



# Tutorial on Deep Learning in Computer Vision

Aykut Erdem, Erkut Erdem, Nazlı İkizler Cinbiş, Seniha Esen Yüksel Erdem



MIT Technology Review

Forbes / Tech

APR 1, 2016 @ 06:47 AM 3,207 VIEWS

What Is Deep Learning And I



Kevin Murnane CONTRIBUTOR

I write about science, technology and the people that connect them.

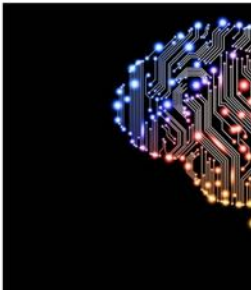


FULL BIO >

Opinions expressed by Forbes Contributors are their own.

TWEET THIS

Deep learning unlocks the treasure trove of unstructu to use it.



Credit: Google

Deep learning recently returned to the he [GOOGLE \(+1.48%\)](#) AlphaGo program crushed Lee ranking Go players in the word. Google ha learning and AlphaGo is just their latest d the news. Google's search engine, voice recognition system and self-driving cars all rely heavily on deep learning. They've used deep learning networks to build a program that picks out an attractive still from a

understand language and then make inferences and decisions on its

Edition: US ▾

Like 7M Follow

FRONT PAGE POLITICS ENTERTAINMENT WHAT'S WORKING HEALTHY LIVING WORLDOPOST HIGHLINE HUFFPOST LIVE ALL SECTIONS

THE BLOG

Artificial Intelligence, Deep Learning, Can It Take Over?

04/15/2016 11:59 am ET



Like 5



Pradeep Aradhya Entrepreneur, Humorist, Fashion Plate, Do Gooder



Bill Gates, Stephen Hawking and Elon Musk first warned us about Artificial Intelligence (AI). Elon Musk then turned around and with other technologists put \$1B into starting a nonprofit research effort - OpenAI just to "keep an eye on it"! Facebook, Google, Amazon, Nvidia, Shopify and others are charging full steam at AI and even open sourcing it! So what is all the AI ruckus about?

FOLLOW HUFFPOST



HuffPost Like 7

Business Like 13

HUFFPOST NEWSLETTERS

Get top stories and blog posts emailed to me each day. Newsletters may offer personalized content or advertisements. [Learn More.](#)

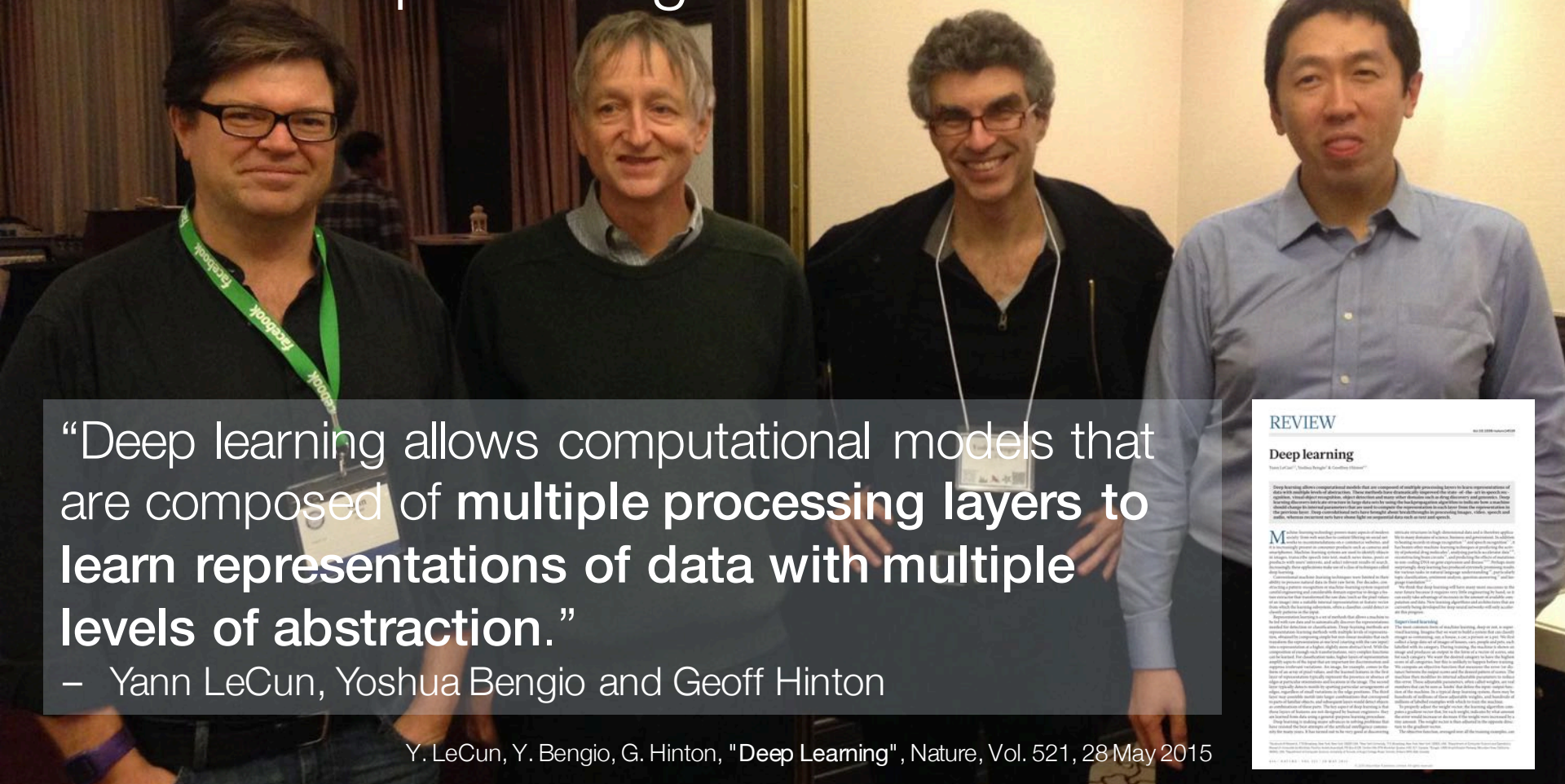
address@email.com

Subscribe!

event for many Silicon Valley companies in the last few years, thanks to the vision in AI. NIPS was where Facebook Chief Executive Officer Mark Zuckerberg > in 2013 to announce the company's plans to form an AI laboratory and where a startup named DeepMind showed off an AI that could learn to play computer games before it was acquired by Google.



