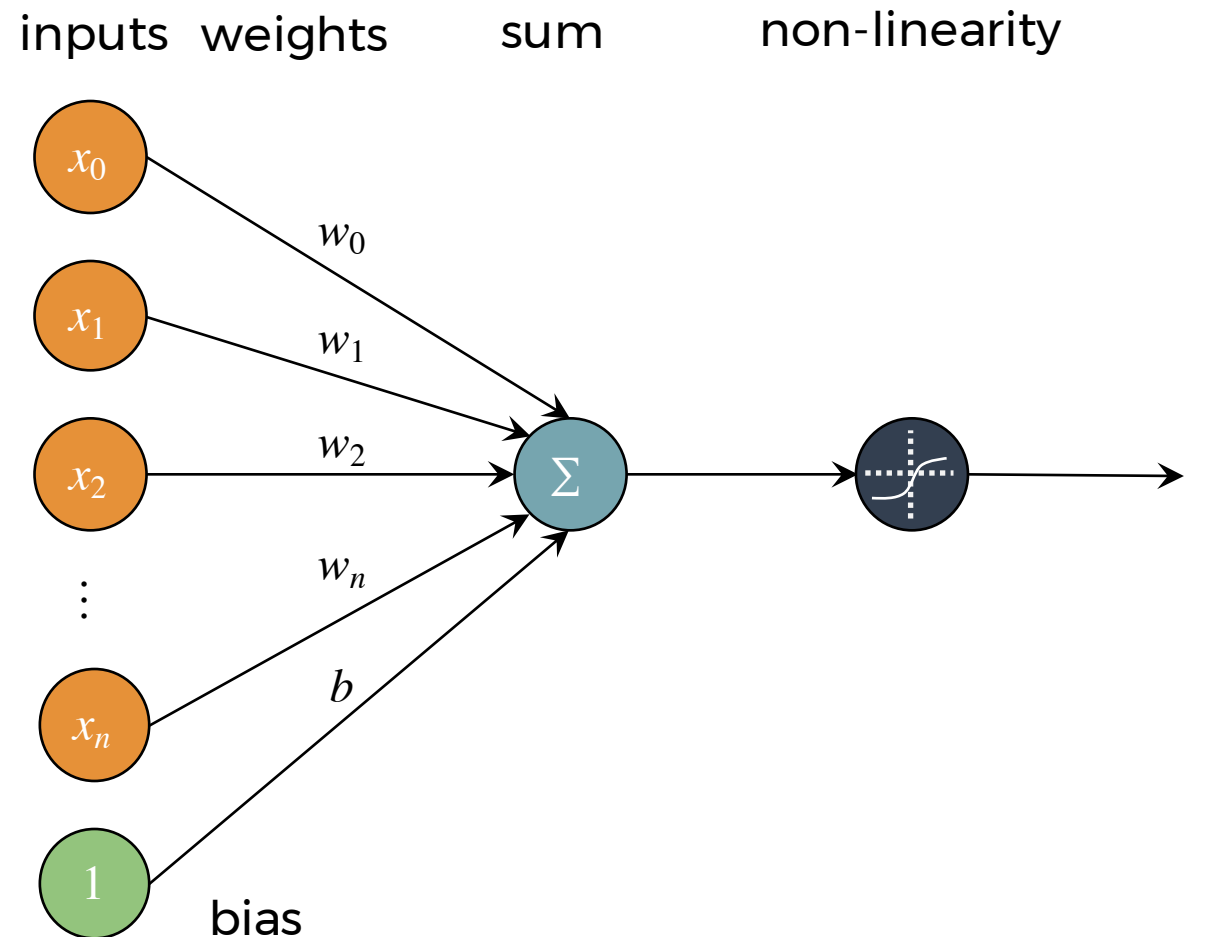
$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

where

$w$  are the weights (parameters)

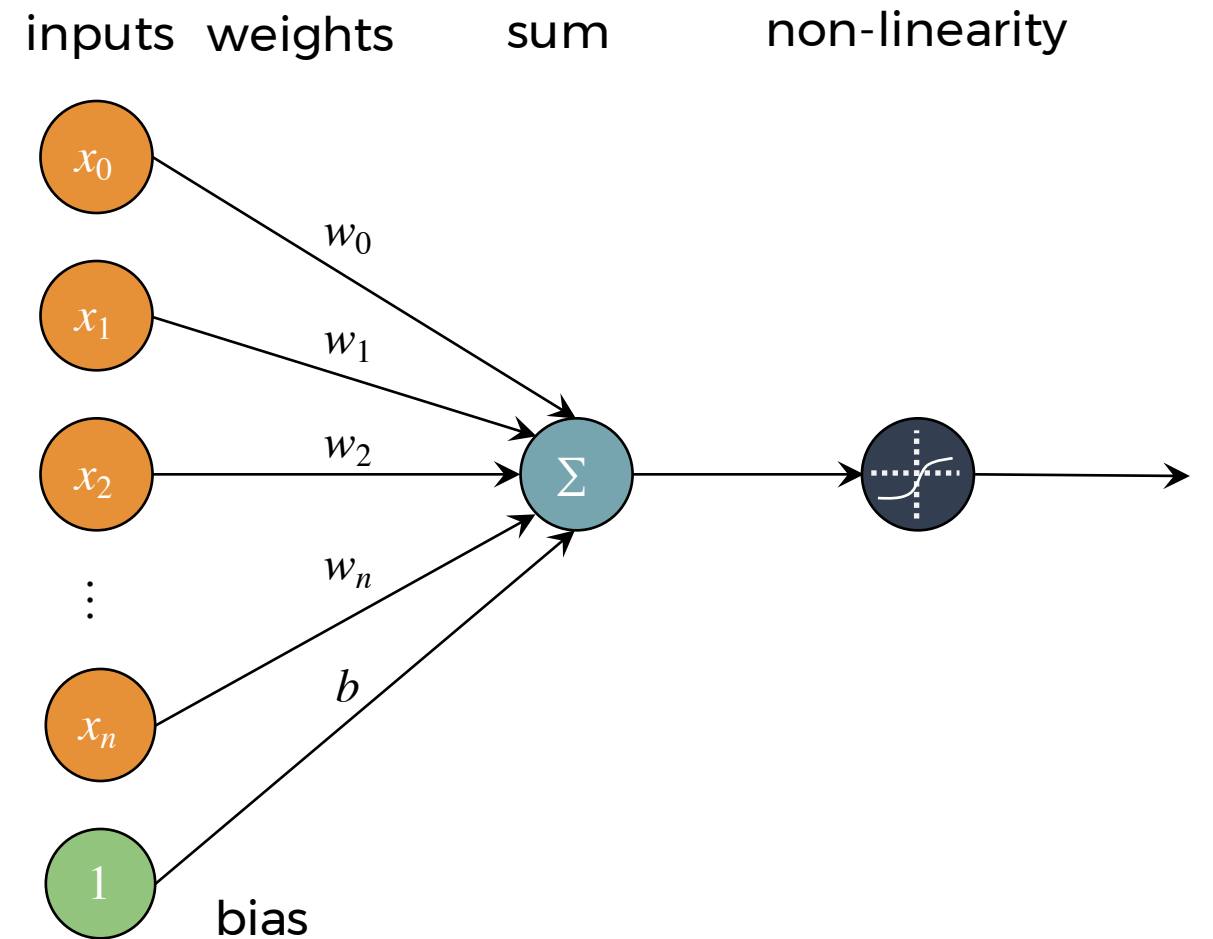
$b$  is the bias term

$g(\cdot)$  is called the activation function

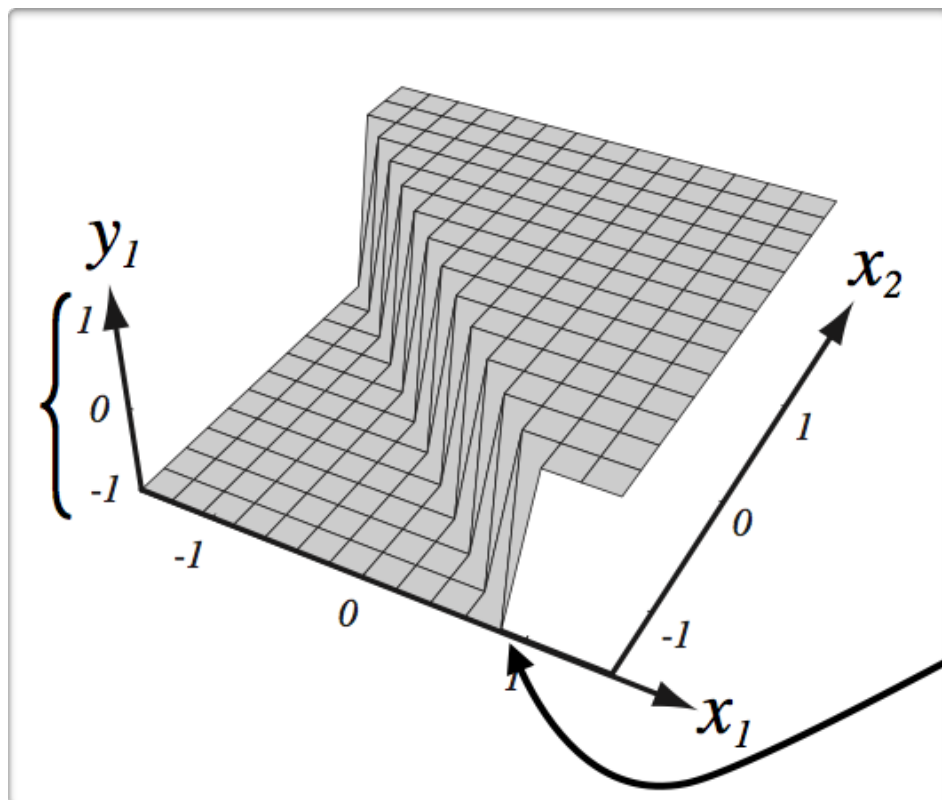


# Output Activation of The Neuron

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$



Range is determined by  $g(\cdot)$

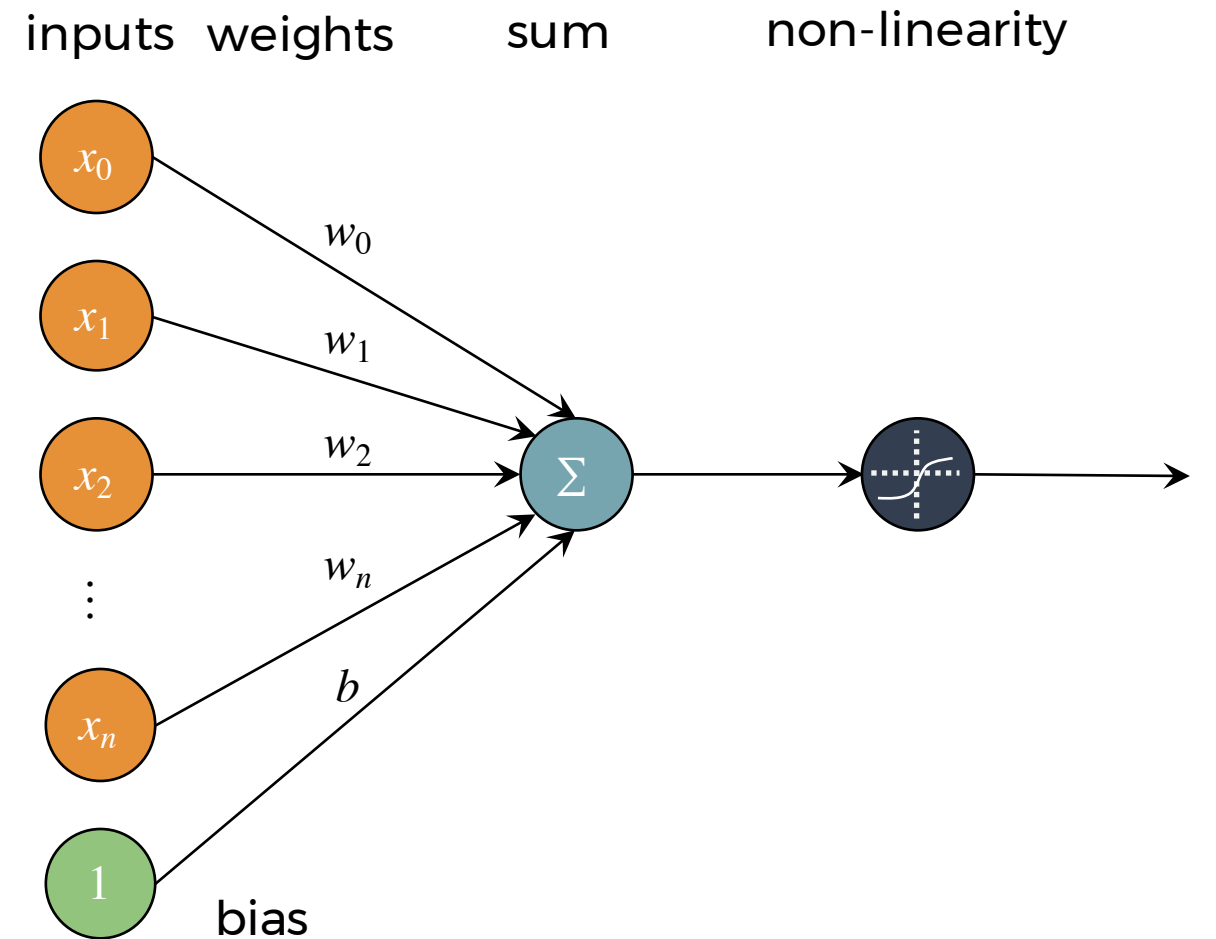
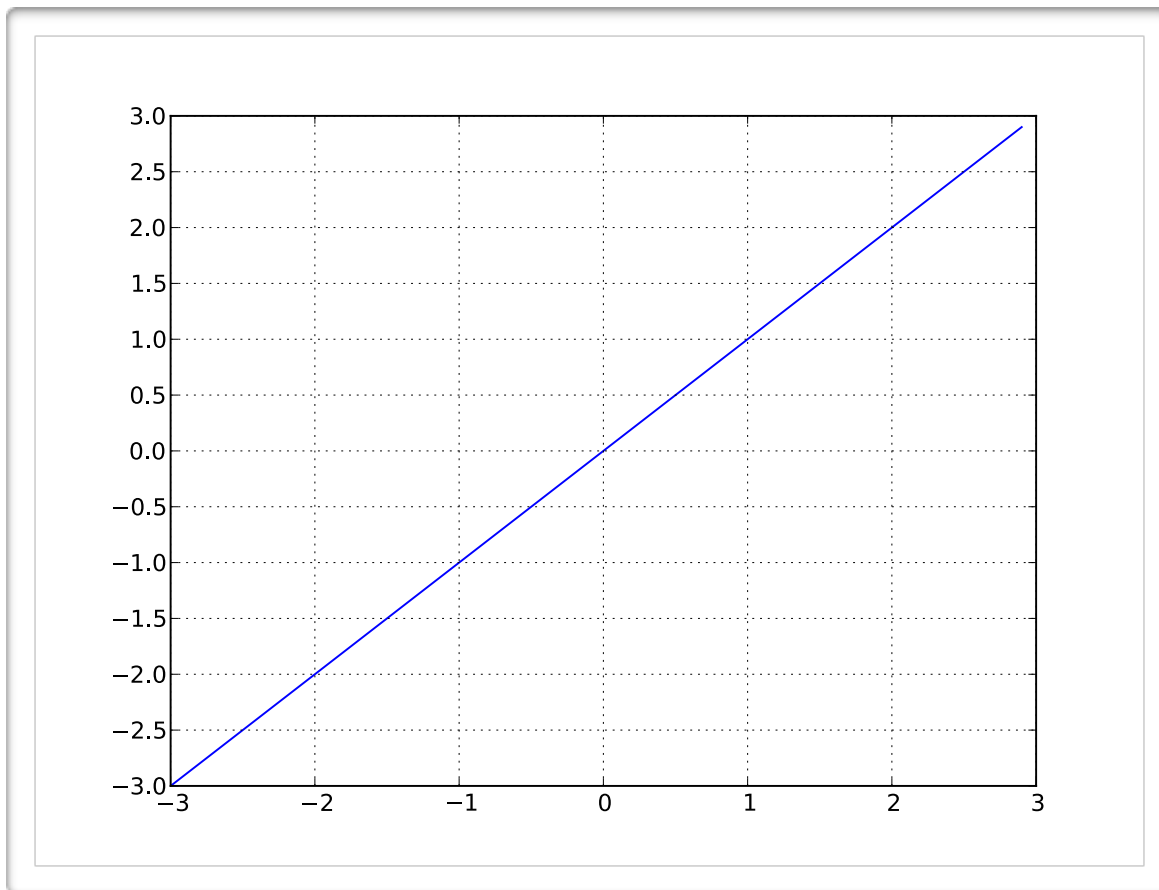


Bias only changes the position of the riff

# Linear Activation Function

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

$$g(a) = a$$

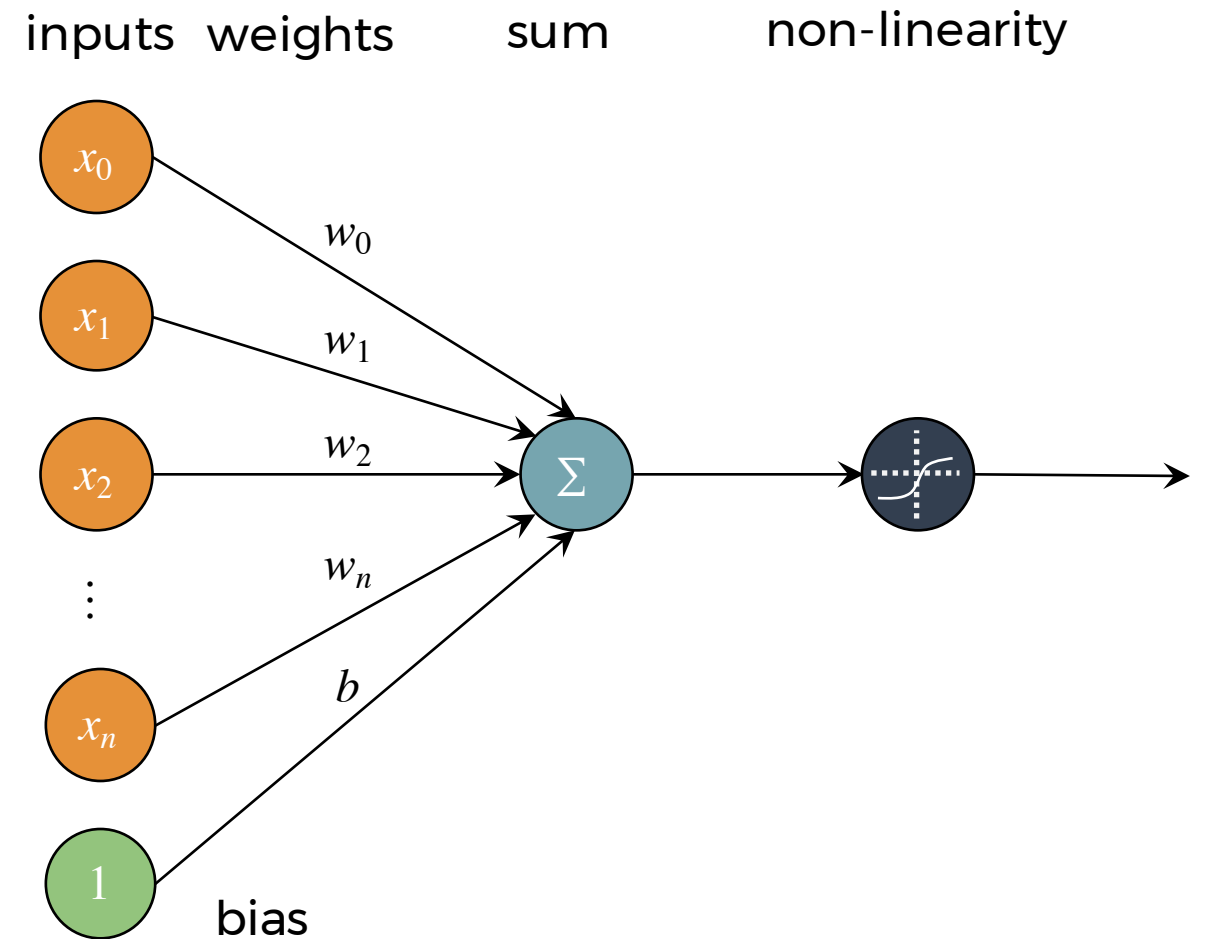
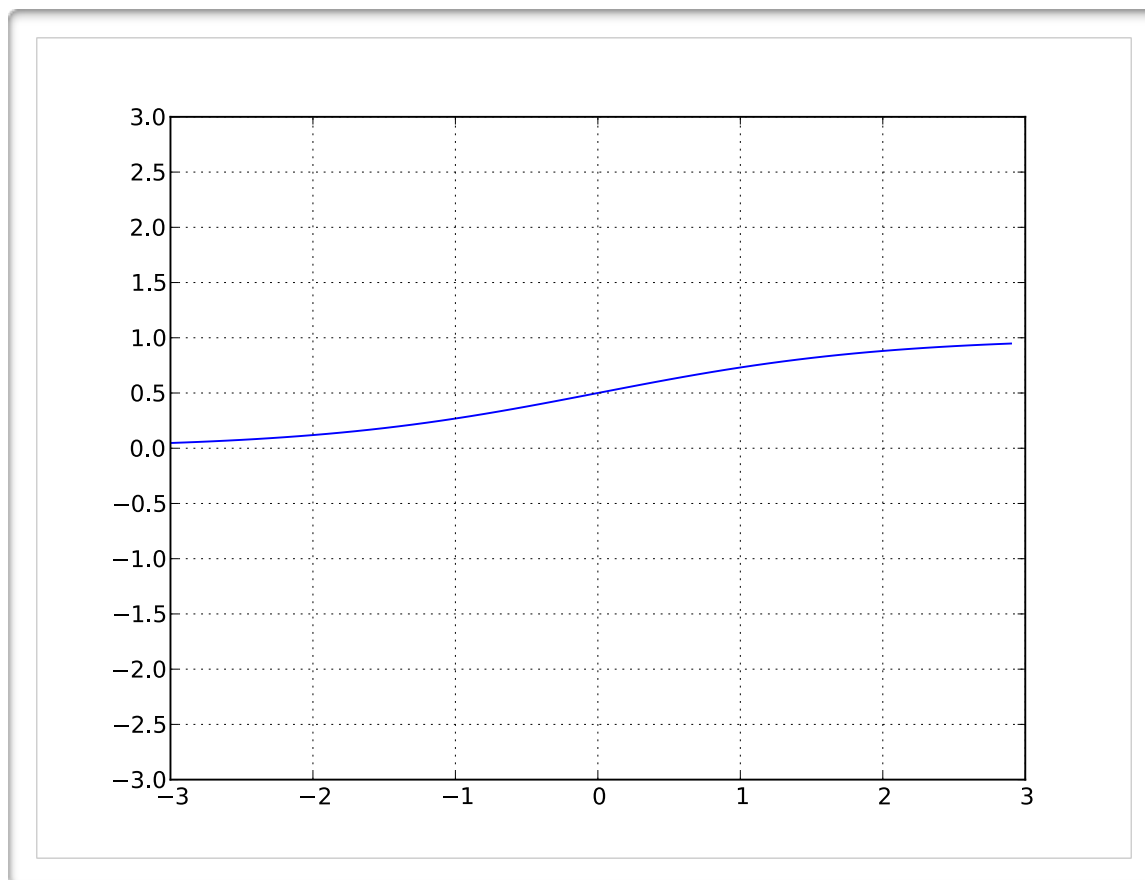


