# I. Mysteries of deep learning

**Yaoyu Zhang**

Institute of Natural Sciences & School of Mathematical Sciences

Shanghai Jiao Tong University

**FAU MoD Course**

**Towards a Mathematical Foundation of Deep Learning: From Phenomena to Theory**

**Date**
Fri. – Thu. May 2 – 8, 2025

**Session Titles**
1. Mysteries of Deep Learning
2. Frequency Principle/Spectral Bias
3. Condensation Phenomenon
4. From Condensation to Loss Landscape Analysis
5. From Condensation to Generalization Theory

# Unimaginable achievements of AI

Krizhevsky, et al, 2012



https://www.linkedin.com/pulse/must-read-path-breaking-papers-image-classification-muktabh-mayank

$p_{\sigma/\rho}(a|s)$

$v_\theta(s')$

$s$

$s'$

Silver, et al, 2017

a — Value network

d — Policy network

Silver, et al, 2017

T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

The Nobel Prize in Chemistry 2024

David Baker

"for computational protein design"



© Nobel Prize Outreach. Photo: Clément Morin

Demis Hassabis

"for protein structure prediction"



© Nobel Prize Outreach. Photo: Clément Morin

John Jumper

"for protein structure prediction"



© Nobel Prize Outreach. Photo: Clément Morin

https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology

Jumper et al., 2021

Turing test?

**Test**

Input: $x_1, \cdots, x_t$ → ChatGPT →

Output: $P(x_{t+1}|x_1, \cdots, x_t)$

**Training**

ChatGPT

**Empirical risk:** $R_S(\boldsymbol{\theta}) = \frac{1}{n}\sum_{i=1}^{n} l(f(\boldsymbol{x_i}, \boldsymbol{\theta}), \boldsymbol{y_i})$

**Model:** $f(\boldsymbol{x}, \boldsymbol{\theta})$

**Data:** $S = \{(x_i, y_i)\}_{i=1}^{n}$

## Just this?

**Common Models:**

Linear models: polynomial models, random feature models, $\cdots$

Neural networks: fully-connected, convolutional, ResNet, Transformer, $\cdots$

**Common loss function:**

Mean-squared error (l2) loss: $l(y, y') = \|y - y'\|_2^2,$

Cross entropy, Hinge loss, …

**Common training algorithm:**

Gradient decent (GD): $\theta^{t+1} = \theta^t - \eta \nabla R_S(\theta^t),$

Stochastic gradient descent (SGD), Adam, …

**Deep neural network**

Input layer    Multiple hidden layers    Output layer

$$\theta := \left( W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]} \right)$$

$$f_\theta^{[l]}(x) := \sigma(W^{[l]} f_\theta^{[l-1]}(x) + b^{[l]})$$

**AlphaGo, AlphaFold, ChatGPT, SORA, ...**

https://www.ibm.com/cloud/learn/neural-networks

- Synthetic diamond

- Atomic bomb

- The Apollo Program

- ChatGPT

- Quantum computer

- 0.1 light-speed spaceship

# Bitter lesson for deep learning theory

**The Bitter Lesson**

**Rich Sutton**

**March 13, 2019**

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of

"The biggest lesson that can be read from **70 years of AI research** is that **general methods that leverage computation** are ultimately the most effective, and by a large margin."

**Leverage computation (learning) instead of human knowledge**

http://www.incompleteideas.net/IncIdeas/BitterLesson.html

- **1969 book of *Perceptrons* (**lead to the first winter)

- **1986 Backpropagation** (emergence of modern deep learning)

- **1989 Universal approximation theorem**

- **1995 Generalization puzzle proposed** (not well solved till now)

  **The Vapnik-Jackel Bet** (witnessed by Yann Lecun)

- **2017 Generalization puzzle demonstrated in SOTA settings**

- **2018 Neural Tangent Kernel** (lead to a surge in DL theory research)

  **Frequency principle/Spectral bias**

**Despite 40 years of effort, framework for its math foundation yet to emerge**

From Lecun's talk

**Intelligent Machines**

# The Dark Secret at the Heart of AI

No one really knows how the most advanced algorithms do what they do. That could be a problem.

by Will Knight     April 11, 2017



**L**ast year, a strange self-driving car was released onto the quiet roads of Monmouth County, New Jersey. The experimental vehicle, developed by researchers at the chip maker Nvidia, didn't look different from other autonomous cars, but it was unlike anything demonstrated by Google, Tesla, or General Motors, and it showed the rising power of artificial intelligence. The car didn't follow a single instruction provided by an engineer or programmer. Instead, it relied entirely on an algorithm that had taught itself to drive by watching a human do it.