# What is deep learning?



“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

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

**REVIEW**

**Deep learning**

David L. Donoho

David L. Donoho, "Deep learning", *Nature*, vol. 521, pp. 435–442, 2015.

**Summary**

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have recently captured the state-of-the-art in speech recognition, visual object recognition, image classification, and many other domains such as drug discovery and genomics. Deep learning architectures have the potential to learn the hierarchical structure of data by using the backpropagation algorithm to tune the weights of each layer. This process is often aided by stochastic gradient descent and other optimization algorithms. The representation of data is learned by the application of the algorithm, which is often aided by stochastic gradient descent and other optimization algorithms. The representation of data is learned by the application of the algorithm, which is often aided by stochastic gradient descent and other optimization algorithms.

**Introduction**

Deep learning is a class of machine learning algorithms that are composed of multiple processing layers. These algorithms are designed to learn representations of data with multiple levels of abstraction. The most common type of deep learning architecture is the feedforward neural network, which consists of an input layer, one or more hidden layers, and an output layer. The weights of the connections between the layers are adjusted during training to minimize the error between the predicted and actual outputs. This process is often aided by stochastic gradient descent and other optimization algorithms.

**Applications**

Deep learning has found applications in a wide range of domains, including speech recognition, image classification, and drug discovery. In speech recognition, deep learning models have achieved state-of-the-art performance by learning to represent the hierarchical structure of speech sounds. In image classification, deep learning models have achieved state-of-the-art performance by learning to represent the hierarchical structure of visual objects. In drug discovery, deep learning models have been used to predict the activity of new compounds and to identify potential drug targets.

**Conclusion**

Deep learning is a powerful tool for learning representations of data with multiple levels of abstraction. It has found applications in a wide range of domains and is expected to continue to play a major role in the development of artificial intelligence.

# Tutorial objectives

- Basics of training deep neural networks
- Good understanding of Convolutional and Recurrent Networks
- A short overview about the future of deep learning
  
- Focus will especially be on computer vision applications
- We expect basic knowledge of machine learning and/or computer vision

# Agenda

- Part I: History and Motivations
- Part II: Training Neural Networks

————— *A short break* —————

- Part III: Convolutional Neural Networks (ConvNets)
- Part IV: Recurrent Neural Networks (RNNs)
- Part V: Concluding remarks

# History and Motivations

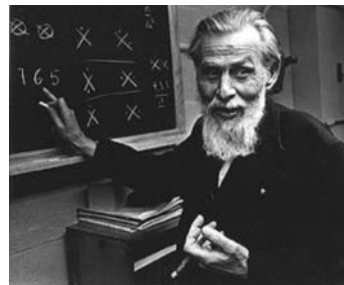


# 1943 – 2006: A Prehistory of Deep Learning



# 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 a 1 if the sum exceeds a certain threshold value, and otherwise outputs a 0



*Bulletin of Mathematical Biology*, Vol. 12, No. 1-2, pp. 165-193, 1946  
Printed in Great Britain

WCC-KOENIGS-1-100  
Program Paper 20  
Library for Mathematical Biology

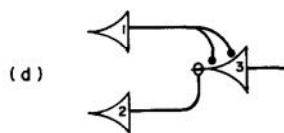
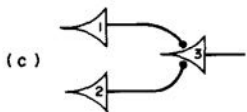
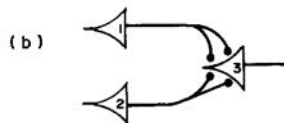
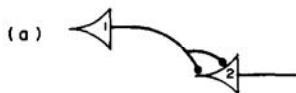
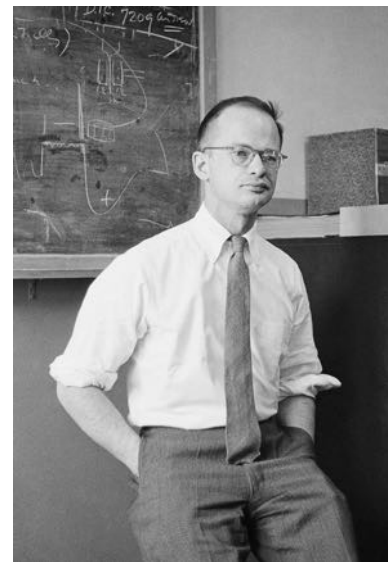
## A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY\*

• WARREN S. McCULLOCH AND WALTER PITTS  
University of Illinois, College of Medicine,  
Department of Psychiatry at the Illinois Neuropsychiatric Institute,  
University of Chicago, Chicago, U.S.A.

Because of the "all-or-none" character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every system can be described in these terms, with the addition of more complicated logical means for nets containing cycles, and that for any logical expression satisfying certain conditions, one can find a net behaving in the fashion it describes. It is shown that many particular choices among possible neurophysiological assumptions are equivalent, in the sense that for every net behaving under one assumption, there exists another net which behaves under the other and gives the same results, although perhaps not in the same time. Various applications of the calculus are discussed.

