

# Maths of AI

## Introduction

Instructor - Simon Lucey



THE UNIVERSITY  
*of* ADELAIDE

AUSTRALIAN  
INSTITUTE FOR  
MACHINE LEARNING

# Today

---

- **About the Course**
- What is AI?
- The Rise of Deep Learning
- Why Math in AI Matters

# Instructors

---



Prof. Simon Lucey



Prof. Matthew Roughan



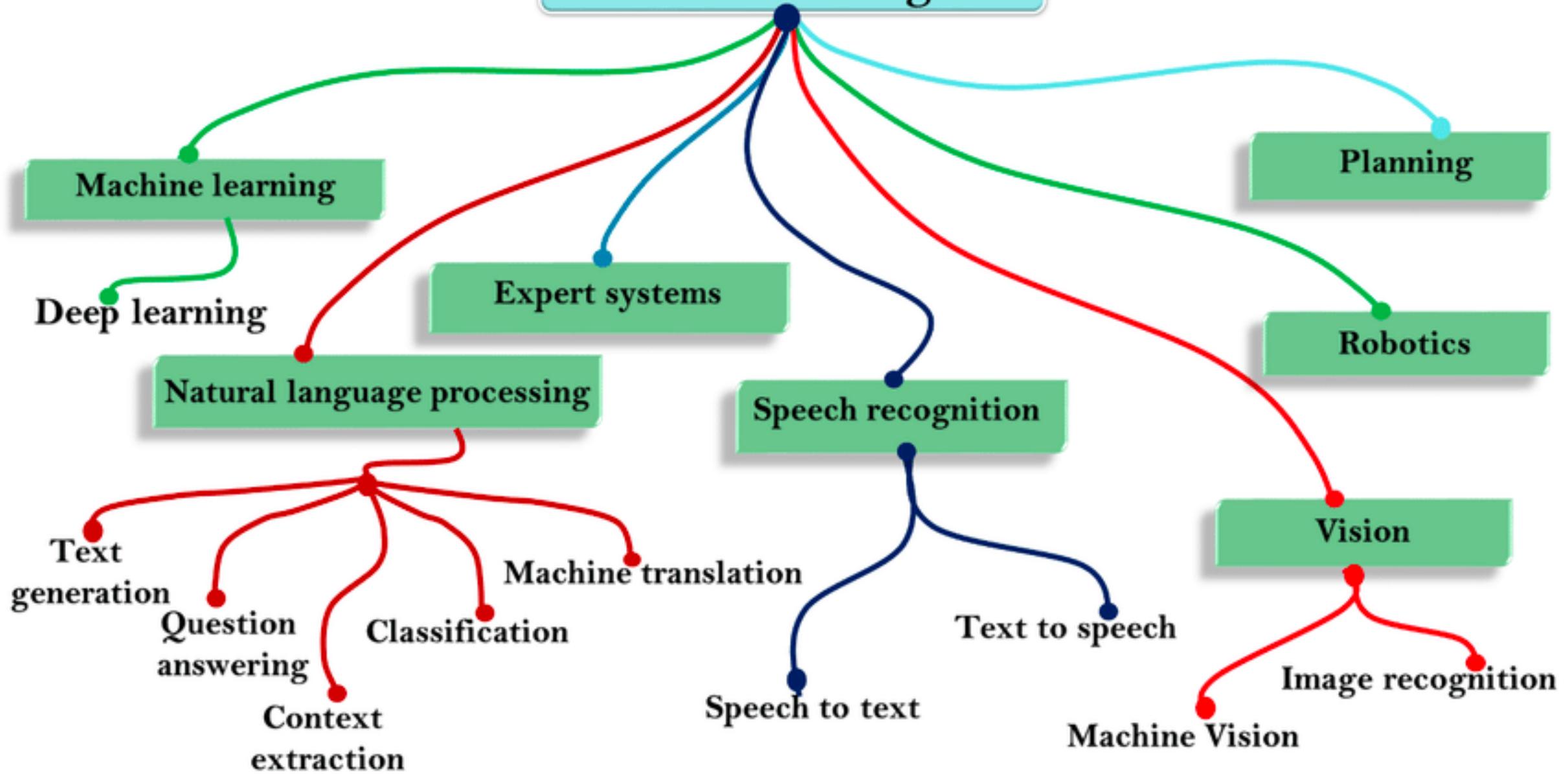
Prof. Lewis Mitchell

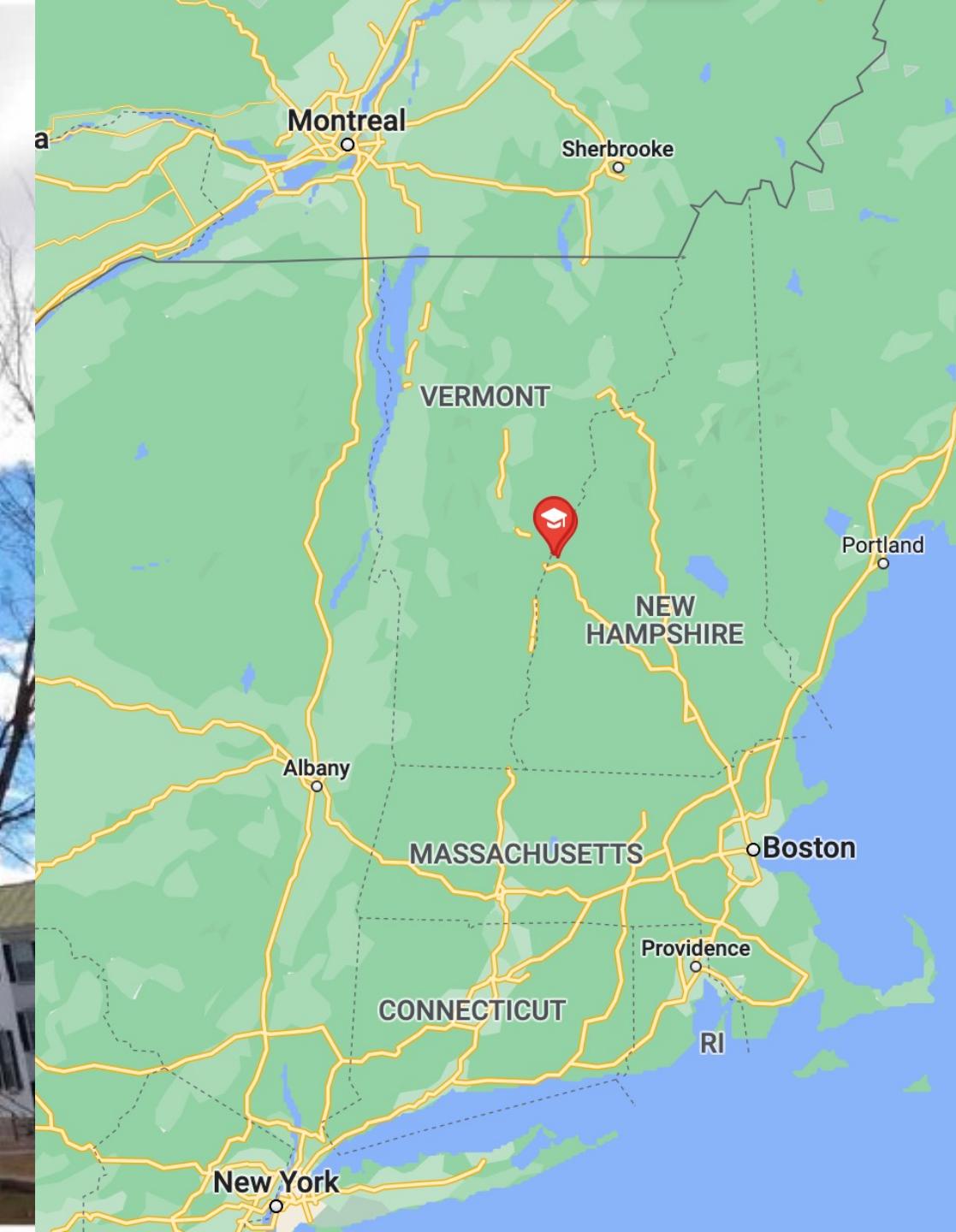
# Today

---

- About the Course
- **What is AI?**
- The Rise of Deep Learning
- Why Math in AI Matters

# Artificial Intelligence

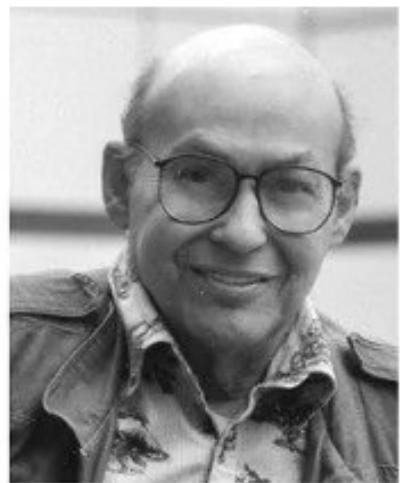




# 1956 Dartmouth Conference: The Founding Fathers of AI



John McCarthy



Marvin Minsky



Claude Shannon



Ray Solomonoff



Alan Newell



Herbert Simon



Arthur Samuel



Oliver Selfridge



Nathaniel Rochester

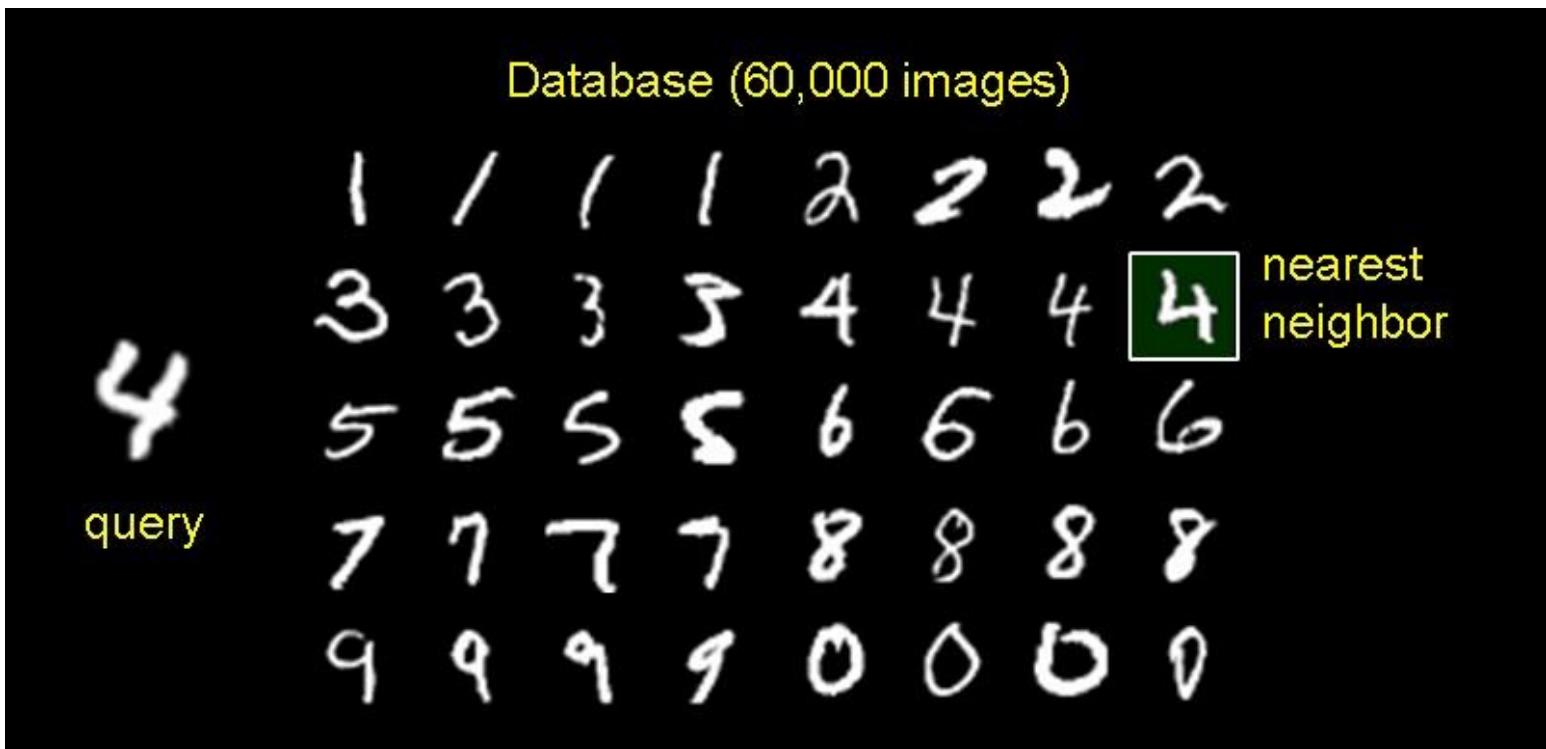


Trenchard More

# Conventional programming

```
2
3 class Room(object):
4     def __init__(self, inventory, desc, short_desc):
5         self.inventory = inventory
6         self._n = None
7         self._s = None
8         self._e = None
9         self._w = None
10        self._desc = desc
11        self._short_desc = short_desc
12        self._gate_n = None
13        self._gate_s = None
14        self._gate_e = None
15        self._gate_w = None
16
17        if not isinstance(desc, str):
18            raise TypeError ("the input provided is not a string.")
19        elif not isinstance(short_desc, str):
20            raise ValueError ("the string provided is empty.")
21
22    # these set the gates
23    # they set the opposite gates, with checks to avoid recursion loops
24    def set_n(self, other):
25        if not isinstance(other, Room) or not other:
26            raise TypeError("Room is not None or an instance of Room")
```

# Machine Learning



# Perceptron

- Rosenblatt simulated the perceptron on a IBM 704 computer at Cornell in 1957.
- Input scene (i.e. printed character) was illuminated by powerful lights and captured on a 20x20 cadmium sulphide photo cells.
- Weights of perceptron were applied using variable rotary resistors.
- Often times referred to as the very first neural network.



