

BBM406

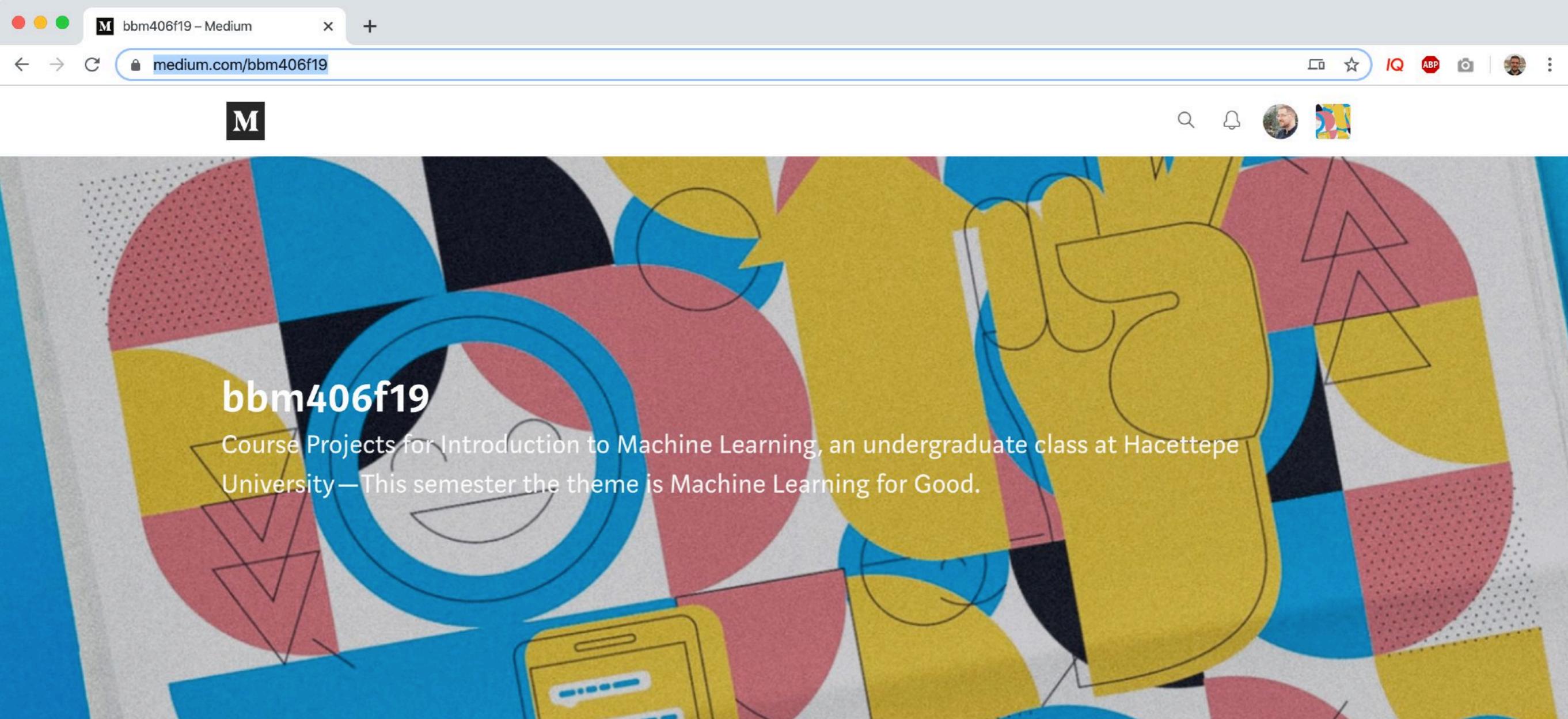
Fundamentals of Machine Learning

Lecture 13:
Introduction to Deep Learning



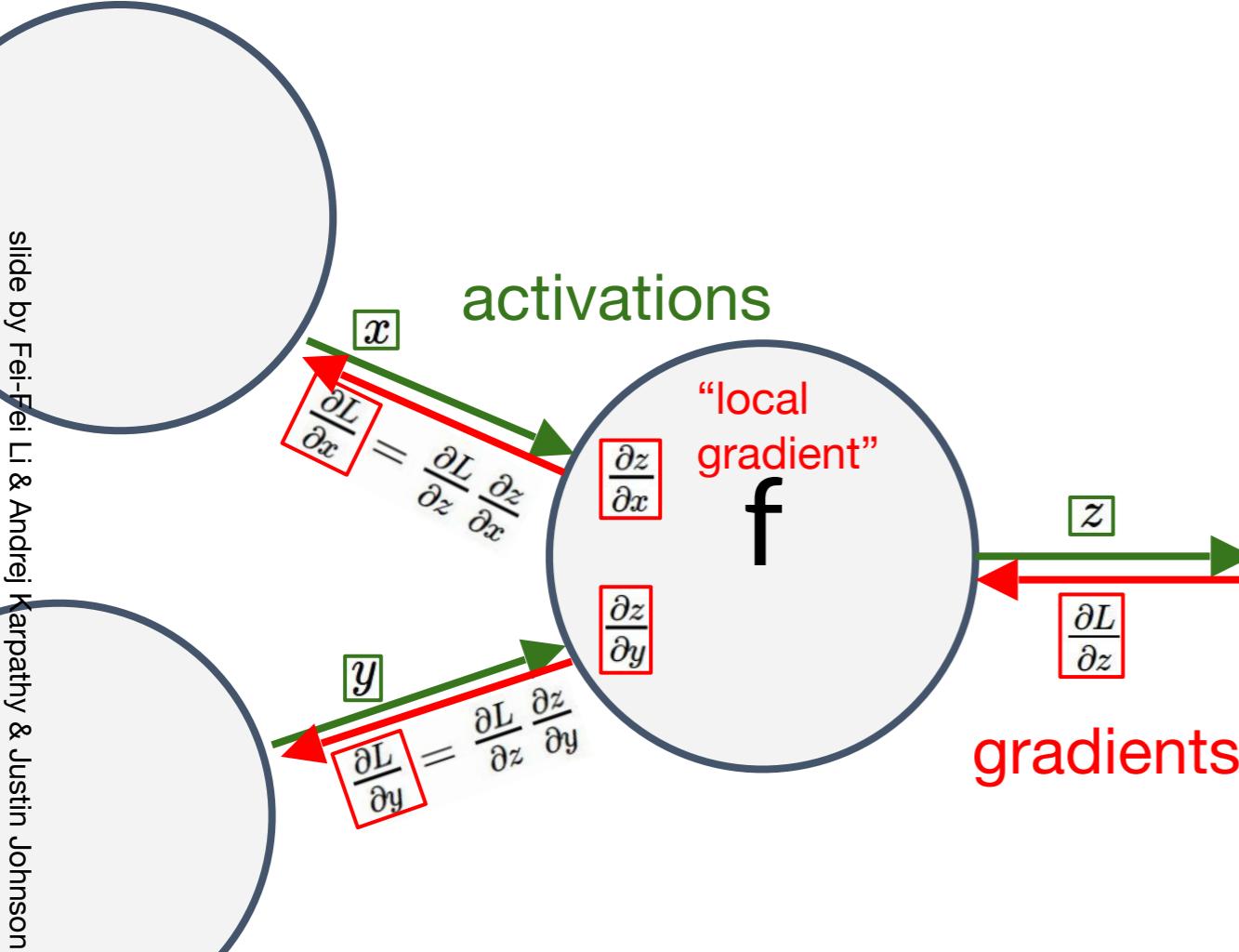
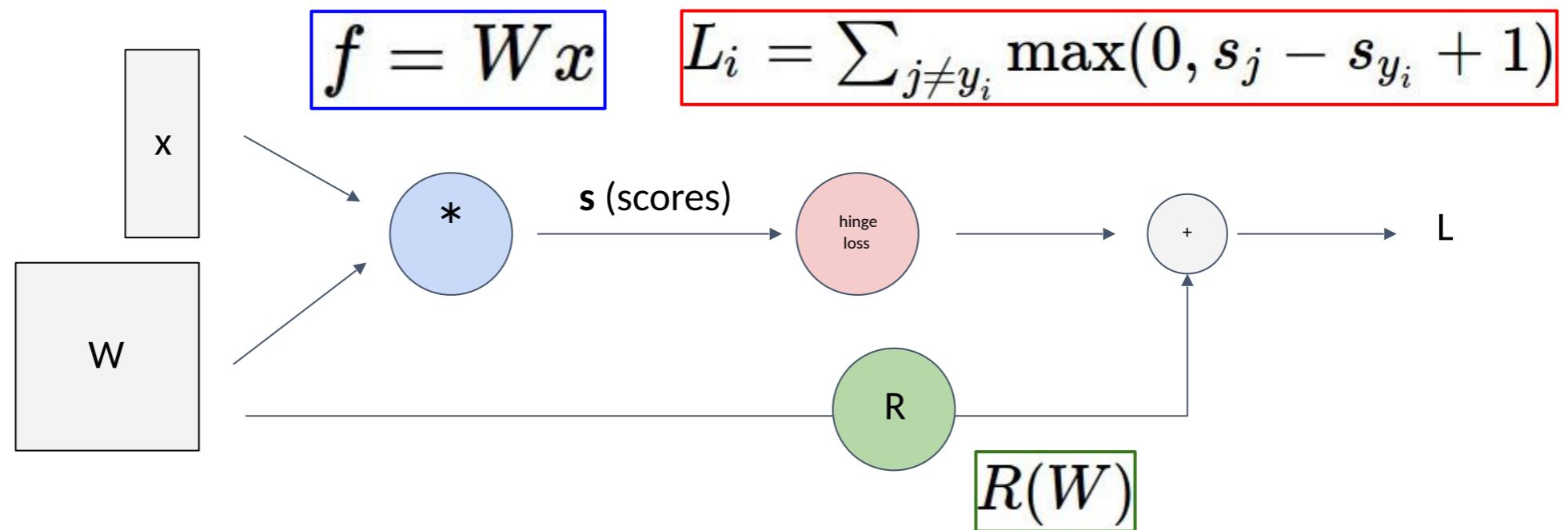
HACETTEPE
UNIVERSITY
COMPUTER
VISION LAB

A reminder about course projects



- From now on, regular (weekly) blog posts about your progress on the course projects!
- We will use medium.com

Last time.. Computational Graph



```
class ComputationalGraph(object):
    ...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

Last time... Training Neural Networks

Mini-batch SGD

Loop:

- 1. Sample** a batch of data
- 2. Forward** prop it through the graph, get loss
- 3. Backprop** to calculate the gradients
- 4. Update** the parameters using the gradient

This week

- Introduction to Deep Learning
- Deep Convolutional Neural Networks



What is deep learning?

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

REVIEW

Deep learning

Vinod Nair¹*, Yoshua Bengio² & Geoffrey Hinton³

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 dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other areas of computer vision. Deep learning discovers intricate structure in large data sets by using the backpropagation algorithm to iteratively change its internal parameters that are used to compute the representation in each layer from the representation in the previous layer. Deep learning has led to significant improvements in processing images, video, speech and audio, whereas recurrent nets have shone light on sequential data such as text and speech.

Machine learning technology powers many aspects of modern society, from web search engines to filtering on social networks to recommendations on e-commerce websites, and it is having records in image recognition^{1–4} and speech recognition^{5–7}. Machine learning systems are used to identify objects in images, transcribe speech into text, match news items, posts or products to users' interests, and much more. In some cases, interestingly, these applications make use of a class of techniques called deep learning.

Traditional machine learning techniques were limited in their ability to process natural data in their raw form. For decades, computers extracted a pattern from raw data using a series of nested, overlapping layers of decision rules. This was a slow and error-prone process that demanded expert knowledge to design a feature extractor that transformed the raw data (such as the pixel values of a face) into features that the learning subsystems, often a classifier, could detect or classify patterns in the input.

Deep learning is a set of methods that allow a machine to be fed with raw data and automatically discover the representations needed for detection or classification. Deep-learning methods are taught to extract features from raw data using a series of nested layers, obtained by compressing simple low-level features (such as edges) into a representation at a higher, slightly more abstracted, level. The combination of enough such transformations, very complex functions can be learned. This task is often referred to as ‘unsupervised learning’ because the input data does not contain explicit labels. A feature extractor that transformed the raw data (such as the pixel values of a face) into features that the learning subsystems, often a classifier, could detect or classify patterns in the input.

