

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



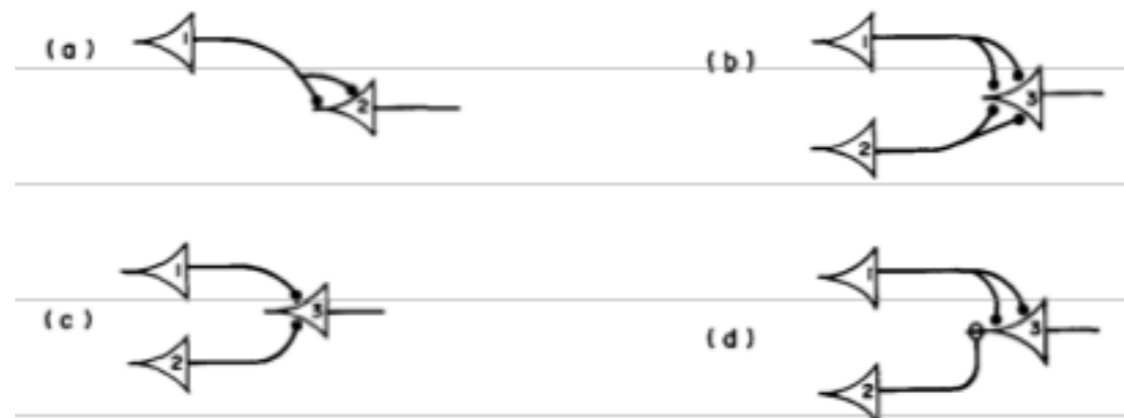
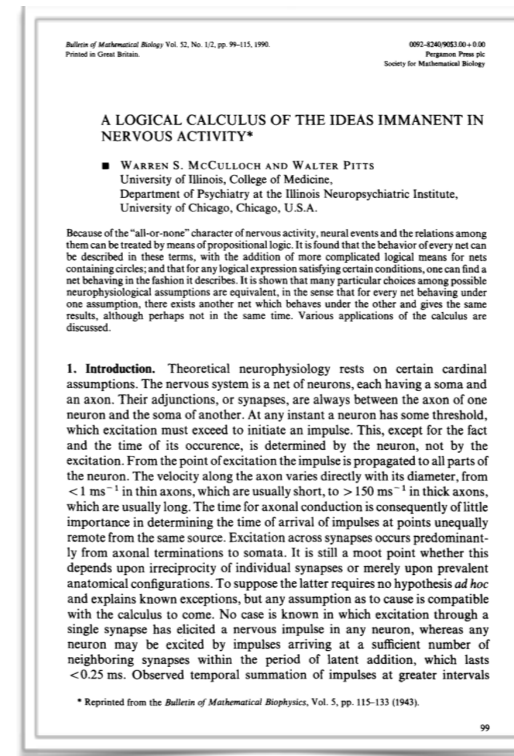
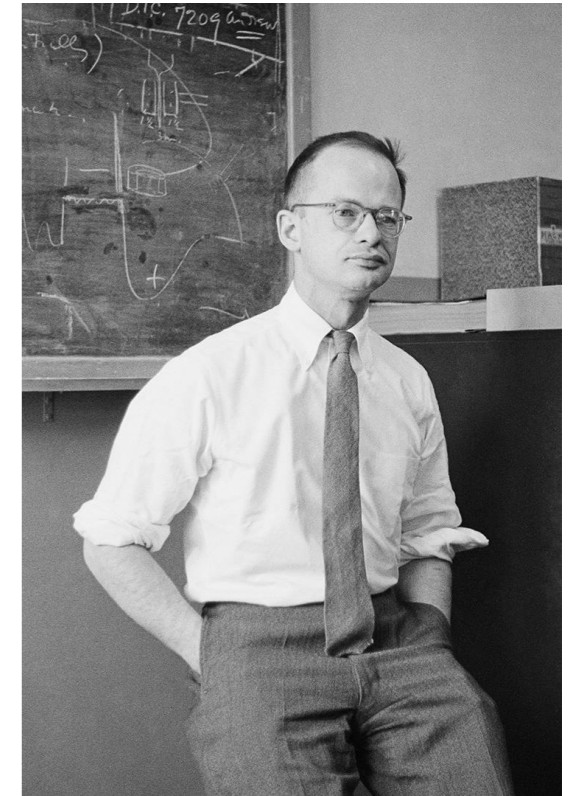
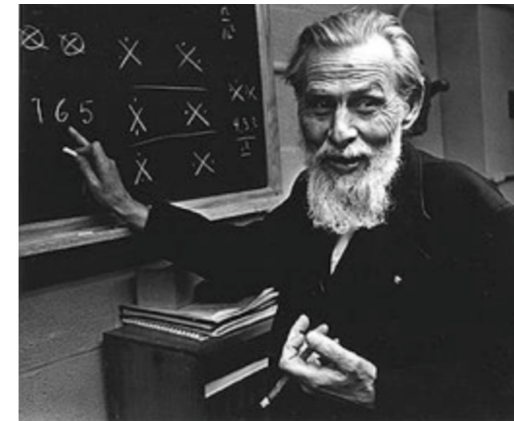
“Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction.”

– Yann LeCun, Yoshua Bengio and Geoff Hinton 2

1943 – 2006:
A Prehistory of Deep Learning

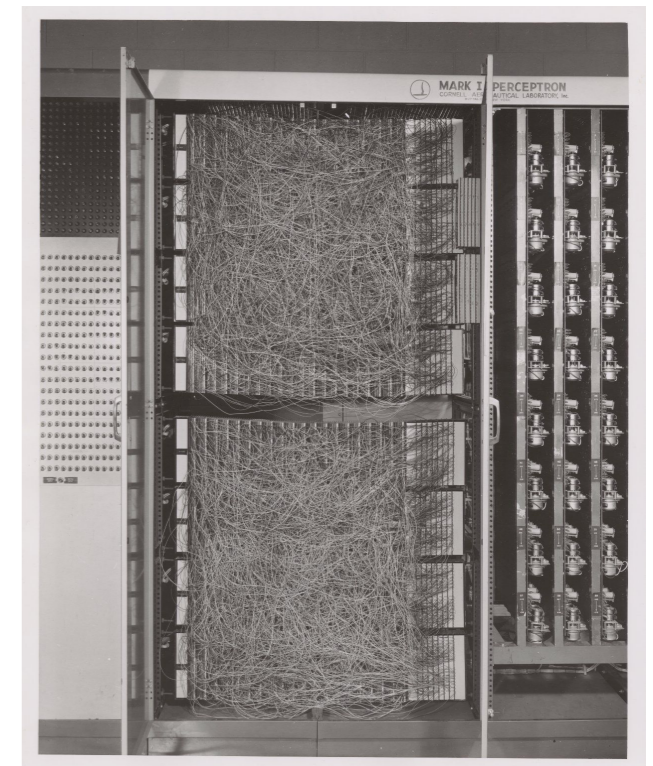
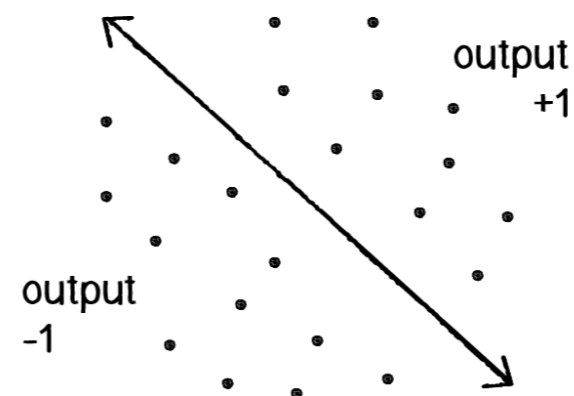
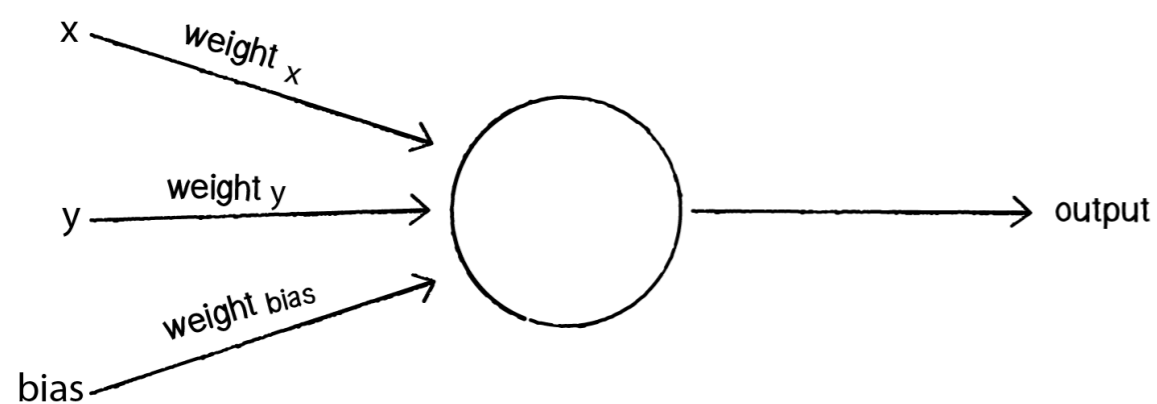
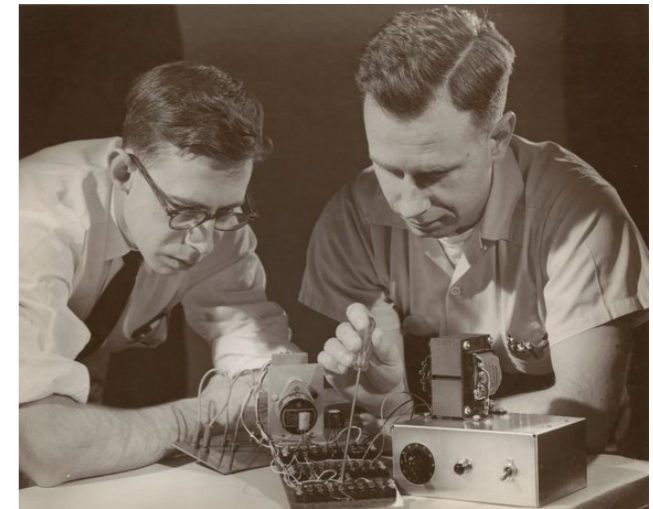
1943: Warren McCulloch and Walter Pitts

- 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

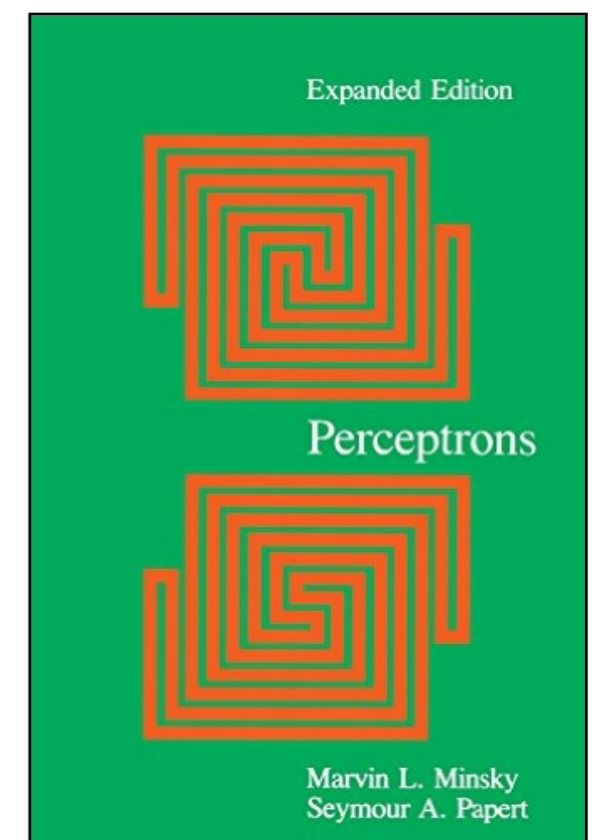
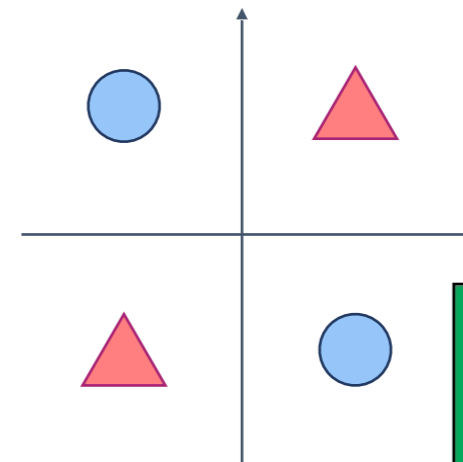


1969: Marvin Minsky and Seymour Papert

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

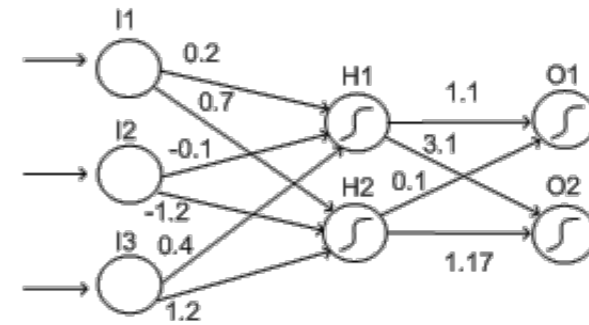


- Perceptrons can only represent linearly separable functions.
 - such as **XOR** Problem
- Wrongly attributed as the reason behind the **AI winter**, a period of reduced funding and interest in AI research



1990s

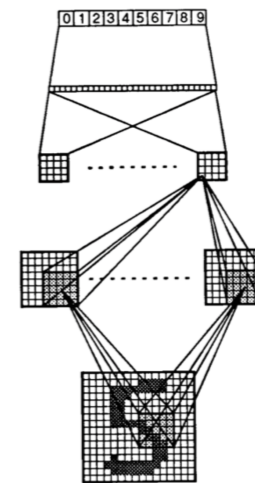
- **Multi-layer perceptrons** can theoretically learn any function (Cybenko, 1989; Hornik, 1991)



- Training multi-layer perceptrons

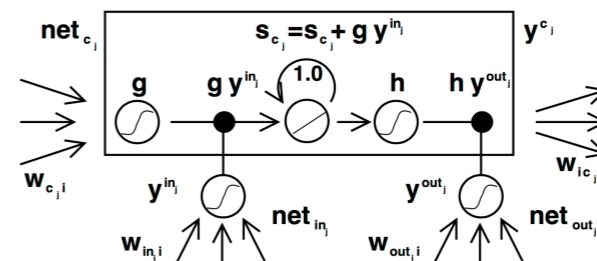
- **Back-propagation** (Rumelhart, Hinton, Williams, 1986)
- **Back-propagation through time (BPTT)** (Werbos, 1988)

10 output units
30 units
12 feature detectors (4 by 4)
12 feature detectors (8 by 8)
16 by 16 input



- New neural architectures

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



Backpropagation Through Time: What It Does and How to Do It
PAUL J. WERBOS

Backpropagation is now the most widely used tool in the field of artificial neural networks. At the same time, backpropagation is a method for calculating derivatives exactly and efficiently in any large system made up of elementary subsystems or calculations which are represented by known, differentiable functions, maps, and operations.

This paper reviews backpropagation, a simple method which is now being widely used in areas like pattern recognition and neural diagnosis. First, it generalizes backpropagation to work with propagation through time, and discusses applications to areas like pattern recognition involving dynamic systems, systems identification, and control. Finally, it describes further extensions of this method, in that with systems other than neural networks, systems and other practical issues which arise with the method. Pseudocode is provided to clarify the algorithm. The clear role of neural networks—the theorem which underlies backpropagation—is briefly discussed.

1. INTRODUCTION
Backpropagation through time is a very powerful tool, with applications to pattern recognition, dynamic modeling, sensitivity analysis, and the control of systems over time, among others. It can be applied to neural networks, to econometric models, to fuzzy logic structures, to fluid dynamics models, and to almost any system built up from elementary subsystems or calculations. The one serious constraint is that the elementary subsystems must be represented by functions known to the user, functions which are both continuous and differentiable (i.e., possess derivatives). For example, the first practical application of backpropagation was for estimating a dynamic model to predict national and social communications in 1987 [1].

Historically, the most general formulation of backpropagation can only be used by those who are willing to work out the mathematics of their particular application. This paper will mainly describe a simpler version of backpropagation, which can be translated into computer code and applied directly by neural network users.

Section II will review the variables and most widely used form of backpropagation, which may be called "basic backpropagation." The concepts here will already be familiar to those who have read the paper by Rumelhart, Hinton, and Williams [2] in the seminal book *Parallel Distributed Processing*, which played a pivotal role in the development of the field. (That book also acknowledged the prior work of Parker [3] and Le Cun [4], and the pivotal role of Charles Smith of the Systems Development Foundation.) This section will use the notation to describe backpropagation in a rigorous manner. (The need for new notation may seem unnecessary to some, but for those who have to apply backpropagation to complex systems, it is essential.)

Section III will use the same notation to describe backpropagation through time. Backpropagation through time has been applied to concrete problems by a number of authors, including, at least, Widrow and Shafer [5], Sussner and Waibel et al. [6], Nguyen and Widrow [7], Jordan [8], Kawato [9], Elman and Zipser, Narendra [10], and myself [11, 12, 13, 14]. Section IV will discuss what is missing in this simplified discussion, and how to do better.

As its core, backpropagation is simply an efficient and exact method for calculating all the derivatives of a single target quantity (such as pattern classification error) with respect to a large set of input quantities (such as the parameters or weights in a classification rule). Backpropagation through time extends this method so that it applies to dynamic systems. This allows one to calculate the derivatives needed when optimizing an iterative analysis procedure, a neural network with memory, or a control system which maximizes performance over time.

II. BASIC BACKPROPAGATION
A. The Supervised Learning Problem
Basic backpropagation is carried the most popular method for training neural networks for supervised learning tasks, which is symbolized in Fig. 1.

In supervised learning, we try to adapt an artificial neural network so that its actual outputs (Y) come close to some target outputs (T) for a training set which contains T patterns. The goal is to adapt the parameters of the network so that it performs well on patterns from outside the training set.

The main use of supervised learning today lies in pattern recognition.

U.S. Government work not protected by U.S. copyright.

Learning representations by back-propagating errors
David E. Rumelhart, Geoffrey E. Hinton, & Ronald J. Williams*

*Institute for Cognitive Sciences, C-111, University of California, San Diego, La Jolla, California 92037, USA.
†Department of Computer Science, Carnegie-Mellon University, Pittsburgh, Pennsylvania 15213, USA.

We describe a new learning procedure, back-propagation, for networks of neurons-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the network and the desired output vector. As a result of the weight adjustments, internal hidden units in the network come to represent features of the task domain, and the output units represent important features of the task domain. This is the first practical application of backpropagation to networks with hidden units. The ability to create useful internal features distinguishes back-propagation from earlier methods such as the perceptron-convergence procedure.

There have been many attempts to design self-organizing networks, but the most successful designs are those in which the modification rule will allow an arbitrarily connected neural network to develop an internal structure that is appropriate for a particular task domain. The task is specified by giving the desired state vector of the output units for each representative of the input units. If the input units are directly connected to the output units, it is relatively easy to find learning rules that modify the weights of the connections so as to progressively reduce the difference between the actual and desired output vectors. Learning becomes more interesting but more difficult when we introduce hidden units whose actual or desired states are not specified by the task. In this procedure, the hidden units are not true hidden units because their output connections are fixed by hand, or their states are randomly determined by the input vector; they do not learn representations. The learning procedure must discover what representations the hidden units should be active in order to help achieve the desired input-output behavior. This discovery is done by back-propagating error signals from the output units through the hidden units. The procedure is powerful enough to construct representations that are more abstract and more useful than those that could be achieved by other methods.

The simplest form of the learning procedure is for layered networks which have a layer of input units at the bottom, one or more layers of intermediate layers, and a layer of output units at the top. Connections within a layer or from layer to layer are all initialized to zero. Connections on each intermediate layer are set to random values. The state of each input unit is presented to the network by setting the state of the input units. Then the state of the units in each layer are determined by applying equations (1) and (2) to the connections coming from lower layers. All units within a layer have their states set in parallel, but different layers have their states set sequentially, starting at the bottom and working upwards.

The total input, x_i , to unit i is a linear function of the outputs, y_j , of the units that are connected by w_{ij} and of the weights, w_{ij} , on these connections.

$$x_i = \sum_j w_{ij} y_j + b_i \quad (1)$$

Units can be given biases by introducing an extra input to each input unit which always has a value of 1. The weight on this extra input is called the bias weight and is denoted by b_i . The bias weight can be fixed just like the other weights.

A unit has a real-valued output, y_i , which is a nonlinear function of its total input:

$$y_i = \frac{1}{1 + e^{-x_i}} \quad (2)$$

Handwritten Digit Recognition with a Back-Propagation Network
Y. Le Cun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel
AT&T Bell Laboratories, Holmdel, N. J. 07733

ABSTRACT
We present an application of back-propagation networks to handwritten digit recognition. Minimal preprocessing of the data was required, but architecture of the network was highly constrained and specifically designed for the task. The input of the network consists of normalized images of isolated digits. The method has a 9% error rate and about a 9% reject rate on spurious digits provided by the U.S. Postal Service.

1 INTRODUCTION
The main point of this paper is to show that large back-propagation (BP) networks can be applied to real image-recognition problems without a large, complex preprocessing stage requiring detailed engineering. Unlike most previous work on the subject (Denker et al., 1989), the learning network is directly fed with images, rather than feature vectors, thus demonstrating the ability of BP networks to deal with large amounts of low level information.

Previous work performed on simple digit images (Le Cun, 1989) showed that the architecture of the network strongly influences the network's generalization ability. Good generalization can only be obtained by designing a network architecture that contains a certain amount of a priori knowledge about the problem. The basic design principle is to minimize the number of free parameters that must be determined by the learning algorithm, without overly reducing the computational power of the network. This principle increases the probability of correct generalization because

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

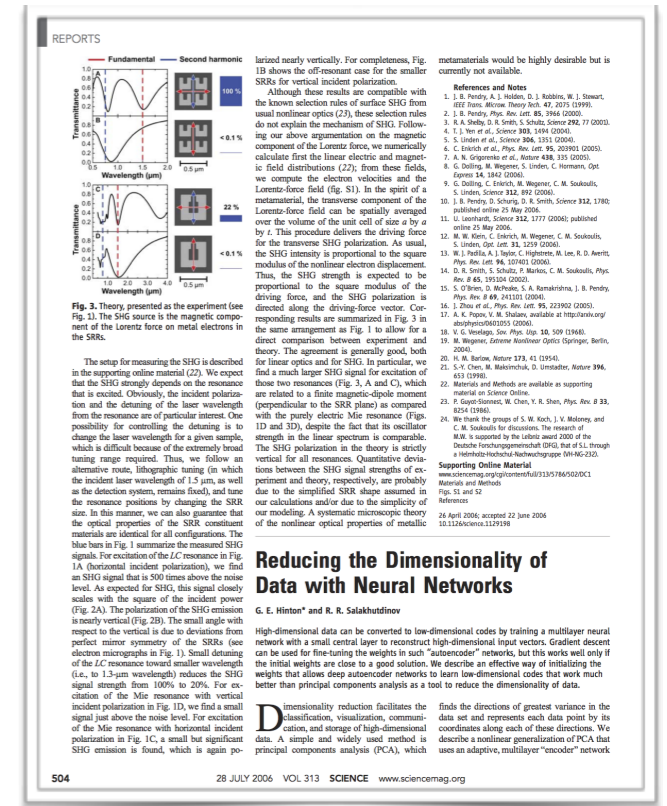
Reducing the Dimensionality of Data with Neural Networks

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

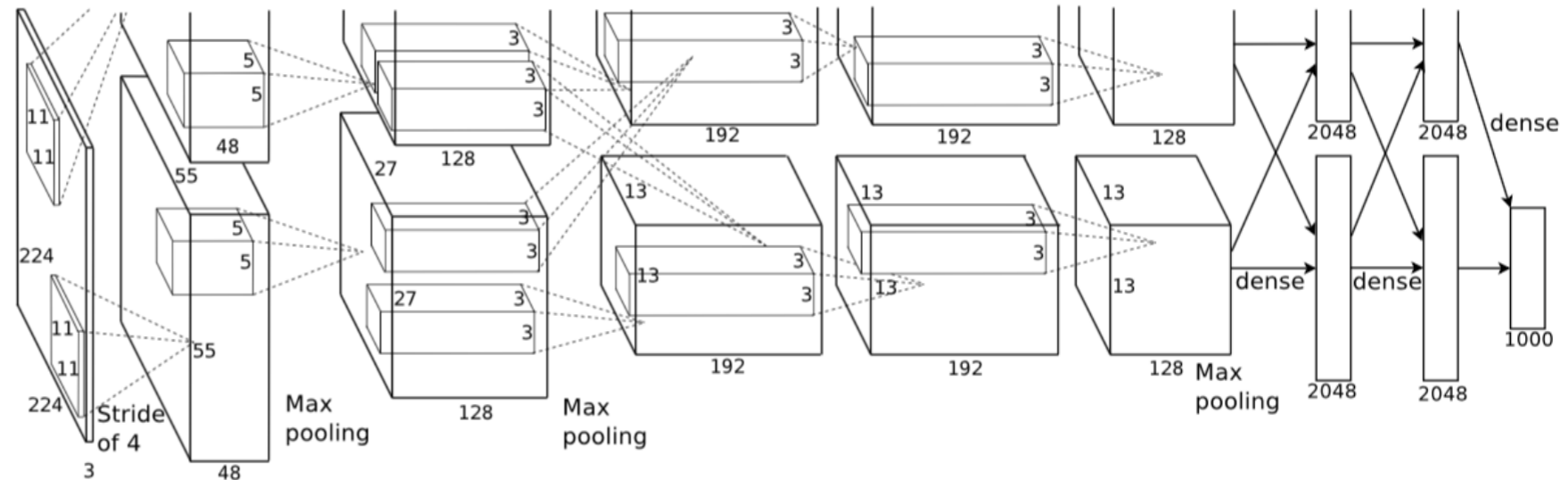


Image classification

ILSVRC 2012 Competition

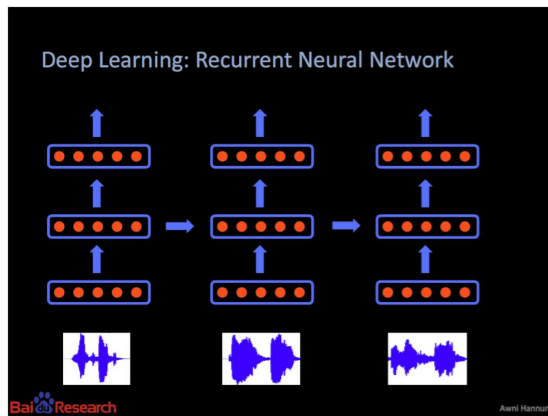
2012 Teams	%Error
Supervision (Toronto)	15.3
ISI (Tokyo)	26.1
VGG (Oxford)	26.9
XRCE/INRIA	27.0
UvA (Amsterdam)	29.6
INRIA/LEAR	33.4