“Frank Rosenblatt”

# The New York Times

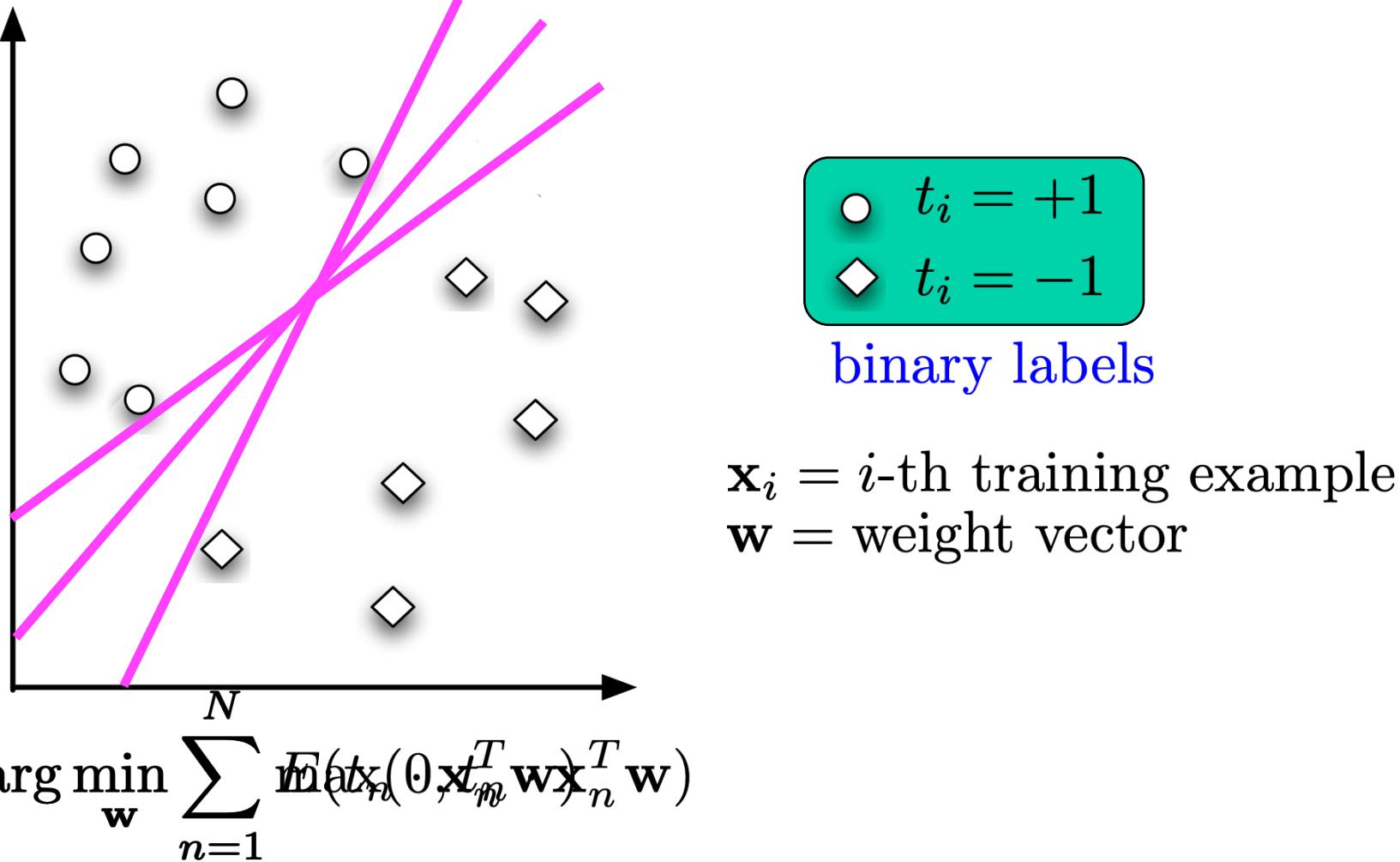
July 8, 1958

**NEW NAVY DEVICE LEARNS BY DOING**

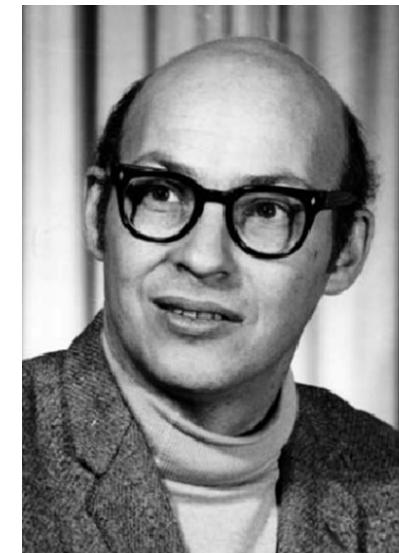
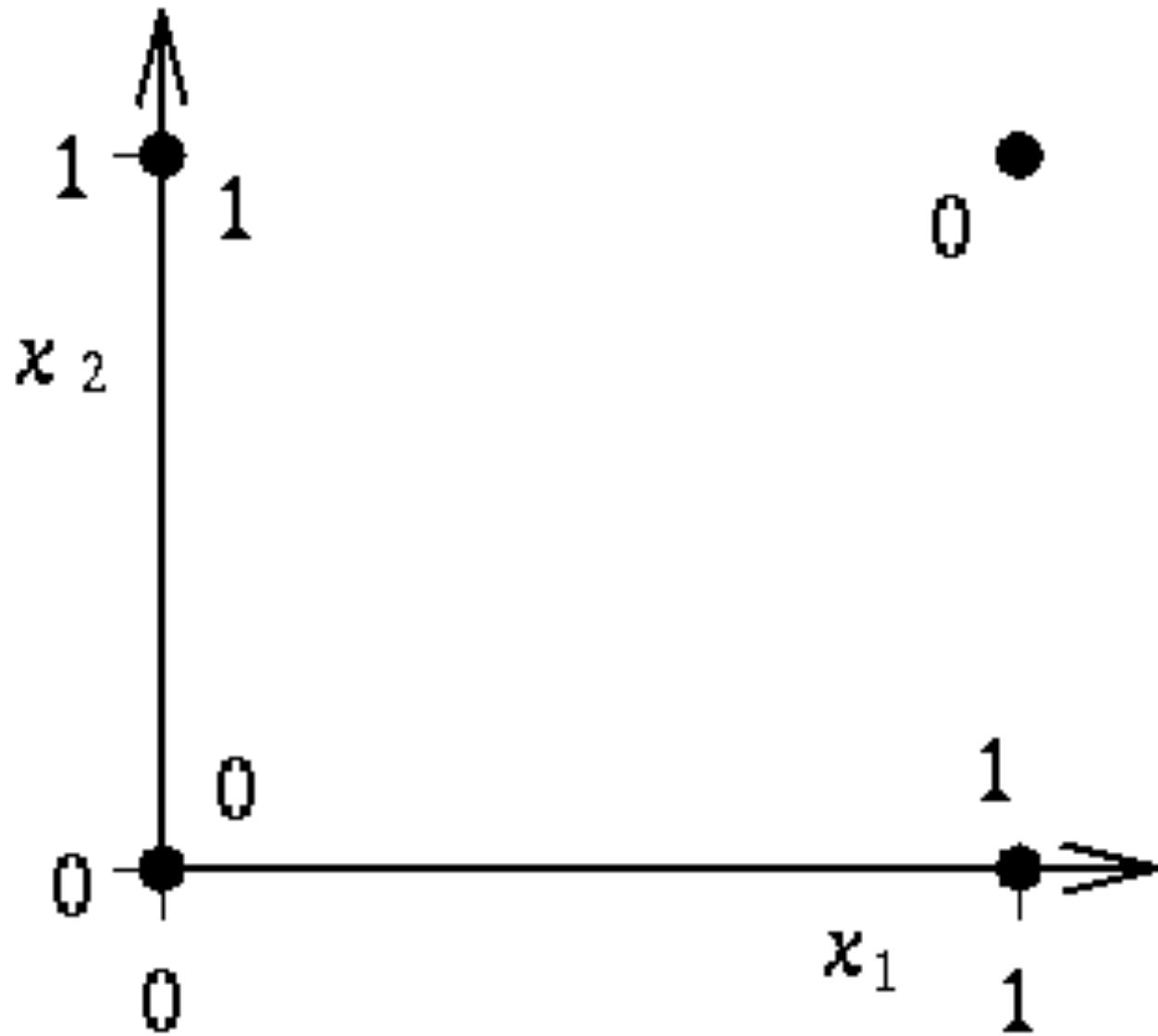
"The Navy revealed the embryo of an electronic computer today that it expects will be able to talk, see, write, reproduce itself and be conscious of its existence... Dr. Frank Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers"



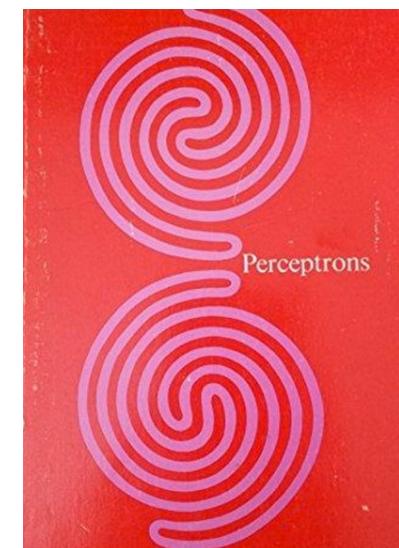
# Perceptron = Linear Discriminant



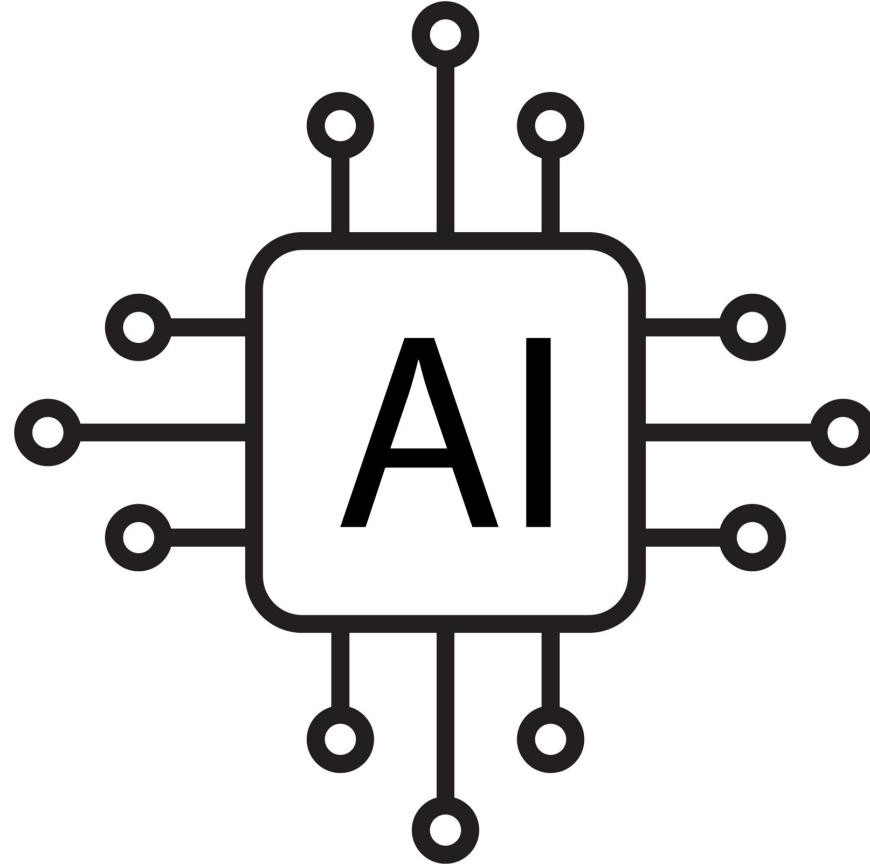
# XORs, Minsky and the AI Winter



“Marvin Minsky”

