

How to reason about Machine Learning

Overview and Introduction

Michael Ngo, April 19th 2025

Hi, I'm Michael Ngo

Pronounced “No”, he/him

- 3rd year undergraduate studying CS
- Research learning theory w/ prof Michael P. Kim
- Teaching assistant for Algorithms (CS 4820)
- Onboarding Chair for Cornell Data Science project team
- Spends way to much time on YouTube



After a really bad haircut

Agenda

1. What is machine learning?
2. Why is machine learning so hot right now?
3. How we can begin to reason about machine learning?
4. What do we have to look out for?

What is machine learning?

Super Mario Bros

How to
program a
computer to
complete this
level?



Super Mario Bros

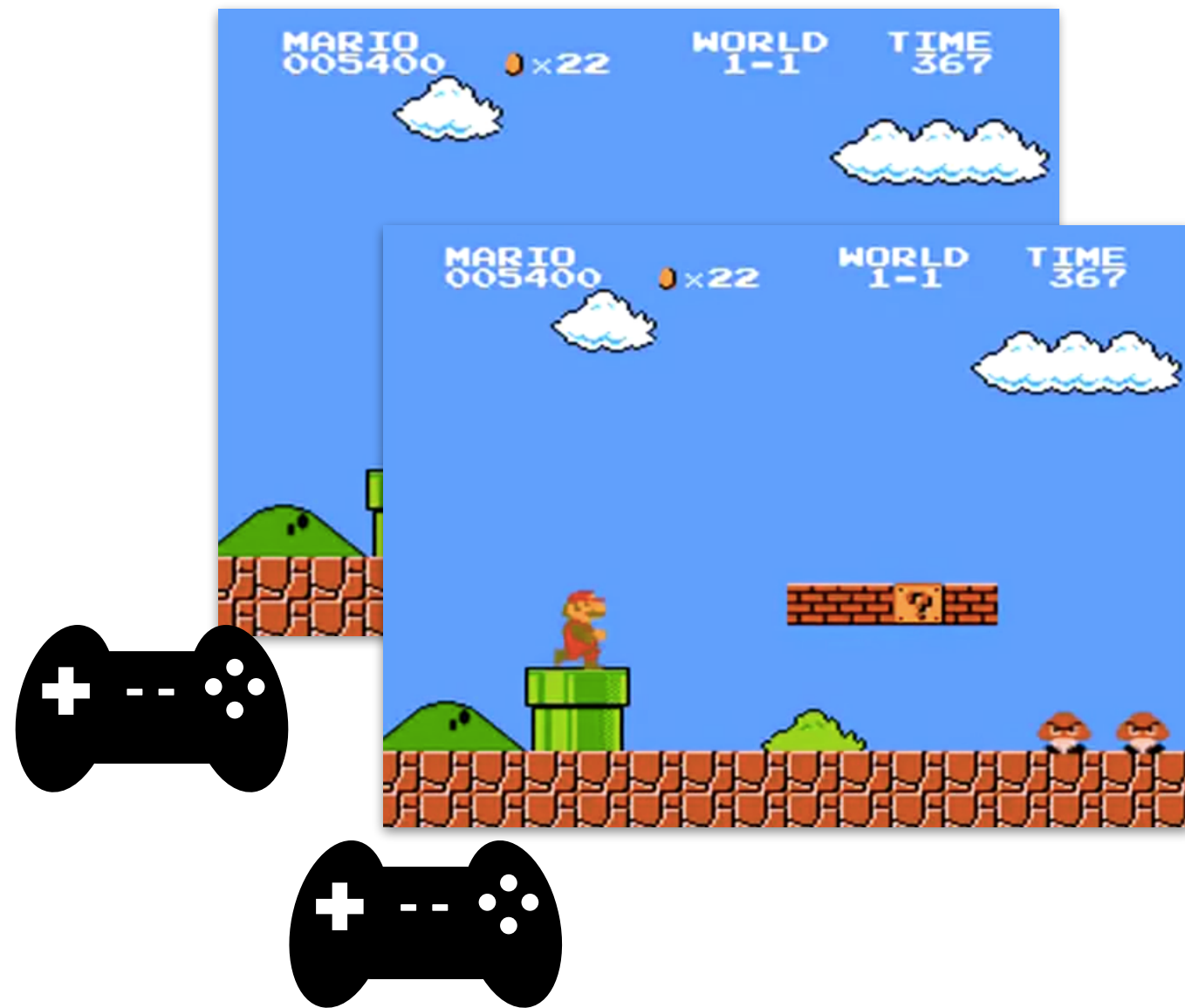
How to
program a
computer to
complete this
level?



Naive:
“Move forward,
jump when see enemy,
Enter pipe”

ML allows computer to learn like humans

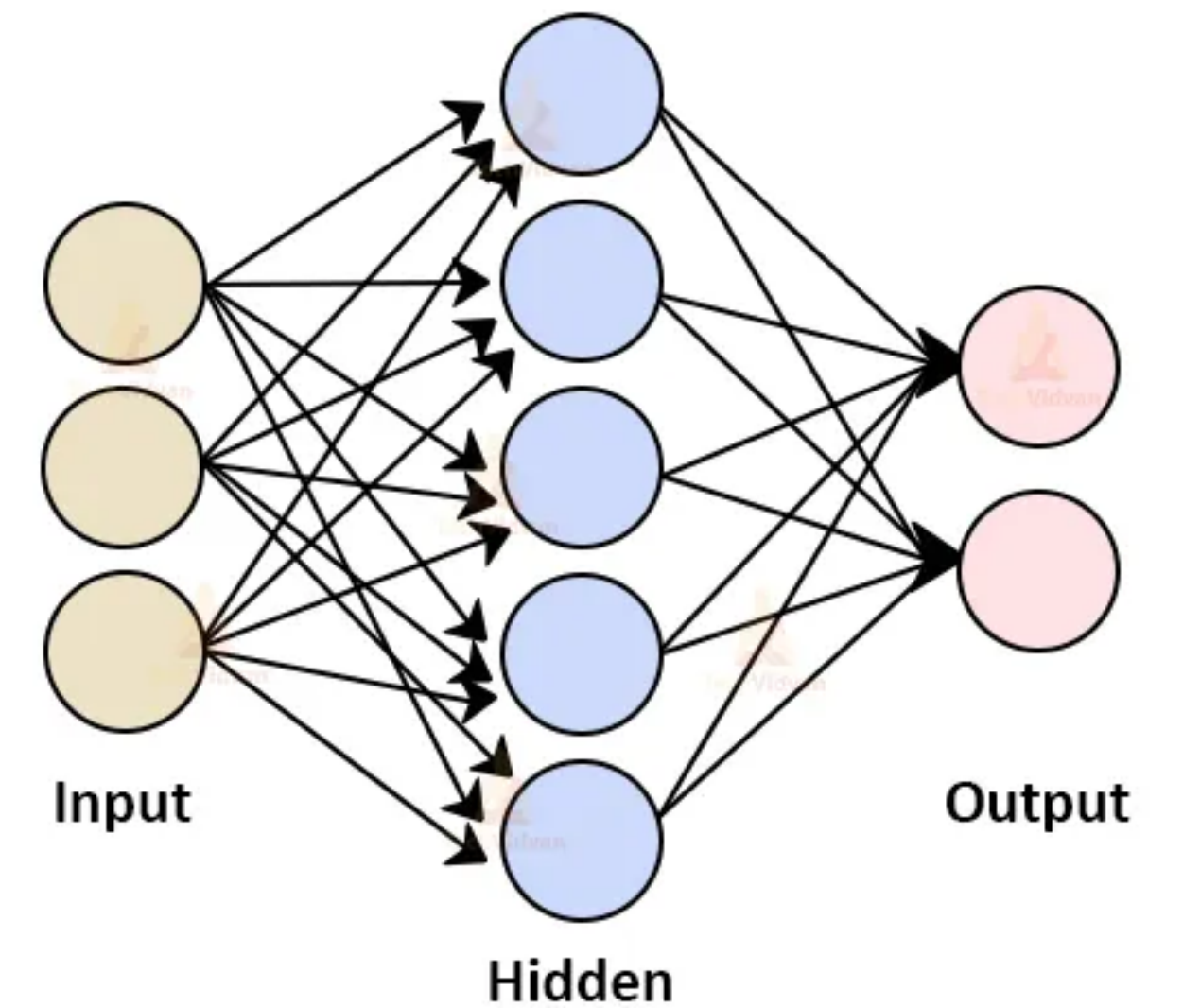
Human Demonstrations



ML
Algorithm



Architecture of
Artificial Neural Network



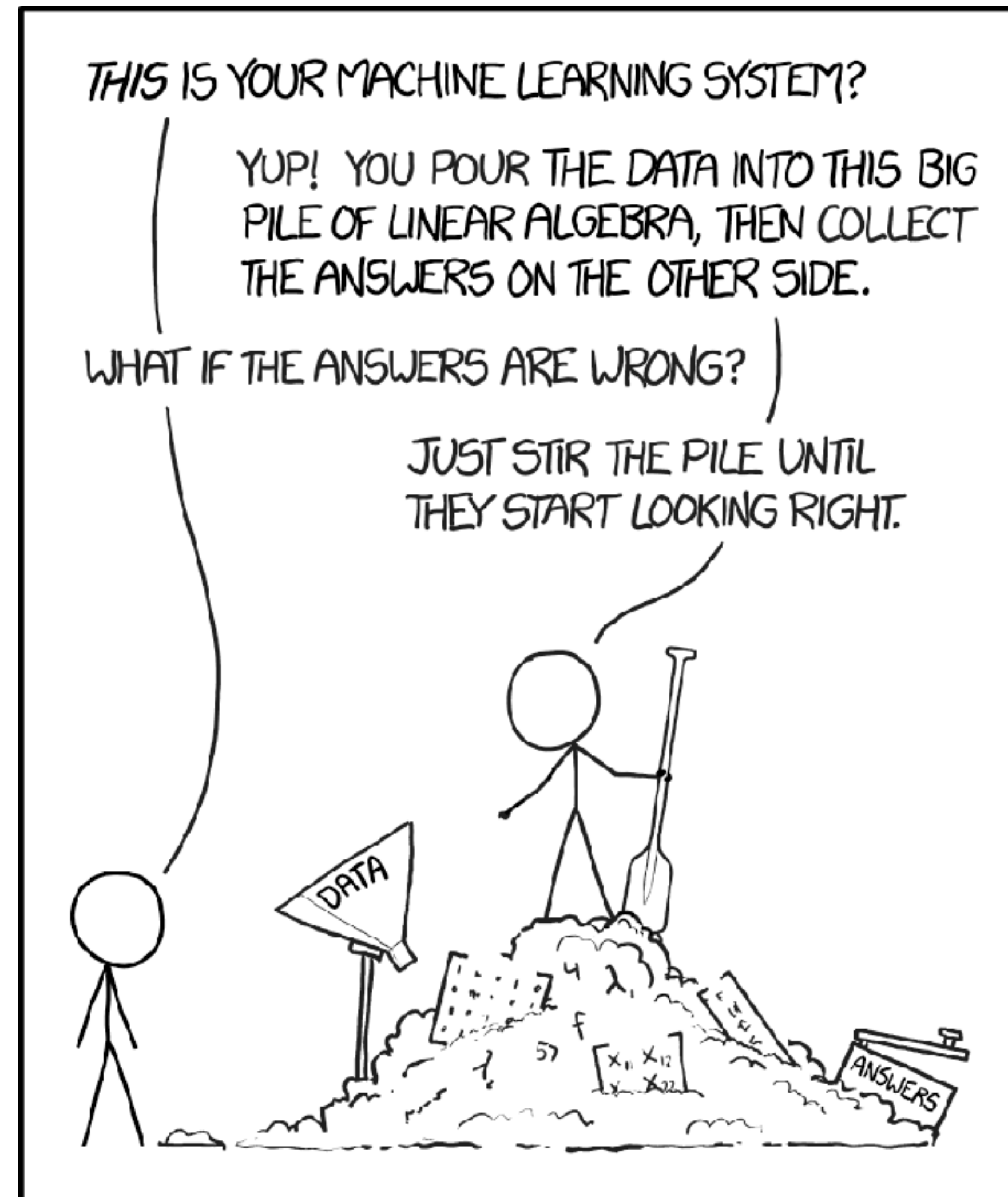
A computer learned this!



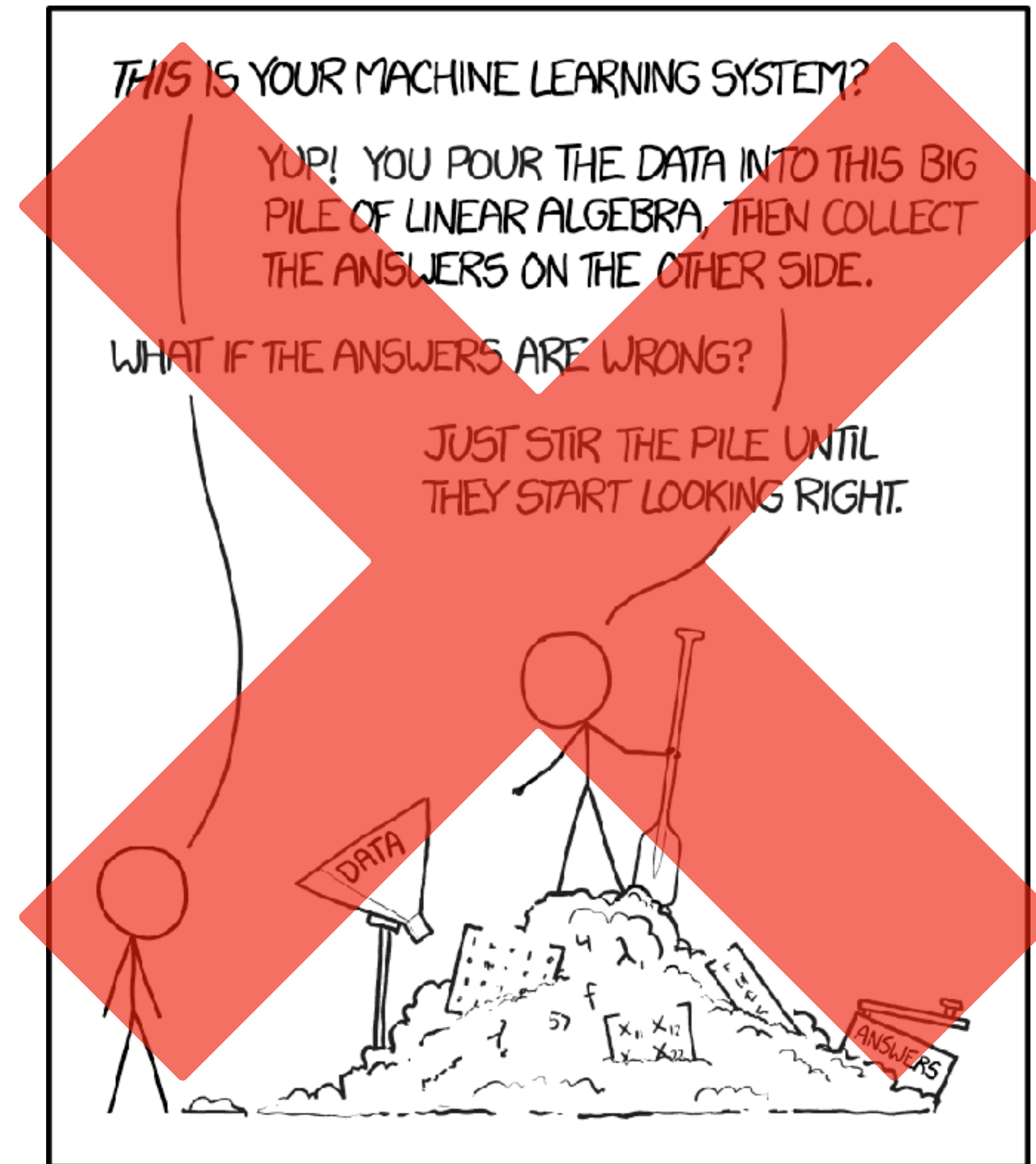
A computer learned this!



