

# Neural Network and Deep Learning

## Lecture 2

Yanwei Fu  
School of Data Science, Fudan University



Fudan-SDS Confidential - Do Not Distribute

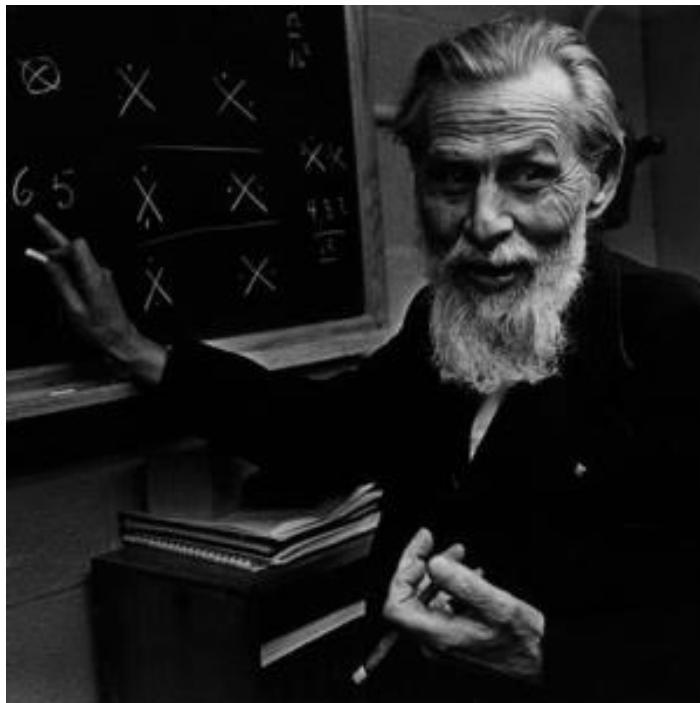




# 神经元进入计算领域的视野



Walter Pitts  
(1923-1969)



Warren McCulloch  
(1898-1969)

BULLETIN OF  
MATHEMATICAL BIOPHYSICS  
VOLUME 5, 1943

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

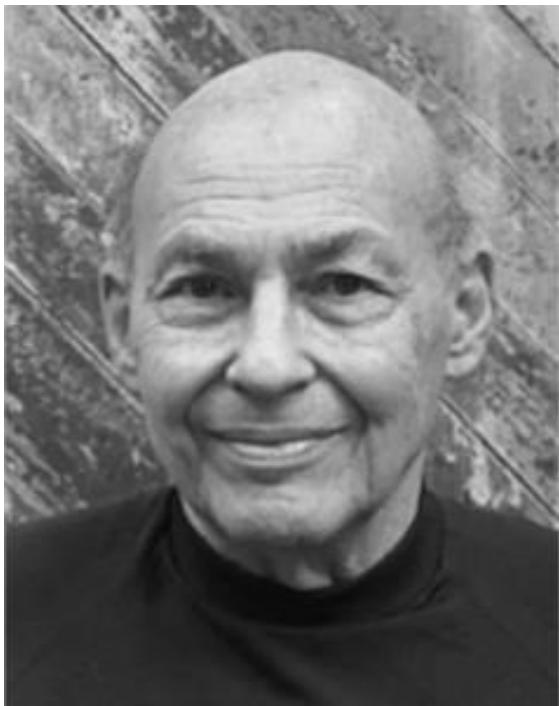
WARREN S. MCCULLOCH AND WALTER PITTS

FROM THE UNIVERSITY OF ILLINOIS, COLLEGE OF MEDICINE,  
DEPARTMENT OF PSYCHIATRY AT THE ILLINOIS NEUROPSYCHIATRIC INSTITUTE,  
AND THE UNIVERSITY OF CHICAGO

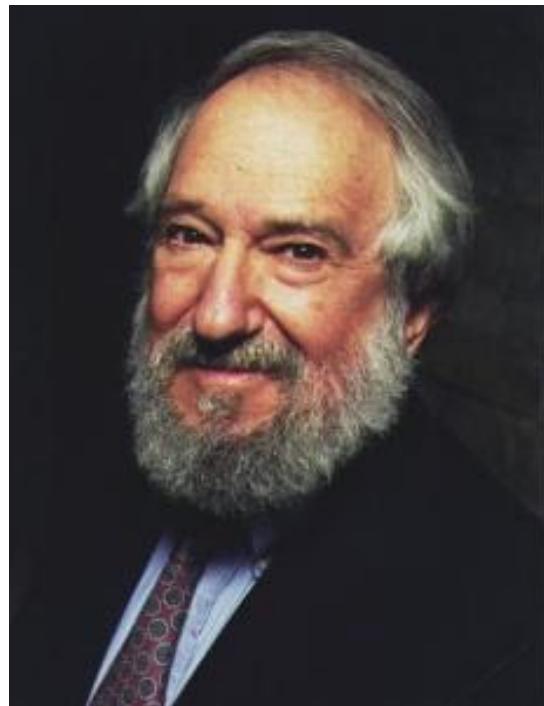
Because of the "all-or-none" character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every net can be described in these terms, with the addition of more complicated logical means for nets containing circles; 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.



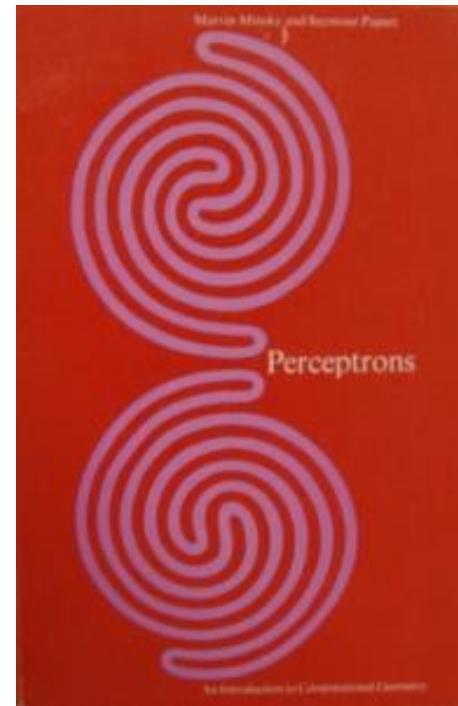
# 神经网络的第一次寒冬



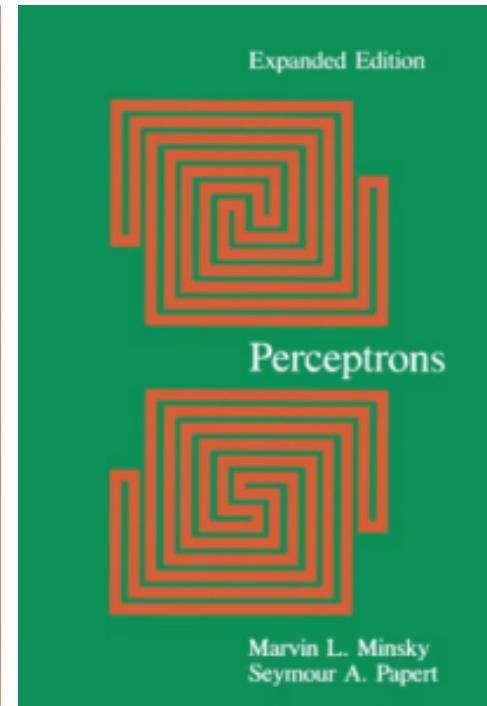
**Marvin Minsky**  
(1927-2016)



**Seymour Papert**  
(1928-)



1969年出版《Perceptrons》一书，认为仅靠局部连接的神经网络无法有效开展训练以及很多被后来的读者们以讹传讹的观点



# Perception: From Rosenblatt, 1962..

Optional subtitle

- "Perception, then, emerges as that relatively primitive, partly autonomous, institutionalized, ratiomorphic subsystem of cognition which achieves prompt and richly detailed orientation habitually concerning the vitally relevant, mostly distal aspects of the environment on the basis of mutually vicarious, relatively restricted and stereotyped, insufficient evidence in uncertainty-geared interaction and compromise, seemingly following the highest probability for smallness of error at the expense of the highest frequency of precision."
  - – From "Perception and the Representative Design of Psychological Experiments," by Egon Brunswik, 1956 (posthumous).
- "That's a simplification. Perception is standing on the sidewalk, watching all the girls go by."
  - – From "The New Yorker", December 19, 1959





# 反向传播算法（BP）的提出



Geoffrey Hinton  
(1947-)



David Rumelhart  
(1942-2011)

## Learning representations by back-propagating errors

David E. Rumelhart\*, Geoffrey E. Hinton†  
& Ronald J. Williams\*

\* Institute for Cognitive Science, C-015, University of California,  
San Diego, La Jolla, California 92093, USA

† Department of Computer Science, Carnegie-Mellon University,  
Pittsburgh, Philadelphia 15213, USA

1986年发表在Nature上的文章将BP算法用于神经网络  
模型，极大降低了计算量 ( $O(n^2) \rightarrow O(n)$ )

当时的计算机普及率、计算能力也远胜60年代



# “首个”重要应用：手写数字识别



Yann LeCun  
(1960-)

---

## *Handwritten Digit Recognition with a Back-Propagation Network*

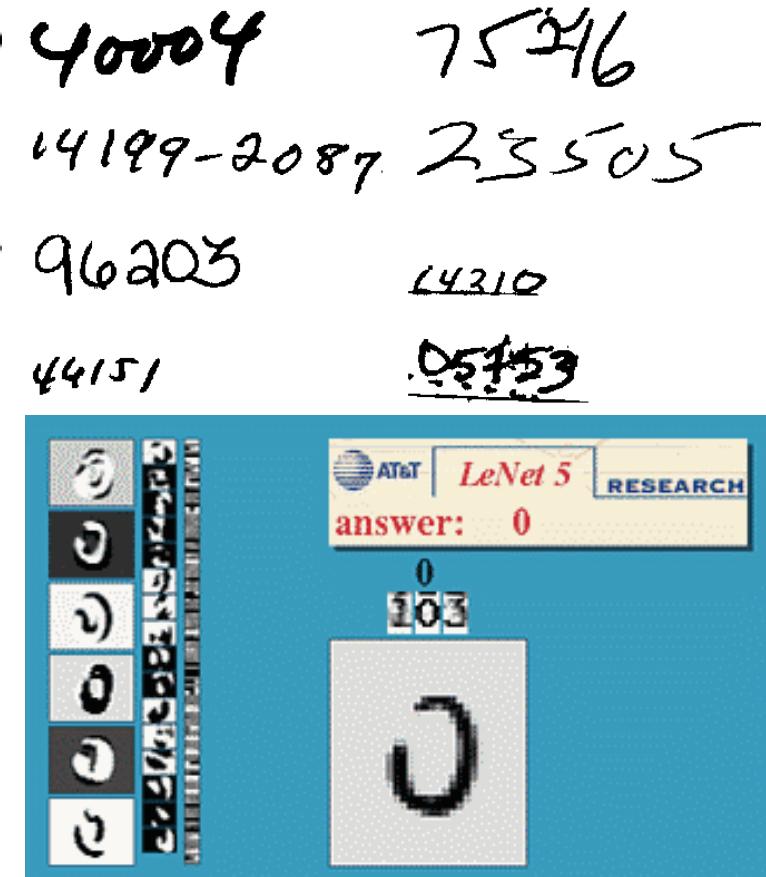
---

Y. Le Cun, B. Boser, J. S. Denker, D. Henderson,  
R. E. Howard, W. Hubbard, and L. D. Jackel  
AT&T Bell Laboratories, Holmdel, N. J. 07733