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

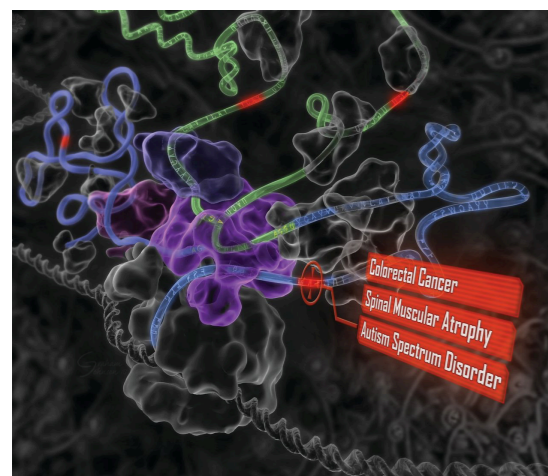
2012 – now
A Cambrian explosion in
deep learning



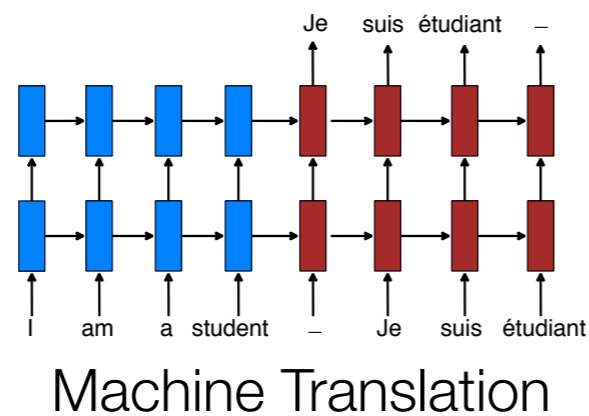
Speech recognition



Robotics



Genomics



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

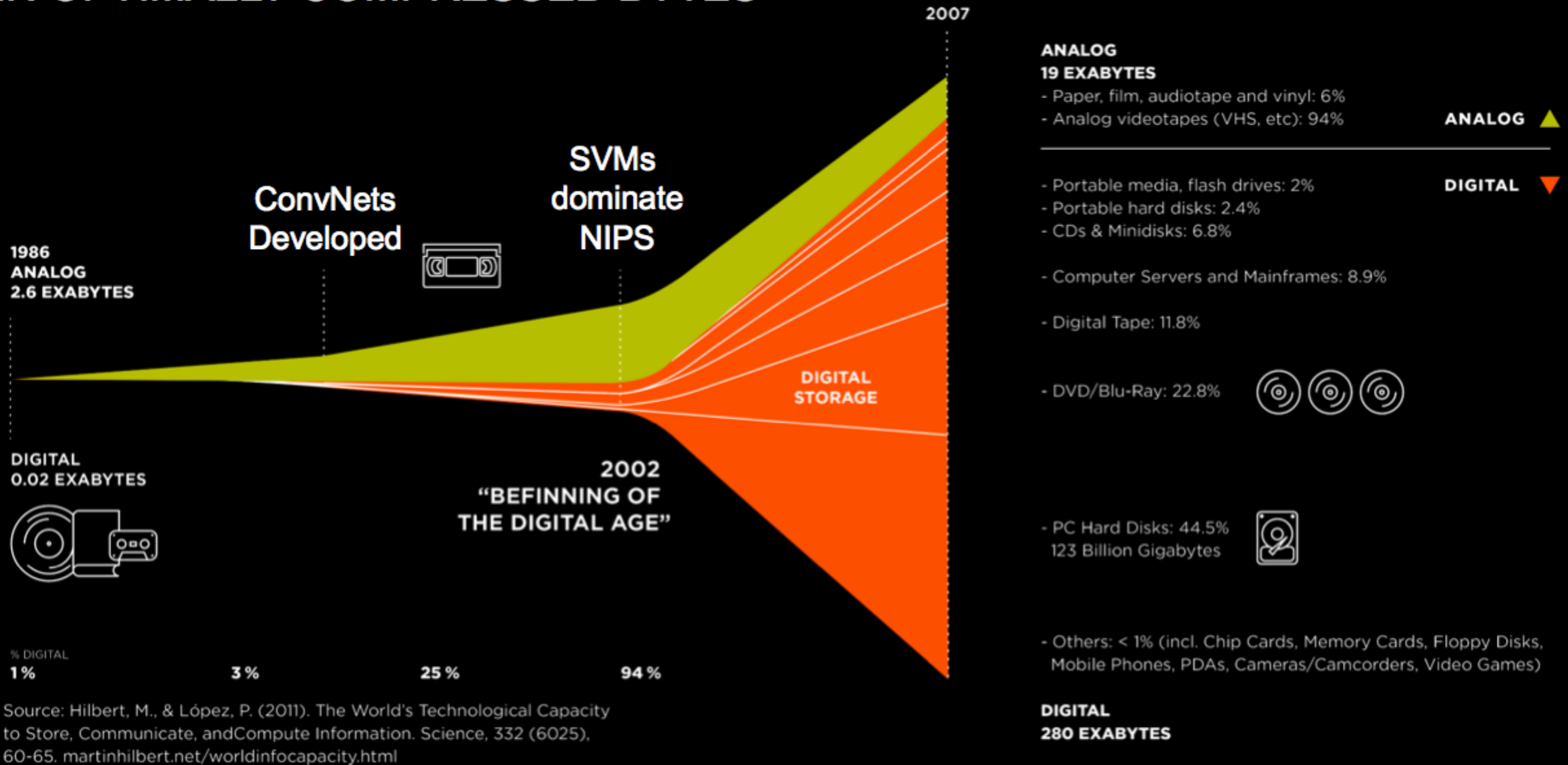
H. Y. Xiong et al., "The human splicing code reveals new insights into the genetic determinants of disease", Science 347, 2015

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)
Average No. of Years to Breakthrough:		3 years	18 years

Powerful Hardware

GOOGLE DATACENTER



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

600 kWatts
\$5,000,000

STANFORD AI LAB

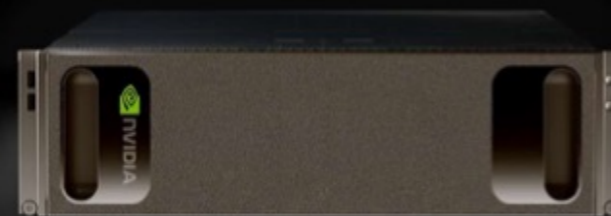


3 GPU-Accelerated Servers
12 GPUs • 18,432 cores

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

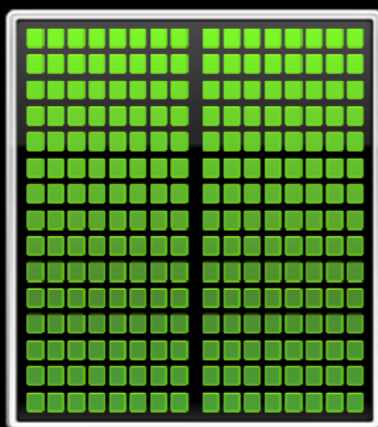
CPU

Optimized for
Serial Tasks



GPU Accelerator

Optimized for
Parallel Tasks



TITAN X

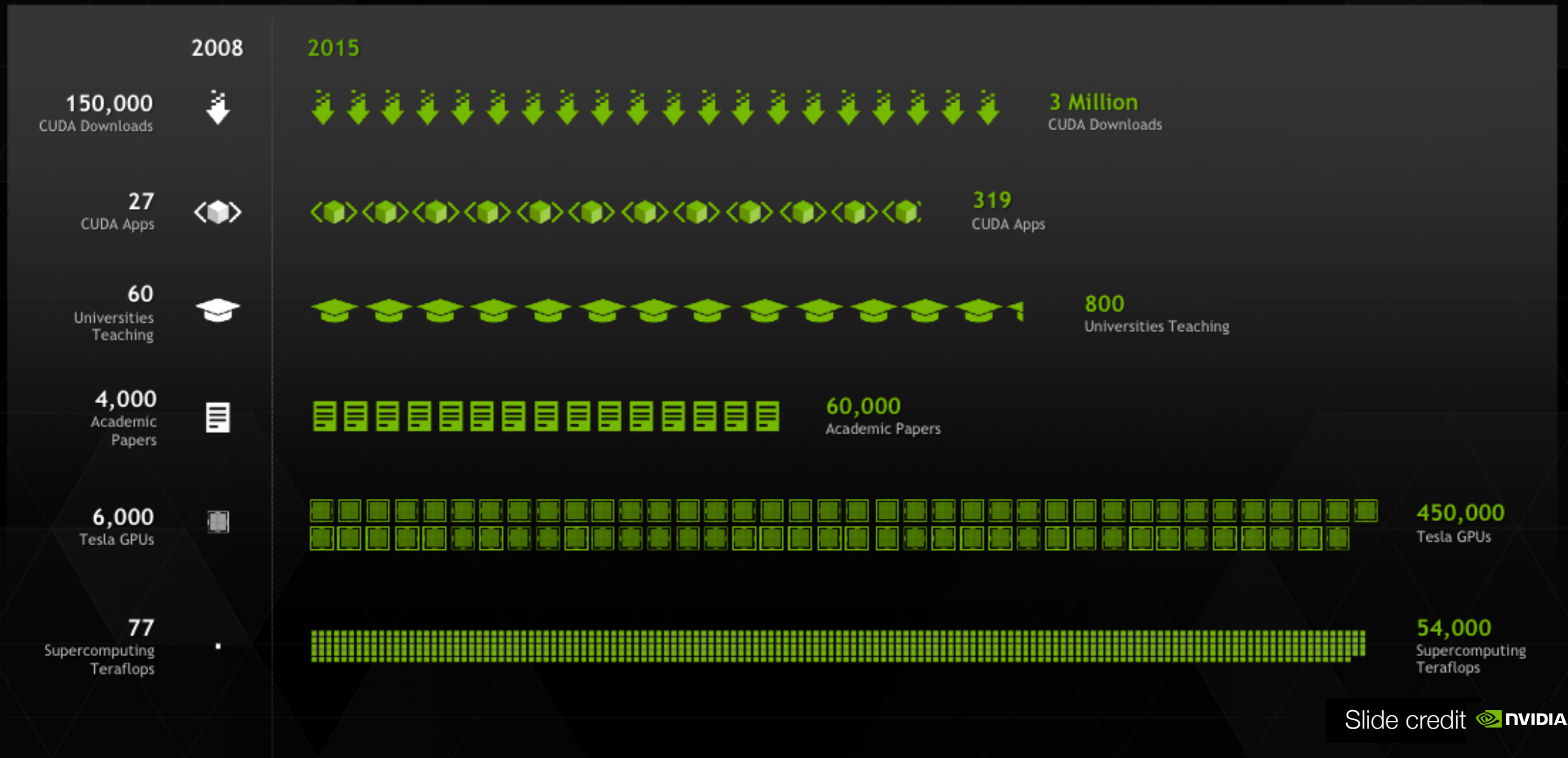
THE WORLD'S FASTEST GPU

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

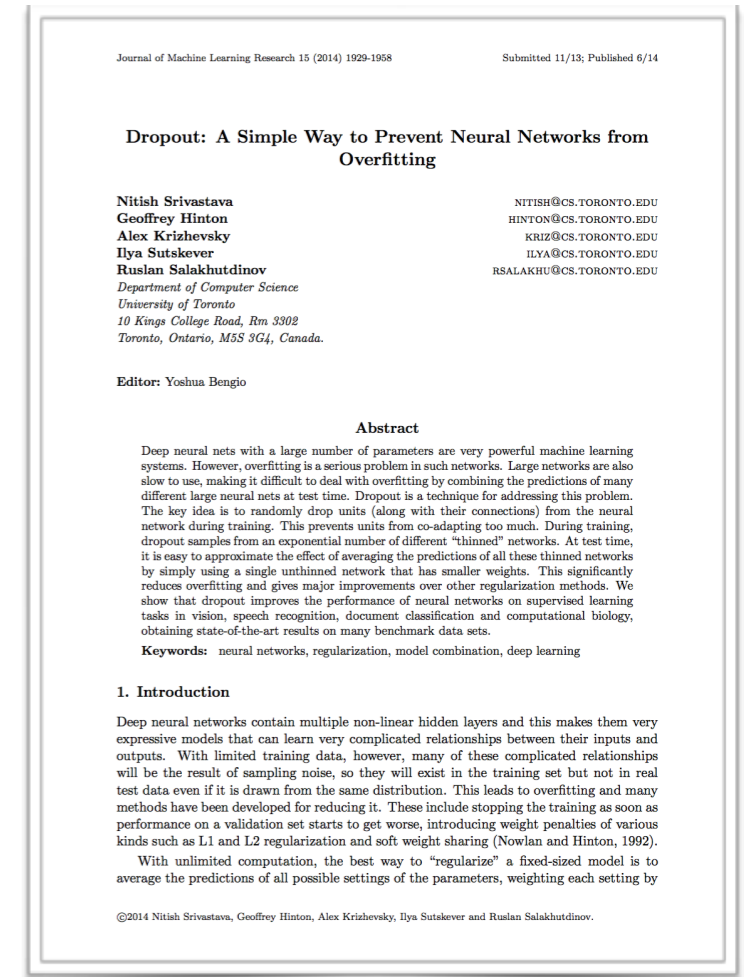
Dropout: A Simple Way to Prevent Neural Networks from Overfitting

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

Christian Szegedy

Google Inc., 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 achieves the same accuracy with 14 times fewer training steps, and beats the original model by a significant margin. Using an ensemble of batch-normalized networks, we improve upon the best published result on ImageNet classification: reaching 4.9% top-5 validation error (and 4.8% test error), exceeding the accuracy of human raters.

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

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

Sergey Ioffe
Google Inc., sioffe@google.com

Christian Szegedy
Google Inc., 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 allows us to use much higher learning rates and be less careful about initialization. It also acts as a regularizer, in some cases eliminating the need for Dropout. Applied to a state-of-the-art image classification model, Batch Normalization achieves the same accuracy with 14 times fewer training steps, and beats the original model by a significant margin. Using an ensemble of batch-normalized networks, we improve upon the best published result on ImageNet classification: reaching 4.9% top-5 validation error (and 4.8% test error), exceeding the accuracy of human raters.

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 are affected by the parameters of all preceding layers – so that small changes to the network parameters amplify as the network becomes deeper.

The change in the distributions of layers' inputs presents a problem because the layers need to continuously adapt to the new distribution. When the input distribution to a learning system changes, it is said to experience *covariate shift* (Shimodaira, 2000). This is typically handled via domain adaptation (Jiang, 2008). However, the notion of covariate shift can be extended beyond the learning system as a whole, to apply to its parts, such as a sub-network or a layer. Consider a network computing

1 Introduction

Deep learning has dramatically advanced the state of the art in vision, speech, and many other areas. Stochastic gradient descent (SGD) has proved to be an effective way of training deep networks, and SGD variants such as momentum (Sutskever et al., 2013) and Adagrad (Duchi et al., 2011) have been used to achieve state of the art performance. SGD optimizes the parameters Θ of the network, so as to minimize the loss

$$\Theta = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(x_i, \Theta)$$

where x_1, \dots, x_N is the training data set. With SGD, the training proceeds in steps, and at each step we consider a *mini-batch* x_1, \dots, x_m of size m . The mini-batch is used to approximate the gradient of the loss function with respect to the parameters, by computing

$$\frac{1}{m} \frac{\partial \ell(x_i, \Theta)}{\partial \Theta}$$

$$\ell = F_2(F_1(u, \Theta_1), \Theta_2)$$

where F_1 and F_2 are arbitrary transformations, and the parameters Θ_1, Θ_2 are to be learned so as to minimize the loss ℓ . Learning Θ_2 can be viewed as if the inputs $x = F_1(u, \Theta_1)$ are fed into the sub-network

$$\ell = F_2(x, \Theta_2).$$

For example, a gradient descent step

$$\Theta_2 \leftarrow \Theta_2 - \frac{\alpha}{m} \sum_{i=1}^m \frac{\partial F_2(x_i, \Theta_2)}{\partial \Theta_2}$$

(for batch size m and learning rate α) is exactly equivalent to that for a stand-alone network F_2 with input x . Therefore, the input distribution properties that make training more efficient – such as having the same distribution between the training and test data – apply to training the sub-network as well. As such it is advantageous for the distribution of x to remain fixed over time. Then, Θ_2 does

1

- 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

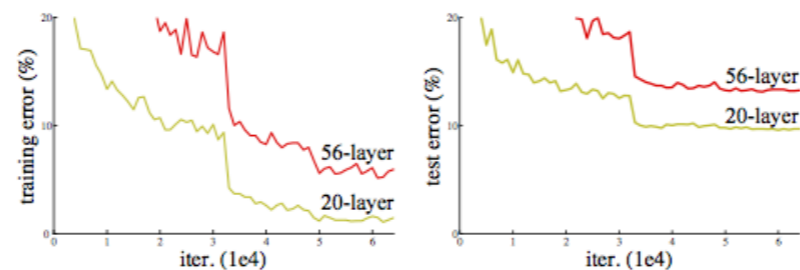


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

greatly benefited from very deep models.

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

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 on the ImageNet test set. This result won the 1st place on the ILSVRC 2015 classification task. We also present analysis on CIFAR-10 with 100 and 1000 layers.

The depth of representations is of central importance for many visual recognition tasks. Solely due to our extremely deep representations, we obtain a 28% relative improvement on the COCO object detection dataset. Deep residual nets are foundations of our submissions to ILSVRC & COCO 2015 competitions¹, where we also won the 1st places on the tasks of ImageNet detection, ImageNet localization, COCO detection, and COCO segmentation.

1. Introduction

Deep convolutional neural networks [22, 21] have led to a series of breakthroughs for image classification [21, 50, 40]. Deep networks naturally integrate low/mid/high-level features [50] and classifiers in an end-to-end multi-layer fashion, and the “levels” of features can be enriched by the number of stacked layers (depth). Recent evidence [41, 44] reveals that network depth is of crucial importance, and the leading results [41, 44, 13, 16] on the challenging ImageNet dataset [36] all exploit “very deep” [41] models, with a depth of sixteen [41] to thirty [16]. Many other non-trivial visual recognition tasks [8, 12, 7, 32, 27] have also

¹<http://image-net.org/challenges/LSVRC/2015/> and <http://mscoco.org/dataset/#detection-challenge2015>.

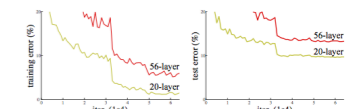


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

greatly benefited from very deep models.

Driven by the significance of depth, a question arises: *Is learning better networks as easy as stacking more layers?* An obstacle to answering this question was the notorious problem of vanishing/exploding gradients [1, 9], which hamper convergence from the beginning. This problem, however, has been largely addressed by normalized initialization [23, 9, 37, 13] and intermediate normalization layers [16], which enable networks with tens of layers to start converging for stochastic gradient descent (SGD) with back-propagation [22].

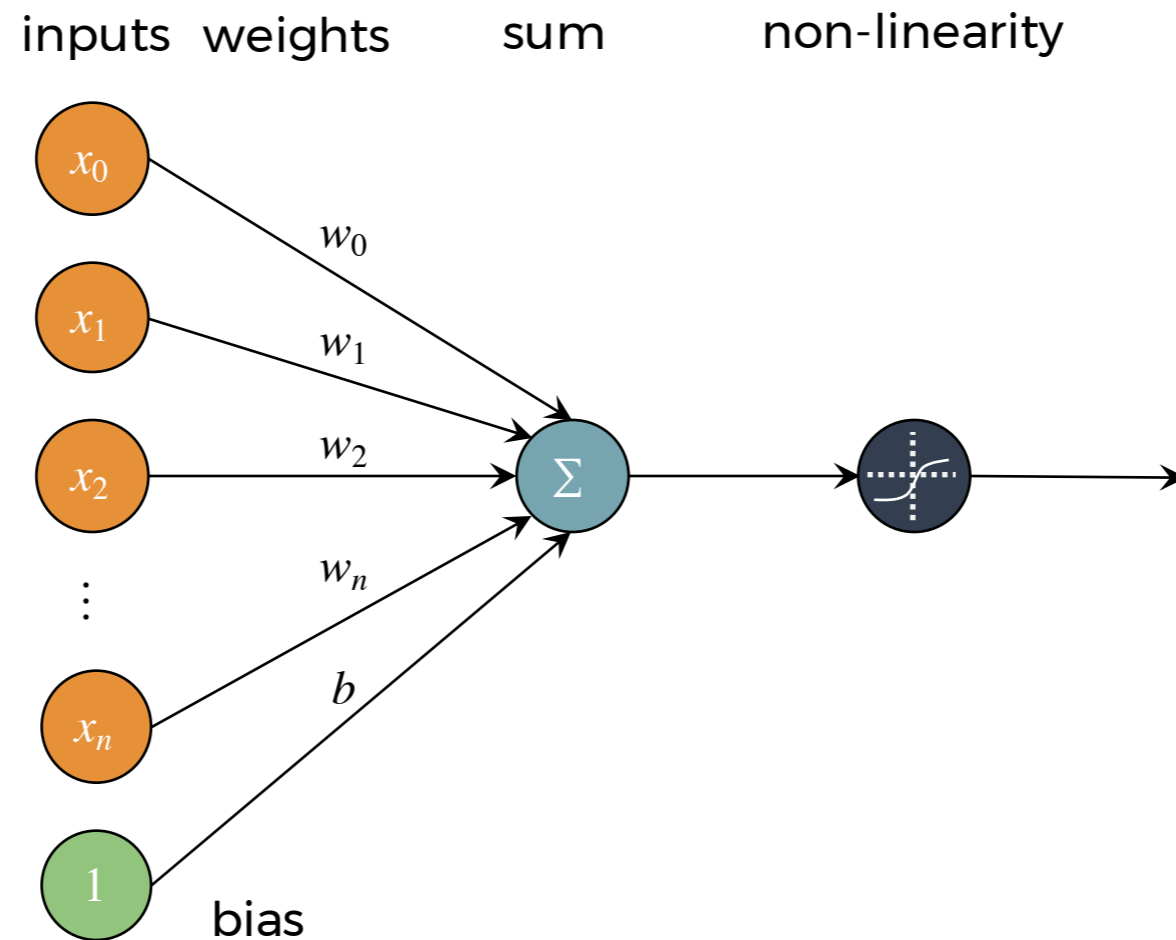
When deeper networks are able to start converging, a degradation problem has been exposed: with the network depth increasing, accuracy gets saturated (which might be unsurprising) and then degrades rapidly. Unexpectedly, such degradation is *not caused by overfitting*, and adding more layers to a suitably deep model leads to *higher training error*, as reported in [11, 42] and thoroughly verified by our experiments. Fig. 1 shows a typical example.