A misconception



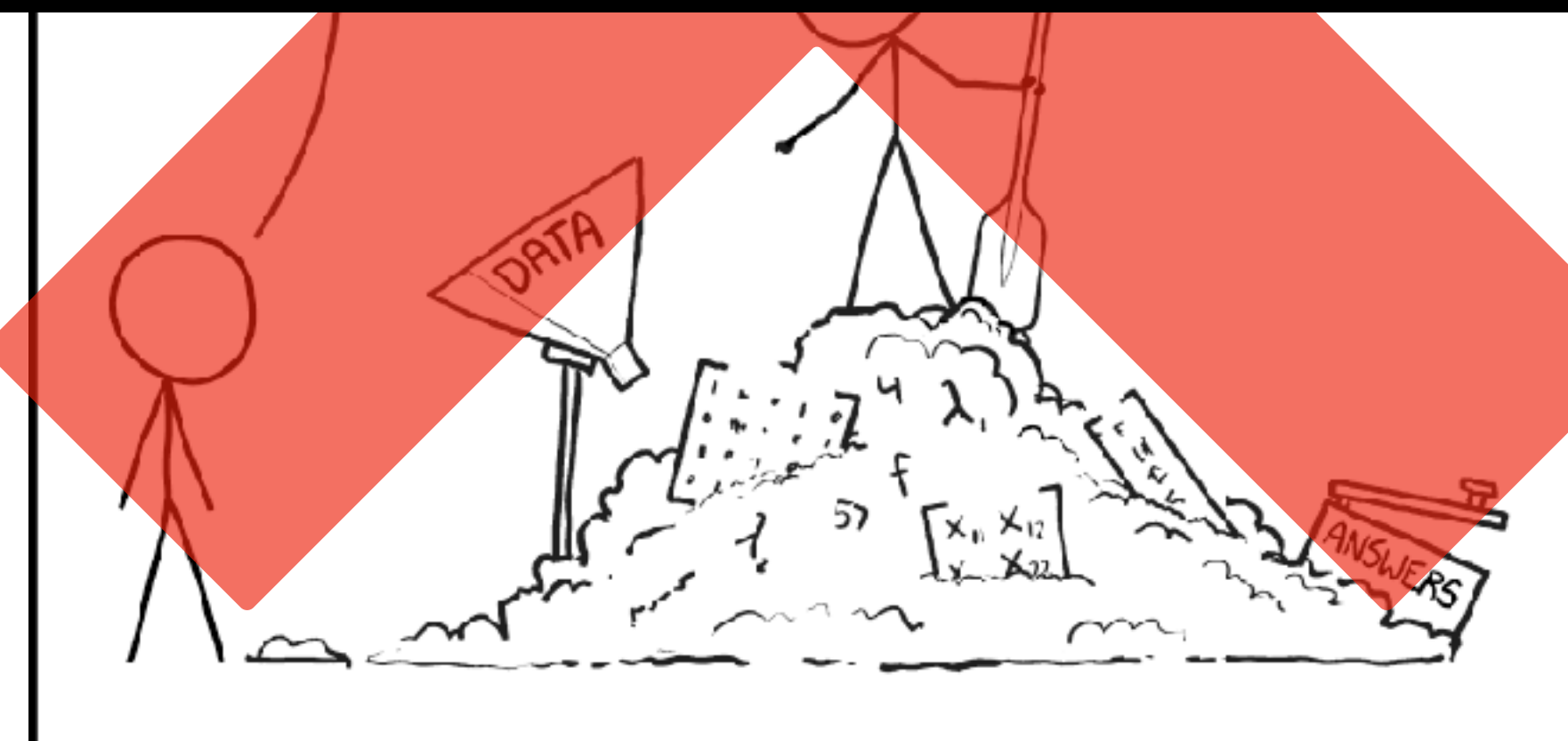
A misconception



A misconception

THIS IS YOUR MACHINE LEARNING SYSTEM?
YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.
WHAT IF THE ANSWERS ARE WRONG?

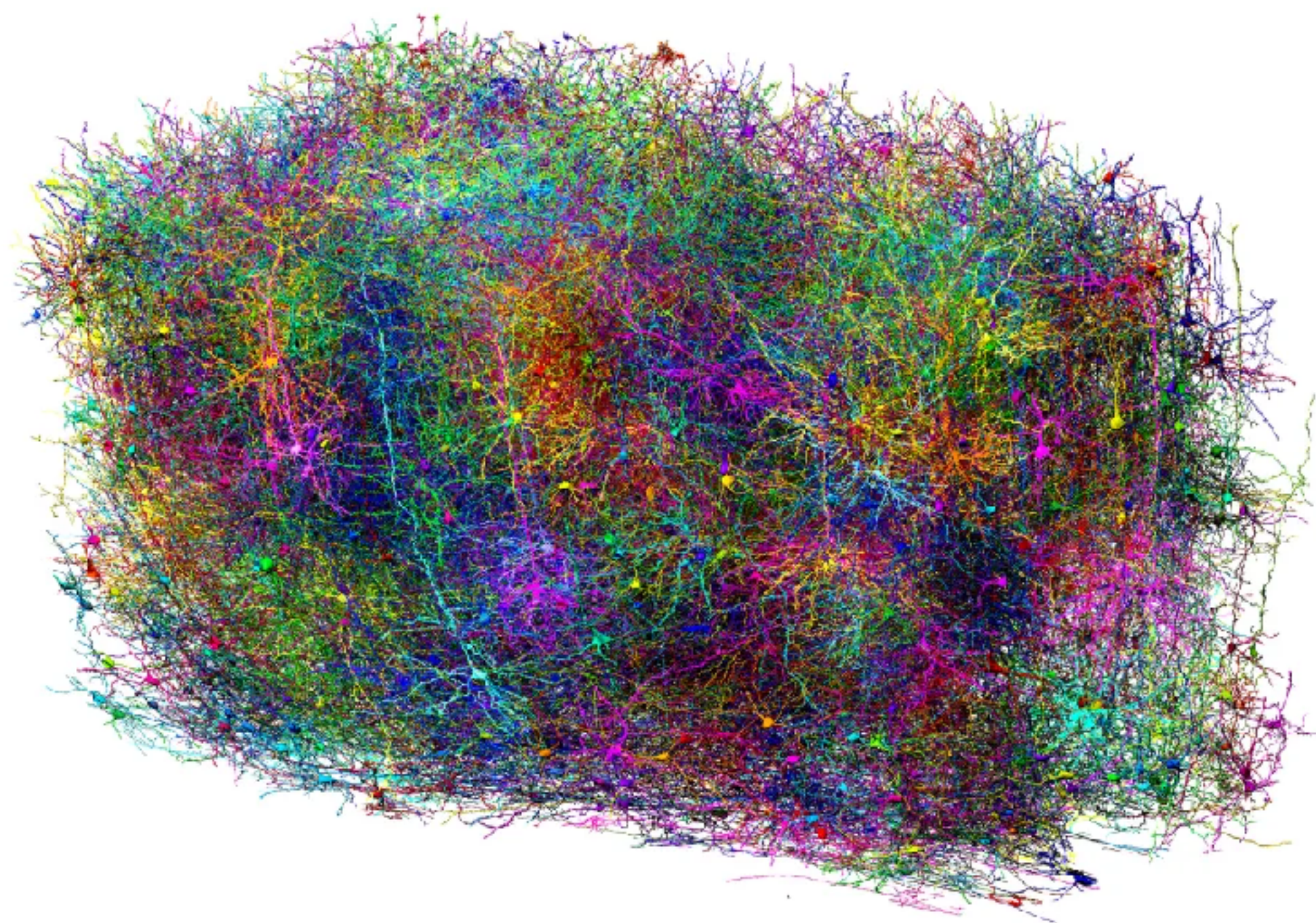
We need a principled way to
reason about Machine Learning!



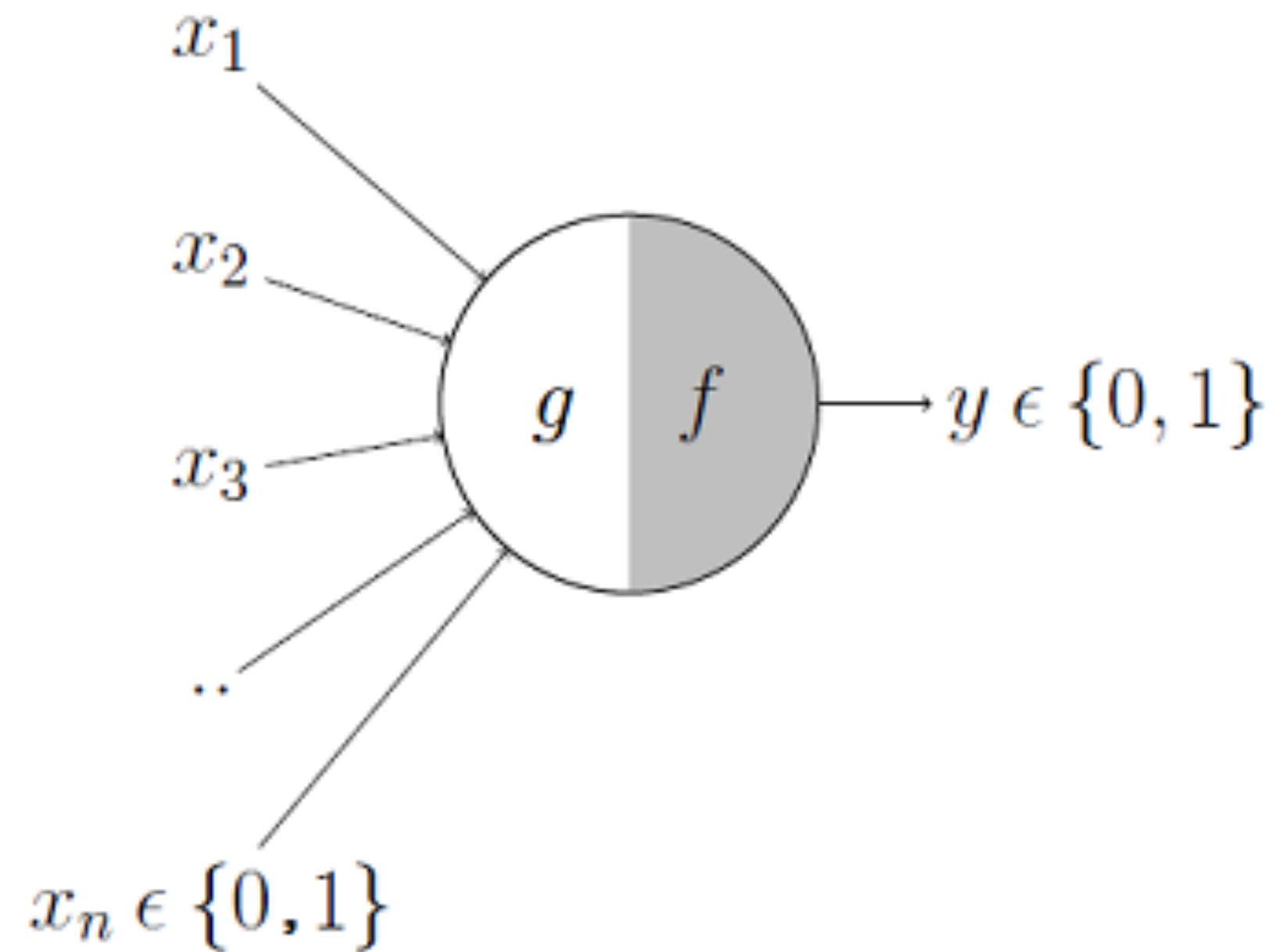
Why should I care about machine learning?

Why should I care about machine learning?

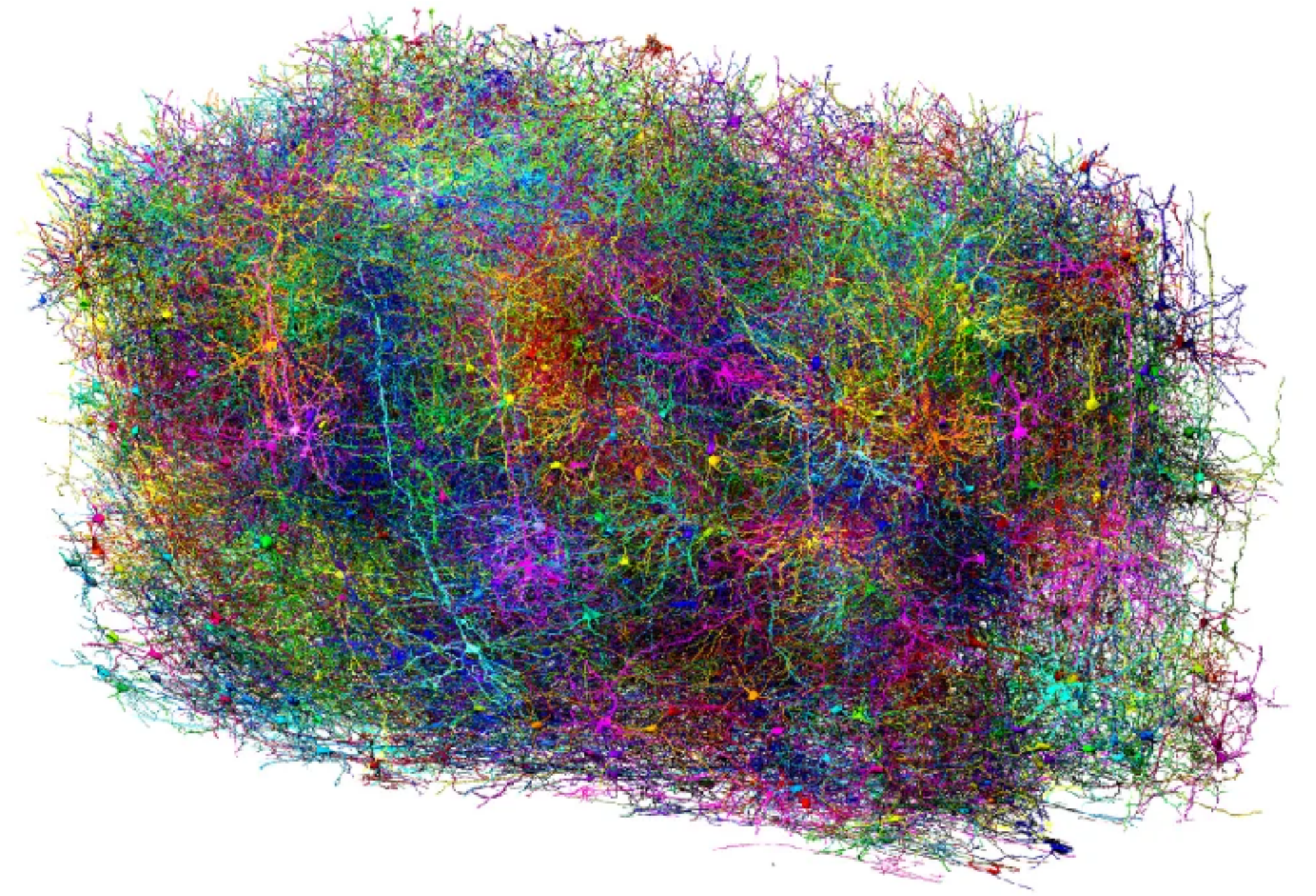
Let me give a *brief* history...