### ABSTRACT

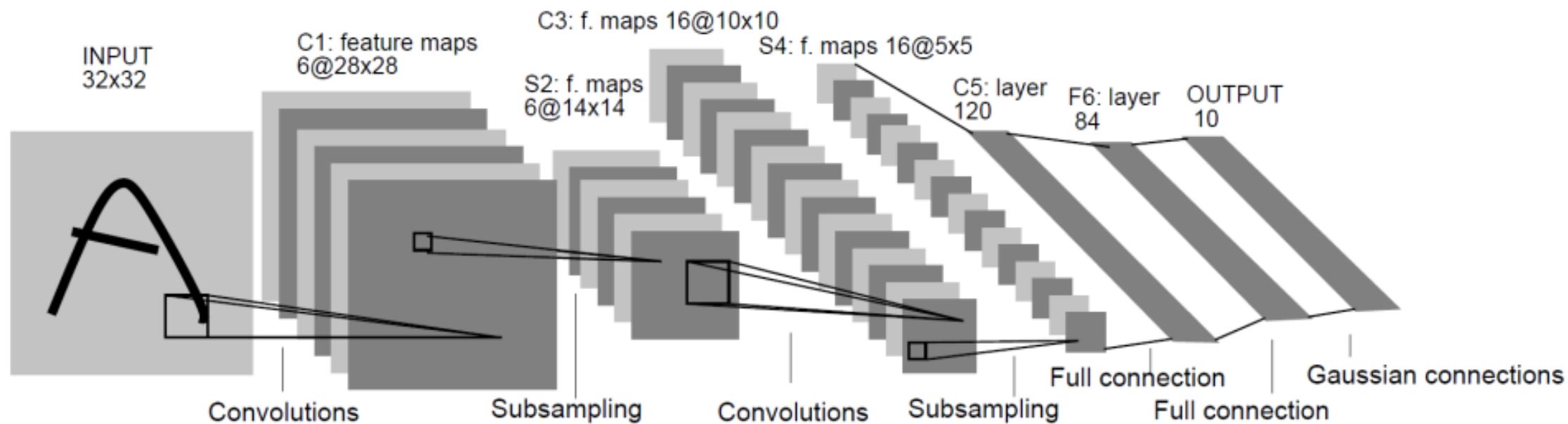
We present an application of back-propagation networks to handwritten digit recognition. Minimal preprocessing of the data was required, but architecture of the network was highly constrained and specifically designed for the task. The input of the network consists of normalized images of isolated digits. The method has 1% error rate and about a 9% reject rate on zipcode digits provided by the U.S. Postal Service.

到上世纪九十年代末，超过10%的美国支票识别采用了相关技术





# LeNet-5



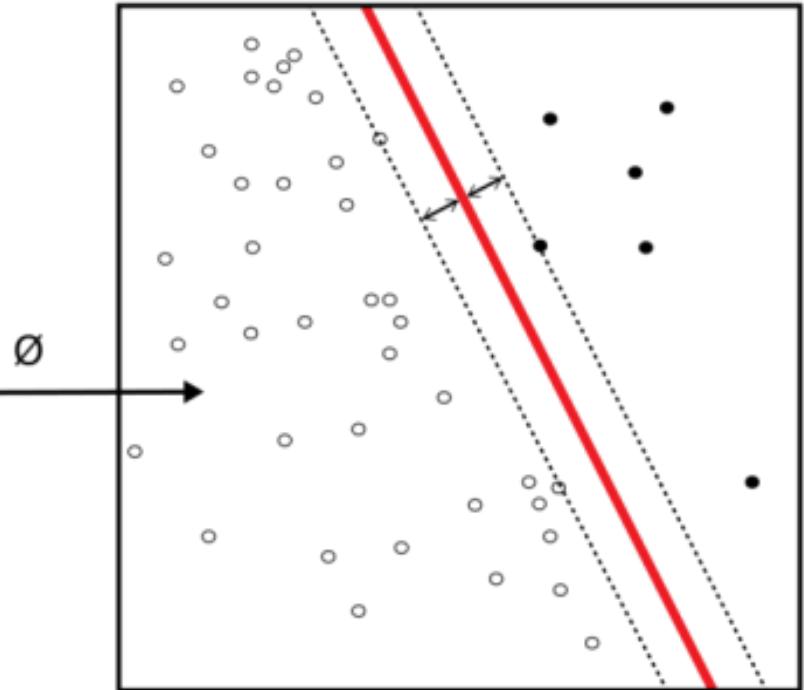
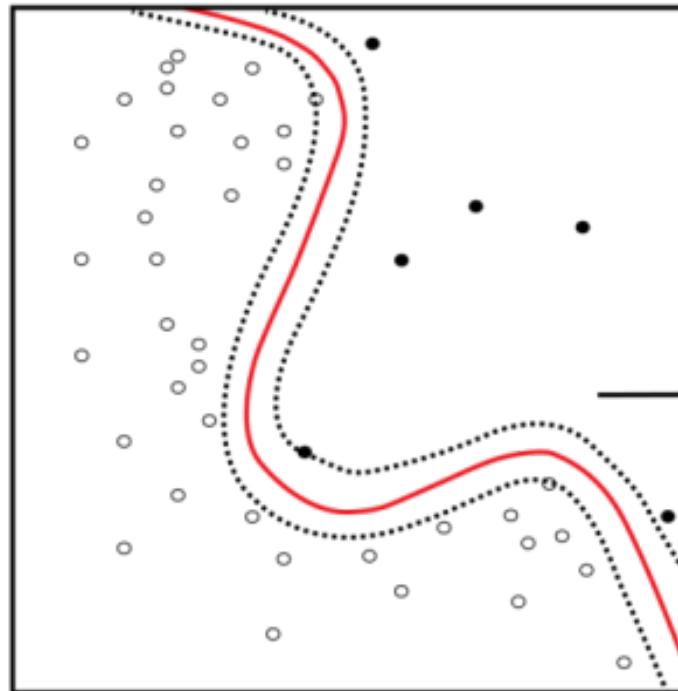
1998年发表的Gradient-based learning applied to Document Recognition一文中提出了LeNet-5  
已是Convolutional Neural Network的基本框架



# 神经网络的第二次寒冬



**Vladimir Vapnik**  
**(1936-)**

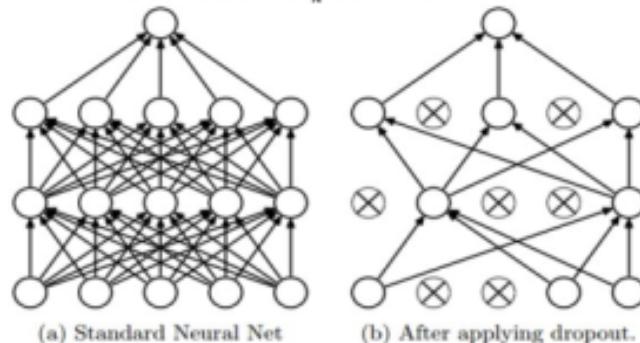
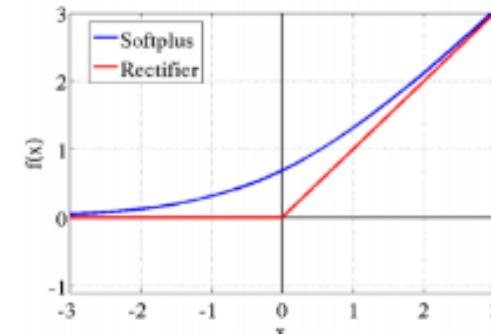


在1963年发表的论文中即提出Support Vector Machine的概念  
在2002年左右，不断改进的SVM方法将手写数字识别的错误率降到了0.56%  
不需要大量样本（事实上也难以支撑大量样本），速度相对可以接受