The degradation (of training accuracy) indicates that not all systems are similarly easy to optimize. Let us consider a shallower architecture and its deeper counterpart that adds more layers onto it. There exists a solution by *construction* to the deeper model: the added layers are *identity* mapping, and the other layers are copied from the learned shallower model. The existence of this constructed solution indicates that a deeper model should produce no higher training error than its shallower counterpart. But experiments show that our current solvers on hand are unable to find solutions that

- 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:

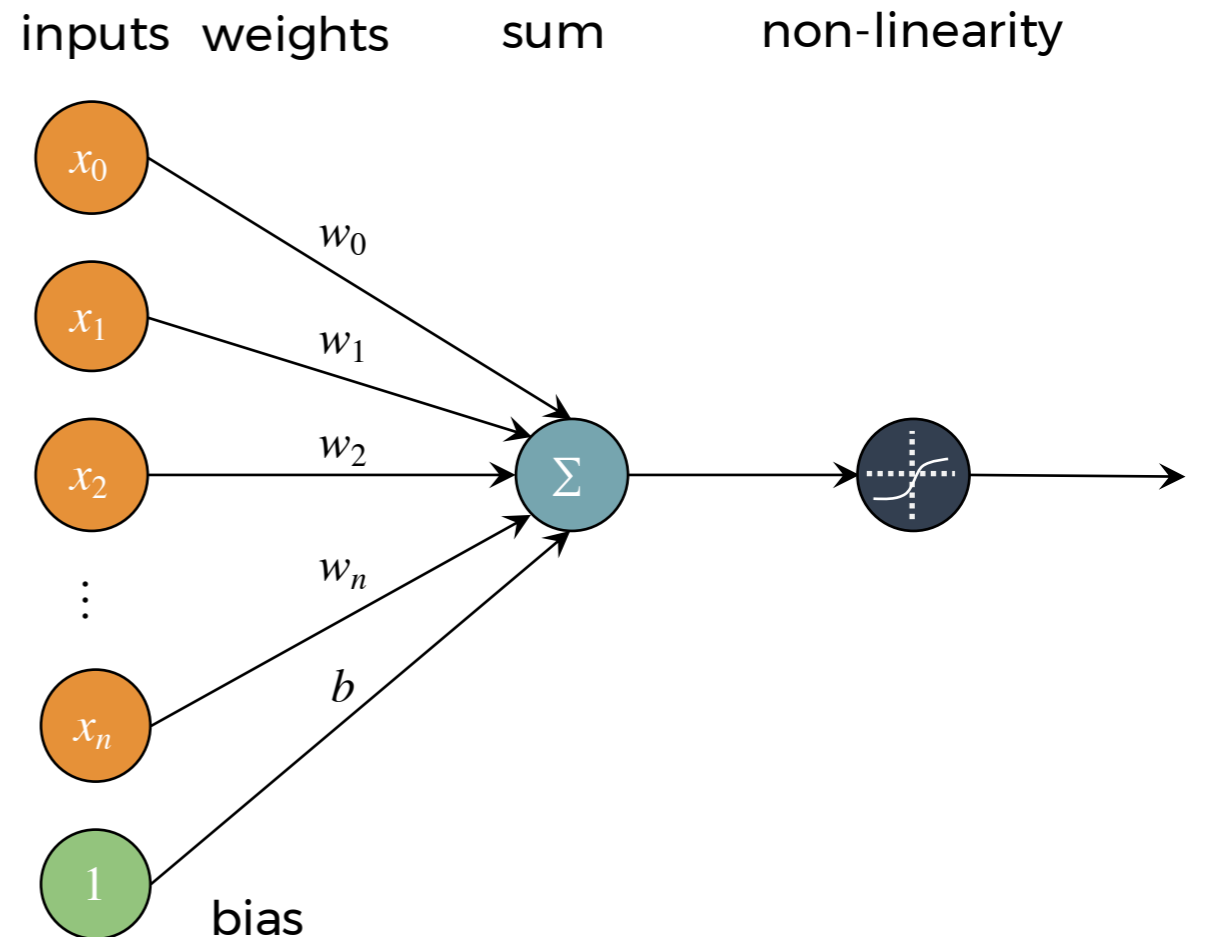
$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

where

w are the weights (parameters)

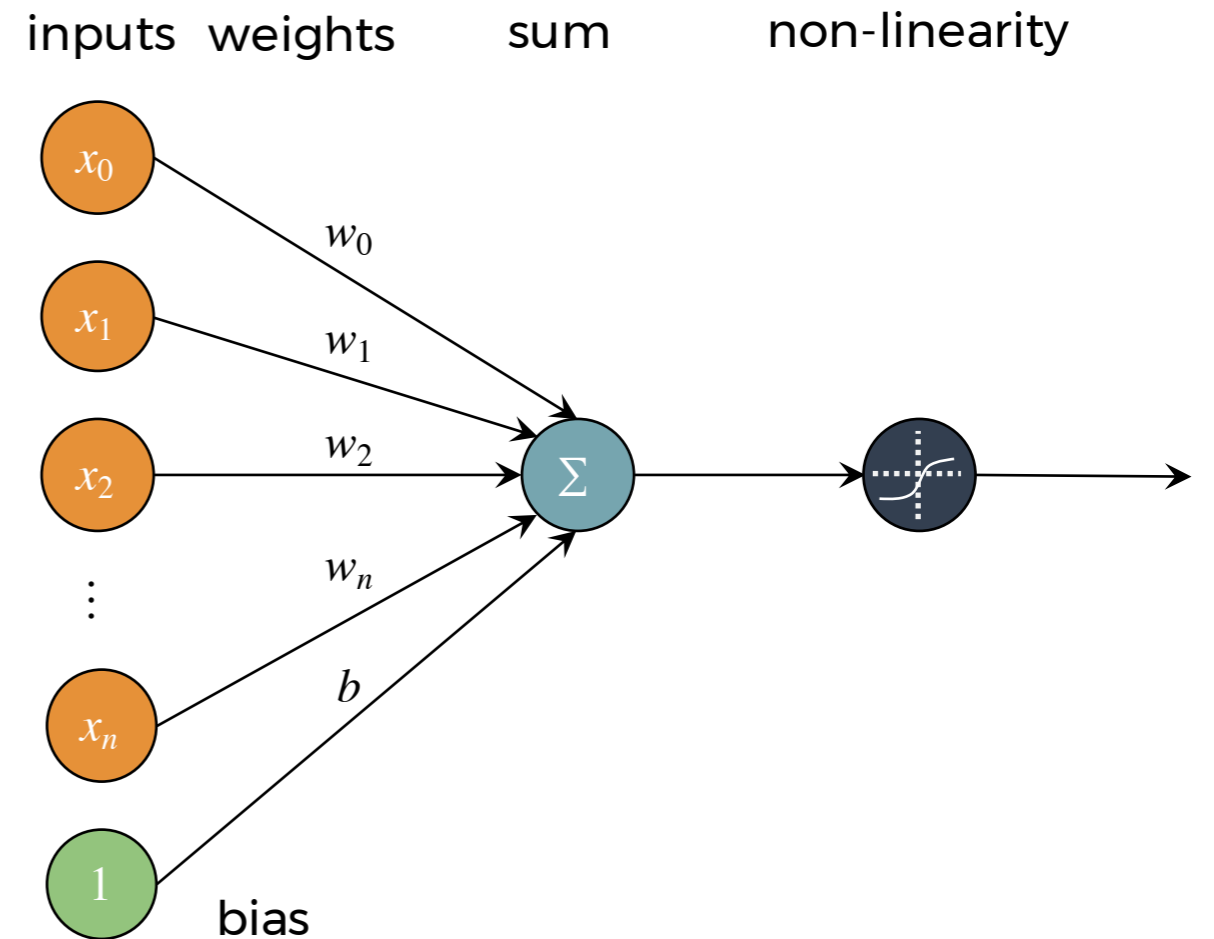
b is the bias term

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

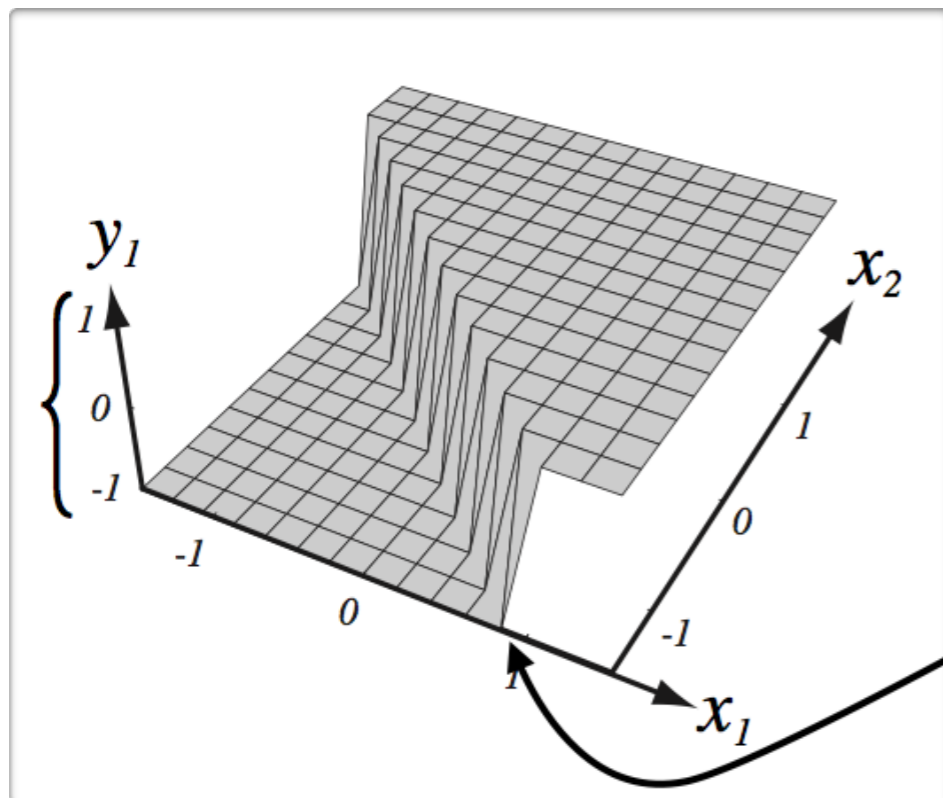


Output Activation of The Neuron

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



Range is determined by $g(\cdot)$

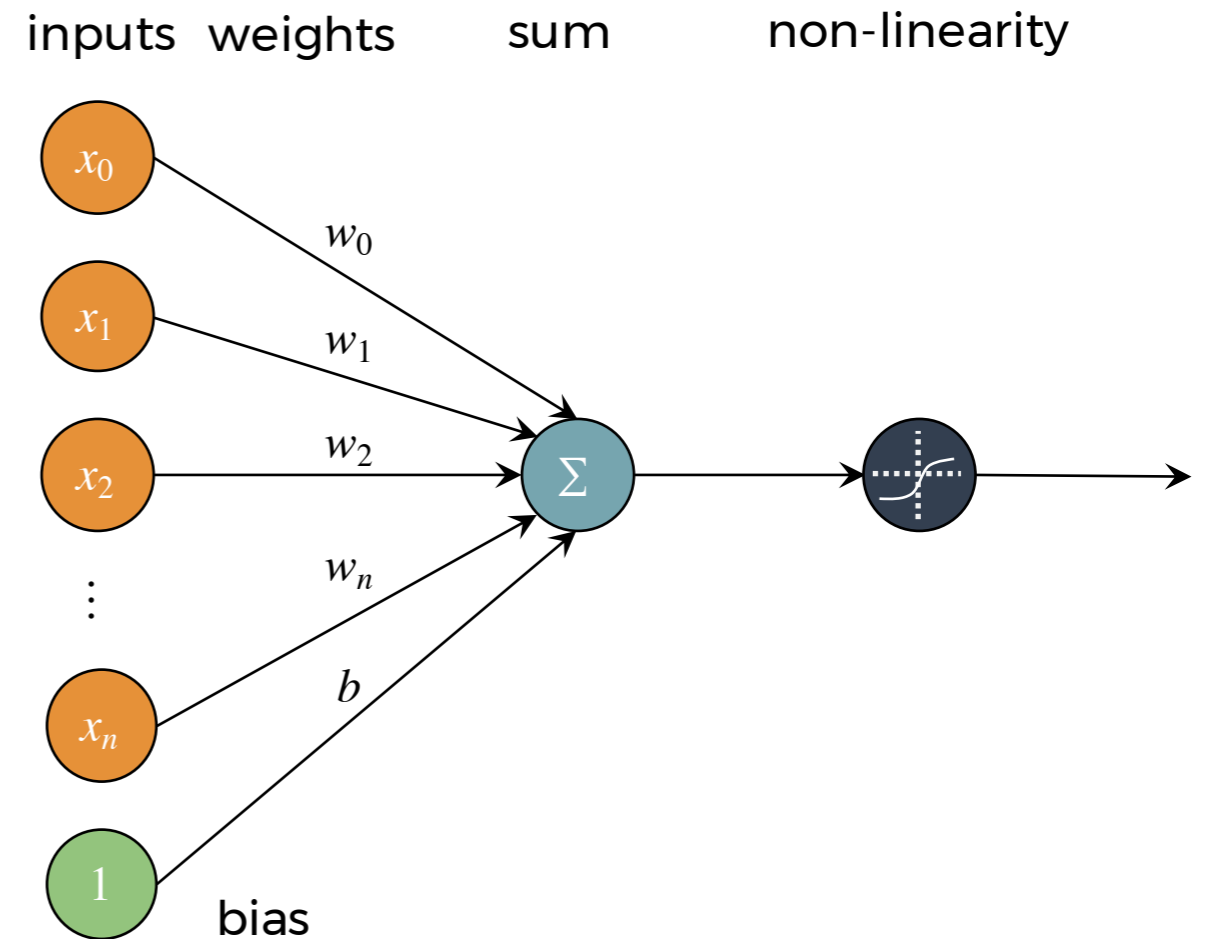
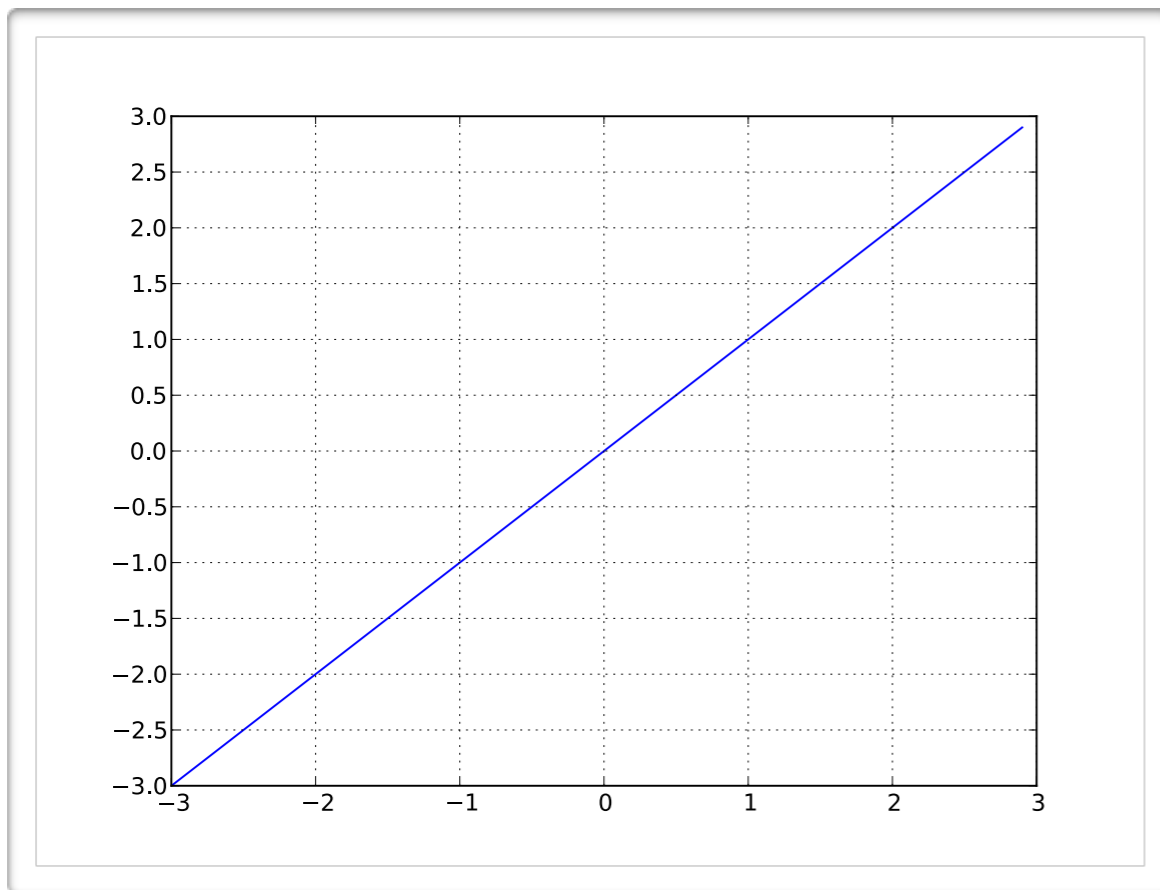


Bias only changes the position of the riff

Linear Activation Function

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

$$g(a) = a$$

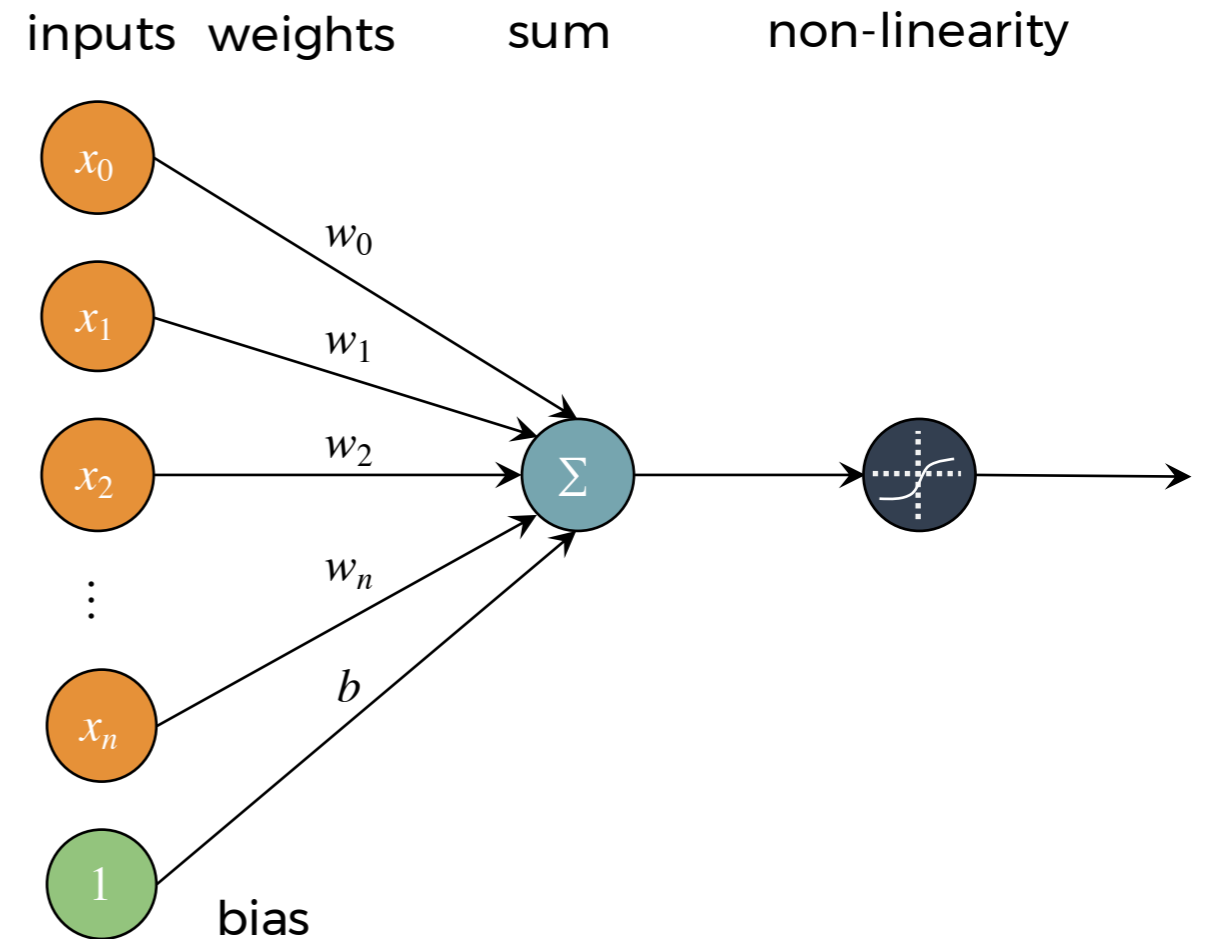
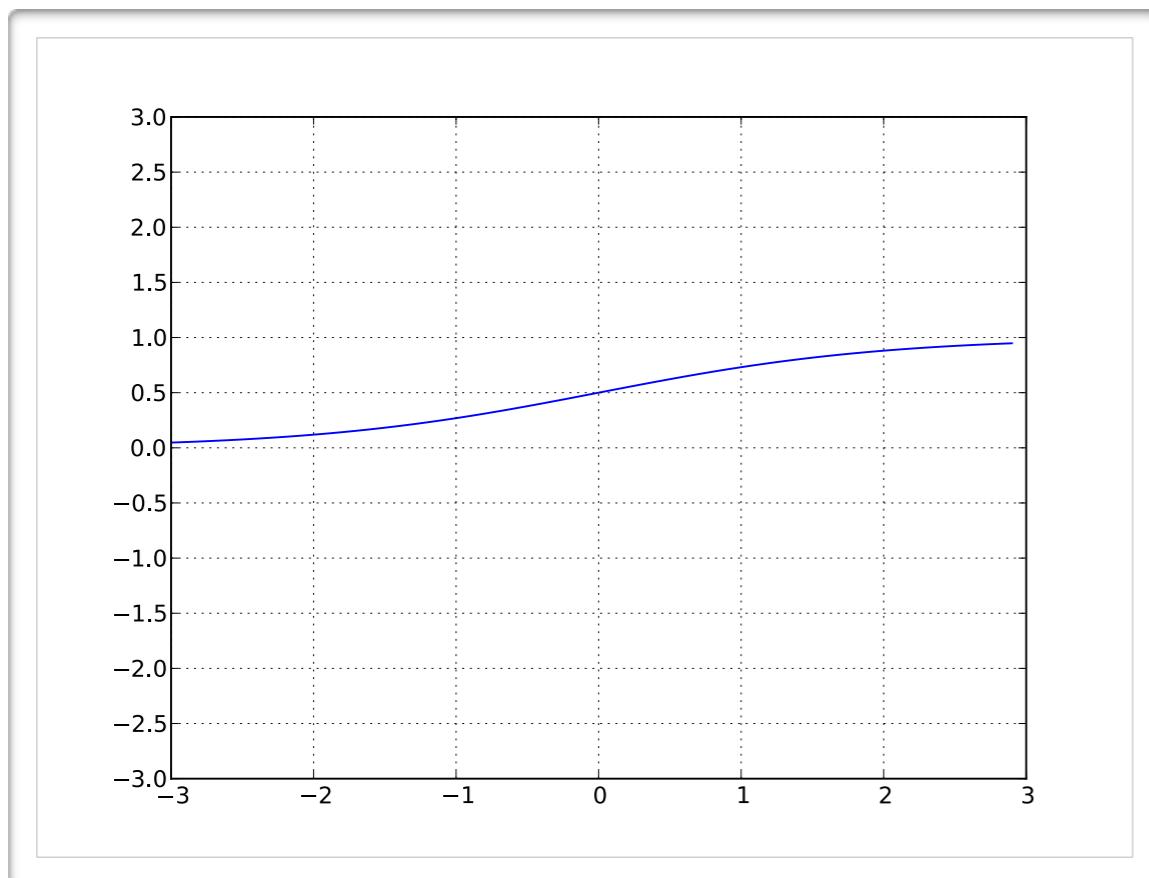