Artificial Neural Network (1943)

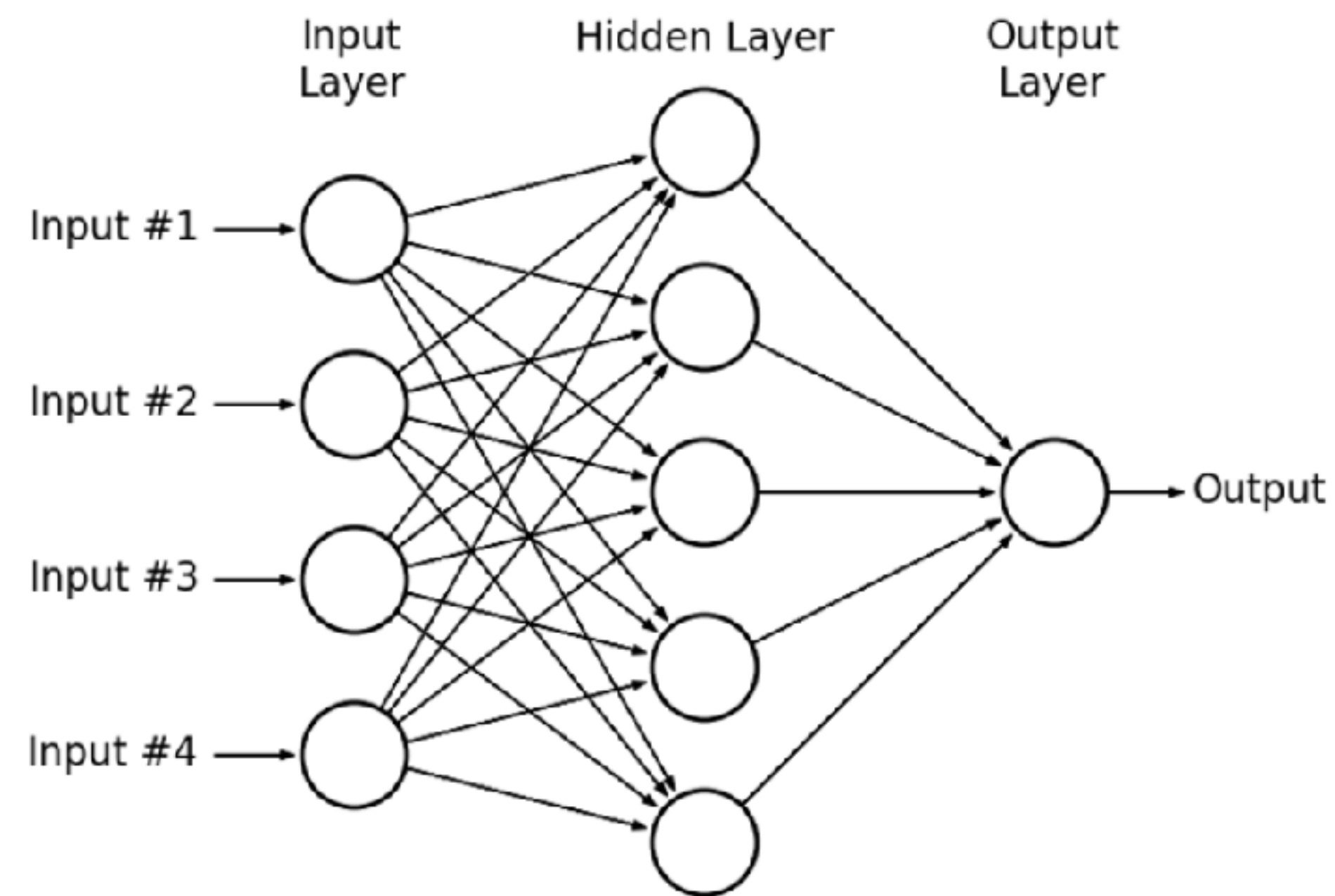


McCulloch-Pitts Neuron



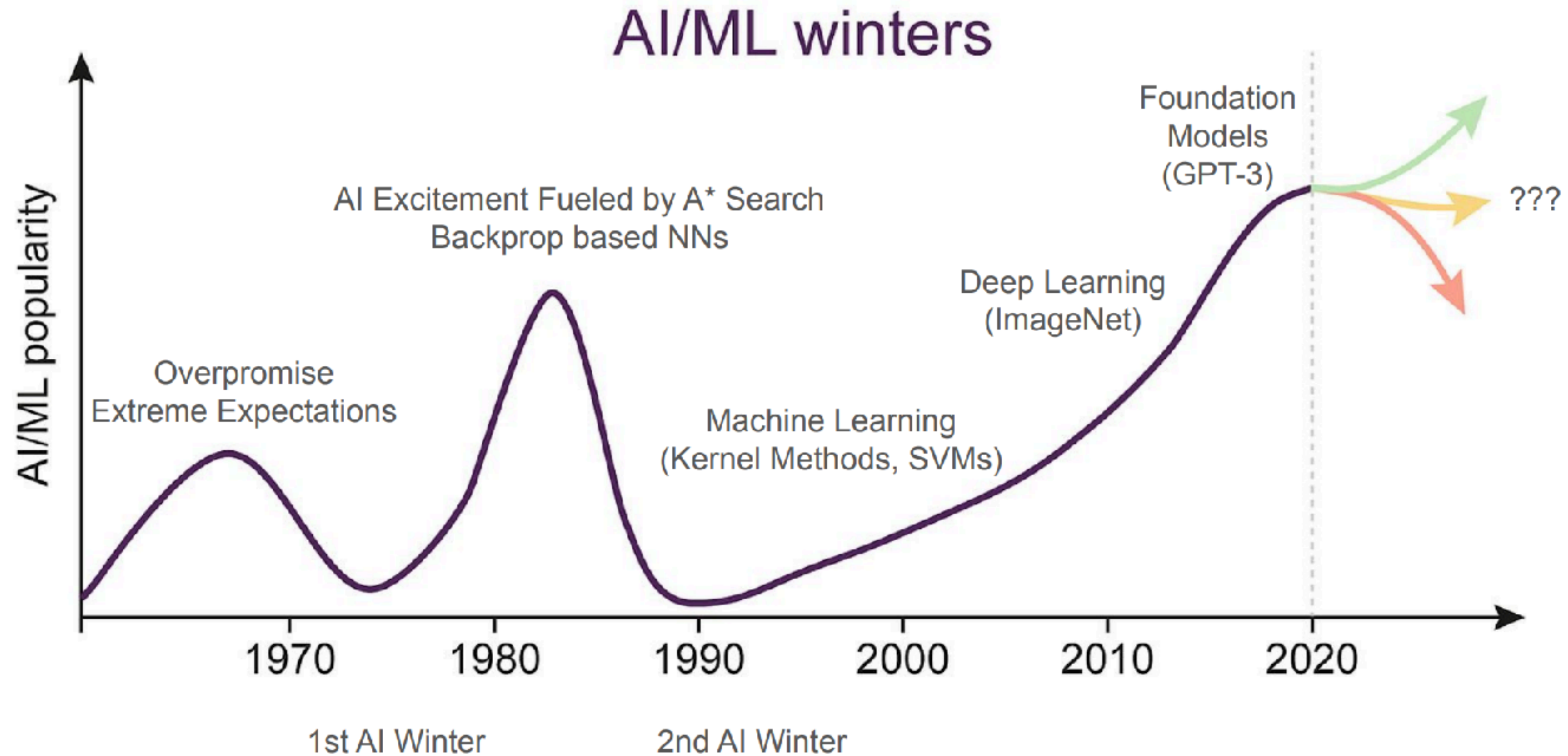
Mouses's brain

Perceptron (1957), Multi-Layer Perceptron (1965)



Rosenblatt's Multi-Layer Perceptron

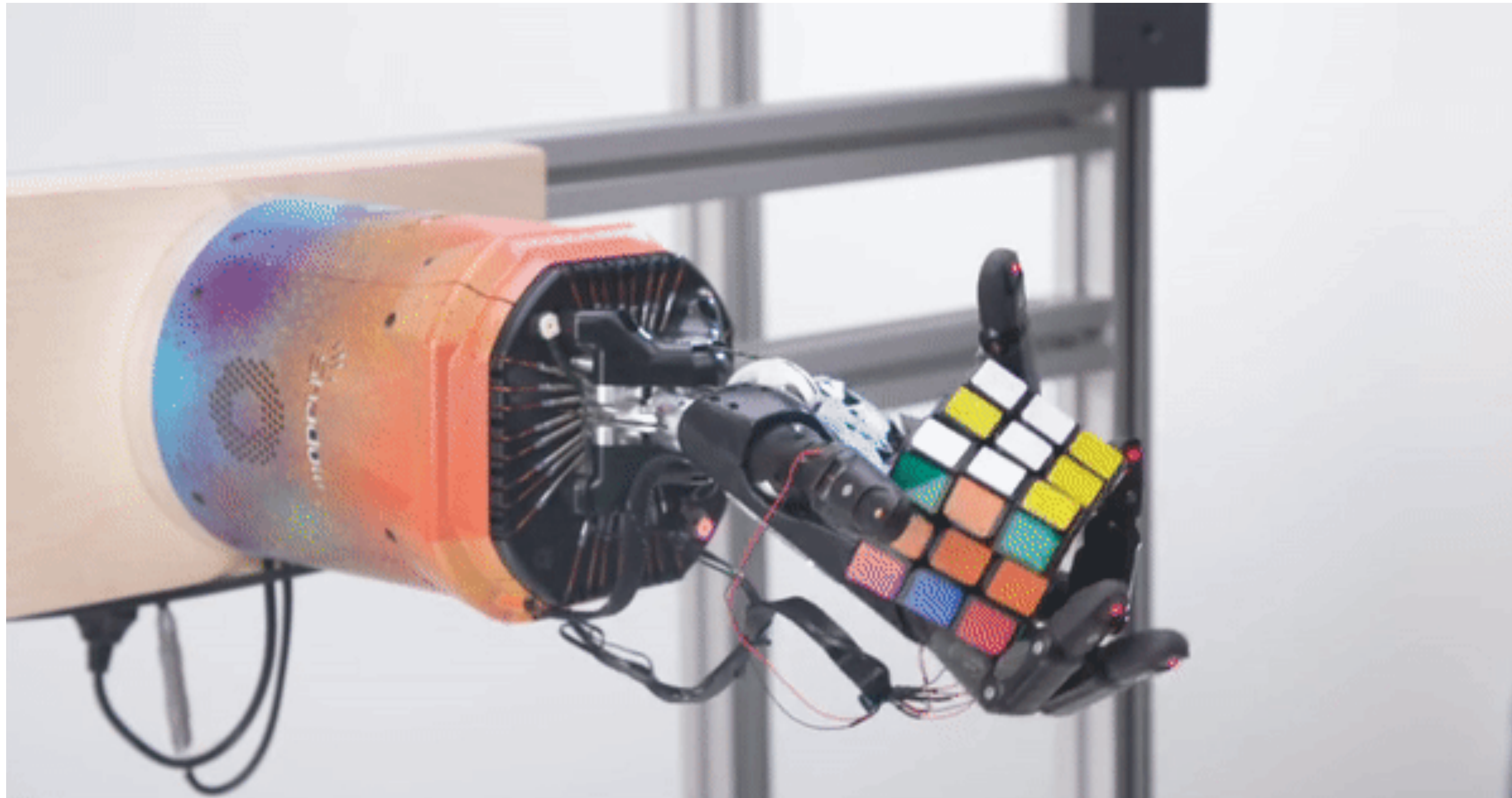
A Brief History of Machine Learning



AlphaGo Beats Lee Se-dol (2016)



Robot Solves Rubik's Cube from Scratch (2019)



DALL-E Image Generation

Prompt

“Tiny potato kings wearing majestic crowns, sitting on thrones, overseeing their vast potato kingdom filled with potato subjects and potato castles.



ChatGPT & Other Large Language Models

how to do inline citations in latex with names like:
Goldwasser, Micali, et. al. (2020)



To format inline citations in LaTeX with author names like "Goldwasser, Micali, et al. (2020)", you should use the `natbib` package with an appropriate bibliography style.

Steps:

1. Include the `natbib` package in your preamble:

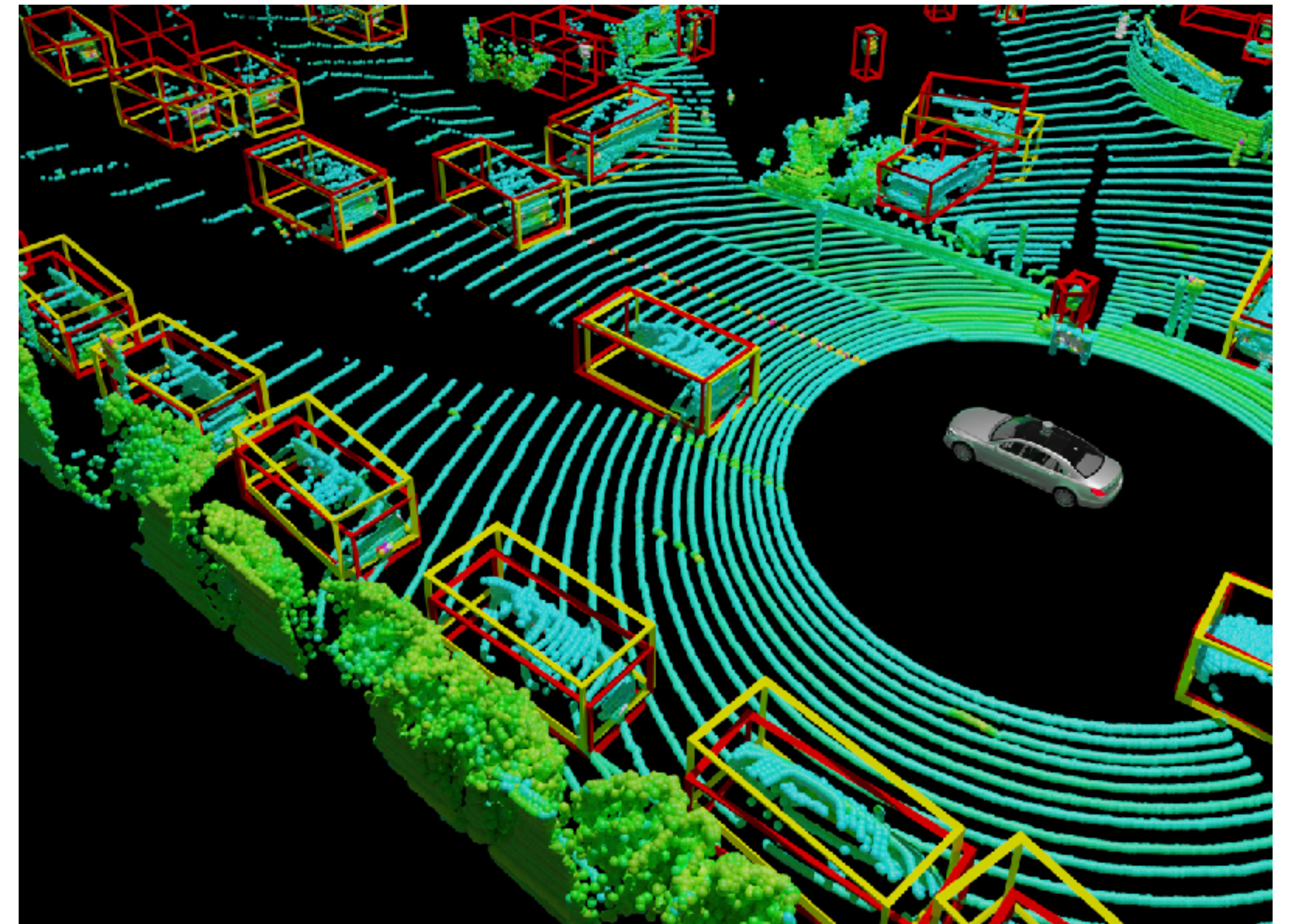
```
latex                                                                    Copy Edit
\usepackage[natbibapa]{apacite} % Alternative: \usepackage{natbib}
\bibliographystyle{apacite}      % Alternative: \bibliographystyle{...}
```