1. Introduction. Theoretical neurophysiology rests on certain cardinal assumptions. The nervous system is a set of neurons, each having a soma and an axon. Their adjunctions, or synapses, are always between the axon of one neuron and the soma of another. At any instant a neuron has some threshold, which excitation must exceed to initiate an impulse. This, except for the fact and the time of its occurrence, is determined by the neuron, not by the excitation. From the point of excitation the impulse is propagated to all parts of the neuron. The velocity along the axon varies directly with its diameter, from  $<1 \text{ ms}^{-1}$  in thin axons, which are usually short, to  $>150 \text{ ms}^{-1}$  in thick axons, which are usually long. The time for axonal conduction is consequently of little importance in determining the time of arrival of impulses at points unequally remote from the same source. Excitation across synapses occurs predominantly from axonal terminations to somata. It is still a moot point whether this depends upon irreversibility of individual synapses or merely upon prevalent anatomical configurations. To suppress the latter requires no hypothesis *ad hoc* and explains known exceptions, but any assumption as to cause is compatible with the calculus so long. No case is known in which excitation through a single synapse has elicited a nervous impulse in any neuron, whereas any neuron may be excited by impulses arriving at a sufficient number of neighboring synapses within the period of latent addition, which lasts  $<0.25 \text{ ms}$ . Observed temporal summation of impulses at greater intervals

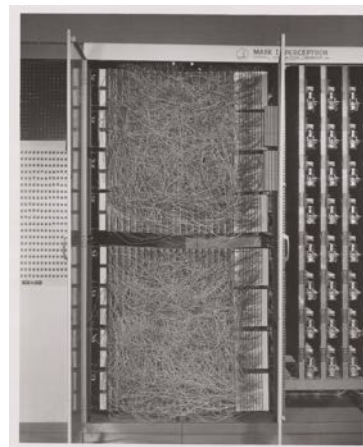
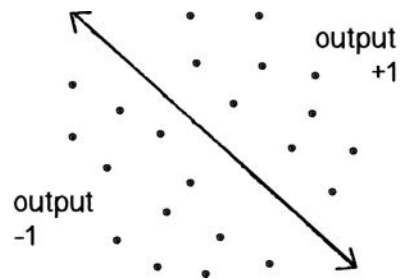
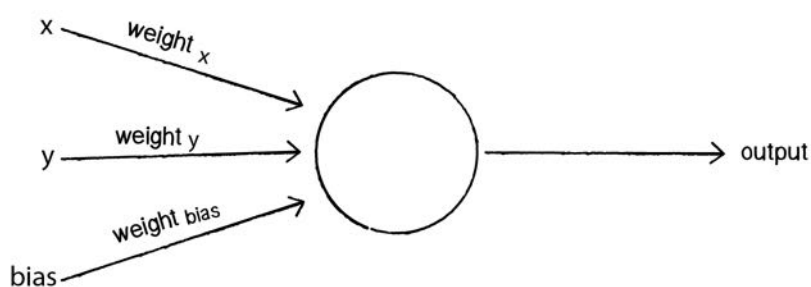
\* Reprinted from the *Bulletin of Mathematical Biophysics*, Vol. 5, pp. 153-157 (1943).





# 1958: Frank Rosenblatt's Perceptron

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

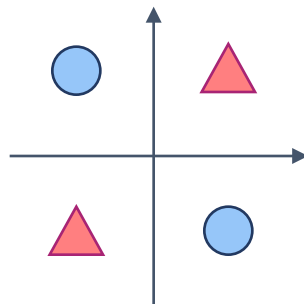


# 1969: Marvin Minsky and Seymour Papert

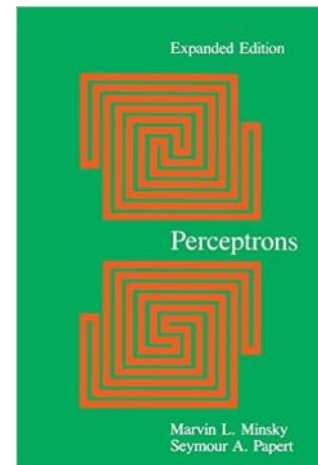
“No machine can learn to recognize X unless it possesses, at least potentially, some scheme for representing X.” (p. xiii)



- Perceptrons can only represent linearly separable functions.
  - such as **XOR** Problem

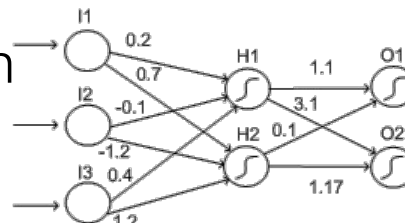


- Wrongly attributed as the reason behind the **AI winter**, a period of reduced funding and interest in AI research



# 1990s

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

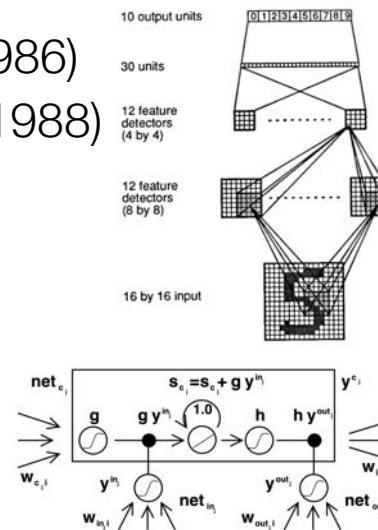


- Training multi-layer perceptrons

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

- 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 through time (BPTT) is a method for training recurrent neural networks (RNNs) using the backpropagation algorithm. It is a generalization of the backpropagation algorithm for feedforward networks to RNNs. The basic idea is to propagate the error gradients backwards through time as well as through the network layers. This allows the network to learn from its own previous actions, which is essential for learning tasks that require temporal dependencies, such as speech recognition, machine translation, and video processing.

© 1990 MIT Press. All rights reserved.