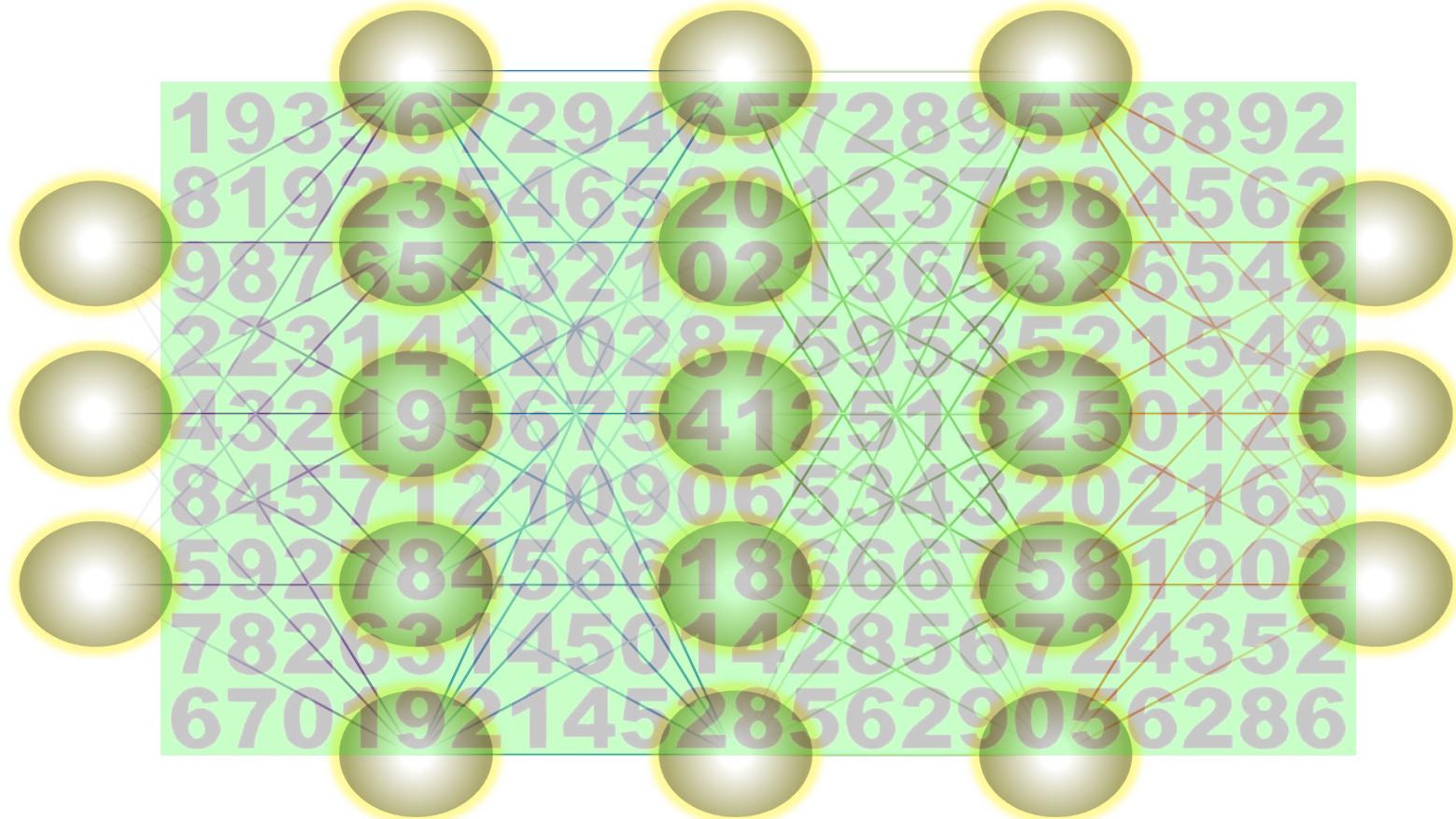


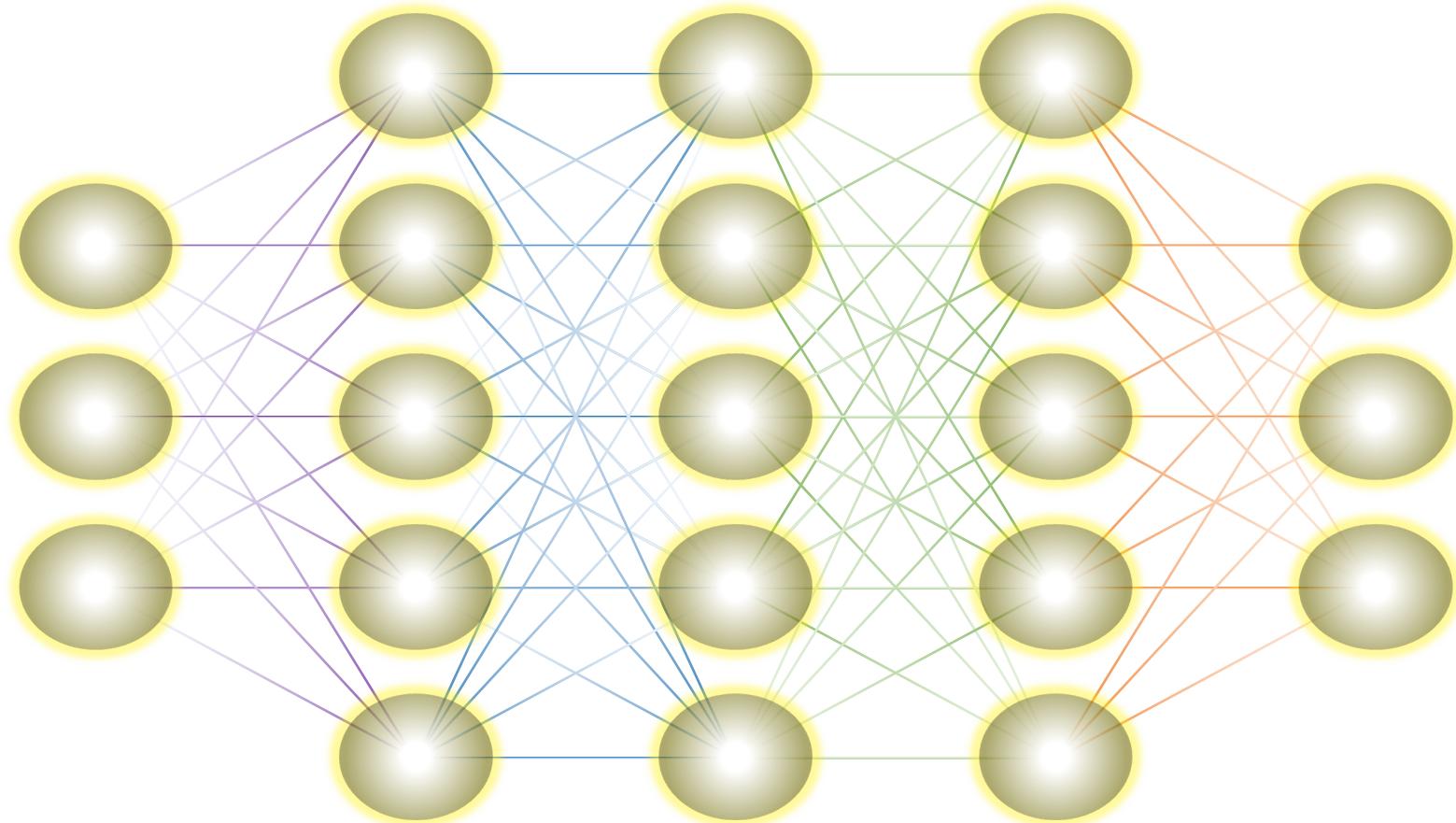
# Artificial Intelligence



**x**

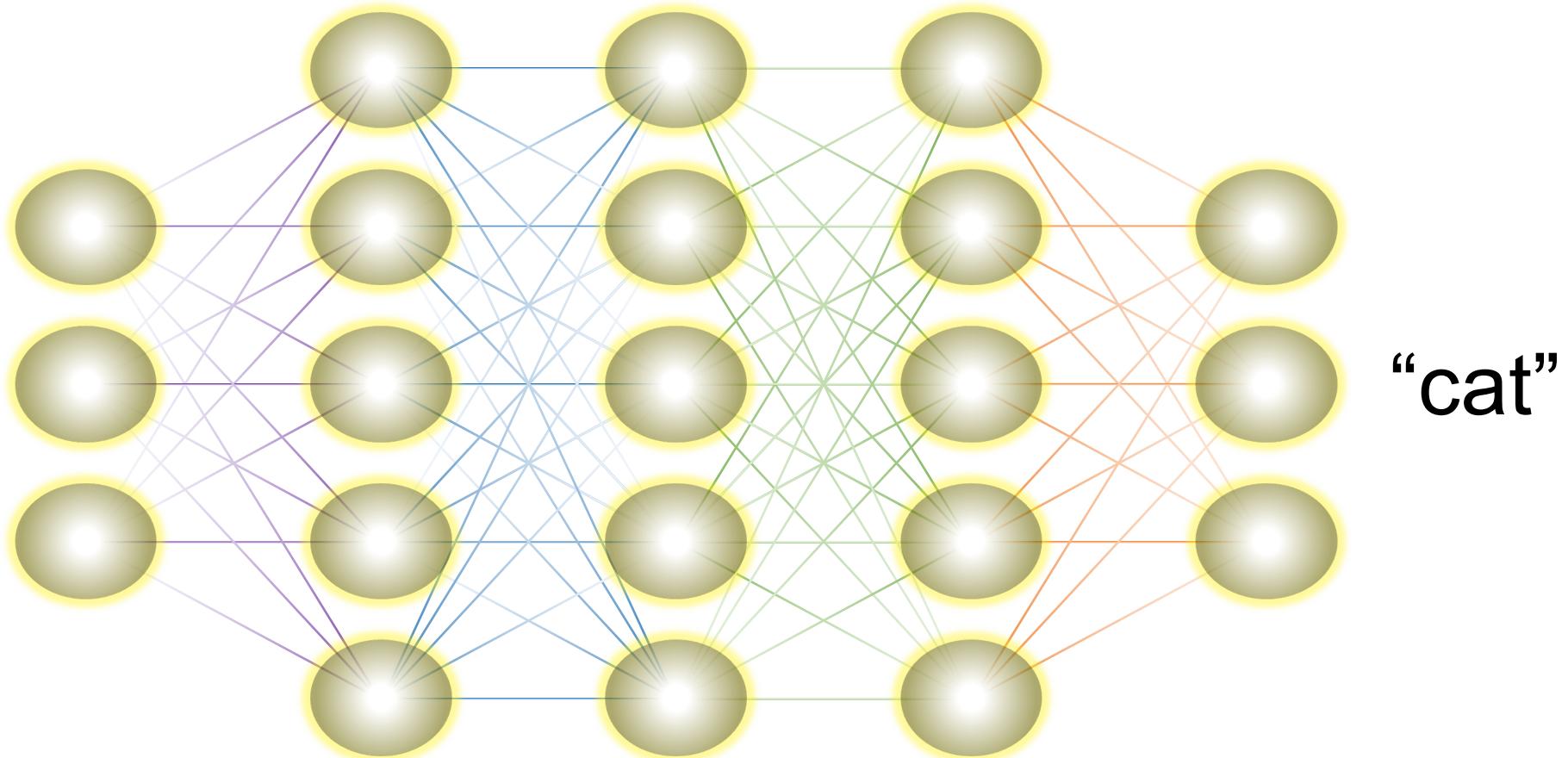
**y**



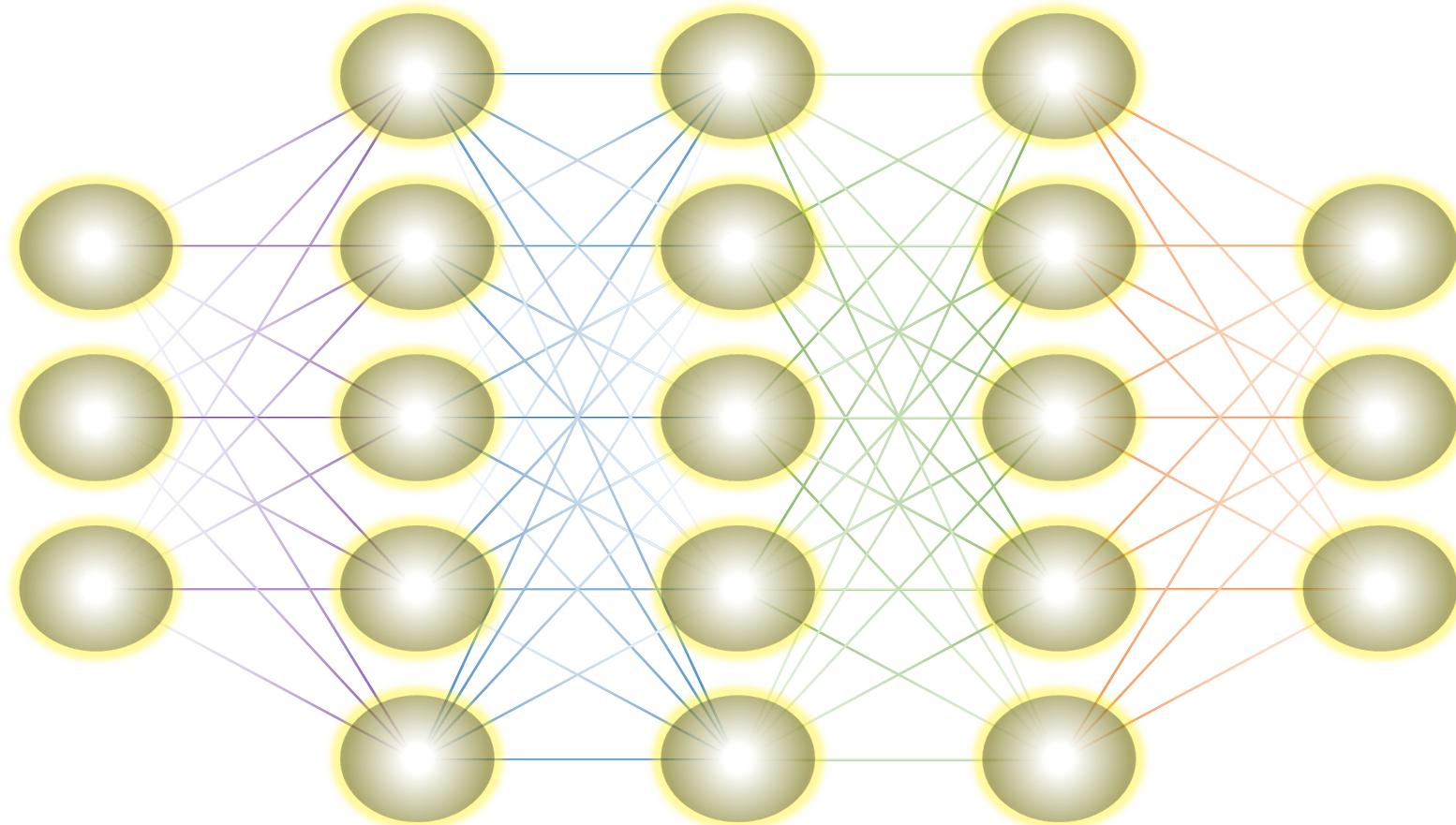


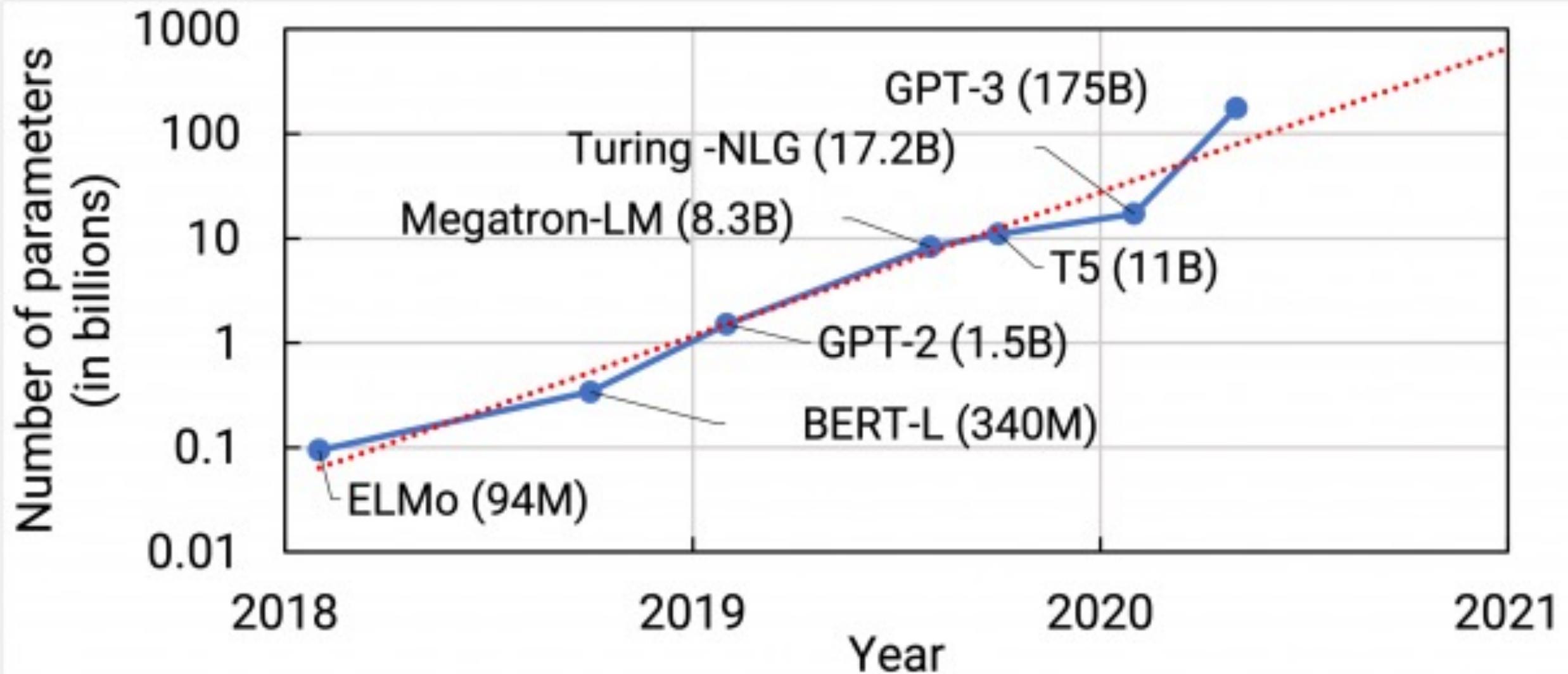
“cat”