Supervised learning

The core component of machine learning, deep or not, is supervised learning. Imagine that we want to build a system that can classify images as containing, say, a chair, a person or a cat. We first need to collect a large number of images of chairs, people and cats, each labeled with its category. During training, the machine is shown one image and produces an output as the form of a vector of scores, one for each category. The machine is trained to produce higher scores for images of categories that it has seen before than for images of categories that it has not seen before. This is done by calculating the difference between the output scores and the desired scores of success. The machine then modifies its internal adjustable parameters to reduce the error. This is done by using a gradient descent algorithm, which computes a gradient vector, for each weight, indicating by what amount the error would increase or decrease if the weight were increased or decreased. The objective function, averaged over all the training examples, can

¹Department of Computer Science, NYU-Princeton, New York, New York 10003, USA, ²Department of Computer Science and Engineering, University of Waterloo, Waterloo, Ontario N2L 3G1, Canada, ³Singapore Institute of Technology, Singapore 573483, Singapore. *e-mail: vnair@cs.nyu.edu

NATURE | VOL. 521 | 28 MAY 2015

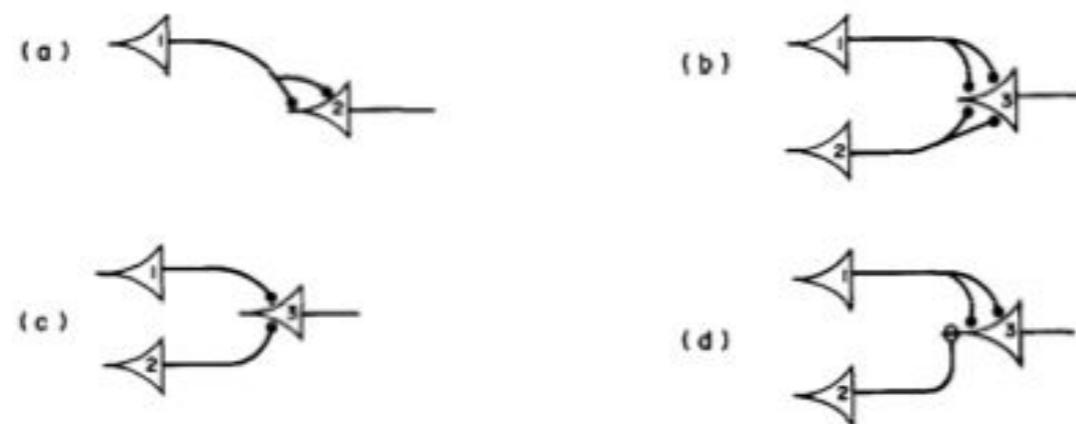
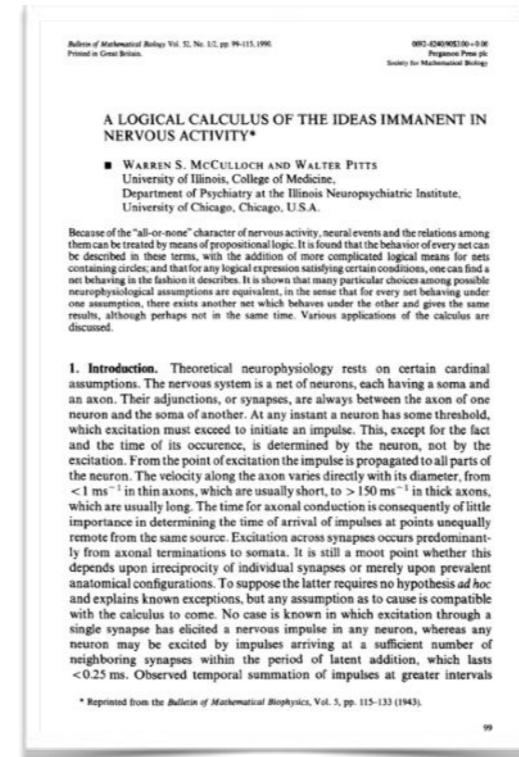
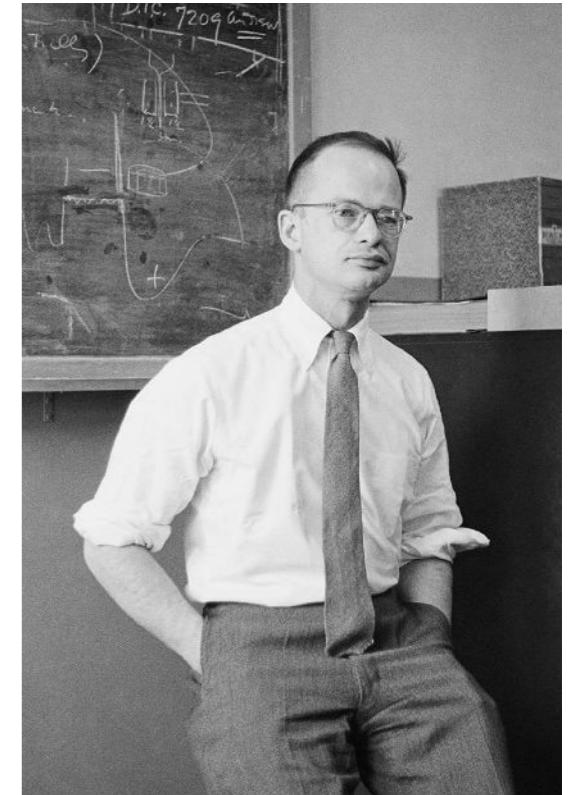
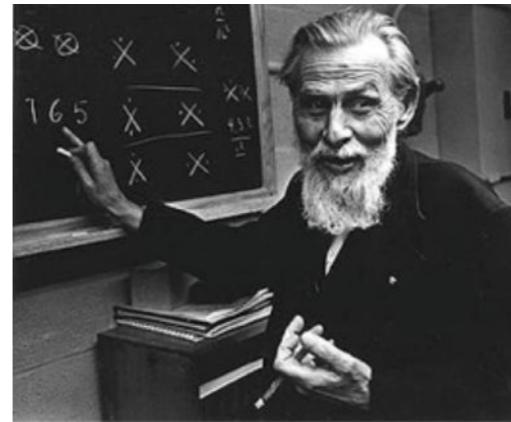
“Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction.”

- Yann LeCun, Yoshua Bengio and Geoff Hinton

1943 – 2006: A Prehistory of Deep Learning

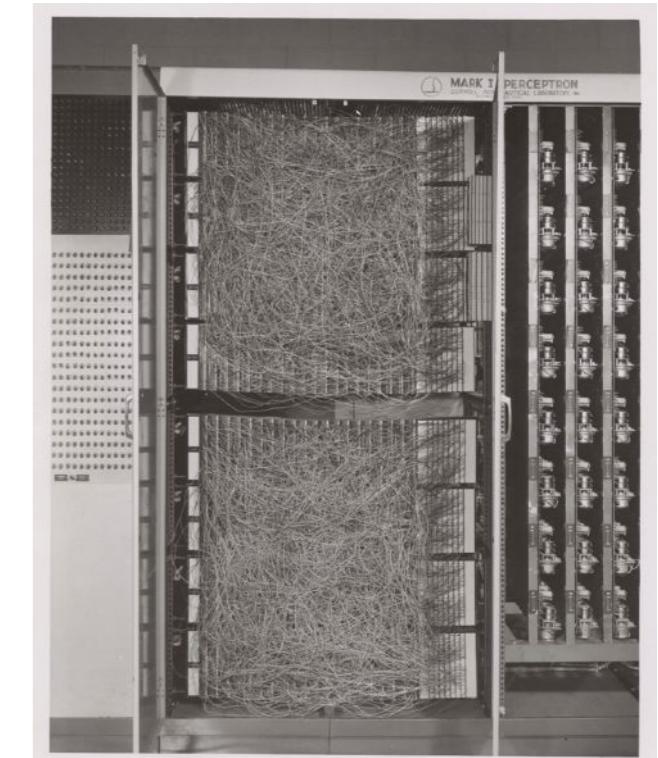
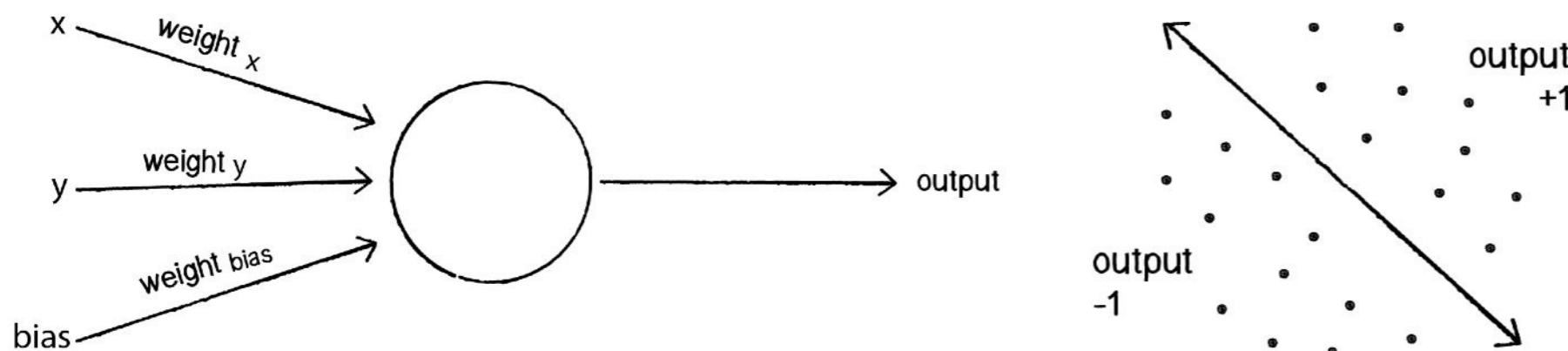
1943: Warren McCulloch and Walter Pitts

- 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

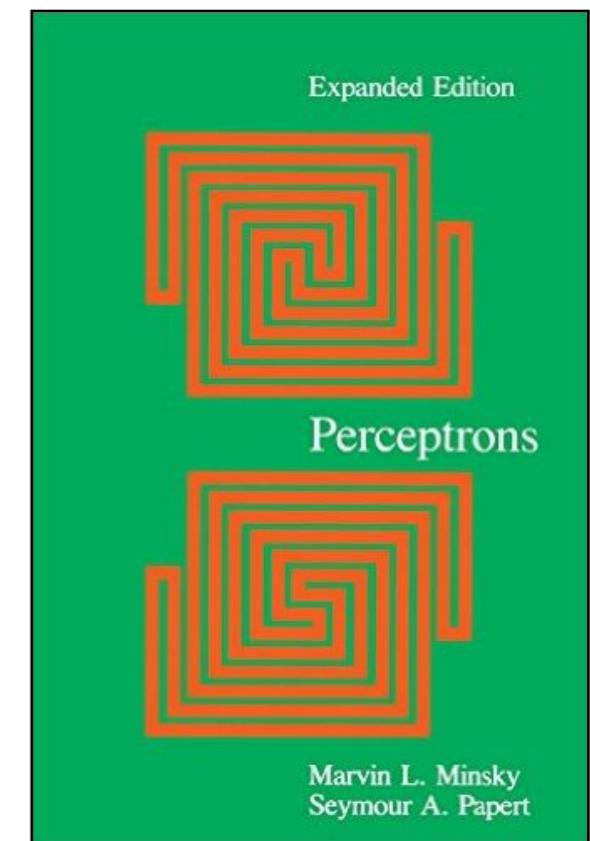
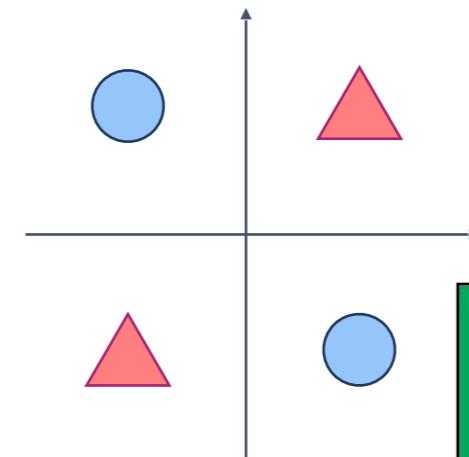


1969: Marvin Minsky and Seymour Papert

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

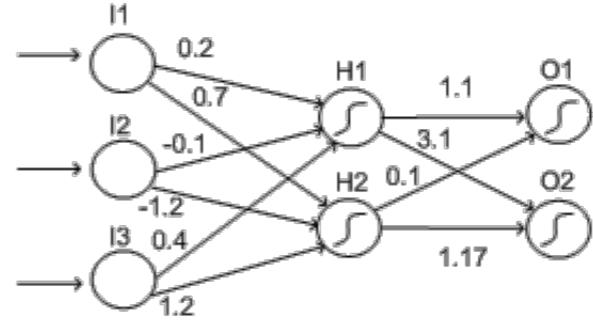


- Perceptrons can only represent linearly separable functions.
 - such as **XOR** Problem
- Wrongly attributed as the reason behind the **AI winter**, a period of reduced funding and interest in AI research



1990s

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



Backpropagation Through Time: What It Does and How to Do It

PAUL J. WERBOS

Abstract: Backpropagation is the most widely used tool in the field of artificial neural networks. At the core of backpropagation is a large system of linear equations mostly and efficiently solved by iterative numerical methods. This paper shows how to implement such a system, and also gives some applications which demonstrate the usefulness of backpropagation.

This paper discusses basic backpropagation, a simple method which is very effective for learning feedforward networks, but which does not always work. It presents the basic equations for backpropagation, and shows how they can be used to learn various pattern-recognition involving dynamic systems, robotics, weather forecasting, and other real-world systems. Other recent researches, extending backpropagation to recurrent networks, are also mentioned, along with practical issues which arise with such methods. Backpropagation is presented as a generalization of the standard backpropagation algorithm, which is the subject which enables backpropagation to work.