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

ARTICLE INFORMATION

# 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



# 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



## Image classification

### Easiest classes



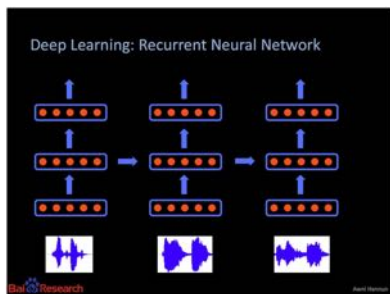
### Hardest classes



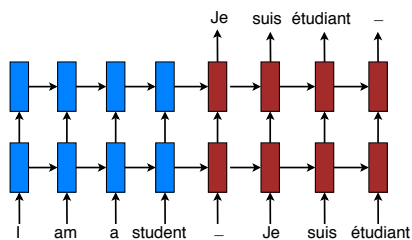


2012 – now

A Cambrian explosion in deep learning



Speech recognition



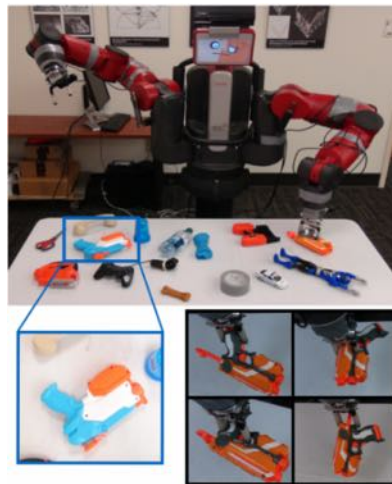
Machine Translation



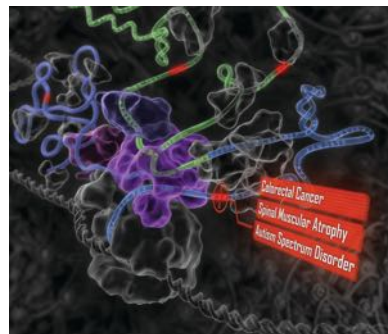
Self-Driving Cars



Game Playing



Robotics



Genomics

"Ode to Joy" harmonized in the style learned from:



Audio Generation

And many more...

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

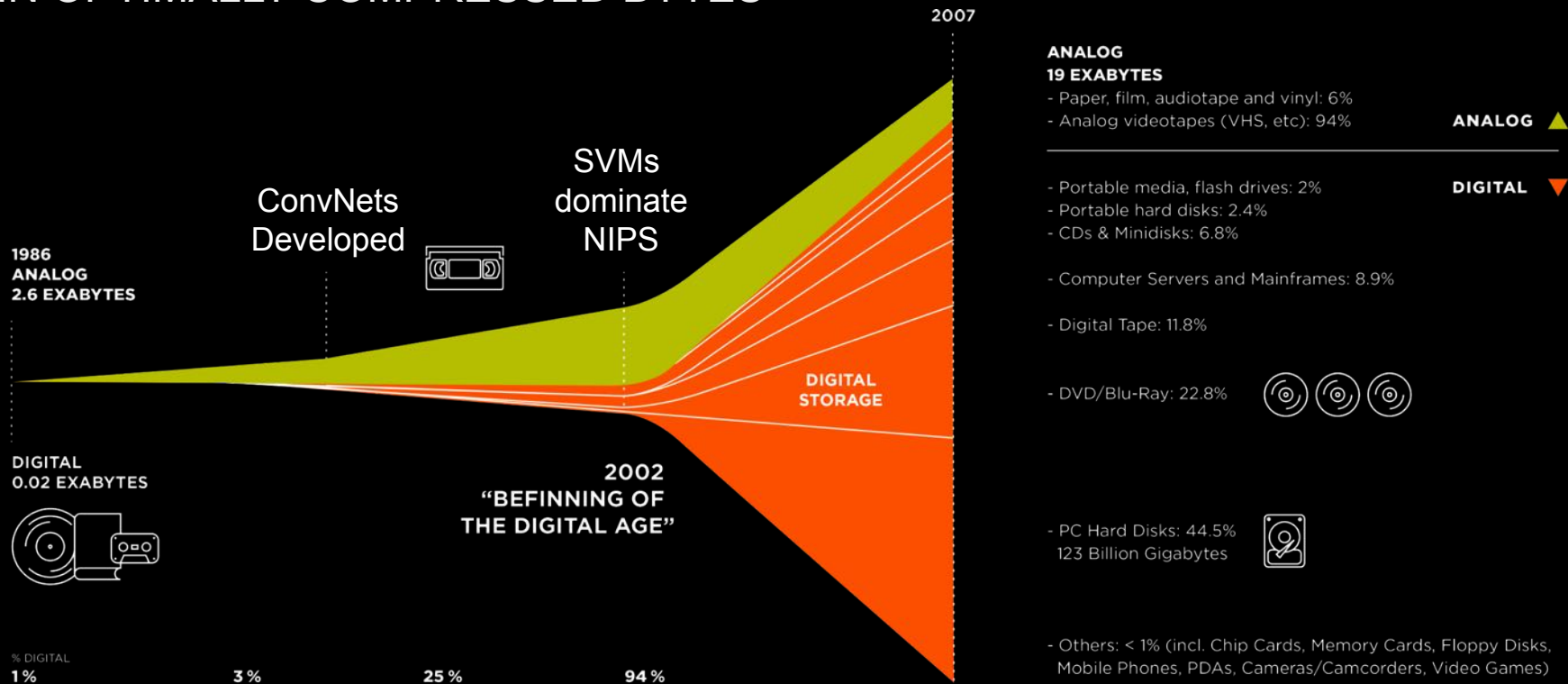
L. Pinto and A. Gupta, "Supersizing Self-supervision: Learning to Grasp from 50K Tries and 700 Robot Hours" ICRA 2015

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

Why now?

# GLOBAL INFORMATION STORAGE CAPACITY IN OPTIMALLY COMPRESSED BYTES



Source: Hilbert, M., & López, P. (2011). The World's Technological Capacity to Store, Communicate, and Compute Information. Science, 332 (6025), 60-65. [martinhilbert.net/worldinfocapacity.html](http://martinhilbert.net/worldinfocapacity.html)

# 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>



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

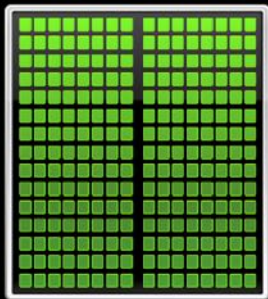
CPU

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



# 10X GROWTH IN GPU COMPUTING

2008

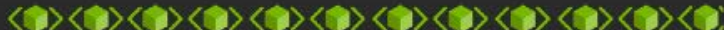
2015

150,000  
CUDA Downloads



3 Million  
CUDA Downloads

27  
CUDA Apps



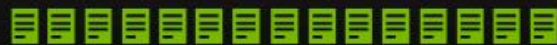
319  
CUDA Apps

60  
Universities  
Teaching



800  
Universities Teaching

4,000  
Academic  
Papers



60,000  
Academic Papers

6,000  
Tesla GPUs



450,000  
Tesla GPUs

77  
Supercomputing  
Teraflops



54,000  
Supercomputing  
Teraflops

# 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,

Journal of Machine Learning Research 15 (2014) 1929-1958

Submitted 11/13; Published 6/14

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

Nitish Srivastava  
Geoffrey Hinton  
Alex Krizhevsky  
Ilya Sutskever  
Ruslan Salakhutdinov  
*Department of Computer Science  
University of Toronto  
19 King's College Road, Rm 3309  
Toronto, Ontario, M5S 3G4, Canada.*

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

Editor: Yueshu Bengio

#### 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, it is easy to approximate the effect of averaging the predictions of all these thinned networks by simply using a single unthinned network that has smaller weights. This significantly reduces overfitting and gives major improvements over other regularization methods. We show that dropout improves the performance of neural networks on supervised learning tasks in vision, speech recognition, document classification and computational biology, obtaining state-of-the-art results on many benchmark data sets.