No nonlinear transformation
No input squashing

Sigmoid Activation Function

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

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

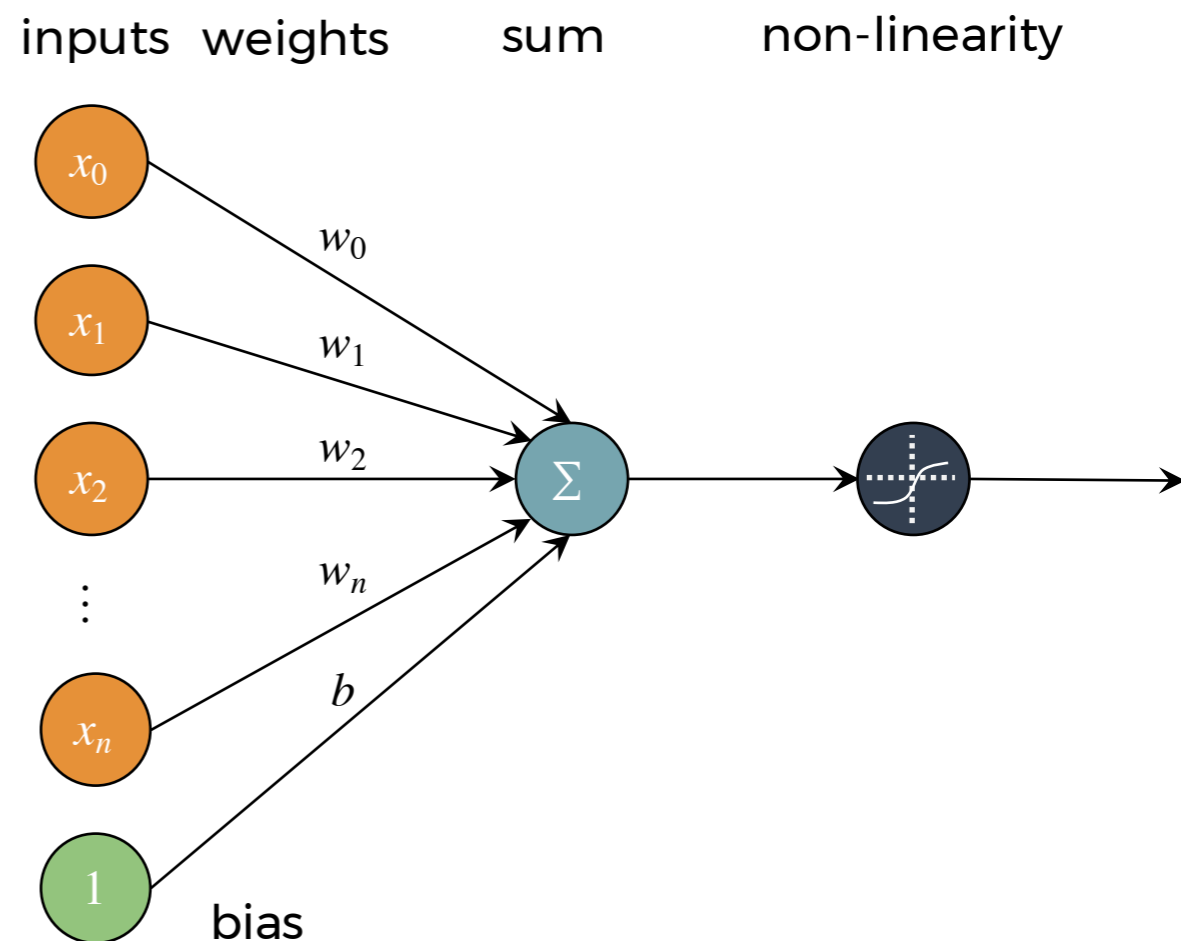
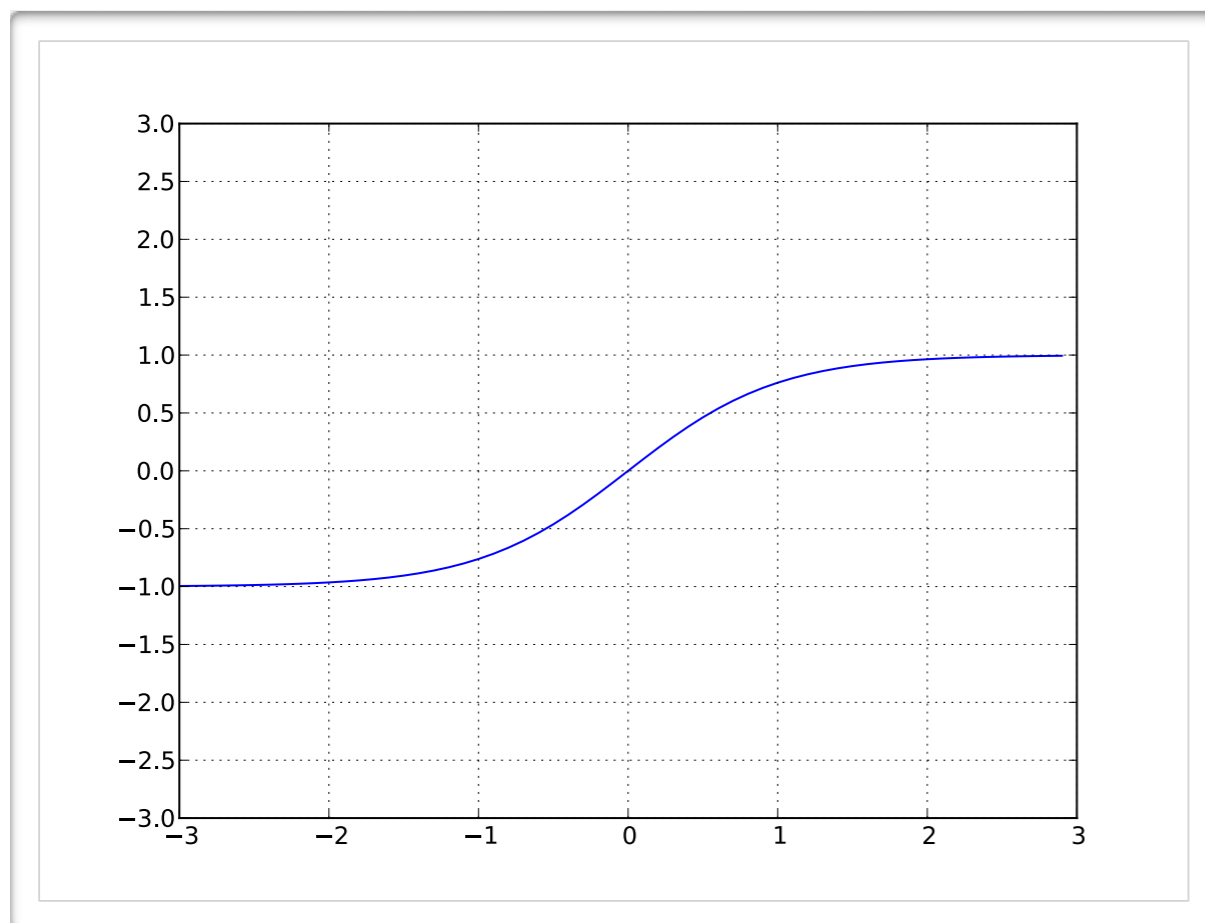


Squashes the neuron's output
between 0 and 1
Always positive
Bounded
Strictly Increasing

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

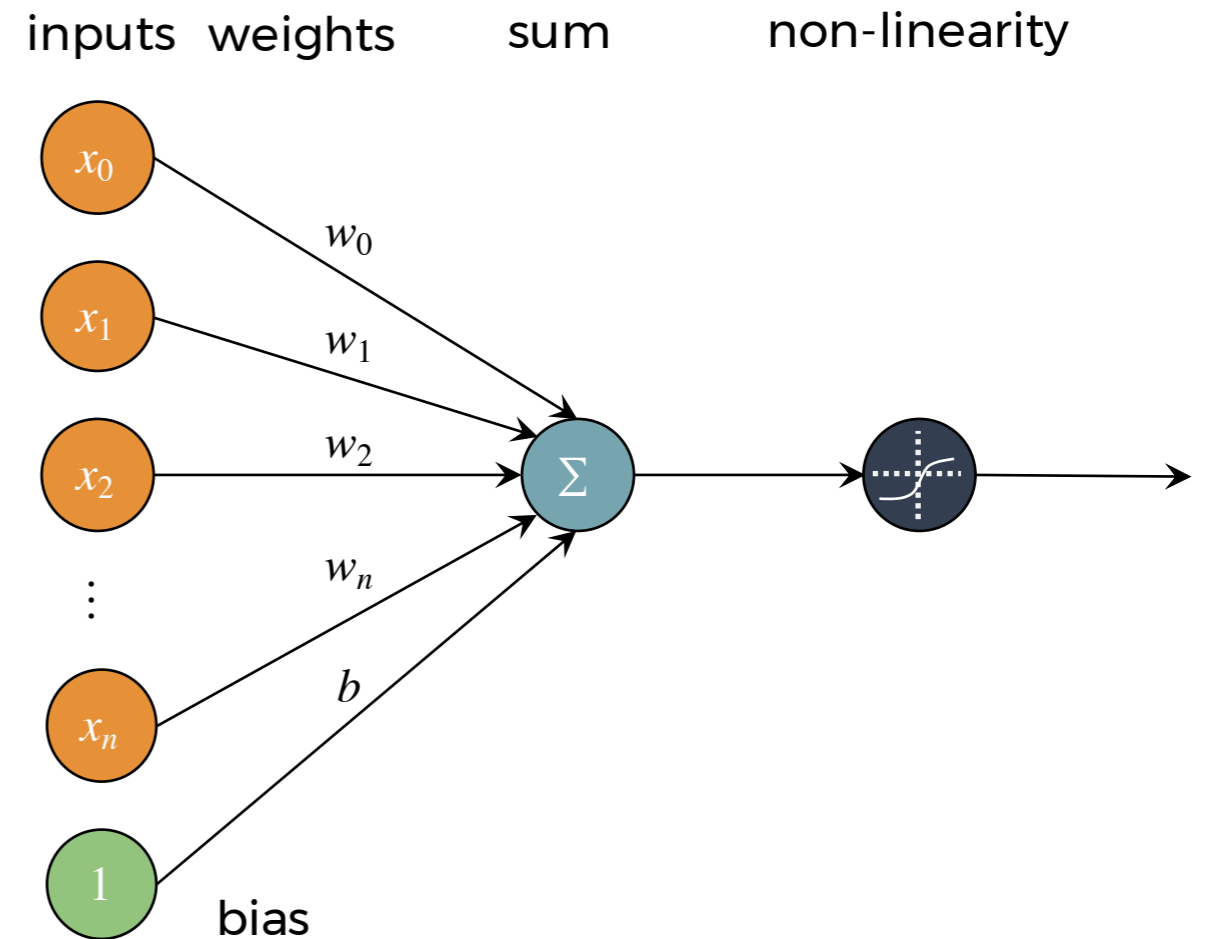
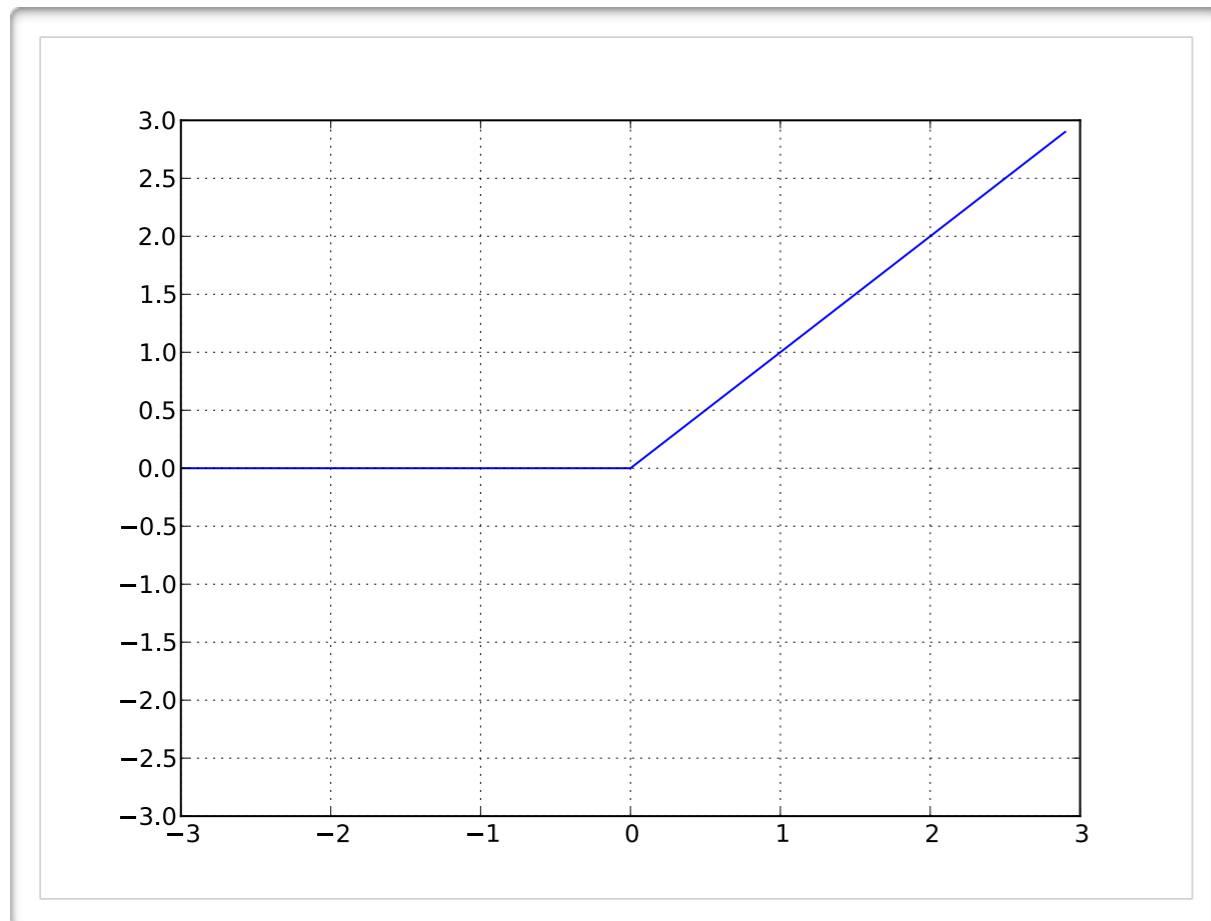
Bounded

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)

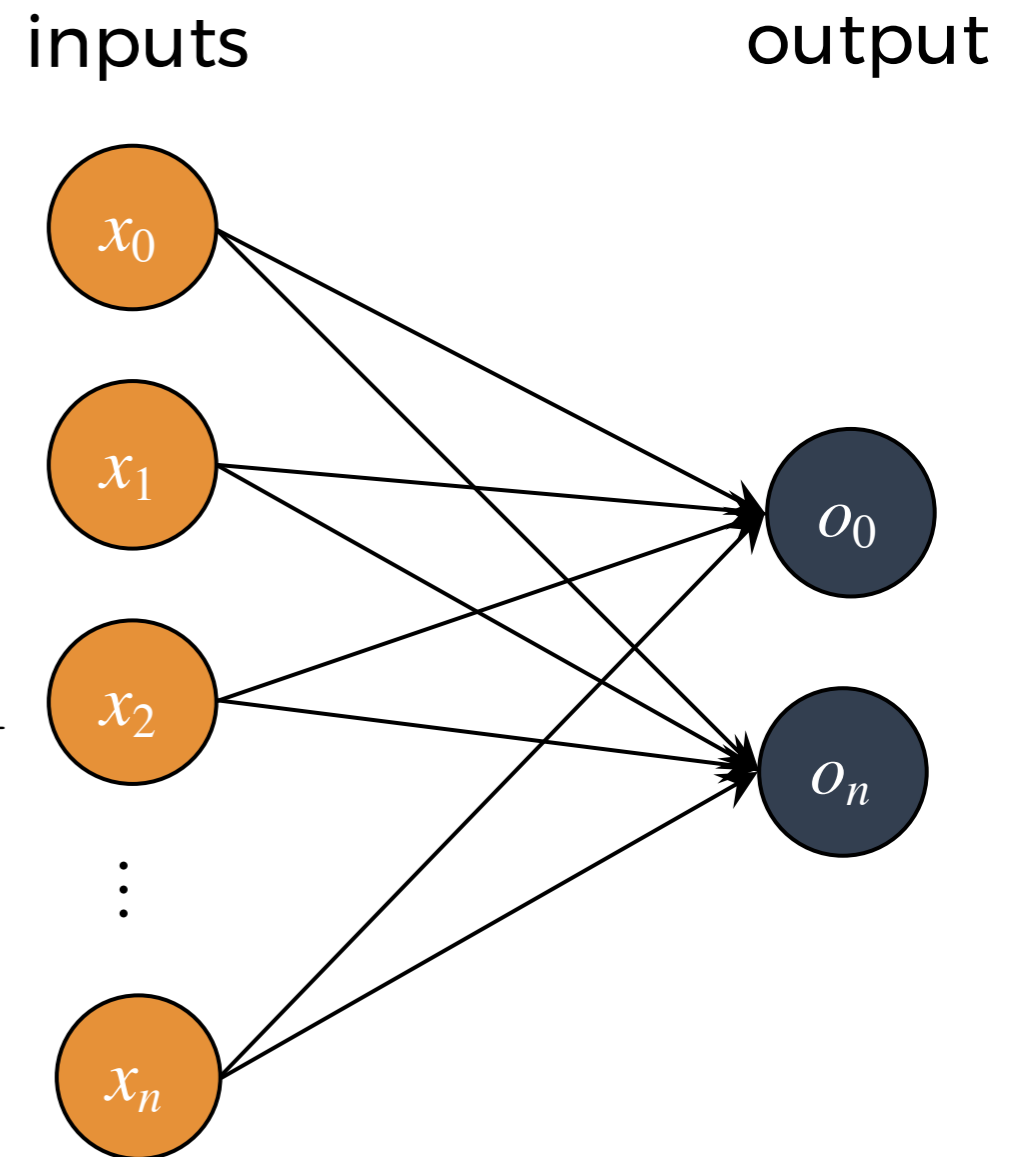
Not upper bounded

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:

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

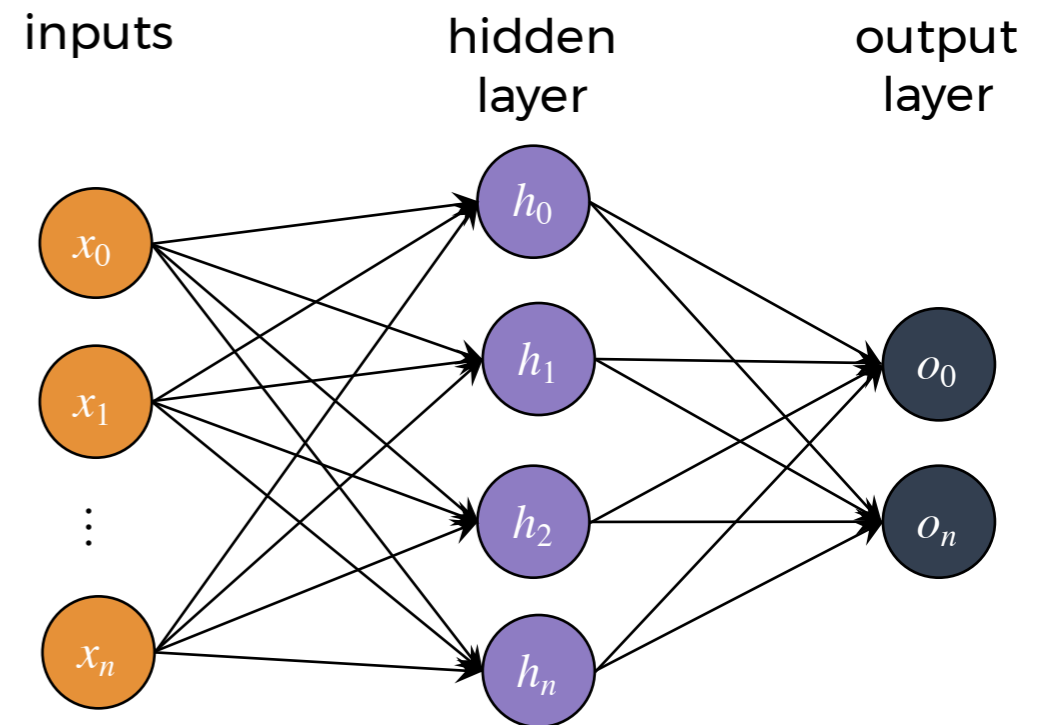
$$\left(a(\mathbf{x})_i = b_i^{(1)} + \sum_j W_{i,j}^{(1)} x_j \right)$$

- 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 $k > 0$

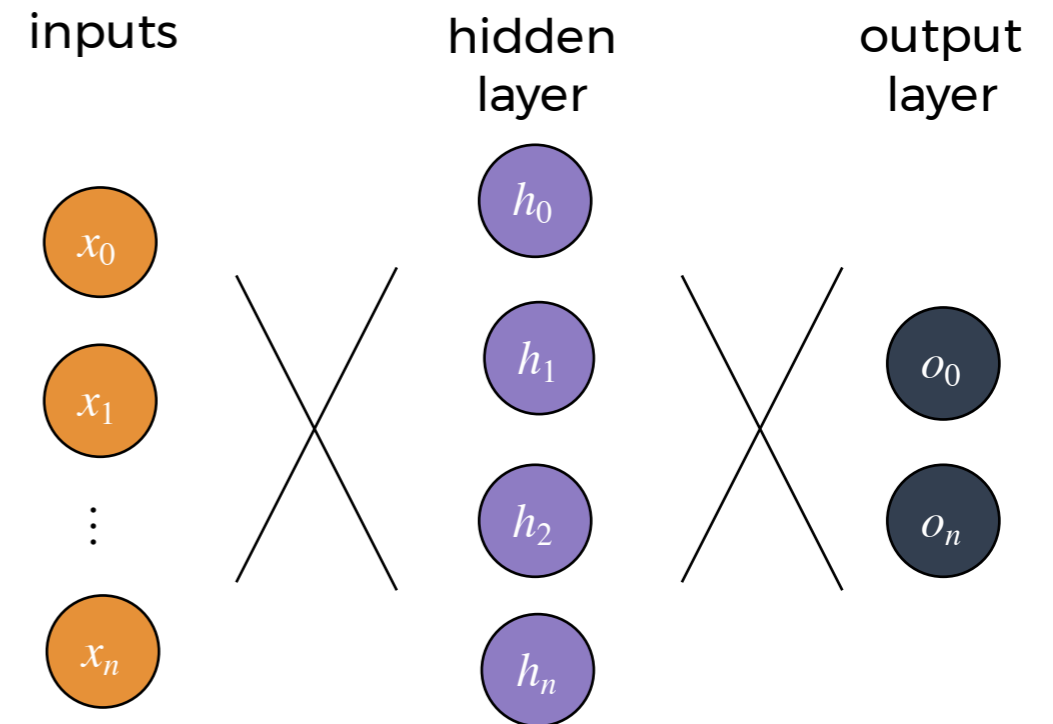
$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$$

- hidden layer activation from 1 to L :

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

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

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



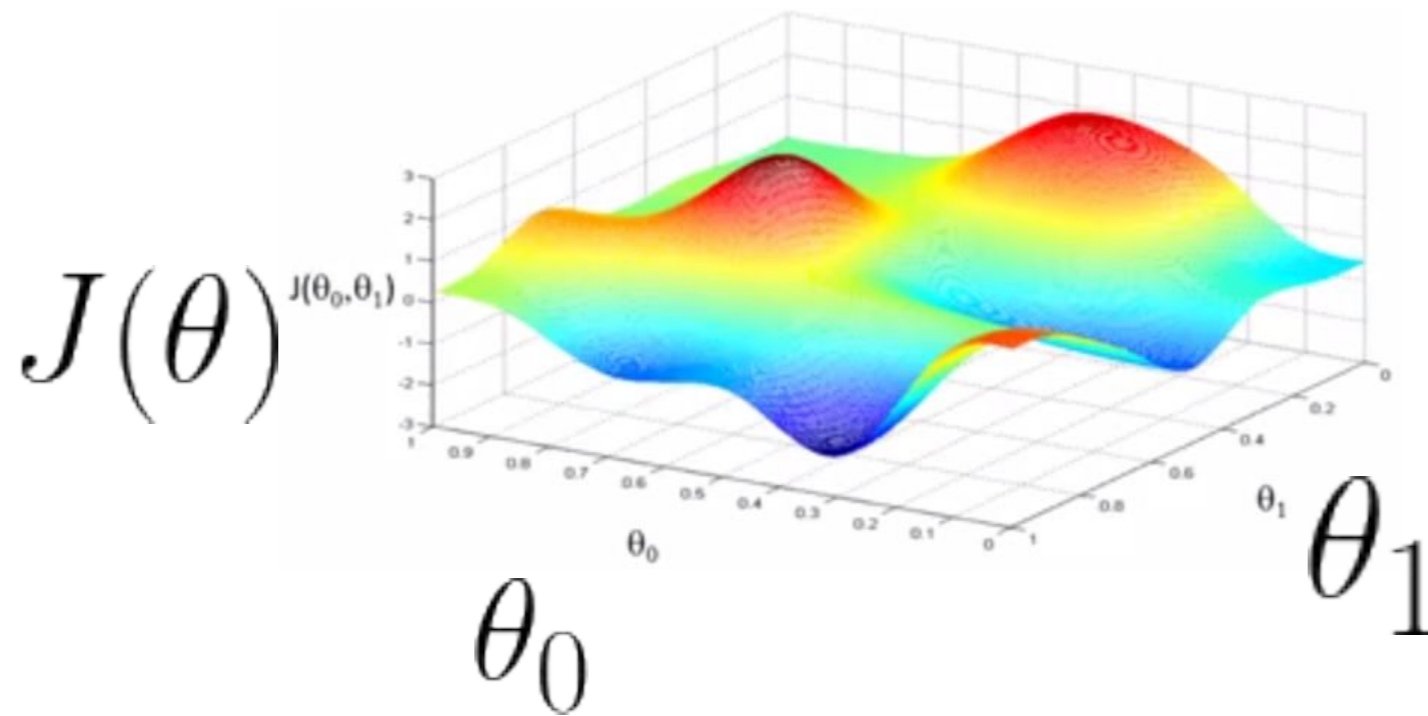
Training

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

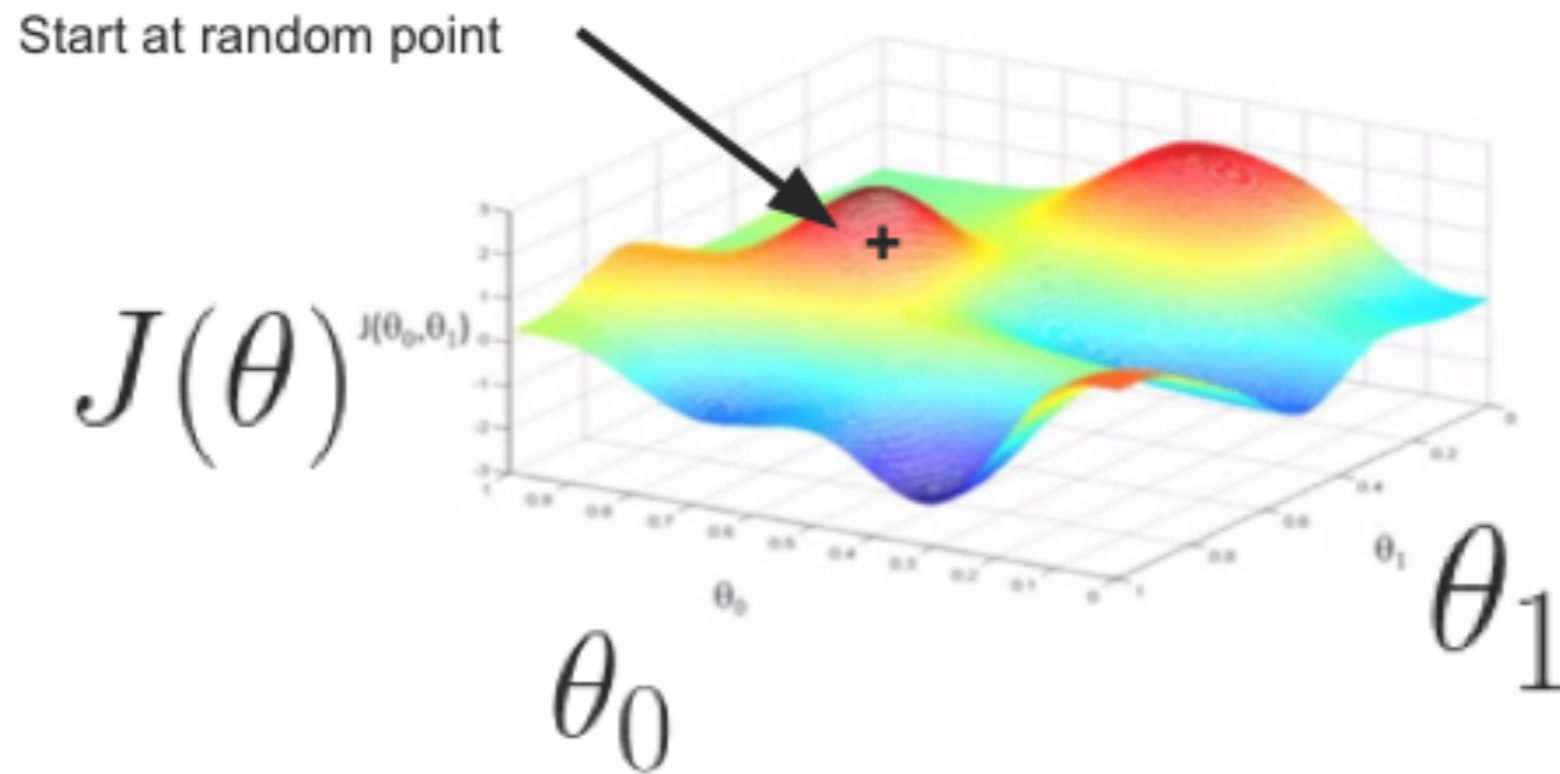
$$\theta = W_1, W_2 \dots W_n$$

- 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)