Donoho's PPT, Stats 385 Stanford

Figure : Every theorist who looks at it see what they wish

## A (personal) bitter lesson:

All previously existing frameworks, irrespective of their origin or demonstrated success, are ineffective for understanding deep learning.

## Existing frameworks:

statistical learning theory, numerical analysis, statistical physics, statistics, optimization, neuroscience, psychology, …

**In face of deep learning, all of us are blind men.**

https://www.sloww.co/blind-men-elephant/

1. **Suspension:** Suspend the prior and belief one may hold and focus on the facts about the object.

2. **Cumulation:** Discover and cumulate all possible facts about the object. Prioritize the more informative ones.

3. **Emergence:** A new framework shall emerge once enough pieces are uncovered.



https://www.sloww.co/blind-men-elephant/

**Suspension**

**Cumulation**

**Emergence**

**Frequency principle/spectral bias**

**Condensation**

Double descent

Edge of stability

Lottery ticket

Neural collapse
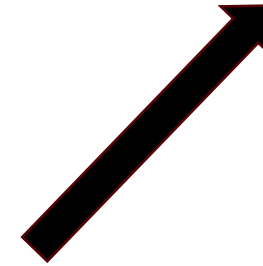
Grokking

……

# Basics of deep learning theory

## Single artificial neuron:

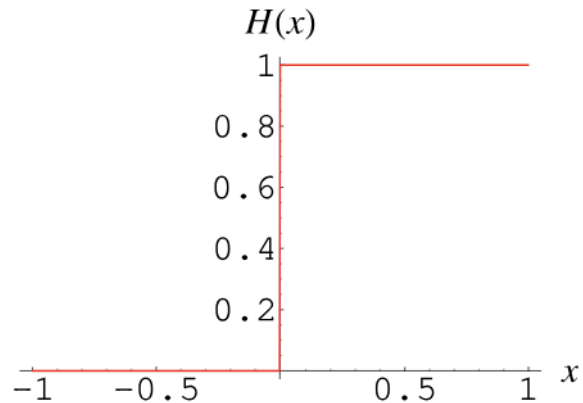$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b)$$

Parameters (weights): $\boldsymbol{\theta} = (\boldsymbol{w}, \boldsymbol{b})$, activation function: $\sigma(\cdot): \mathbb{R} \to \mathbb{R}$

## Illustration:

# Deep neural networks:

**Deep neural network**

Input layer    Multiple hidden layers    Output layer



$$\boldsymbol{\theta} := \left( \boldsymbol{W}^{[1]}, \boldsymbol{b}^{[1]}, \ldots, \boldsymbol{W}^{[L]}, \boldsymbol{b}^{[L]} \right)$$

$$\boldsymbol{f}_{\boldsymbol{\theta}}^{[l]}(\boldsymbol{x}) := \sigma \left( \boldsymbol{W}^{[l]} \boldsymbol{f}_{\boldsymbol{\theta}}^{[l-1]}(\boldsymbol{x}) + \boldsymbol{b}^{[l]} \right)$$

https://www.ibm.com/cloud/learn/neural-networks

> **Neural networks with a single hidden layer can be used to approximate any continuous function to any desired precision.**

Cybenko 89, Hornik 89, Hornik 91, Barron 93

Requirement for transfer function:

$\sigma(z)$ is well-defined as $z \to -\infty$ and $z \to \infty$

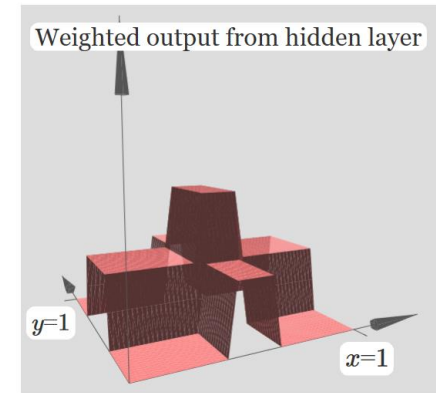$$\left| f(x) - \sum_j k_j \sigma(w_{ij} x_i + b_j) \right| < \epsilon$$
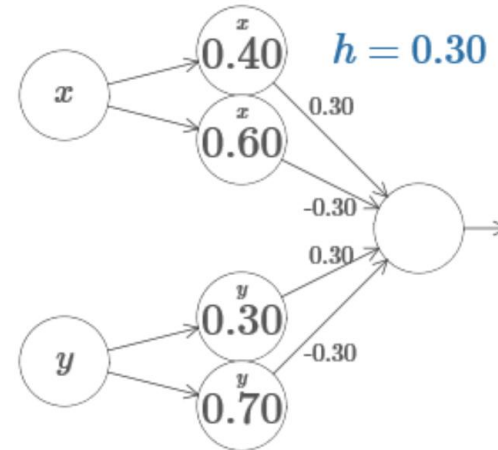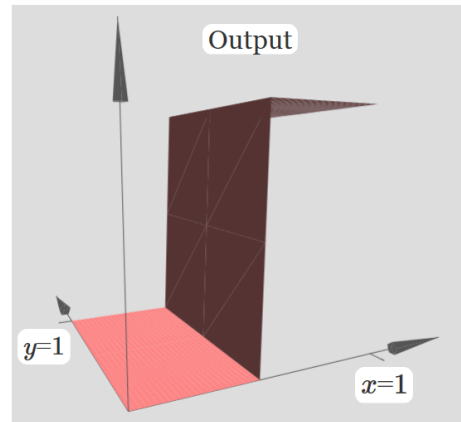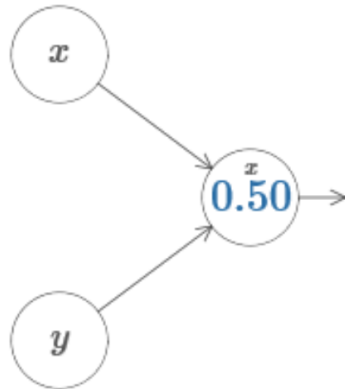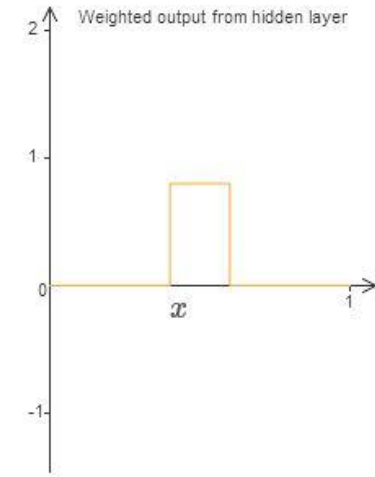
**Sketch of a constructive proof:**

1. Construct Heaviside function from the given transfer function
2. Construct "bump" function (1-d) or "tower" function (2-d)
3. Approximate the target continuous function with "bump" or "tower" functions

$s = 0.40$

Output from top hidden neuron

$s_1 = 0.40$

$w_1 = 0.8$

$x$

$s_2 = 0.60$

$w_2 = \text{-}0.8$

Weighted output from hidden layer

$x$

Output

$x$

$0.50$

$y$

$y=1$

$x=1$

$x$

$0.40$

$x$

$0.60$

$h = 0.30$

0.30

-0.30

0.30

$y$

$0.30$

-0.30

$y$

$0.70$

Weighted output from hidden layer

$y=1$

$x=1$

**Theorem**—Given a finite set $V$ and a finite set $S$ of real numbers, assume that $f : V \rightarrow S$ is chosen at random according to uniform distribution on the set $S^V$ of all possible functions from $V$ to $S$. For the problem of optimizing $f$ over the set $V$, then no algorithm performs better than blind search.