# 卷土重来的（深度）神经网络

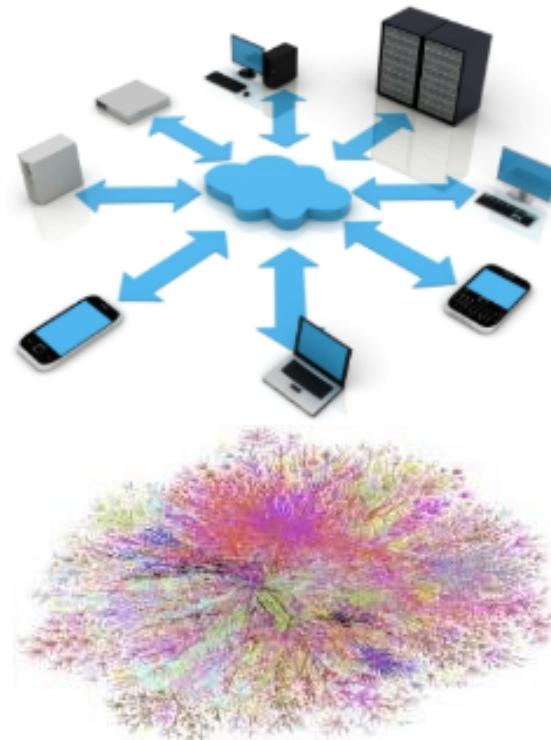
## 优化策略

梯度弥散问题通过ReLU、Dropout、Deep Residual Learning等方法得到缓解



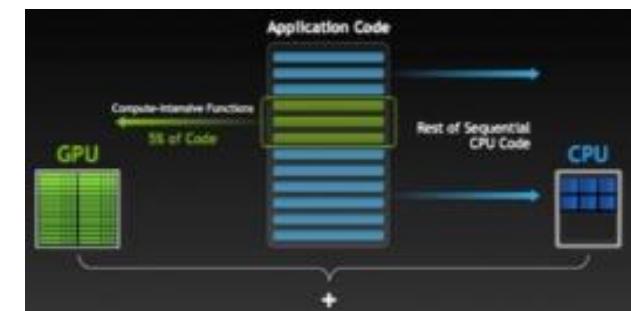
## 数据规模

真正意义上的大数据出现  
社会标注机制



## 计算能力

CPU、GPU的长足进步  
共享权值



# AlexNet

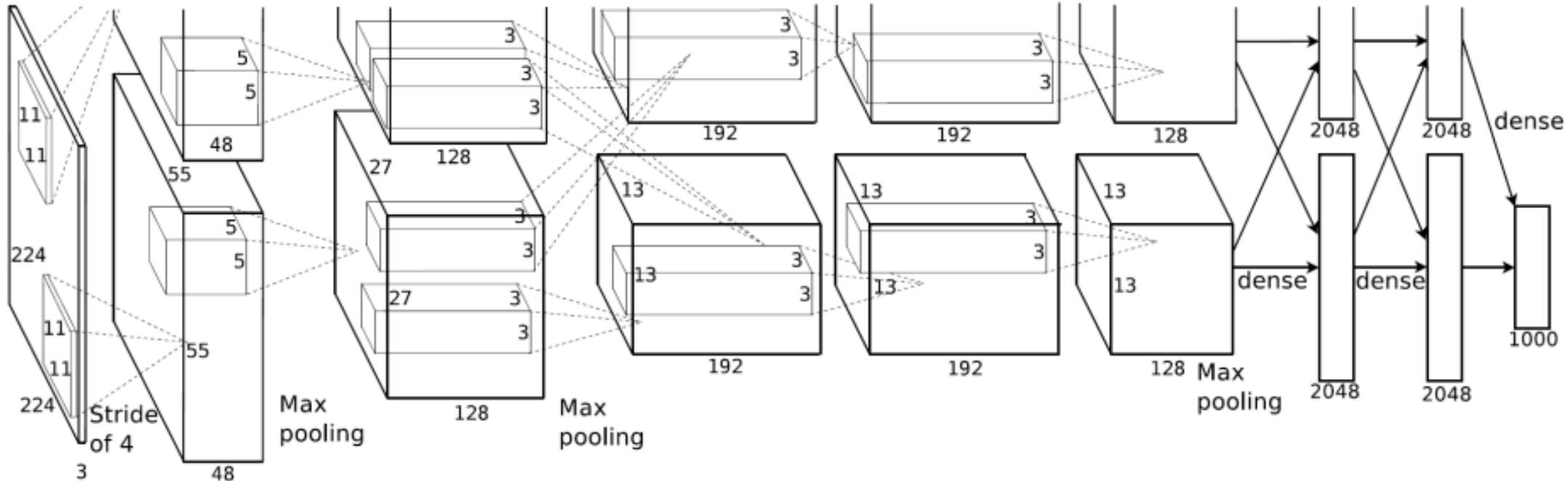
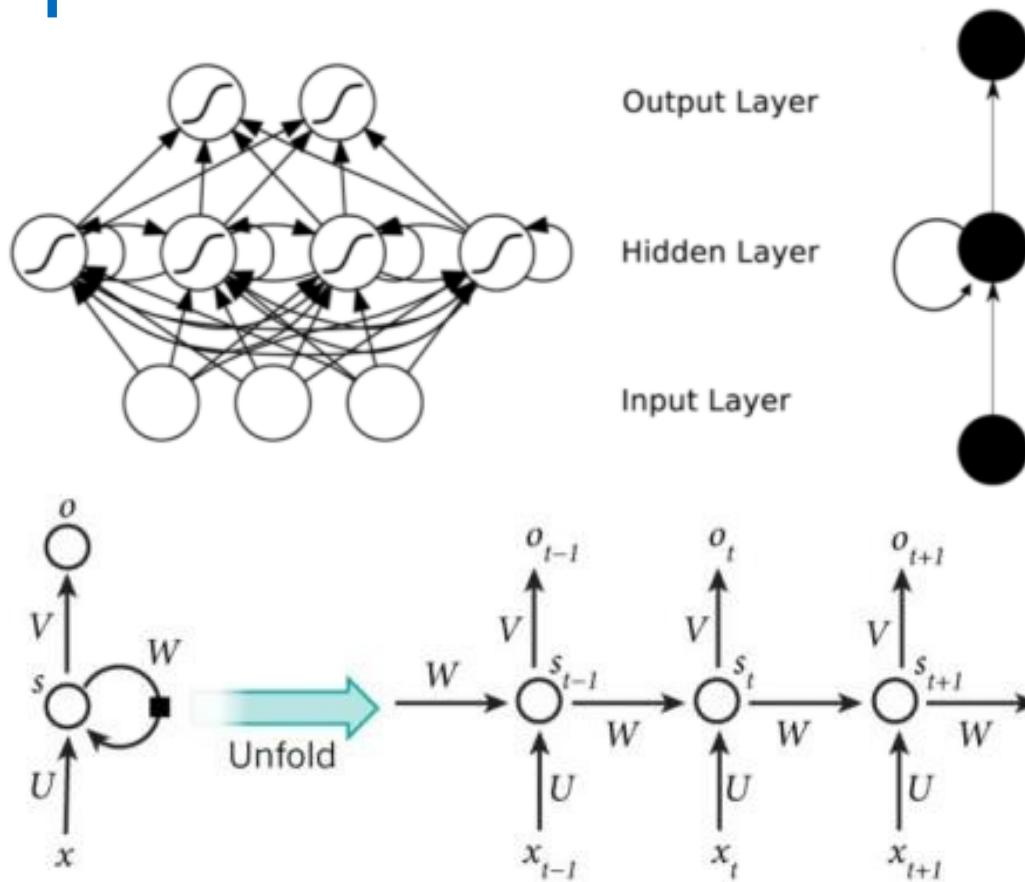


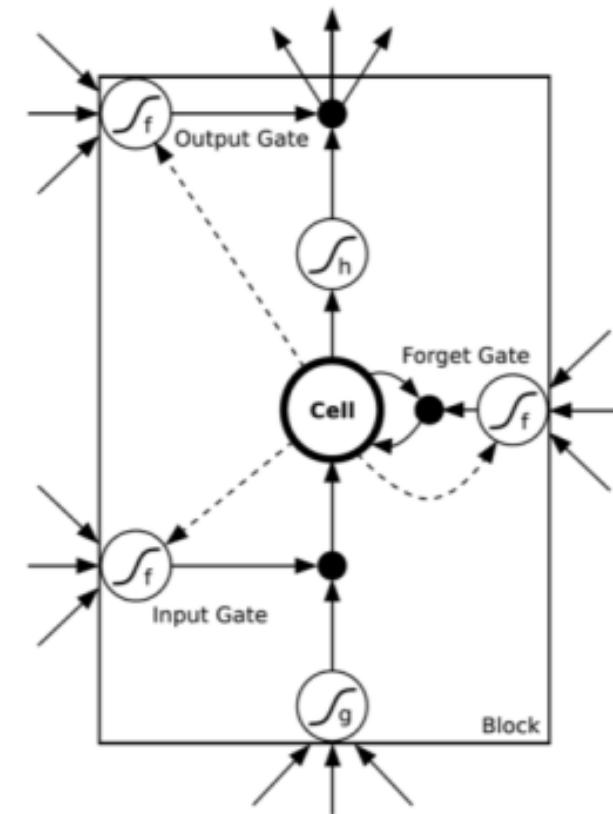
Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network’s input is 150,528-dimensional, and the number of neurons in the network’s remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

# 可处理时序数据的深度学习机制

## | Recursive Neural Networks



## | Long Short Term Memory





# 大型企业不同于以往的密切关注





# 开源框架

| 名称              | 基本语言  | 是否支持多GPU | 速度   | 应用领域 |
|-----------------|---|----------|------|------|
| TensorFlow      |  Python、C++  | 是        | ★★☆  | 通用   |
| Caffe           | C++   | 是        | ★★★☆ | 图像分类 |
| Torch           |  Lua         | 是        | ★★☆  | 通用   |
| Theano          |  Python    | 默认为否     | ★★★☆ | 通用   |
| IBM DL Platform |  SystemML等 | 是        | ?    | 通用   |

# Breakthroughs with neural networks

Optional subtitle

TECHNEWSWORLD      EMERGING TECH

SEARCH

Computing Internet IT Mobile Tech Reviews Security Technology Tech Blog Reader Services

## Microsoft AI Beats Humans at Speech Recognition

By Richard Adhikari  
Oct 20, 2016 11:40 AM PT

G+ 5  
Twitter 25  
Facebook Share 45  
LinkedIn Share 11  
Reddit Share 0  
StumbleUpon share 104



Image: Adobe Stock

Print Email

Most Popular Newsletters News Alerts

How do you feel about Black Friday and Cyber Monday?

- They're great -- I get a lot of bargains!
- The deals are too spread out -- I'd prefer just one day.
- They're a fun way to kick off the holiday season.
- I don't like the commercialization of Thanksgiving Day.
- They're crucial for the retail industry and the economy.
- The deals typically aren't that good.

Vote to See Results

### E-Commerce Times

Black Friday Shoppers Hungry for New Experiences, New Tech

Pay TV's Newest Innovation: Giving Users Control

Apple Celebrates Itself in \$300 Coffee Table Tome

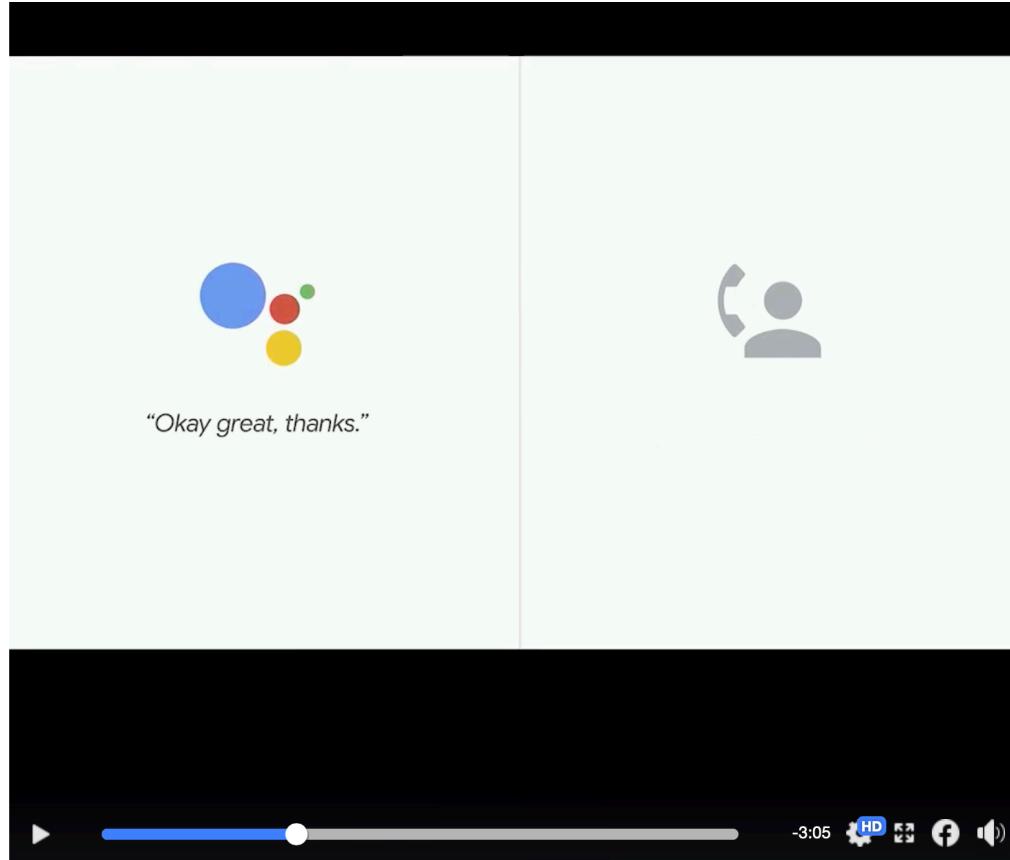
AWS Enjoys Top Perch in IaaS, PaaS Markets

US Comptroller Gears Up for Blockchain and



# Google just gave a stunning demo of Assistant making an actual phone call

<https://www.theverge.com/2018/5/8/17332070/google-assistant-makes-phone-call-demo-duplex-io-2018>



Onstage at I/O 2018, Google showed off a jaw-dropping new capability of Google Assistant: in the not too distant future, it's going to make phone calls on your behalf. CEO Sundar Pichai played back a phone call recording that he said was placed by the Assistant to a