**Keywords:** neural networks, regularization, model combination, deep learning

#### 1. Introduction

Deep neural networks contain multiple non-linear hidden layers and this makes them very expressive models that can learn very complicated relationships between their inputs and outputs. With limited training data, however, many of these complicated relationships will be the result of sampling noise, so they will exist in the training set but not in real test data even if it is drawn from the same distribution. This leads to overfitting and many methods have been developed for reducing it. These include stopping the training as soon as performance on a validation set starts to get worse, introducing weight penalties of various kinds such as L1 and L2 regularization and soft weight sharing (Nowlan and Hinton, 1992).

With unlimited computation, the best way to “regularize” a fixed-sized model is to average the predictions of all possible settings of the parameters, weighting each setting by

©2014 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever and Ruslan Salakhutdinov.



# 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 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

## 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 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.

### 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  $\theta$  of the network, so as to minimize the loss

$$\theta = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n f(x_i, \theta)$$

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

$$\frac{1}{m} \frac{\partial \ell(x_1, \theta_2)}{\partial \theta}$$

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 (Elzinga, 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

$$f = F_2(F_1(x, \theta_1), \theta_2)$$

where  $F_1$  and  $F_2$  are arbitrary transformations, and the parameters  $\theta_1, \theta_2$  are to be learned so as to minimize the loss  $\ell$ . Learning  $\theta_1$  can be viewed as if the inputs  $x = F_1(x, \theta_1)$  are fed into the sub-network

$$f = F_2(x, \theta_2)$$

For example, a gradient descent step

$$\theta_2 \leftarrow \theta_2 - \frac{\alpha}{m} \sum_{i=1}^m \frac{\partial f(x_i, \theta_2)}{\partial \theta_2}$$

(for batch size  $m$  and learning rate  $\alpha$ ) 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,  $\theta_1$  does

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

# 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 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.

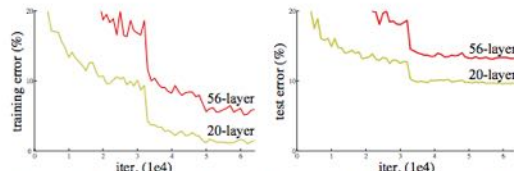


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

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 20% relative improvement on the COCO object detection dataset. Deep residual nets are foundations of our submissions to ILSVRC & COCO 2015 competitions<sup>1</sup>, where we also won 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-to-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 [11, 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] layers. Other notable visual recognition tasks [8, 12, 7, 32, 27] have also

<sup>1</sup>[https://image-net.org/show\\_results.cgi?det=1&top=5&range=ILSVRC2015](https://image-net.org/show_results.cgi?det=1&top=5&range=ILSVRC2015) and [http://mscoco.org/dataset/object\\_detection\\_coco2015/](http://mscoco.org/dataset/object_detection_coco2015/).

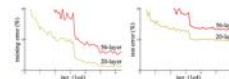


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

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 [11, 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**)

So what is deep learning?

# Three key ideas

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

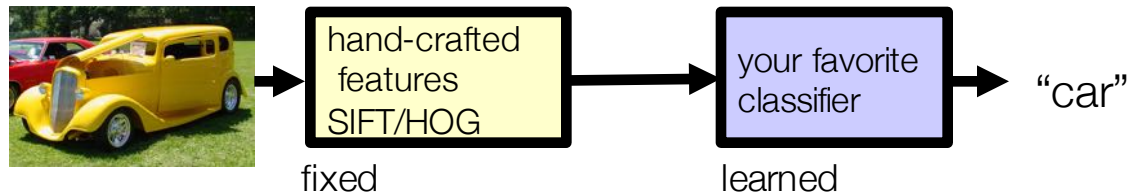


# Three key ideas

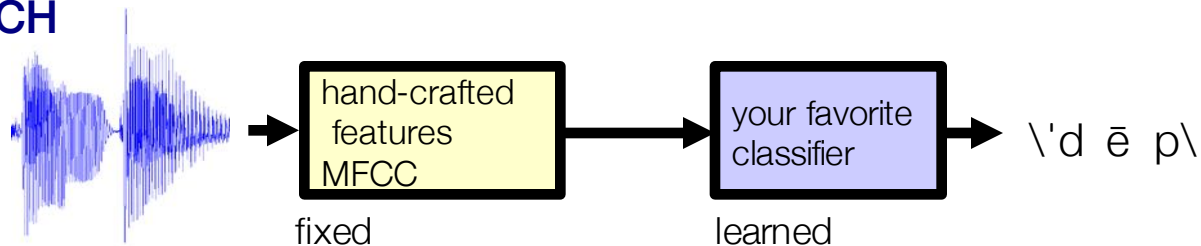
- **(Hierarchical) Compositionality**
  - Cascade of non-linear transformations
  - Multiple layers of representations
- End-to-End Learning
  - Learning (goal-driven) representations
  - Learning to feature extract
- Distributed Representations
  - No single neuron “encodes” everything
  - Groups of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