No nonlinear transformation  
No input squashing

# Sigmoid Activation Function

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

$$g(a) = \text{sigm}(a) = \frac{1}{1 + \exp(-a)}$$



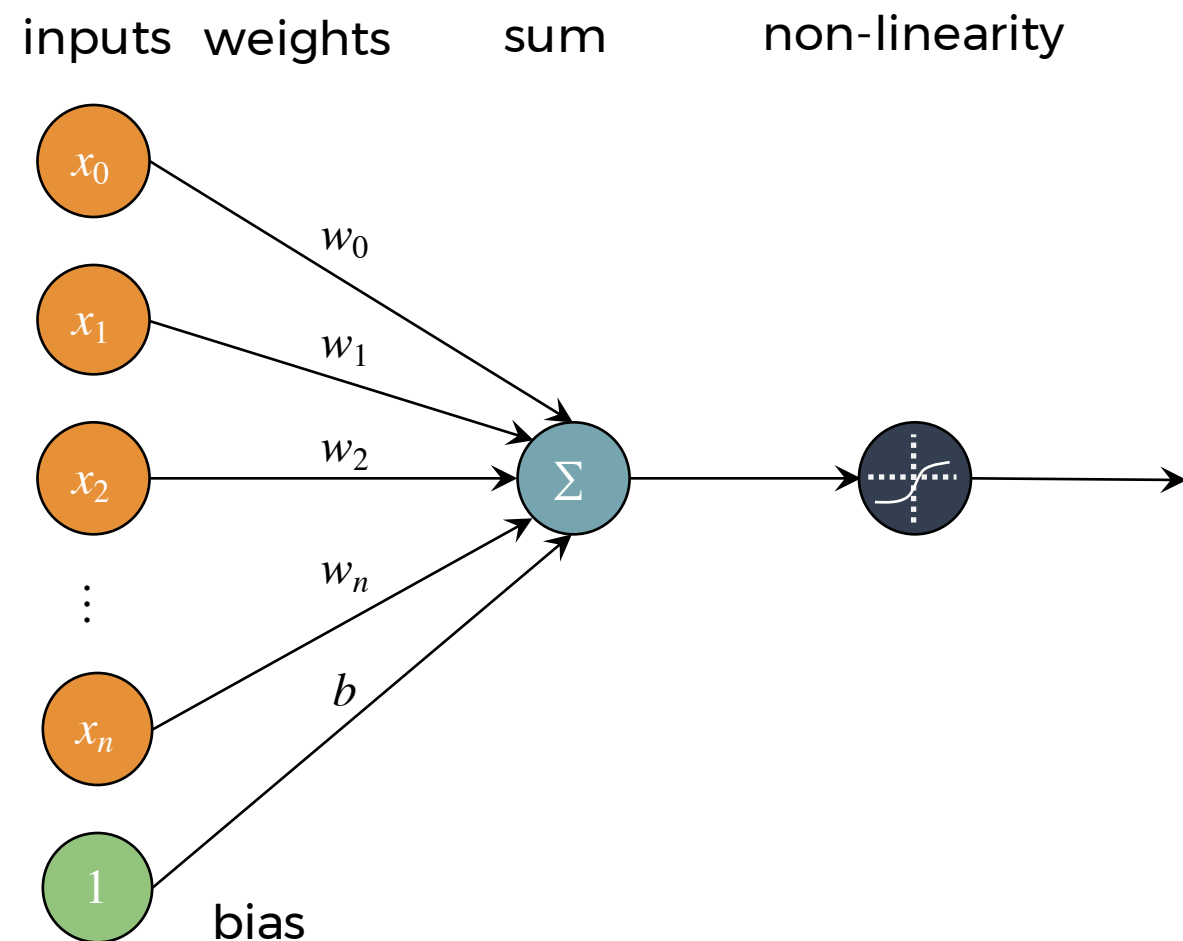
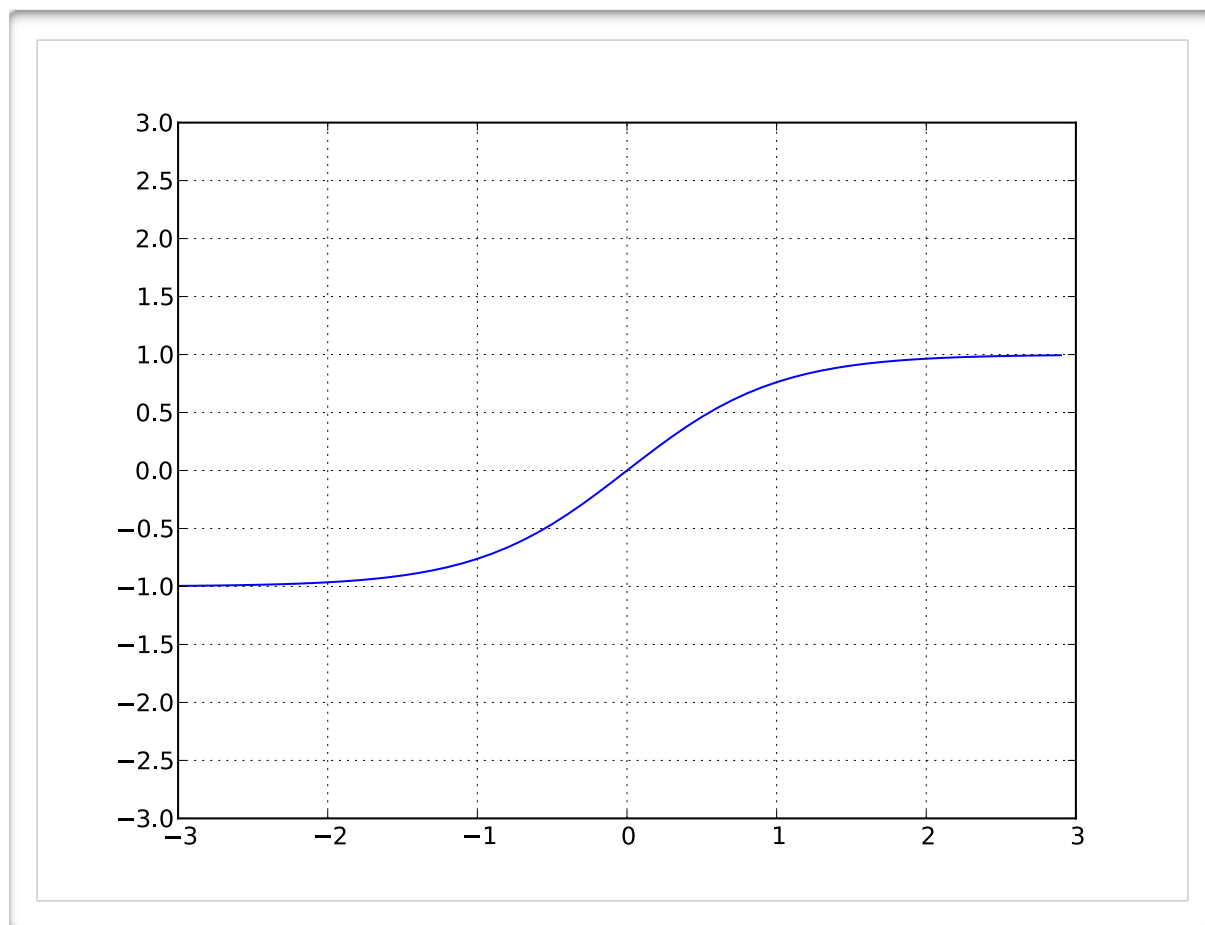
Squashes the neuron's output  
between 0 and 1  
Always positive  
Bounded  
Strictly Increasing

# Hyperbolic Tangent (tanh) Activation Function

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

$$g(a) = \tanh(a) =$$

$$= \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} = \frac{\exp(2a) - 1}{\exp(2a) + 1}$$



Squashes the neuron's output  
between

-1 and 1

Can be positive or negative

Bounded

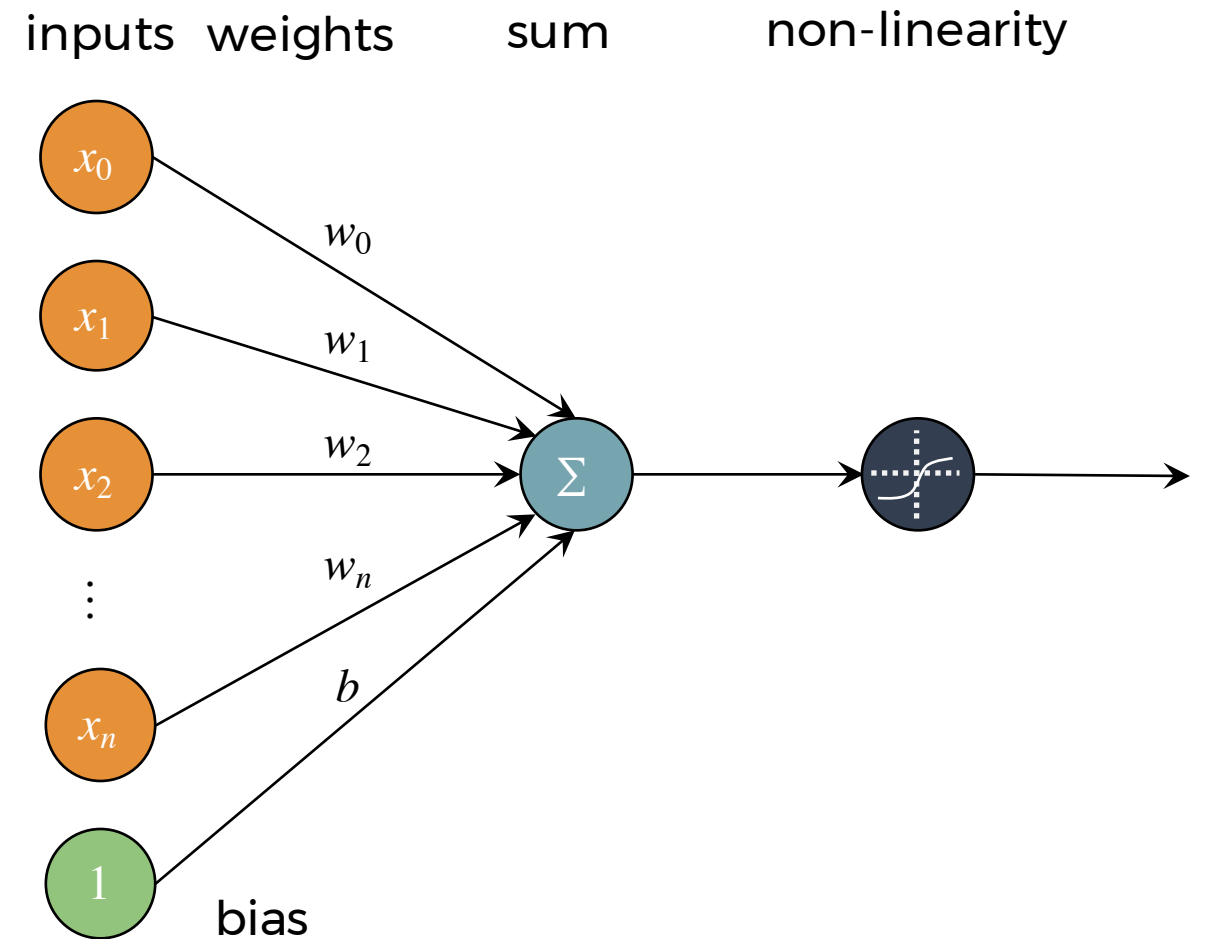
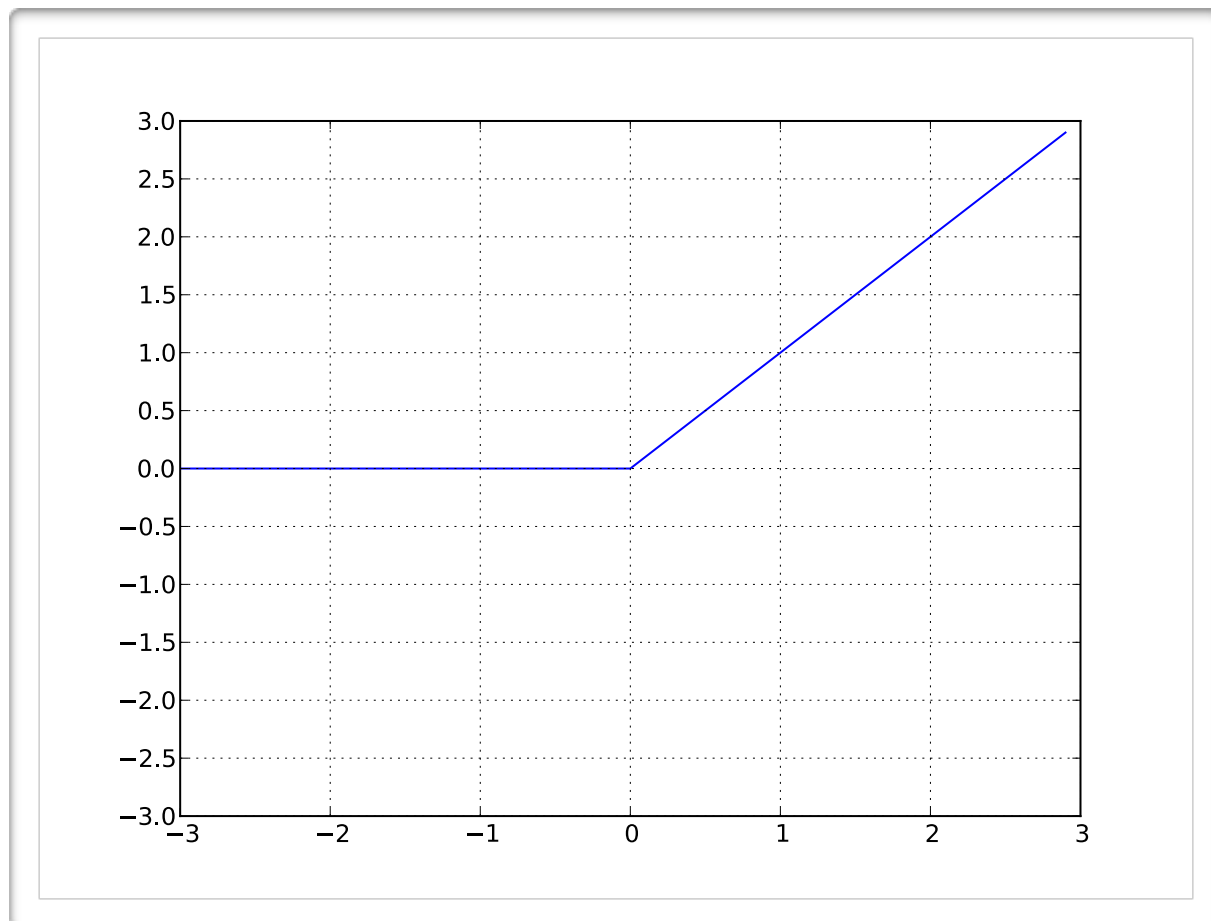
Strictly Increasing



# Rectified Linear (ReLU) Activation Function

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

$$g(a) = \text{reclin}(a) = \max(0, a)$$



Bounded below by 0 (always non-negative)

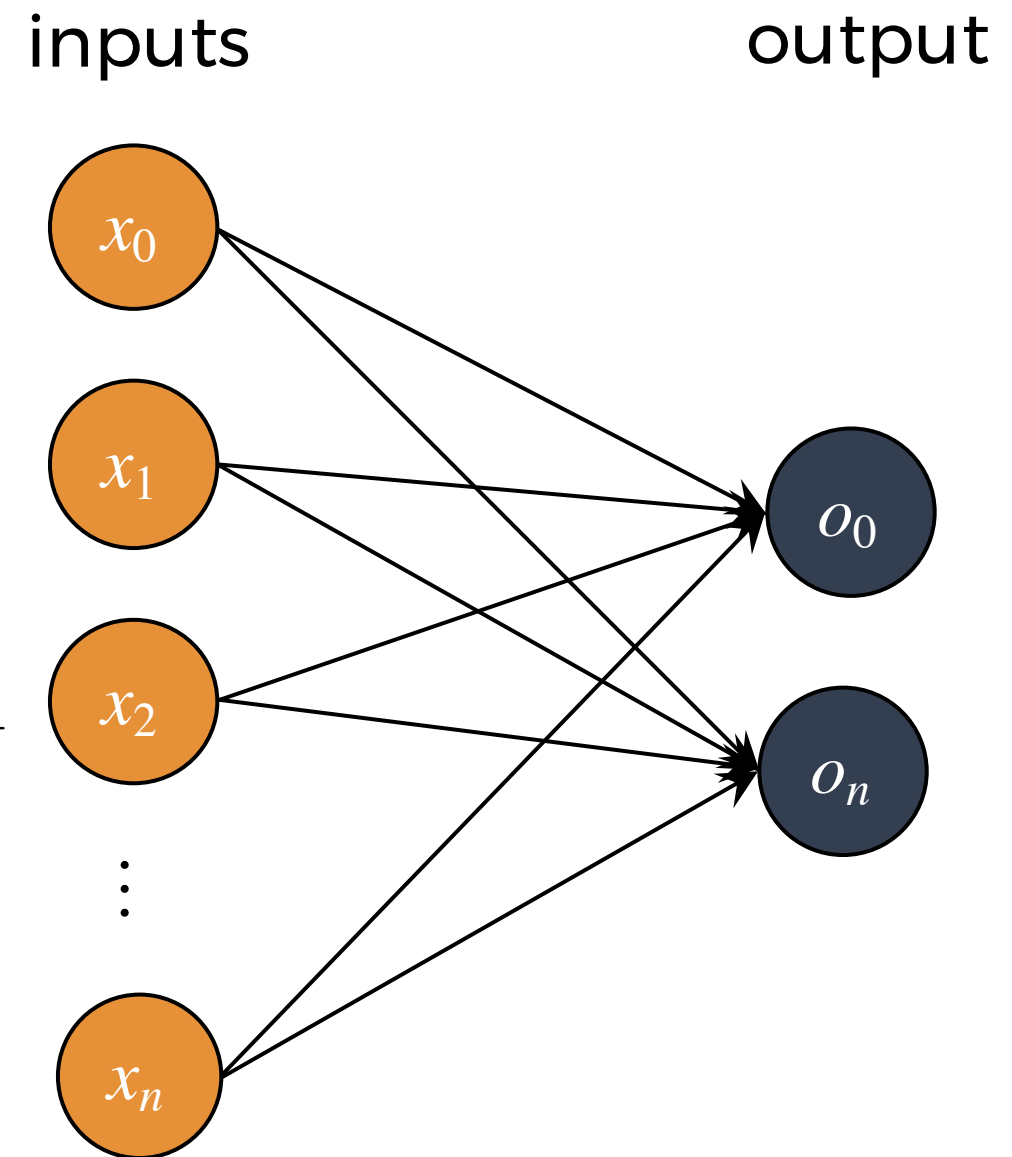
Not upper bounded

Strictly increasing

Tends to produce units with sparse activities

# Multi-Output Perceptron

- We need multiple outputs  
(1 output per class)  $i, j$
- We need to estimate conditional probability  
 $p(y = c|x)$
- Discriminative Learning
- Softmax activation function at the output
  - Strictly positive
  - sums to one
- Predict class with the highest estimated class conditional probability.