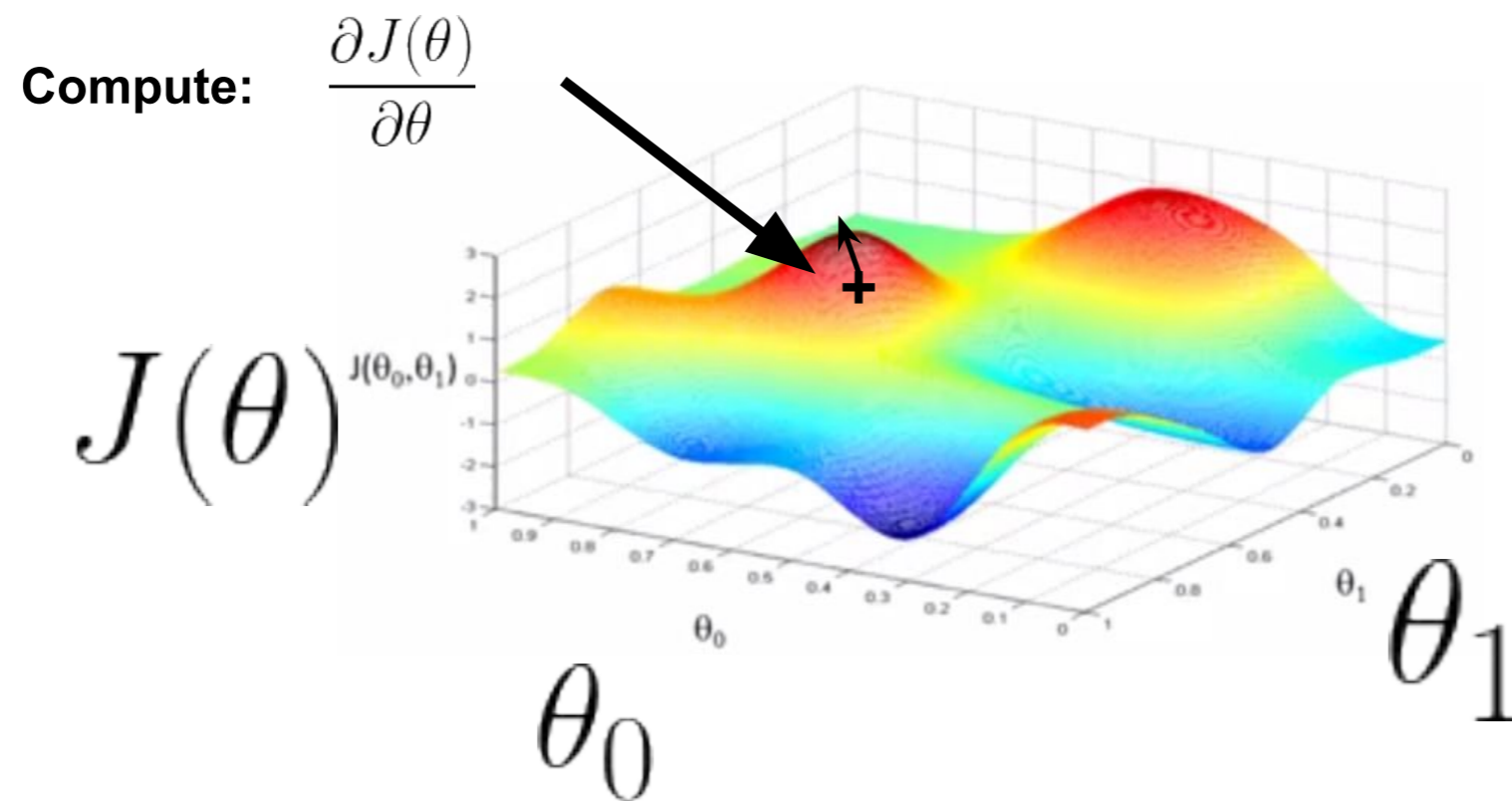
Loss is a **function** of the model's parameters



How to minimize loss?

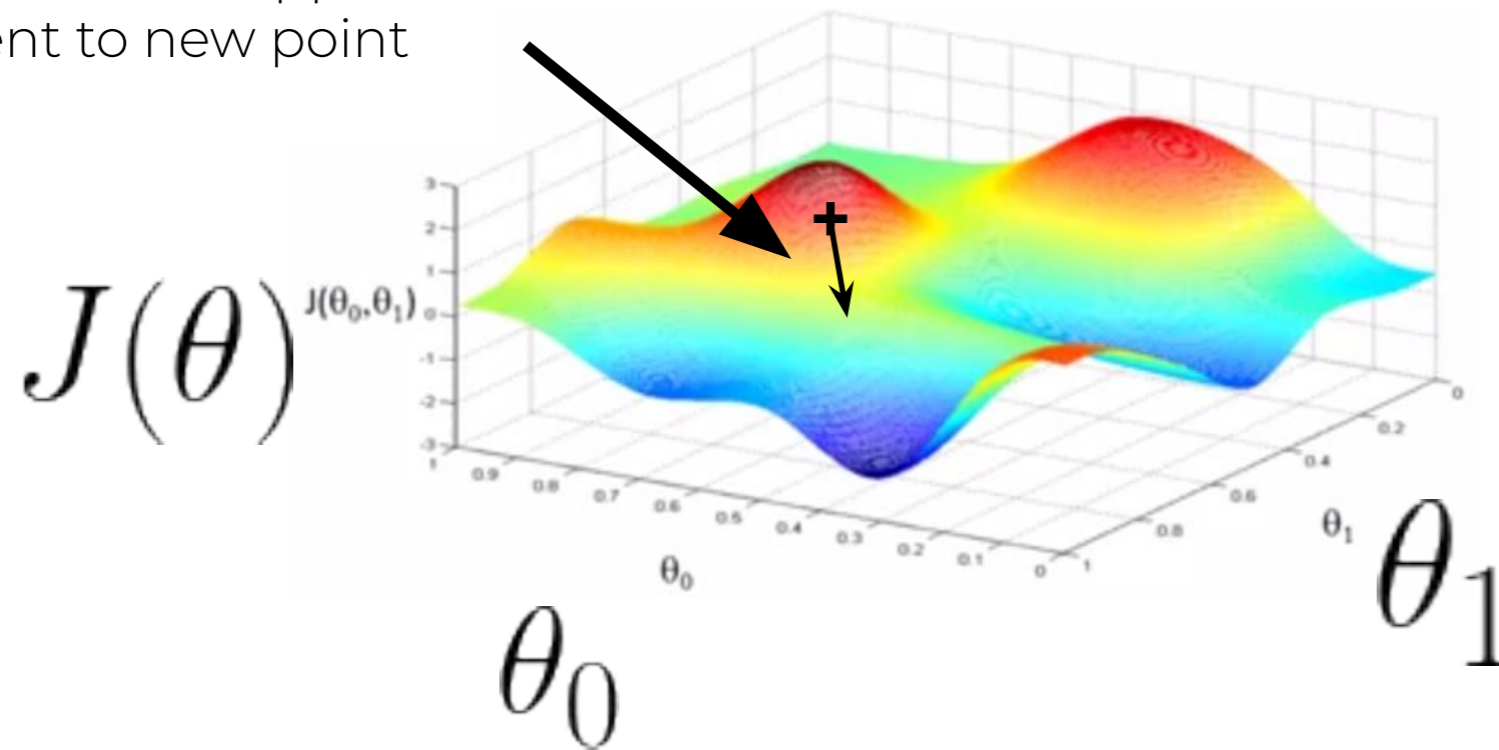


How to minimize loss?



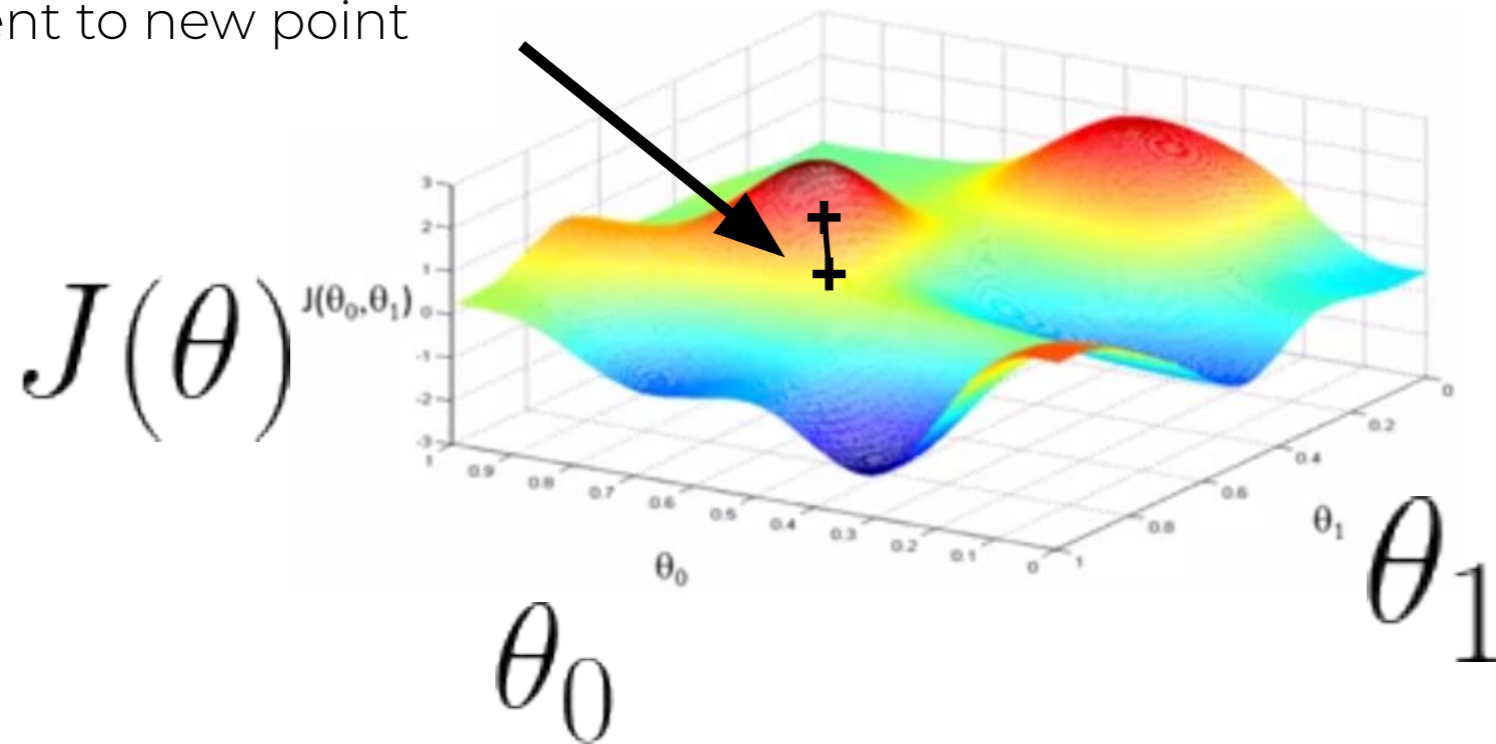
How to minimize loss?

Move in direction opposite of gradient to new point



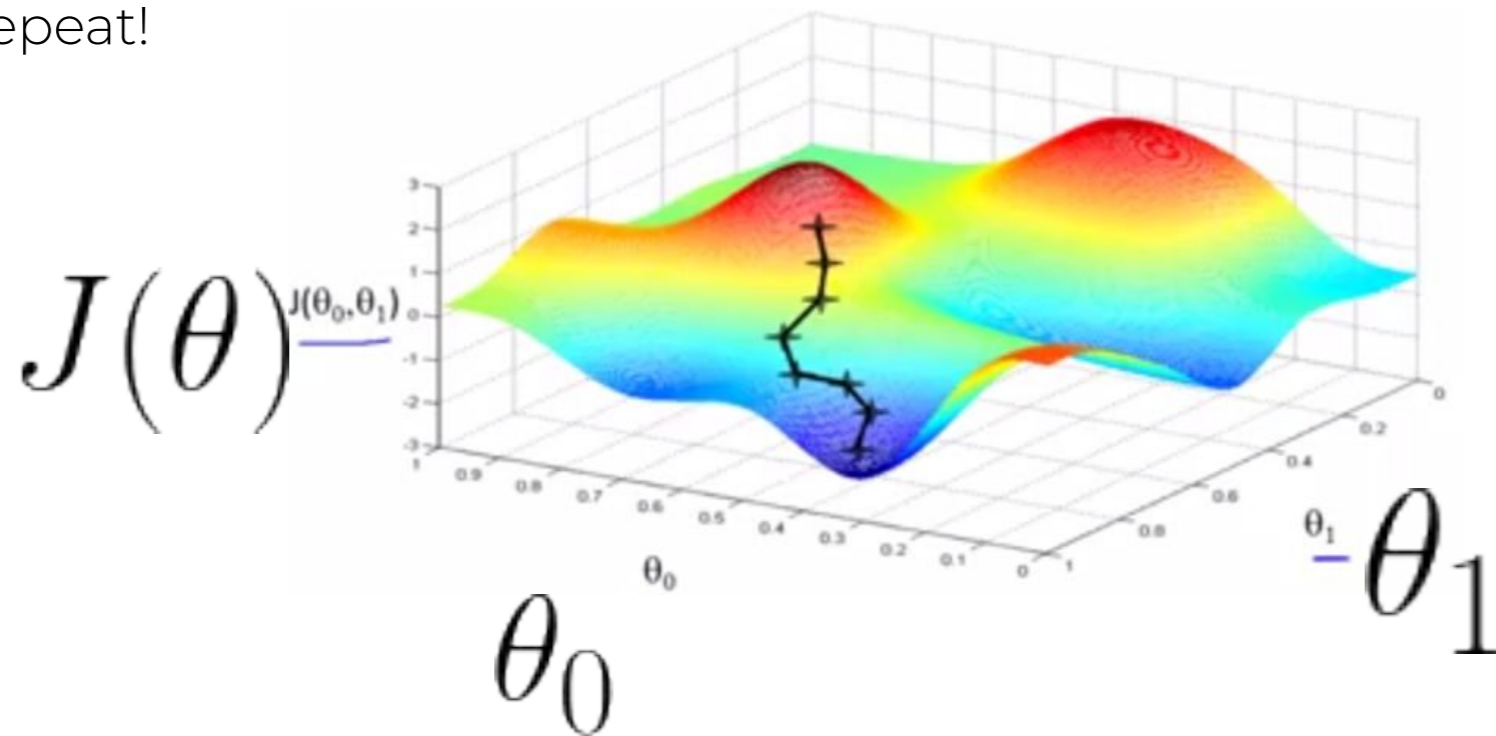
How to minimize loss?

Move in direction opposite of gradient to new point



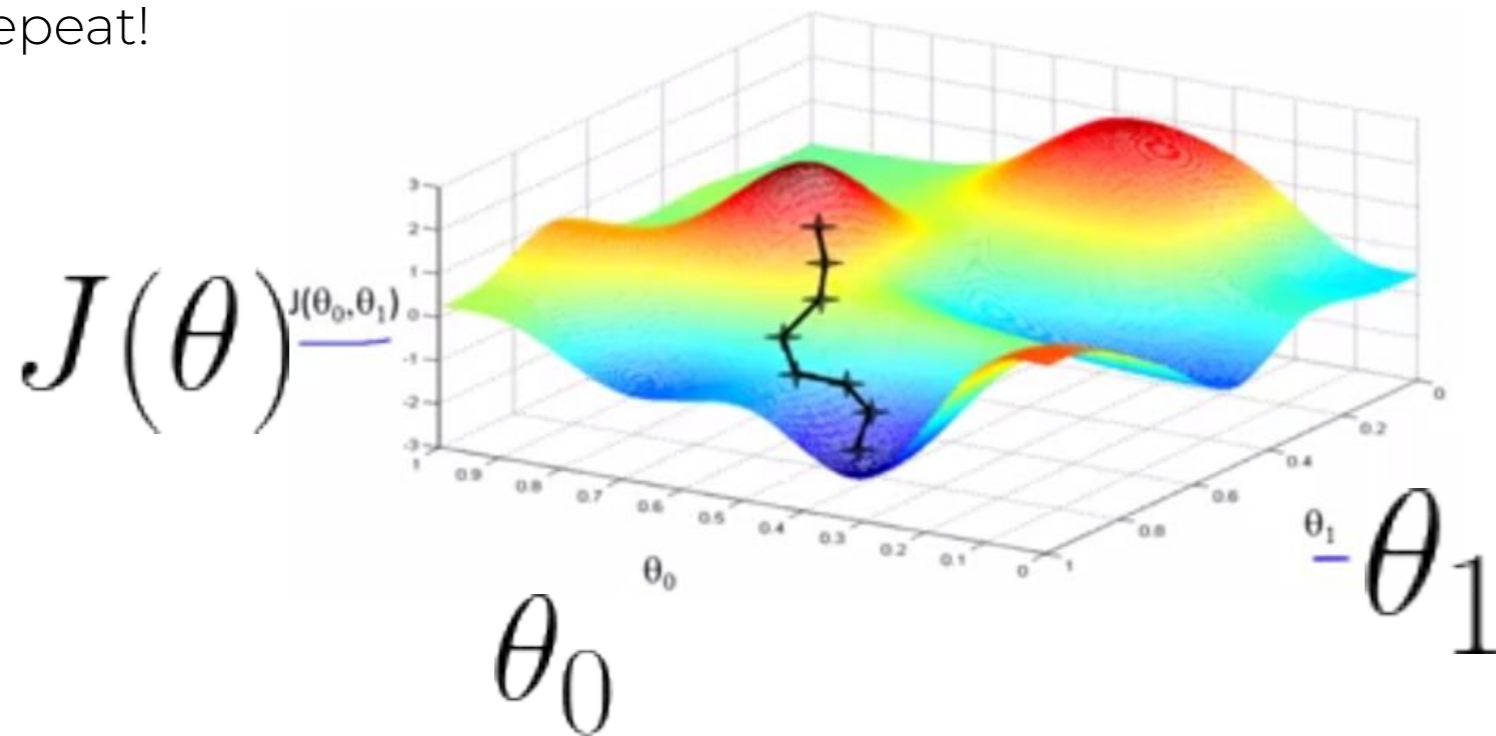
How to minimize loss?

Repeat!



This is called Stochastic Gradient Descent (SGD)

Repeat!



Stochastic Gradient Descent (SGD)

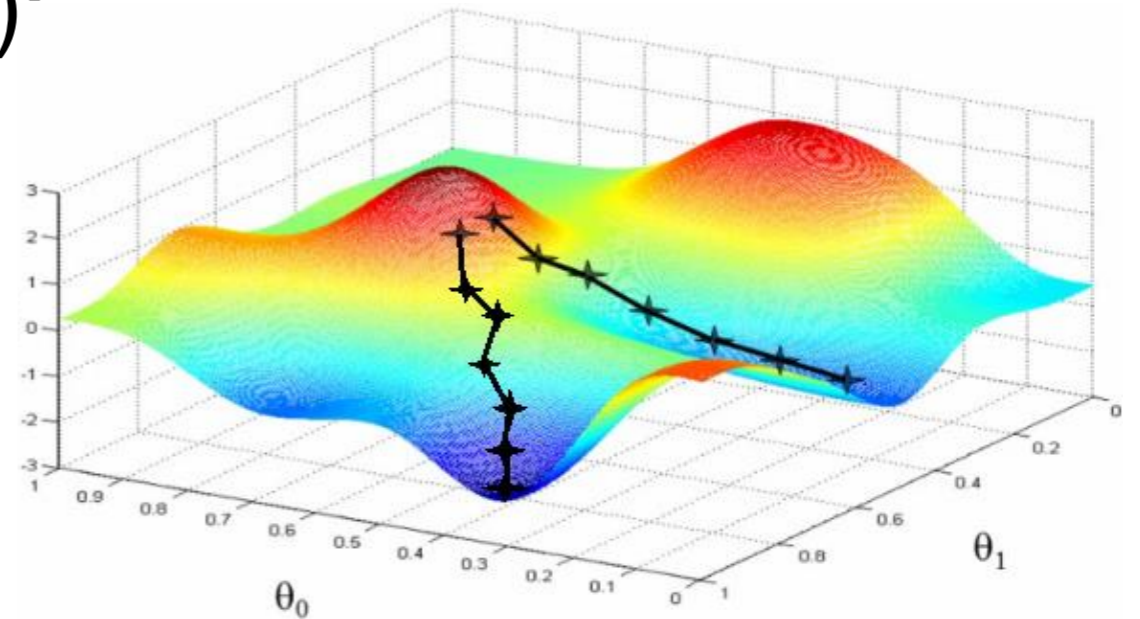
- Initialize θ randomly
- For N Epochs
- For each training example (x, y)

- Compute Loss Gradient:

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

- Update θ with update rule:

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



Why is it **Stochastic** Gradient Descent?

- Initialize θ randomly
- For N Epochs
- For each training example (x, y)

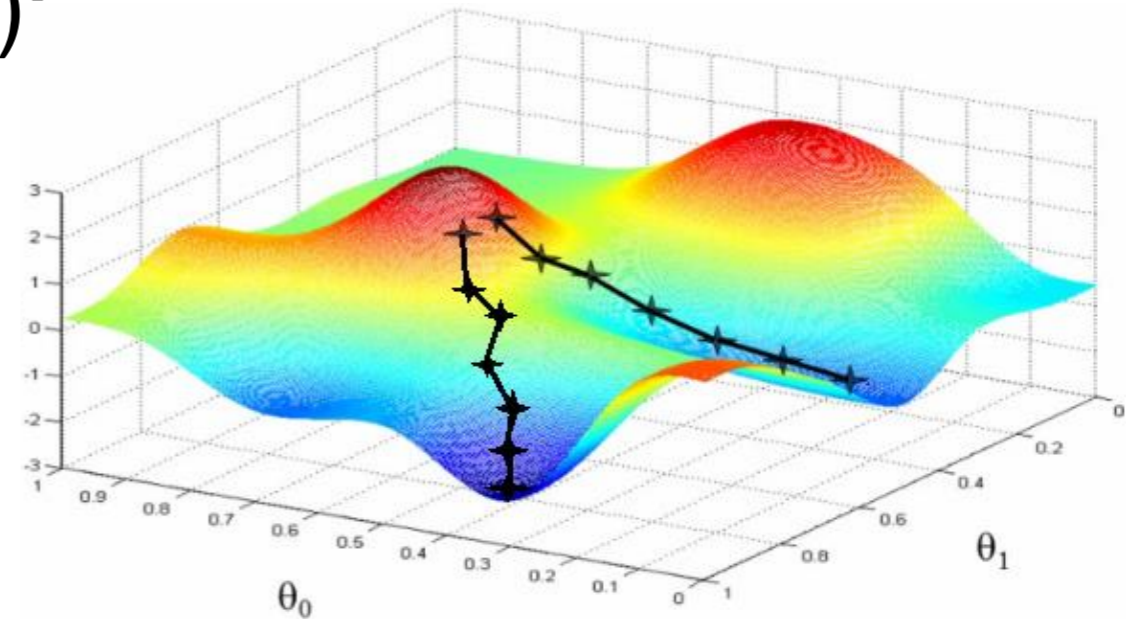
- Compute Loss Gradient:

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

← Only an estimate of true gradient!

- Update θ 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 θ randomly
- For N Epochs
- For each training example (x, y)

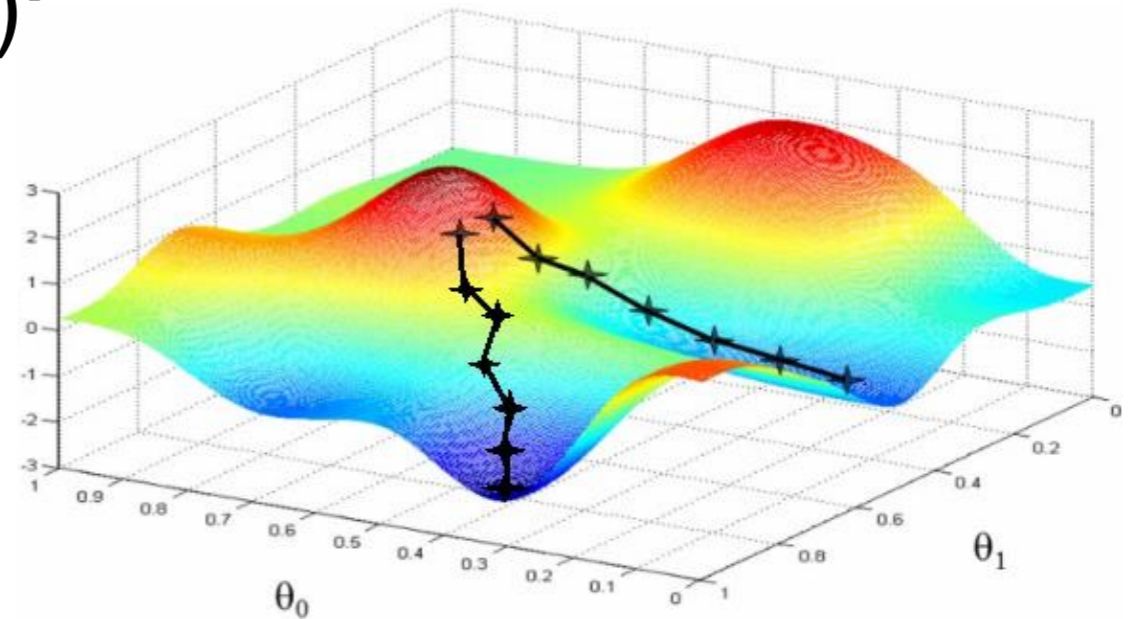
- Compute Loss Gradient:

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

More accurate estimate!