**How to infer the missing spot?**

https://en.wikipedia.org/wiki/No_free_lunch_theorem

**Generalization**

**Optimization**

**Approximation**

**Robustness**

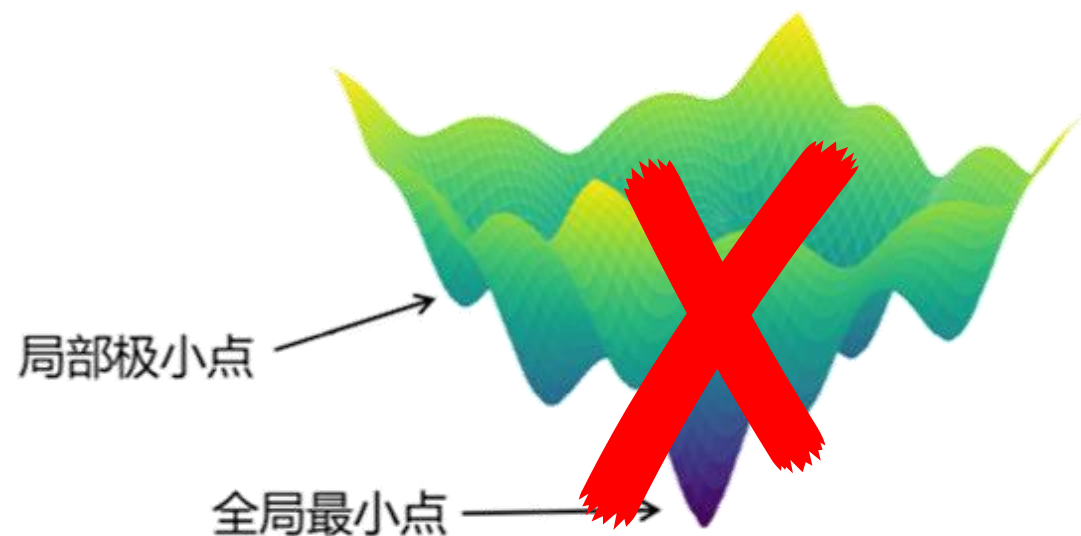**Interpretability**

**…**

**Phenomenon**

Despite strongly nonconvex loss landscape, gradient-based training of large DNNs often find global minima.



局部极小点

全局最小点

**Problem**

What is the geometry of loss landscape?

**Phenomenon**

Some architectures are more parameter efficient than others regarding particular class of tasks.

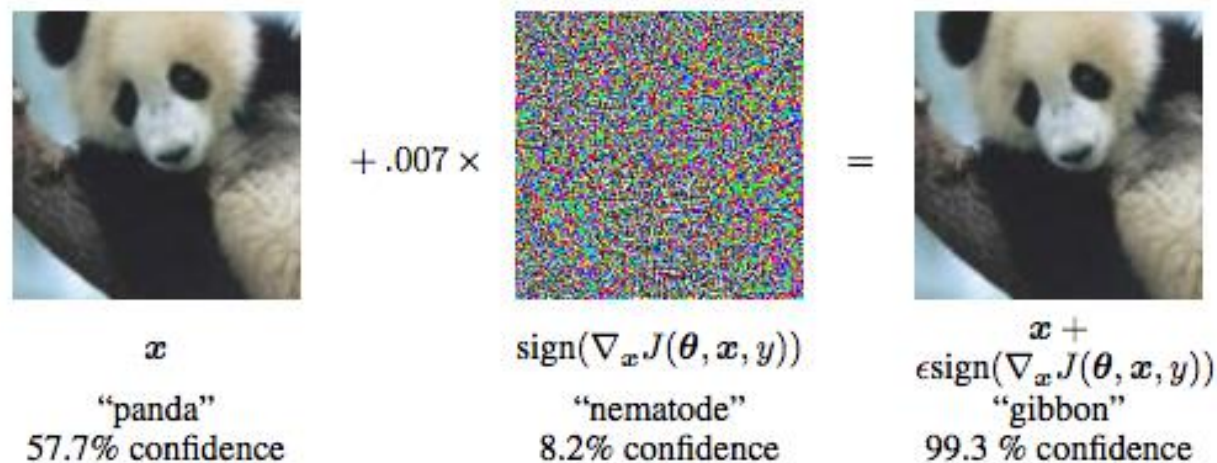**Ex:**CNN vs. FNN for image, Transformer vs. LSTM for language

**Problem**

How to quantify the difference in parameter efficiency between architectures?