$$\mathbf{o}(\mathbf{a}) = \text{softmax}(\mathbf{a}) = \left[ \frac{\exp(a_1)}{\sum_c \exp(a_c)} \cdots \frac{\exp(a_C)}{\sum_c \exp(a_c)} \right]^T$$

# Single Hidden Layer Neural Network

- Hidden layer pre-activation:

$$\mathbf{a}(\mathbf{x}) = \mathbf{b}^{(1)} + \mathbf{W}^{(1)}\mathbf{x}$$

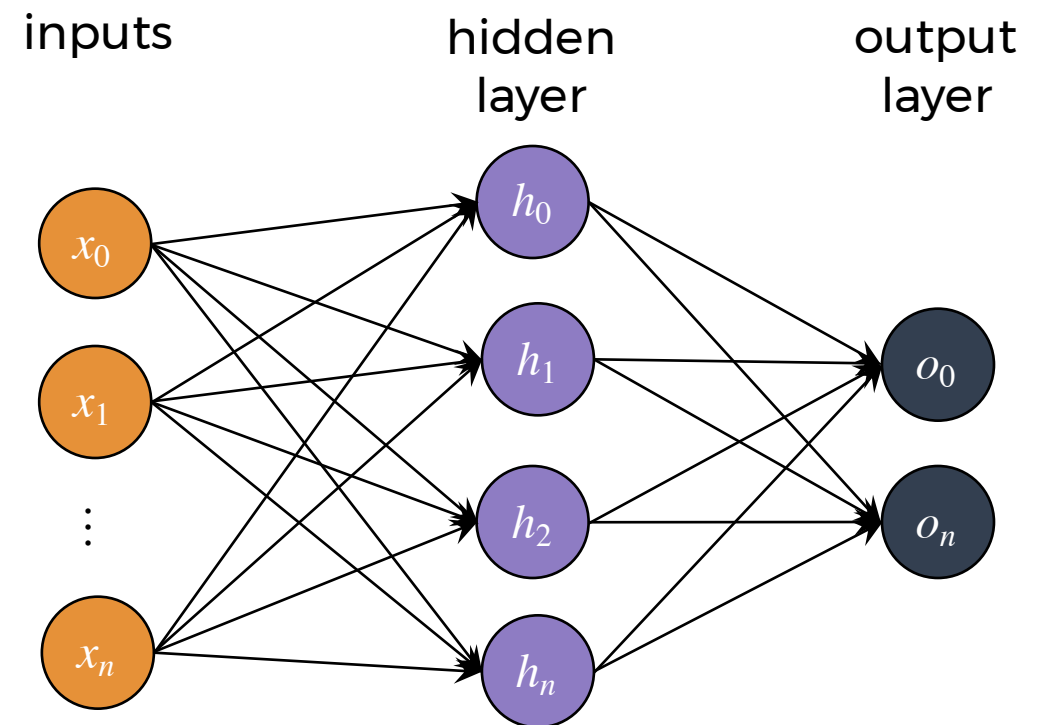
$$\left( a(\mathbf{x})_i = b_i^{(1)} + \sum_j W_{i,j}^{(1)} x_j \right)$$

- Hidden layer activation:

$$\mathbf{h}(\mathbf{x}) = \mathbf{g}(\mathbf{a}(\mathbf{x}))$$

- Output layer activation:

$$\mathbf{o}(\mathbf{x}) = \mathbf{o} \left( b^{(2)} + \mathbf{w}^{(2)} \mathbf{h}^{(1)} \mathbf{x} \right)$$



# Multi-Layer Perceptron (MLP)

Consider a network with L hidden layers.

- layer pre-activation for  $k > 0$

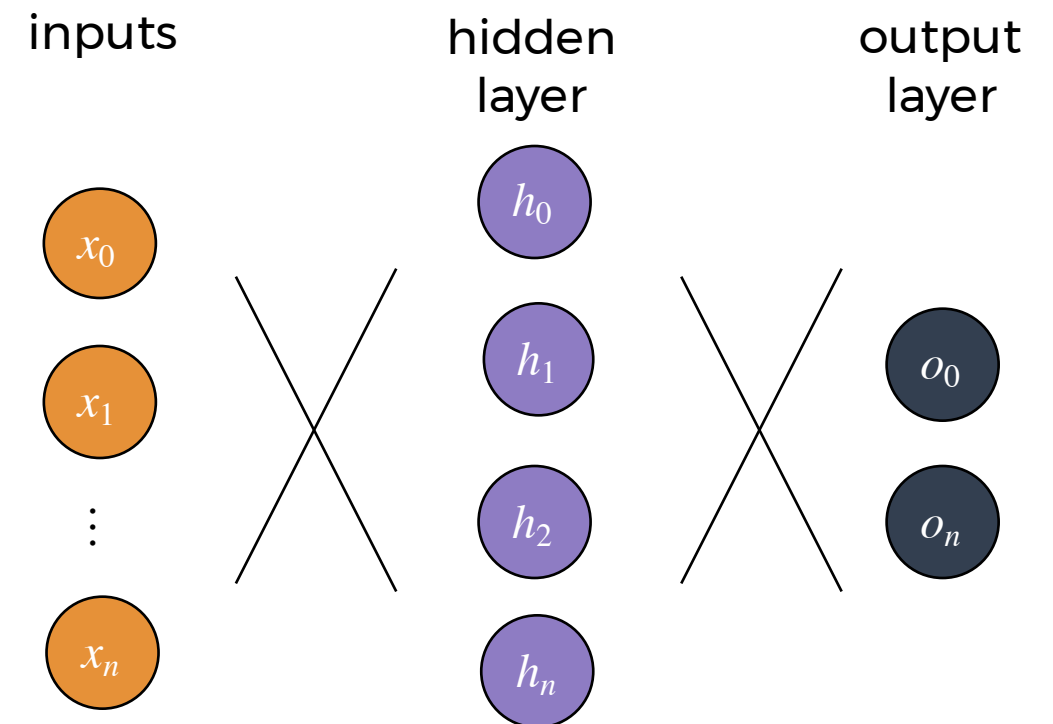
$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$$

- hidden layer activation from 1 to L:

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

- output layer activation ( $k=L+1$ )

$$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$



# Training

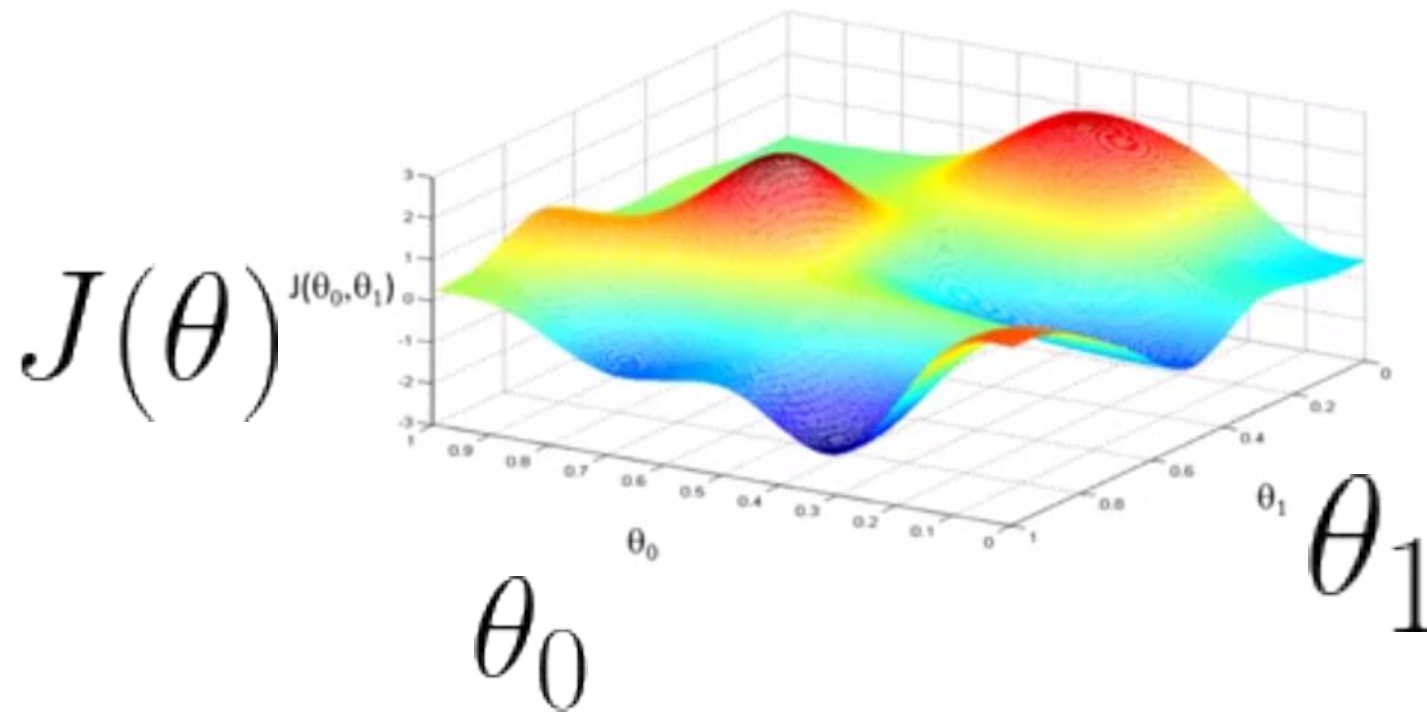
$$J(\theta) = \arg \min_{\theta} \frac{1}{T} \sum_t \underbrace{l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})}_{\text{Loss function}} + \underbrace{\lambda \Omega(\theta)}_{\text{Regularizer}}$$

$$\theta = W_1, W_2 \dots W_n$$

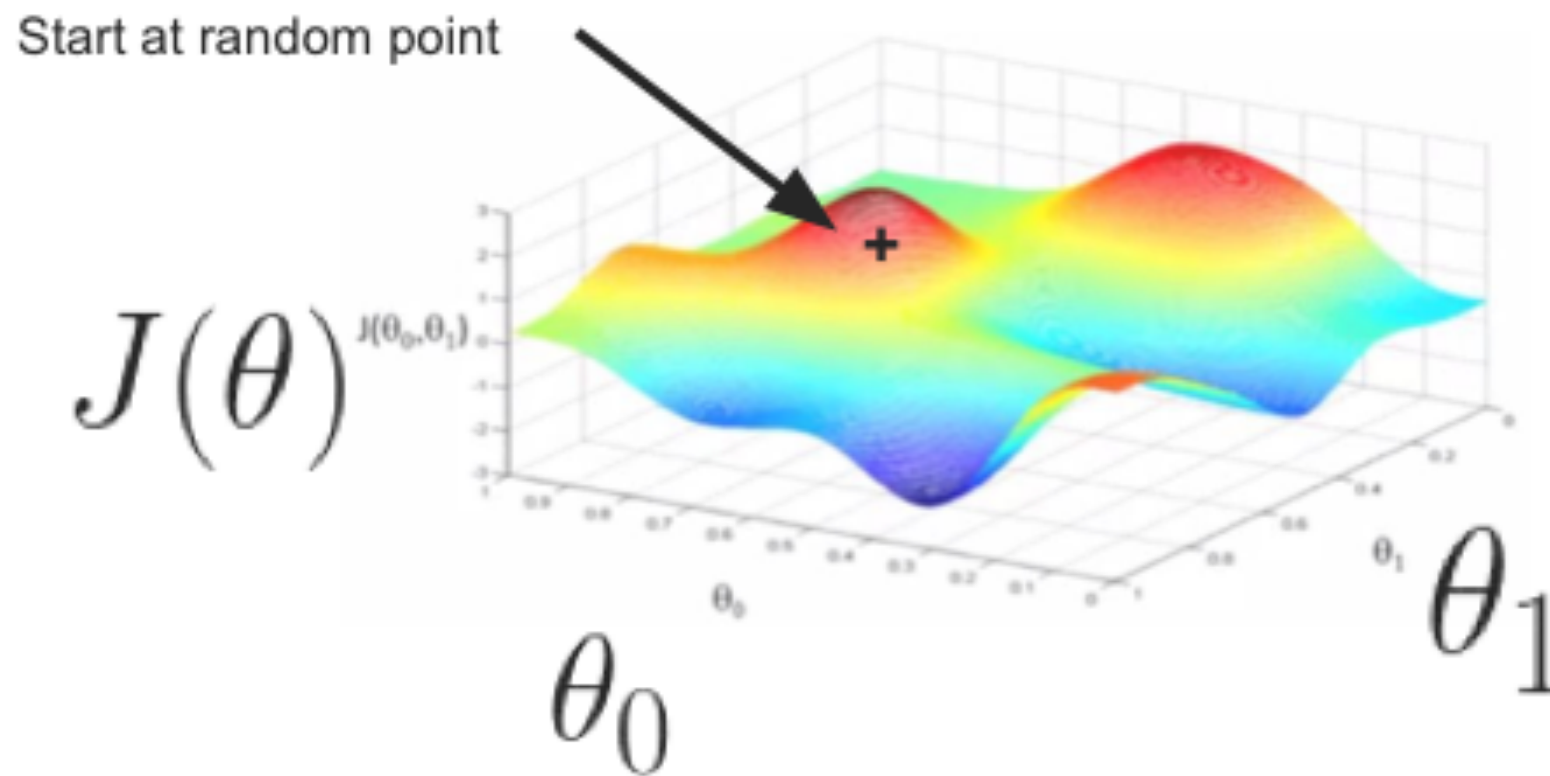
- Learning is cast as optimization.
- For classification problems, we would like to minimize classification error
- Loss function can sometimes be viewed as a surrogate for what we want to optimize (e.g. upper bound)



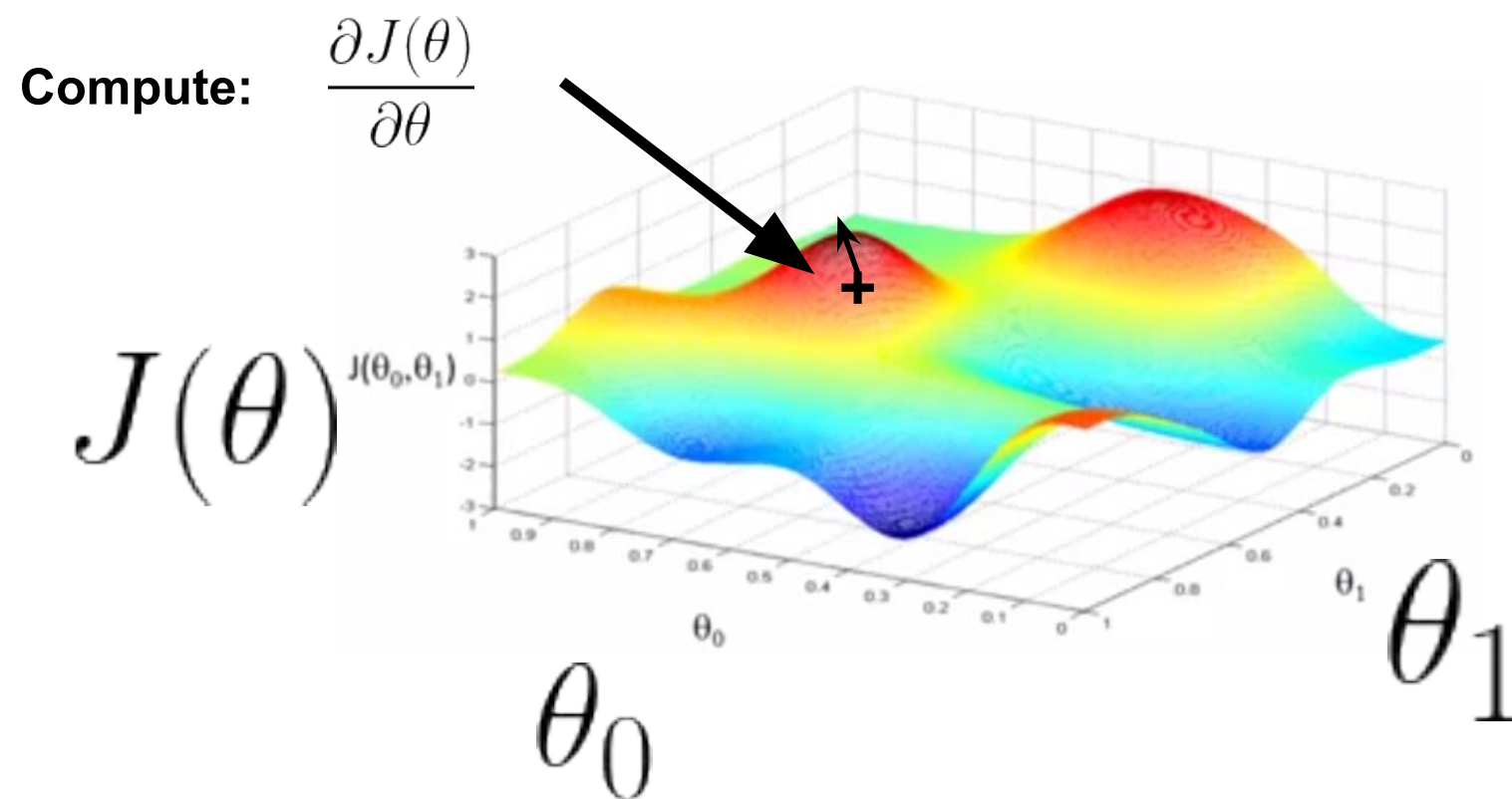
Loss is a **function** of the model's parameters