Section III will give some variation to describe backpropagation, and show how it has been applied to various problems by a number of authors, including, at least, Wasserman and Shannahan [1], Seiden [2], and LeCun et al. [3]. Section IV will discuss the work of Kawato [4], Elman and Zipser [5], and Seiden [6]. This section will also mention what needs to be done to extend backpropagation to recurrent networks. Section V will discuss the simple discussion, and how to do better.

At its core, backpropagation is simply an efficient and effective way to learn the weights of a feedforward network to work out the mathematics of their particular application. This paper will show how to implement backpropagation, which can be translated into computer code and applied directly to neural networks.

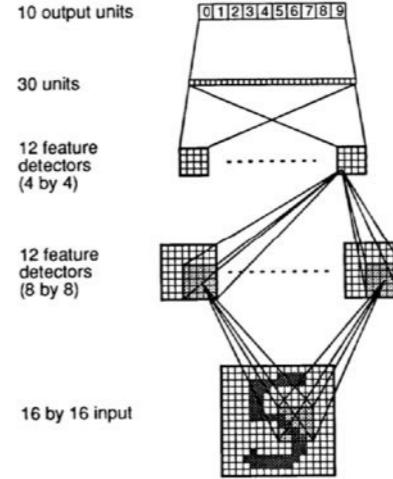
Some of the applications of backpropagation have already used terms of backpropagation, which may be called "basic backpropagation".

Manuscript received September 12, 1989; revised March 11, 1990.
The author is with the National Science Foundation, 1801 G St., NW, Washington, DC 20546.
GSS Log Number W9002.

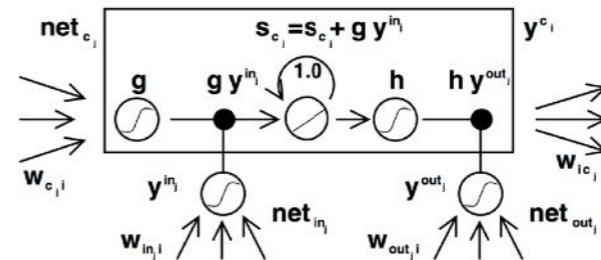
U.S. Government work not protected by U.S. copyright.

1190 PROCEEDINGS OF THE IEEE, VOL. 78, NO. 10, OCTOBER 1990

- Training multi-layer perceptrons
 - Back-propagation
(Rumelhart, Hinton, Williams, 1986)
 - Back-propagation through time (BPTT)
(Werbos, 1988)



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



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. 07732

ABSTRACT

We present an application of back-propagation networks to handwritten digit recognition. Minimal preprocessing of the data was performed, no feature extraction or feature selection was done, and the network was specifically designed for the task. The input of the network consists of normalized images of digits. The method has 18 neurons and about a 9% reject rate on degraded digits provided by the U.S. Postal Service.

INTRODUCTION

The main point of this paper is to show that large back-propagation (BP) networks can be trained to perform complex tasks, such as handwritten digit recognition, with little or no prior knowledge of the problem. Unlike most previous work on the subject (Fischer et al., 1989), the learning network is fed with images, rather than features, thus improving the ability of BP networks to deal with large amounts of low-level information.

Prior work performed on simple digit images (LeCun, 1989) showed that the addition of the network significantly influences the network's generalization ability. Good generalization, however, requires the network to have knowledge 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.

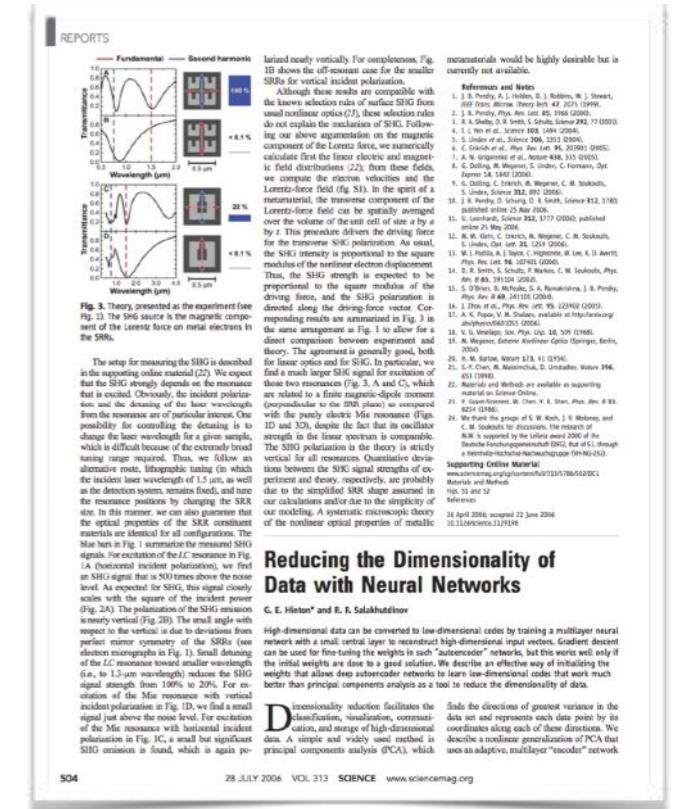


Fig. 13. Theory, presented as the experiment (see Fig. 13). The SHG source is the magnetic component of the Lorentz force on metal electrons in the SRR.

The setup for measuring the SHG is described in the supporting online material (27). We expect that the SHG strength depends on the resonance that is excited. Obviously, the incident polarization does not have to be vertical to excite resonance. Even the resonance of zero polarizations interest. One possibility for controlling the detailing is to change the laser wavelength for a given example, which is the most simple and directly having range required. Thus, we follow an alternative route, lithographic tuning (in which the incident laser wavelength of 1.5 μm , as well as the size of the SRR, are varied simultaneously with the resonance position by changing the SRR size). In this manner, we can also guarantee that the optical properties of the SRR constituent materials are identical for all configurations. The idea here is to measure the unpolarized SHG signals. For excitation of the LC resonance in Fig. 1A (normal incident polarization), we first use SHG signal that is SOTR shows the lowest level. After excitation of the SOTR, the SHG signal scales with the square of the incident polarization (Fig. 2A). The polarization of the SHG emission is nearly vertical (Fig. 2D). The small angle with respect to the vertical is due to the finite periodical symmetry of the SRRs (see electron microscopy in Fig. 1). Small detailing of the LC resonance toward smaller wavelengths (i.e., to 1.3 μm) is needed to reach the required signal level from 100% to 20%. For excitation of the Mie resonance with vertical incident polarization in Fig. 1D, we find a small signal just above the noise level. For excitation of the SRR with vertical incident polarization in Fig. 1C, a small but significant SHG emission is found, which is again polarized nearly vertically. For completeness, Fig. 1B shows off-resonance results for the smaller SRRs for vertical incident polarization.

Instrumentation would be highly desirable but is currently not available.

References and Notes

1. J. B. Pendry, M. J. Holden, D. J. Robbins, W. J. Stewart, *IEEE Trans. Microw. Theory Tech.* **45**, 1297 (1997).
2. J. B. Pendry, *Phys. Rev. Lett.* **85**, 3964 (2000).
3. J. B. Pendry, *Nature* **425**, 77 (2003).
4. I. I. Smirnov et al., *Science* **308**, 1494 (2005).
5. G. Gralak et al., *Phys. Rev. Lett.* **95**, 203901 (2005).
6. A. N. Grigorenko et al., *Nature* **436**, 515 (2005).
7. L. V. Butkov, *Electromagnetic Theory* (Pergamon, Oxford, 1962).
8. G. Dohle, C. Trichet, M. Weiland, C. M. Soukoulis, *J. Appl. Phys.* **91**, 7070 (2002).
9. J. B. Pendry, D. Schurig, D. R. Smith, *Science* **312**, 1780 (2006).
10. G. Leonhardt, *Science* **312**, 1777 (2006); published online 29 May 2006.
11. M. L. Brongersma, R. P. van der Heijden, C. M. Soukoulis, S. Linden, *Opt. Lett.* **31**, 1203 (2006).
12. M. L. Brongersma, R. P. van der Heijden, C. M. Soukoulis, *Phys. Rev. Lett.* **96**, 017401 (2006).
13. S. D. Quiles, D. McPhale, S. A. Ramachandran, J. B. Pendry, *Phys. Rev. B* **49**, 24115 (1994).
14. A. C. Pipkin, *Phys. Rev. B* **50**, 223942 (1994).
15. A. C. Pipkin, V. M. Shalaev, available at <http://arxiv.org/abs/physics/0605025>.
16. A. C. Pipkin, *Z. Phys. Chem.* **16**, 309 (1985).
17. M. Wegener, *Extreme Wavefield Optics* (Springer, Berlin, 2005).
18. H. K. Button, *Nature* **173**, 11 (1954).
19. S.-K. Choi, M. Maksimchuk, C. Umstadter, *Nature* **364**, 525 (2000).
20. Materials and Methods are available in supporting information on the web.
21. M. L. Brongersma, R. P. van der Heijden, C. M. Soukoulis, *Phys. Rev. Lett.* **96**, 017401 (2006).
22. Materials and Methods are available in supporting information on the web.
23. A. C. Pipkin, V. M. Shalaev, *Phys. Rev. B* **53**, 13472 (1996).
24. C. M. Soukoulis for discussions, the research of R. R. S. was supported by the Leibniz award 2006 of the German Research Foundation (DFG) and the DFG through a Helmholtz-National-Macromodel (HNM-2).

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.