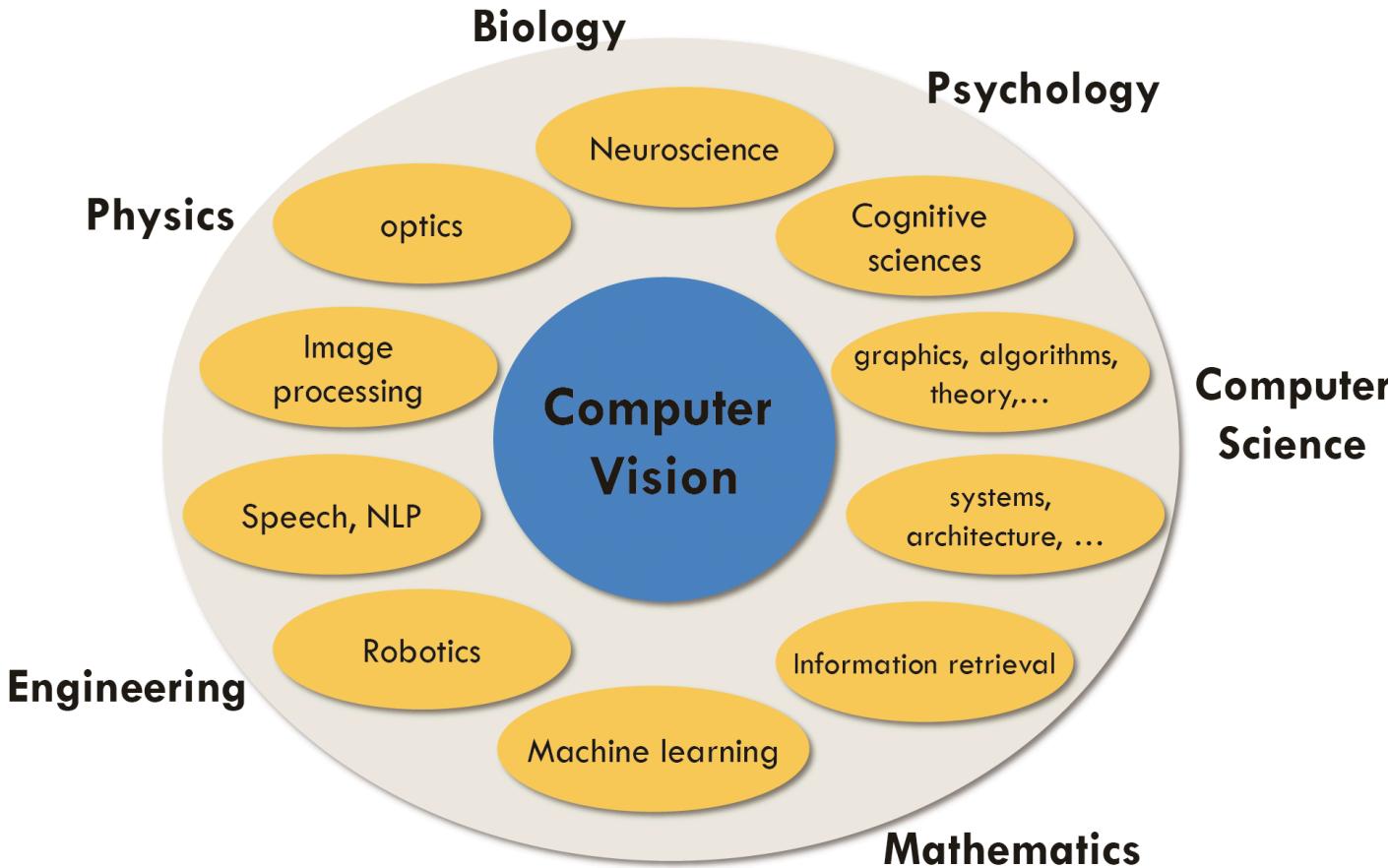
Fudan-SDS Confidential - Do Not Distribute

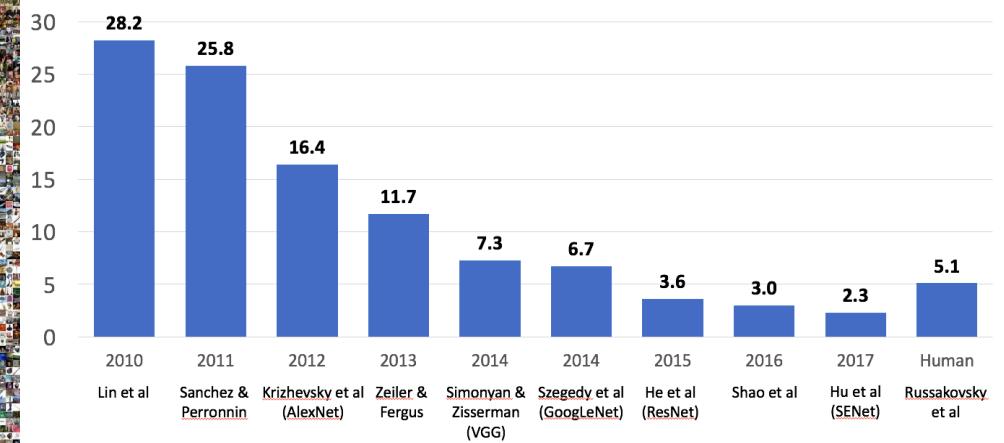
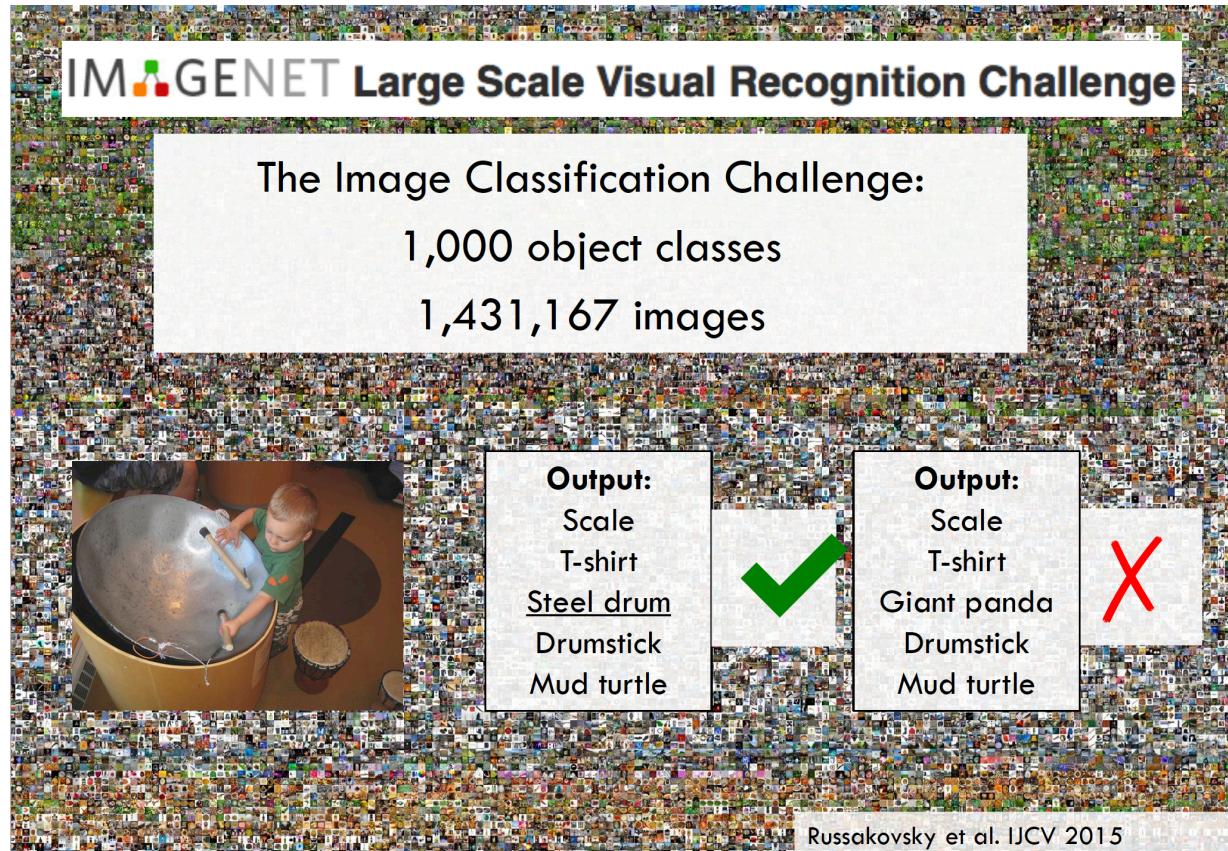
# Breakthroughs with neural networks

Optional subtitle

The screenshot shows a news article from 'The Keyword'. The top navigation bar includes the Google logo, 'The Keyword', 'Latest Stories', 'Product News', 'Topics', a search icon, and a more options icon. Below the header, there's a 'TRANSLATE' button and the date 'NOV 15, 2016'. The main title of the article is 'Found in translation: More accurate, fluent sentences in Google Translate'. Below the title, the author is listed as 'Barak Turovsky' and 'PRODUCT LEAD, GOOGLE TRANSLATE'. A short summary at the bottom states: 'In 10 years, Google Translate has gone from supporting just a few languages to 103, connecting strangers, reaching across language barriers and even helping'. There is a blue circular arrow icon in the bottom right corner.



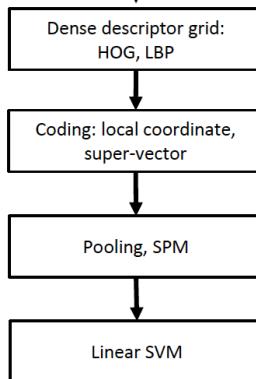




# IMAGENET Large Scale Visual Recognition Challenge

## Year 2010

NEC-UIUC

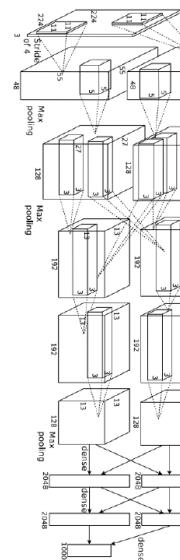


[Lin CVPR 2011]

[Lion image](#) by Swissfrog is licensed under [CC BY 3.0](#)

## Year 2012

SuperVision



[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

## Year 2014

GoogLeNet

- Pooling
- Convolution
- Softmax
- Other



[Szegedy arxiv 2014]

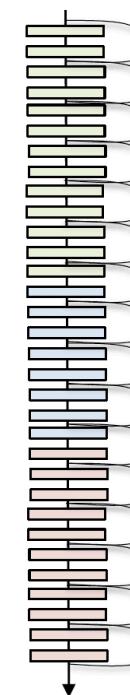
VGG



[Simonyan arxiv 2014]

## Year 2015

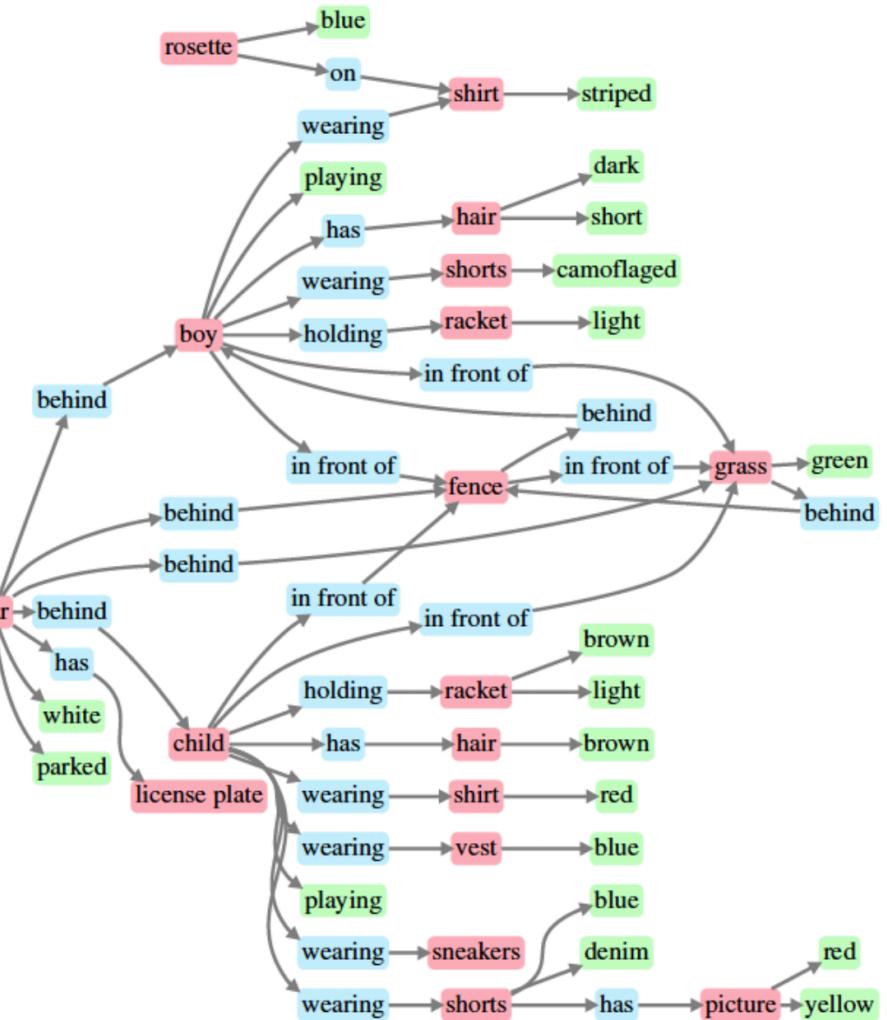
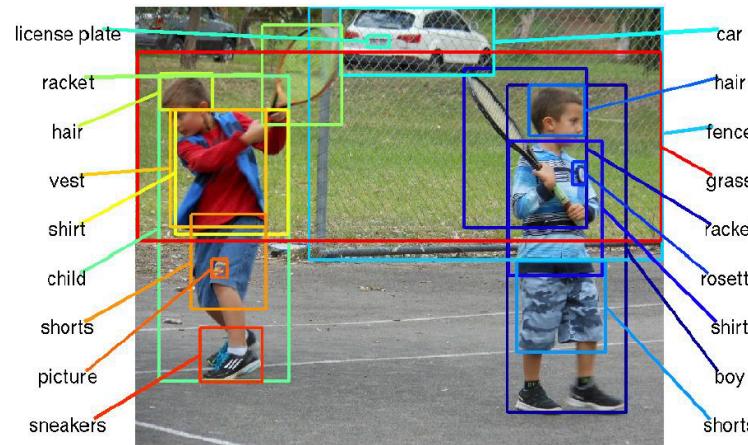
MSRA