- Update θ 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

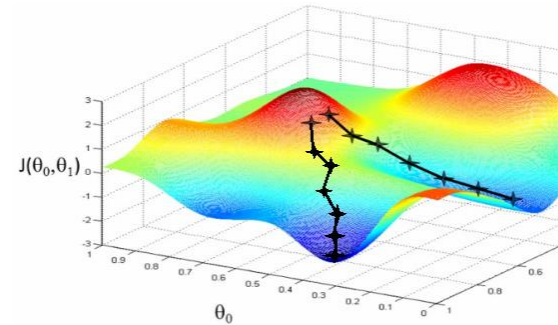
$$\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta)$$

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

Training epoch

=

Iteration over **all** examples



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

- the loss function $l(\mathbf{f}(\mathbf{x}^{(t)}; \theta), y^{(t)})$

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

- 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

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

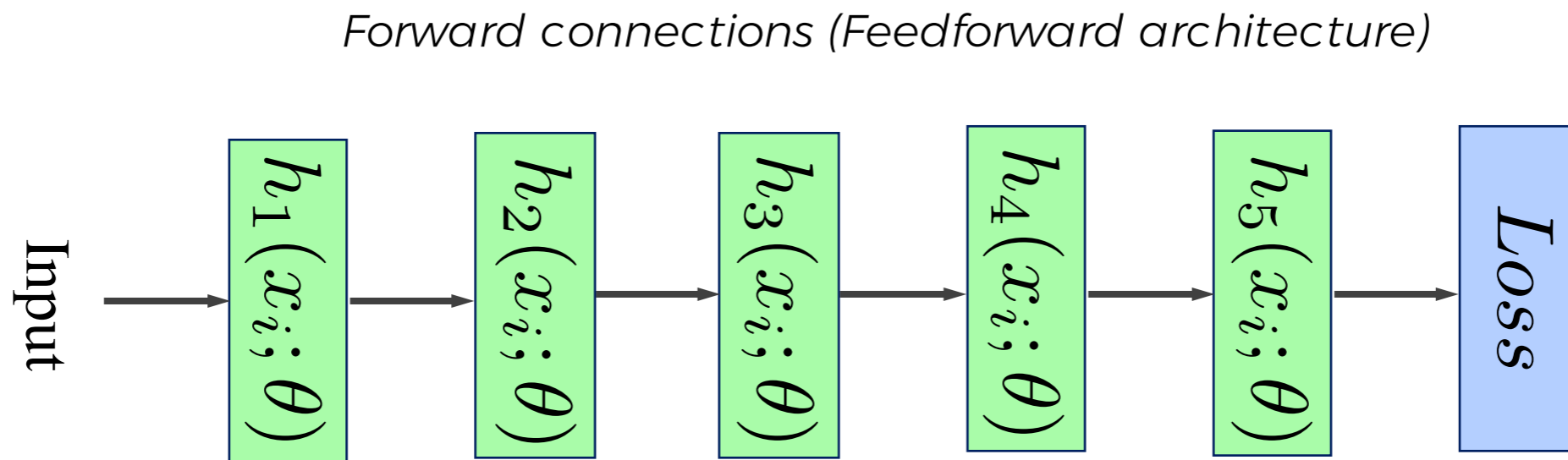
- x : input, θ_l : parameters for layer l , $a_l = h_l(x, \theta_l)$: (non-)linear function

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

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

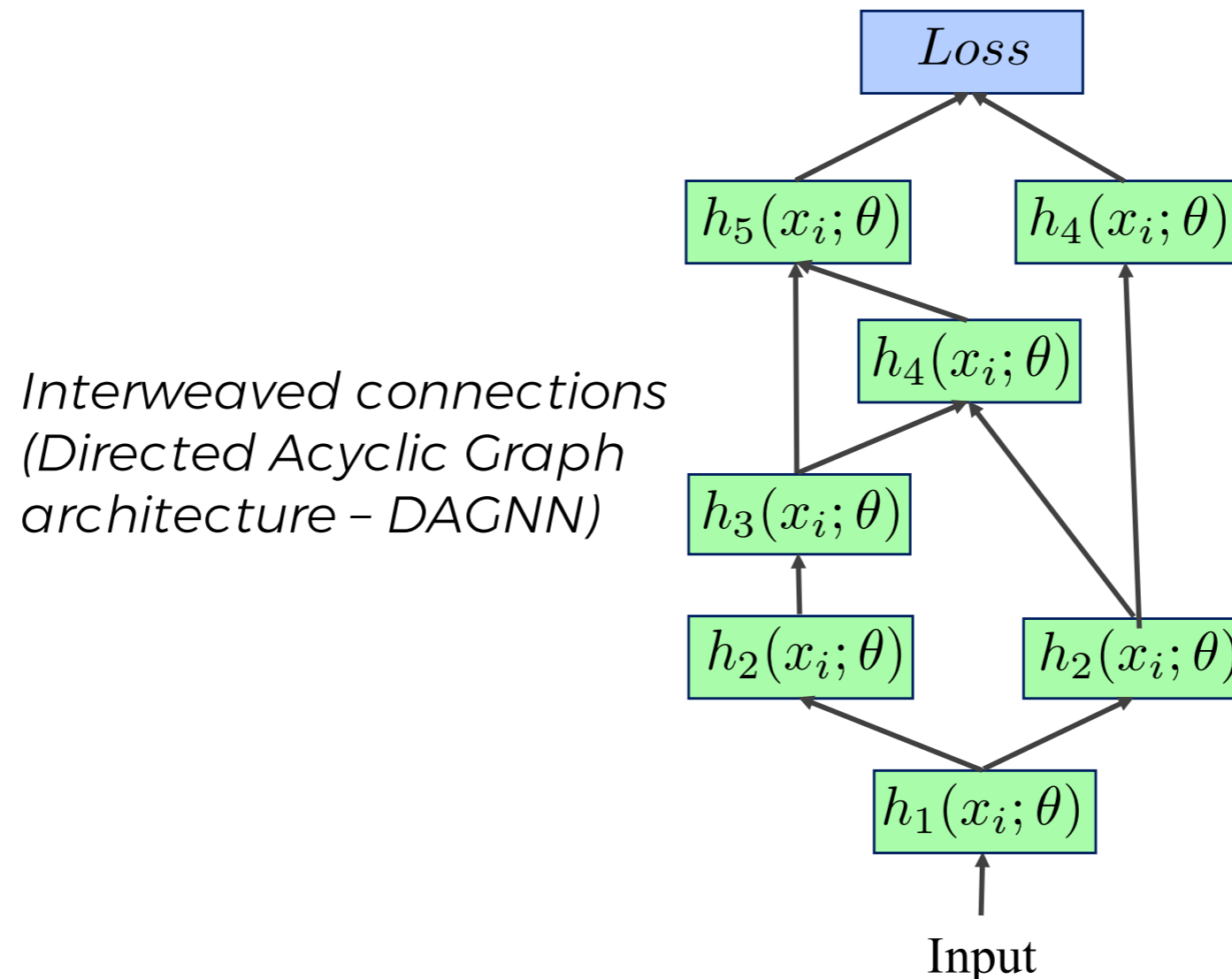
Neural network models

- 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



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, θ_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
- we need the gradients $\left(\theta^{t+1} = \theta^t - \eta_t \frac{\partial \mathcal{L}}{\partial \theta^t} \right)$
 $\frac{\partial \mathcal{L}}{\partial \theta^l}, l = 1, \dots, L$
- 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 $y = g(x)$

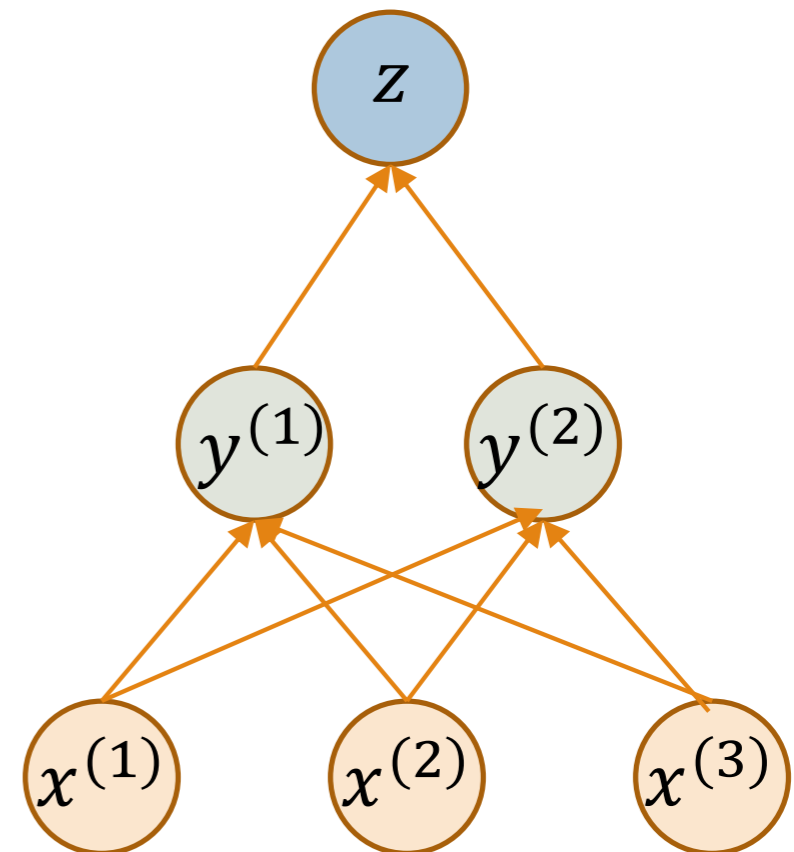
- Chain Rule for scalars x, y, z

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

- 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 $y = g(x)$

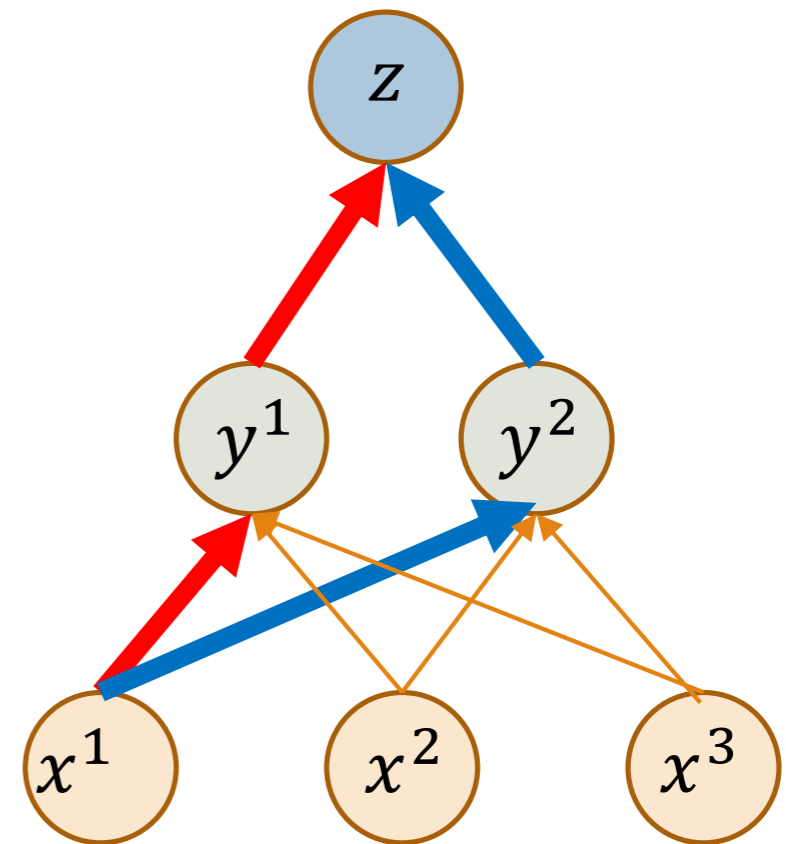
- Chain Rule for scalars x, y, z

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

- 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 $y = g(x)$

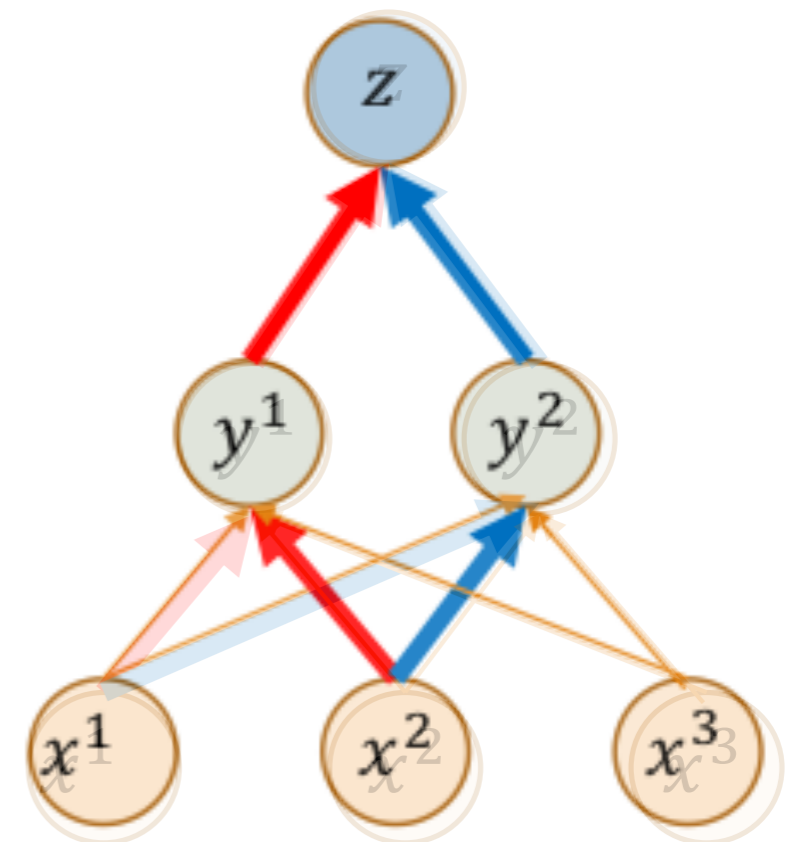
- Chain Rule for scalars x, y, z

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

- 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 $y = g(x)$

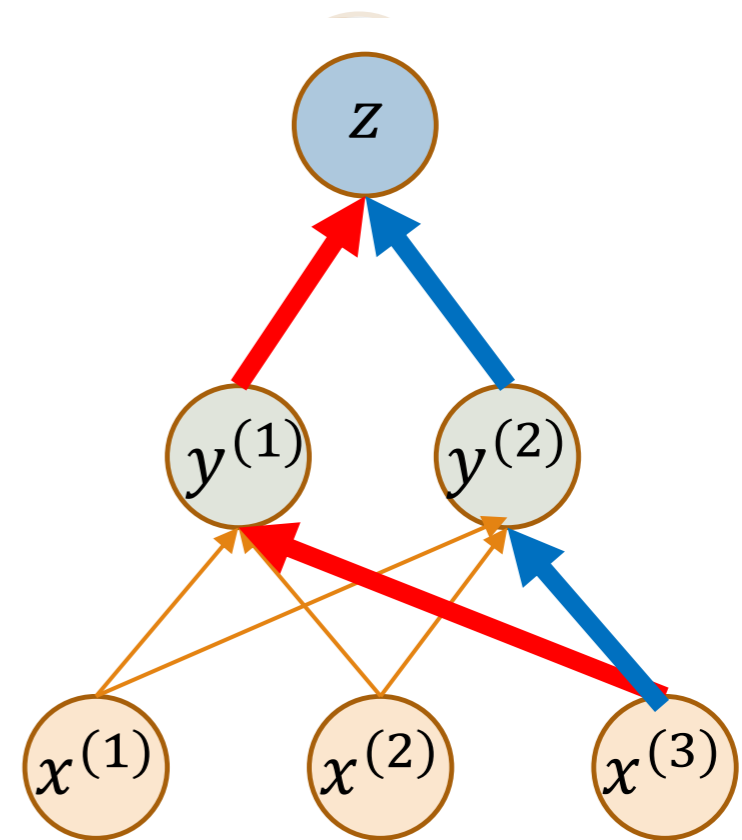
- Chain Rule for scalars x, y, z

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

- 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 \iff Chain rule!!!

- The loss function $\mathcal{L}(y, a_L)$ depends on a_L , which depends on a_{L-1} , ..., which depends on a_2 :

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

- 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 two quantities

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

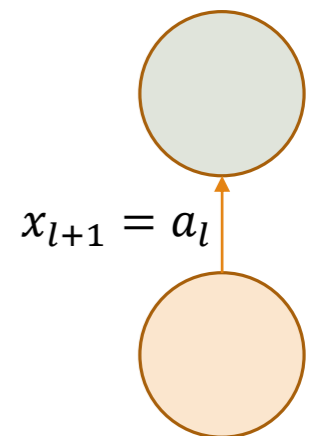
Gradient of a module w.r.t. its parameters Gradient of loss w.r.t. the module output

Backpropagation \iff 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

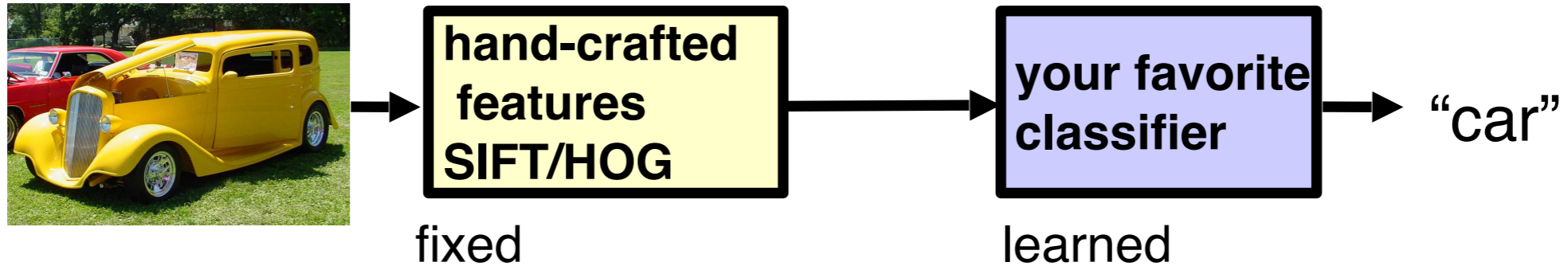
- (Hierarchical) Compositionality
- End-to-End Learning
- Distributed Representations

Three key ideas

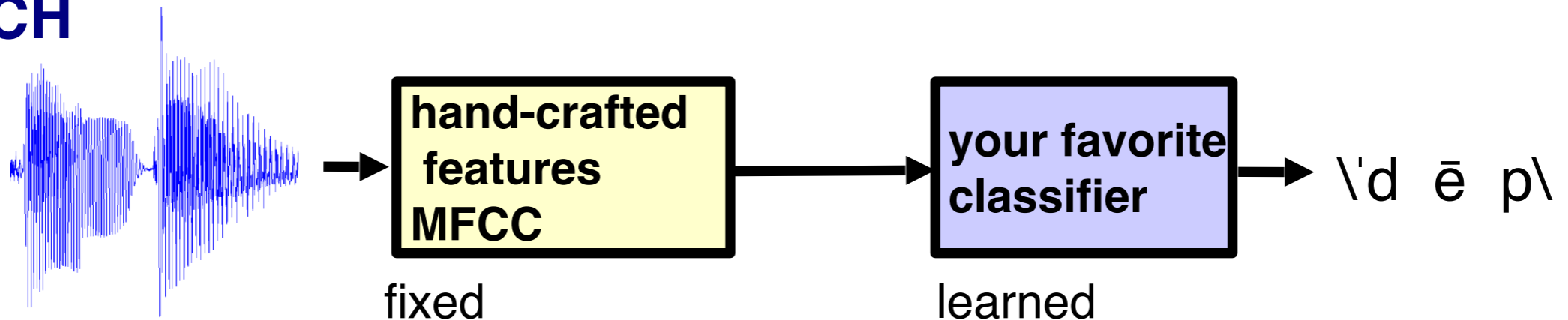
- **(Hierarchical) Compositionality**
 - Cascade of non-linear transformations
 - Multiple layers of representations
- End-to-End Learning
 - Learning (goal-driven) representations
 - Learning to feature extract
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons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

Traditional Machine Learning

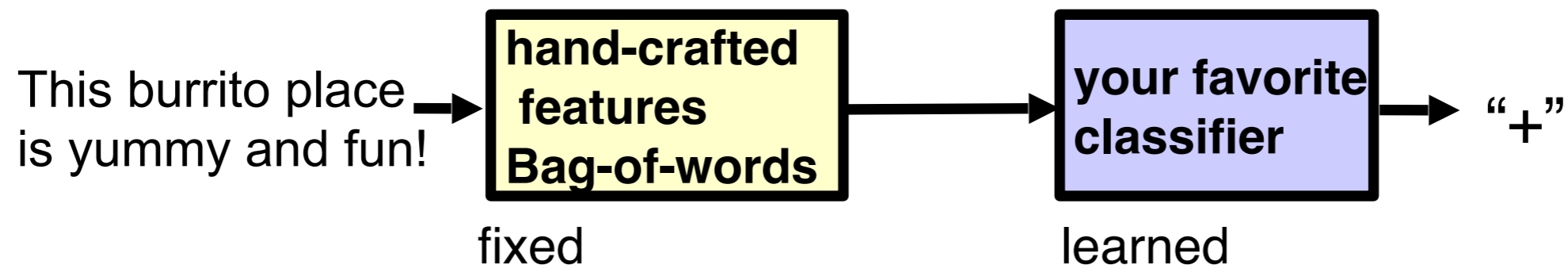
VISION



SPEECH

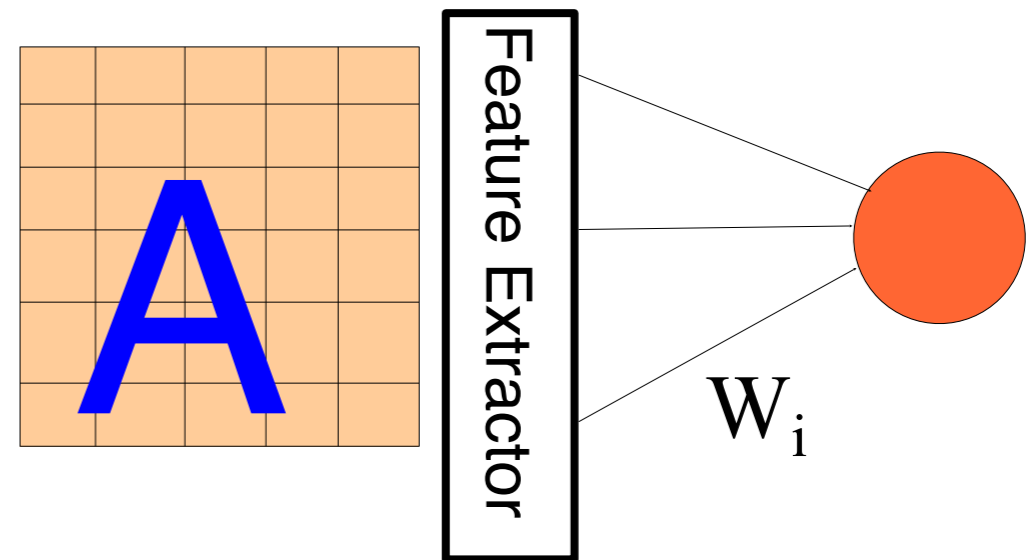


NLP

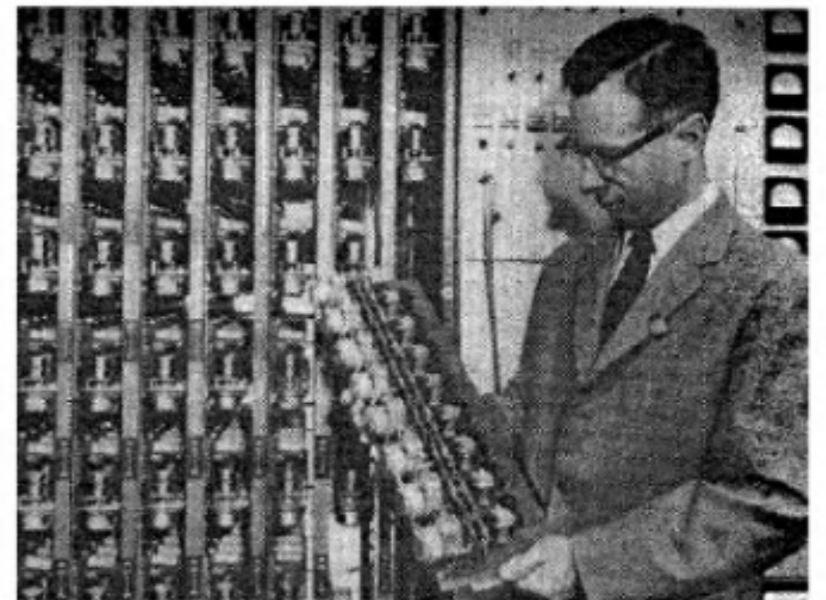
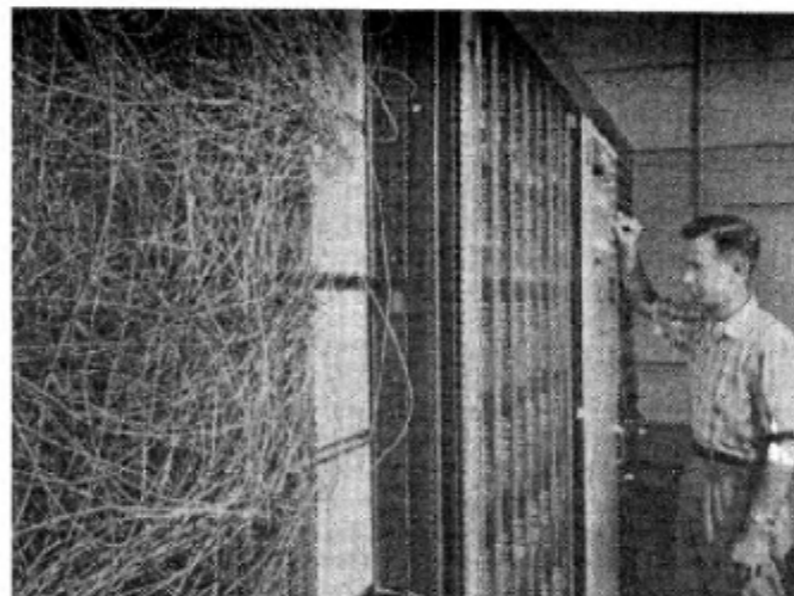
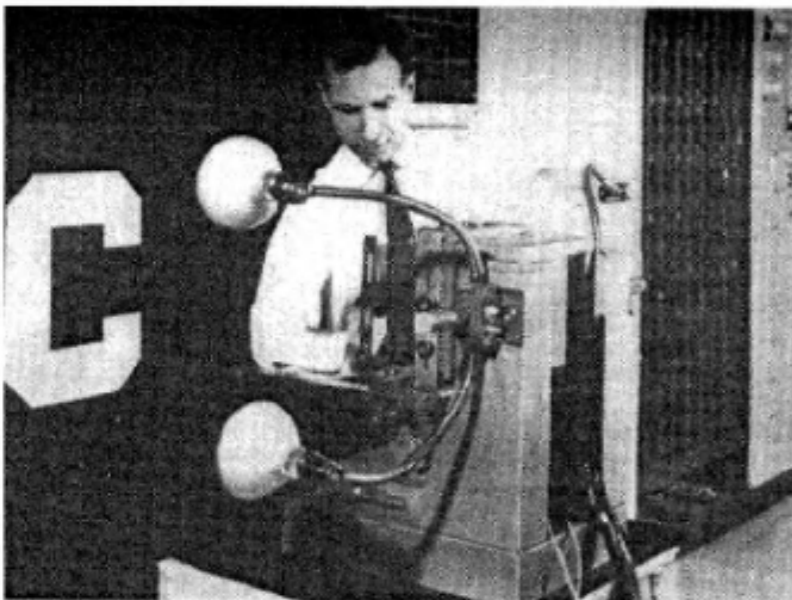