# How to minimize loss?

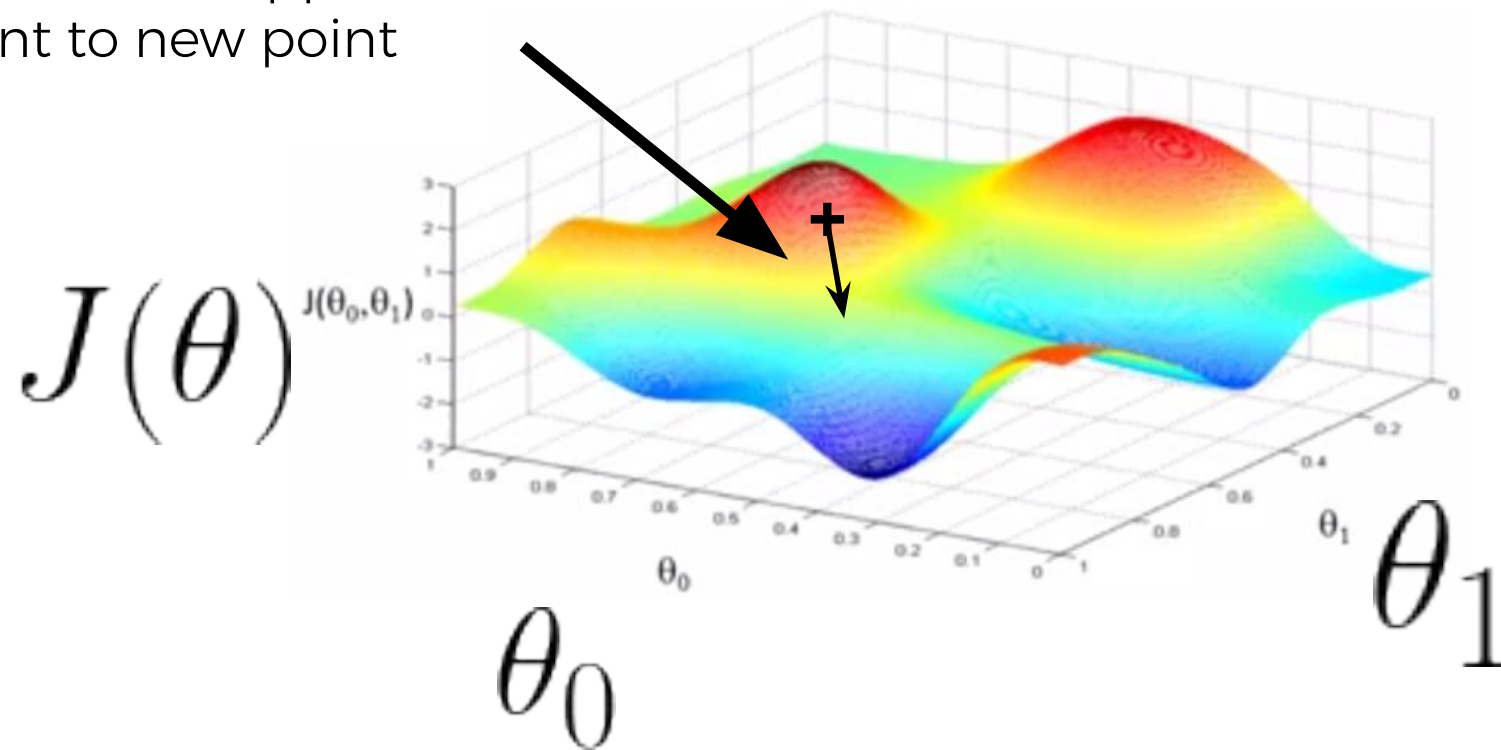


# How to minimize loss?



# How to minimize loss?

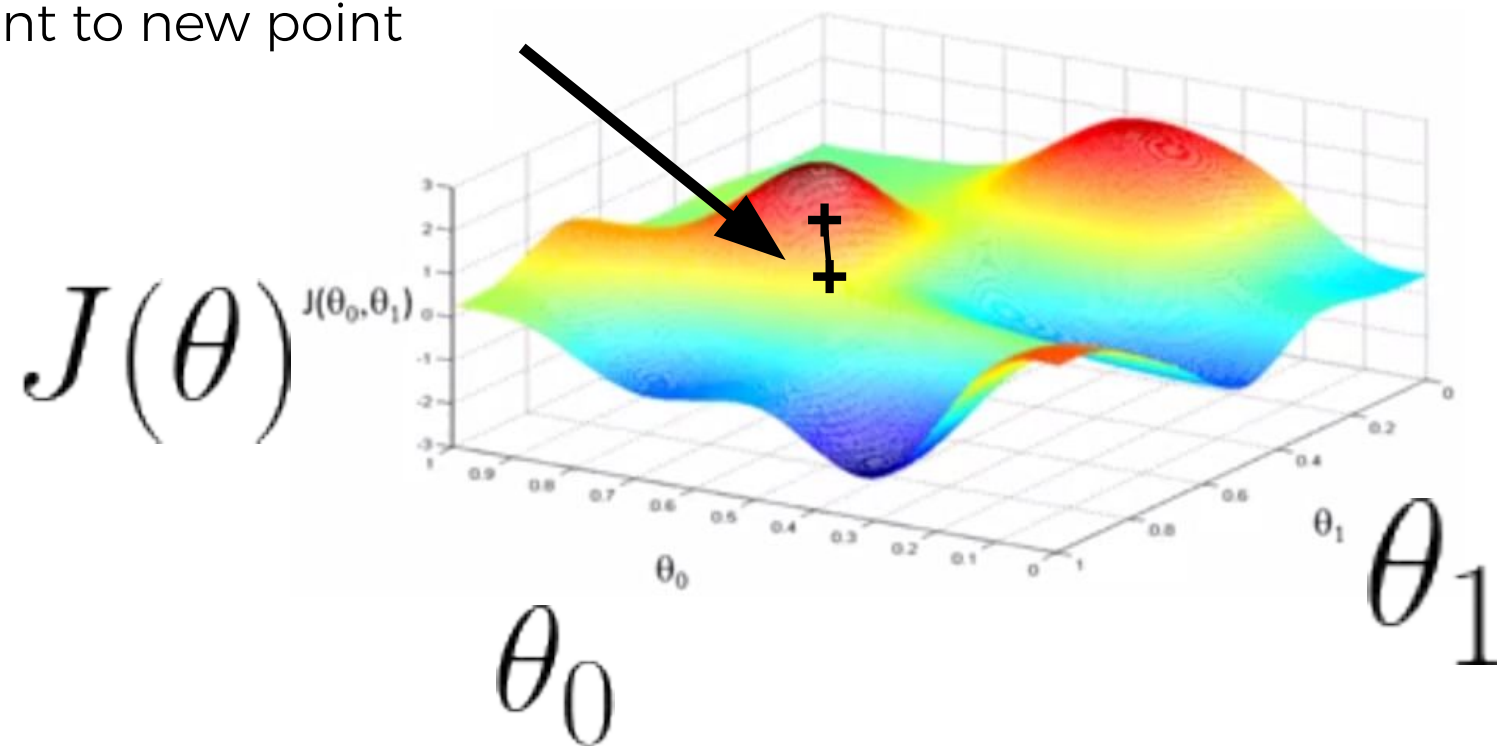
Move in direction opposite of gradient to new point





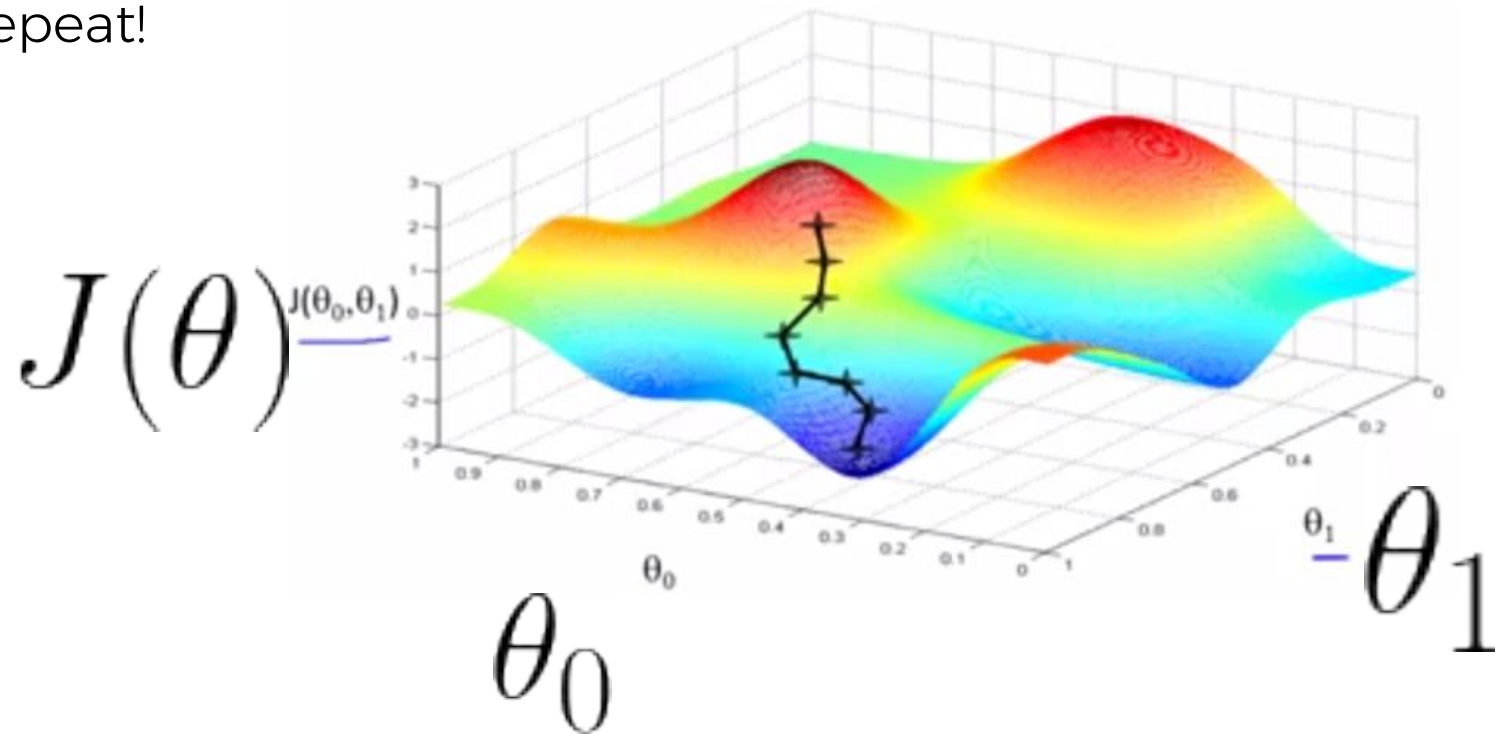
# How to minimize loss?

Move in direction opposite of gradient to new point



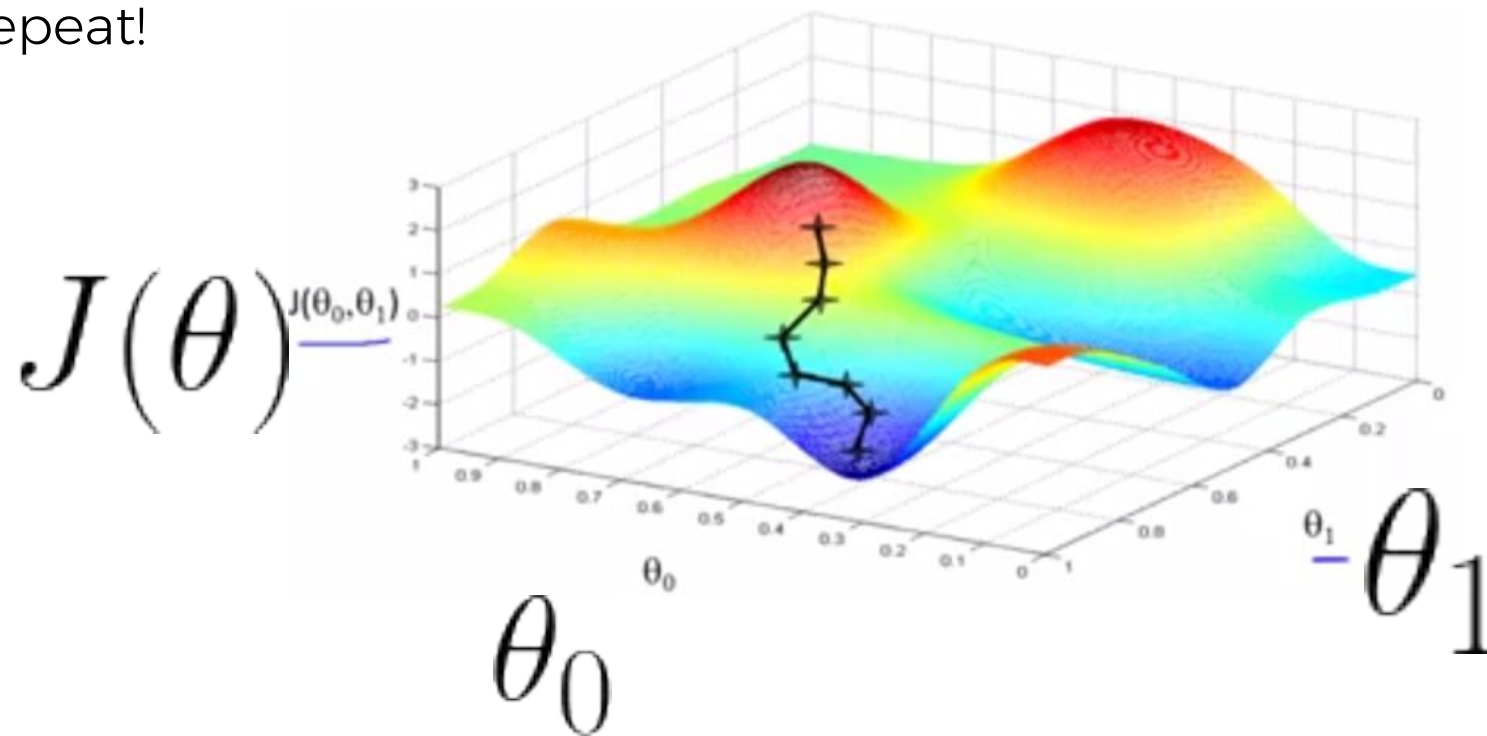
# How to minimize loss?

Repeat!



# This is called Stochastic Gradient Descent (SGD)

Repeat!



# Stochastic Gradient Descent (SGD)

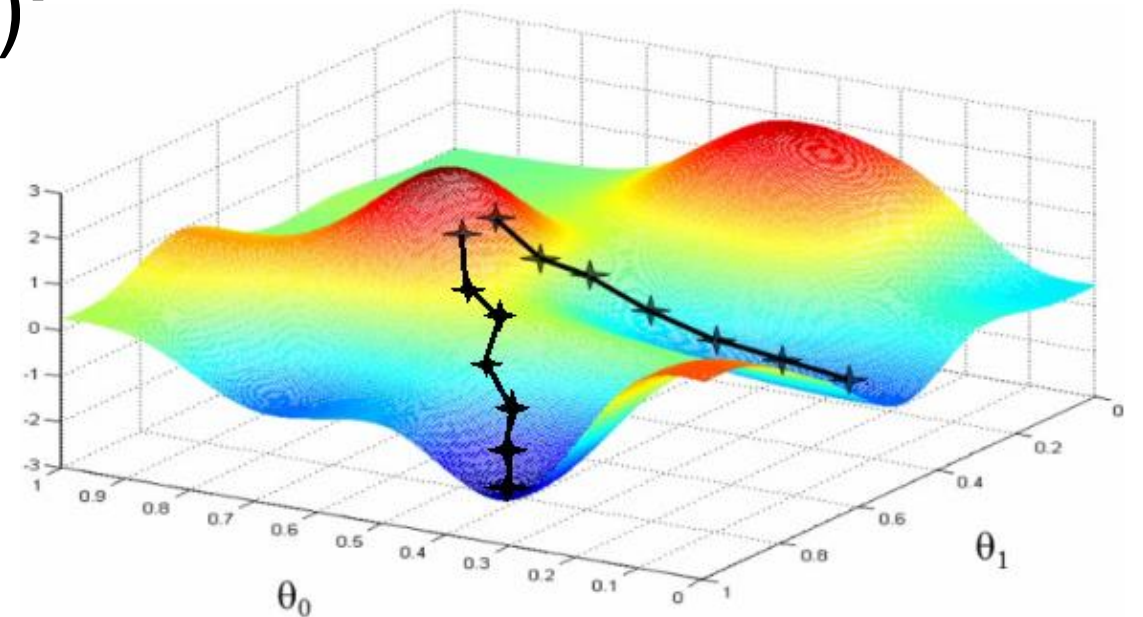
- Initialize  $\theta$  randomly
- For N Epochs
- For each training example  $(x, y)$

- Compute Loss Gradient:

$$\frac{\partial J(\theta)}{\partial \theta}$$

- Update  $\theta$  with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$





# Why is it **Stochastic** Gradient Descent?

- Initialize  $\theta$  randomly
- For N Epochs
- For each training example  $(x, y)$

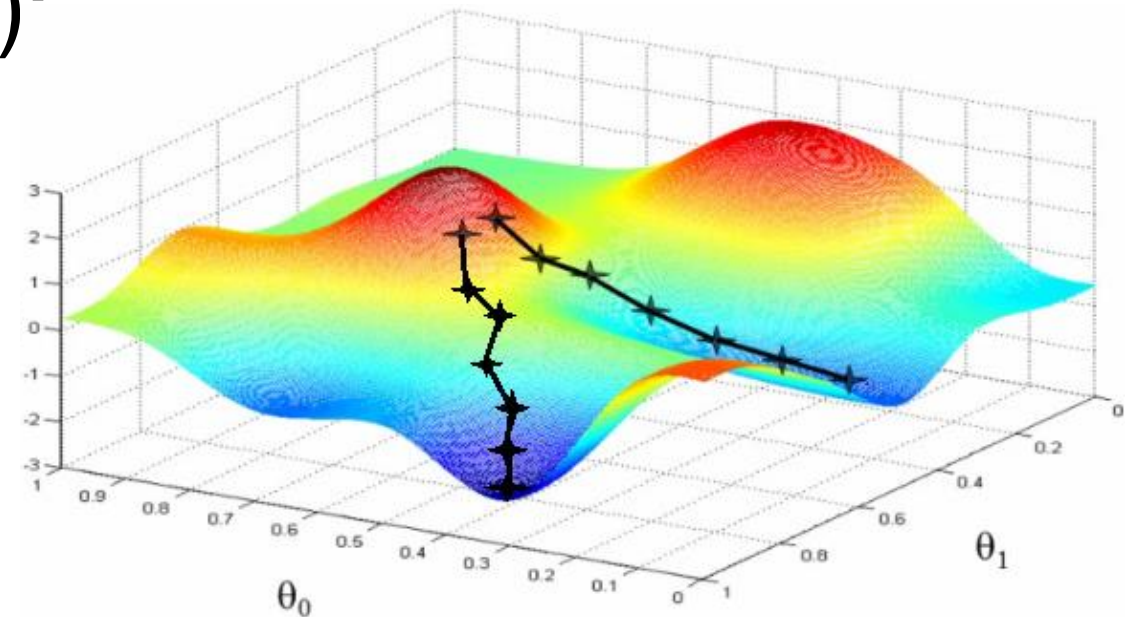
- Compute Loss Gradient:

$$\frac{\partial J(\theta)}{\partial \theta}$$

← Only an estimate of true gradient!

- Update  $\theta$  with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$



## Advantages:

- More accurate estimation of gradient
  - Smoother convergence
  - Allows for larger learning rates
- Minibatches lead to fast training!
  - Can parallelize computation + achieve significant speed increases on GPU's

# Why is it **Stochastic** Gradient Descent?

- Initialize  $\theta$  randomly
- For N Epochs
- For each training example  $(x, y)$

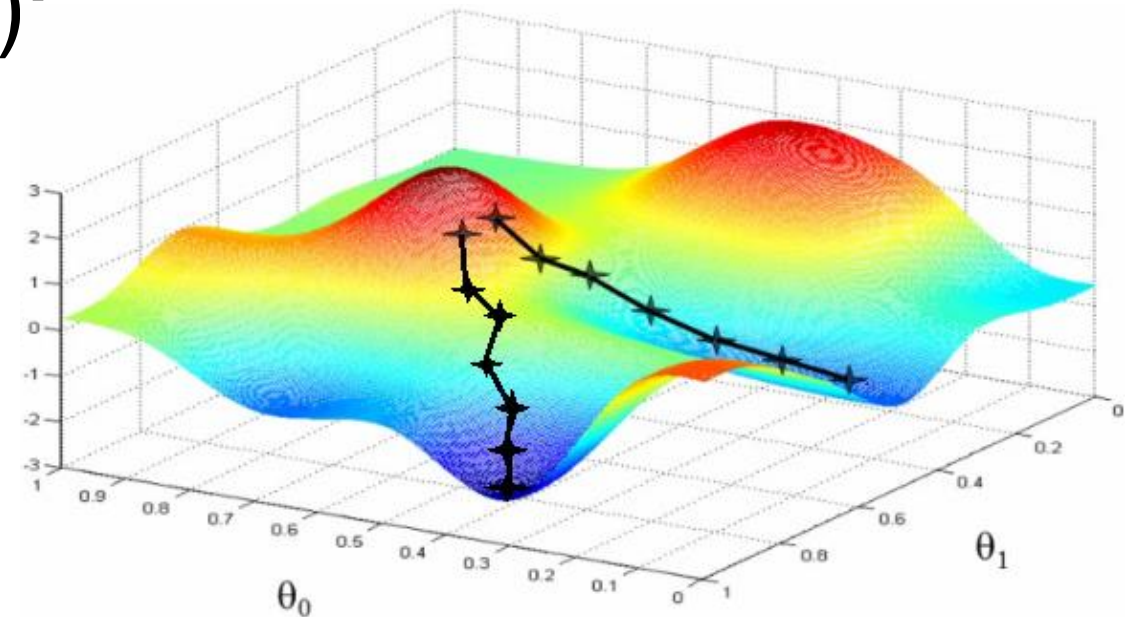
- Compute Loss Gradient:

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_i^B \frac{\partial J_i(\theta)}{\partial \theta}$$

More accurate estimate!