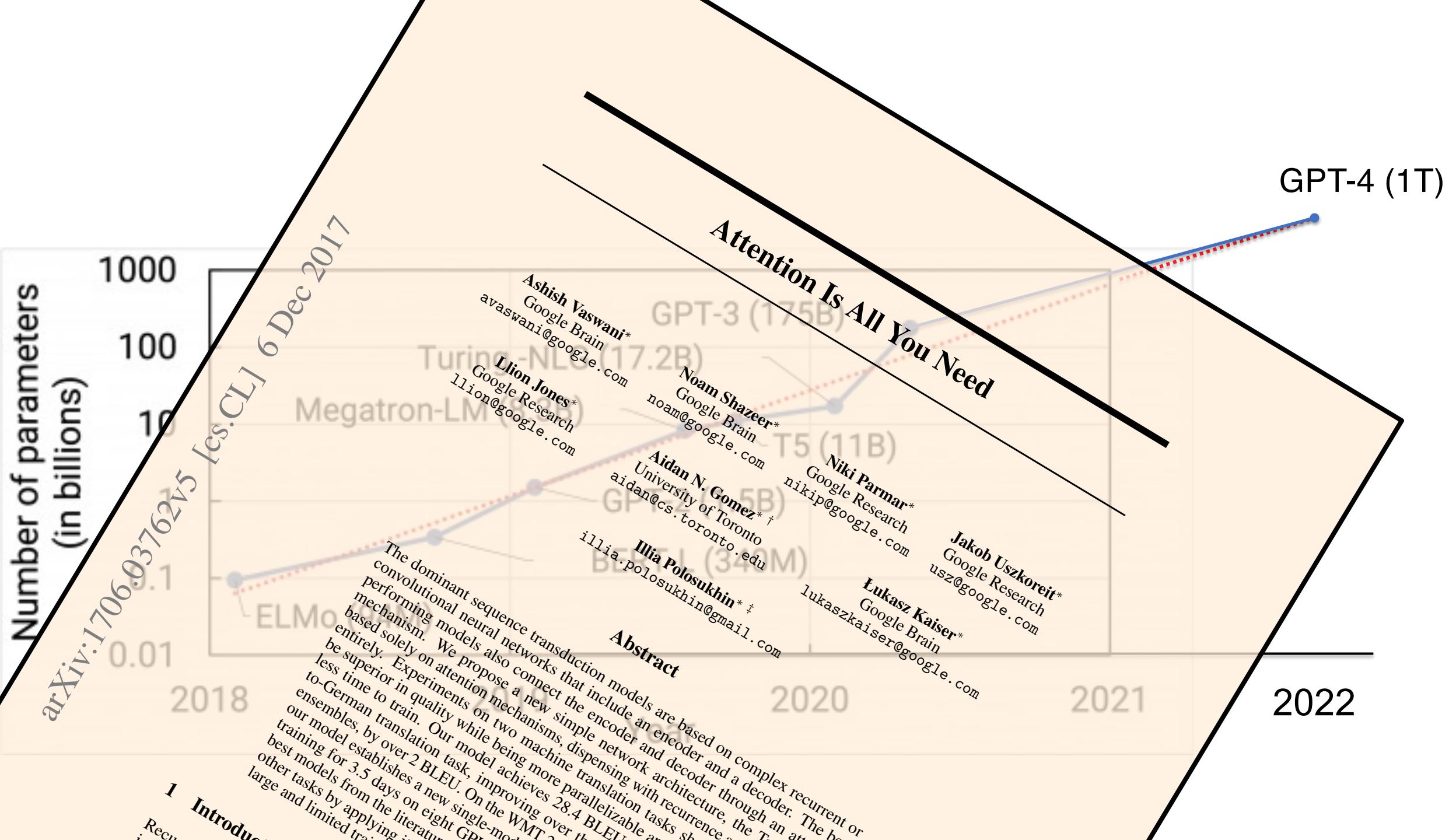
# Artificial Neural Networks



“cat”





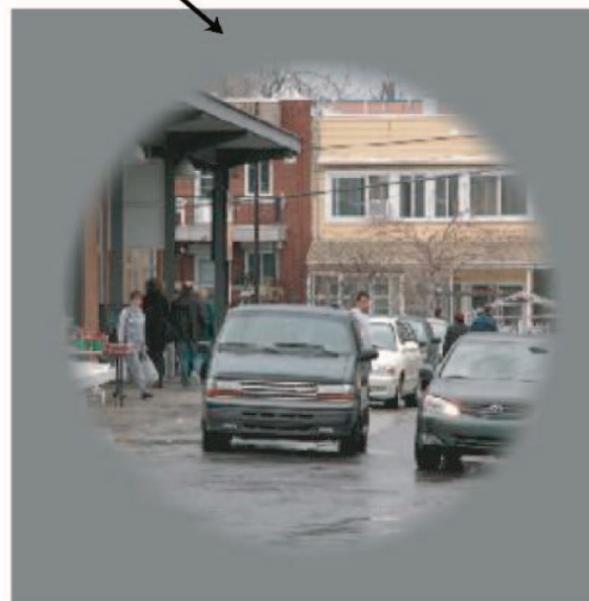
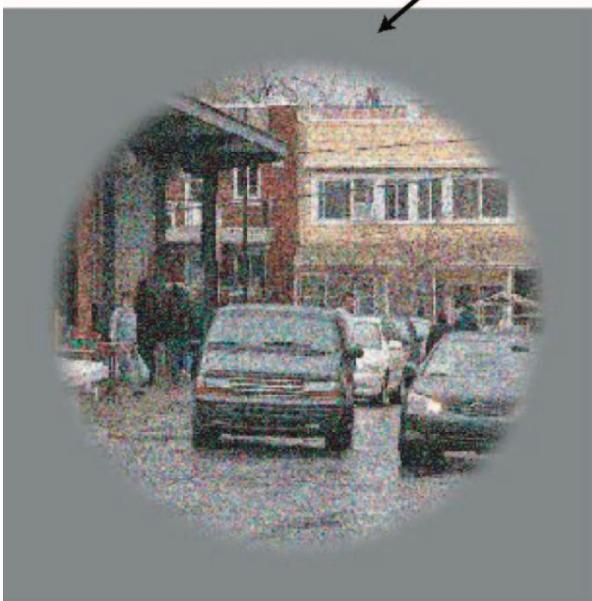
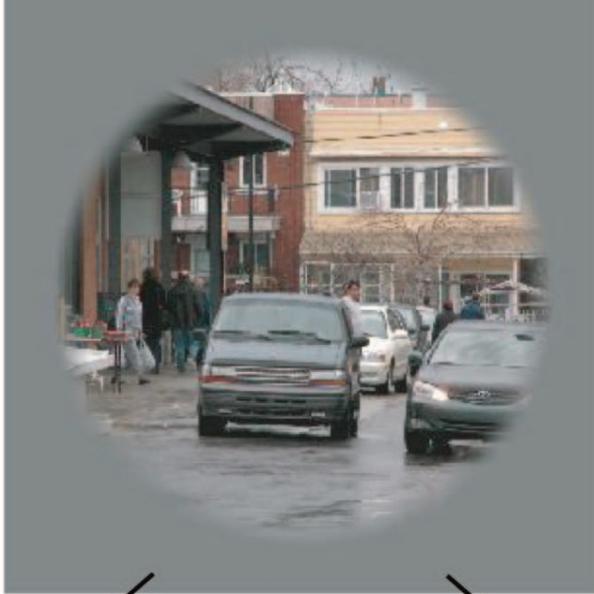


# Today

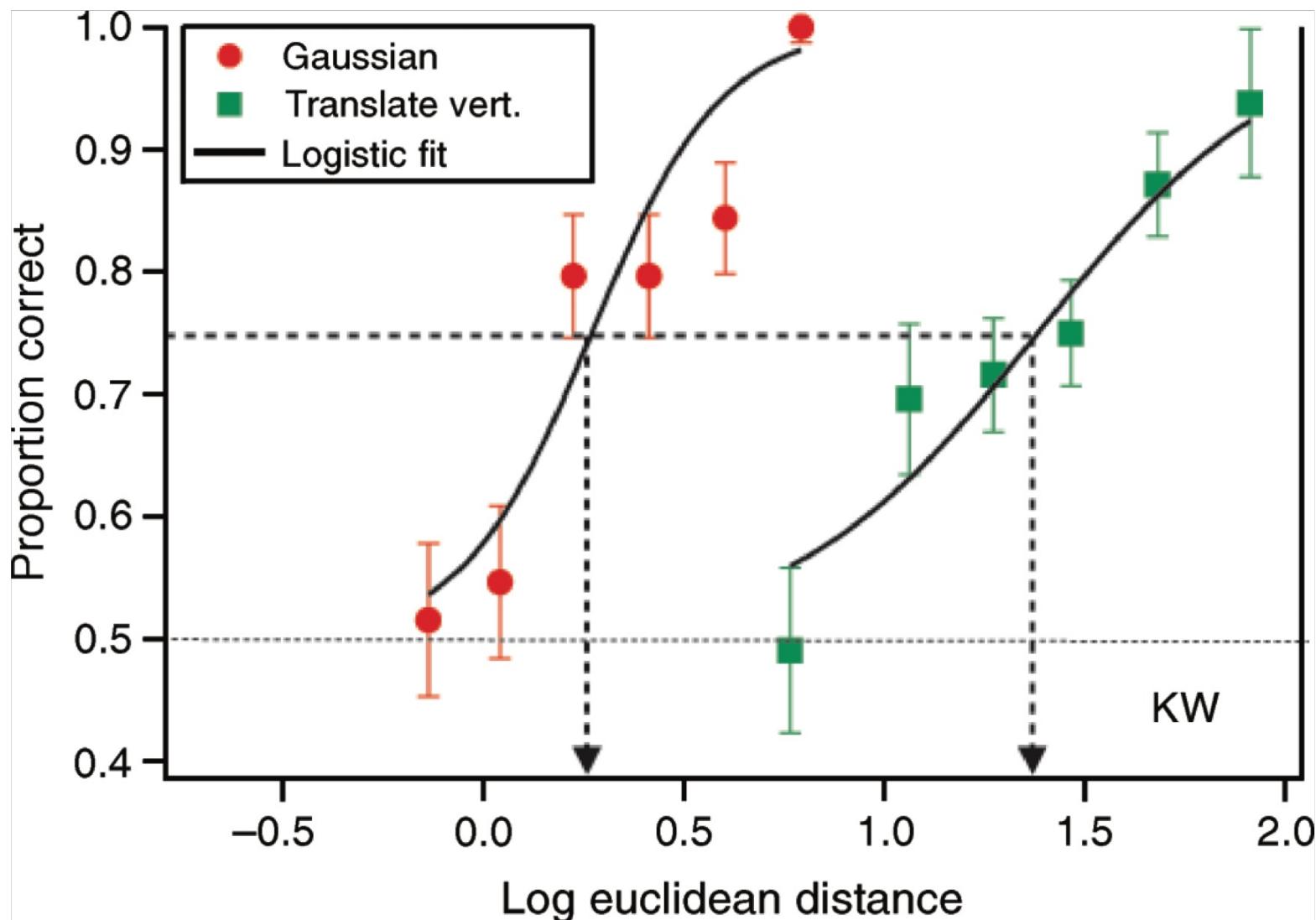
---

- About the Course
- What is AI?
- **The Rise of Deep Learning**
- Why Math in AI Matters