2. Use the citation command:



Use in Critical Systems

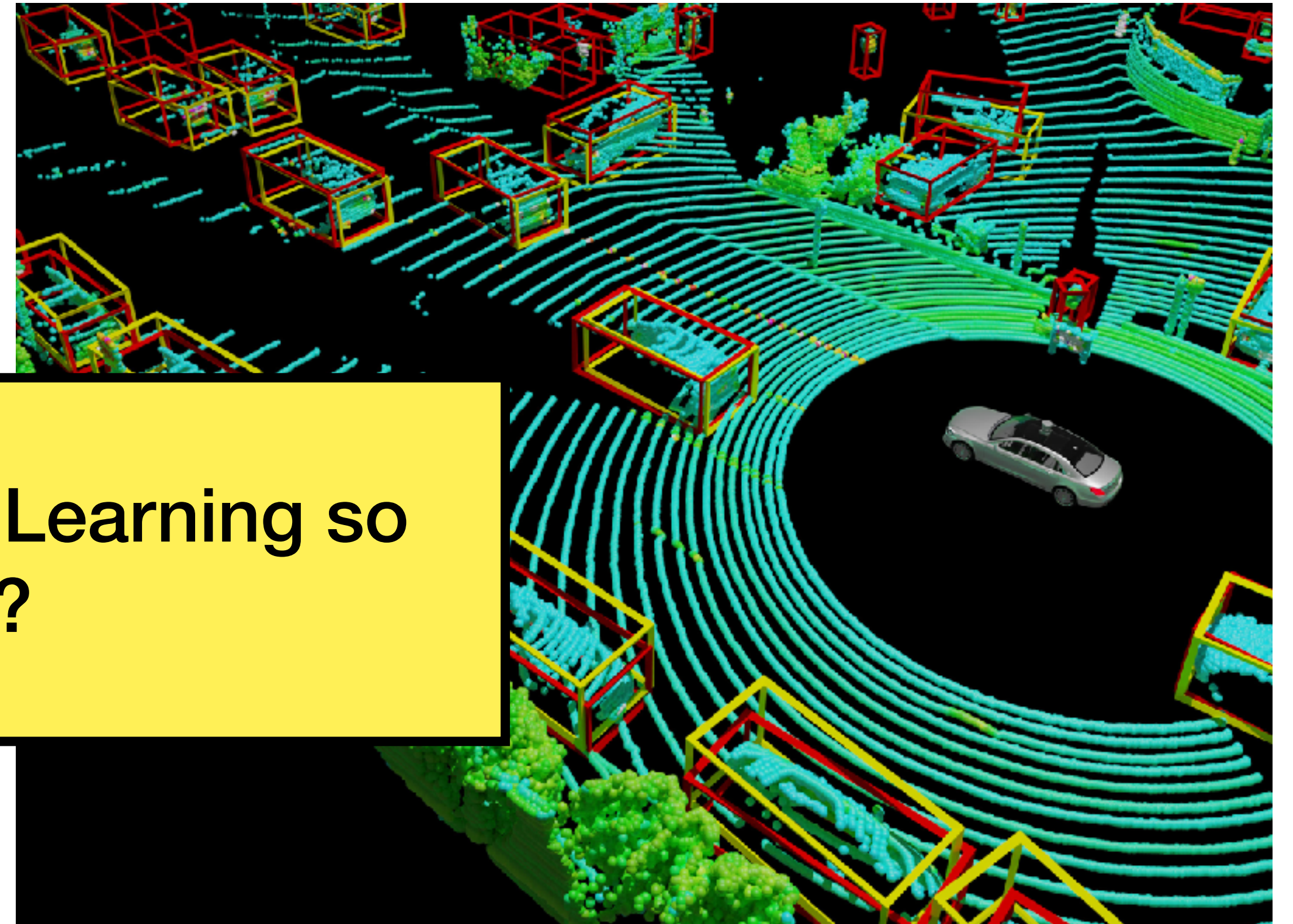
- Wildfire & earthquake prediction
- Self-driving cars



Use in Critical Systems

- Wildfire & earthquake prediction
- Self-driving cars

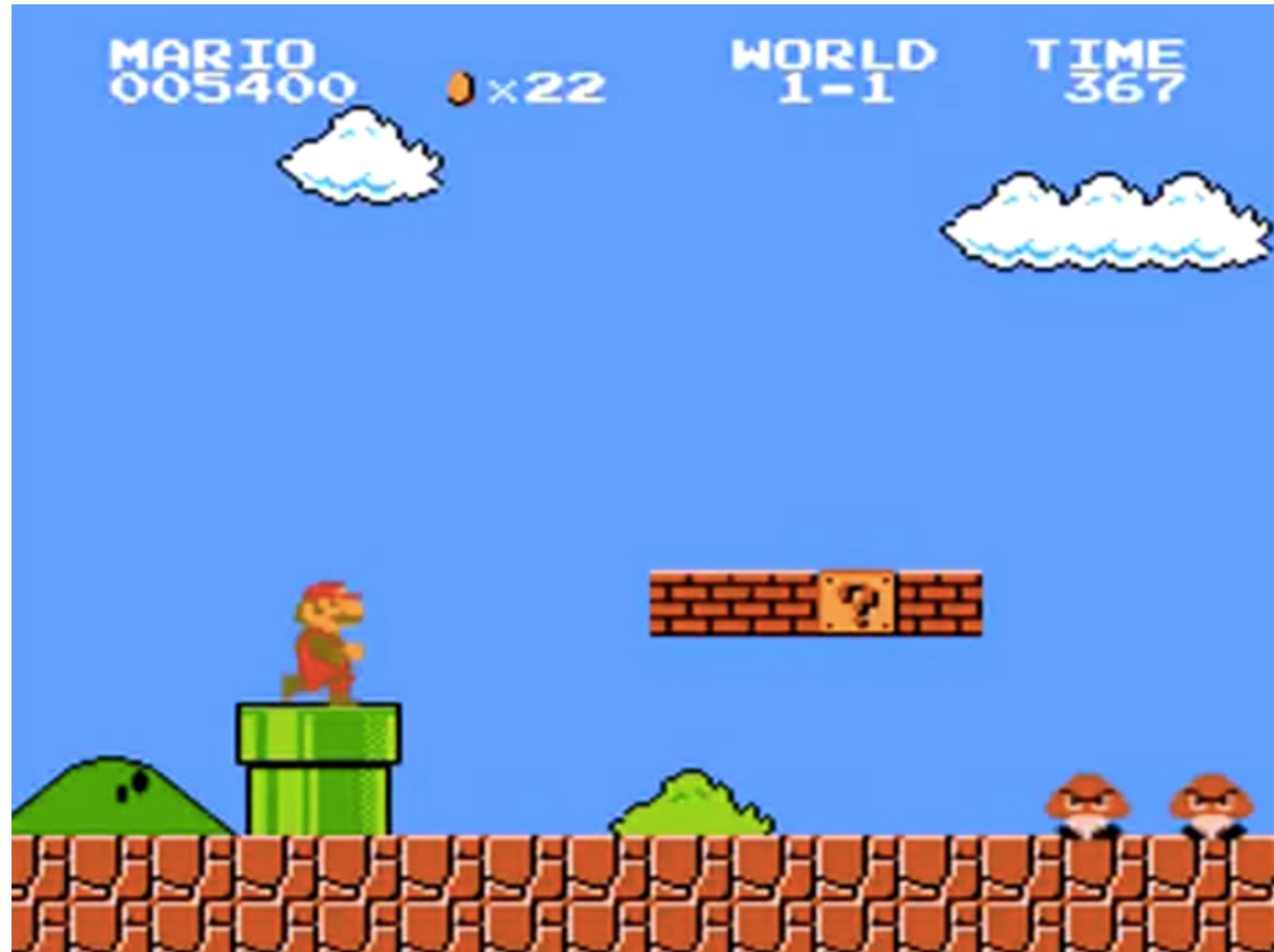
Why is Machine Learning so good?



How to reason about machine learning

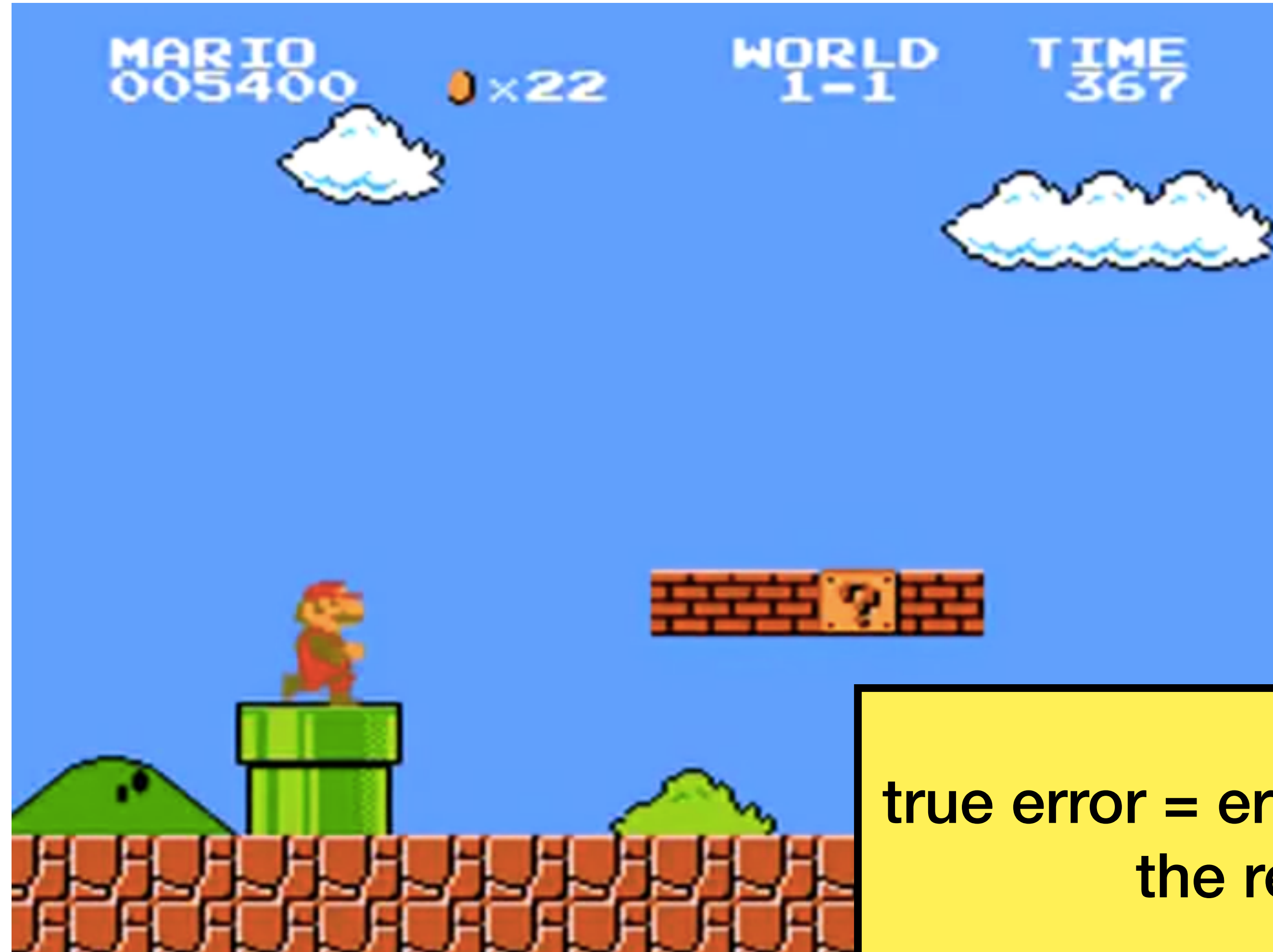
Super Mario Bros

How to
measure
success?



Super Mario Bros

How to
measure
success?



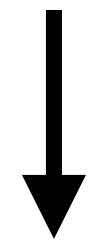
true error = error of the model in
the real world!

Bias-variance decomposition

$$\text{true error} = \text{bias} + \text{variance}$$

Bias-variance decomposition

Error of model in
the real world



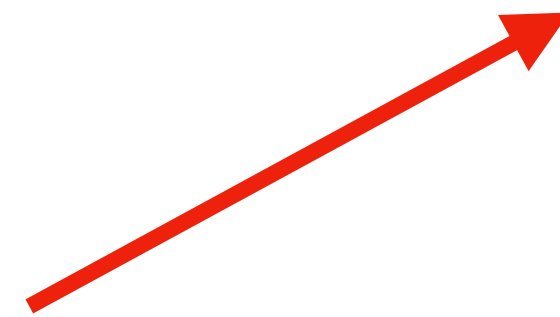
$$\text{true error} = \text{bias} + \text{variance}$$

Bias-variance decomposition

Error of model in
the real world



$$\text{true error} = \text{bias} + \text{variance}$$



Absolute best our
ML algorithm can
do with unlimited
data and
computation

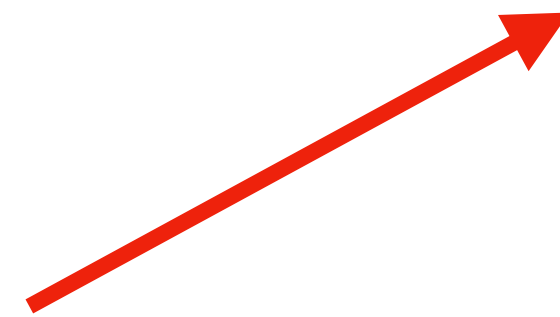
Bias-variance decomposition

Error of model in
the real world

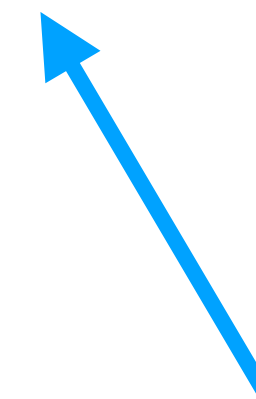


$$\text{true error} = \text{bias} + \text{variance}$$

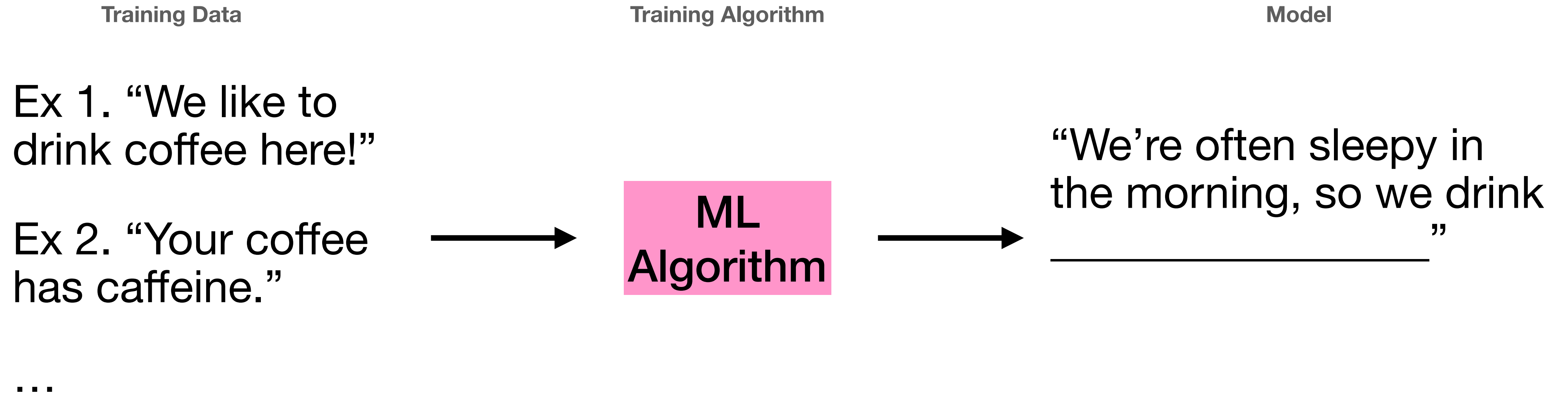
Absolute best our
ML algorithm can
do with unlimited
data and
computation