## Phenomenon

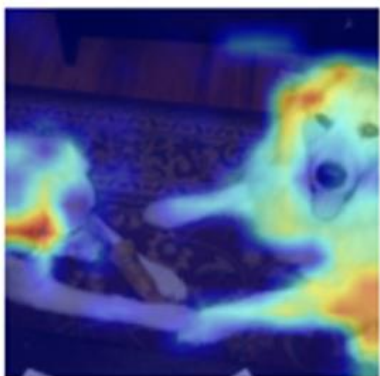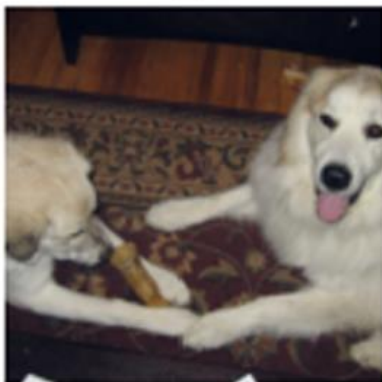Output of well-trained DNNs are often susceptible to tiny adversarial perturbation.



$$+.007 \times$$

$$x$$
"panda"
57.7% confidence

$$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"nematode"
8.2% confidence

$$x + \epsilon \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"gibbon"
99.3 % confidence

Goodfellow et al.

## Problem

Why is that? How to improve robustness?

## Phenomenon

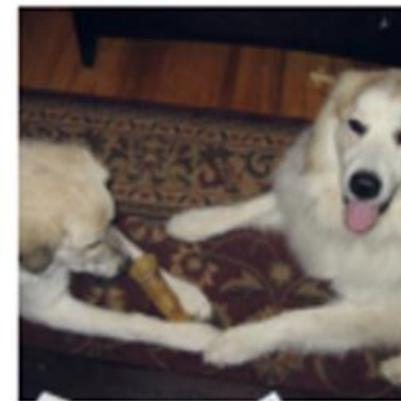One can hardly obtain an explanation with prediction power.