- Update  $\theta$  with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$



# Stochastic Gradient Descent (SGD)

- Algorithm that performs updates after each example

- initialize  $\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$
- for N iterations
  - for each training example  $(\mathbf{x}^{(t)}, y^{(t)})$  or batch

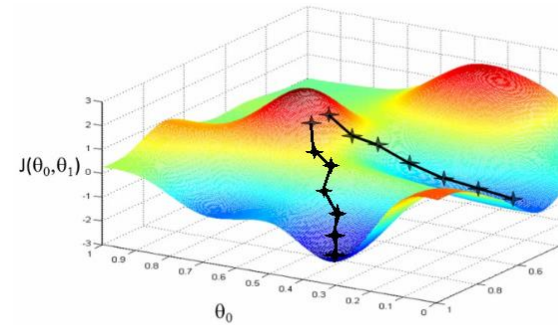
$$\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta)$$

$$\theta \leftarrow \theta + \alpha \Delta$$

Training epoch

=

Iteration over **all** examples



- To apply this algorithm to neural network training, we need:

- the loss function  $l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$
- a procedure to compute the parameter gradients:  $\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$
- the regularizer  $\Omega(\theta)$  (and the gradient  $\nabla_{\theta} \Omega(\theta)$ )

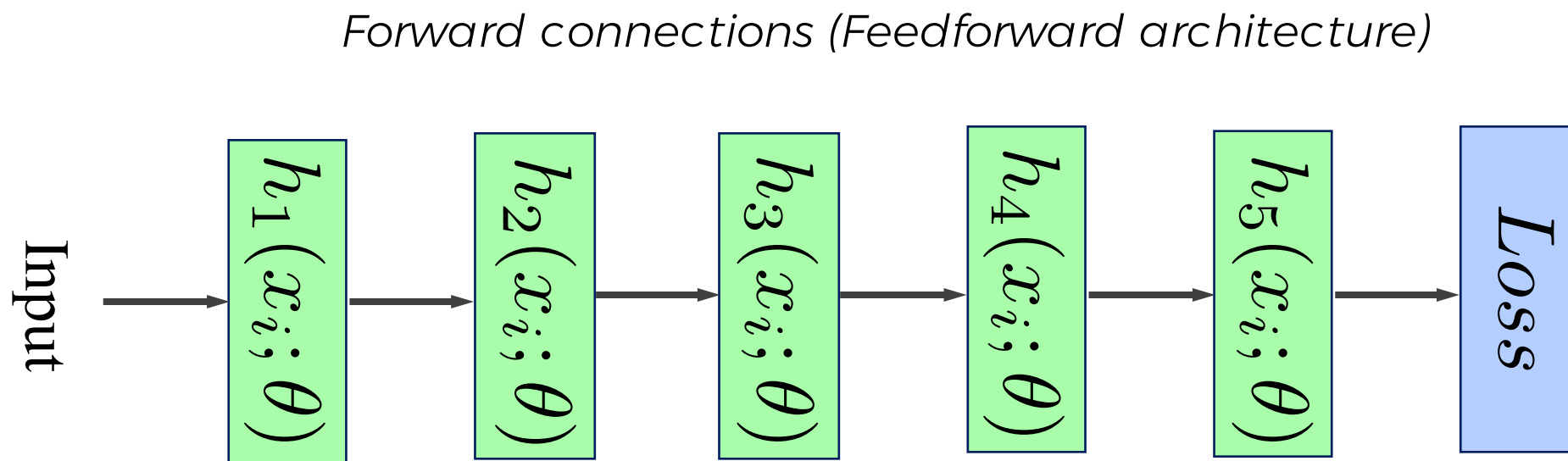
# What is a neural network again?

- A family of parametric, non-linear and hierarchical representation learning functions
- $a_L(x; \theta_1, \dots, \theta_L) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$
- $x$ : input,  $\theta_l$ : parameters for layer  $l$ ,  $a_l = h_l(x, \theta_l)$ : (non-)linear function
- Given training corpus  $\{X, Y\}$  find optimal parameters

$$\theta^* \leftarrow \arg \min_{\theta} \sum_{(x, y) \in (X, Y)} \ell(y, a_L(x; \theta_1, \dots, \theta_L))$$

# Neural network models

- A neural network model is a series of hierarchically connected functions
- The hierarchy can be very, very complex

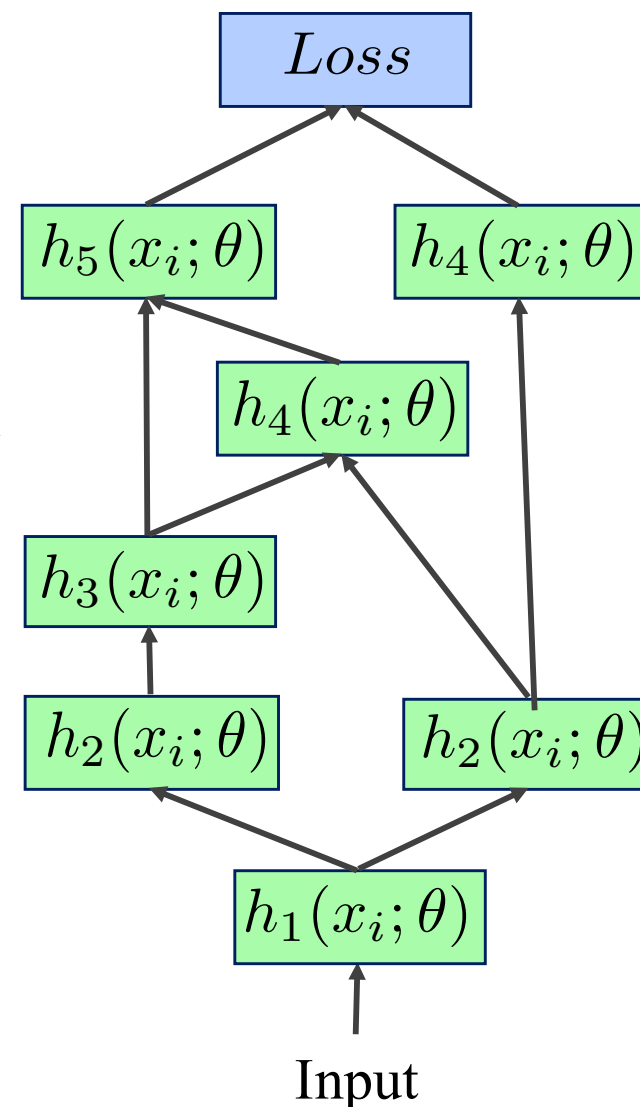




# Neural network models

- A neural network model is a series of hierarchically connected functions
- The hierarchy can be very, very complex

*Interweaved connections  
(Directed Acyclic Graph  
architecture - DAGNN)*



# Again, what is a neural network again?

- $a_L(x; \theta_1, \dots, \theta_L) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$
- $x$ : input,  $\theta_l$ : parameters for layer  $l$ ,  $a_l = h_l(x, \theta_l)$ : (non-)linear function

- Given training corpus  $\{X, Y\}$  find optimal parameters

$$\theta^* \leftarrow \arg \min_{\theta} \sum_{(x,y) \in (X,Y)} \ell(y, a_L(x; \theta_1, \dots, \theta_L))$$

- To use gradient descent optimization  $\left( \theta^{t+1} = \theta^t - \eta_t \frac{\partial \mathcal{L}}{\partial \theta^t} \right)$
- we need the gradients  $\frac{\partial \mathcal{L}}{\partial \theta_l}, l = 1, \dots, L$
- How to compute the gradients for such a complicated function enclosing other functions, like  $a_L(\dots)$ ?

# Chain rule

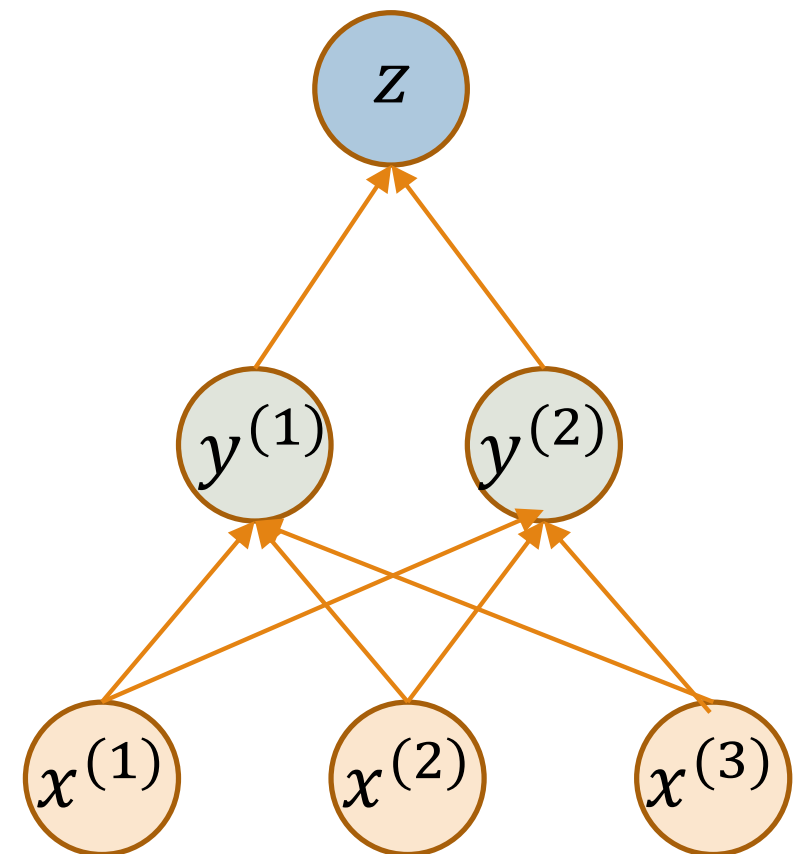
- Assume a nested function,  $z = f(y)$  and  $y = g(x)$
- Chain Rule for scalars  $x, y, z$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

- When

$$\frac{dz}{dx^i} = \sum_j \frac{dz}{dy^j} \frac{dy^j}{dx^i}$$

gradients from all possible paths



# Chain rule

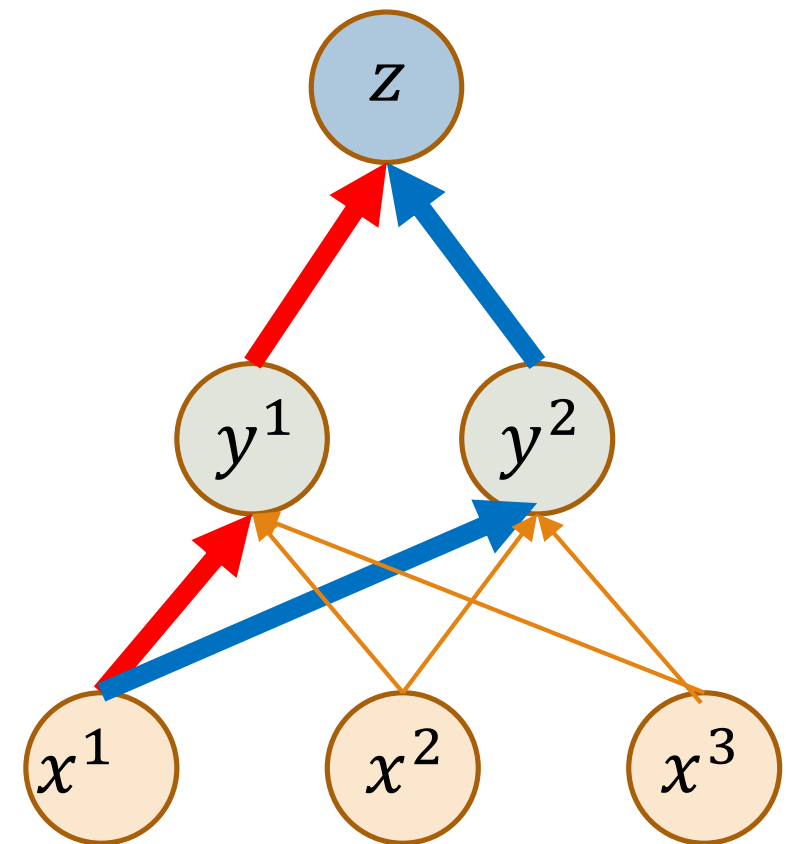
- Assume a nested function,  $z = f(y)$  and  $y = g(x)$
- Chain Rule for scalars  $x, y, z$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

- When

$$\frac{dz}{dx^i} = \sum_j \frac{dz}{dy^j} \frac{dy^j}{dx^i}$$

gradients from all possible paths



$$\frac{dz}{dx^1} = \frac{dz}{dy^1} \frac{dy^1}{dx^1} + \frac{dz}{dy^2} \frac{dy^2}{dx^1}$$

# Chain rule

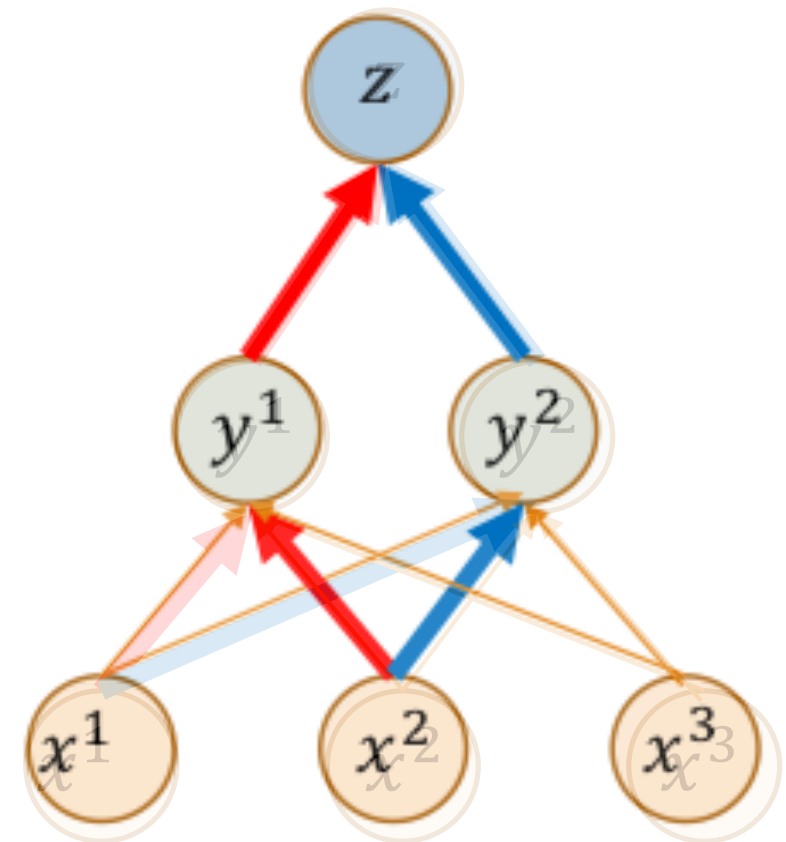
- Assume a nested function,  $z = f(y)$  and  $y = g(x)$
- Chain Rule for scalars  $x, y, z$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

- When

$$\frac{dz}{dx^i} = \sum_j \frac{dz}{dy^j} \frac{dy^j}{dx^i}$$

gradients from all possible paths



$$\frac{dz}{dx^2} = \frac{dz}{dy^1} \frac{dy^1}{dx^2} + \frac{dz}{dy^2} \frac{dy^2}{dx^2}$$



# Chain rule

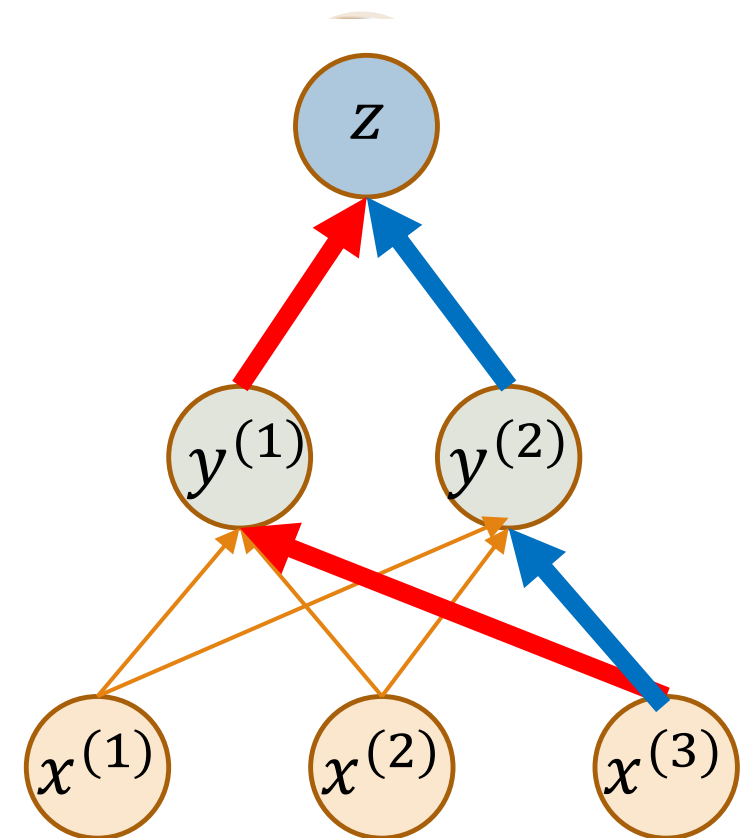
- Assume a nested function,  $z = f(y)$  and  $y = g(x)$
- Chain Rule for scalars  $x, y, z$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

- When

$$\frac{dz}{dx^i} = \sum_j \frac{dz}{dy^j} \frac{dy^j}{dx^i}$$

gradients from all possible paths



$$\frac{dz}{dx^3} = \frac{dz}{dy^1} \frac{dy^1}{dx^3} + \frac{dz}{dy^2} \frac{dy^2}{dx^3}$$

# Backpropagation $\Leftrightarrow$ Chain rule!!!

- The loss function  $\mathcal{L}(y, a_L)$  depends on  $a_L$ , which depends on  $a_{L-1}$ , ..., which depends on  $a_1$ :

$$a_L(x; \theta_{1, \dots, L}) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$

- Gradients of parameters of layer  $l \rightarrow$  Chain rule

$$\frac{\partial \mathcal{L}}{\partial \theta_l} = \frac{\partial \mathcal{L}}{\partial a_l} \cdot \frac{\partial a_l}{\partial a_{l-1}} \cdot \frac{\partial a_{l-1}}{\partial a_{l-2}} \dots$$

- When shortened, we need to two quantities

$$\frac{\partial \mathcal{L}}{\partial \theta_l} = \left( \frac{\partial a_l}{\partial \theta_l} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial a_l}$$

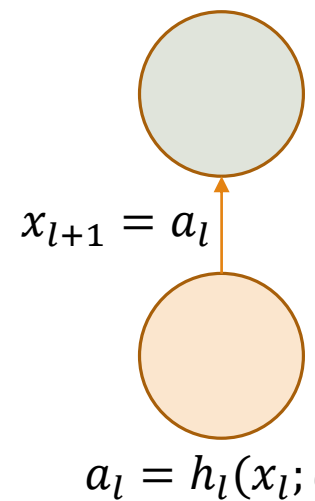
Gradient of a module w.r.t. its parameters      Gradient of loss w.r.t. the module output

# Backpropagation $\Leftrightarrow$ Chain rule!!!

- For  $\frac{\partial \mathcal{L}}{\partial a_l}$  in  $\frac{\partial \mathcal{L}}{\partial \theta_l} = \left( \frac{\partial a_l}{\partial \theta_l} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial a_l}$  we apply chain rule again

$$\frac{\partial \mathcal{L}}{\partial a_l} = \left( \frac{\partial a_{l+1}}{\partial a_l} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial a_{l+1}}$$

$$a_{l+1} = h_{l+1}(x_{l+1}; \theta_{l+1})$$



- We can rewrite  $\frac{\partial a_{l+1}}{\partial a_l}$  as gradient of module w.r.t. to input

- Remember, the output of a module is the input for the next one:

$$\frac{\partial \mathcal{L}}{\partial a_l} = \left( \frac{\partial a_{l+1}}{\partial x_{l+1}} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial a_{l+1}}$$

Recursive rule

So what is deep learning?