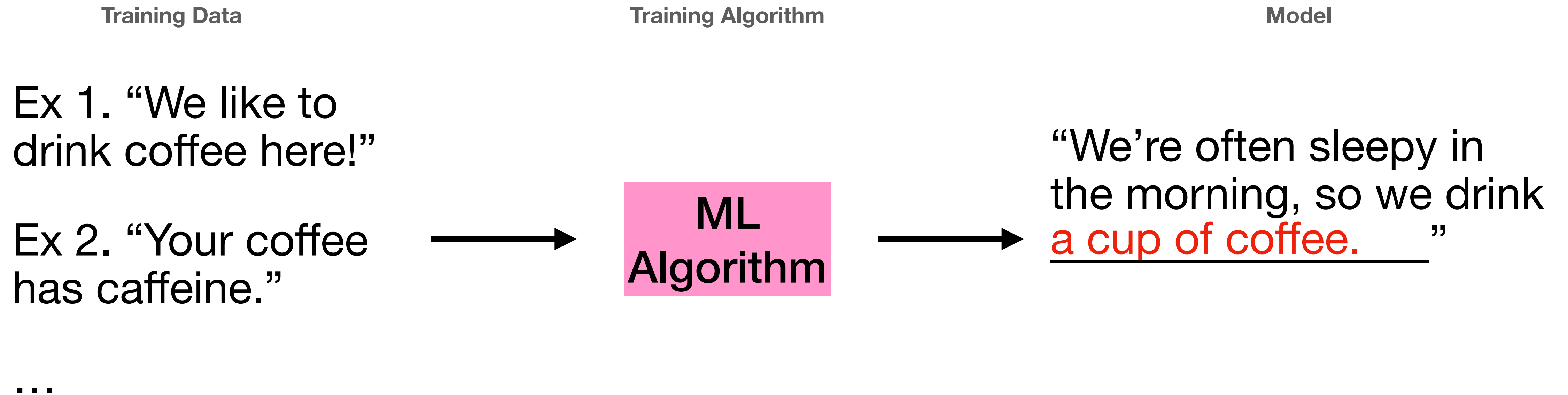
Variance from optimal
model. Approximation
error from not having
enough data or compute
power.



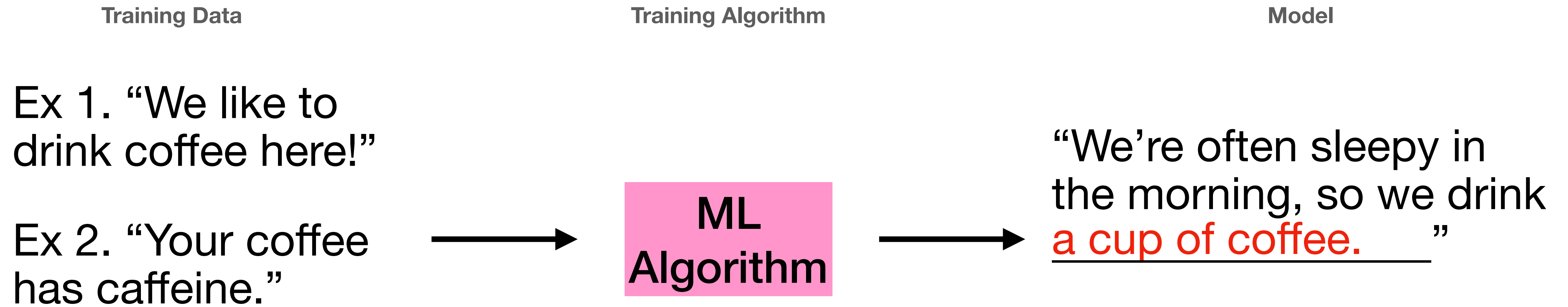
Learning to finish sentences



Learning to finish sentences



Learning to finish sentences



...

This is how ChatGPT works! It learns how to finish the user's prompt.

Machine learning on bi-grams

Ex 1. “We like to
drink coffee here!” → (“We”, “like”), (“like”, “to”), (“to”, “drink”),
 (“drink”, “coffee”), (“coffee”, “here”), (“here”, “!”)

Ex 2. “Your coffee
has caffeine.” → (“Your”, “coffee”), (“coffee”, “has”), (“has”,
 “caffeine”), (“caffeine”, “.”)

Model

converts each sentence into a bi-gram and use each
bi-gram to predict the next word.

“We’re tired in the morning, so we

”

Machine learning on bi-grams

Ex 1. “We like to
drink coffee here!” → (“We”, “like”), (“like”, “to”), (“to”, “drink”),
 (“drink”, “coffee”), (“coffee”, “here”), (“here”, “!”)

Ex 2. “Your coffee
has caffeine.” → (“Your”, “coffee”), (“coffee”, “has”), (“has”,
 “caffeine”), (“caffeine”, “.”)

Model

converts each sentence into a bi-gram and use each
bi-gram to predict the next word.

“We’re tired in the morning, so we like to drink coffee has caffeine.”

Machine learning on bi-grams

Ex 1. “We like to drink coffee here!” → (“We”, “like”), (“like”, “to”), (“to”, “drink”), (“drink”, “coffee”), (“coffee”, “here”), (“here”, “!”)

Ex 2. “Your coffee has caffeine.” → (“Your”, “coffee”), (“coffee”, “has”), (“has”, “caffeine”), (“caffeine”, “.”)

Model

converts each sentence into a bi-gram and use each bi-gram to predict the next word.

“We’re tired in the morning, so we like to drink coffee has caffeine.”



Bi-grams have high bias

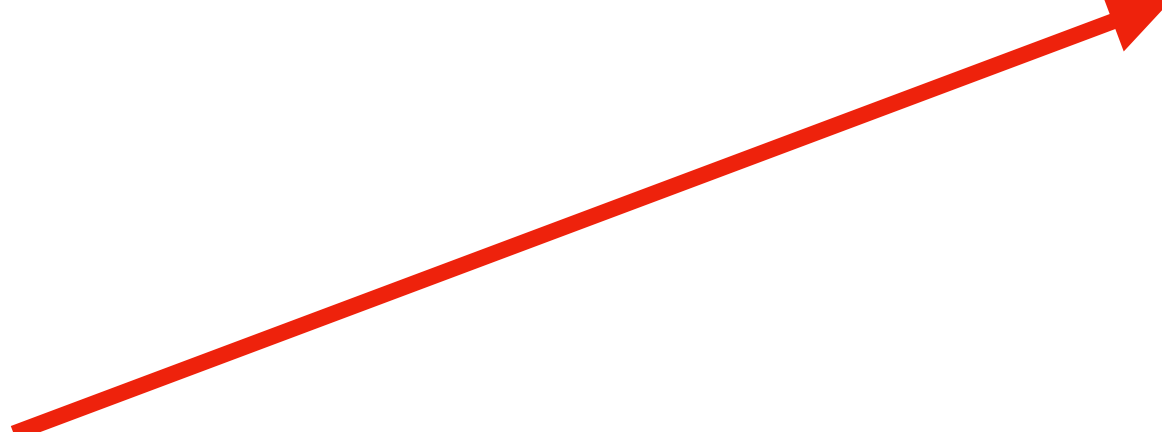
“We’re tired in the morning so... we drink coffee we work hard during the morning so we drink coffee we are always tired in the evening we work we're happy.”

$$\text{true error} = \text{bias} + \text{variance}$$

Bi-grams have high bias

“We’re tired in the morning so... we drink coffee we work hard during the morning so we drink coffee we are always tired in the evening we work we're happy.”

$$\text{true error} = \text{bias} + \text{variance}$$

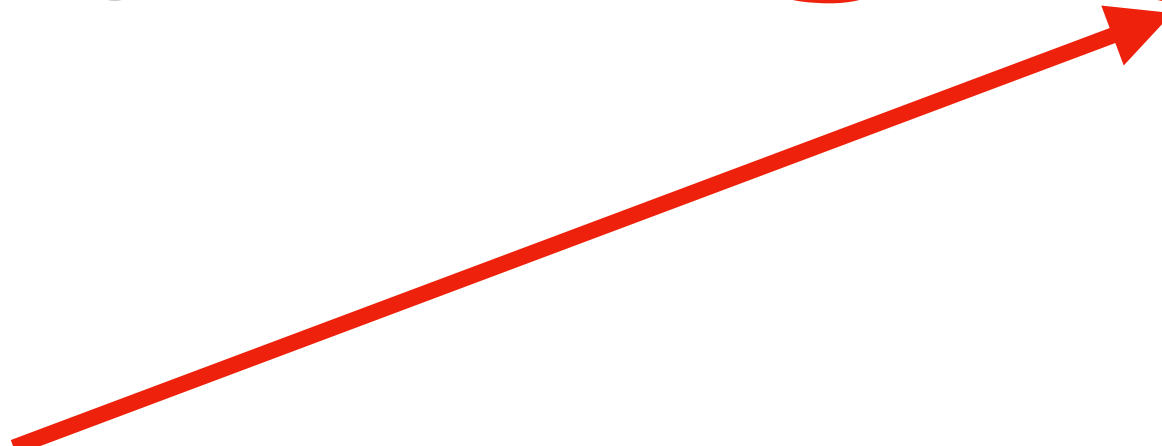


Finishing a sentence one word at a time will almost never make sense.


Bi-grams have high bias

“We’re tired in the morning so... we drink coffee we work hard during the morning so we drink coffee we are always tired in the evening we work we're happy.”

$$\text{true error} = \text{bias} + \text{variance}$$



Finishing a sentence one word at a time will almost never make sense.



Over enough sentences, the actually most frequent of bigrams will always be used. So there's not much variance from the optimal.

Machine learning by memorization

Ex 1. “We like to drink coffee here!”

Ex 2. “Your coffee has caffeine.”

Model

Checks if the prompt matches any of the sentences in the training set, and finishes accordingly.

“We like to drink....”

“We’re tired in the morning, so we”

Machine learning by memorization

Ex 1. “We like to drink coffee here!”

Ex 2. “Your coffee has caffeine.”

Model

Checks if the prompt matches any of the sentences in the training set, and finishes accordingly.

“We like to drink.... **coffee here!**”

“We’re tired in the morning, so we”

Machine learning by memorization

Ex 1. “We like to drink coffee here!”

Ex 2. “Your coffee has caffeine.”

Model

Checks if the prompt matches any of the sentences in the training set, and finishes accordingly.

“We like to drink.... **coffee here!**”

“We’re tired in the morning, so we **[NO OUTPUT]**”

Memorization has high variance

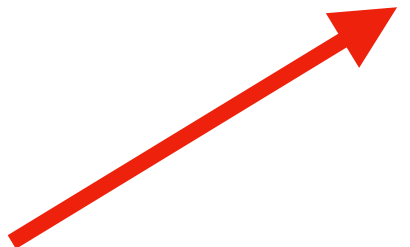
“Four score and seven years ago... our fathers brought forth on this continent, a new nation, conceived in Liberty, and dedicated to the proposition that all men are created equal.”

$$\text{true error} = \text{bias} + \text{variance}$$

Memorization has high variance

“Four score and seven years ago... our fathers brought forth on this continent, a new nation, conceived in Liberty, and dedicated to the proposition that all men are created equal.”

$$\text{true error} = \text{bias} + \text{variance}$$

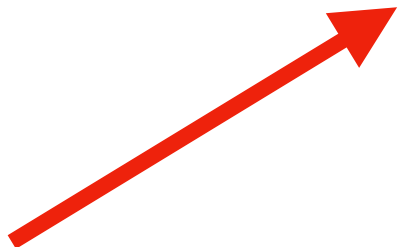


With access to every possible conceivable sentence and the memory to store it, it will always finish correctly!


Memorization has high variance

“Four score and seven years ago... our fathers brought forth on this continent, a new nation, conceived in Liberty, and dedicated to the proposition that all men are created equal.”

$$\text{true error} = \text{bias} + \text{variance}$$

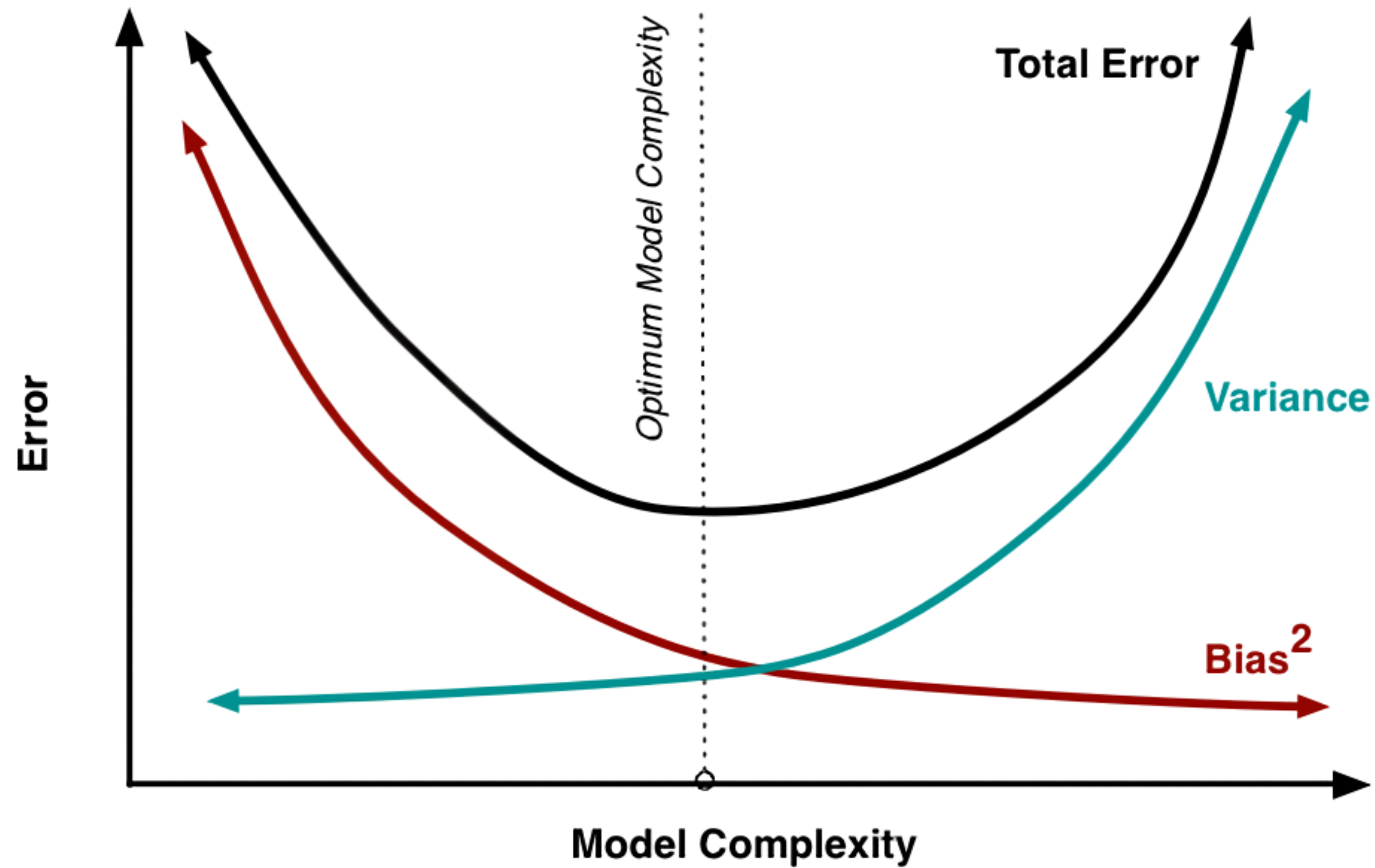


With access to every possible conceivable sentence and the memory to store it, it will always finish correctly!