[He ICCV 2015]

# Ingredients for Deep Learning



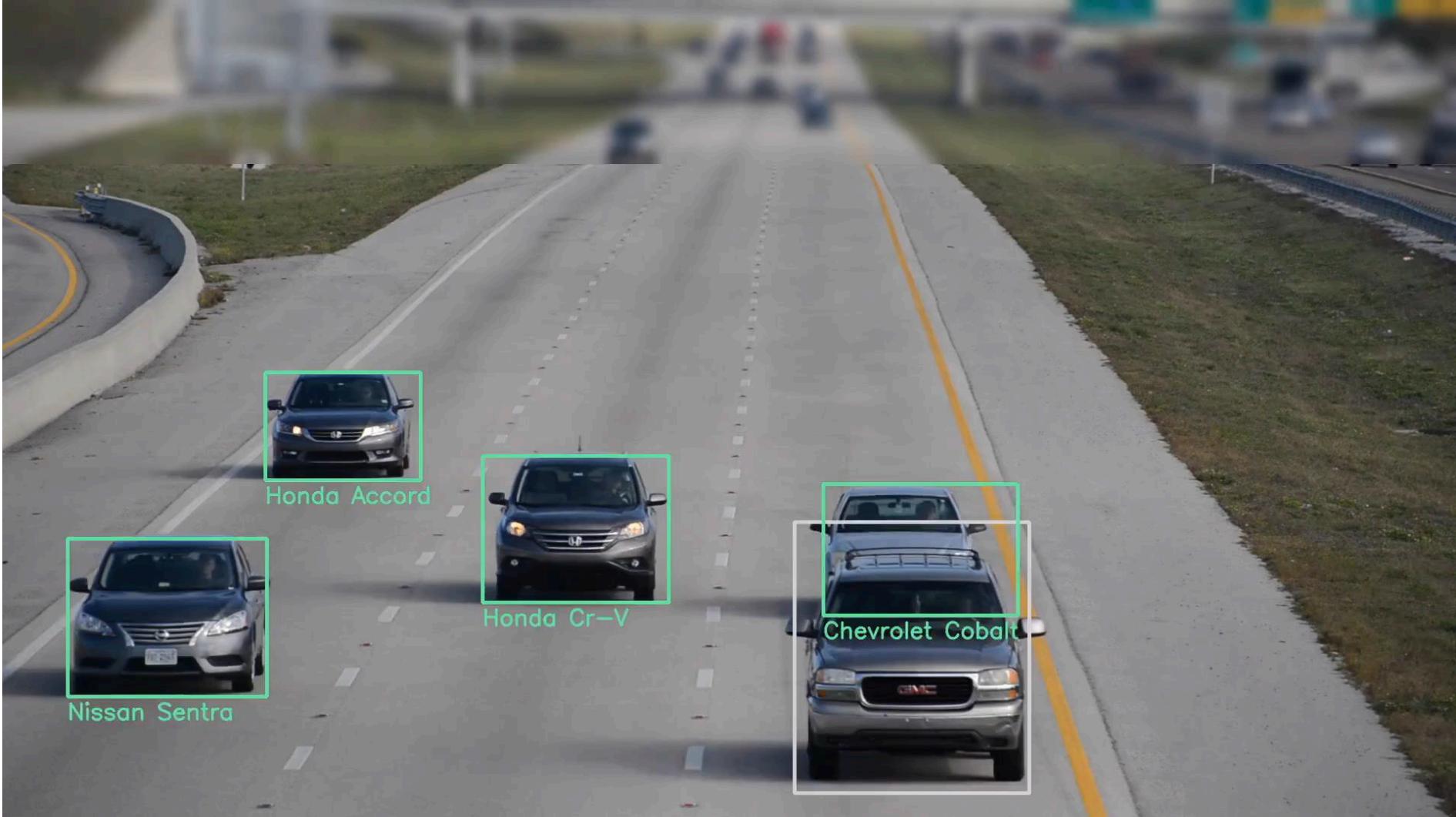


Johnson *et al.*, “Image Retrieval using Scene Graphs”, CVPR 2015

Figures copyright IEEE, 2015. Reproduced for educational purposes

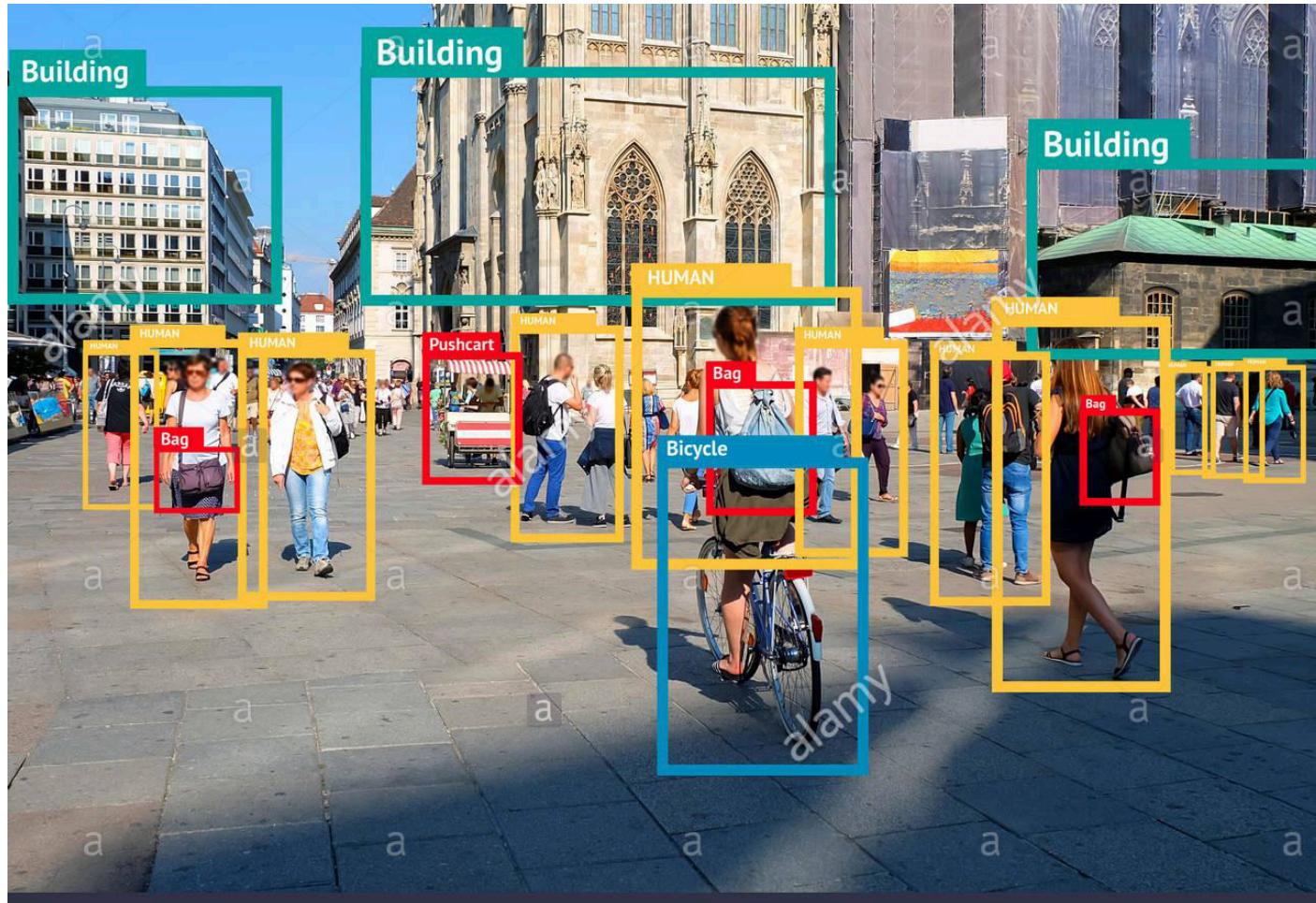
# Image Recognition

Optional subtitle



# Image segmentation & recognition

Optional subtitle



# Breakthroughs with neural networks

Optional subtitle

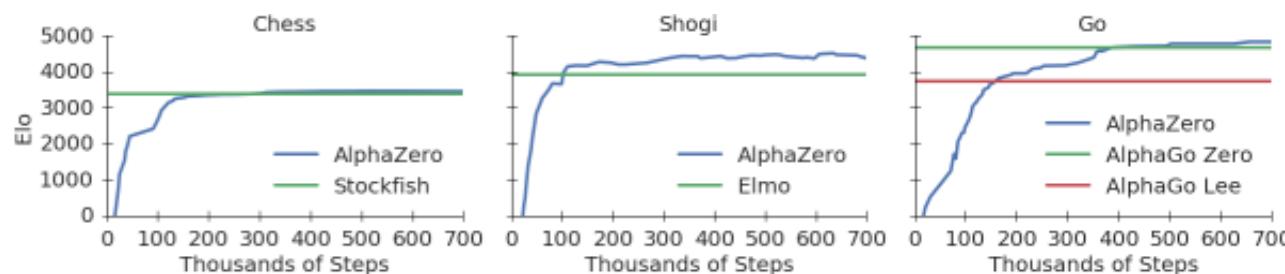


Figure 1: Training *AlphaZero* for 700,000 steps. Elo ratings were computed from evaluation games between different players when given one second per move. **a** Performance of *AlphaZero* in chess, compared to 2016 TCEC world-champion program *Stockfish*. **b** Performance of *AlphaZero* in shogi, compared to 2017 CSA world-champion program *Elmo*. **c** Performance of *AlphaZero* in Go, compared to *AlphaGo Lee* and *AlphaGo Zero* (20 block / 3 day) (29).

# Breakthroughs with neural networks

Optional subtitle



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."

# Neural Network In General



## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

*Face*



*Face*



*Not a face*



Examples: Linear regression,  
Logistic regression, Neural  
Network

Examples: Mean-squared error,  
Cross Entropy

## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps  
opposite the gradient)

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

## Background

1. Given training data

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

– Loss function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

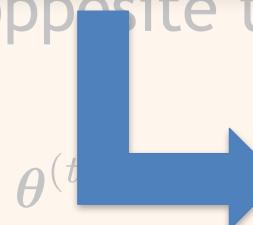
## A Recipe for

# Gradients

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!