Dimensionality reduction facilitates the classification, visualization, communication, and interpretation of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which finds the directions of greatest variance in the data set and represents each data point by its projection onto these axes. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer “encoder” network

504 28 JULY 2006 VOL 313 SCIENCE www.sciencemag.org

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

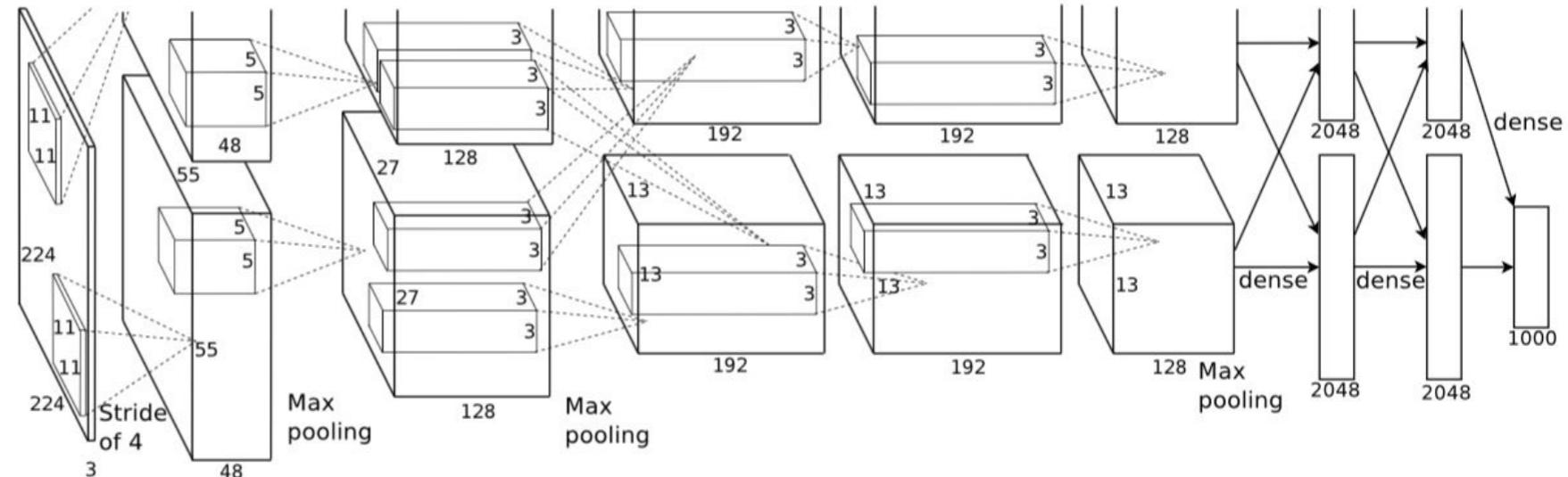


Image classification

ILSVRC 2012 Competition

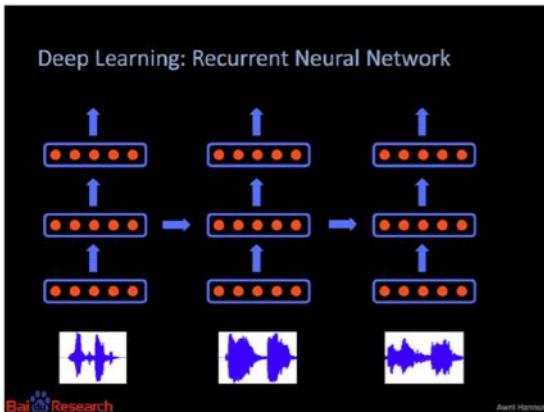
2012 Teams	%Error
Supervision (Toronto)	15.3
ISI (Tokyo)	26.1
VGG (Oxford)	26.9
XRCE/INRIA	27.0
UvA (Amsterdam)	29.6
INRIA/LEAR	33.4

CNN based,
non-CNN based

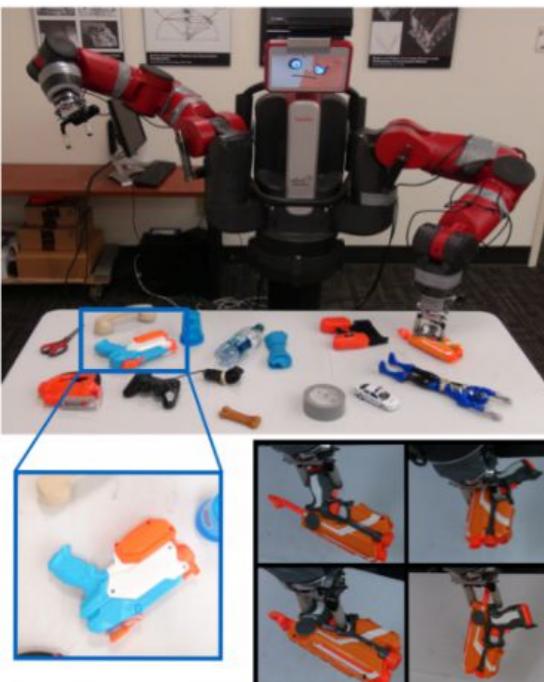


- The success of AlexNet, a deep convolutional network
 - 7 hidden layers (not counting some max pooling layers)
 - 60M parameters
- Combined several tricks
 - ReLU activation function, data augmentation, dropout

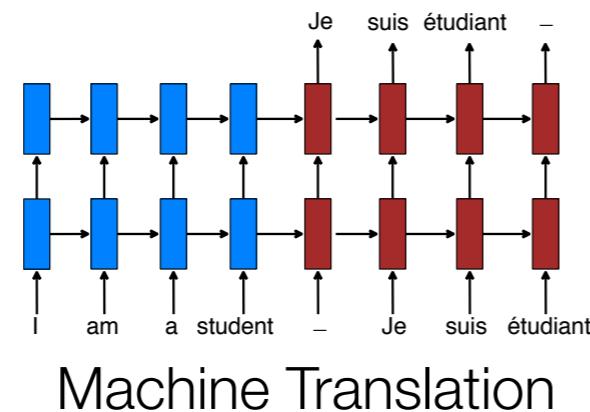
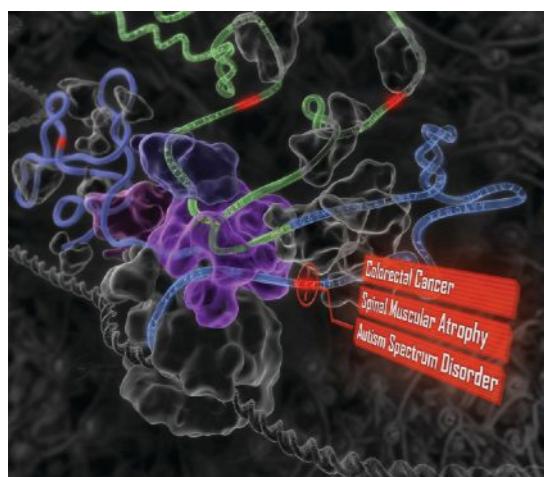
2012 – now
Deep Learning Era



Speech recognition



Robotics



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

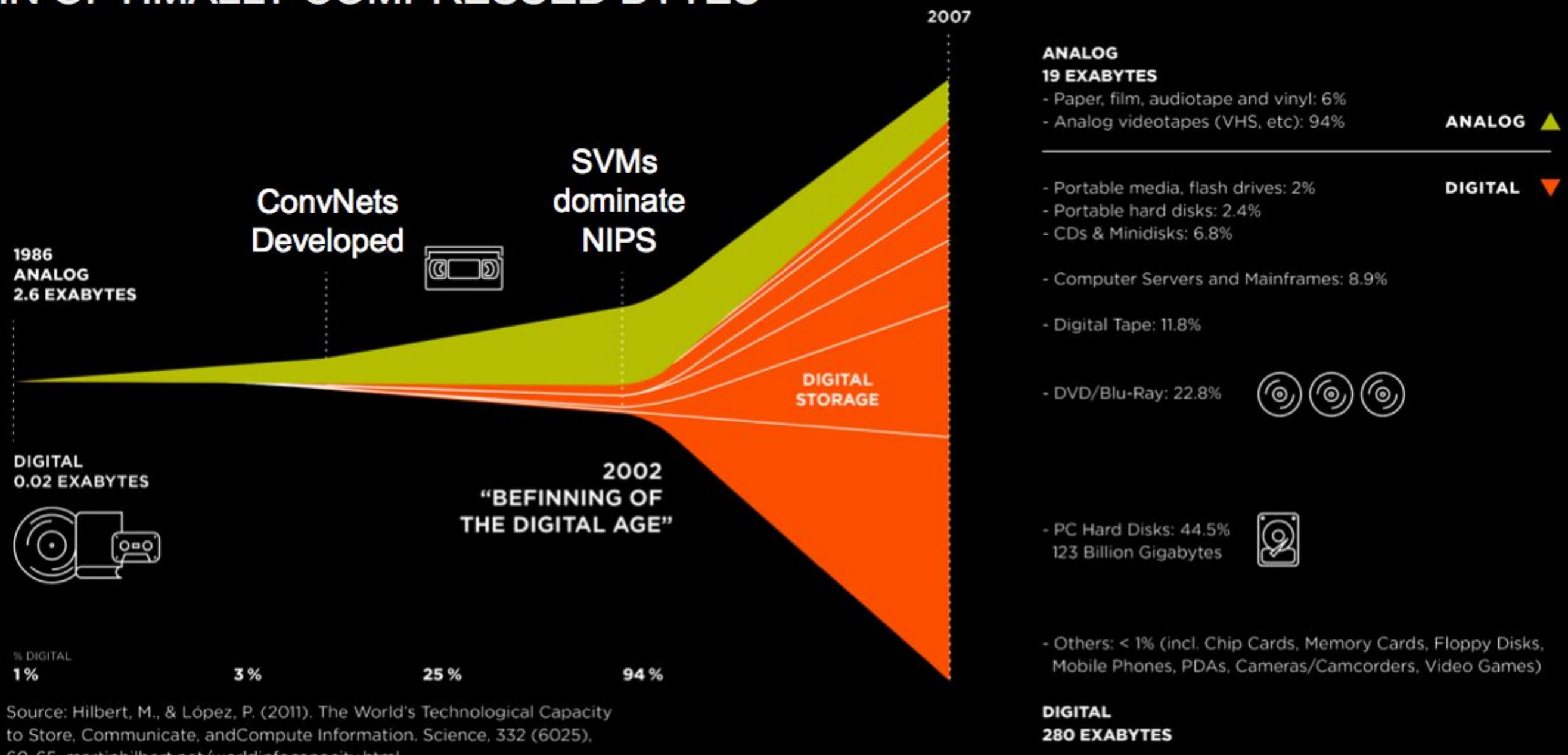
H. Y. Xiong et al., "The human splicing code reveals new insights into the genetic determinants of disease", Science 347, 2015

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

And many more... 19

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

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



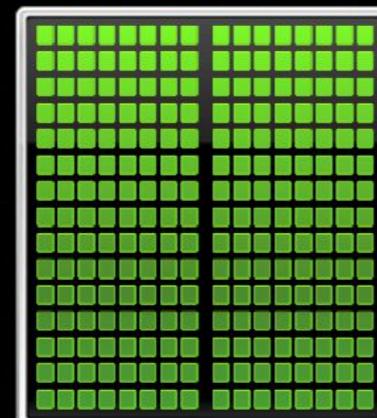
170 TFLOPS FP16
8x Tesla P100 16GB
NVLink Hybrid Cube Mesh
Accelerates Major AI Frameworks
Dual Xeon
7 TB SSD Deep Learning Cache
Dual 10GbE, Quad IB 100Gb
3RU - 3200W

CPU
Optimized for
Serial Tasks



GPU Accelerator

Optimized for
Parallel Tasks



TITAN X

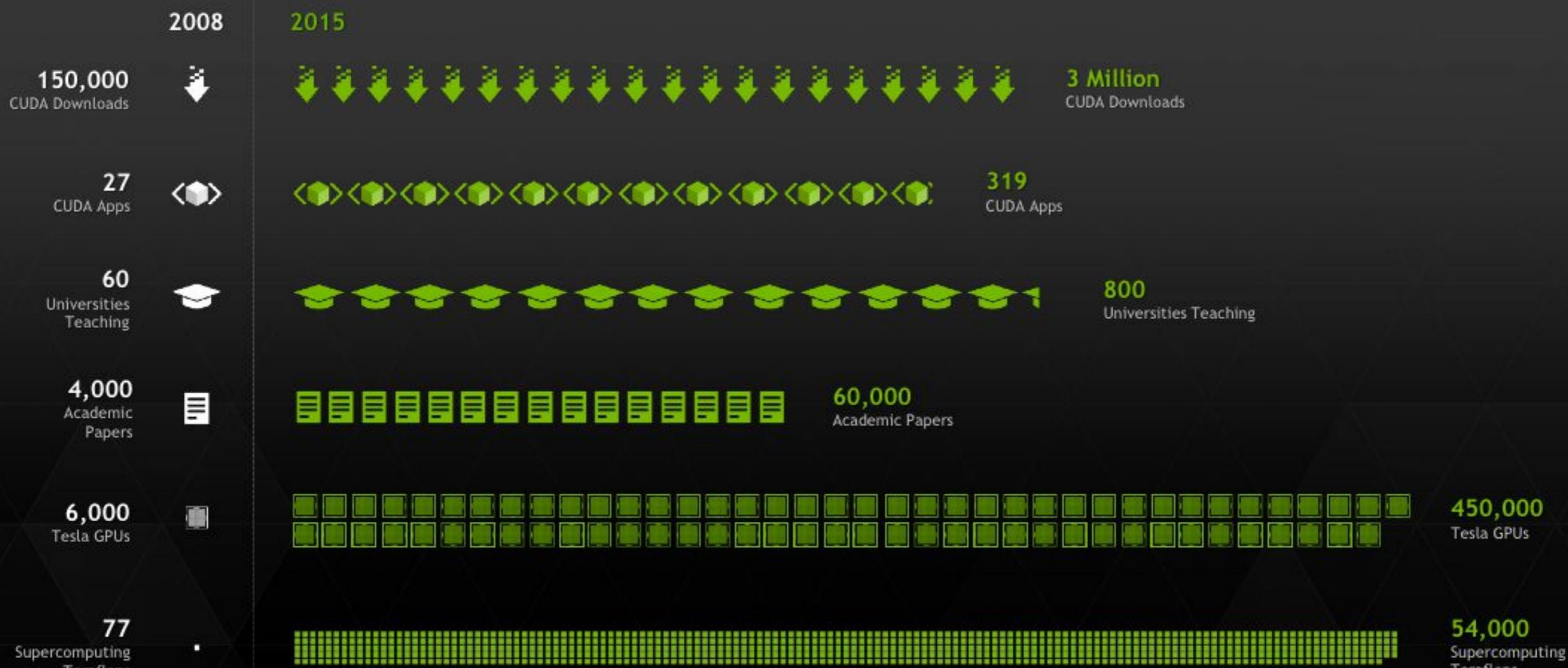
THE WORLD'S FASTEST GPU

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,

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
10 Kings College Road, Rm 3302
Toronto, Ontario, M5S 3G4, Canada.*

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

Editor: Yoshua 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 untrained 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.

- 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 Nor-

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.

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 $x_{1\dots N}$ is the training data set. With SGD, the training proceeds in steps, and at each step we consider a *mini-batch* $x_{1\dots m}$ of size m . The mini-batch is used to approximate the gradient of the loss function with respect to the parameters, by computing

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

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

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

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

- 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 unreference 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

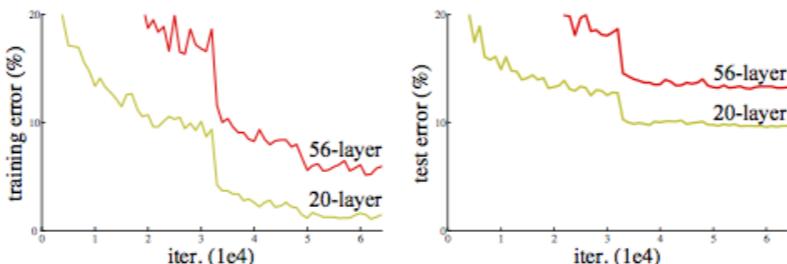


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

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 unreference 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/datasets/coco2015challenge2015>

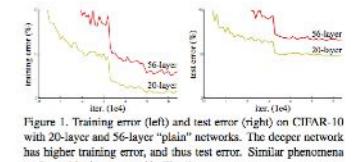


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

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)

So what is deep learning?