# Traditional Machine Learning

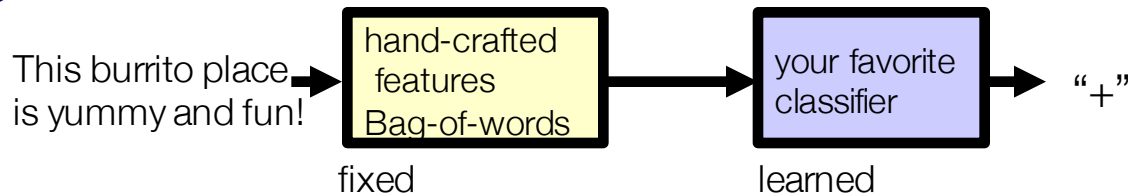
## VISION



## SPEECH



## NLP



# Hierarchical Compositionality

## VISION

pixels → edge → texton → motif → part → object

## SPEECH

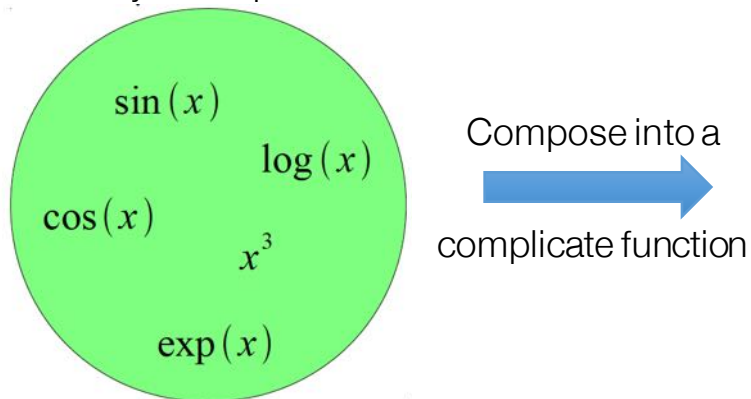
sample → spectral  
band → formant → motif → phone → word

## NLP

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

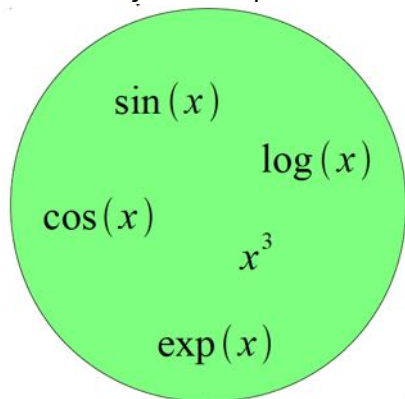
# Building A Complicated Function


Given a library of simple functions



# Building A Complicated Function

Given a library of simple functions

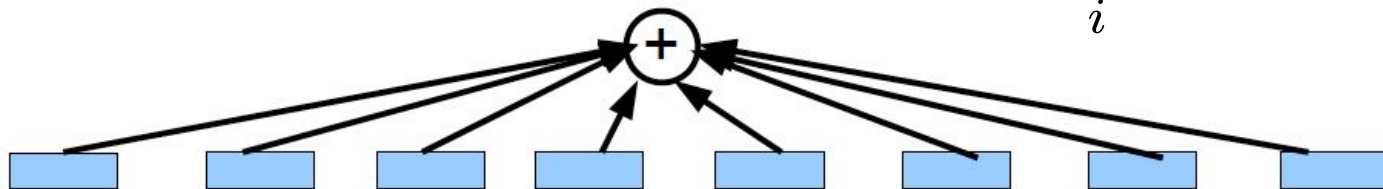


Compose into a  
  
 complicate function

## Idea 1: Linear Combinations

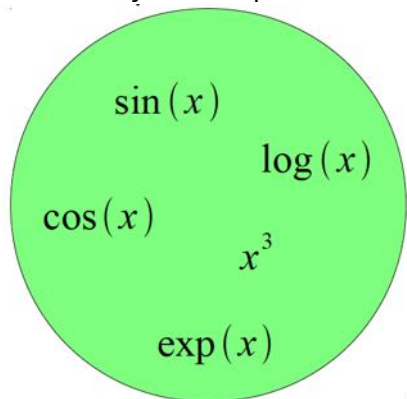
- Boosting
- Kernels
- ...