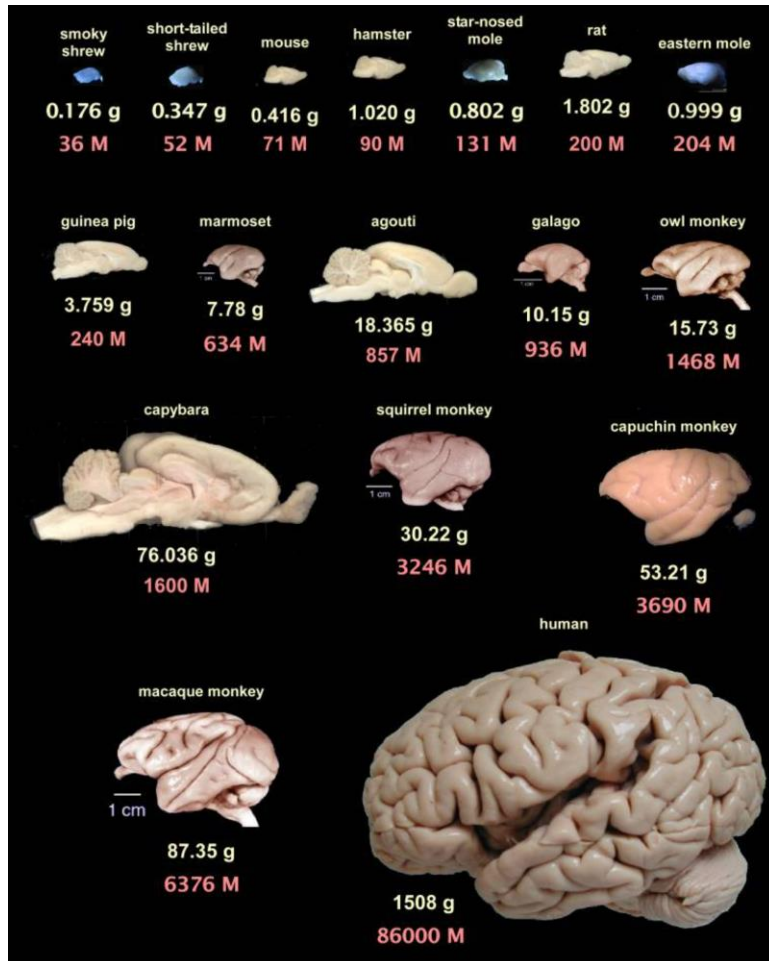## Problem

When is it possible to obtain explanations with prediction power?