It's an old paradigm

- 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

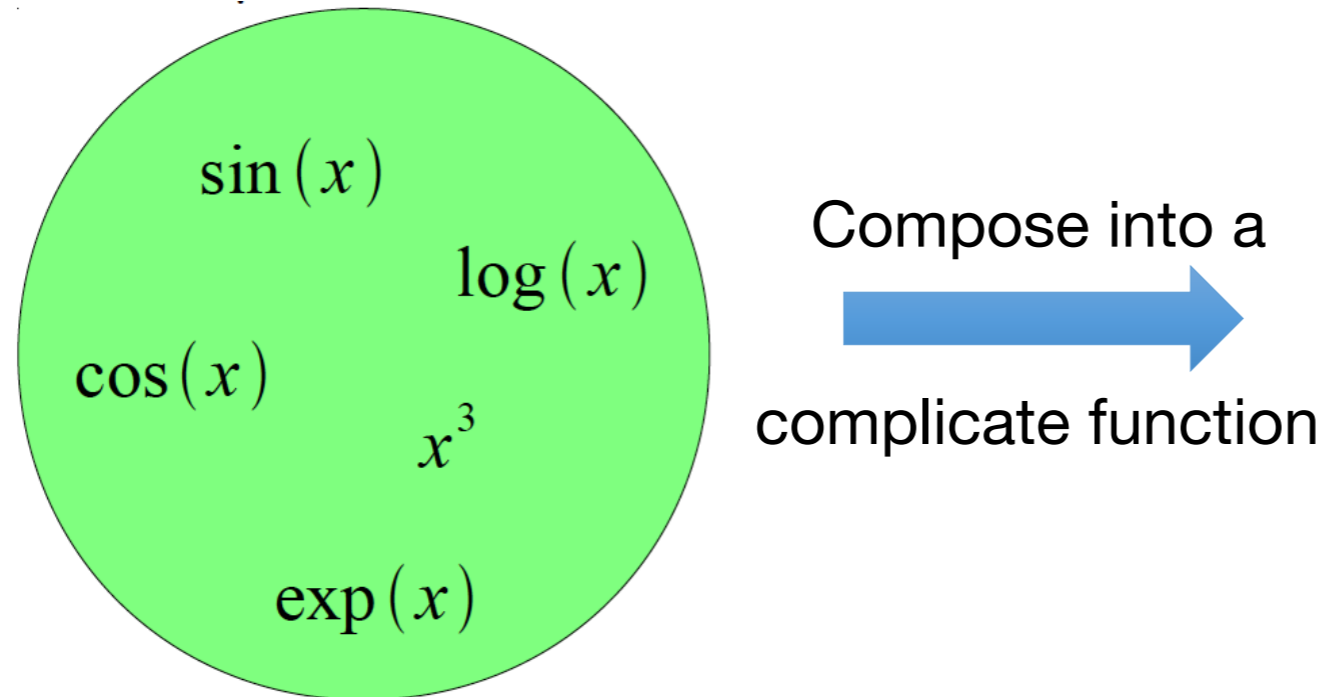
sample → spectral band → formant → motif → phone → word

NLP

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

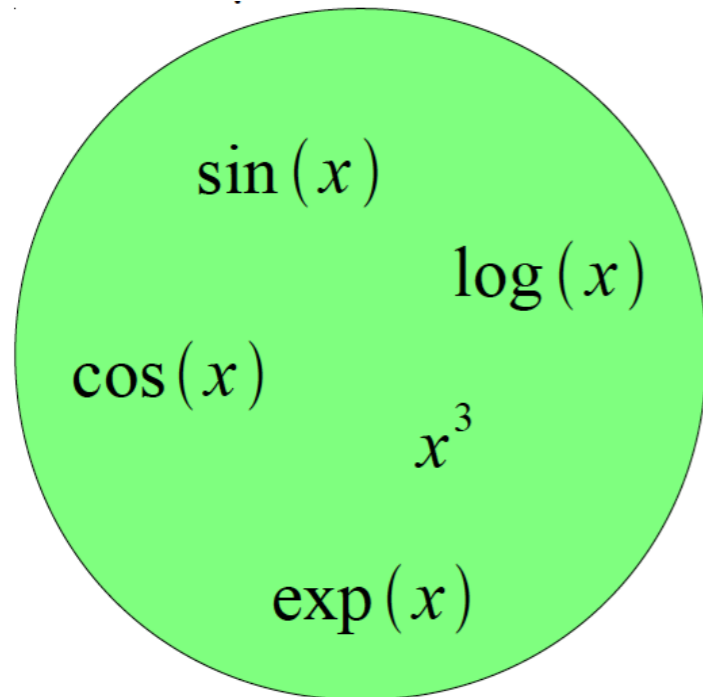
Building A Complicated Function

Given a library of simple functions



Building A Complicated Function

Given a library of simple functions

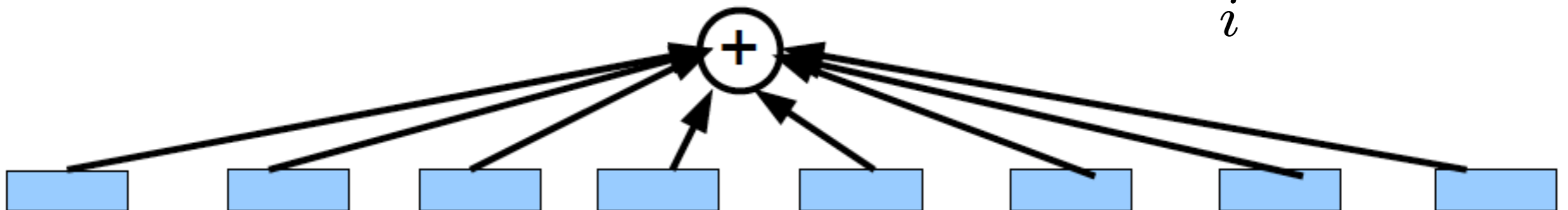


Compose into a
→
complicate function

Idea 1: Linear Combinations

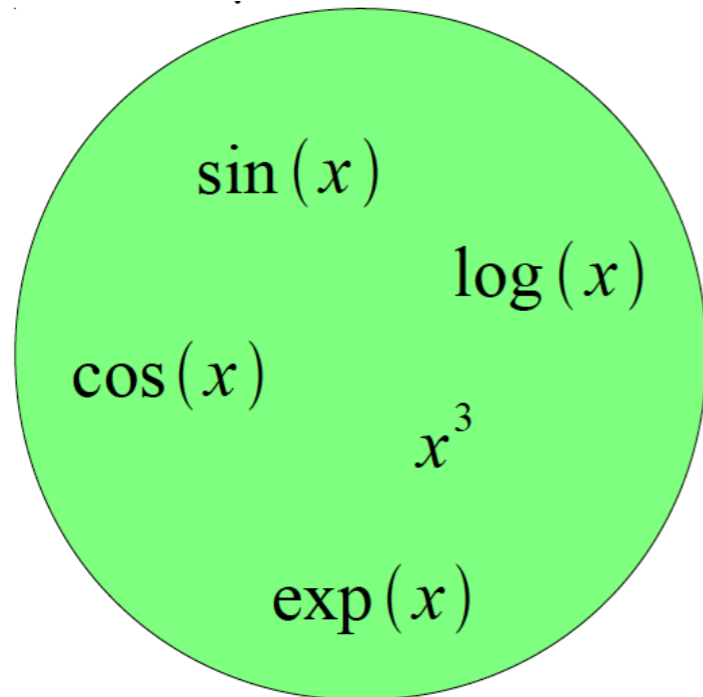
- Boosting
- Kernels
- ...

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



Building A Complicated Function

Given a library of simple functions

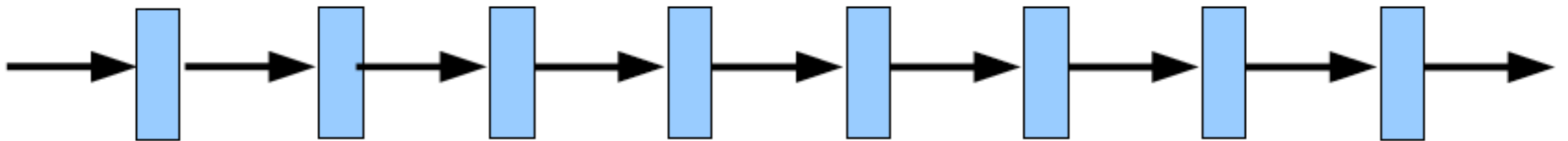


Compose into a
→
complicate function

Idea 2: Compositions

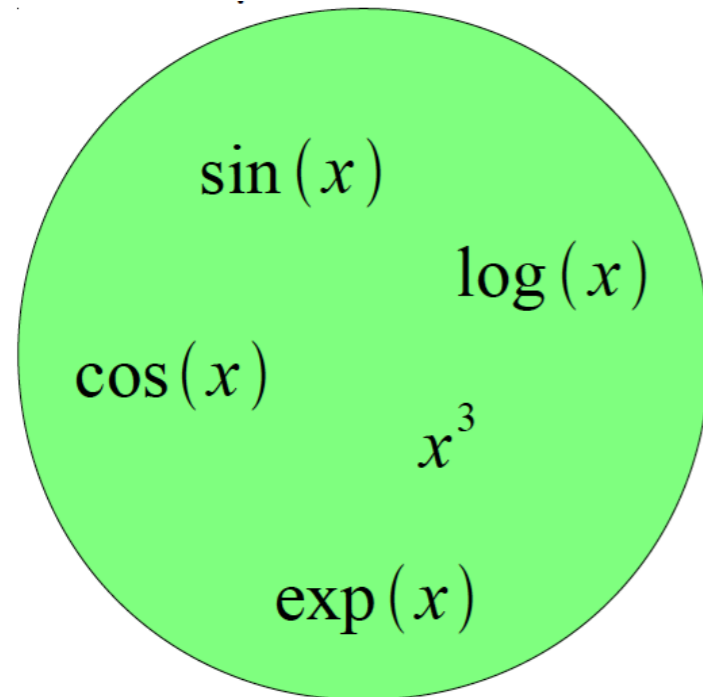
- Deep Learning
- Grammar models
- Scattering transforms...

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



Building A Complicated Function

Given a library of simple functions

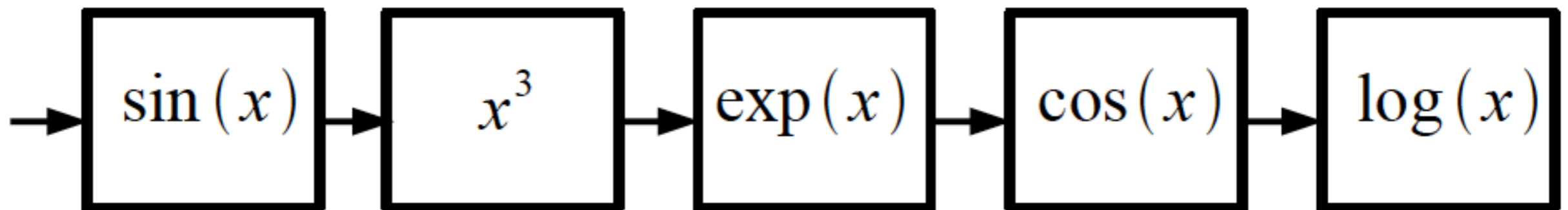


Compose into a
→
complicate function

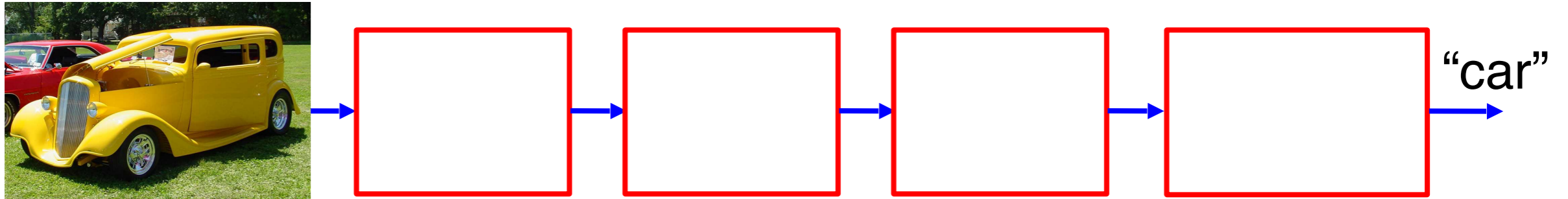
Idea 2: Compositions

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

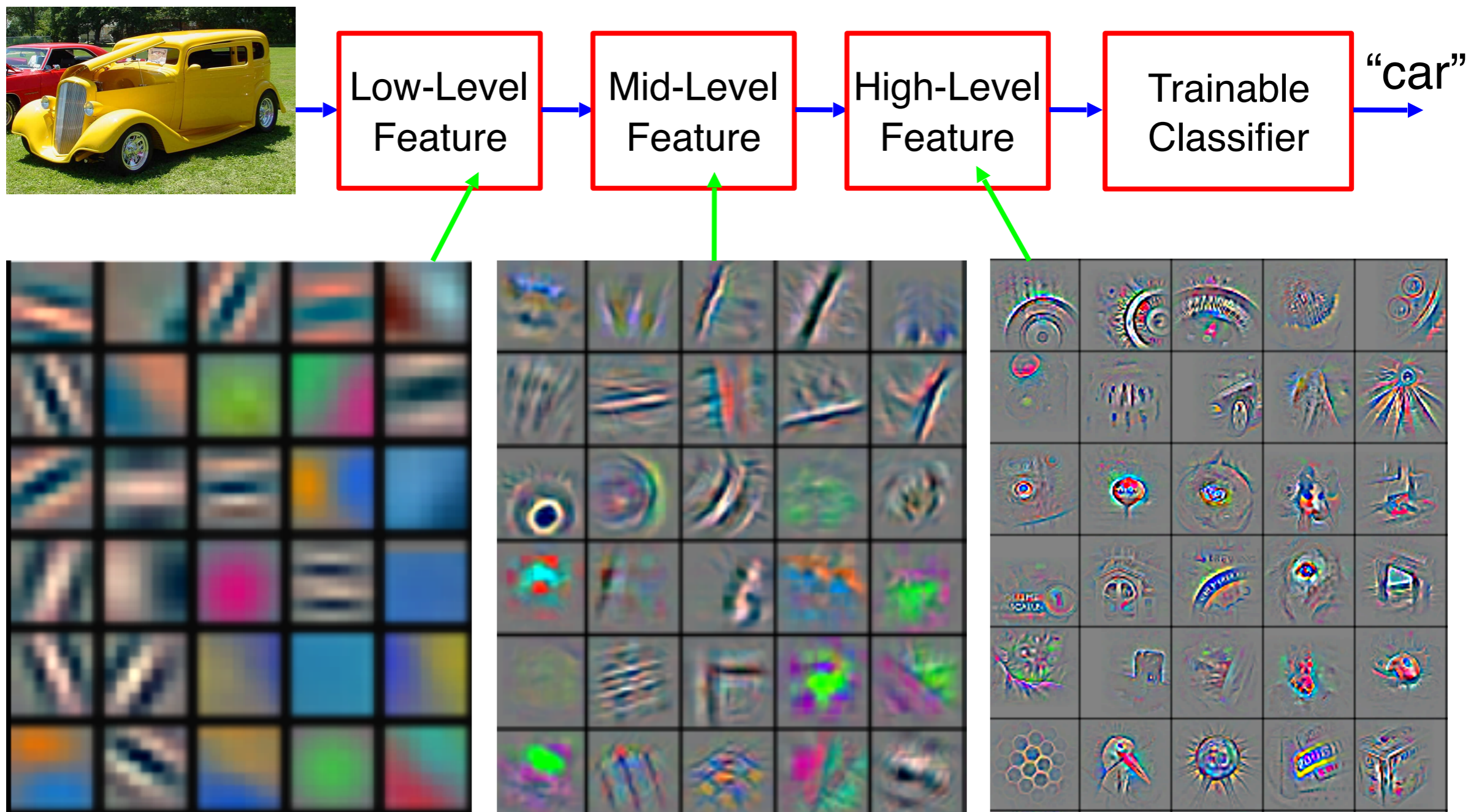
$$f(x) = \log(\cos(\exp(\sin^3(x))))$$



Deep Learning = Hierarchical Compositionality



Deep Learning = Hierarchical Compositionality



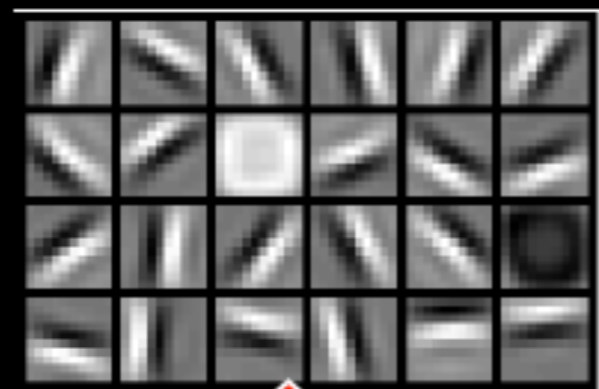
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



Face detectors



Face parts
(combination
of edges)



edges



pixels

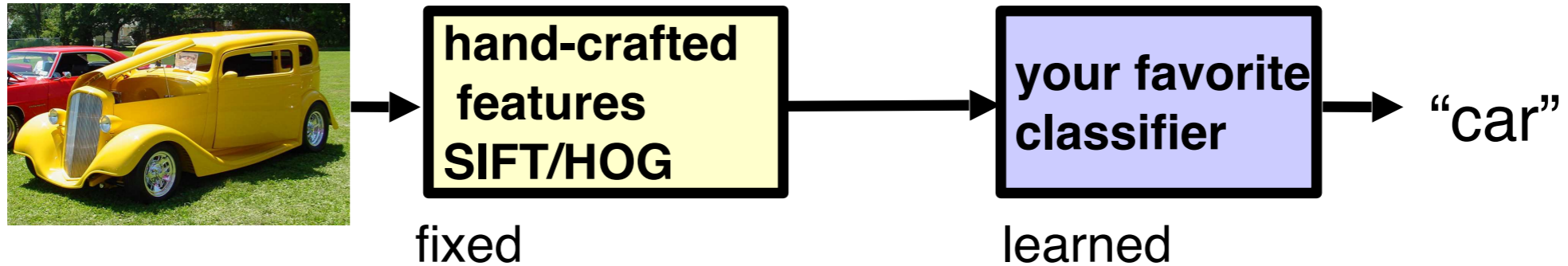
Sparse DBNs
[Lee et al. ICML '09]
Figure courtesy: Quoc Le

Three key ideas

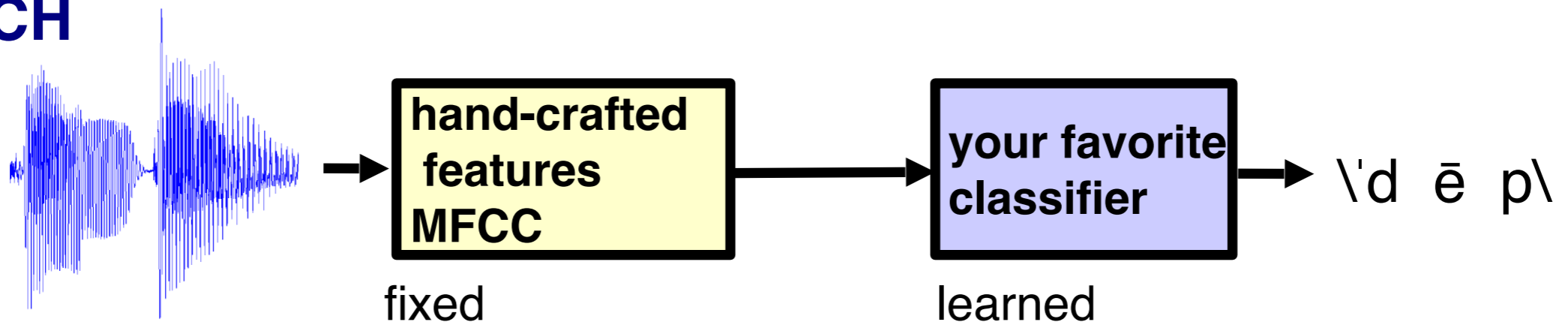
- (Hierarchical) Compositionality
 - Cascade of non-linear transformations
 - Multiple layers of representations
- **End-to-End Learning**
 - Learning (goal-driven) representations
 - Learning to feature extract
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons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

Traditional Machine Learning

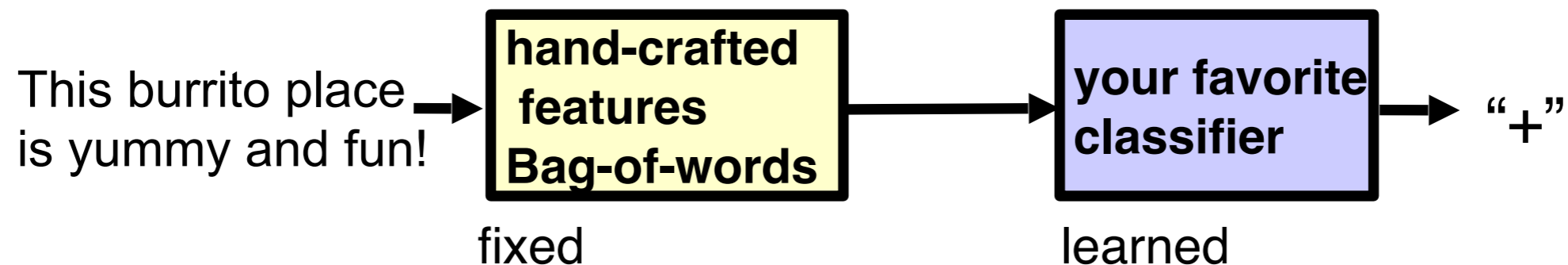
VISION



SPEECH

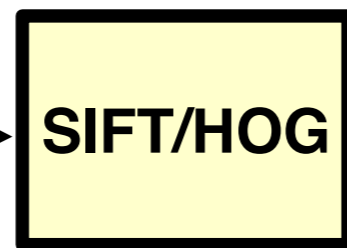


NLP

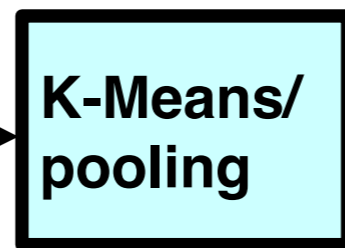


Traditional Machine Learning (more accurately)