Three key ideas

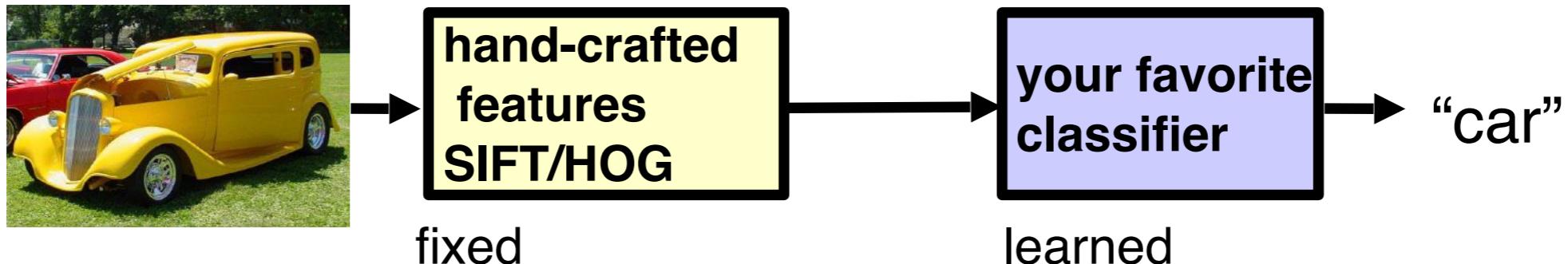
- (Hierarchical) Compositionality
- End-to-End Learning
- Distributed Representations

Three key ideas

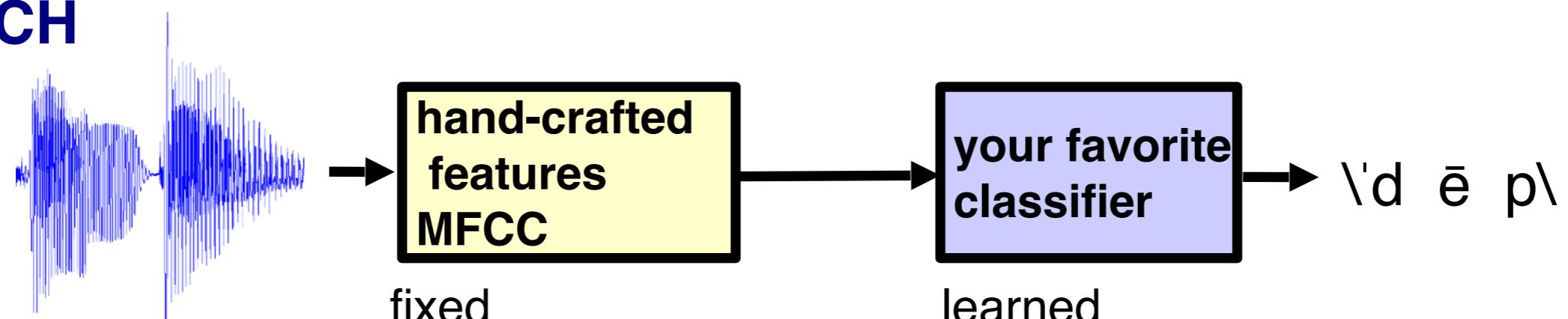
- **(Hierarchical) Compositionality**
 - Cascade of non-linear transformations
 - Multiple layers of representations
- End-to-End Learning
 - Learning (goal-driven) representations
 - Learning to feature extract
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons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

Traditional Machine Learning

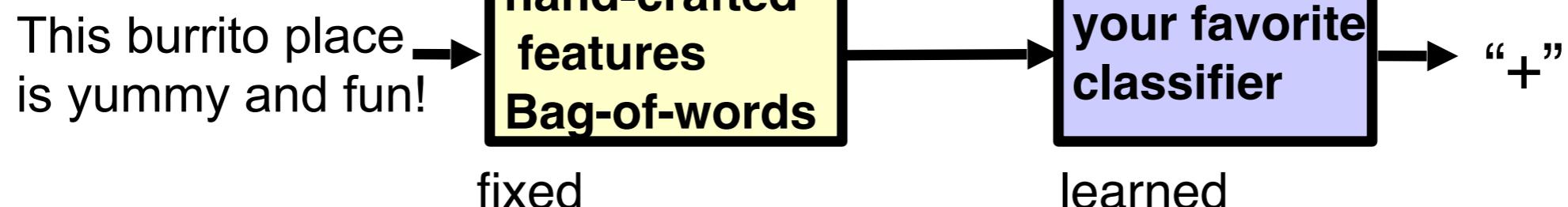
VISION



SPEECH

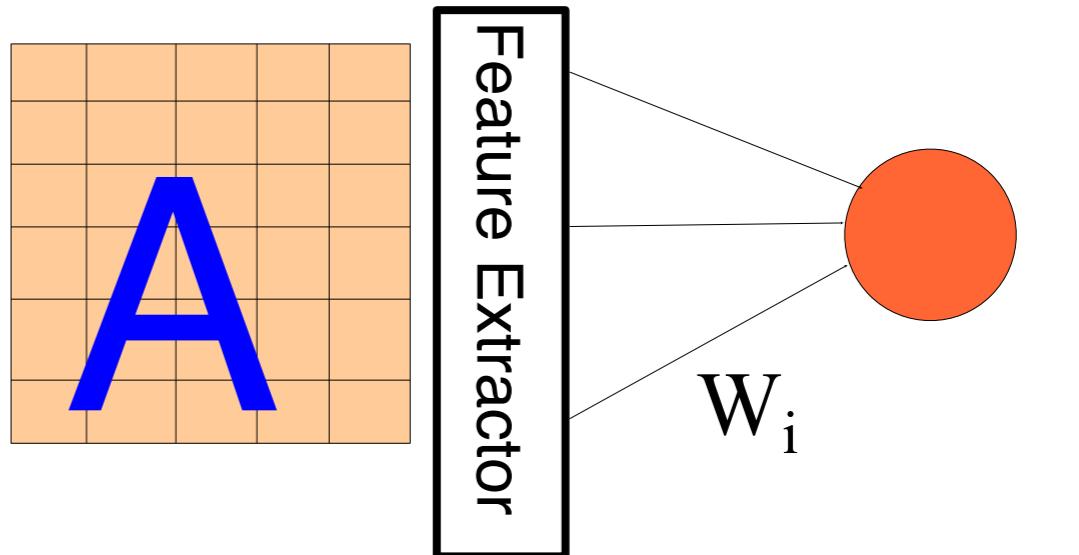


NLP

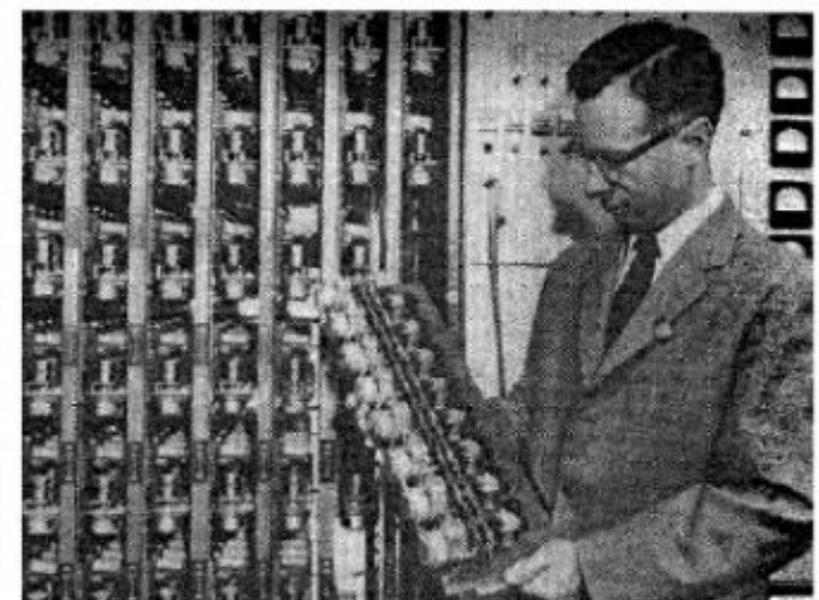
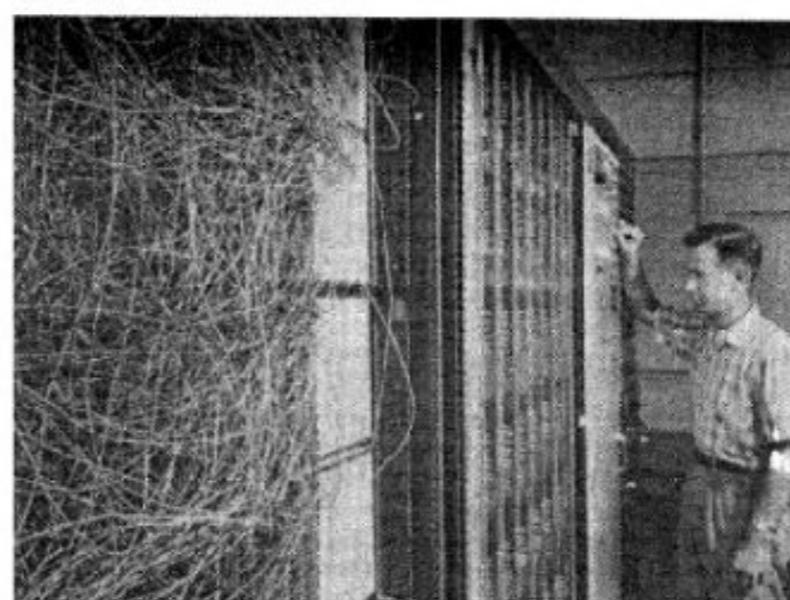
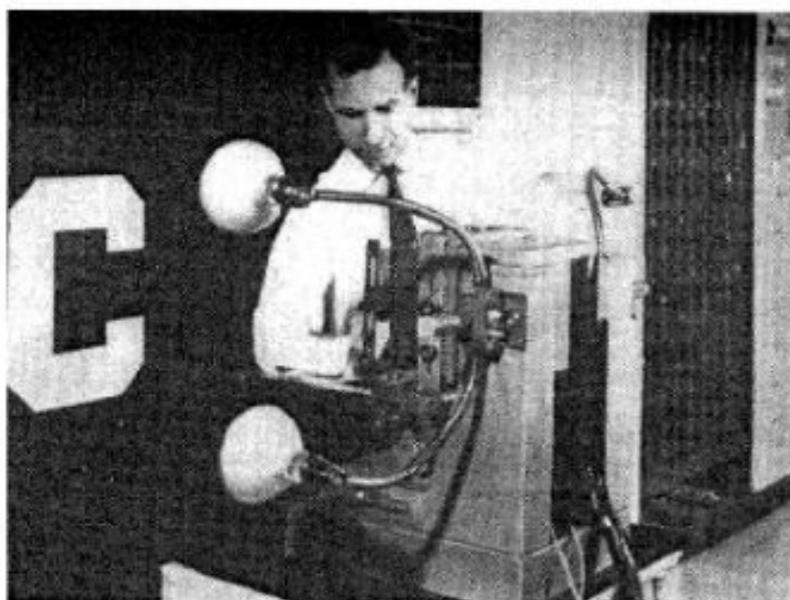


It's an old paradigm

- 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 → texton → motif → part → object

SPEECH

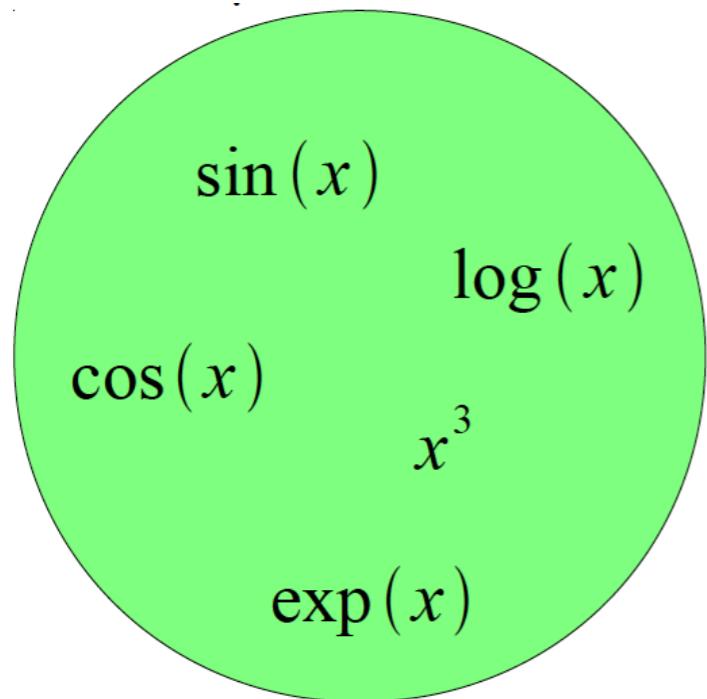
sample → spectral band → formant → motif → phone → word

NLP

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