# Three key ideas

- (Hierarchical) Compositionality
- End-to-End Learning
- Distributed Representations

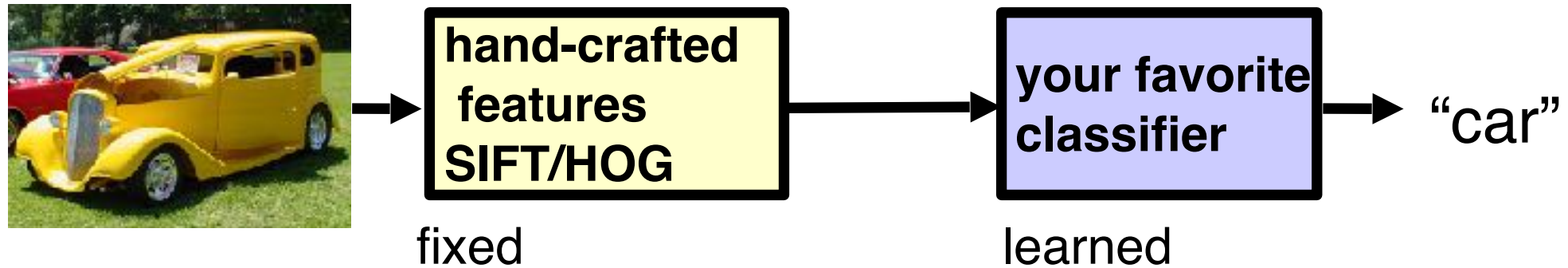


# Three key ideas

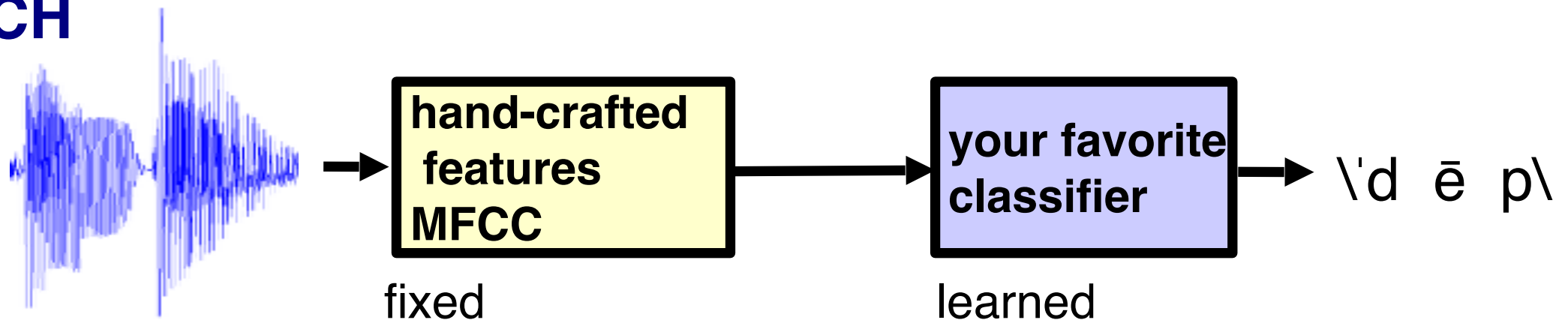
- **(Hierarchical) Compositionality**
  - Cascade of non-linear transformations
  - Multiple layers of representations
- End-to-End Learning
  - Learning (goal-driven) representations
  - Learning to feature extract
- Distributed Representations
  - No single neuron “encodes” everything
  - Groups of neurons work together

# Traditional Machine Learning

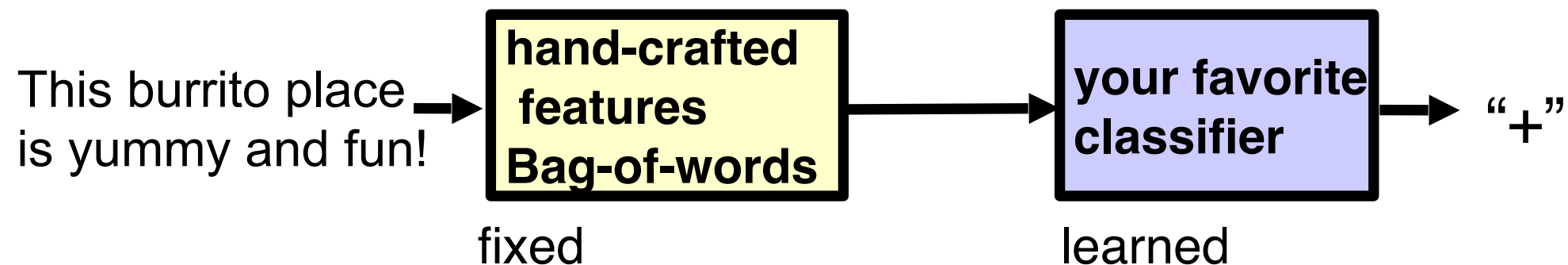
## VISION



## SPEECH

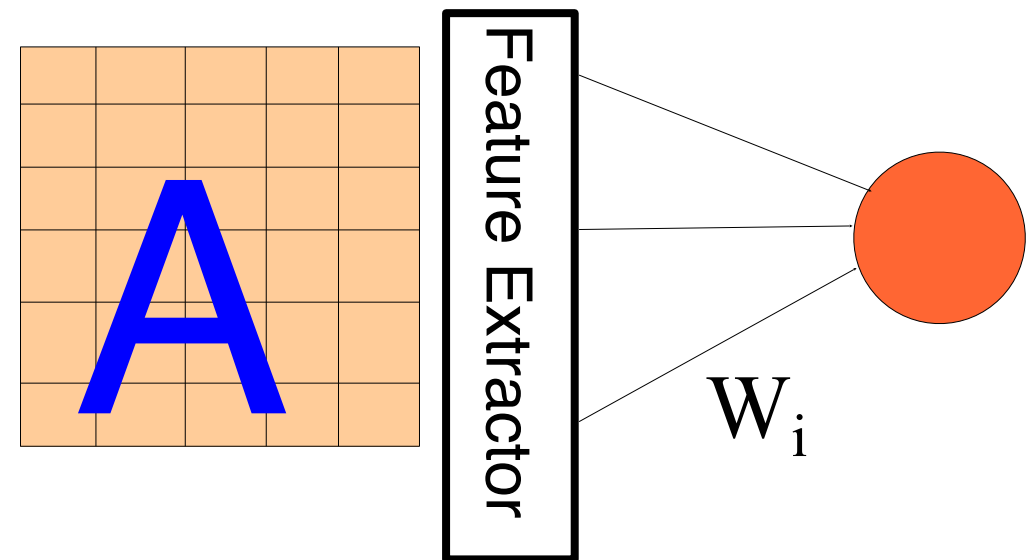


## NLP

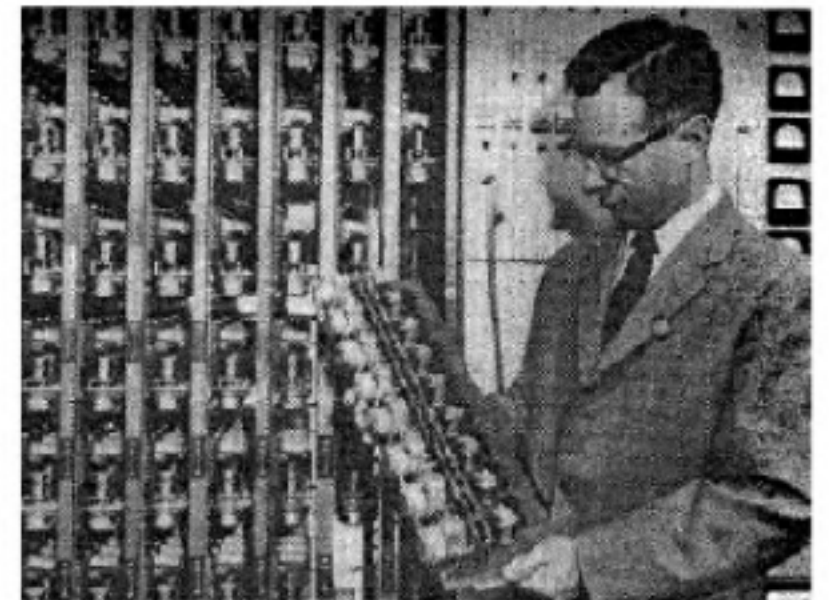
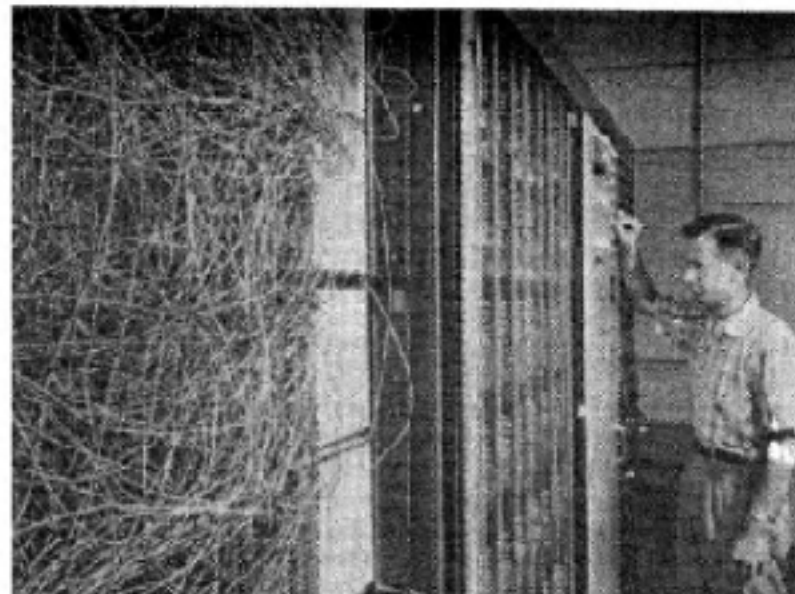
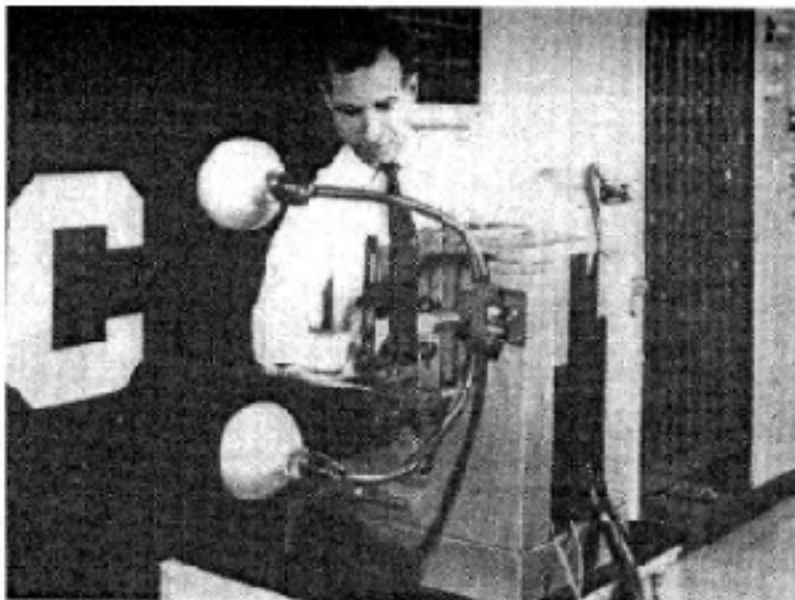


# It's an old paradigm

- The first learning machine: the **Perceptron**
  - Built at Cornell in 1960
- The Perceptron was a **linear classifier** on top of a simple **feature extractor**
- The vast majority of practical applications of ML today use glorified **linear classifiers** or glorified template matching.
- Designing a feature extractor requires considerable efforts by experts.



$$y = \text{sign} \left( \sum_{i=1}^N W_i F_i(X) + b \right)$$



# Hierarchical Compositionality

## VISION

pixels → edge → texture → motif → part → object

## SPEECH

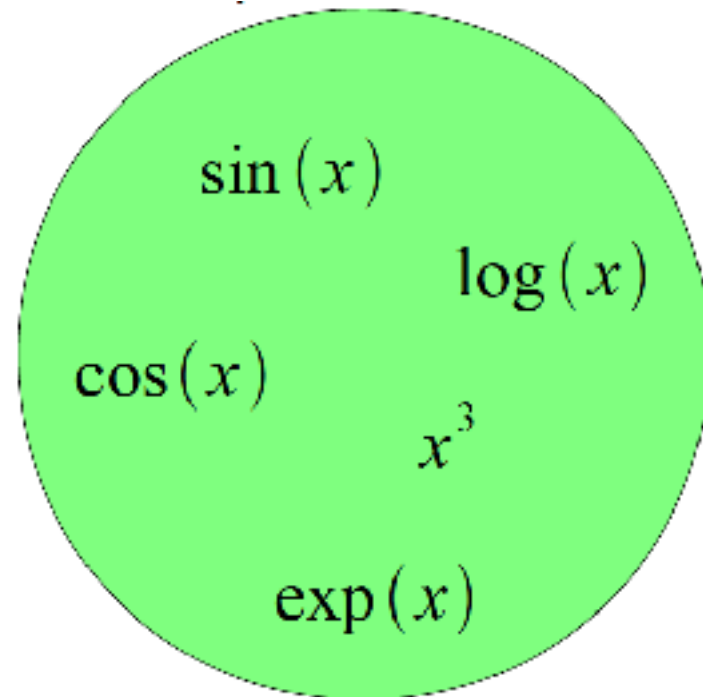
sample → spectral band → formant → motif → phone → word


## NLP

character → word → NP/VP/.. → clause → sentence → story

# Building A Complicated Function

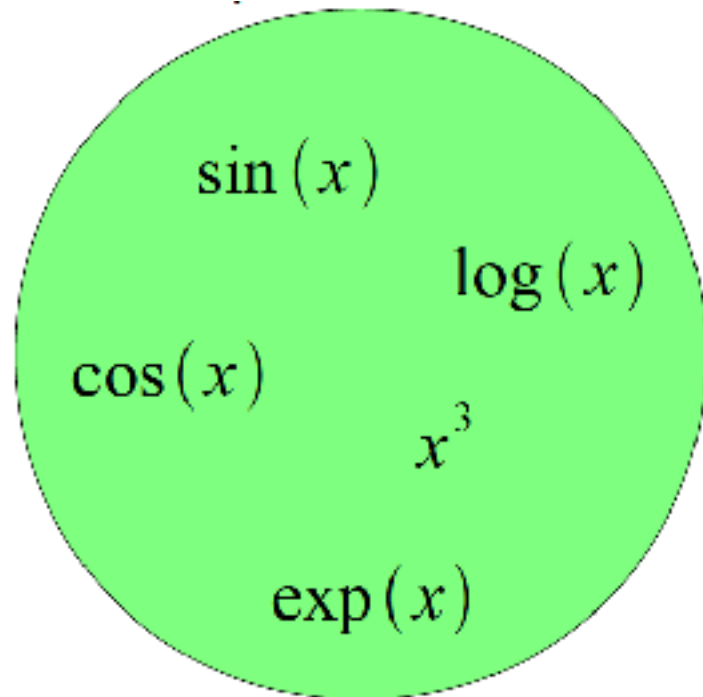
Given a library of simple functions



Compose into a  
  
complicate function

# Building A Complicated Function

Given a library of simple functions

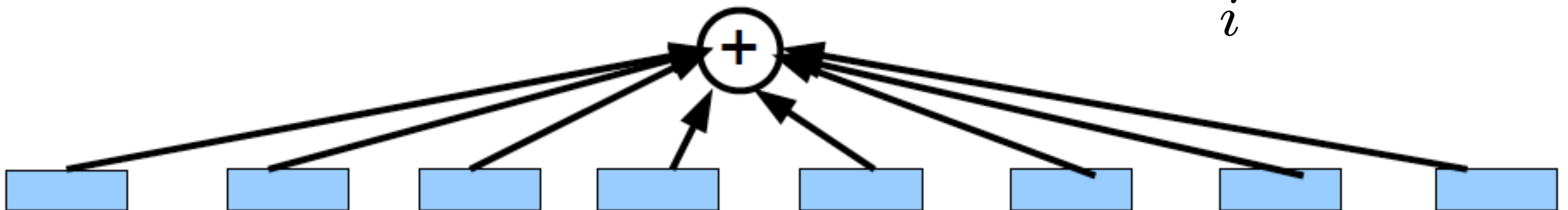


Compose into a  
→  
complicate function

## Idea 1: Linear Combinations

- Boosting
- Kernels
- ...

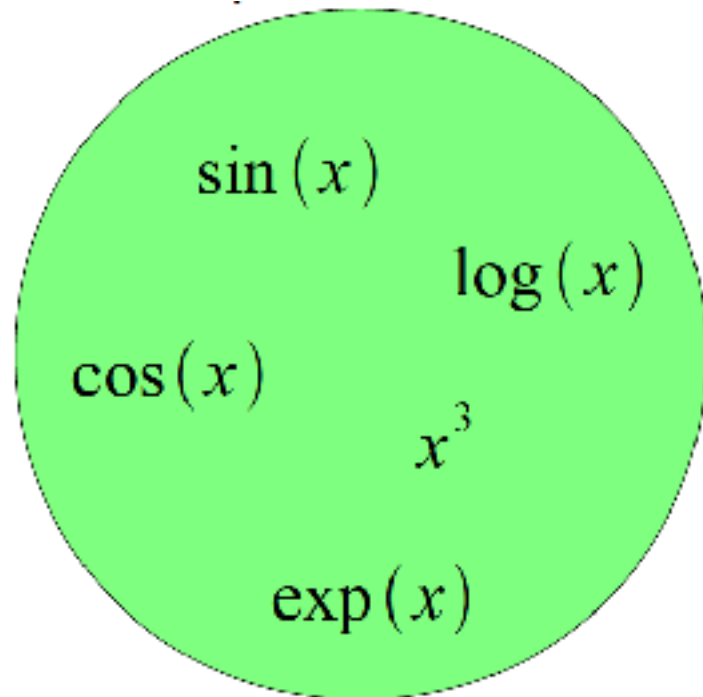
$$f(x) = \sum_i \alpha_i g_i(x)$$





# Building A Complicated Function

Given a library of simple functions

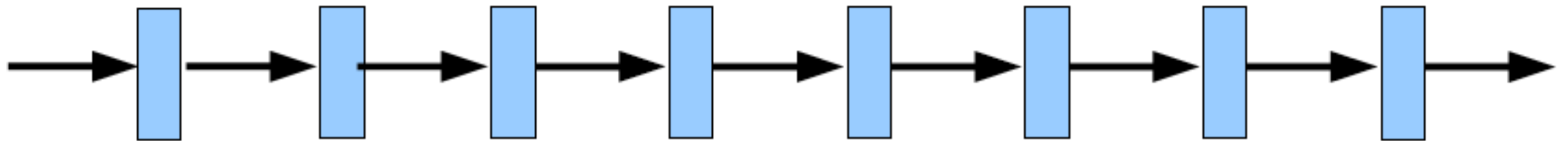


Compose into a  
→  
complicate function

## Idea 2: Compositions

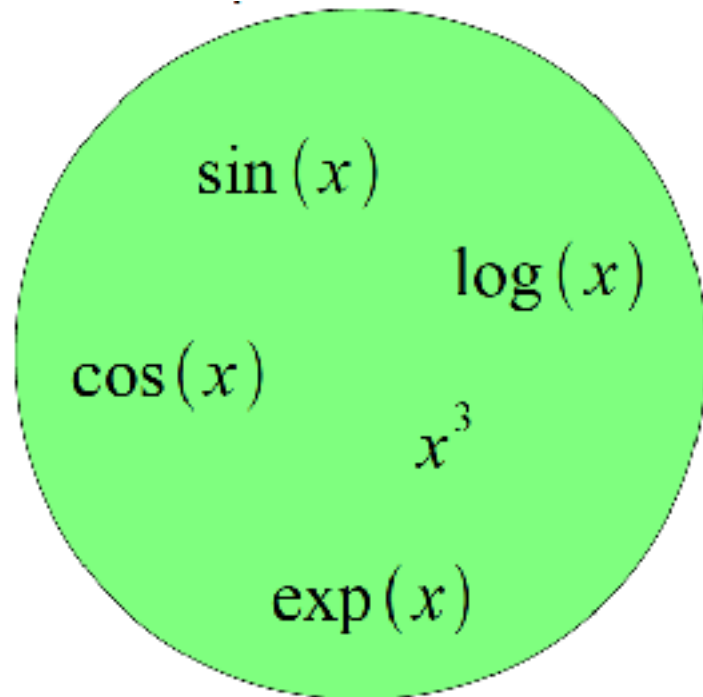
- Deep Learning
- Grammar models
- Scattering transforms...

$$f(x) = g_1(g_2(\dots(g_n(x)\dots)))$$



# Building A Complicated Function

Given a library of simple functions

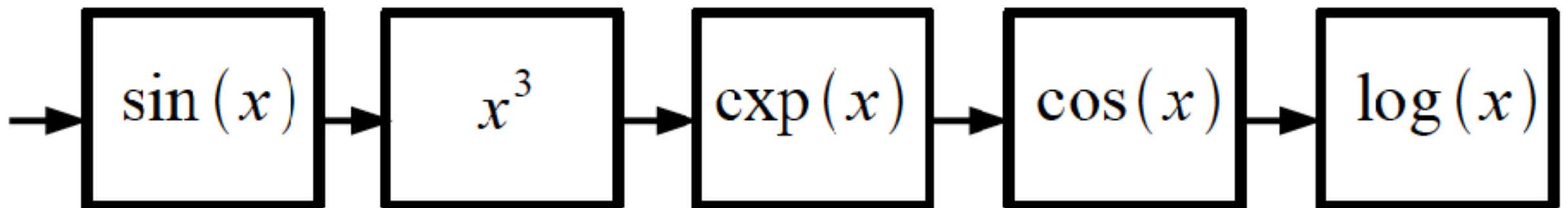


Compose into a  
→  
complicate function

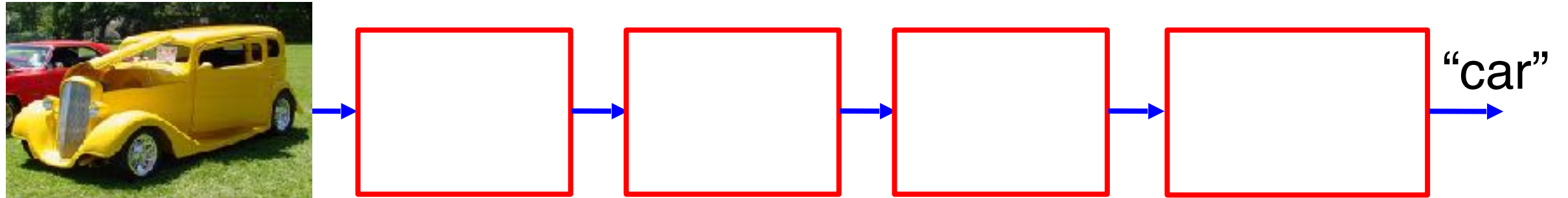
## Idea 2: Compositions

- Deep Learning
- Grammar models
- Scattering transforms...