# Question?

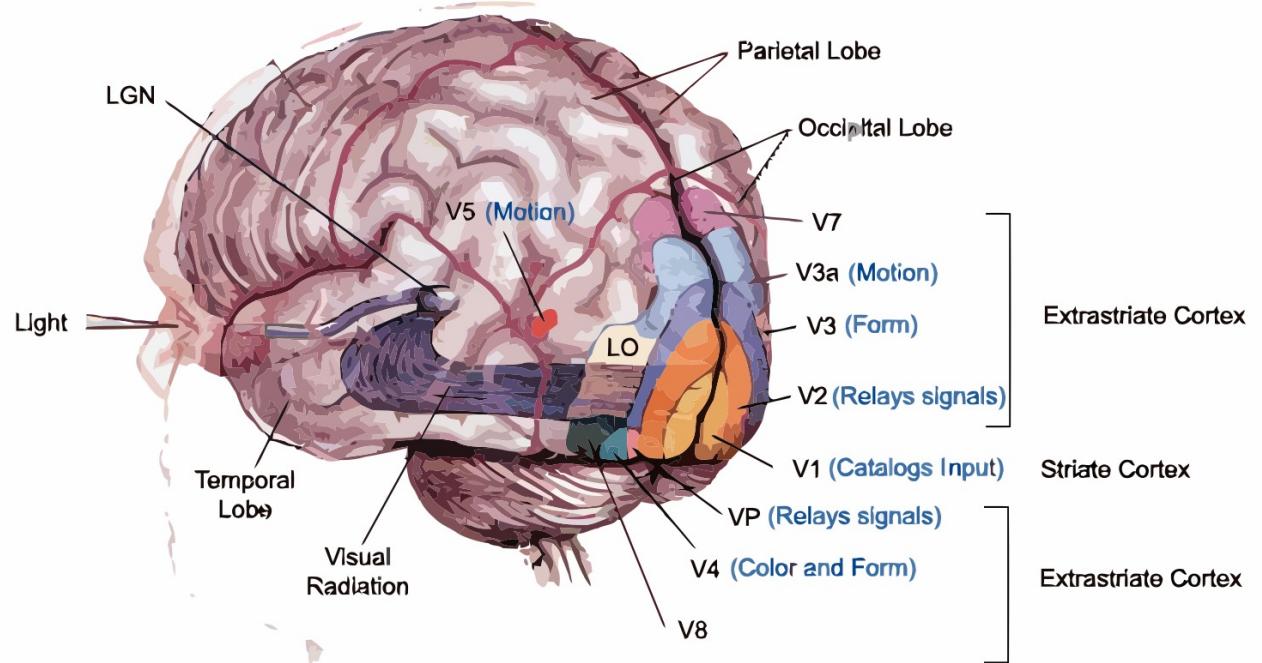


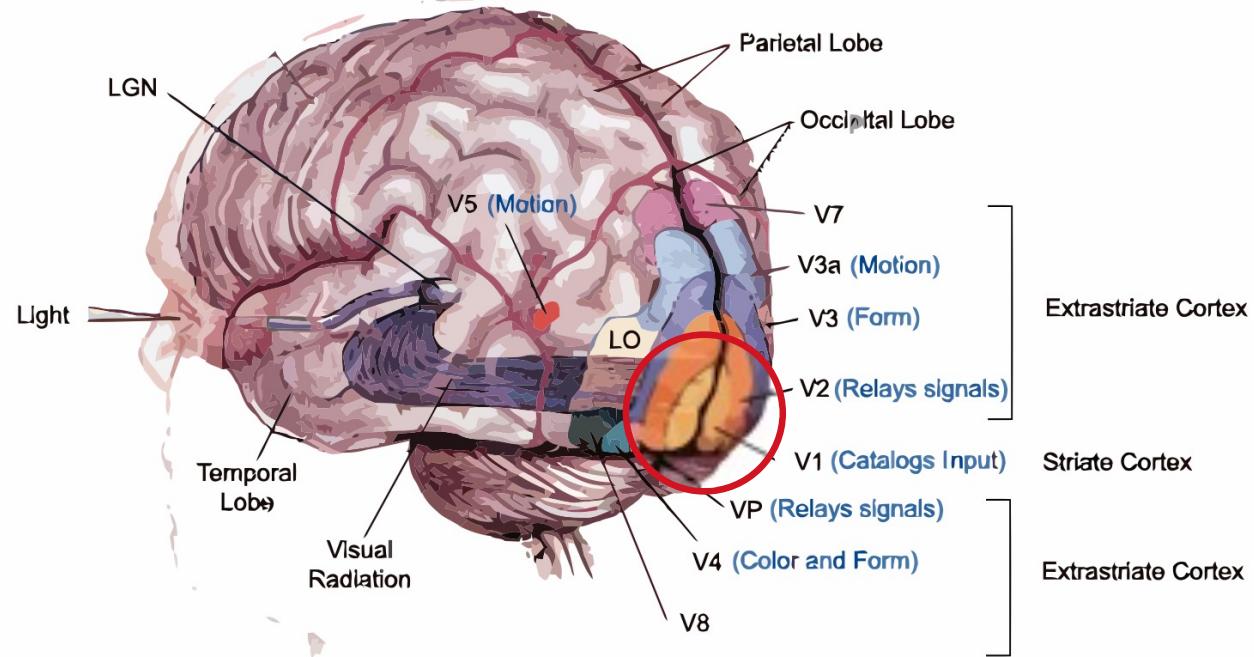
# Answer

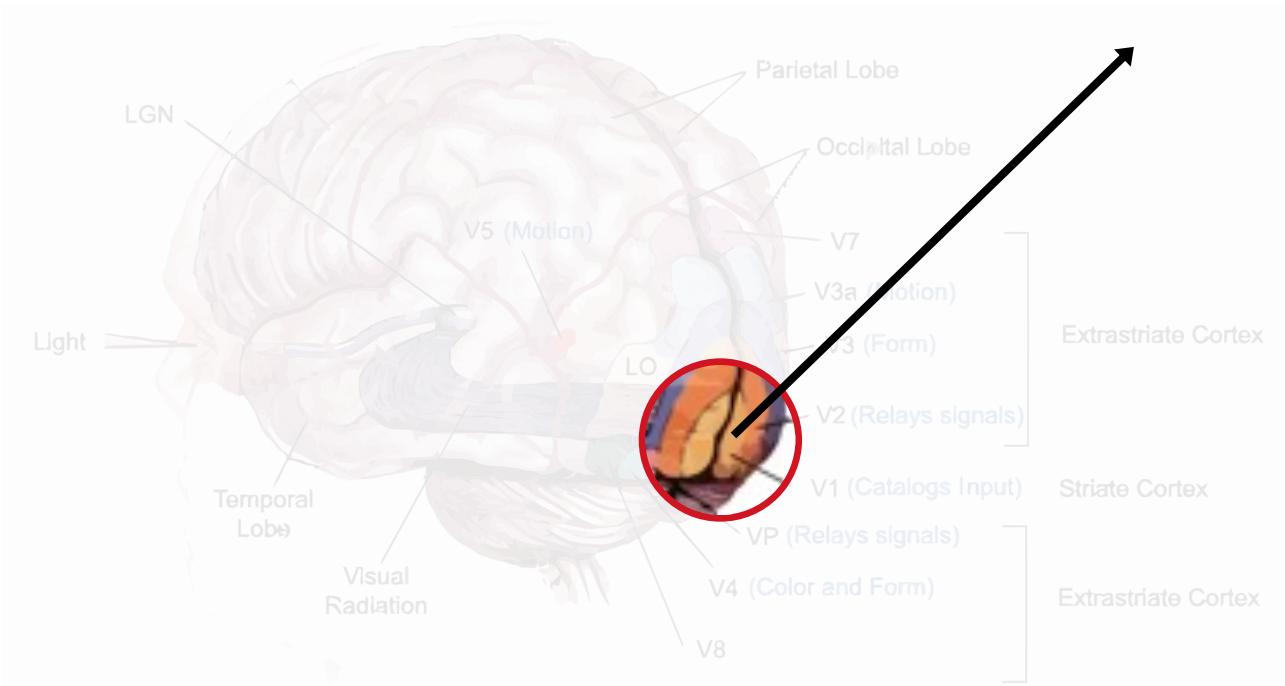


# Why?

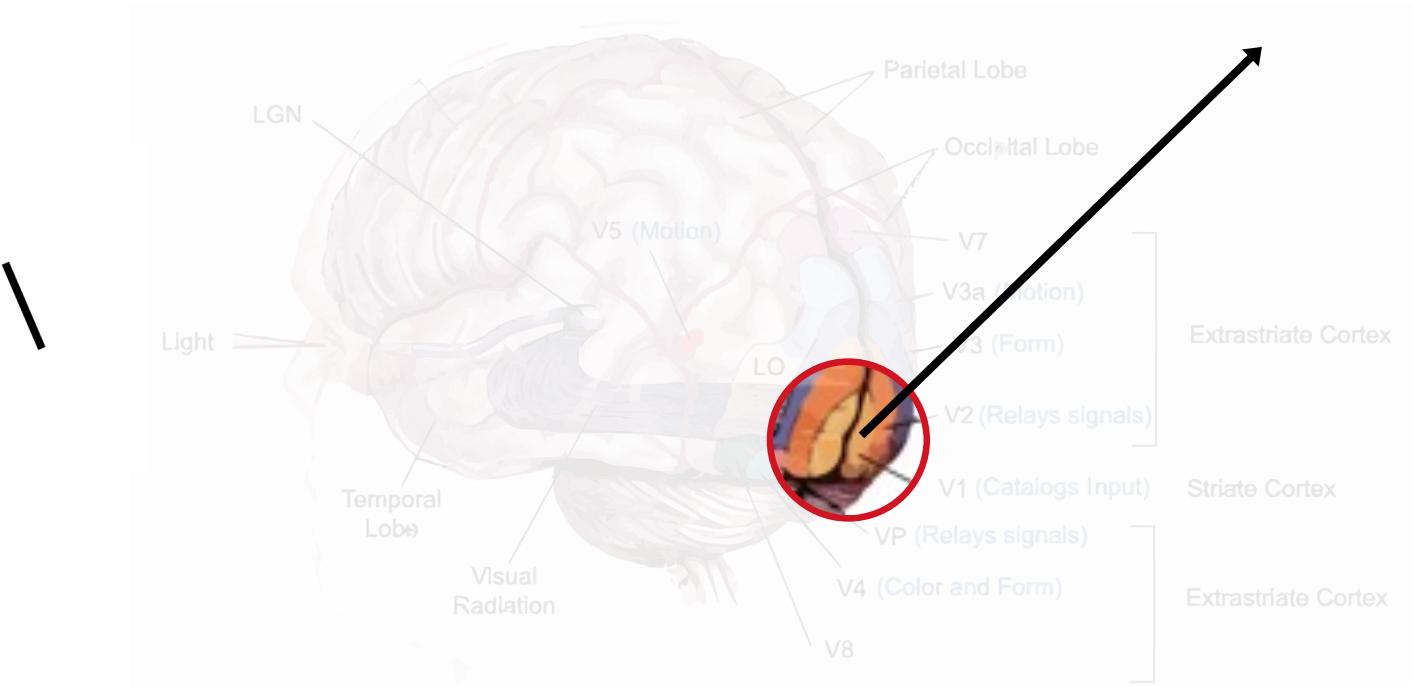
207	245	77	21	247	211	240	1
219	41	58	179	161	154	184	98
215	145	187	71	251	249	65	100
192	2	189	247	166	63	232	213
105	94	66	190	156	61	89	145
159	154	87	184	101	105	72	71
192	111	6	94	60	70	65	226
175	120	210	226	80	183	168	184
134	56	36	240	159	178	76	135
239	244	199	9	132	104	188	185
245	210	78	199	0	92	9	246
5	121	187	122	107	47	12	119
230	171	135	36	82	54	65	37
61	140	79	19	161	96	127	187
56	223	46	6	180	186	142	244
28	20	61	2	178	187	98	220