Building A Complicated Function

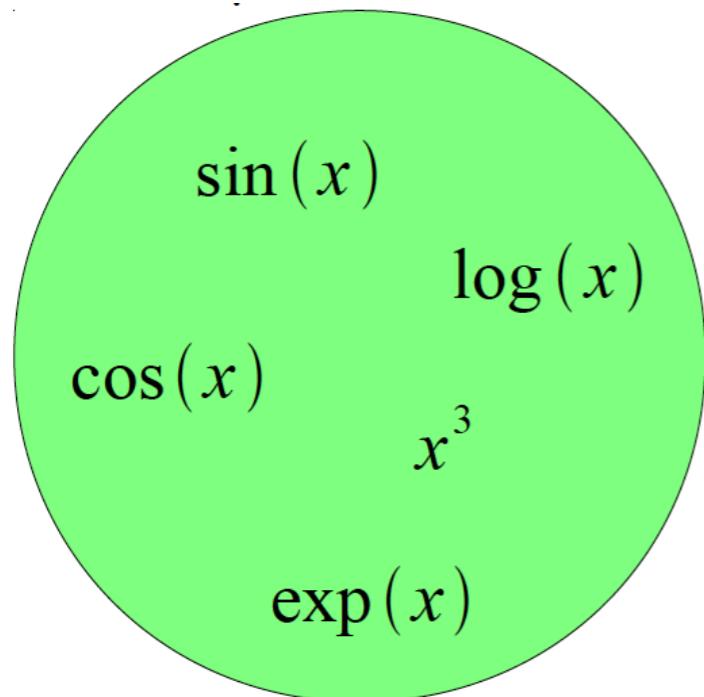
Given a library of simple functions



Compose into a
complicate function

Building A Complicated Function

Given a library of simple functions

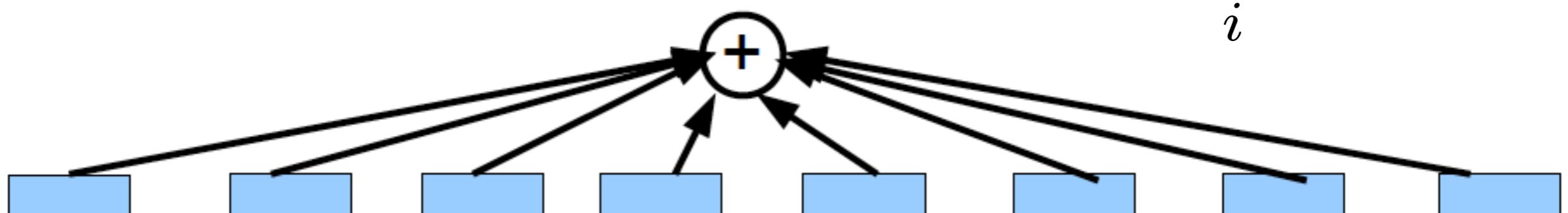


Compose into a
complicate function

Idea 1: Linear Combinations

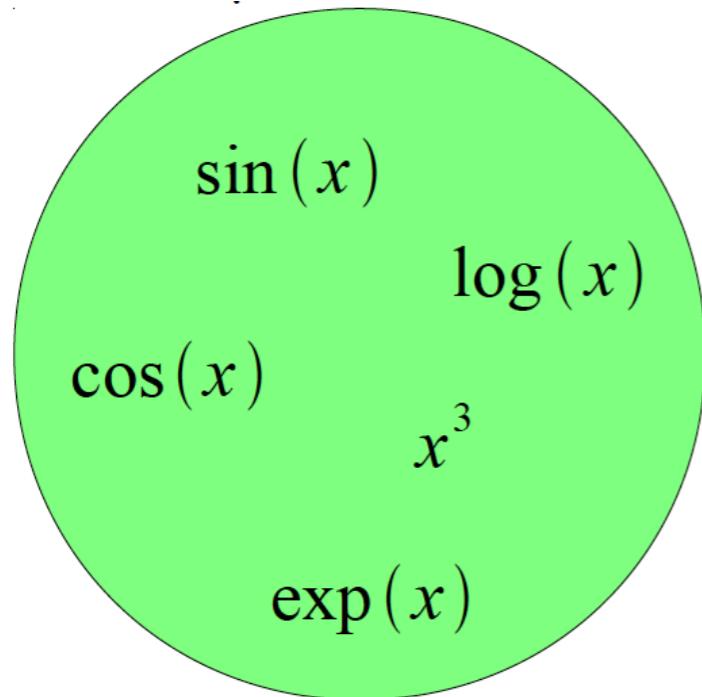
- Boosting
- Kernels
- ...

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



Building A Complicated Function

Given a library of simple functions

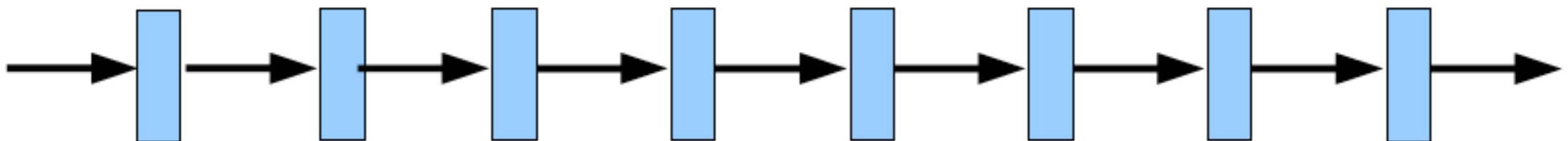


Compose into a
complicate function

Idea 2: Compositions

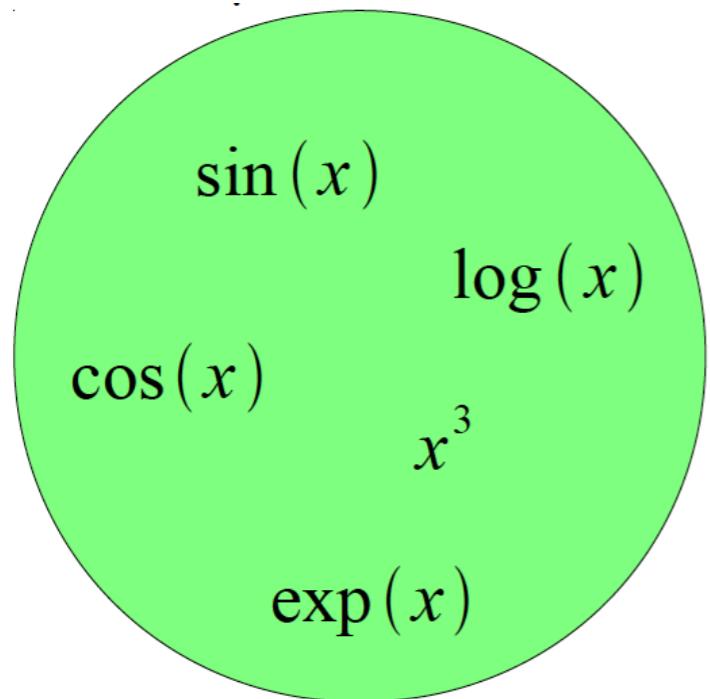
- Deep Learning
- Grammar models
- Scattering transforms...

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



Building A Complicated Function

Given a library of simple functions

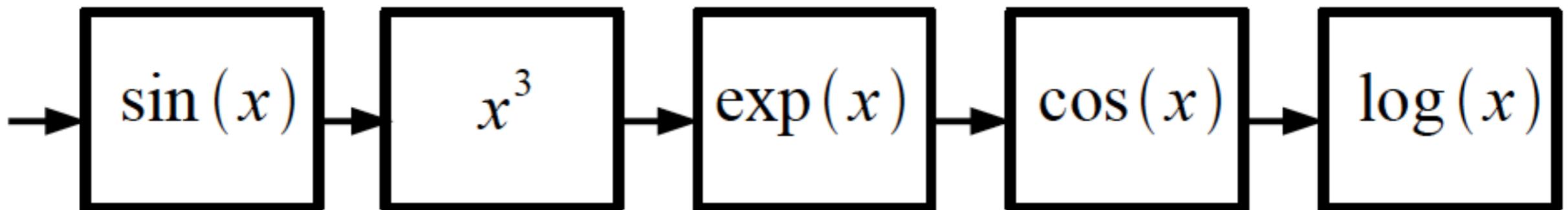


Compose into a
complicate function

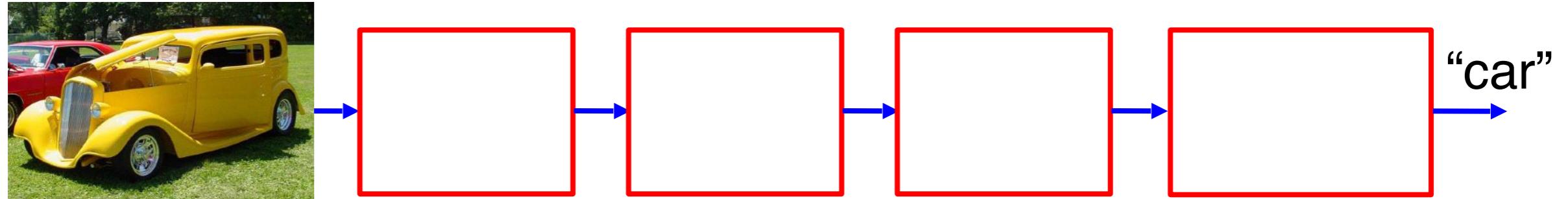
Idea 2: Compositions

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