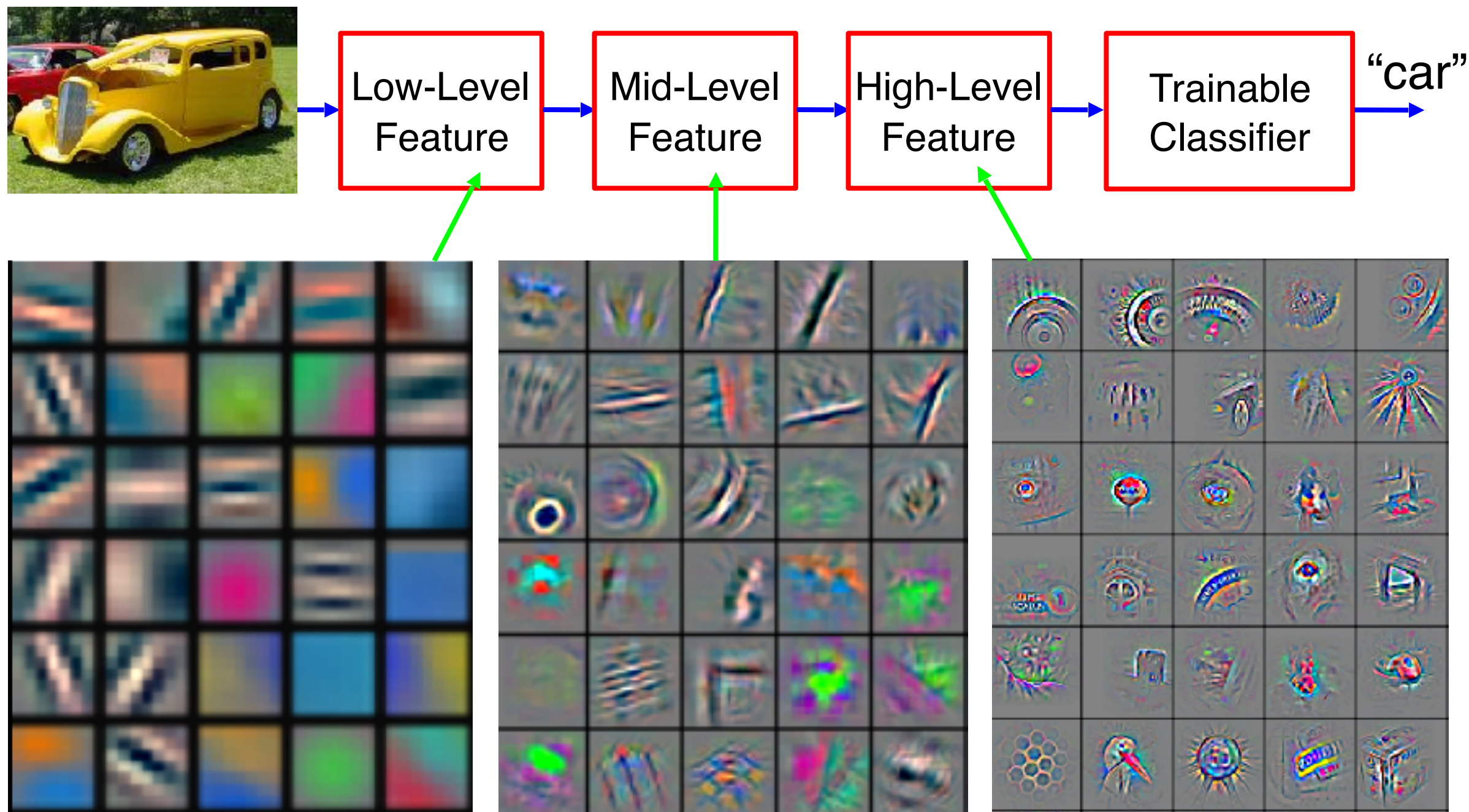
$$f(x) = \log(\cos(\exp(\sin^3(x))))$$



# Deep Learning = Hierarchical Compositionality



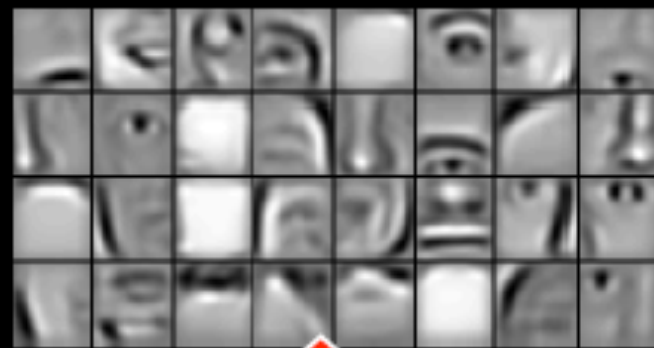
# Deep Learning = Hierarchical Compositionality



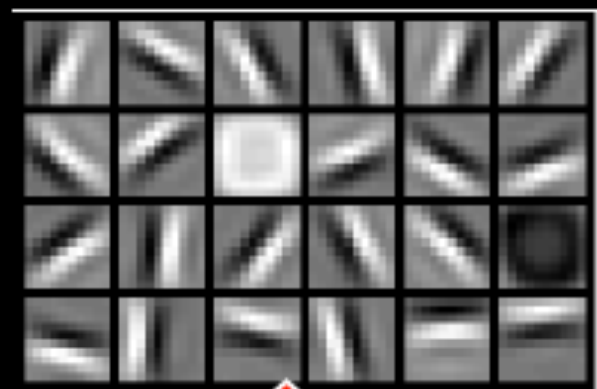
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



Face detectors



Face parts  
(combination  
of edges)



edges



pixels

Sparse DBNs  
[Lee et al. ICML '09]  
Figure courtesy: Quoc Le

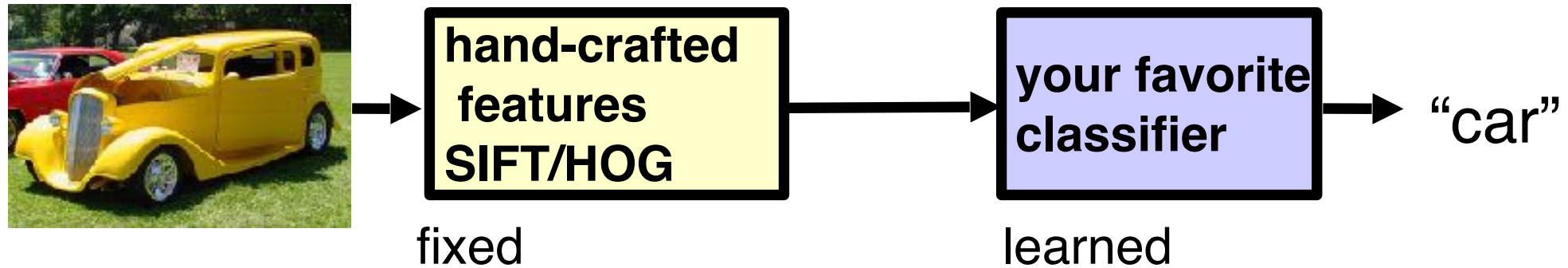
# Three key ideas

- (Hierarchical) Compositionality
  - Cascade of non-linear transformations
  - Multiple layers of representations
- **End-to-End Learning**
  - Learning (goal-driven) representations
  - Learning to feature extract
- Distributed Representations
  - No single neuron “encodes” everything
  - Groups of neurons work together