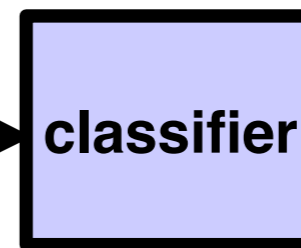
VISION



fixed



unsupervised

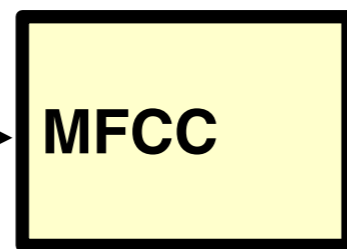
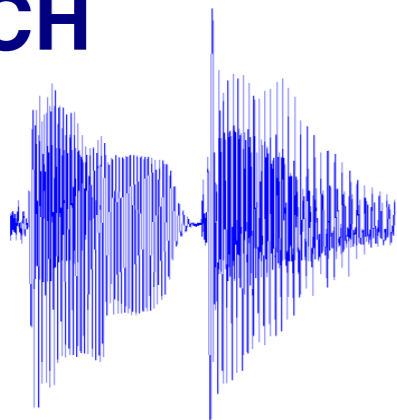


supervised

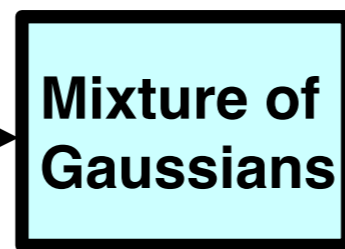
"car"

"Learned"

SPEECH



fixed



unsupervised

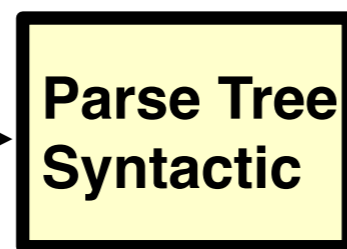


supervised

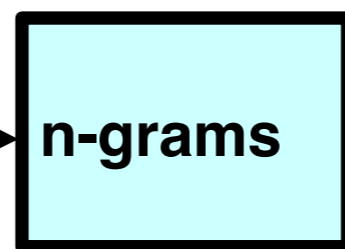
\ 'd ē p \

NLP

This burrito place
is yummy and fun!



fixed



unsupervised

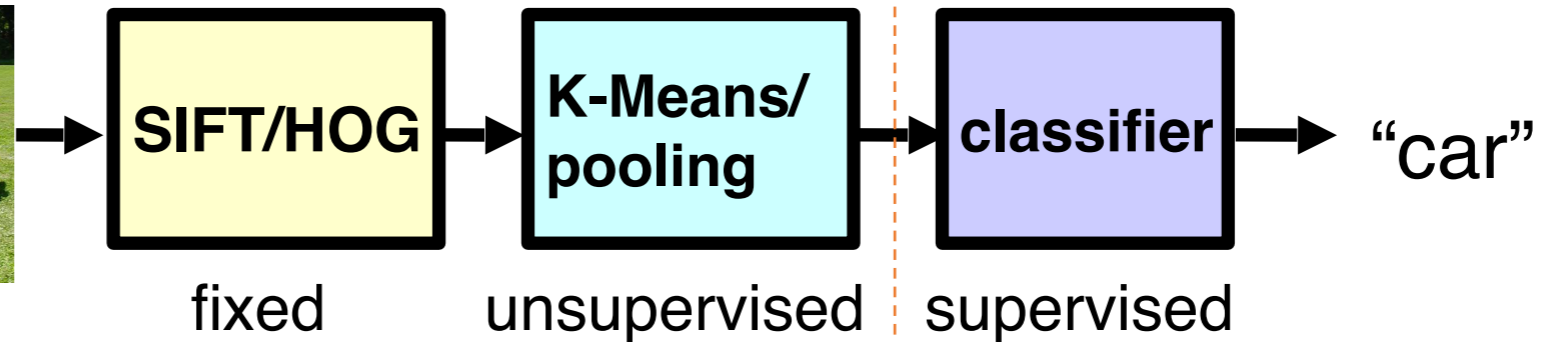


supervised

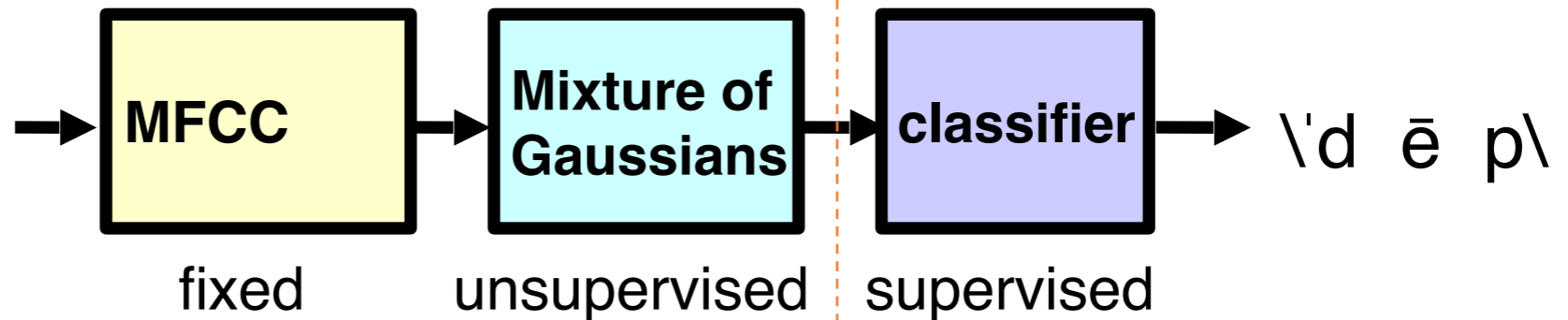
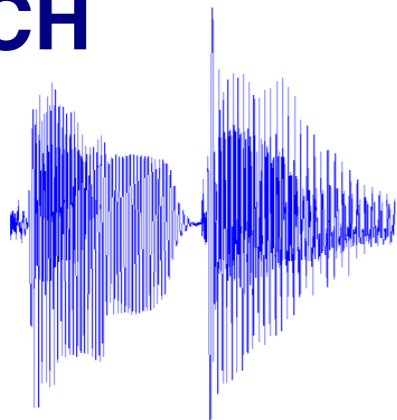
"+"

Deep Learning = End-to-End Learning

VISION

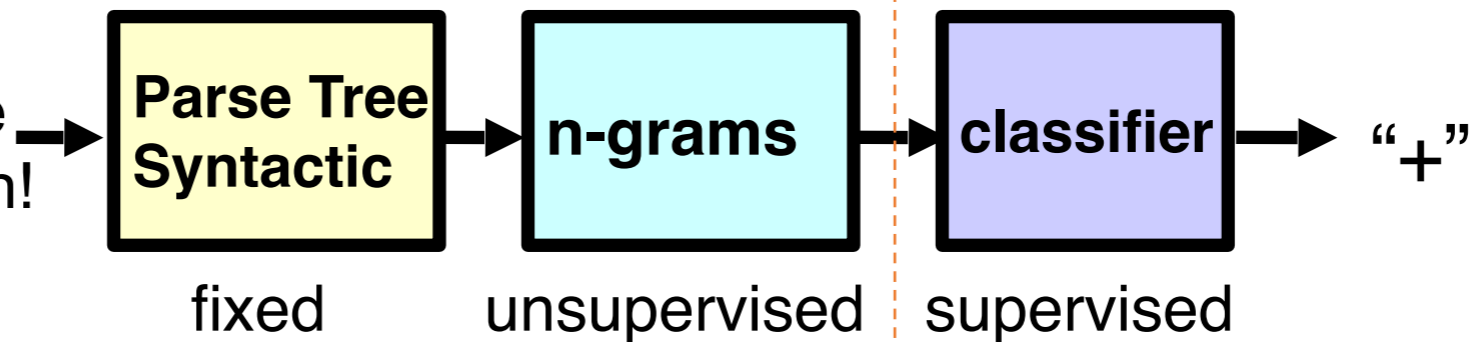


SPEECH



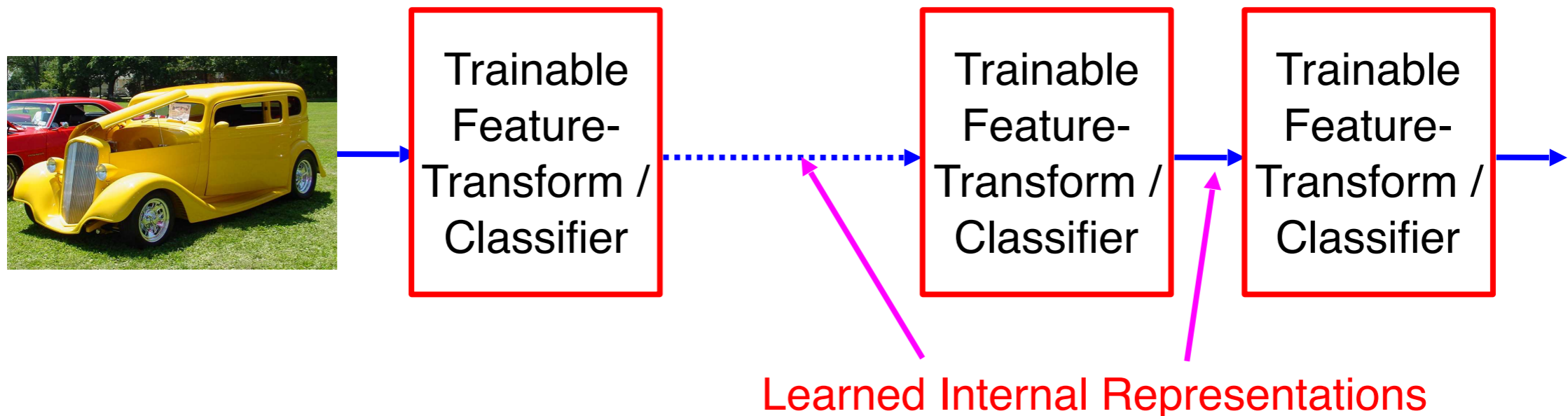
NLP

This burrito place is yummy and fun!



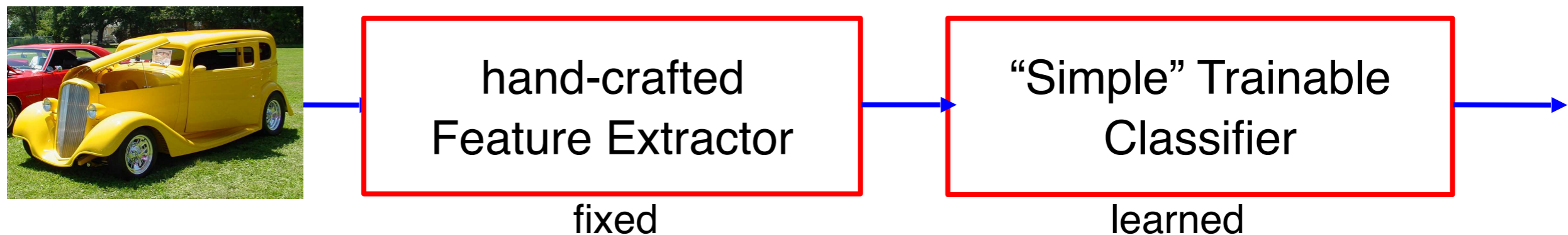
Deep Learning = End-to-End Learning

- 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

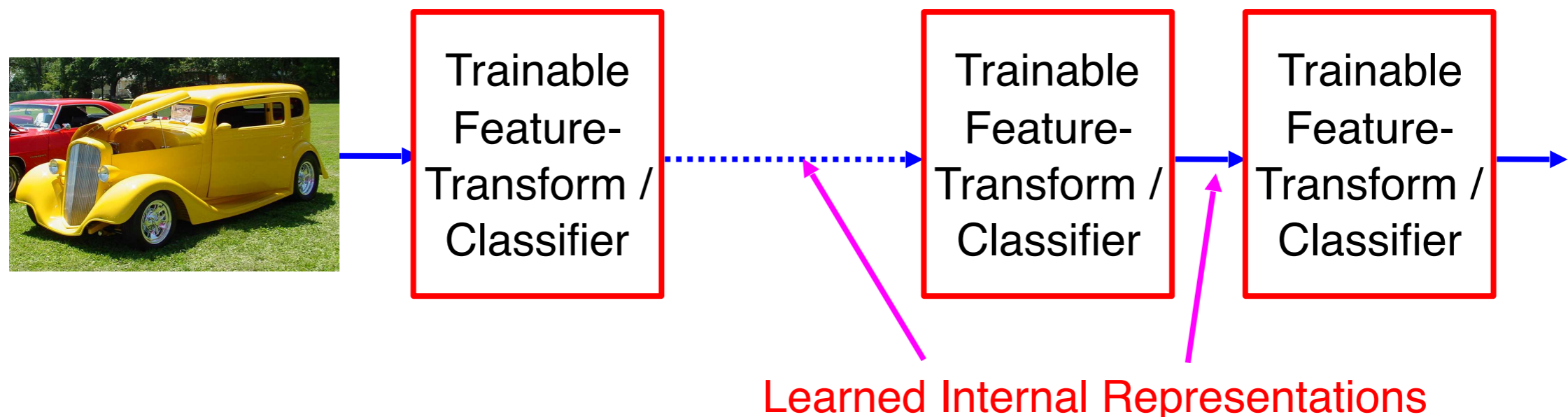


“Shallow” vs Deep Learning

- “Shallow” models



- Deep models



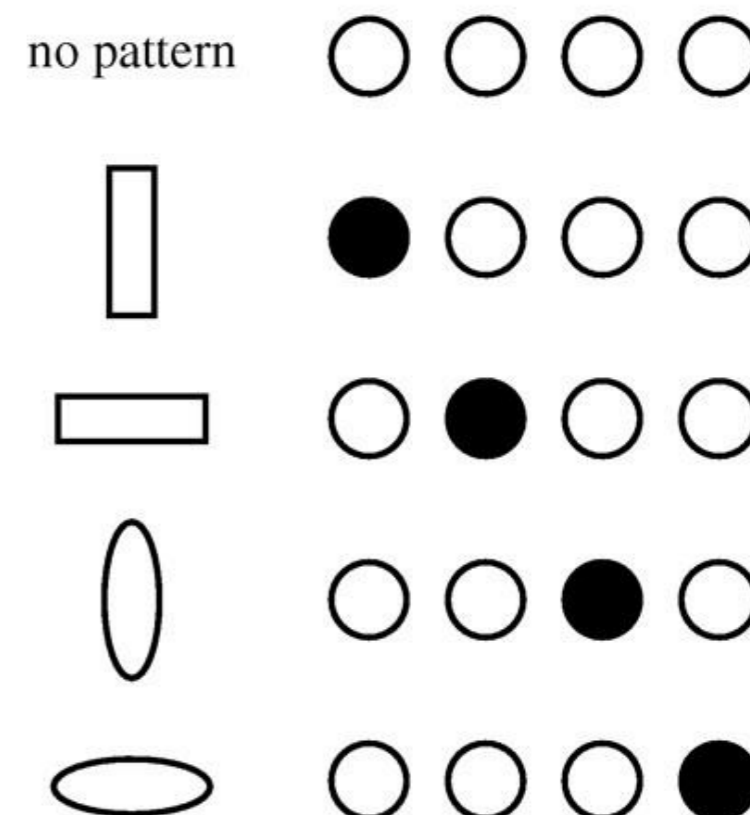
Three key ideas

- (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)

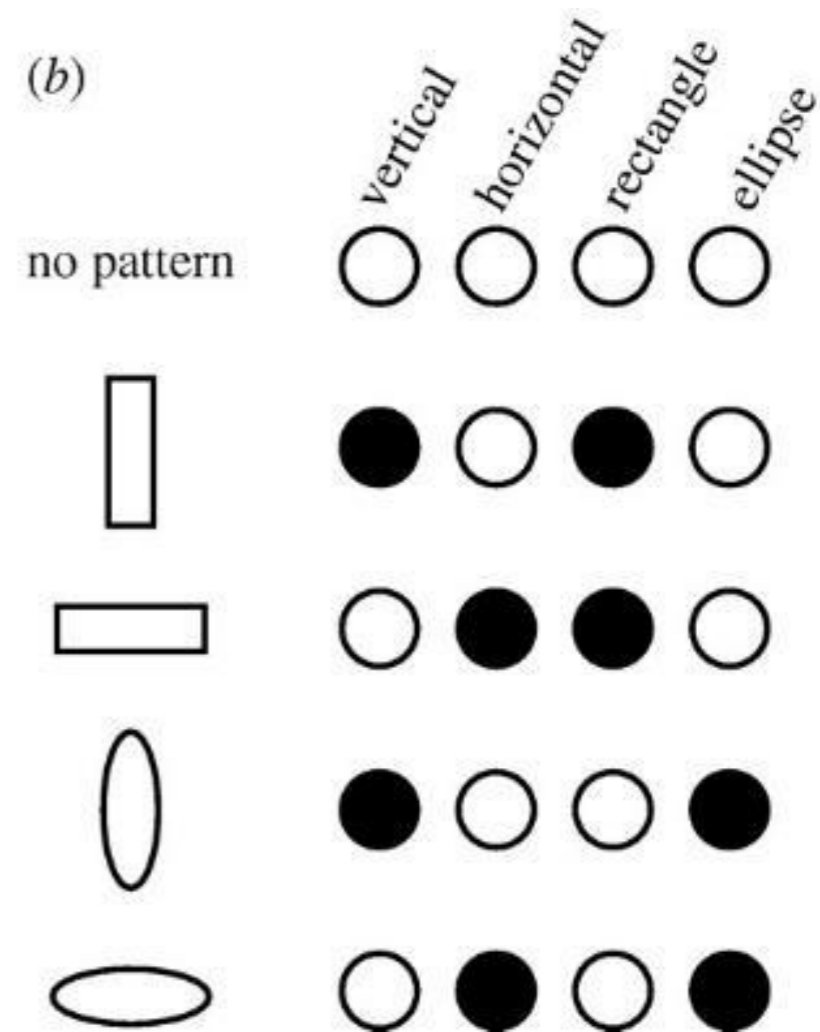
- Easy to understand.
- Easy to code by hand
 - Often used to represent inputs to a net
- Easy to learn
 - This is what mixture models do.
 - Each cluster corresponds to one neuron
- Easy to associate with other representations or responses.



- 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

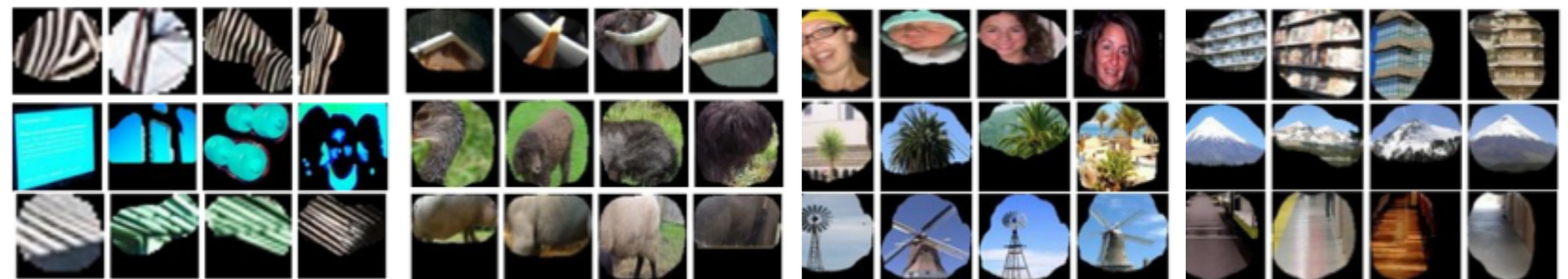


mountain



- 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:



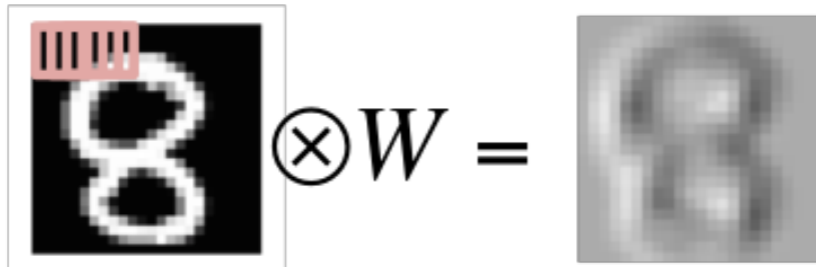
- Apply convolution to local sliding windows

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

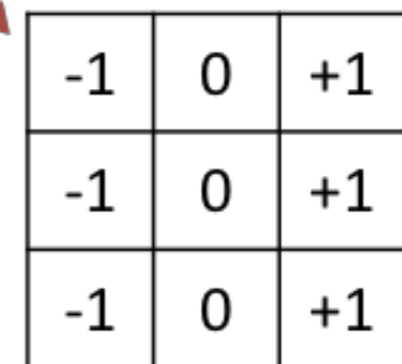
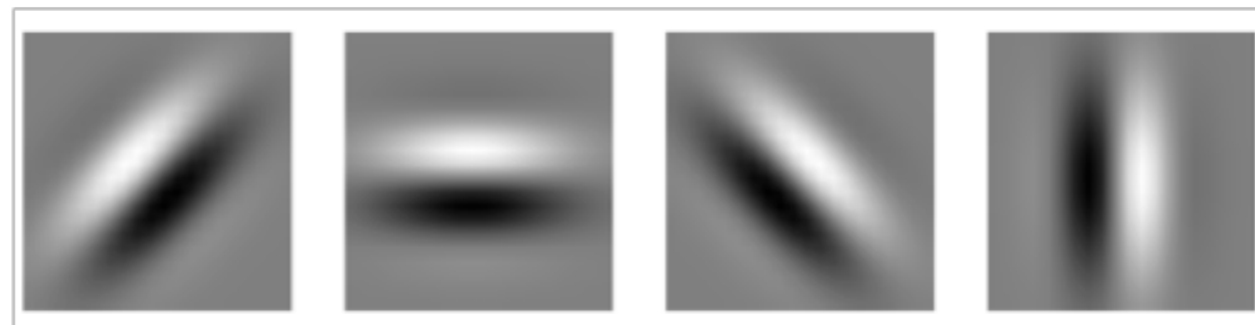
-1	0	+1
-1	0	+1
-1	0	+1

W

Gabor Filters

- Most common low-level convolutions for computer vision

Example W:



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

W

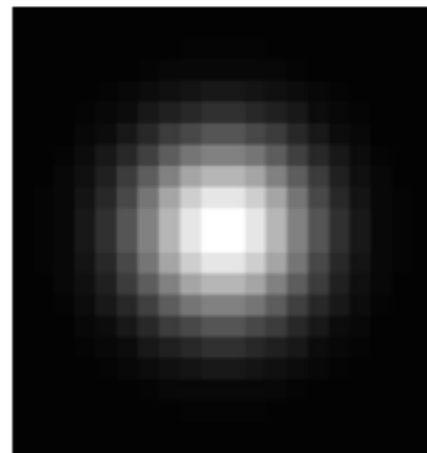
http://en.wikipedia.org/wiki/Gabor_filter

Gaussian Blur Filters

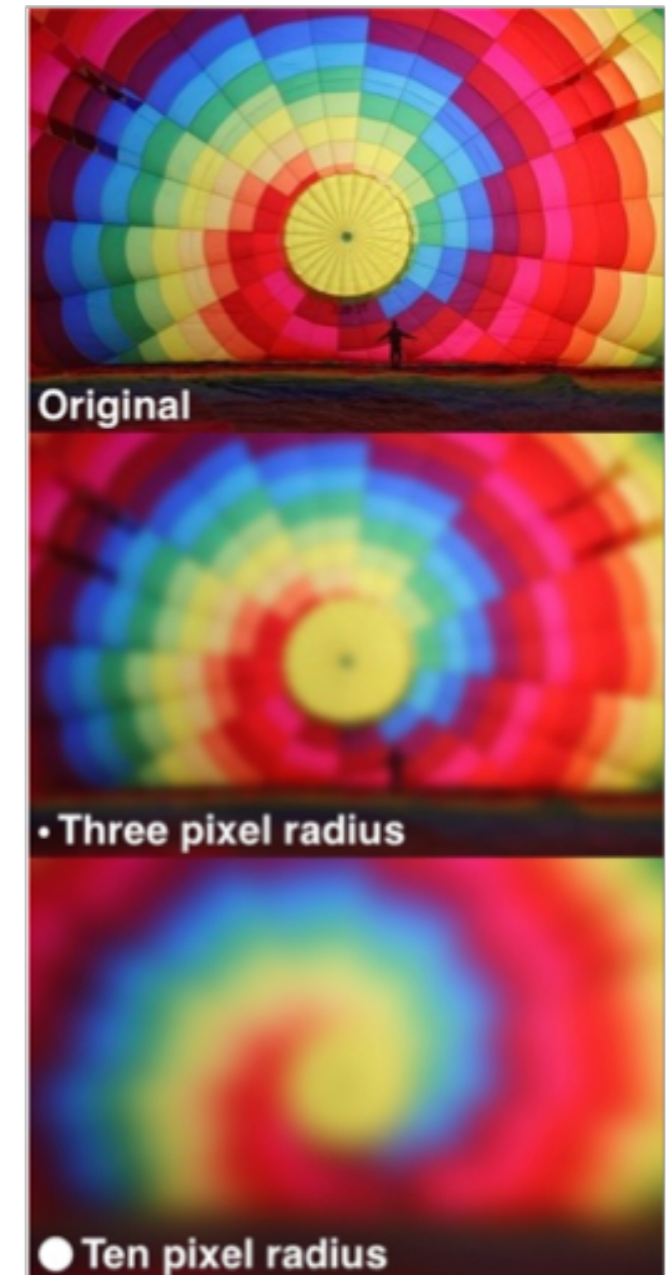
- Weights decay according to Gaussian Distribution
 - Variance term controls radius

Example W:

Apply per RGB Channel



- Black = 0
- White = Positive



http://en.wikipedia.org/wiki/Gaussian_blur