(opposite the gradient)



$$\theta^{(t)} \leftarrow \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

B

A Recipe for

## Goals for This Course

1. **Explore a new class of decision functions  
(Neural Networks)**
2. Consider variants of this recipe for training

2. Choose each of these:

– Decision function

$$\hat{y} = f_{\theta}(x_i)$$

– Loss function

$$\ell(\hat{y}, y_i) \in \mathbb{R}$$

4. Train with SGD:  
(take small steps  
opposite the gradient)

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(x_i), y_i)$$

# Representation Learning



# Causal Inference in Deep Learning?

Optional subtitle



Judea Pearl, 2011年,  
因通过概率和因果推理的算法研发  
在人工智能取得的杰出贡献而获得图灵奖

“无人问津”的贝叶斯网络之父Judea Pearl在NIPS 2017上到底报告了啥



Fudan-SDS Confidential - Do Not Distribute



# Representations Matter

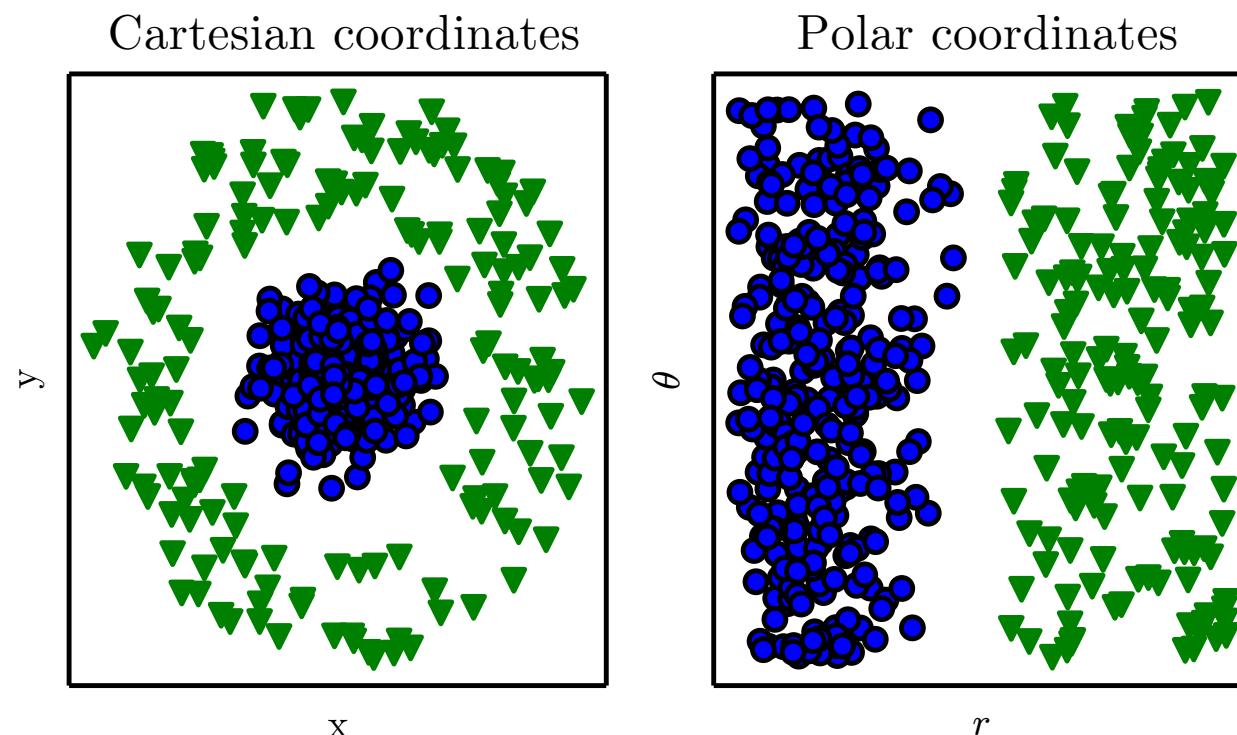


Figure 1.1

(Goodfellow 2016)

# Depth: Repeated Composition

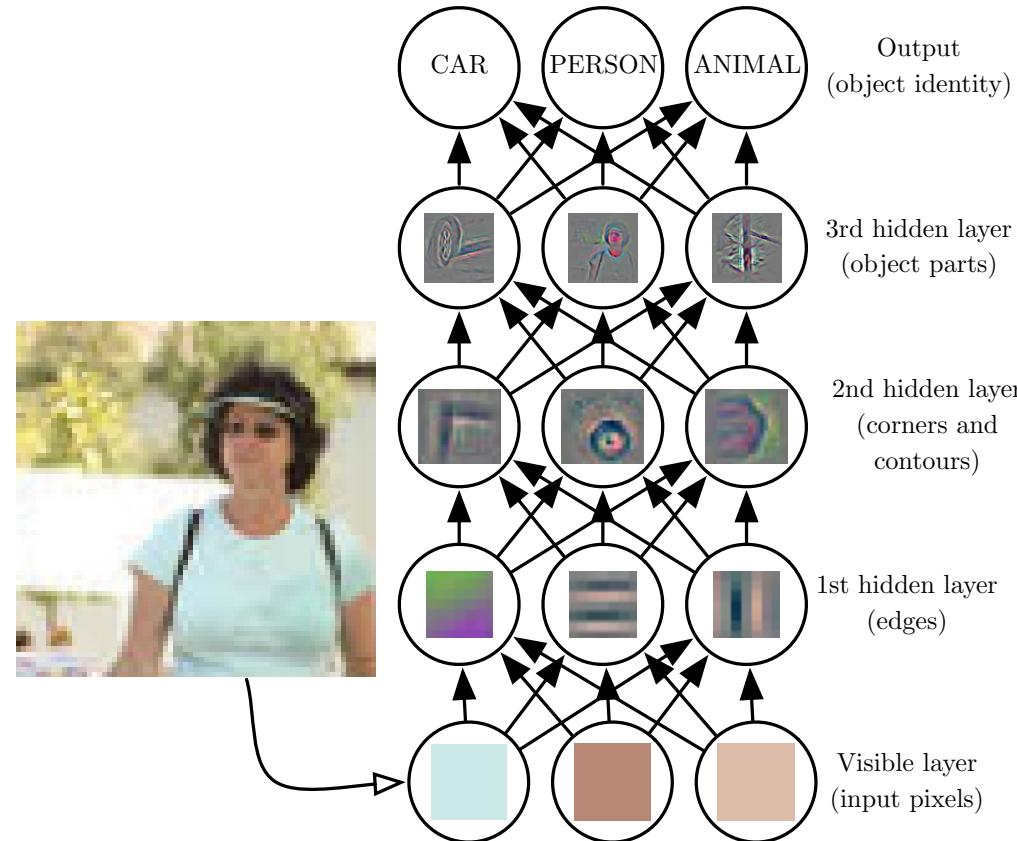


Figure 1.2

(Goodfellow 2016)

# Computational Graphs

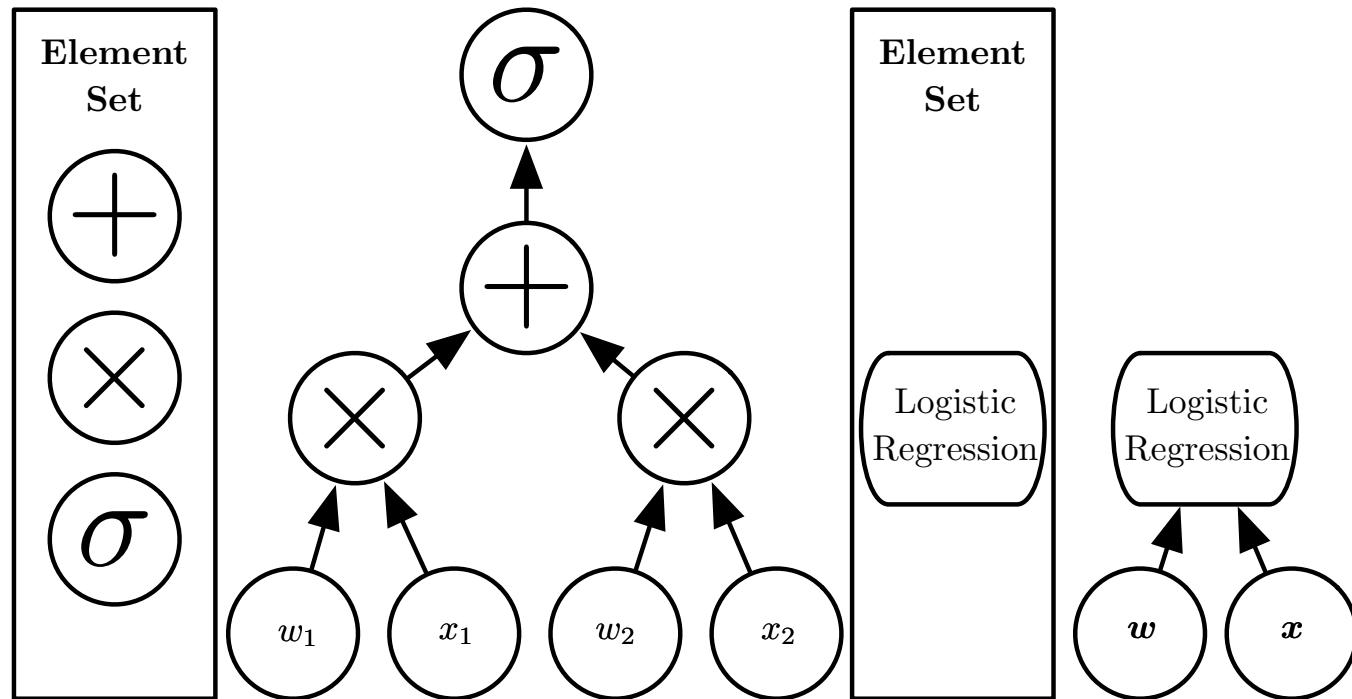


Figure 1.3

(Goodfellow 2016)