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



# Building A Complicated Function

Given a library of simple functions

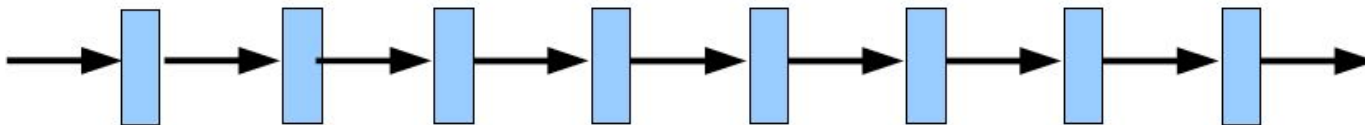


Compose into a  
  
 complicate function

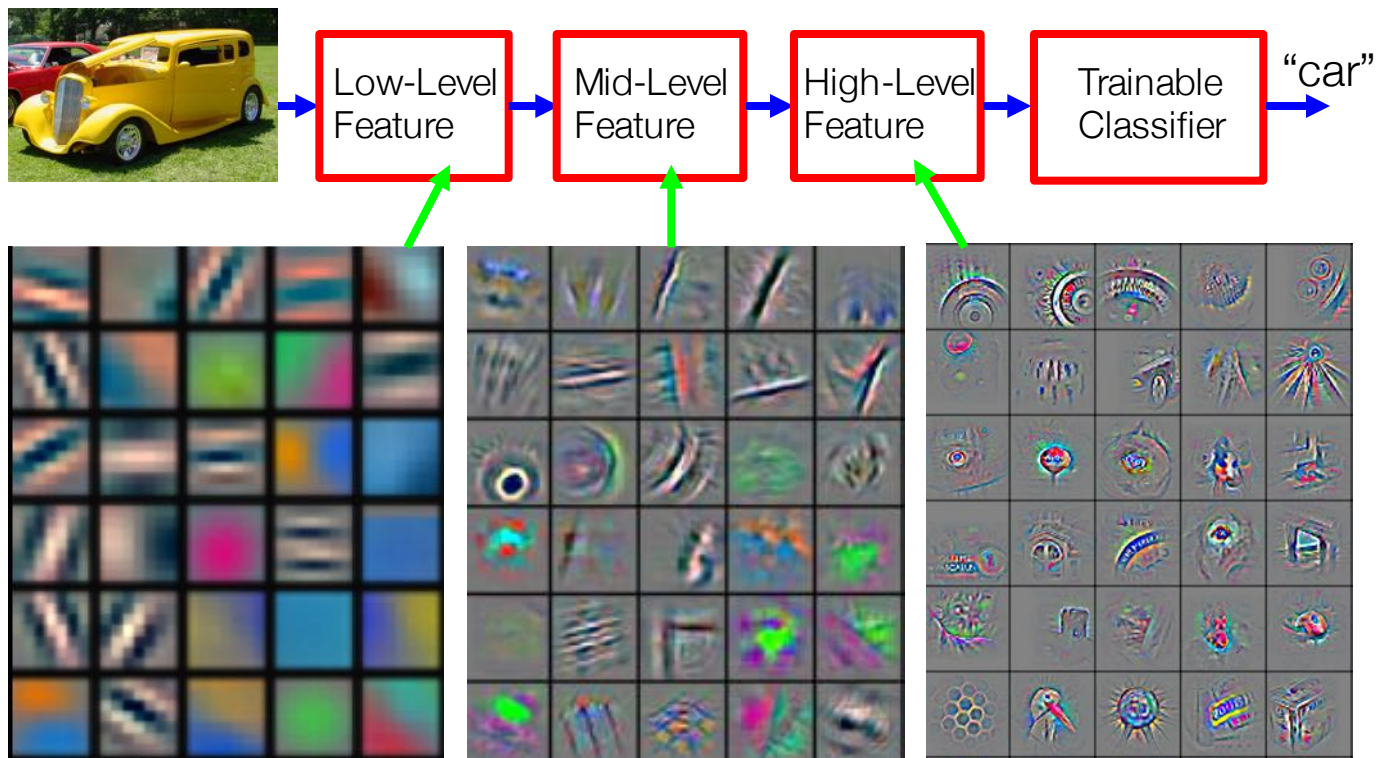
## Idea 2: Compositions

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

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



# Deep Learning = Hierarchical Compositionality



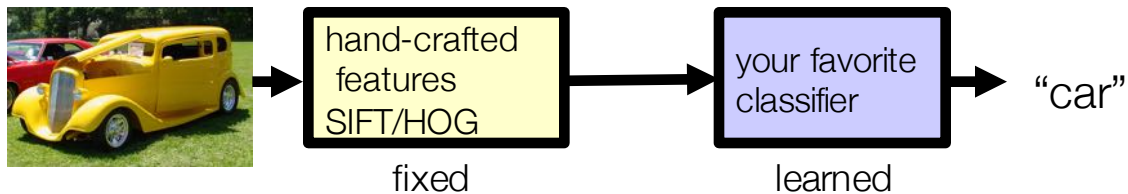
# Three key ideas

- (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

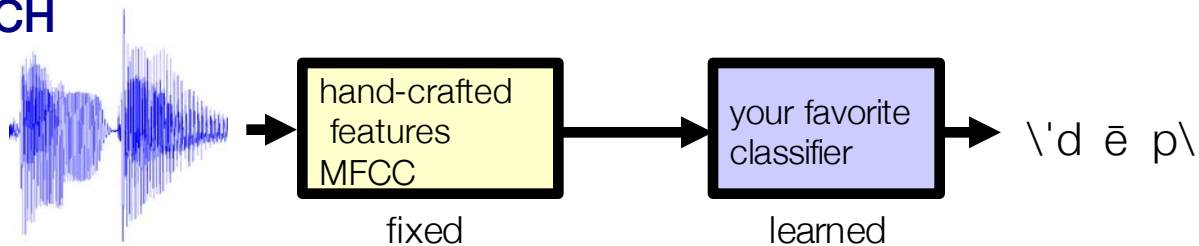


# Traditional Machine Learning

## VISION

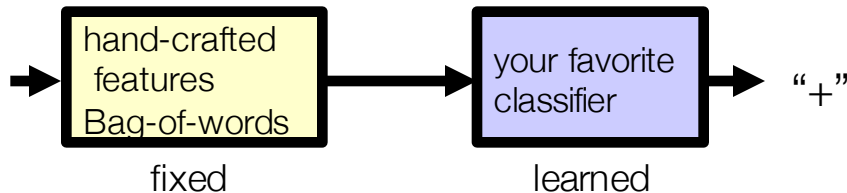


## SPEECH



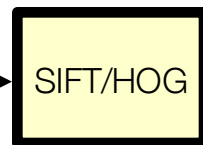
## NLP

This burrito place  
is yummy and fun!

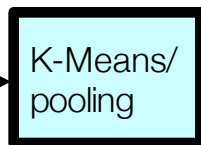


# More accurate version

## VISION



fixed



unsupervised



supervised

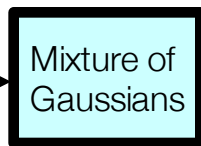
“car”

“Learned”

## SPEECH



fixed



unsupervised

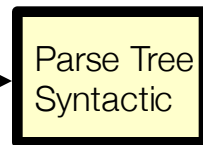


supervised

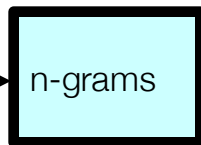
\ 'd ē p \

## NLP

This burrito place  
is yummy and fun!