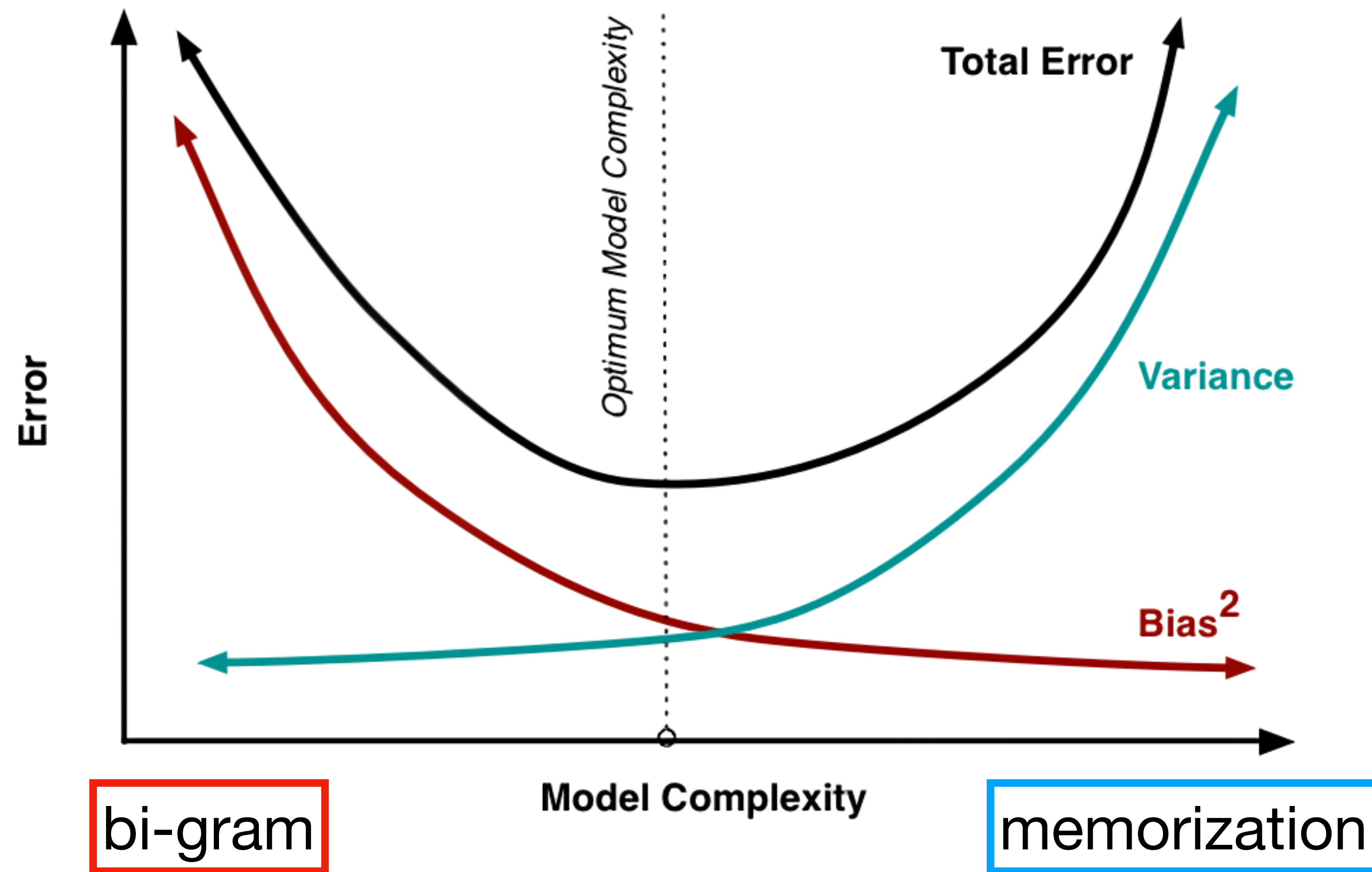


Can only finish sentences in dataset. With limited memory, it will not even be close to finishing every possible sentence.