# Machine Learning and AI

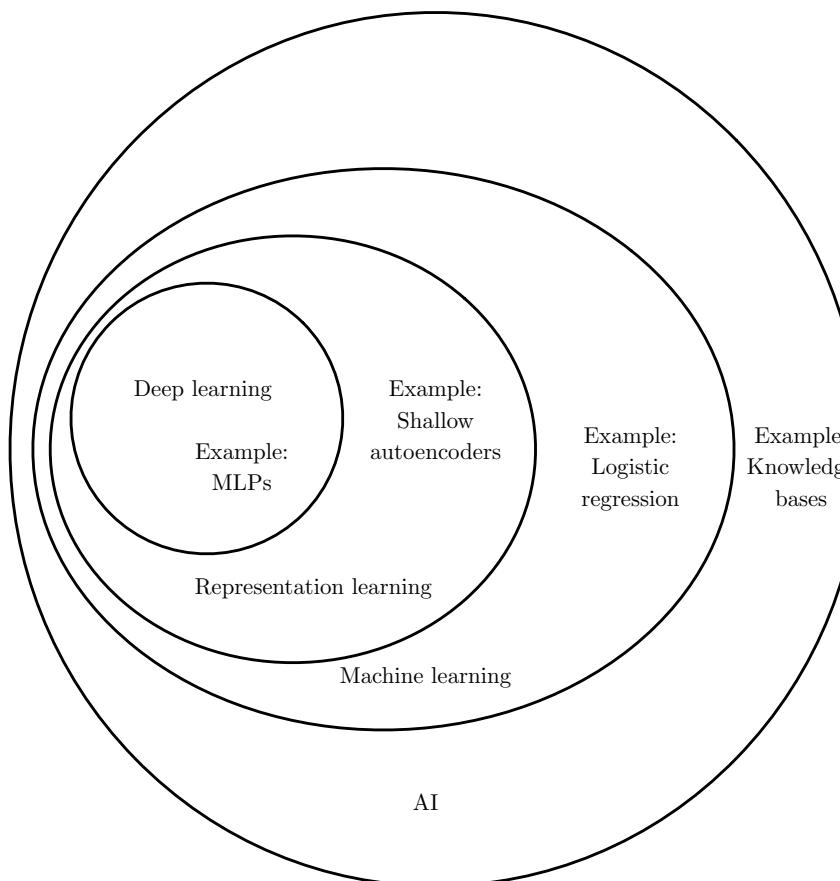
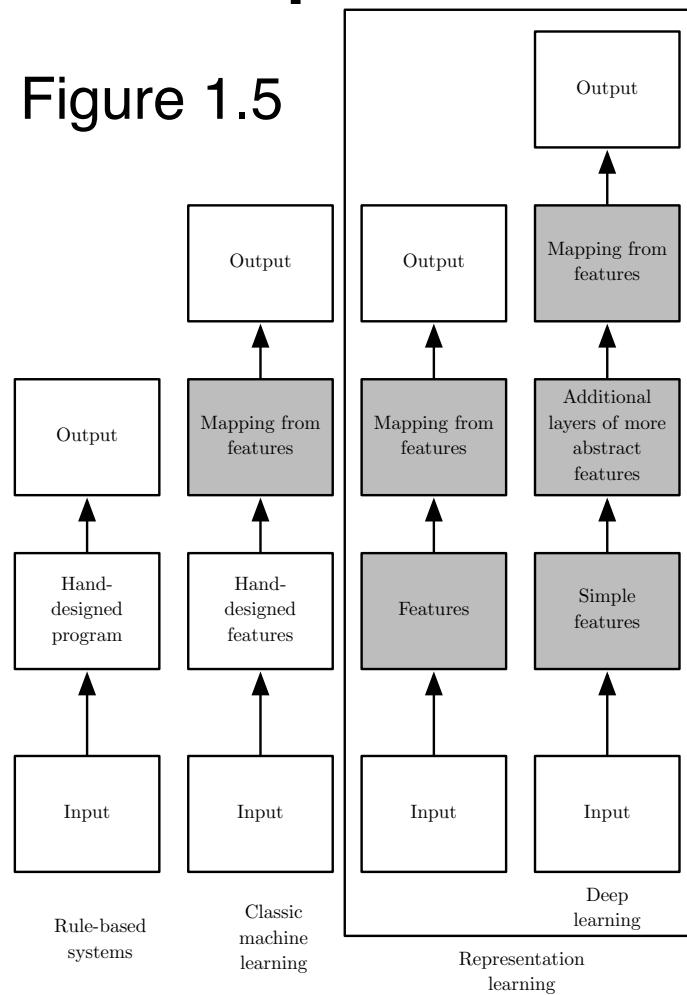


Figure 1.4

(Goodfellow 2016)

# Learning Multiple Components

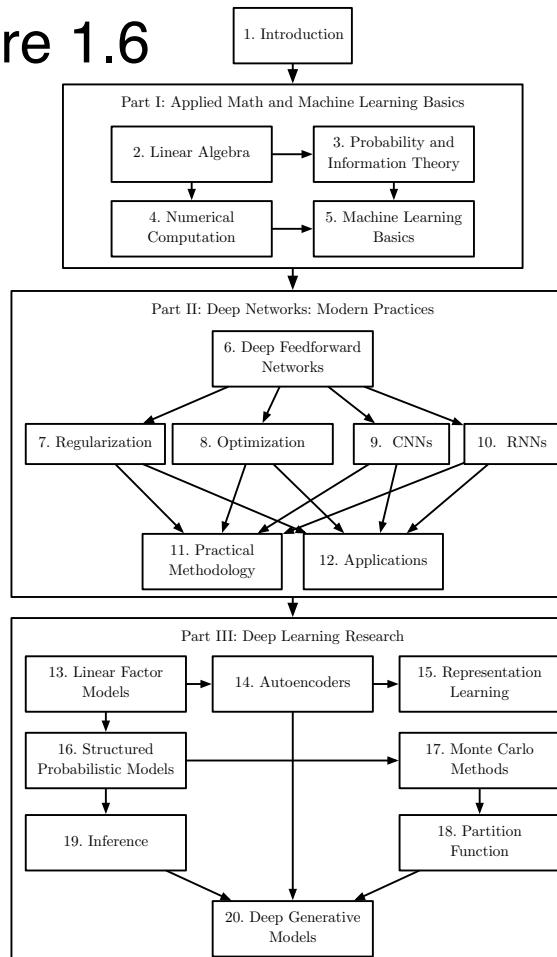
Figure 1.5



(Goodfellow 2016)

# Organization of the Book

Figure 1.6



(Goodfellow 2016)

# Historical Waves

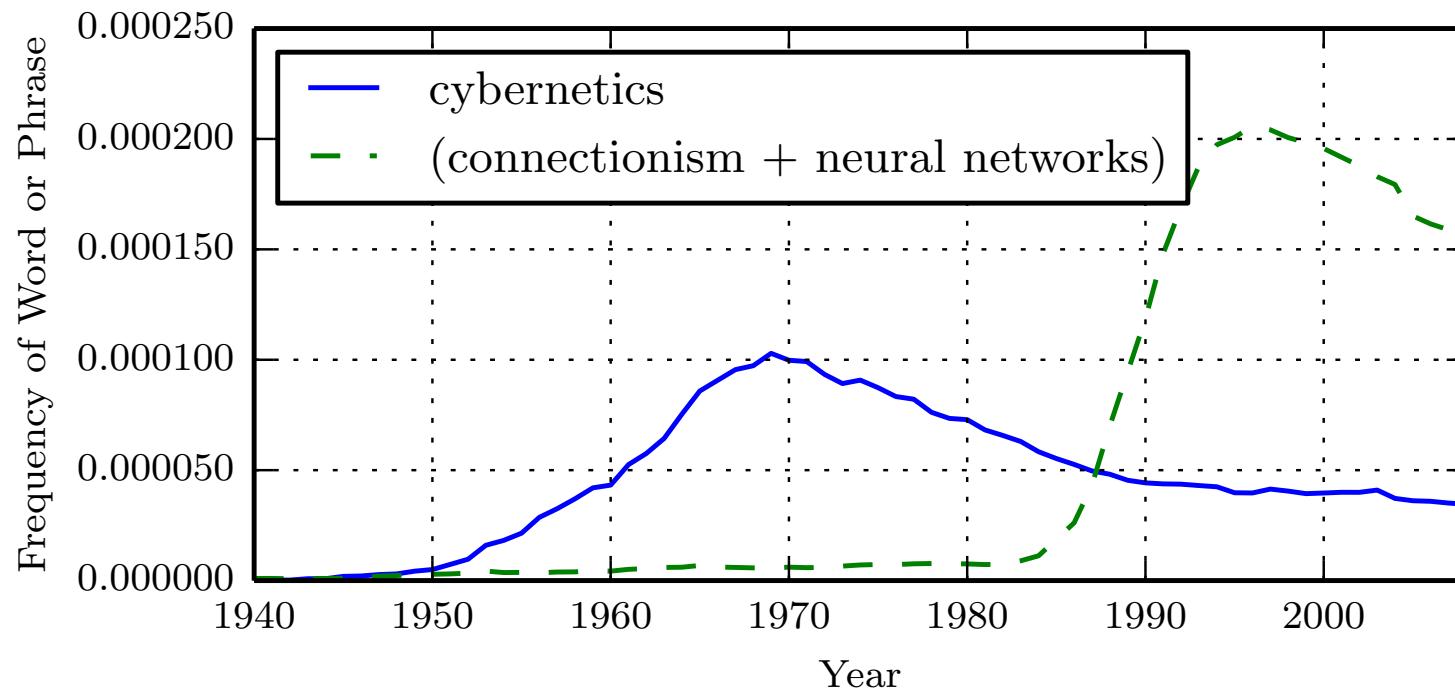


Figure 1.7

(Goodfellow 2016)

# Historical Trends: Growing Datasets

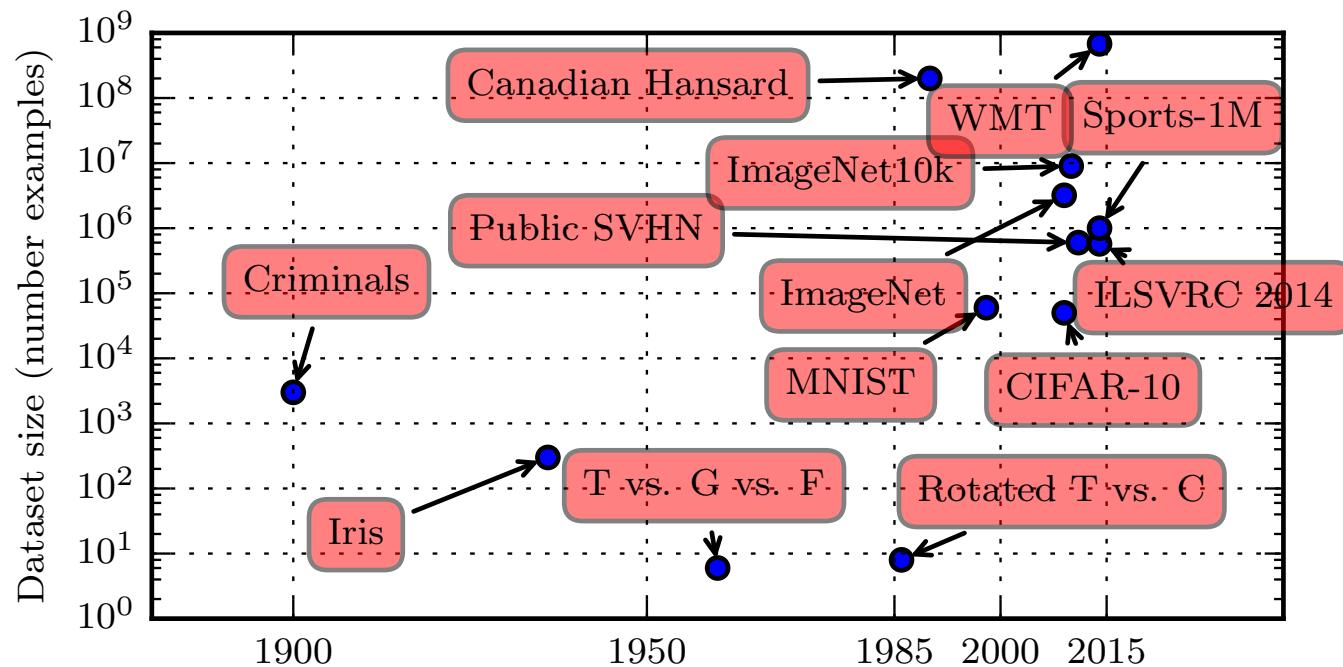


Figure 1.8

(Goodfellow 2016)