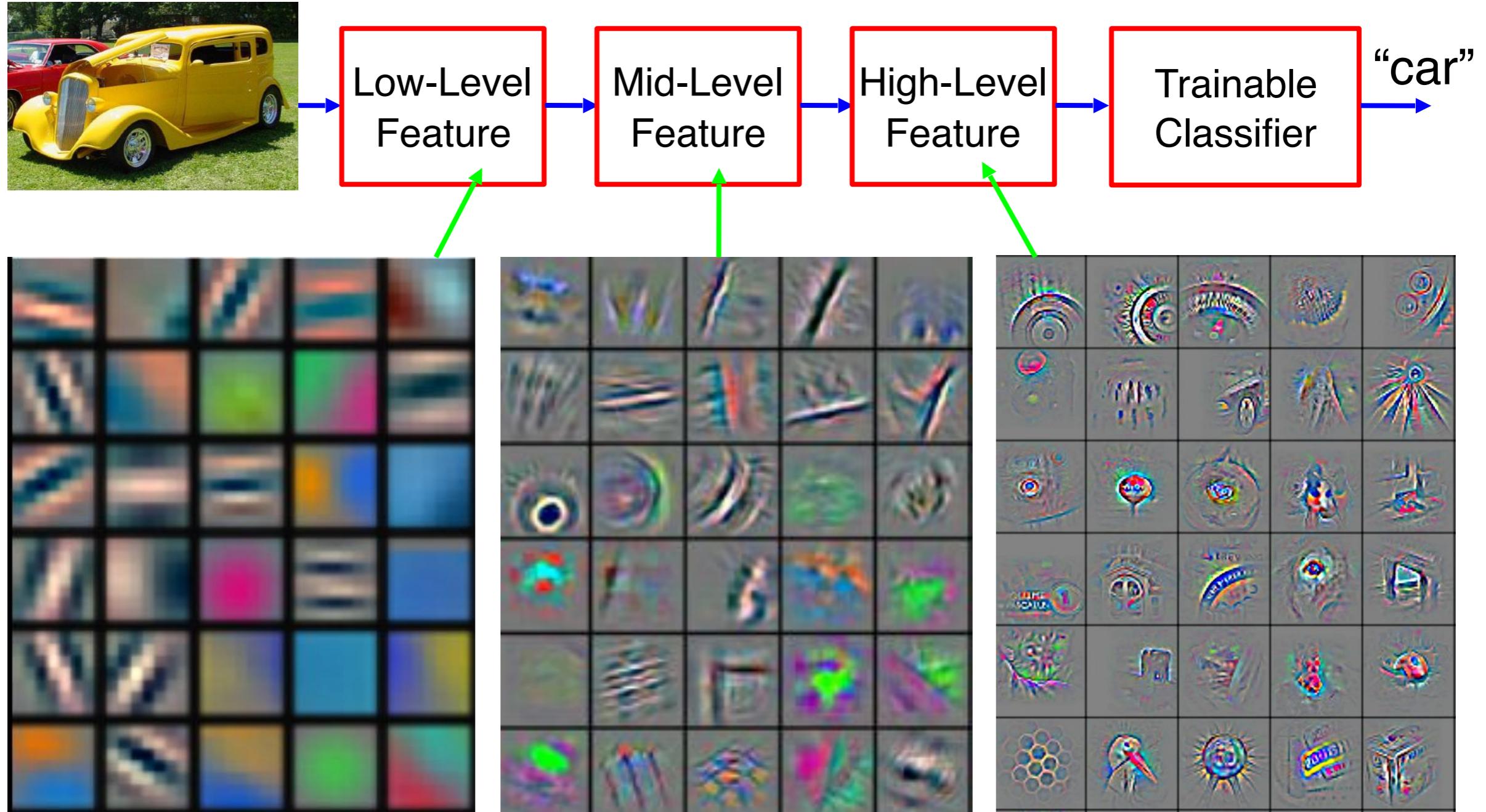
$$f(x) = \log(\cos(\exp(\sin^3(x))))$$



Deep Learning = Hierarchical Compositionality



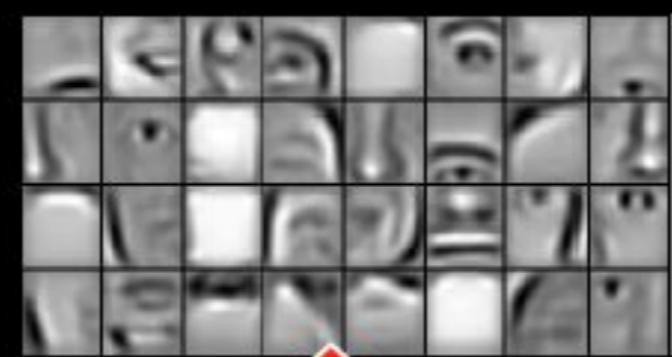
Deep Learning = Hierarchical Compositionality



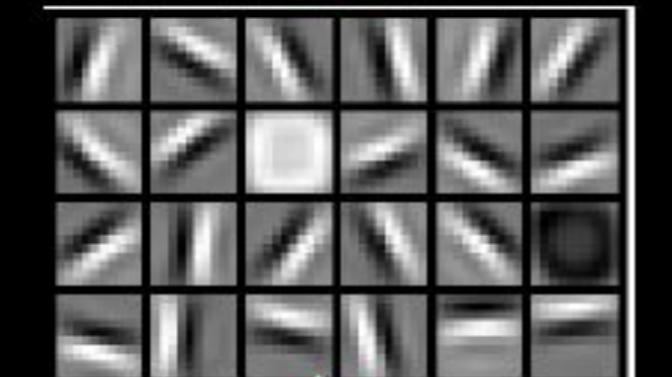
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



Face detectors



Face parts
(combination
of edges)



edges



pixels

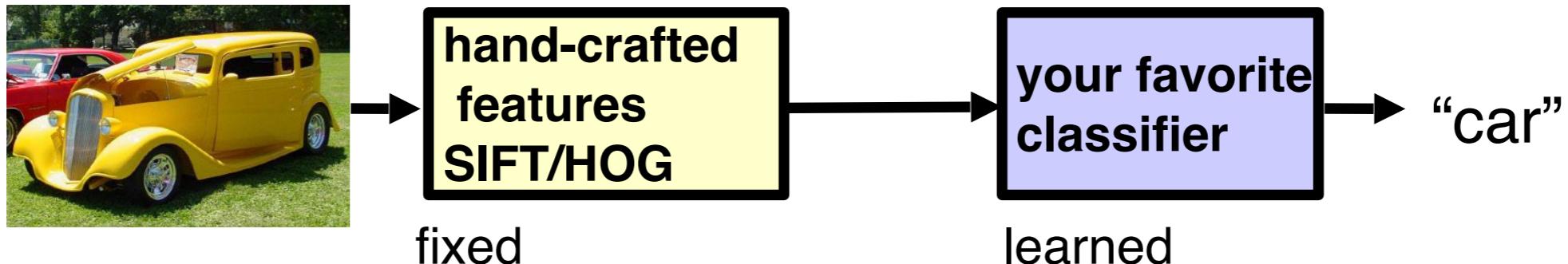
Sparse DBNs
[Lee et al. ICML '09]
Figure courtesy: Quoc Le

Three key ideas

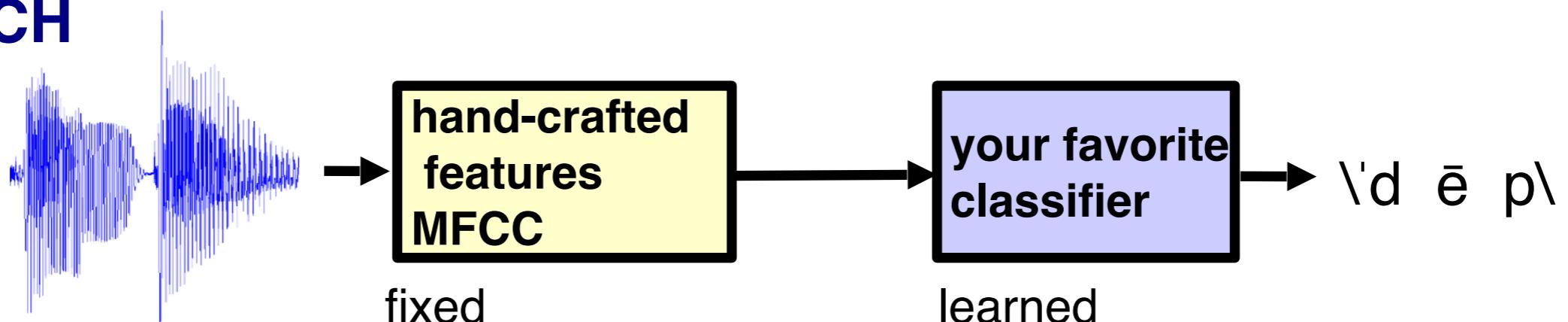
- (Hierarchical) Compositionality
 - Cascade of non-linear transformations
 - Multiple layers of representations
- **End-to-End Learning**
 - Learning (goal-driven) representations
 - Learning to feature extract
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons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

Traditional Machine Learning

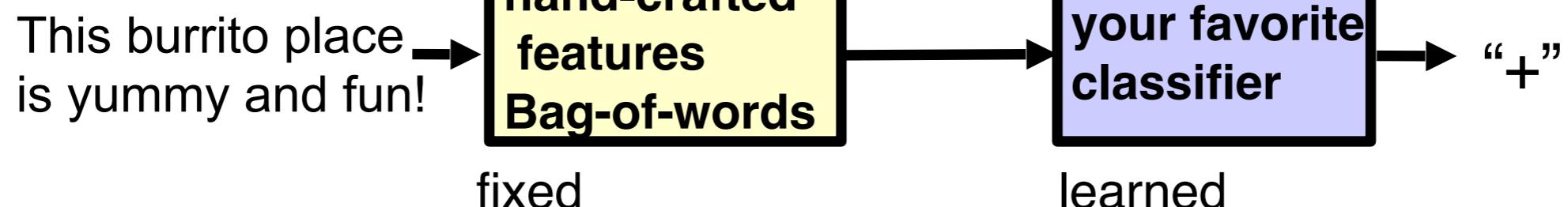
VISION



SPEECH

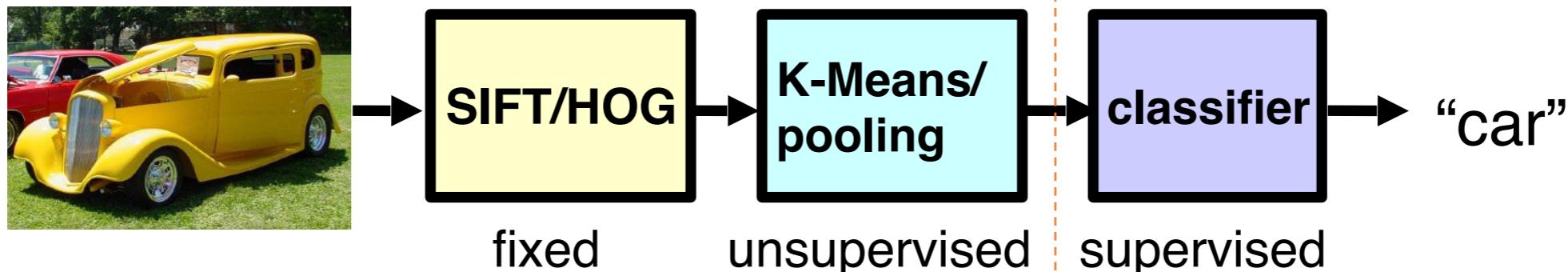


NLP

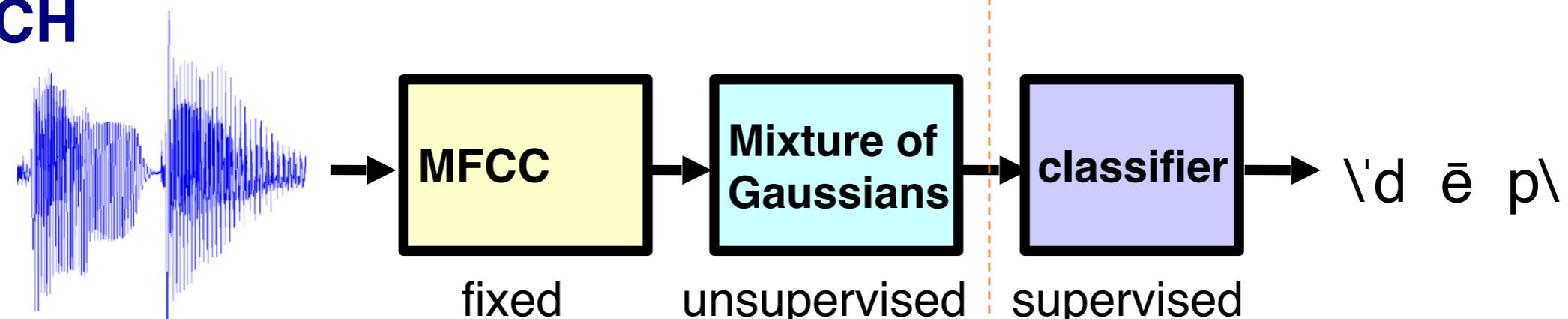