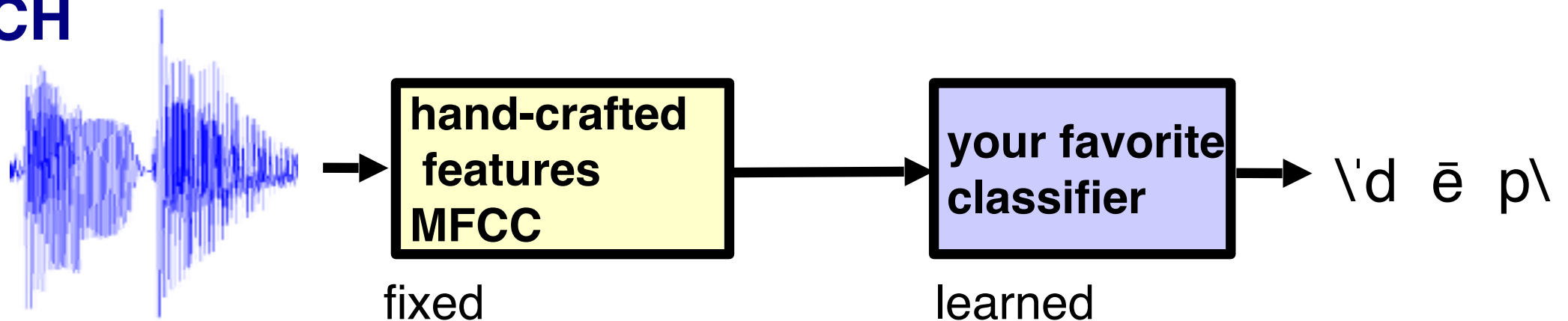


# Traditional Machine Learning

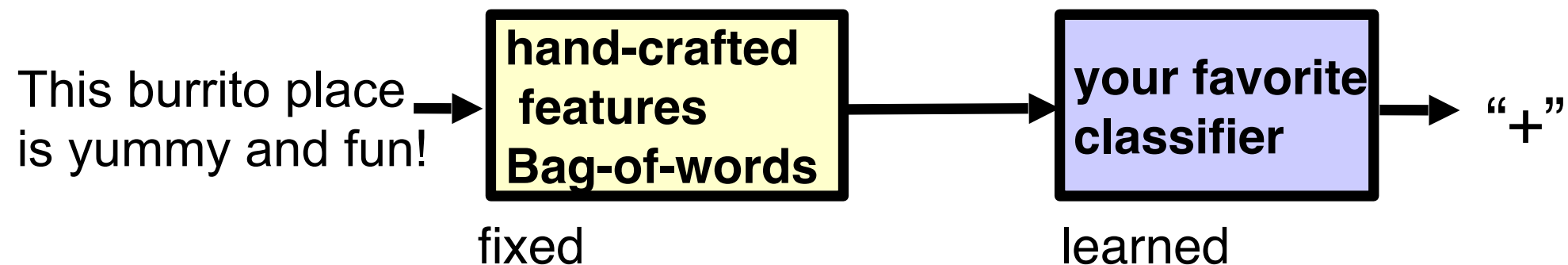
## VISION



## SPEECH

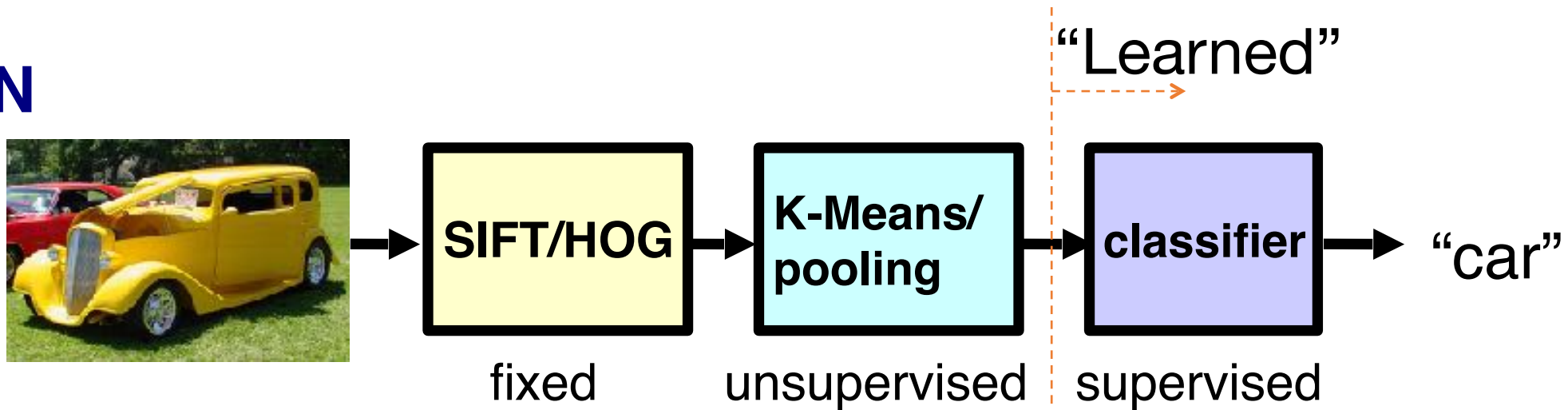


## NLP

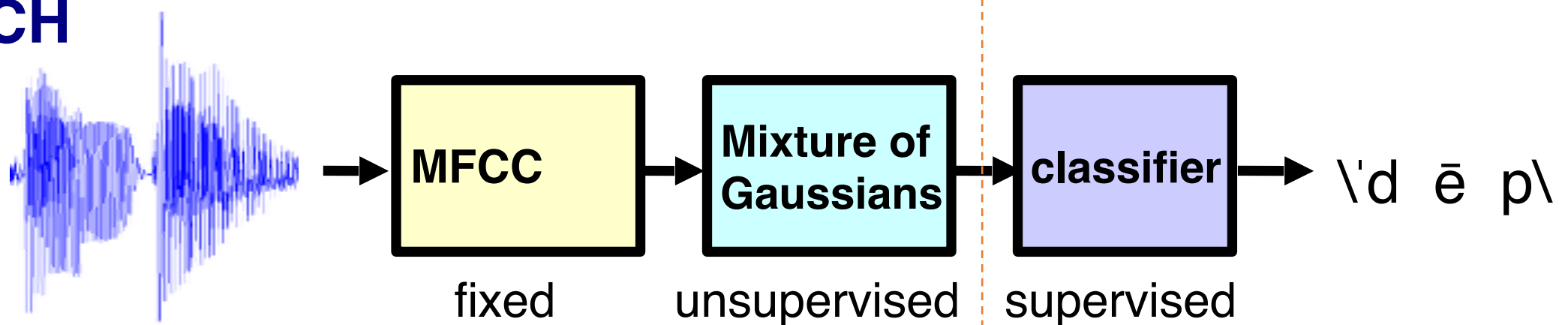


# Traditional Machine Learning (more accurately)

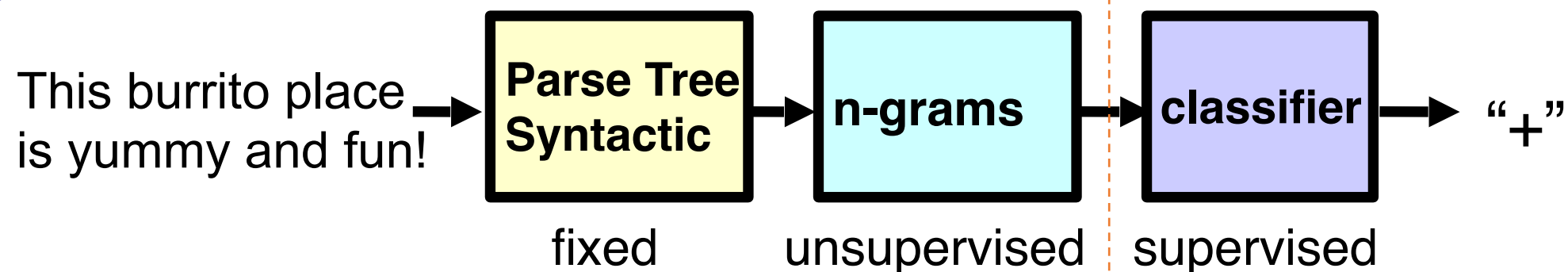
## VISION



## SPEECH

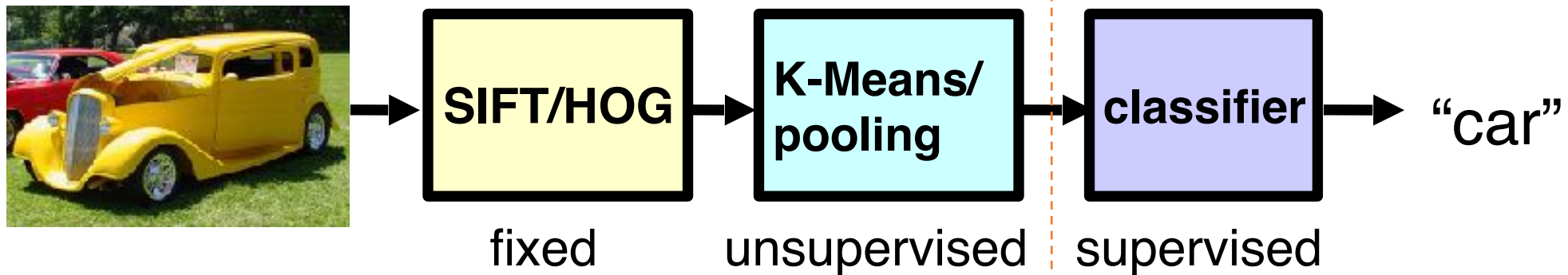


## NLP

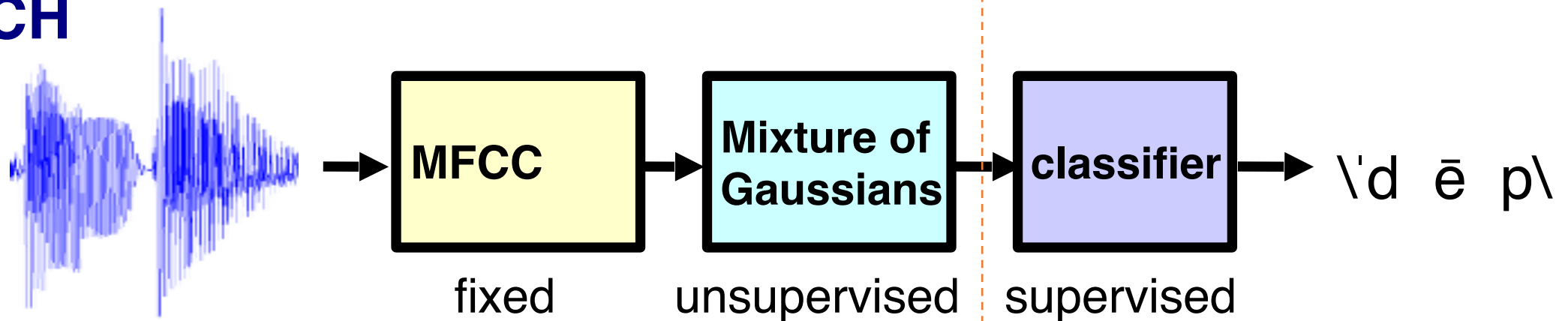


# Deep Learning = End-to-End Learning

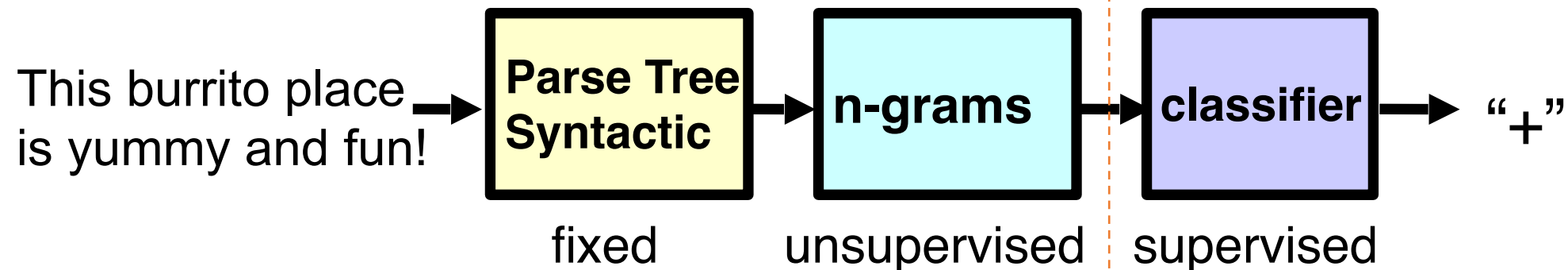
## VISION



## SPEECH

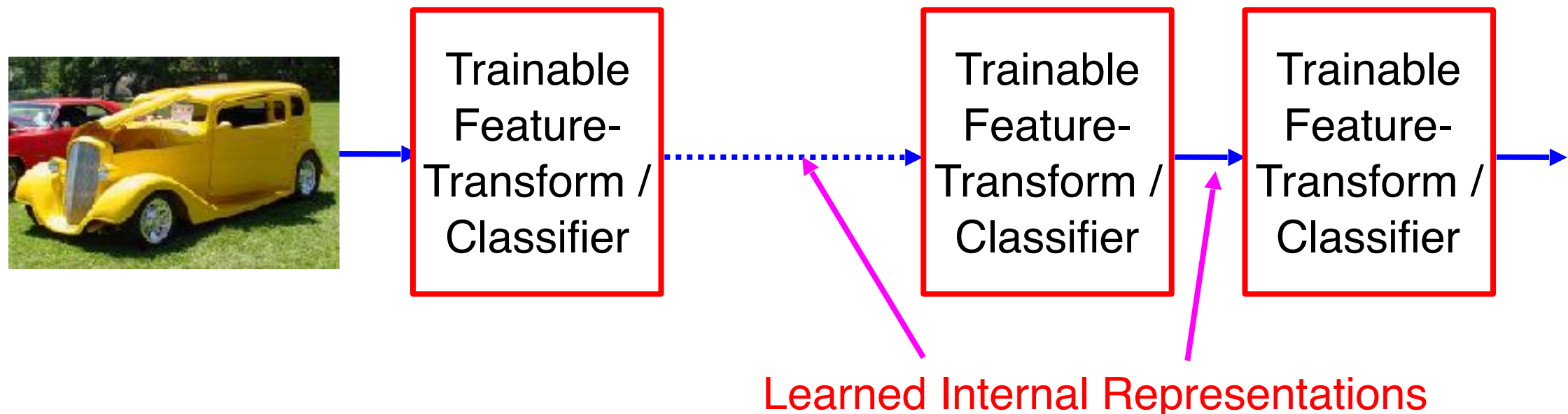


## NLP



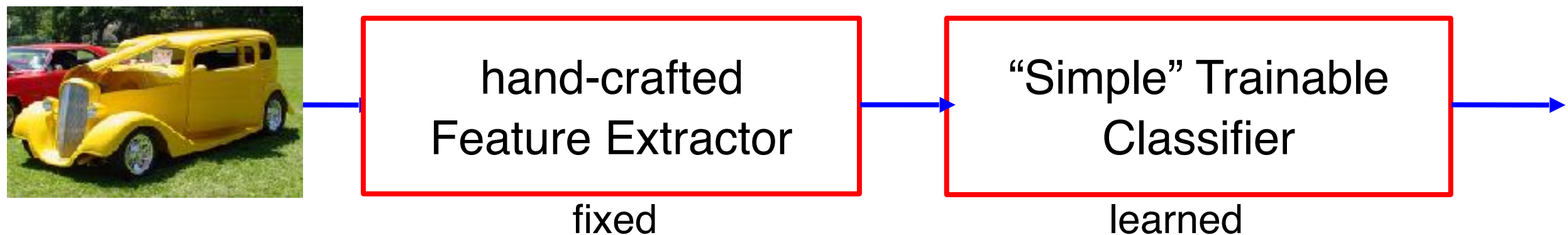
# Deep Learning = End-to-End Learning

- A hierarchy of trainable feature transforms
  - Each module transforms its input representation into a higher-level one.
  - High-level features are more global and more invariant
  - Low-level features are shared among categories

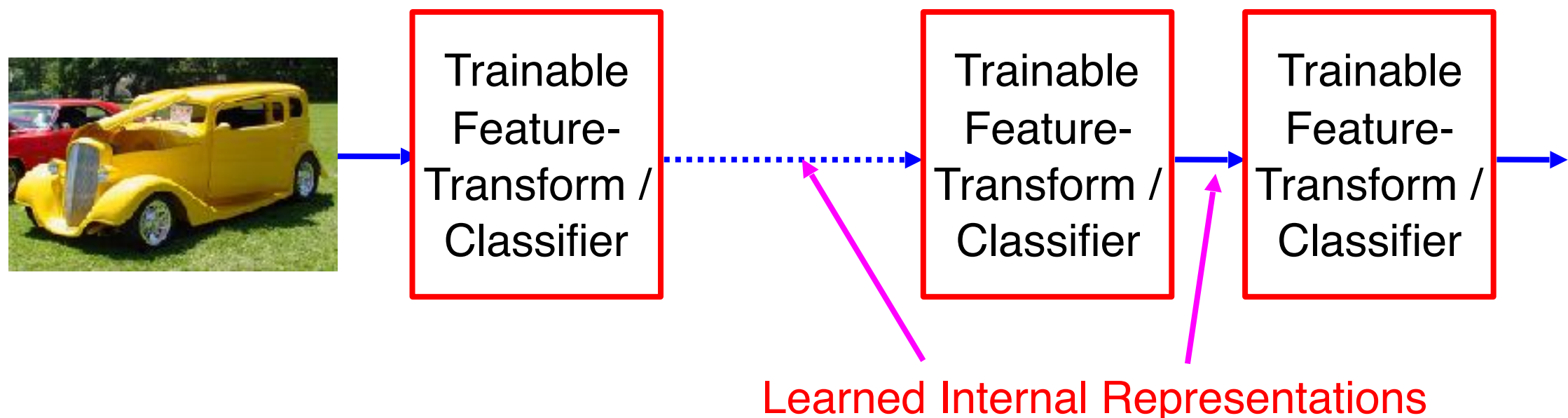


# “Shallow” vs Deep Learning

- “Shallow” models



- Deep models



# Three key ideas

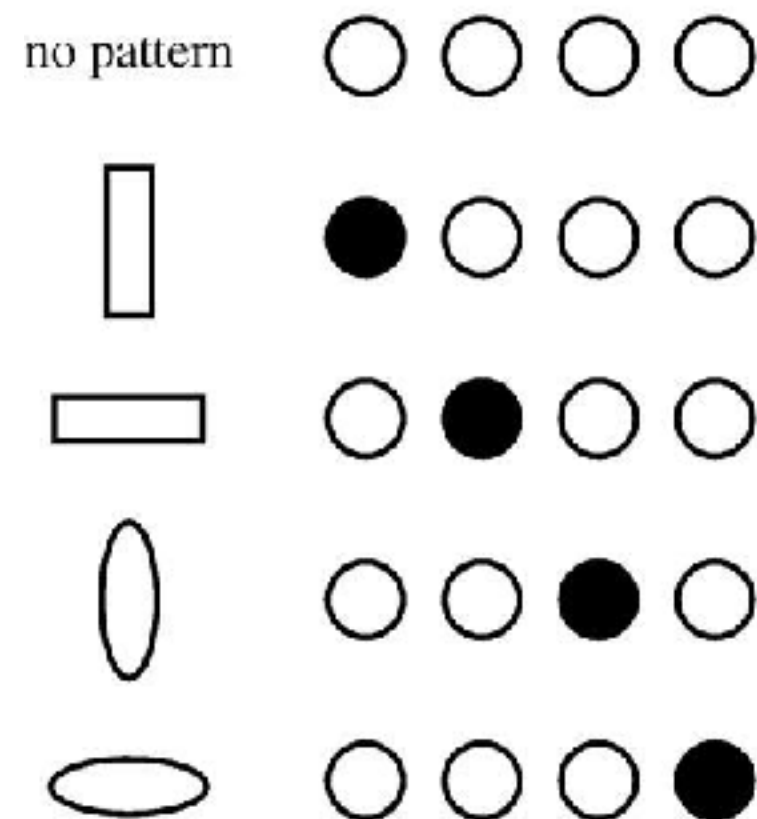
- (Hierarchical) Compositionality
  - Cascade of non-linear transformations
  - Multiple layers of representations
- End-to-End Learning
  - Learning (goal-driven) representations
  - Learning to feature extract
- **Distributed Representations**
  - No single neuron “encodes” everything
  - Groups of neurons work together



# Localist representations

- The simplest way to represent things with neural networks is to **dedicate one neuron to each thing**. (a)

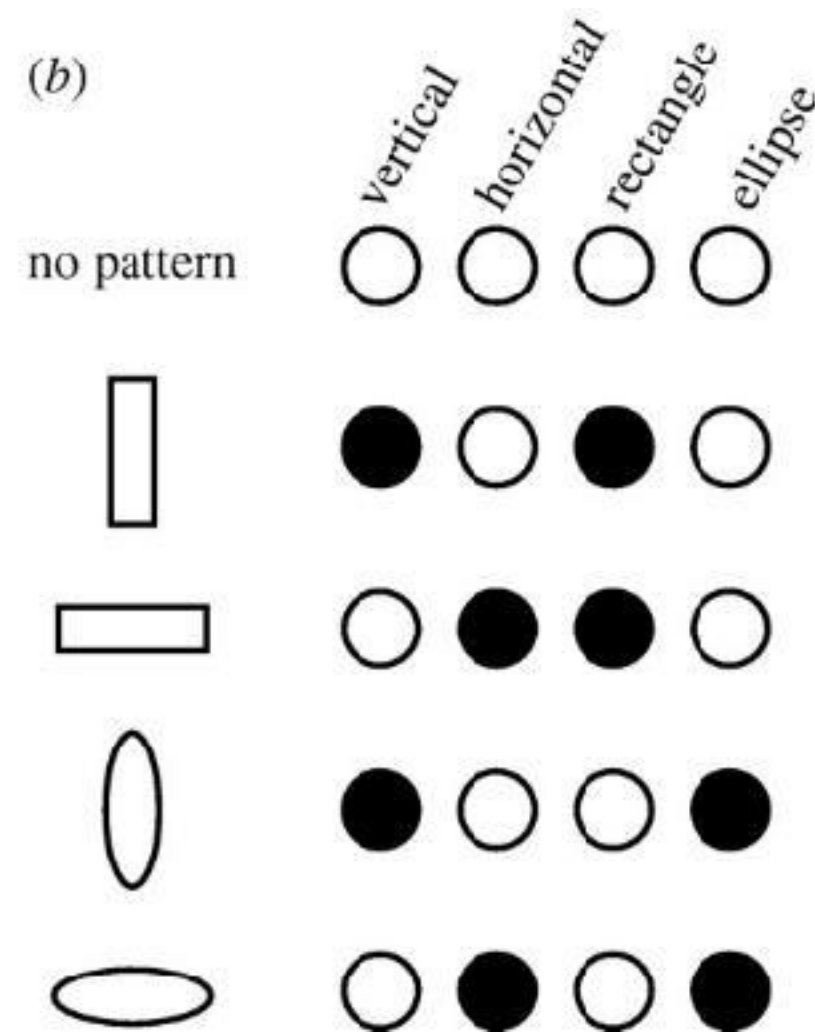
- Easy to understand.
- Easy to code by hand
  - Often used to represent inputs to a net
- Easy to learn
  - This is what mixture models do.
  - Each cluster corresponds to one neuron
- Easy to associate with other representations or responses.



- But localist models are very inefficient whenever the data has componential structure.

# Distributed Representations

- Each neuron must represent something, so this must be a local representation.
- **Distributed representation** means a many-to-many relationship between two types of representation (such as concepts and neurons).
  - Each concept is represented by many neurons
  - Each neuron participates in the representation of many concepts



Local ● ● ○ ● = VR + HR + HE = ?

Distributed ● ● ○ ● = V + H + E ≈ ○

# Power of distributed representations!

## Scene Classification

bedroom

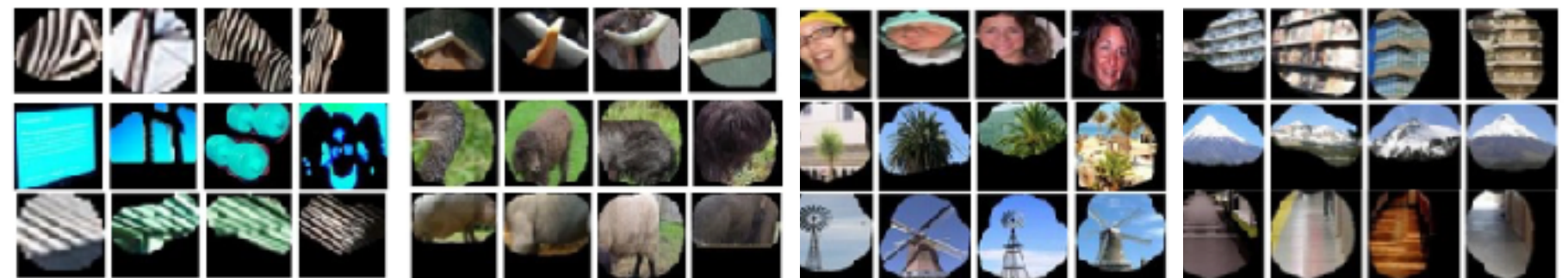


mountain



- Possible internal representations:

- Objects
- Scene attributes
- Object parts
- Textures



Simple elements & colors

Object part

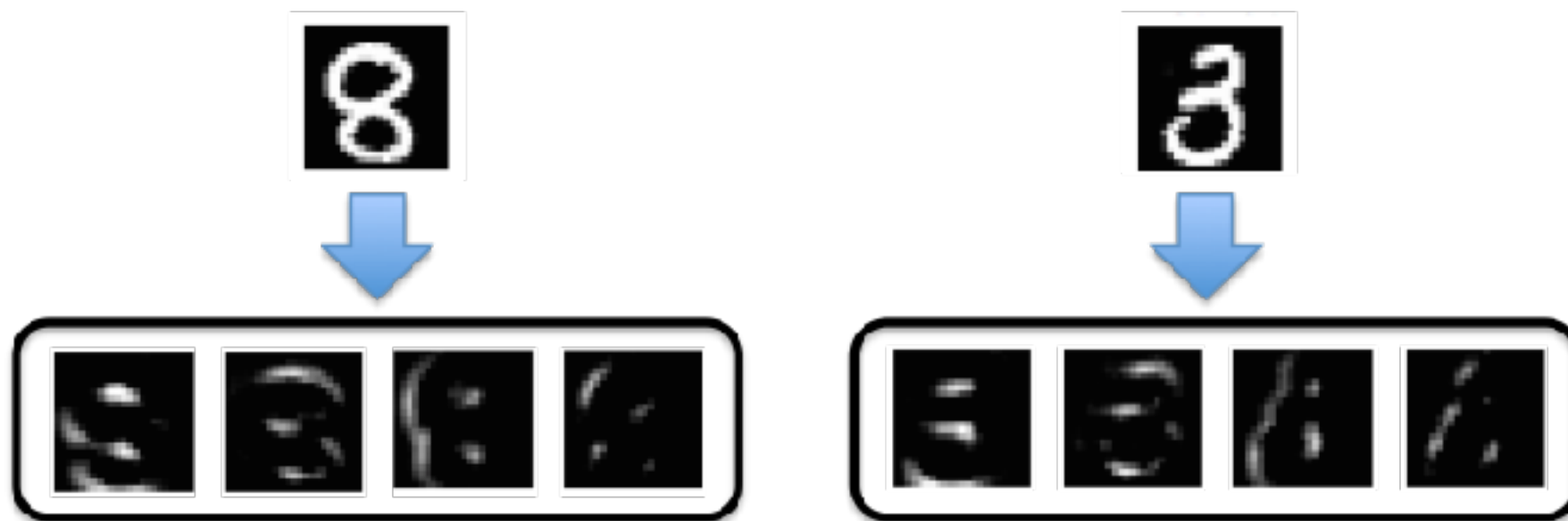
Object

Scene

# Deep Convolutional Neural Networks

# Convolutions

- Images typically have invariant patterns
  - E.g., directional gradients are translational invariant:



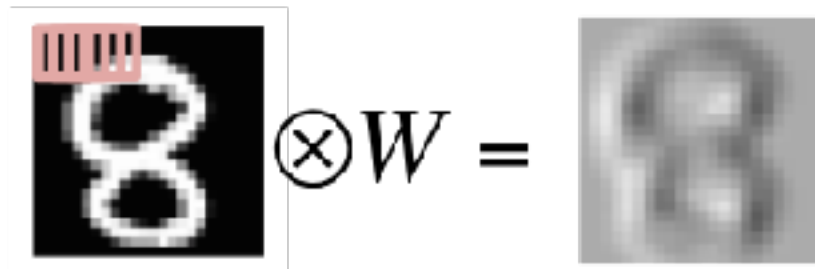
- Apply convolution to local sliding windows

# Convolution Filters

- Applies to an image patch  $x$ 
  - Converts local window into single value
  - Slide across image

$$x \otimes W = \sum_{ij} W_{ij} x_{ij}$$

↑  
Local Image Patch



Left-to-Right  
Edge Detector

-1	0	+1
-1	0	+1
-1	0	+1


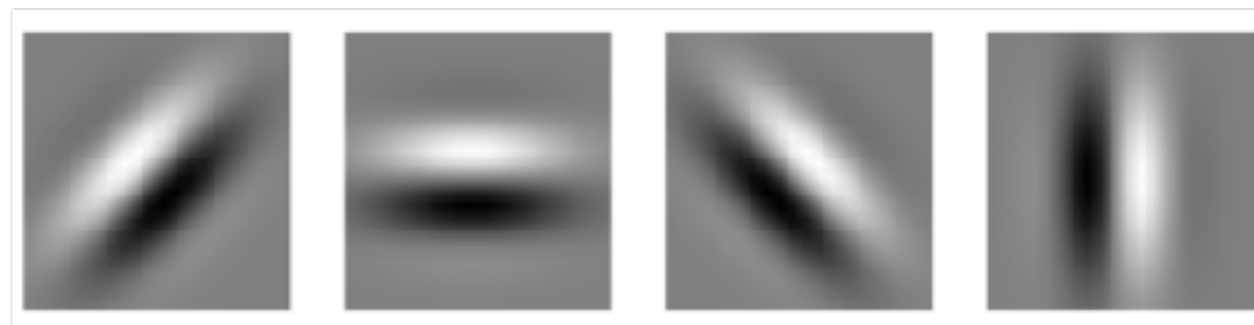
$W$



# Gabor Filters

- Most common low-level convolutions for computer vision

Example  $W$ :



-1	0	+1
-1	0	+1
-1	0	+1

$W$

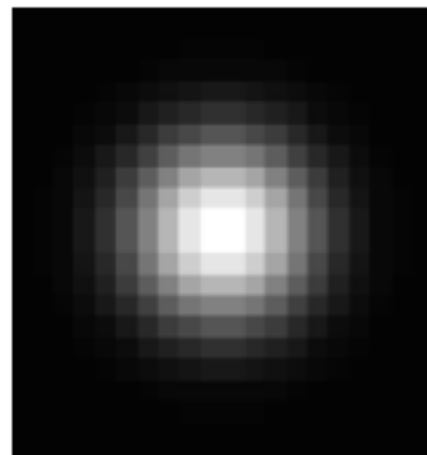
[http://en.wikipedia.org/wiki/Gabor\\_filter](http://en.wikipedia.org/wiki/Gabor_filter)

# Gaussian Blur Filters

- Weights decay according to Gaussian Distribution
  - Variance term controls radius

Example W:

Apply per RGB Channel



- Black = 0
- White = Positive



[http://en.wikipedia.org/wiki/Gaussian\\_blur](http://en.wikipedia.org/wiki/Gaussian_blur)