# **Generalization puzzle of deep learning**

Suzana Herculano-Houzel, 2009
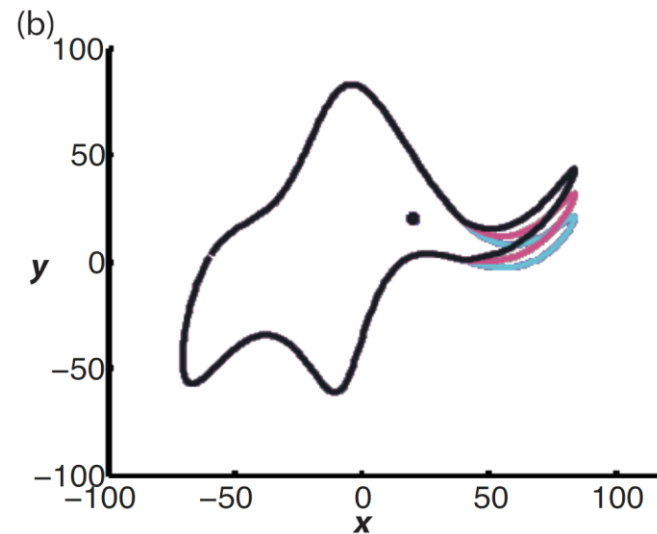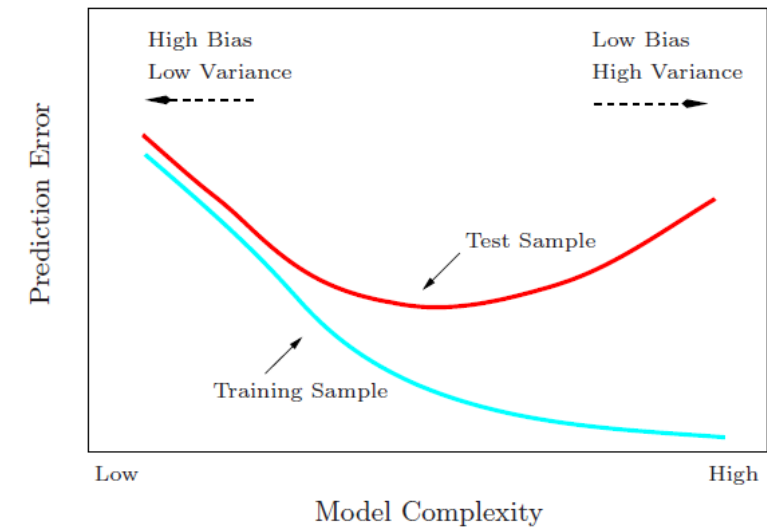
**"With four parameters you can fit an elephant to a curve; with five you can make him wiggle his trunk."**
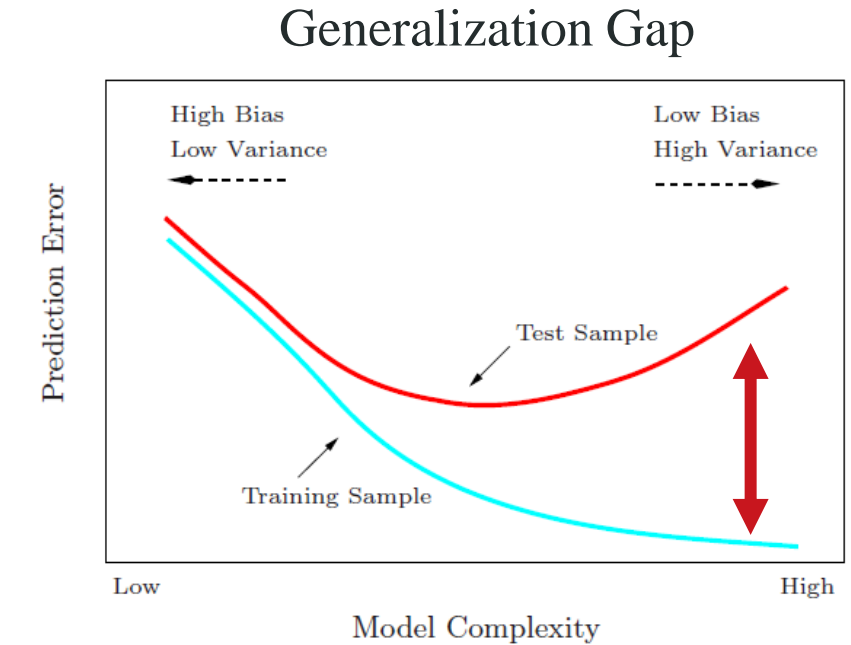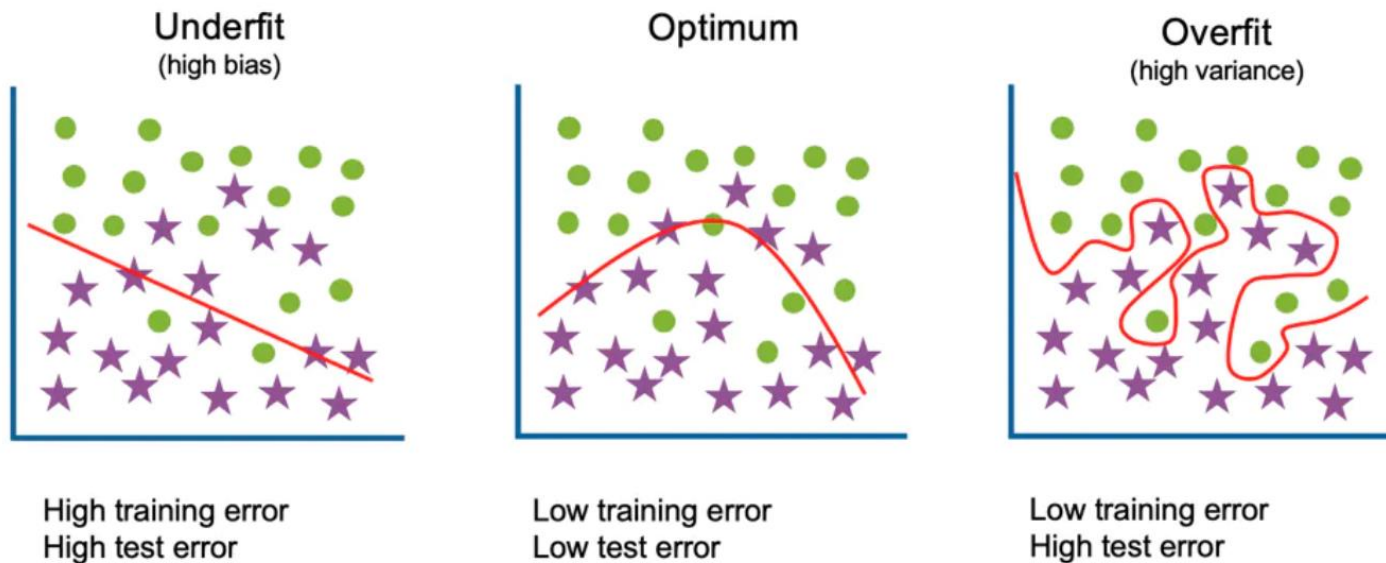
-- John von Neumann



Mayer et al., 2010

**Complex models easily overfit.**

**Large complexity → Large generalization gap**



**Occam Razor:** Entities should not be multiplied unnecessarily

**1995**

**Leo Breiman**

Statistics Department, University of California, Berkeley, CA 94305;
e-mail: leo@stat.berkeley.edu

# Reflections After Refereeing Papers for NIPS

Our fields would be better off with far fewer theorems, less emphasis on faddish stuff, and much more scientific inquiry and engineering. But the latter requires real thinking.