Traditional Machine Learning (more accurately)

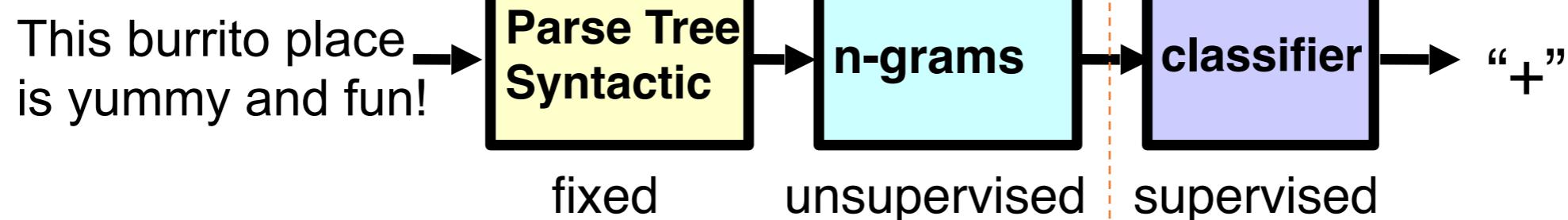
VISION



SPEECH

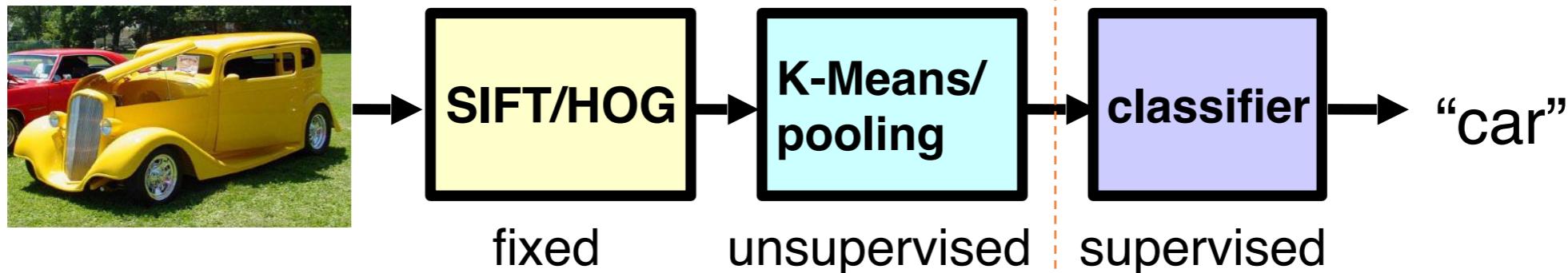


NLP

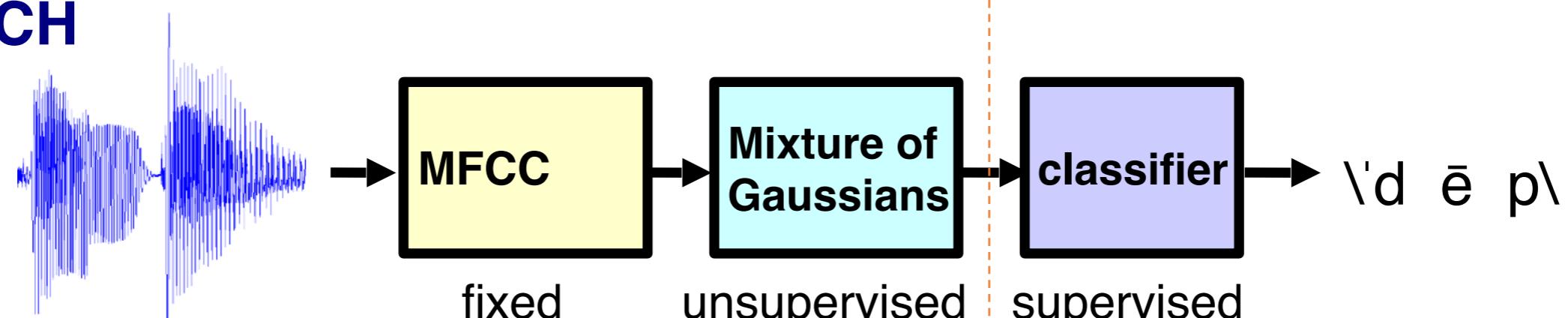


Deep Learning = End-to-End Learning

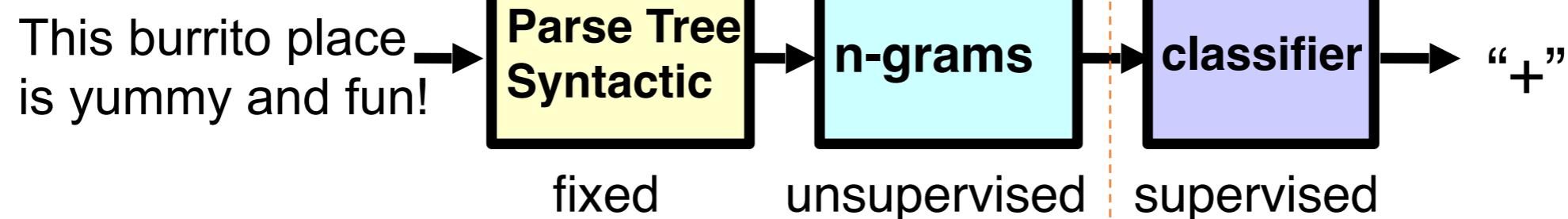
VISION



SPEECH

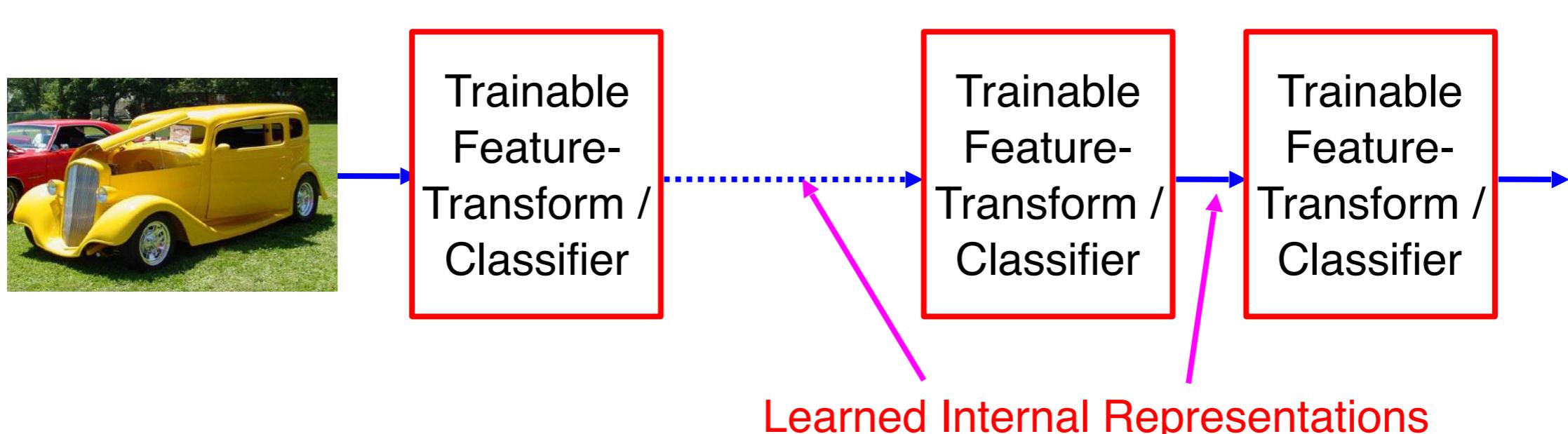


NLP



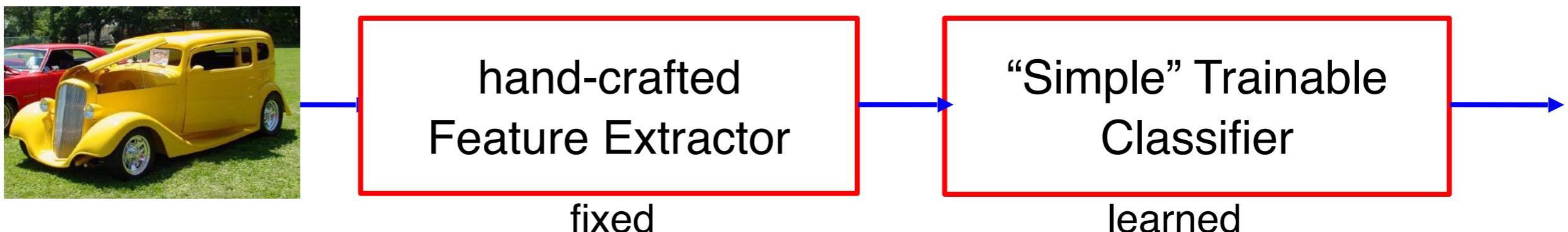
Deep Learning = End-to-End Learning

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

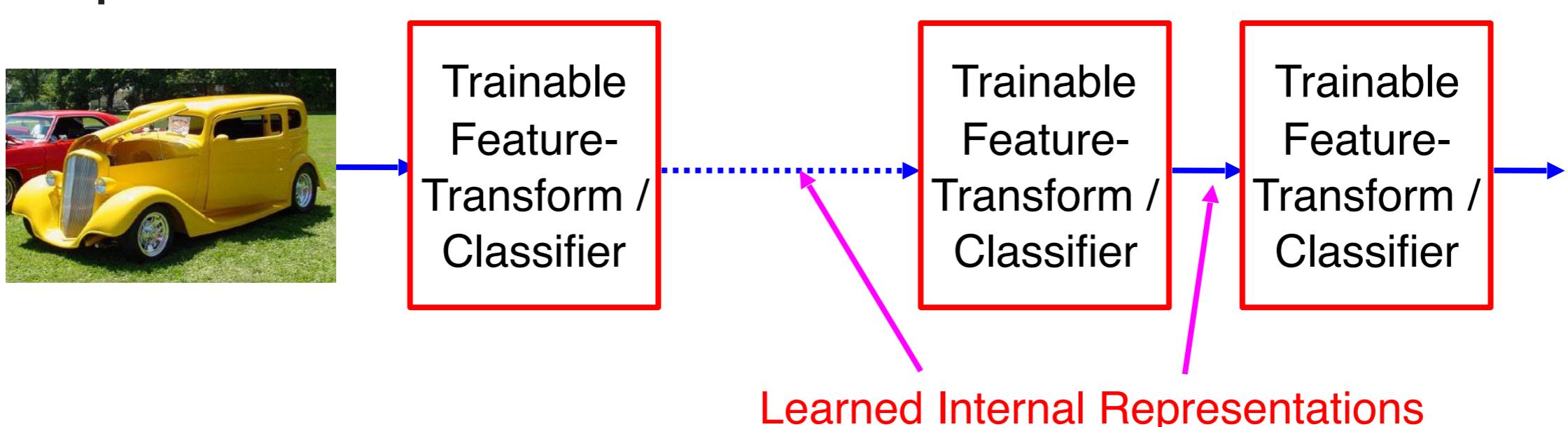


“Shallow” vs Deep Learning

- “Shallow” models



- Deep models



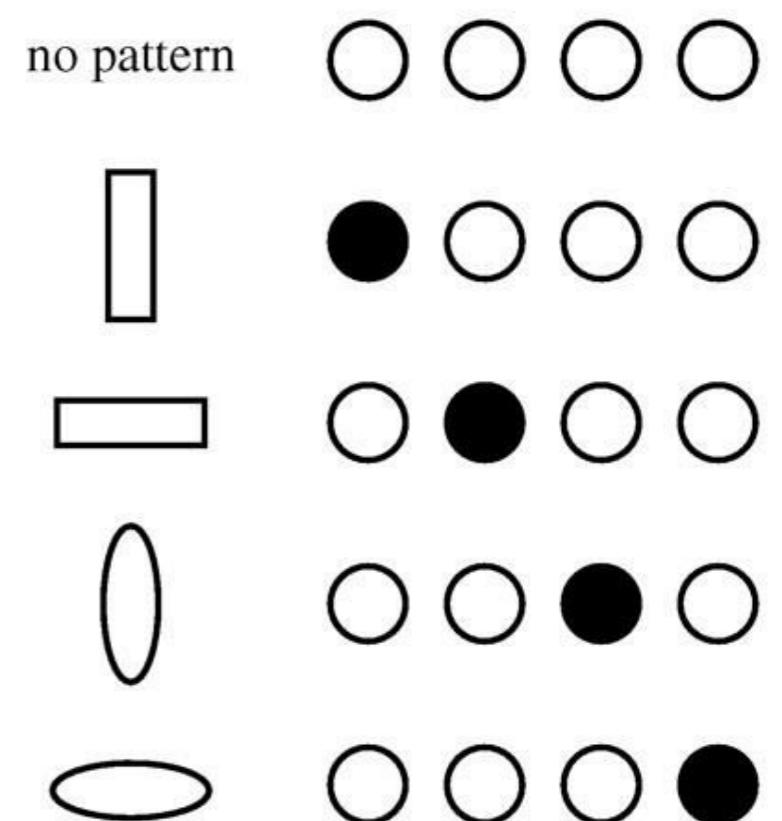
Three key ideas

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

Localist representations

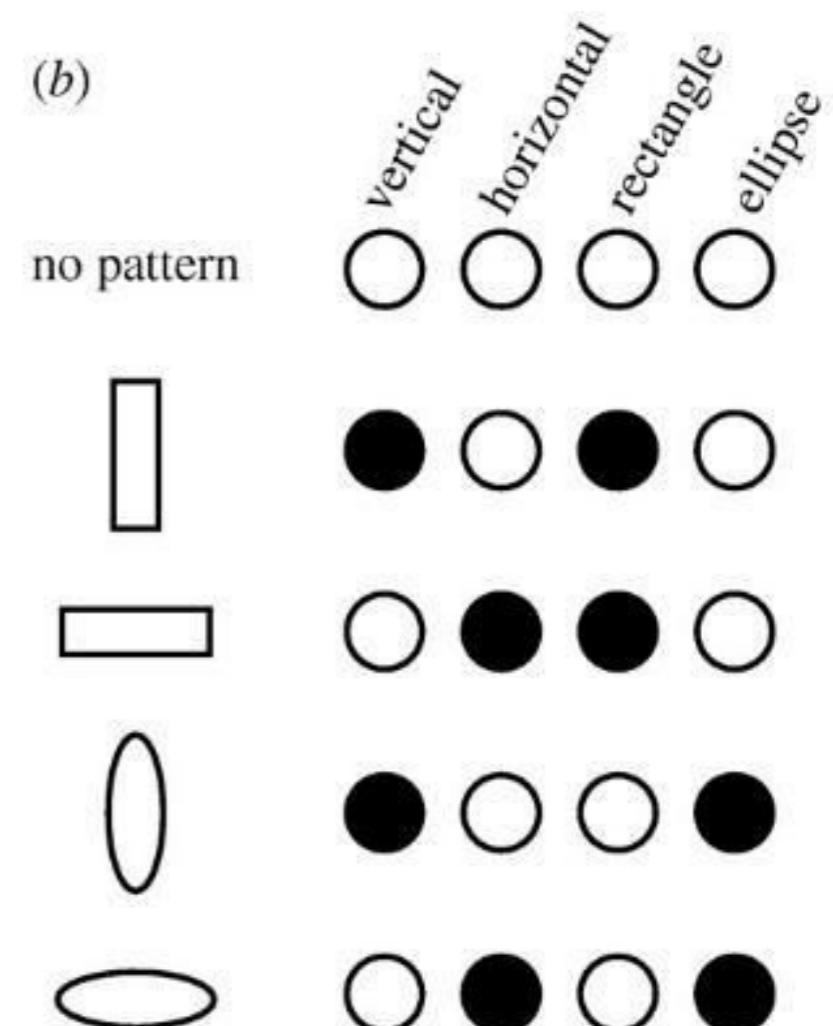
- The simplest way to represent things with neural networks is to **dedicate one neuron to each thing.** (a)

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



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

Next Lecture:

Convolutional Neural Networks