fixed



unsupervised

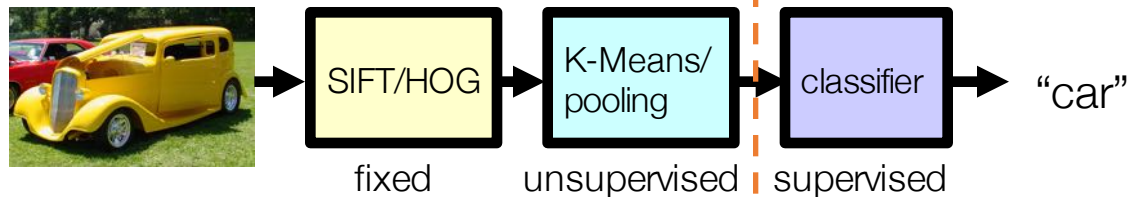


supervised

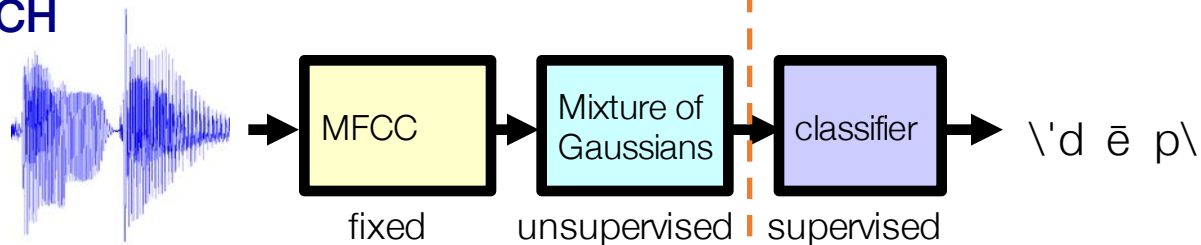
“+”

# Deep Learning = End-to-End Learning

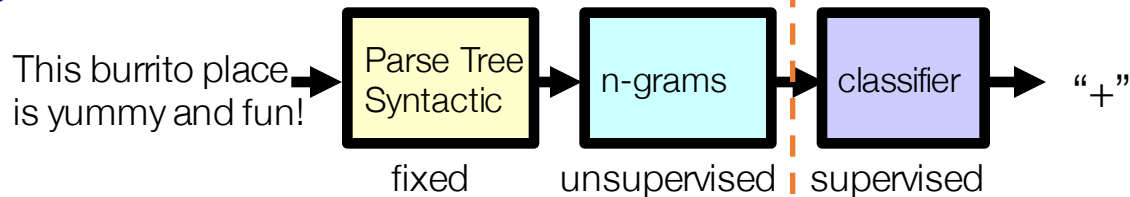
## VISION



## SPEECH



## NLP

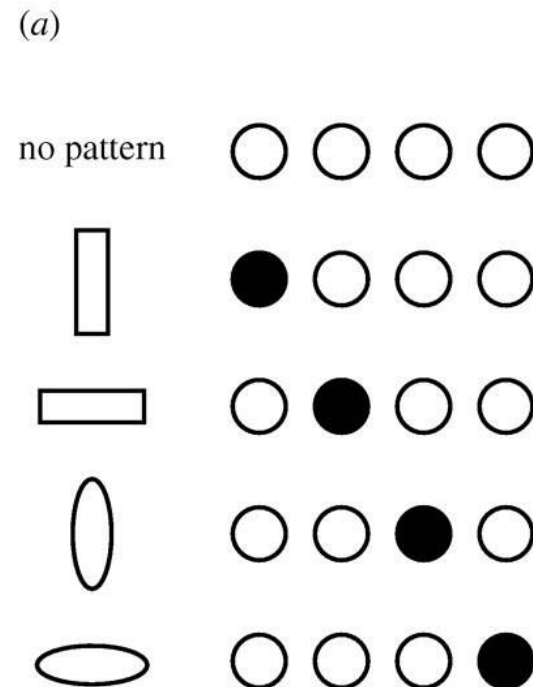


# Three key ideas

- (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

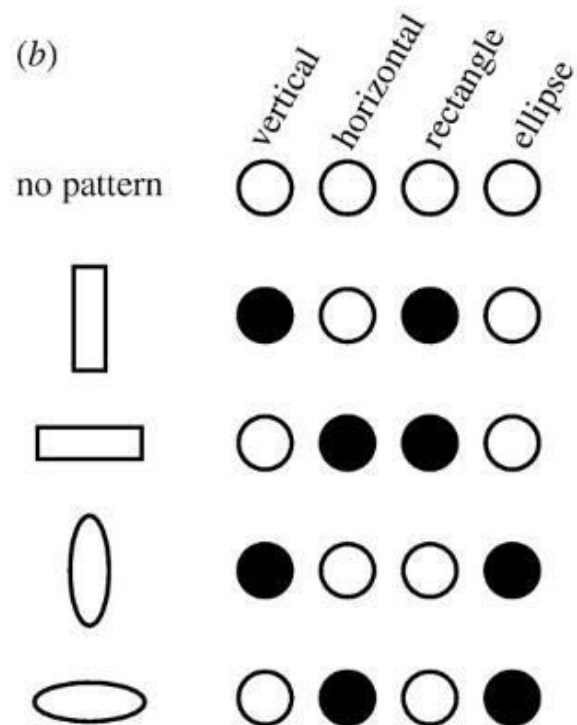
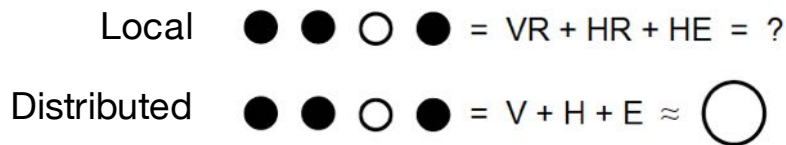
# Localist representations

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



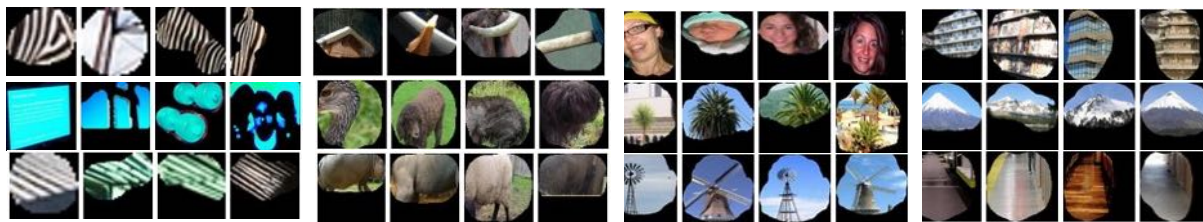
# Power of distributed representations!

## Scene Classification



- Possible internal representations:

- Objects
- Scene attributes
- Object parts
- Textures



Simple elements & colors

Object part

Object

Scene



# Three key ideas of deep learning

- **(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