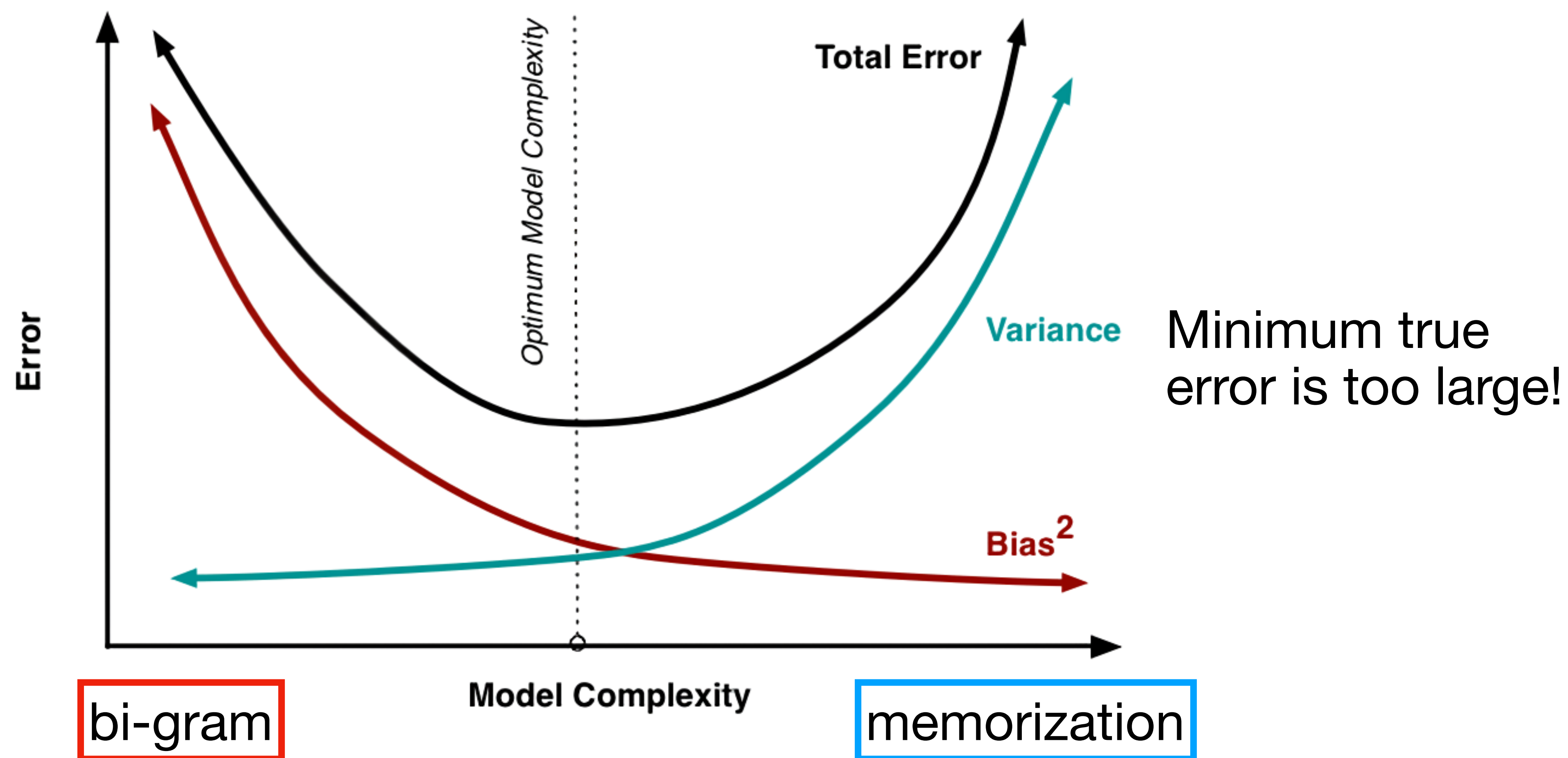
Bias-variance tradeoff



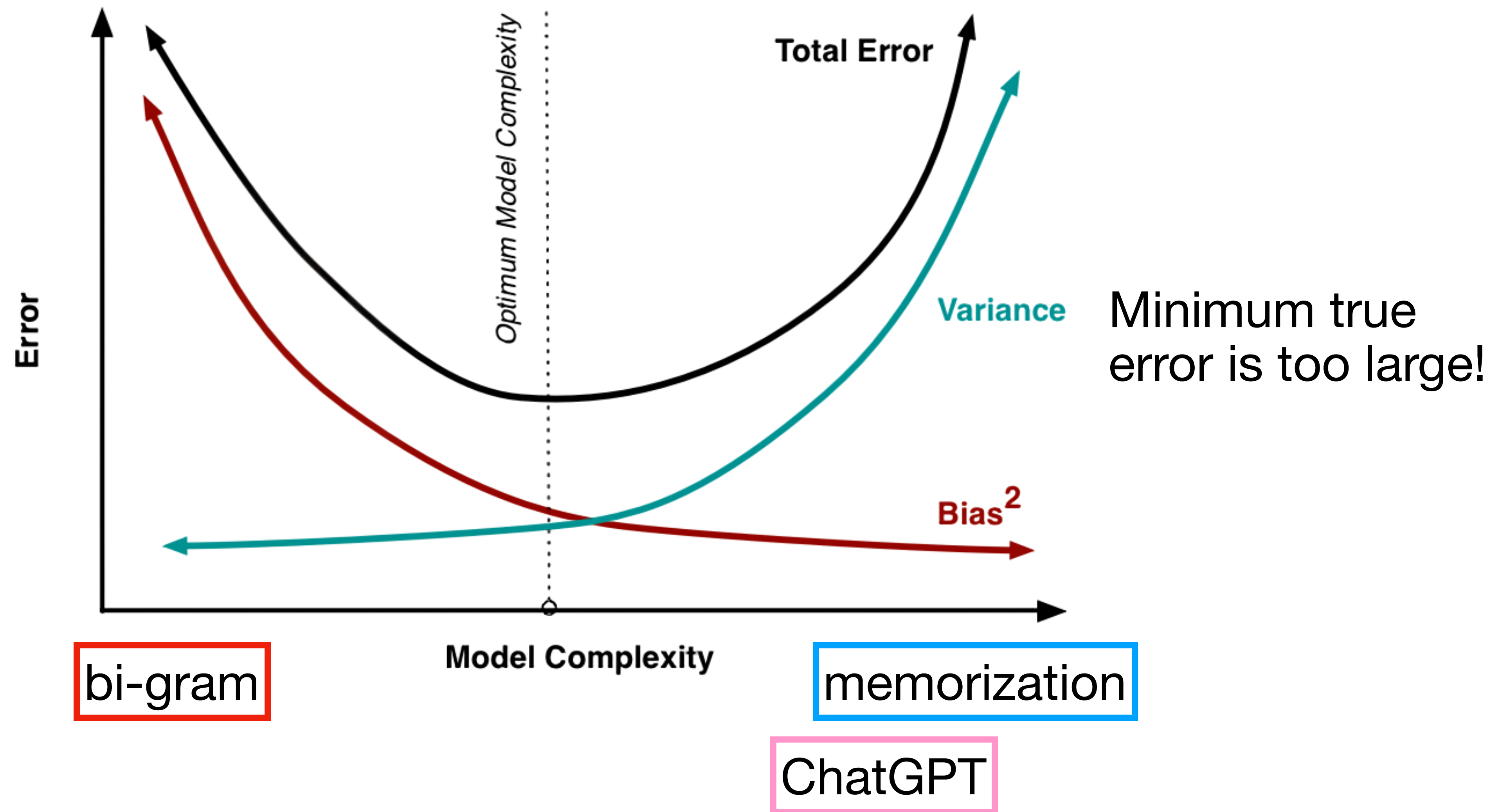
Bias-variance tradeoff



Bias-variance tradeoff

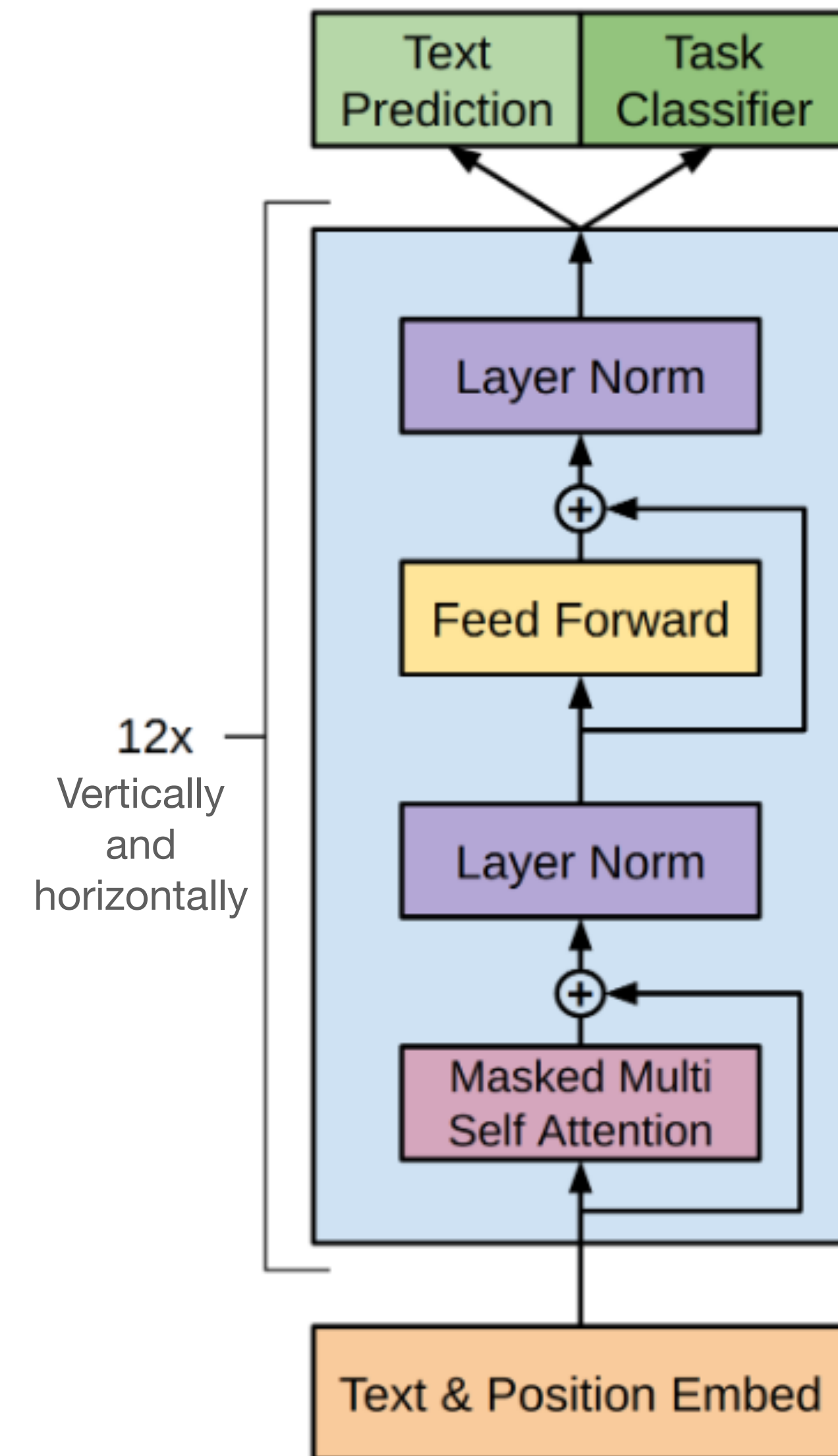


Bias-variance tradeoff



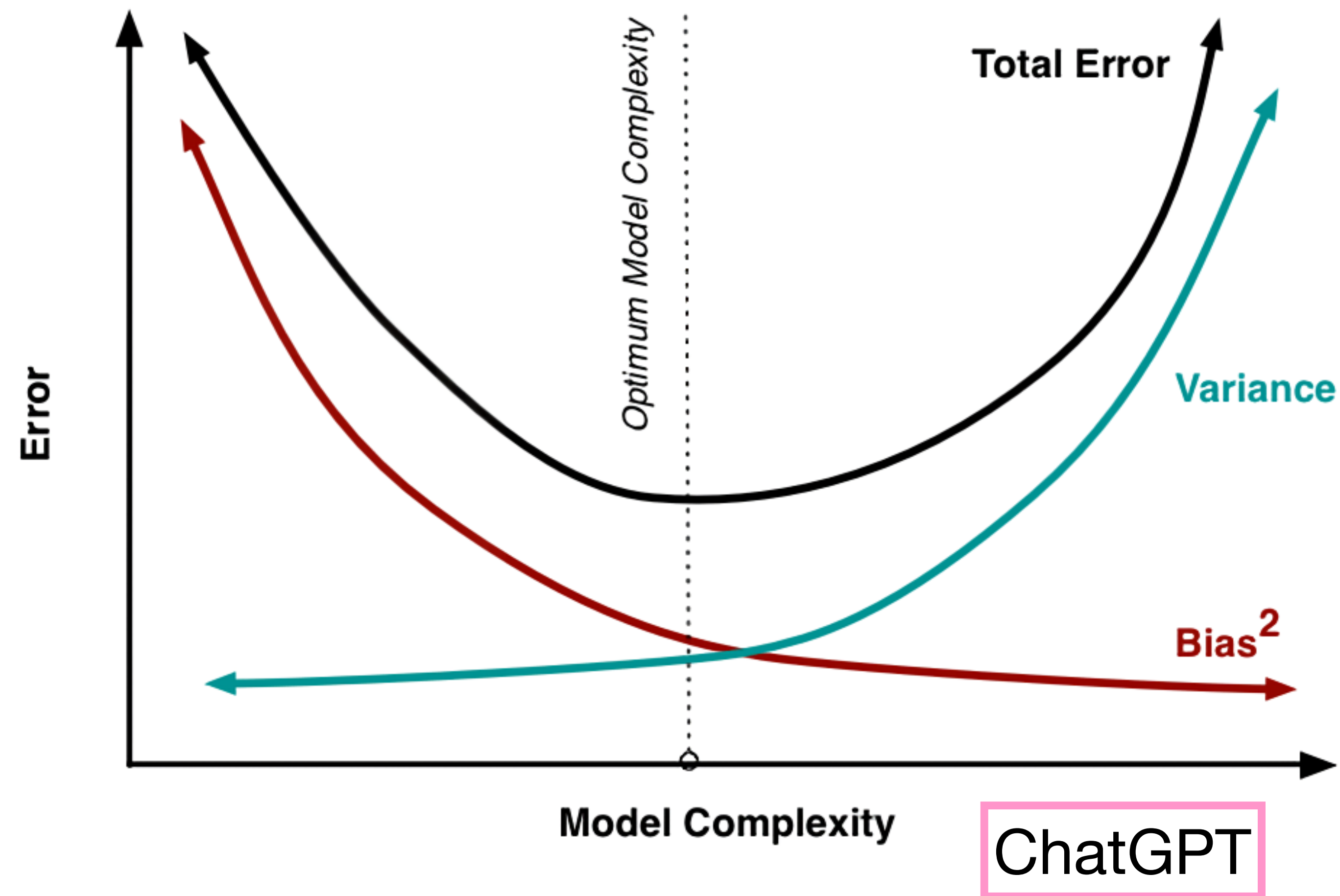
How ChatGPT cheats

- GPT-4 has **~1 trillion** parameters
- Push variance to the right via
 - Train on the a massive dataset (the internet)
 - Use a transformer-based architectures which allows for really good parallelization with GPUs.
 - Spend > **\$100 million**

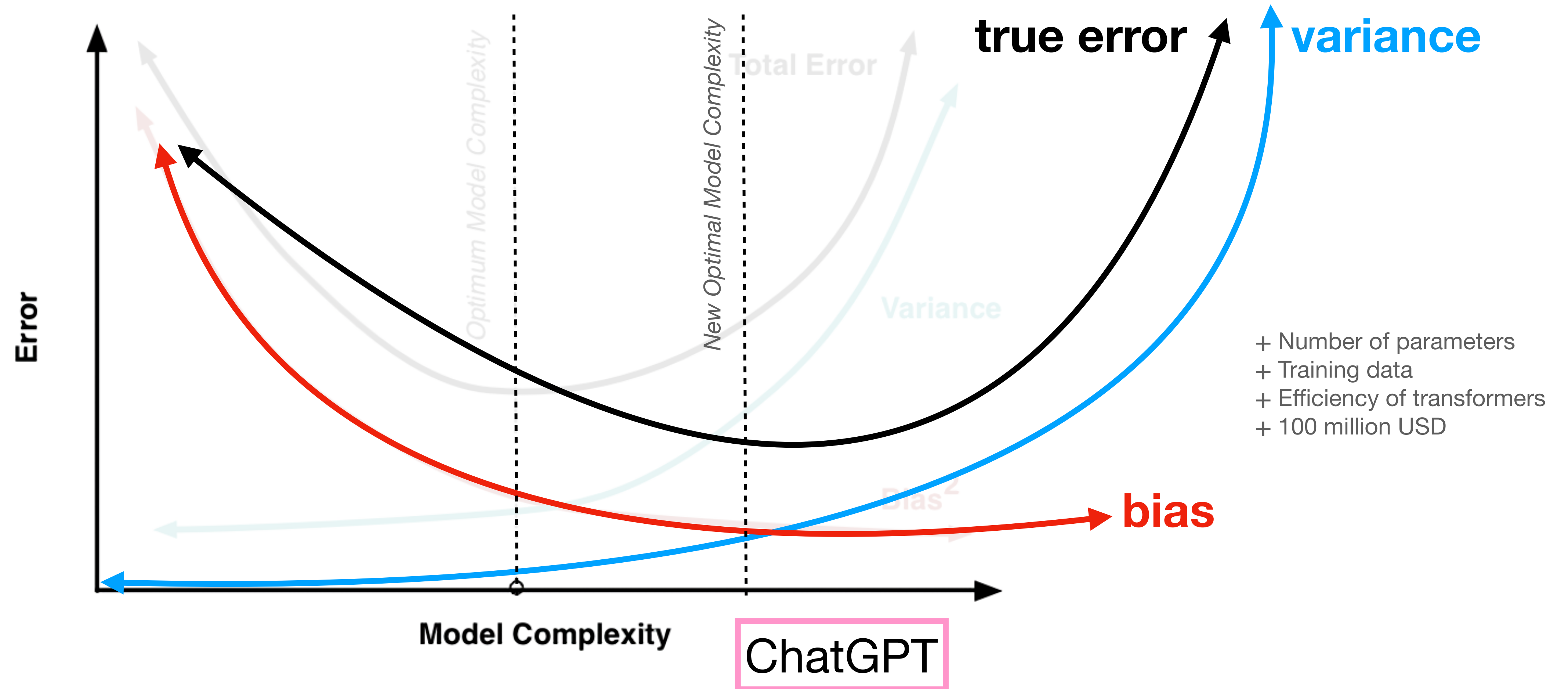


GPT-3 architecture

Bias-variance tradeoff for GPT



Bias-variance tradeoff for GPT



Scaling laws

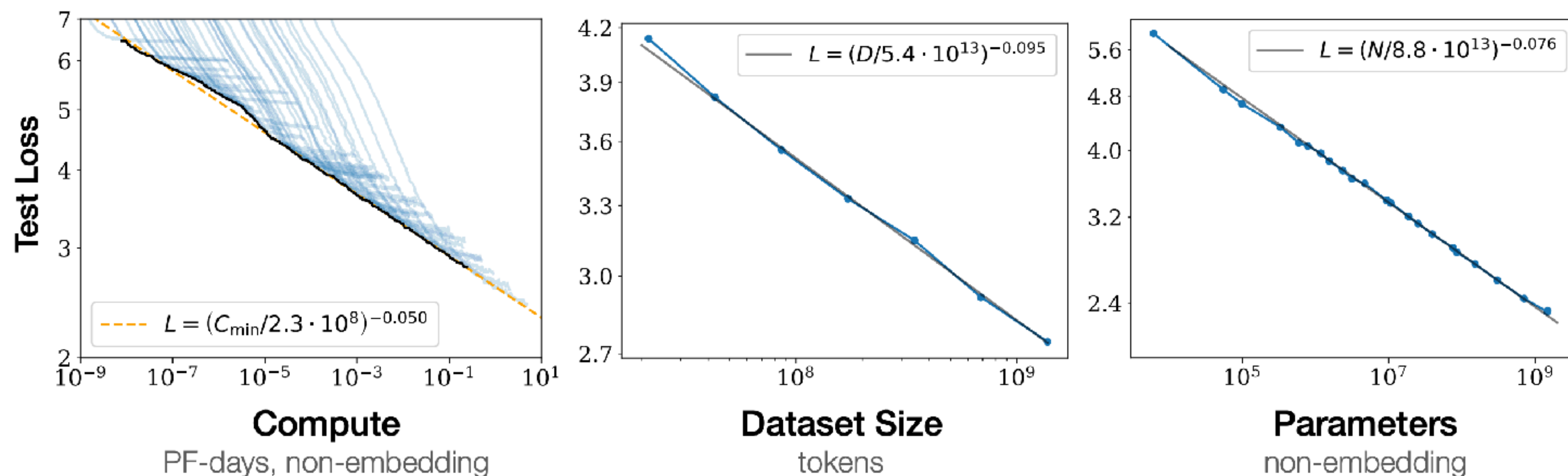
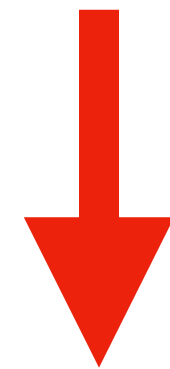


Figure 1 Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute² used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

ML is powerful!

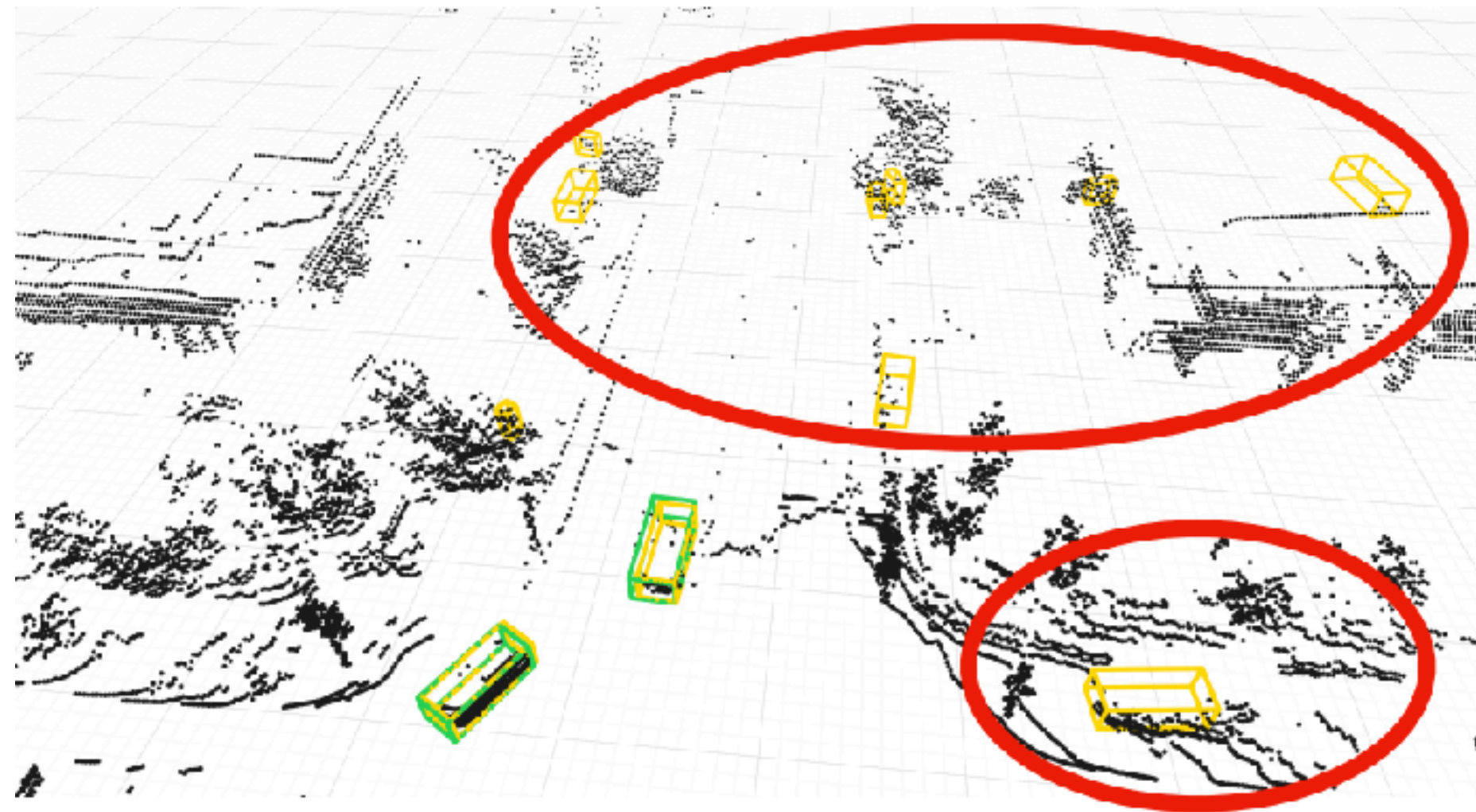
ML is powerful!