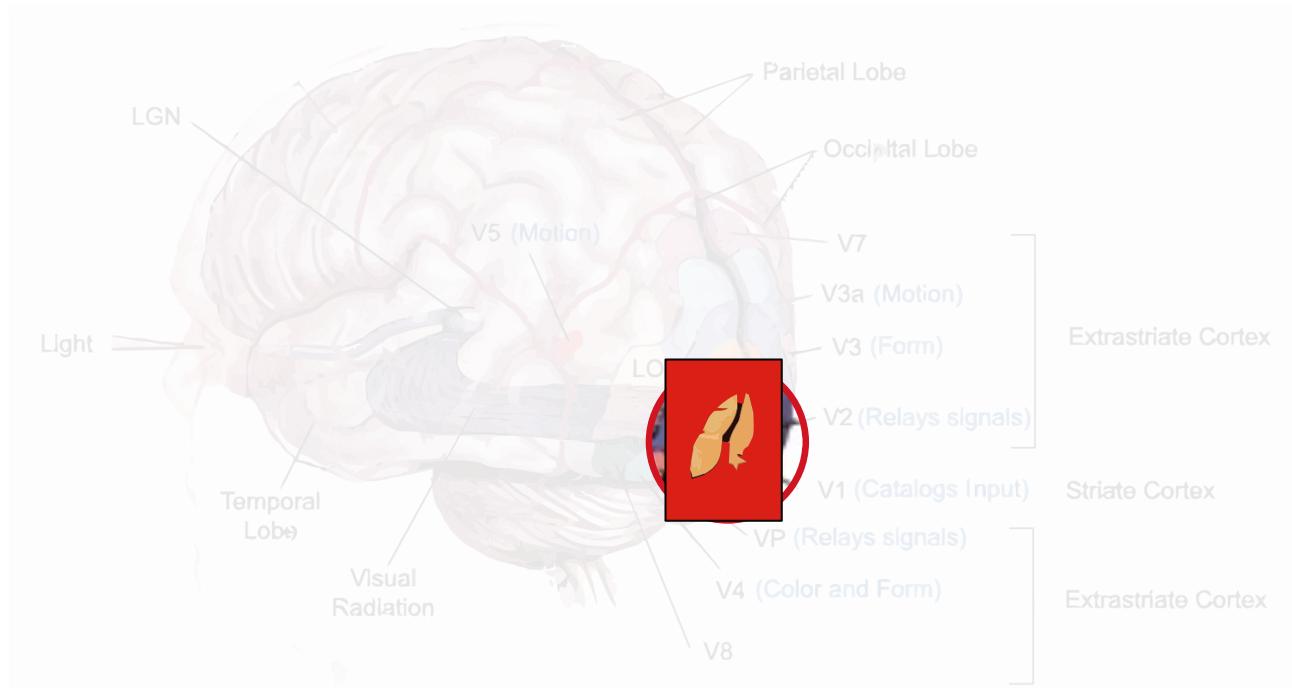




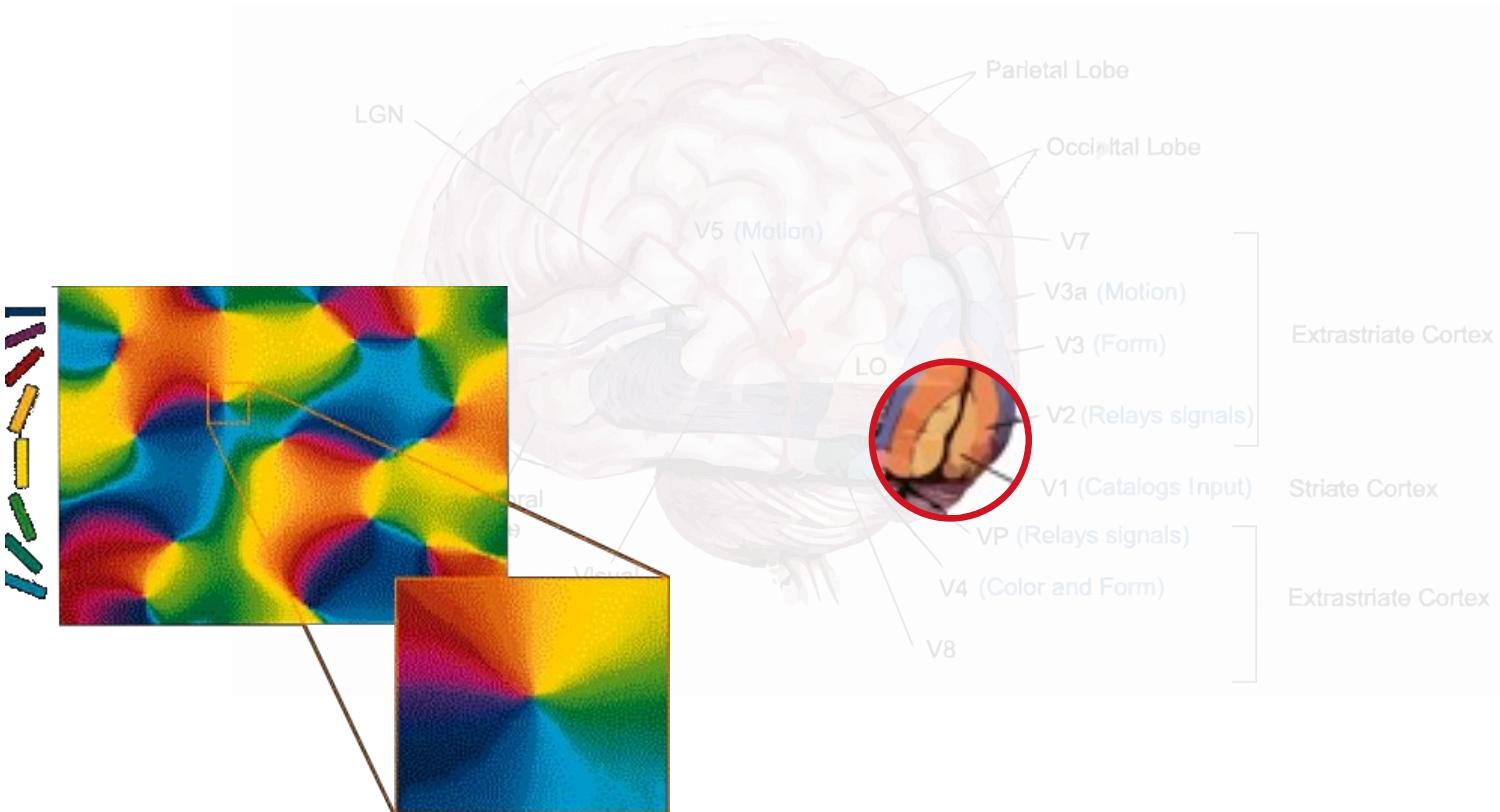
D.H. Hubel & T.N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of Physiology*, 160(1):106, 1962.



D.H. Hubel & T.N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of Physiology*, 160(1):106, 1962.



D.H. Hubel & T.N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of Physiology*, 160(1):106, 1962.



D.H. Hubel & T.N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of Physiology*, 160(1):106, 1962.

# Emergence of simple-cell receptive field properties by learning a sparse code for natural images

Bruno A. Olshausen\* & David J. Field

Department of Psychology, Uris Hall, Cornell University, Ithaca,  
New York 14853, USA

THE receptive fields of simple cells in mammalian primary visual cortex can be characterized as being spatially localized, oriented<sup>1–4</sup> and bandpass (selective to structure at different spatial scales), comparable to the basis functions of wavelet transforms<sup>5,6</sup>. One approach to understanding such response properties of visual neurons has been to consider their relationship to the statistical structure of natural images in terms of efficient coding<sup>7–12</sup>. Along these lines, a number of studies have attempted to train unsupervised learning algorithms on natural images in the hope of developing receptive fields with similar properties<sup>13–18</sup>, but none has succeeded in producing a full set that spans the image space and contains all three of the above properties. Here we investigate the proposal<sup>8,12</sup> that a coding strategy that maximizes sparseness is sufficient to account for

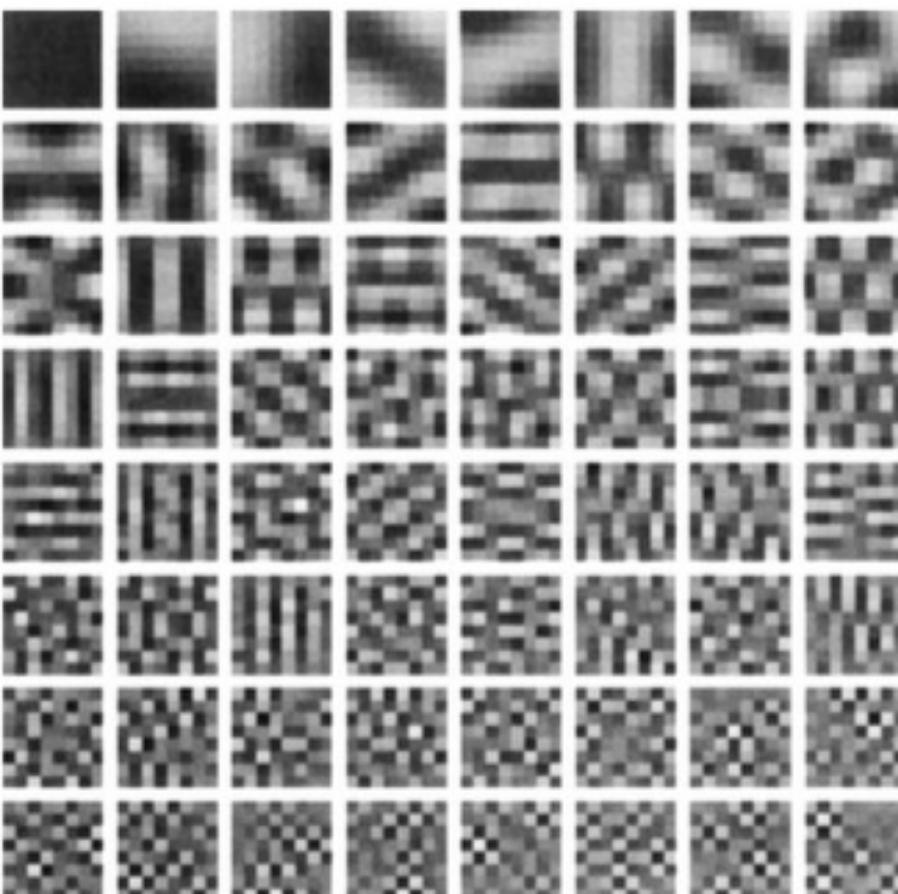
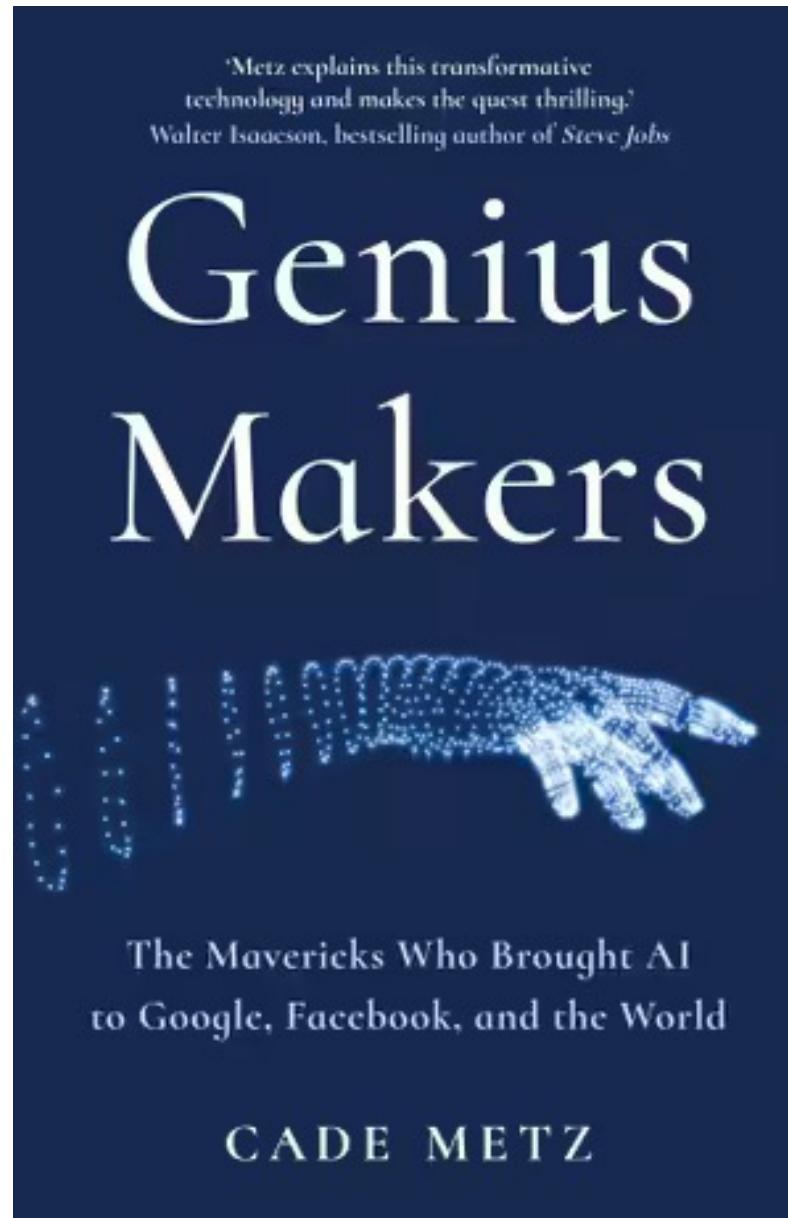


FIG. 1 Principal components calculated on  $8 \times 8$  image patches extracted from natural scenes by using Sanger's rule<sup>14</sup>. The full set of 64 components is shown, ordered by their variance (by columns, then by rows). The oriented structure of the first few principal components does not arise as a result of the oriented structures in natural images, but rather because these functions are composed of a small number of low-frequency components.

# The Debate Rages On!!!!



Yoshua Bengio



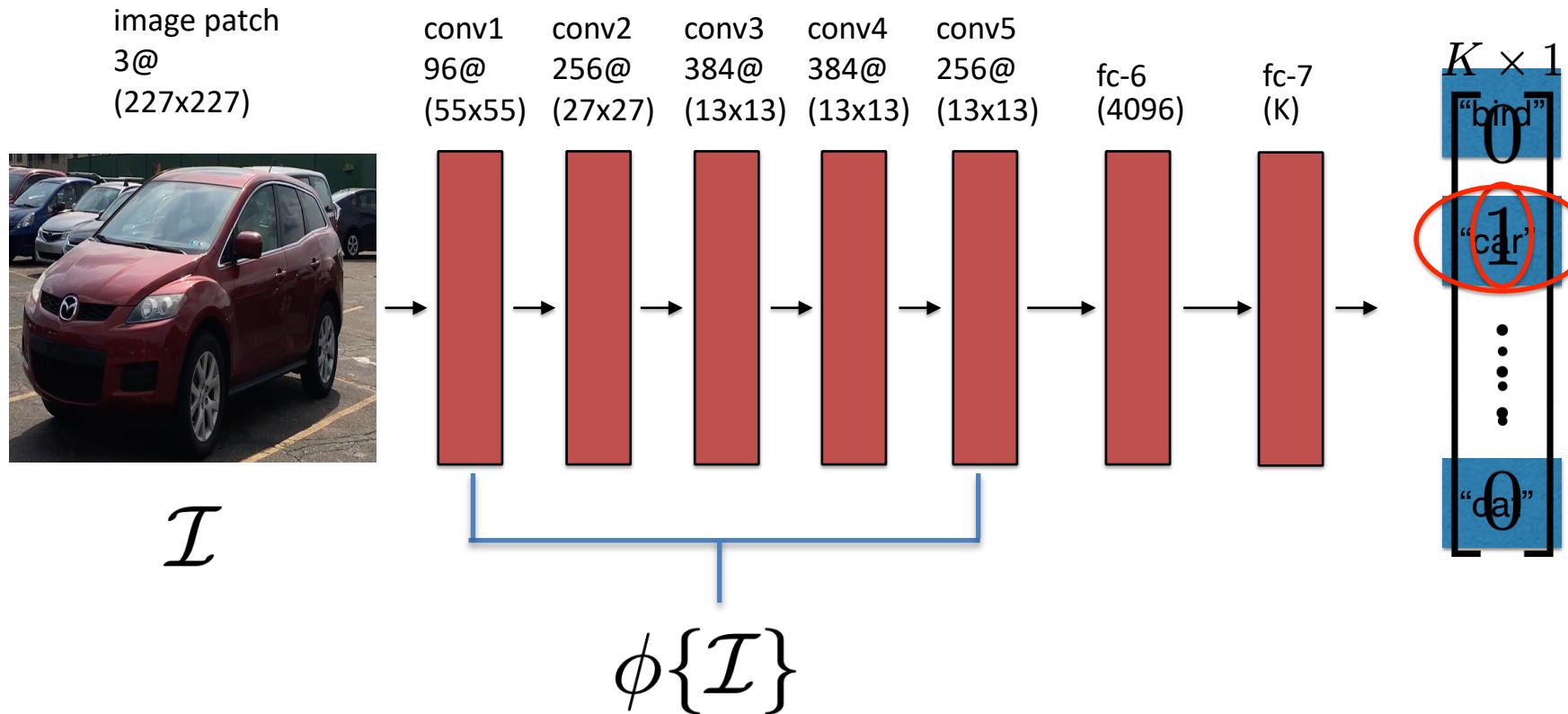
Gary Marcus

# Challenge

207	245	77	21	247	211	240	1
219	41	58	179	161	154	184	98
215	145	187	71	251	249	65	100
192	2	189	247	166	63	232	213
105	94	66	190	156	61	89	145
159	154	87	184	101	105	72	71
192	111	6	94	60	70	65	226
175	120	210	226	80	183	168	184
134	56	36	240	159	178	76	135
239	244	199	9	132	104	188	185
245	210	78	199	0	92	9	246
5	121	187	122	107	47	12	119
230	171	135	36	82	54	65	37
61	140	79	19	161	96	127	187
56	223	46	6	180	186	142	244
28	20	61	2	178	187	98	220