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?

UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

**Chiyuan Zhang***
Massachusetts Institute of Technology
chiyuan@mit.edu

**Samy Bengio**
Google Brain
bengio@google.com

**Moritz Hardt**
Google Brain
mrtz@google.com

**Benjamin Recht**[†]
University of California, Berkeley
brecht@berkeley.edu

**Oriol Vinyals**
Google DeepMind
vinyals@google.com

**Cifar10: 60,000 training data**

| model | # params | random crop | weight decay | train accuracy | test accuracy |
|---|---|---|---|---|---|
| Inception | 1,649,402 | yes | yes | 100.0 | 89.05 |
| | | yes | no | 100.0 | 89.31 |
| | | no | yes | 100.0 | 86.03 |
| | | no | no | 100.0 | 85.75 |
| (fitting random labels) | | no | no | 100.0 | 9.78 |

Zhang et al., 2017

Find an interpolation of $\mathcal{D}: \{(x_i, y_i)\}_{i=1}^n$ in $\mathcal{H}: \{h(\cdot; \Theta) | \Theta \in \mathbb{R}^m\}$

Example:
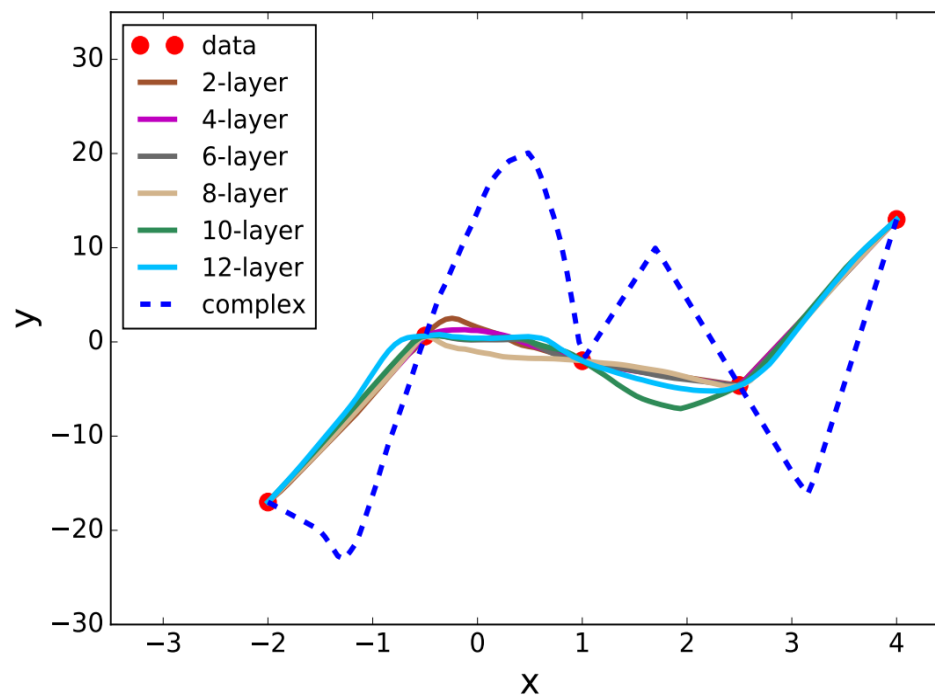$h(x; \Theta) = \theta_1 + \theta_2 x + \cdots + \theta_M x^{m-1}$ with $m = n$



Traditional wisdom: $m < n$.

Modern wisdom?
Using neural network with $m \gg n$.

Lei Wu, Zhanxing Zhu, Weinan E, 2017