# The MNIST Dataset

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 8 | 9 | 0 | 1 | 2 | 3 | 4 | 7 | 8 | 9 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 6 |
| 4 | 2 | 6 | 4 | 7 | 5 | 5 | 4 | 7 | 8 | 9 | 2 | 9 | 3 | 9 | 3 | 8 | 2 | 0 | 5 |
| 0 | 1 | 0 | 4 | 2 | 6 | 5 | 3 | 5 | 3 | 8 | 0 | 0 | 3 | 4 | 1 | 5 | 3 | 0 | 8 |
| 3 | 0 | 6 | 2 | 7 | 1 | 1 | 8 | 1 | 7 | 1 | 3 | 8 | 9 | 7 | 6 | 7 | 4 | 1 | 6 |
| 7 | 5 | 1 | 7 | 1 | 9 | 8 | 0 | 6 | 9 | 4 | 9 | 9 | 3 | 7 | 1 | 9 | 2 | 2 | 5 |
| 3 | 7 | 8 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 8 | 1 | 0 | 5 | 5 | 1 | 9 | 0 | 4 | 1 | 9 |
| 3 | 8 | 4 | 7 | 7 | 8 | 5 | 0 | 6 | 5 | 5 | 3 | 3 | 3 | 9 | 8 | 1 | 4 | 0 | 6 |
| 1 | 0 | 0 | 6 | 2 | 1 | 1 | 3 | 2 | 8 | 8 | 7 | 8 | 4 | 6 | 0 | 2 | 0 | 3 | 6 |
| 8 | 7 | 1 | 5 | 9 | 9 | 3 | 2 | 4 | 9 | 4 | 6 | 5 | 3 | 2 | 8 | 5 | 9 | 4 | 1 |
| 6 | 5 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 8 | 9 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 6 | 4 | 2 | 6 | 4 | 7 | 5 | 5 |
| 4 | 7 | 8 | 9 | 2 | 9 | 3 | 9 | 3 | 8 | 2 | 0 | 9 | 8 | 0 | 5 | 6 | 0 | 1 | 0 |
| 4 | 2 | 6 | 5 | 5 | 5 | 4 | 3 | 4 | 1 | 5 | 3 | 0 | 8 | 3 | 0 | 6 | 2 | 7 | 1 |
| 1 | 8 | 1 | 7 | 1 | 3 | 8 | 5 | 4 | 2 | 0 | 9 | 7 | 6 | 7 | 4 | 1 | 6 | 8 | 4 |
| 7 | 5 | 1 | 2 | 6 | 7 | 1 | 9 | 8 | 0 | 6 | 9 | 4 | 9 | 9 | 6 | 2 | 3 | 7 | 1 |
| 9 | 2 | 2 | 5 | 3 | 7 | 8 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 | 8 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 2 | 1 | 2 | 1 | 3 |
| 9 | 9 | 8 | 5 | 3 | 7 | 0 | 7 | 7 | 5 | 7 | 9 | 9 | 4 | 7 | 0 | 3 | 4 | 1 | 4 |
| 4 | 7 | 5 | 8 | 1 | 4 | 8 | 4 | 1 | 8 | 6 | 6 | 4 | 6 | 3 | 5 | 7 | 2 | 5 | 9 |

Figure 1.9

(Goodfellow 2016)

# Connections per Neuron

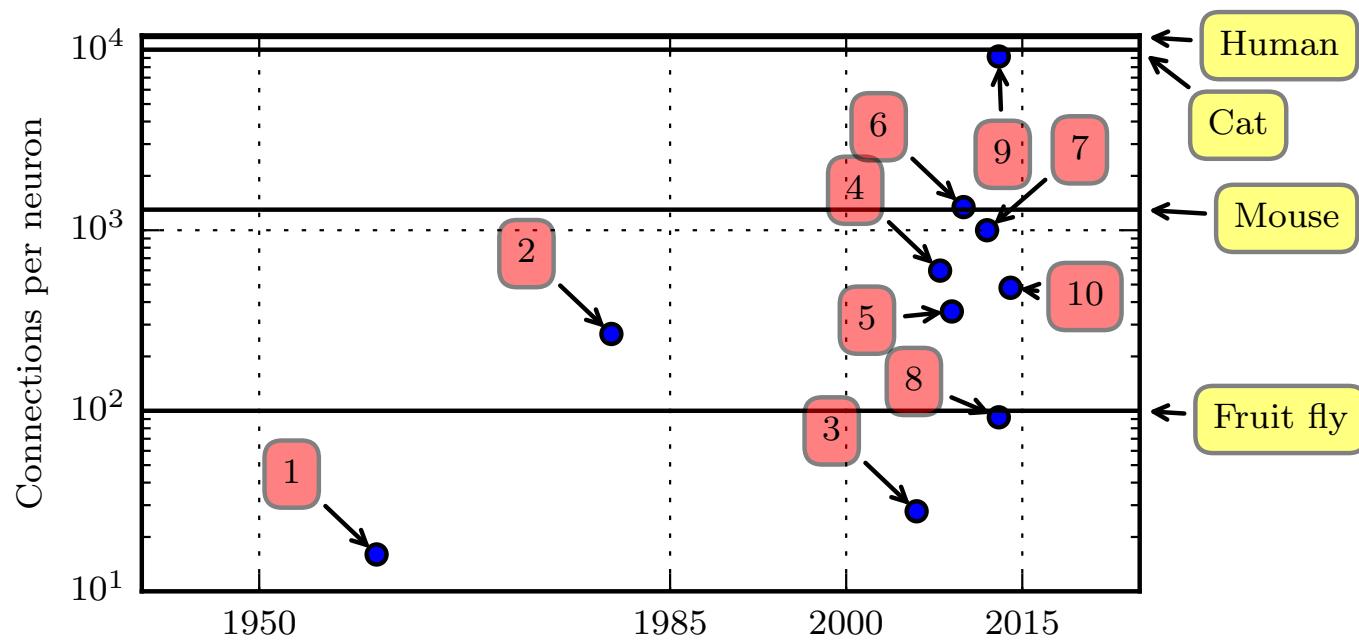


Figure 1.10

(Goodfellow 2016)

# Number of Neurons

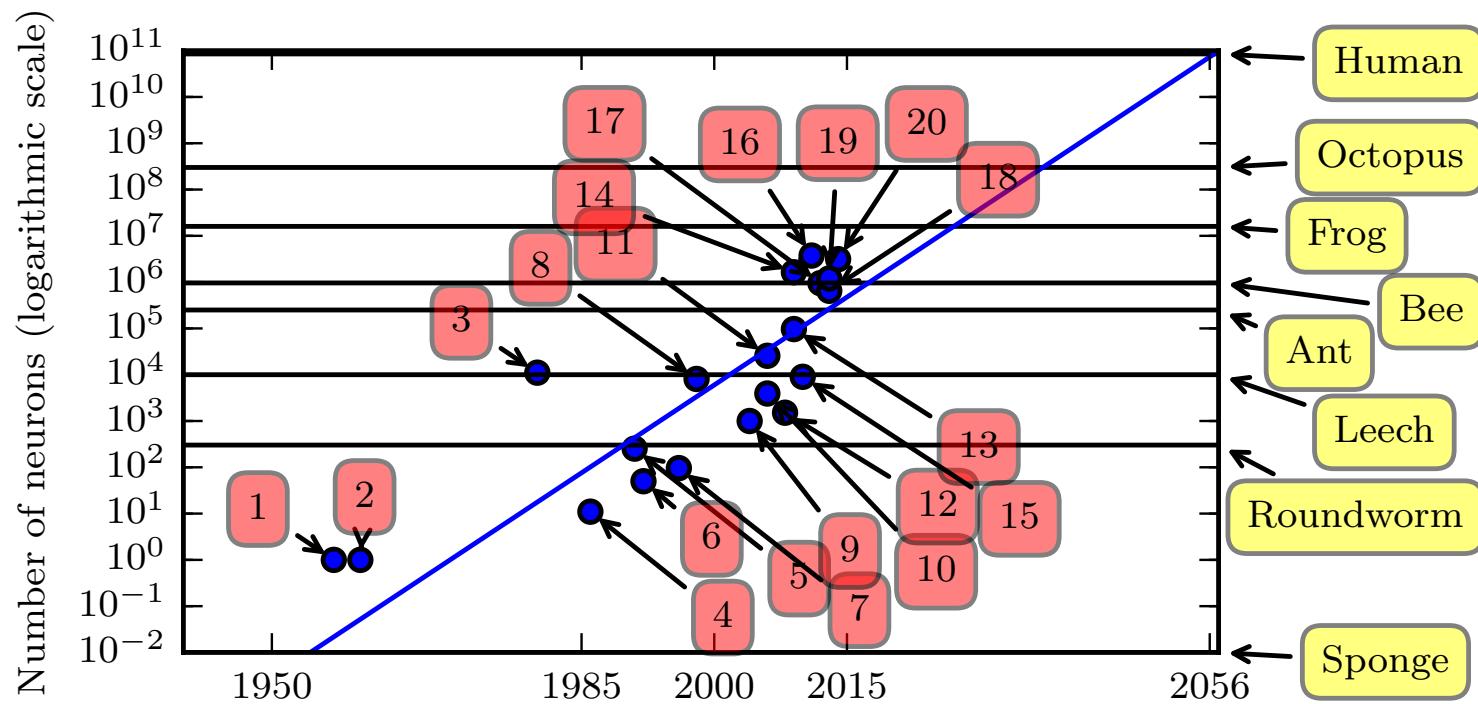


Figure 1.11

(Goodfellow 2016)

# Solving Object Recognition

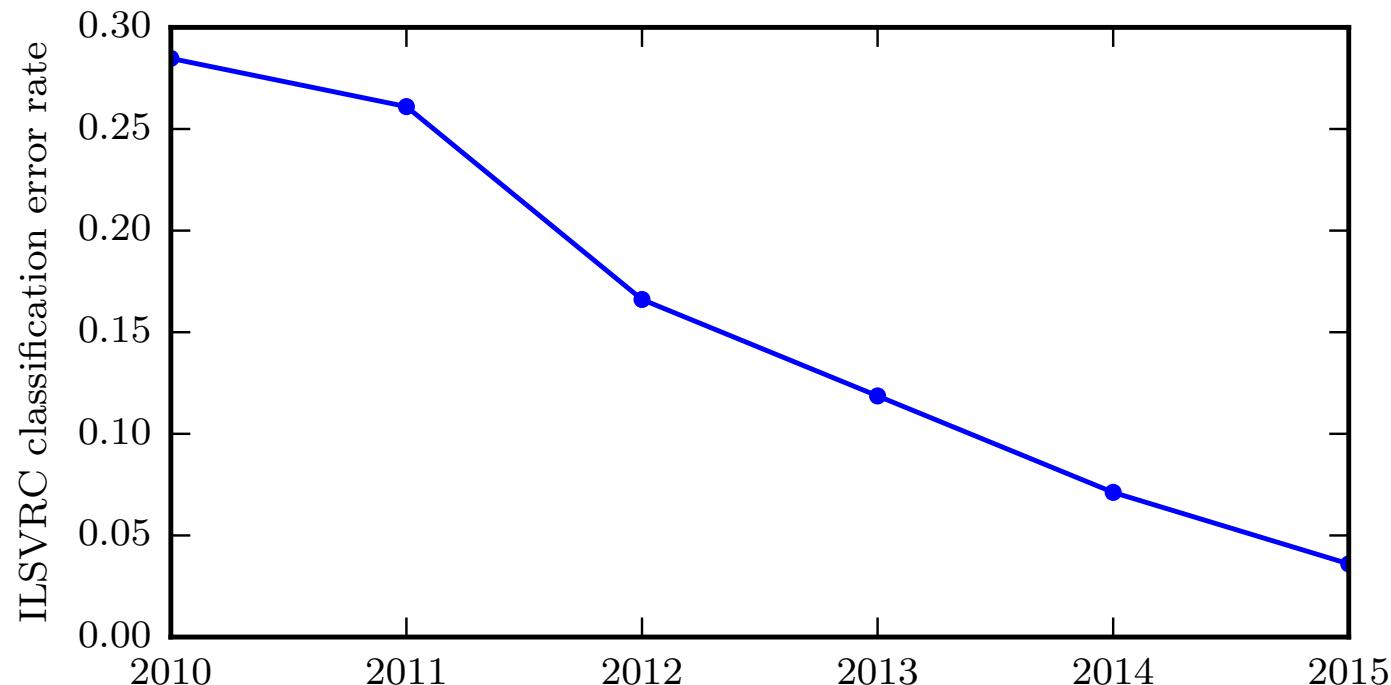


Figure 1.12

(Goodfellow 2016)

# Thanks

This is the end of this lecture!



# Neurons