$$\phi\{\mathcal{I}\}$$

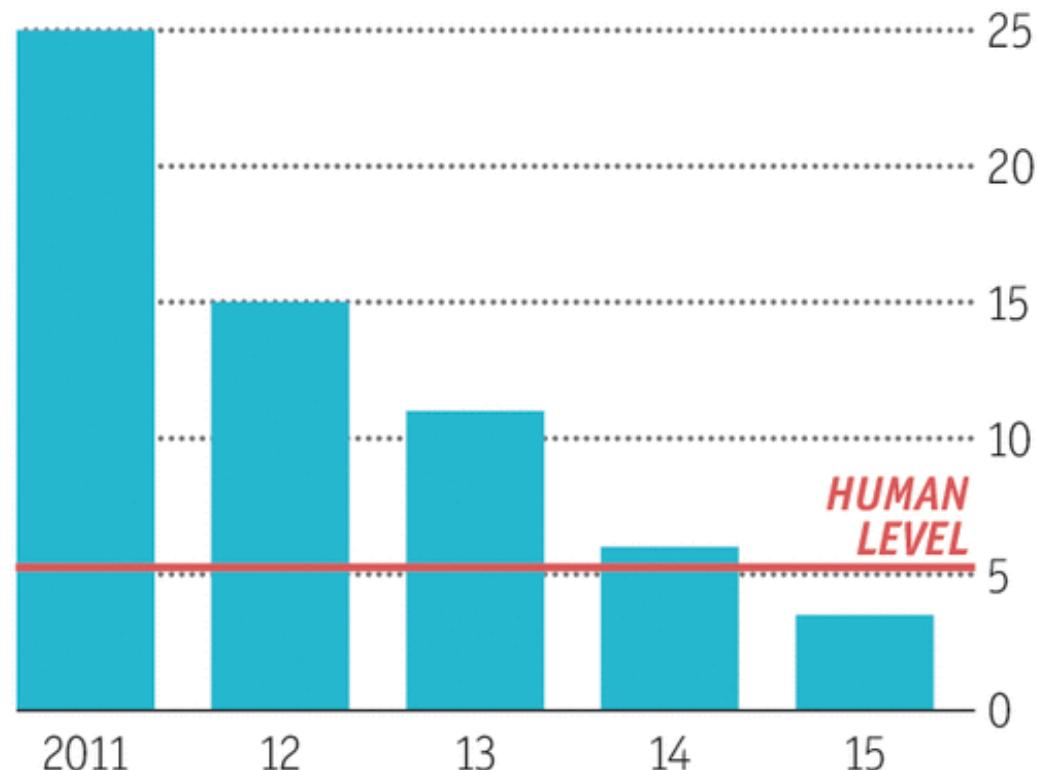
# Deep Learning - A Breakthrough!!!



# Deep Learning

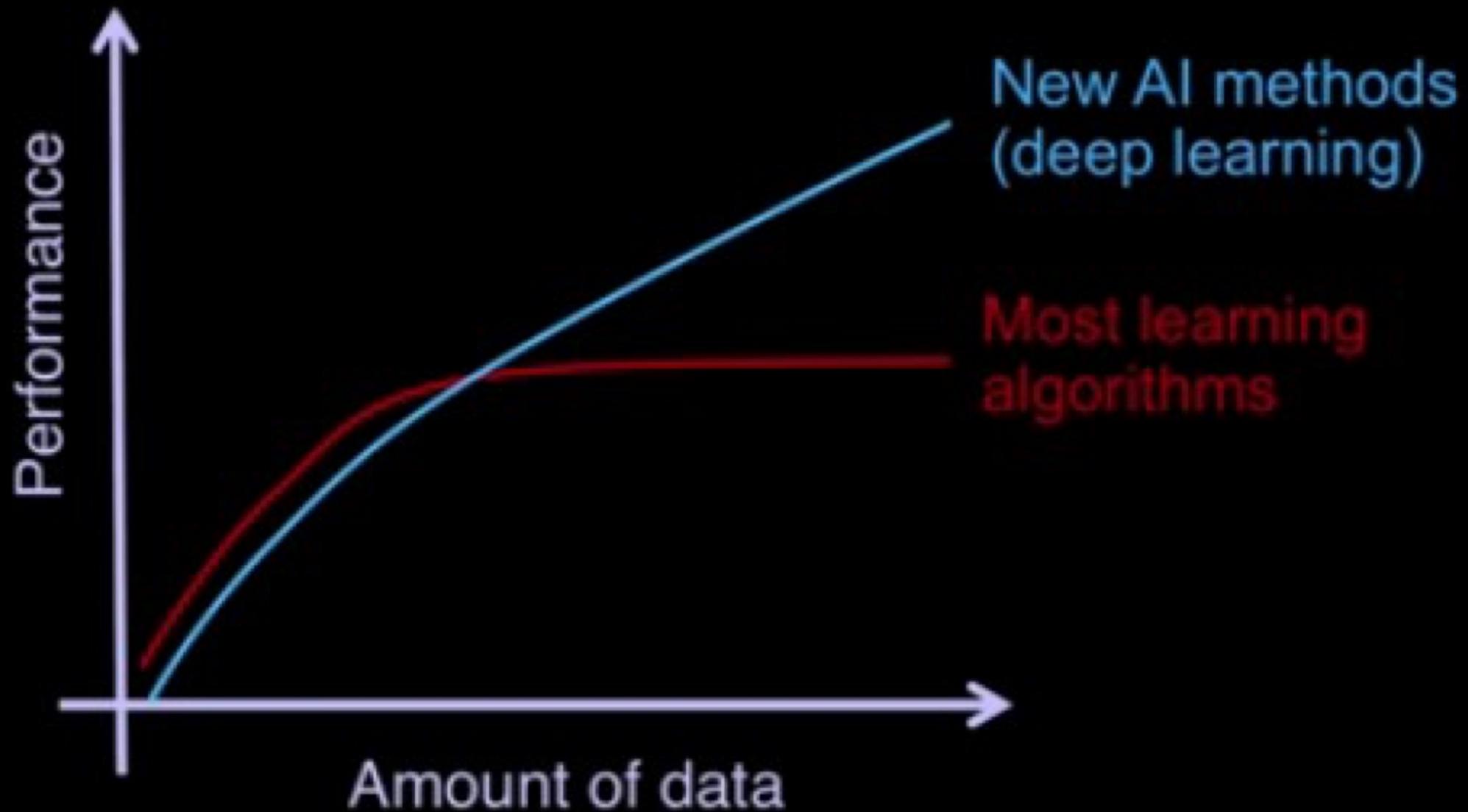
## Ever cleverer

Error rates on ImageNet Visual Recognition Challenge, %

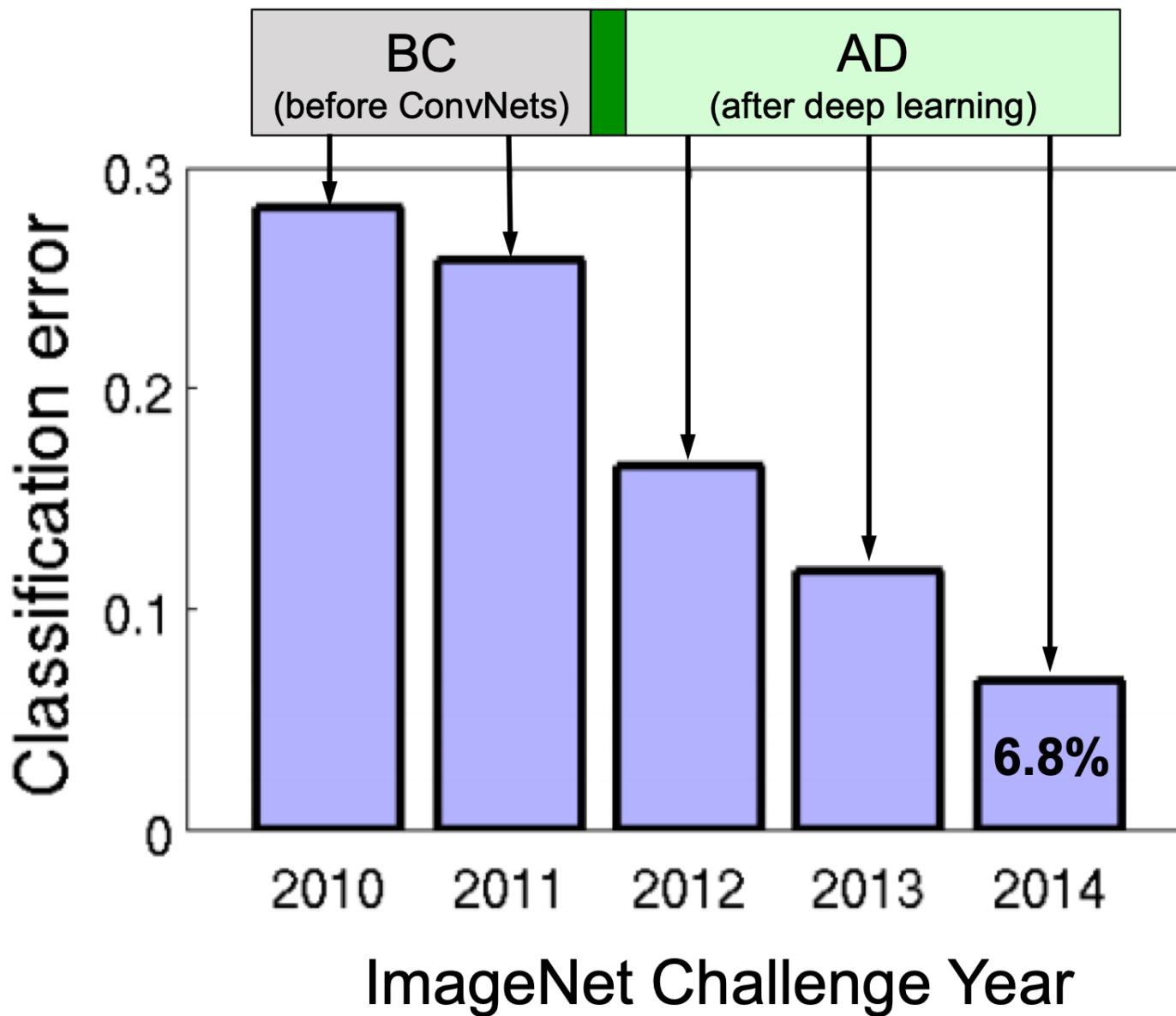


Sources: ImageNet; Stanford Vision Lab

# Data and machine learning

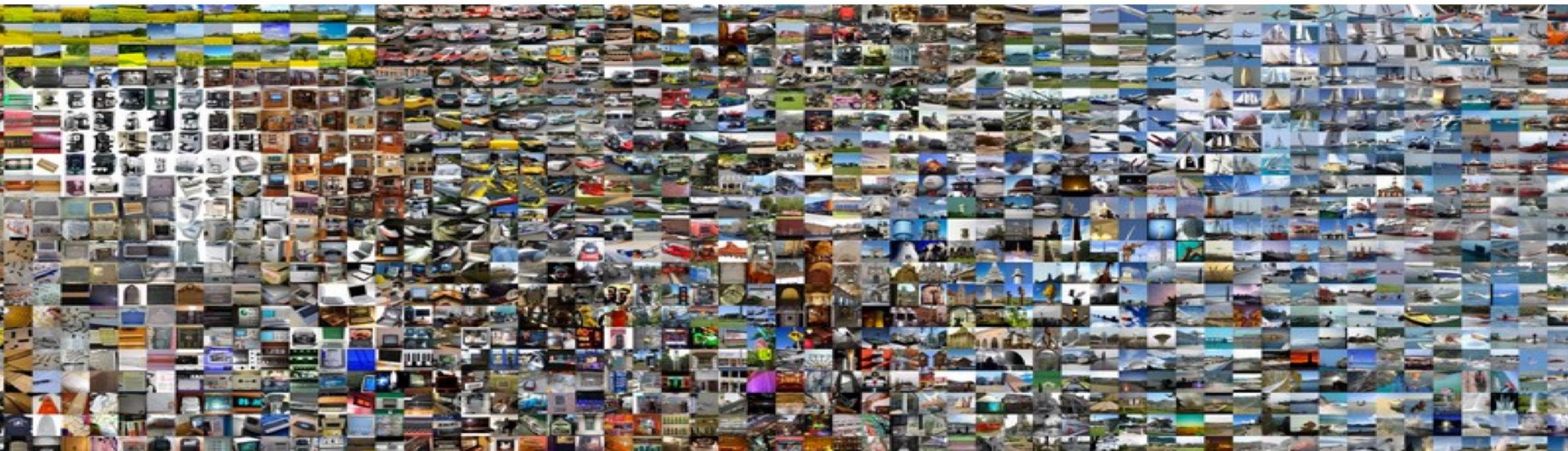


# Impact on Object Recognition



# ImageNet

- Over 15M labeled high resolution images.
- Roughly 22K categories.
- Collected from web and labeled by Amazon Mechanical Turk.



# Machines are starting to act like us!





## Audio

TIMIT Phone classification	Accuracy	TIMIT Speaker identification	Accuracy
Prior art (Clarkson et al., 1999)	79.6%	Prior art (Reynolds, 1995)	99.7%
Feature learning	<b>80.3%</b>	Feature learning	<b>100.0%</b>

## Images

CIFAR Object classification	Accuracy	NORB Object classification	Accuracy
Prior art (Ciresan et al., 2011)	80.5%	Prior art (Scherer et al., 2010)	94.4%
Feature learning	<b>82.0%</b>	Feature learning	<b>95.0%</b>

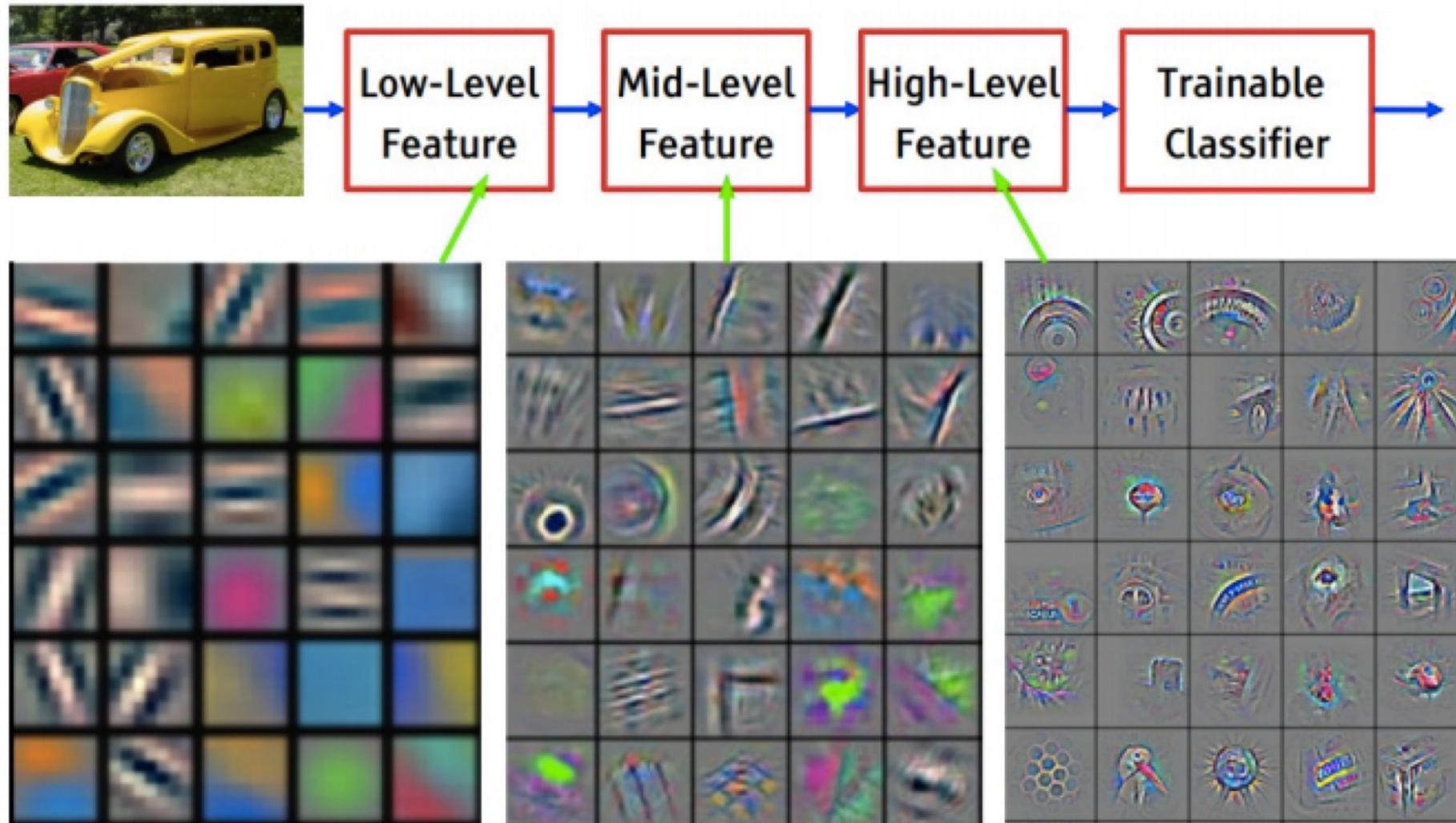
## Video

Hollywood2 Classification	Accuracy	YouTube	Accuracy
Prior art (Laptev et al., 2004)	48%	Prior art (Liu et al., 2009)	71.2%
Feature learning	<b>53%</b>	Feature learning	<b>75.8%</b>
KTH	Accuracy	UCF	Accuracy
Prior art (Wang et al., 2010)	92.1%	Prior art (Wang et al., 2010)	85.6%
Feature learning	<b>93.9%</b>	Feature learning	<b>86.5%</b>

## Text/NLP

Paraphrase detection	Accuracy	Sentiment (MR/MPQA data)	Accuracy
Prior art (Das & Smith, 2009)	76.1%	Prior art (Nakagawa et al., 2010)	77.3%
Feature learning	<b>76.4%</b>	Feature learning	<b>77.7%</b>

# Visualizing CNNs



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

# Mimicking ≠ Understanding

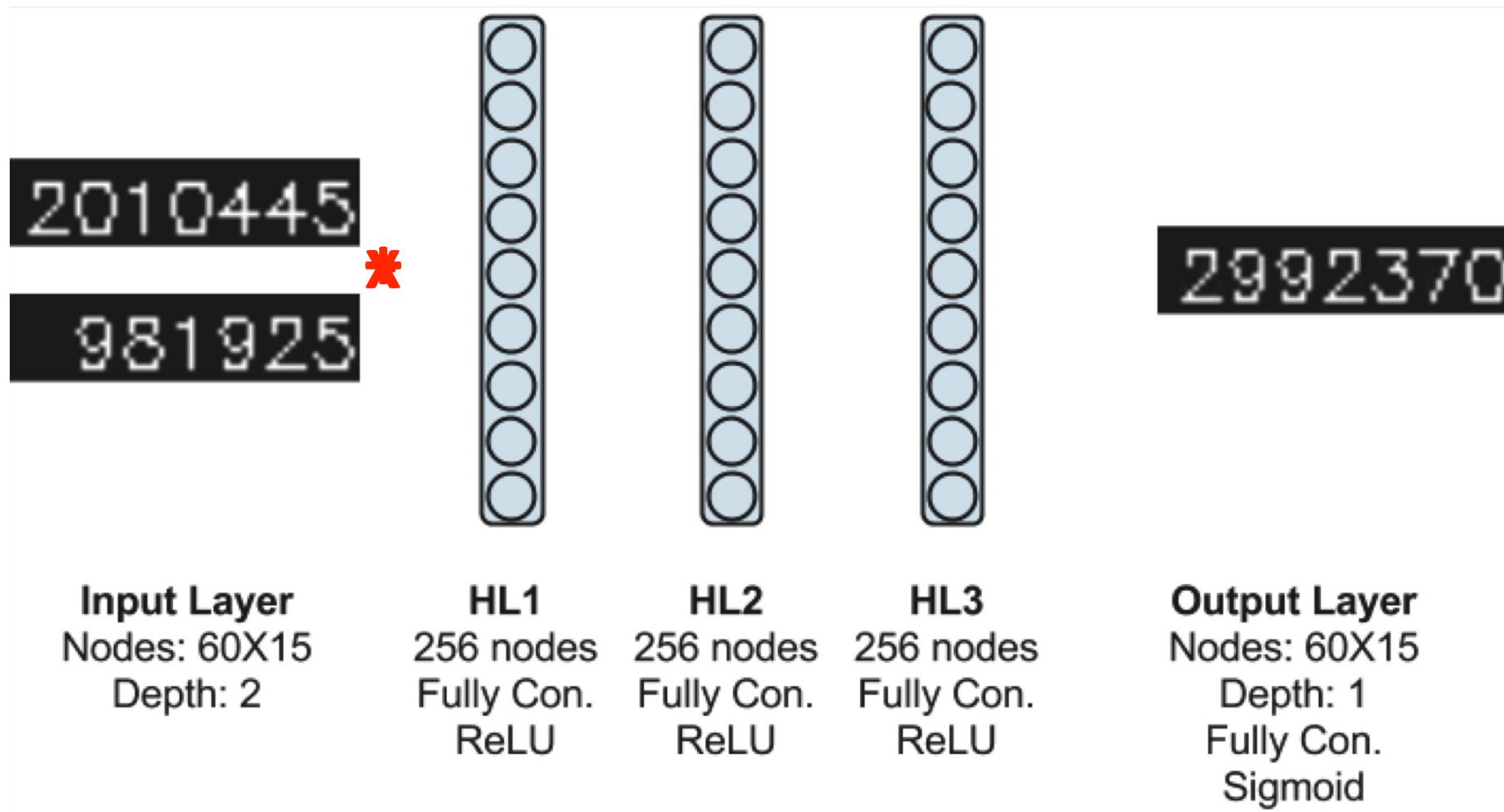
---



# Deep Learning & Visual Arithmetic

$$\begin{array}{r} 981925 \\ + 2010445 \\ \hline 2992370 \end{array}$$

# Deep Learning & Visual Arithmetic



# Deep Learning & Visual Arithmetic

	Example A	Example B
Input Picture 1	981925	3570002
Input Picture 2	2010445	1216536
Network Output Picture	2992370	4786538
Ground Truth Picture	2992370	4786538

# Deep Learning & Visual Arithmetic

Operation	Pictures		1-hot Vectors	
	No. Layers	% Error	No. Layers	% Error
Add	3	1.9%	1	1.7%
Subtract	3	3.2%	1	2.1%
Multiply	5	71.5%	3	37.6%
Roman Addition	5	74.3 %	3	0.7 %



# Multiplying Limits of GPT-4

# Faith and Fate: formers on

# Faith and Fate: Transformers on Compositionality

## Abstract

The correct answer is 866,133

**Abstract**

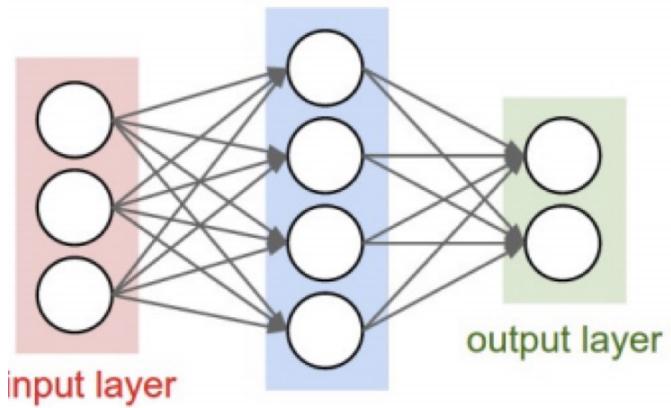
The correct answer is accuracy. Transformer large language models (LLMs) have sparked admiration for their exceptional performance on tasks that demand intricate multi-step reasoning. Yet these models simultaneously show failures on surprisingly trivial problems. This begs the question: Are these errors incidental, or do they signal more substantial limitations? In an attempt to demystify Transformers, we investigate the limits of these models across three representative compositional tasks—multi-digit multiplication, logic grid puzzles, and a classic dynamic programming problem. These tasks require breaking problems down into sub-steps and synthesizing them into a precise answer. We formulate compositional tasks as computational sub-procedures. Our empirical findings suggest a systematic way to quantify the level of complexity, and break down intermediate sub-tasks by reducing them to subgraph matching, without necessarily need for compositional matching. Our empirical study shows that subgraph matching can be used to quickly identify the correct answer. To round off our empirical study, we conducted a multi-step reasoning problem with increasing complexity. The results show that the model's performance decays rapidly as the number of steps increases.

# I Introduction

## *"It was the epoch of h-*

# Shallow vs. Deep Layers

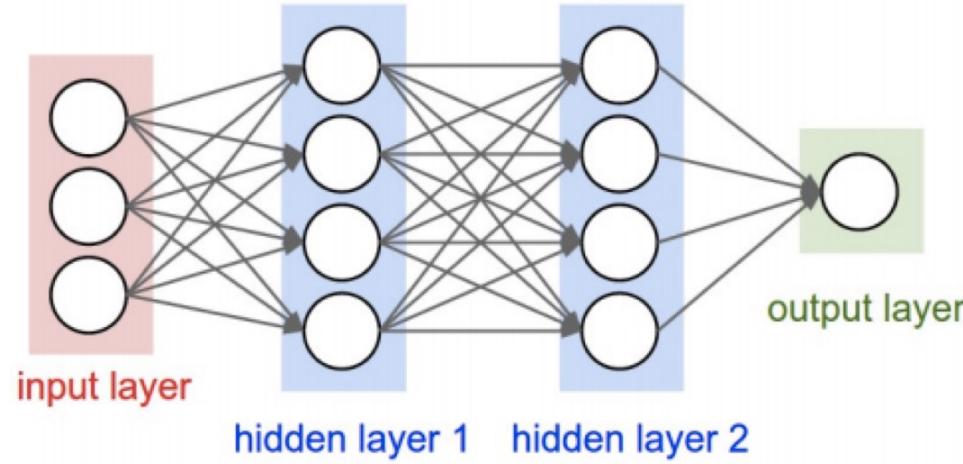
“shallow”



“2-layer neural net,” or  
“1-hidden-layer neural net”

**“Fully-connected” layers**

“deep”



“3-layer neural net,” or  
“2-hidden-layer neural net”

## Why deepness?

# Shallow vs. Deep Learning

---

- A shallow learner has only one hidden unit,

$$\mathbf{W}^{(2)} h(\mathbf{W}^{(1)} \mathbf{x})$$

- A deep learner has more than one hidden unit,

$$\mathbf{W}^{(3)} h(\mathbf{W}^{(2)} h(\mathbf{W}^{(1)} \mathbf{x}))$$



“Geoffrey Hinton”

# Today

---

- About the Course
- What is AI?
- The Rise of Deep Learning
- **Why Math in AI Matters**

# AI: Bakery vs. Chemistry?



# Many mysteries still remain?

DOI:10.1145/3446776

## Understanding Deep Learning (Still) Requires Rethinking Generalization

By Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals

### Abstract

Despite their massive size, successful deep artificial neural networks can exhibit a remarkably small gap between training and test performance. Conventional wisdom attributes small generalization error either to properties of the model family or to the regularization techniques used during training.

Through extensive systematic experiments, we show how these traditional approaches fail to explain why large neural networks generalize well in practice. Specifically, our experiments establish that state-of-the-art convolutional networks for image classification trained with stochastic gradient methods easily fit a random labeling of the training data. This phenomenon is qualitatively unaffected by explicit regularization and occurs even if we replace the true images by completely unstructured random noise. We corroborate these experimental findings with a theoretical construction showing that simple depth two neural networks already have perfect finite sample expressivity as soon as the number of parameters exceeds the number of data points as it usually does in practice.

We interpret our experimental findings by comparison with traditional models.

We supplement this republication with a new section at the end summarizing recent progresses in the field since the original version of this paper.

underwrites the generalization ability of a model has occupied the machine learning research community for decades.

There are a variety of theories proposed to explain generalization.

Uniform convergence, margin theory, and algorithmic stability are but a few of the important conceptual tools to reason about generalization. Central to much theory are different notions of *model complexity*. Corresponding generalization bounds quantify how much data is needed as a function of a particular complexity measure. Despite much significant theoretical work, the prescriptive and descriptive value of these theories remains debated.

This work takes a step back. We do not offer any new theory of generalization. Rather, we offer a few simple experiments to interrogate the empirical import of different purported theories of generalization. With these experiments at hand, we broadly investigate what practices do and do not promote generalization, what does and does not measure generalization?

### 1.1. The randomization test

In our primary experiment, we create a copy of the training data where we replace each label independently by a random label chosen from the set of valid labels. A dog picture labeled “dog” might thus become a dog picture labeled “air-

# More to read...