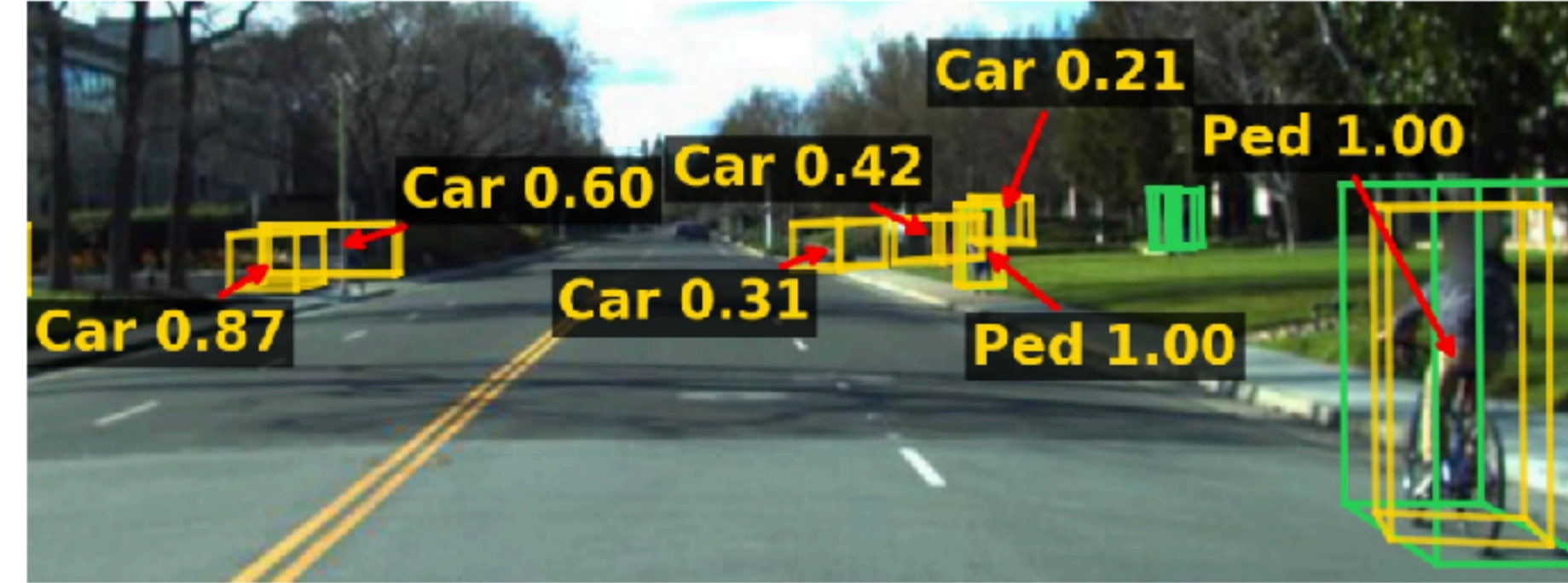
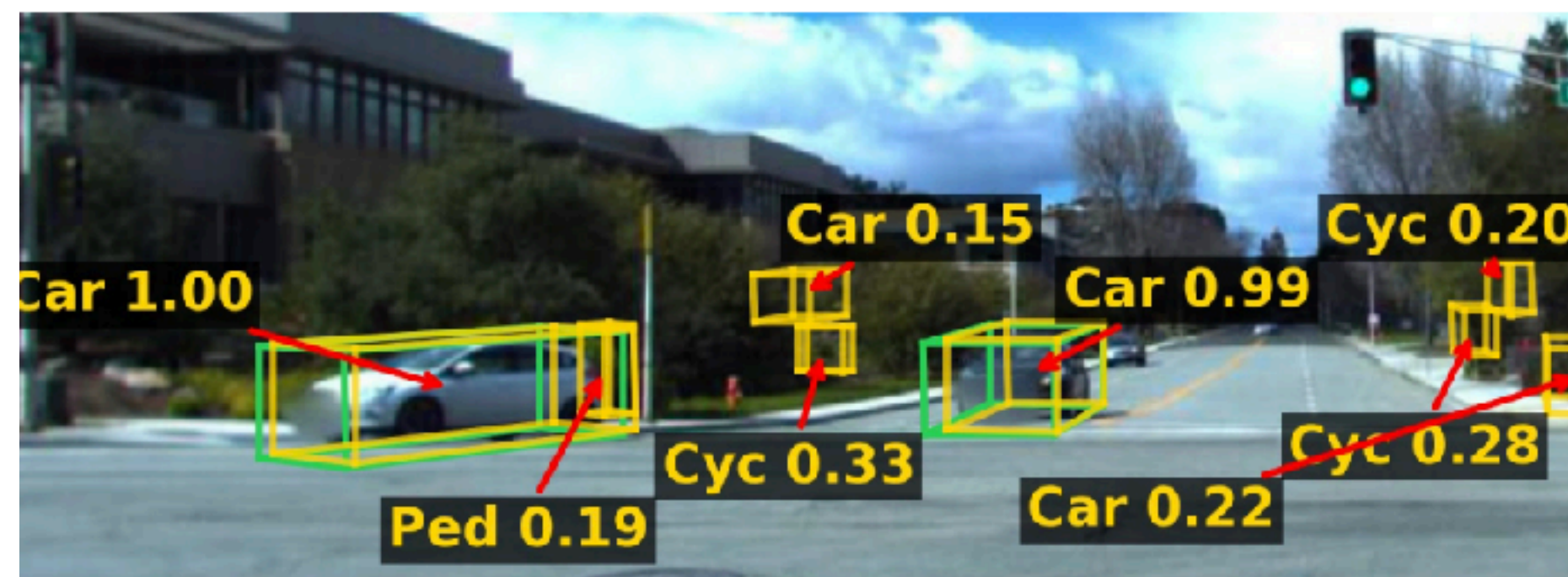
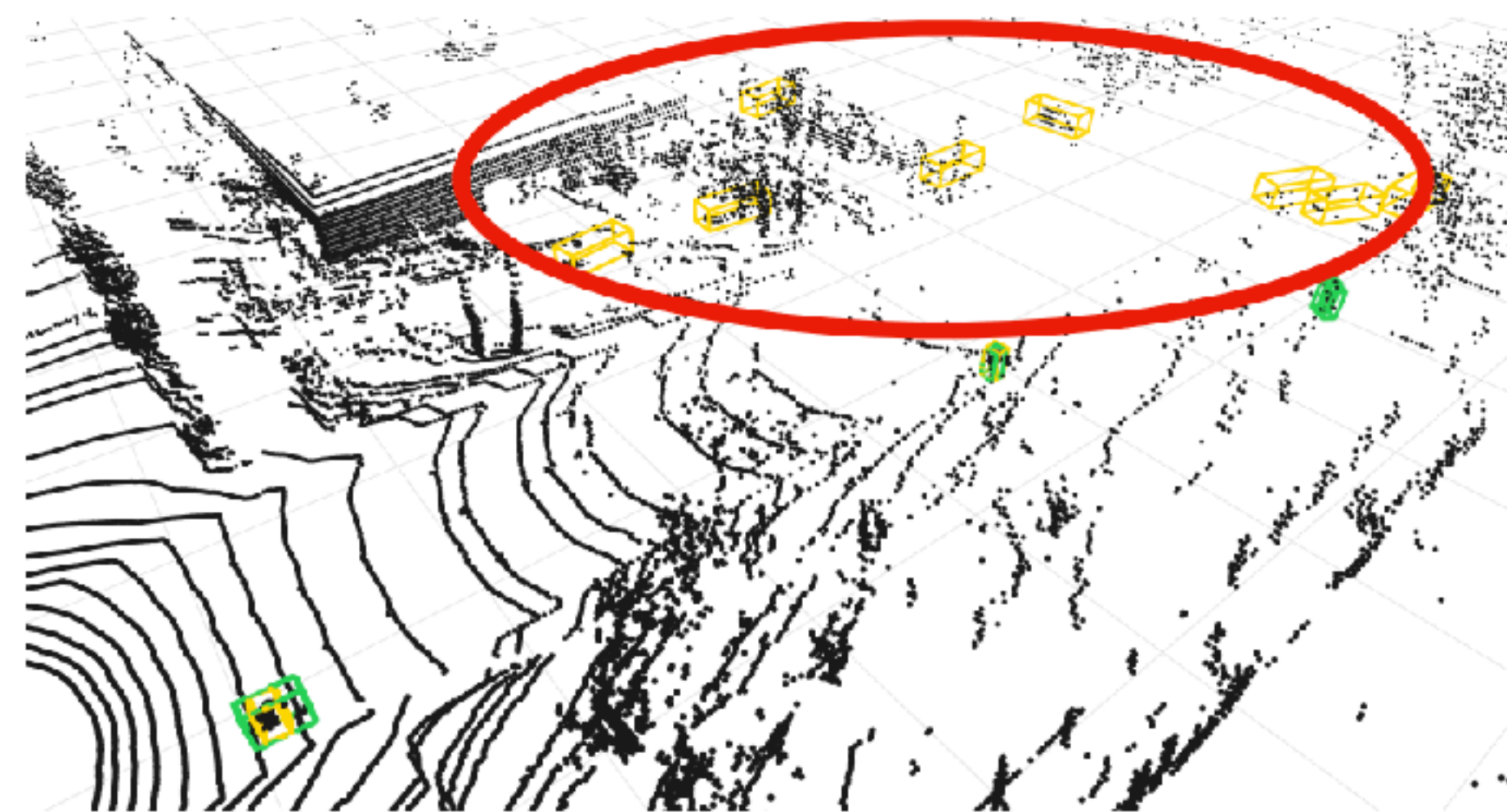
ML is dangerous!

ML can make mistakes

Example Scene 1



Example Scene 2



ML aren't smart



ML algorithm thinks: *brake light on* => *brake*

ML aren't smart



ML algorithm thinks: *brake light on* => *brake*

Correlation \neq Causation
This is called **overfitting**.

ML aren't smart



ML algorithm thinks: *brake light on* => *brake*

Correlation \neq Causation

This is called **overfitting**.

Try talking to ChatGPT! Does it feel like a person?

Even GPT-4o struggles at reasoning

Solid choice — you placed your **X** in position **8**. Here's the updated board:

markdownCopy code

```
0 | X | 0
-----
0 | X | 6
-----
7 | X | X
```

I'll go with **6**:

markdownCopy code

```
0 | X | 0
-----
0 | X | 0
-----
7 | X | X
```

Your turn! Only positions **7** are left.



**AI is nowhere close to human intelligence.
But it doesn't have to be smarter than us to
be harmful.**

<https://www.cnn.com/2025/04/02/tech/ai-future-of-humanity-2035-report/index.html>

<https://www.ibm.com/think/news/agi-right-goal>

<https://venturebeat.com/ai/agi-is-coming-faster-than-we-think-we-must-get-ready-now/>

Biased data implies biased machine learning

The Best Algorithms Struggle to Recognize Black Faces Equally

Google's algorithm shows prestigious job ads to men, but not to women. Here's why that should worry you.

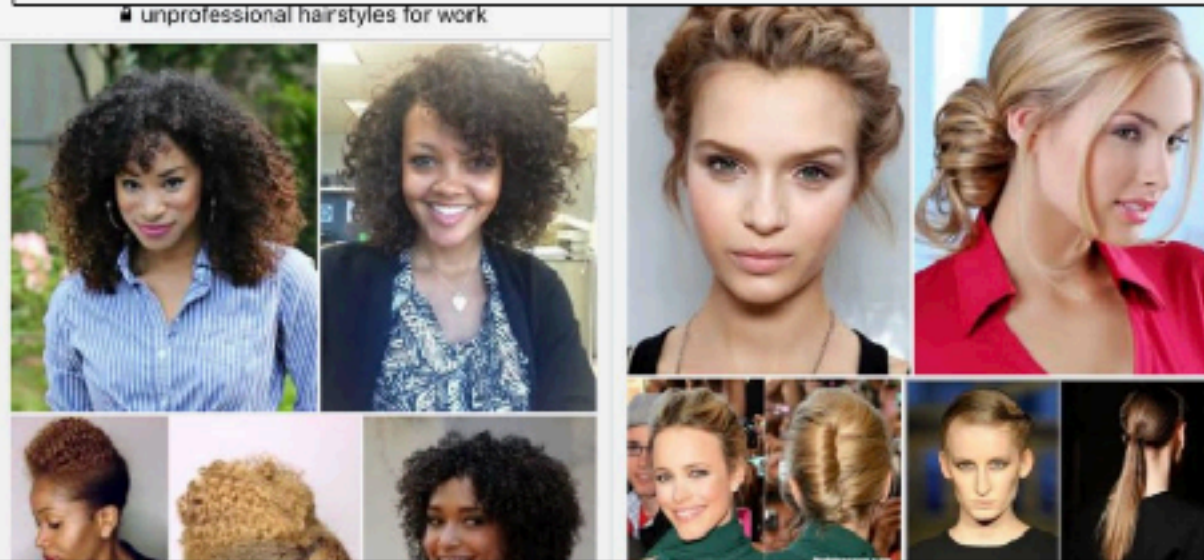
Gender and racial bias found in Amazon's facial recognition technology (again)

How Amazon Accidentally Invented a Sexist Hiring Algorithm

A company experiment to use artificial intelligence in hiring inadvertently favored male candidates.


Do Google's 'unprofessional hair' results show it is racist?

unprofessional hairstyles for work



When an Algorithm Helps Send You to Prison

By Ellora Thadaney Israni



ML can be used for bad



Can you tell the difference? Jake Tapper uses his own deepfake to show how powerful AI is

The Lead

How AI deepfakes polluted elections in 2024

DECEMBER 21, 2024 • 5:00 AM ET

HEARD ON **ALL THINGS CONSIDERED**



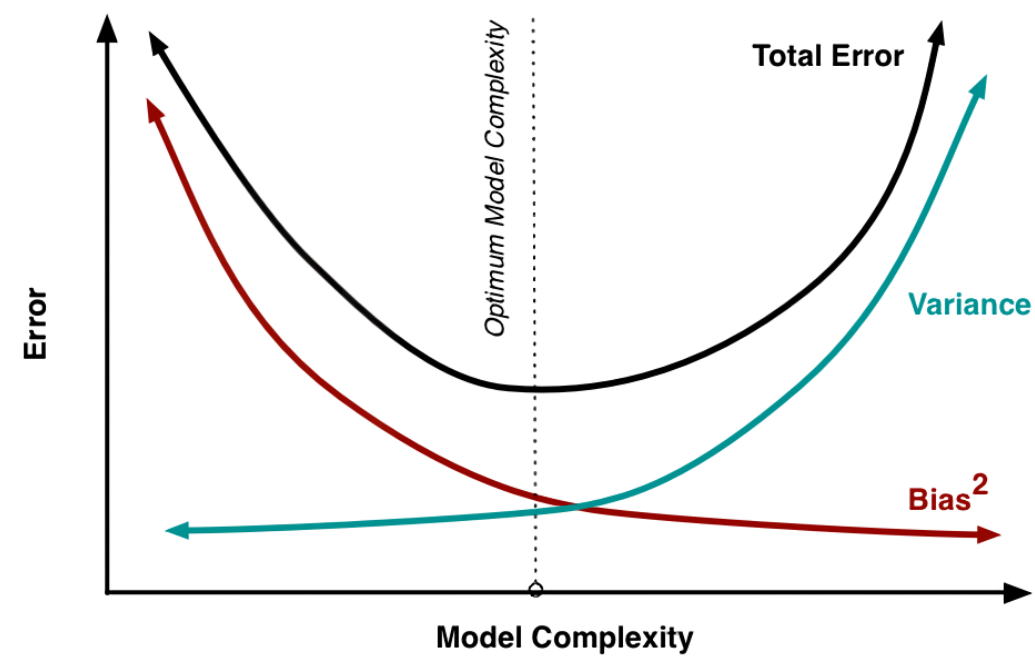
Shannon Bond

<https://www.cnn.com/2024/10/04/politics/video/ai-elections-jake-tapper-lead-digvid>

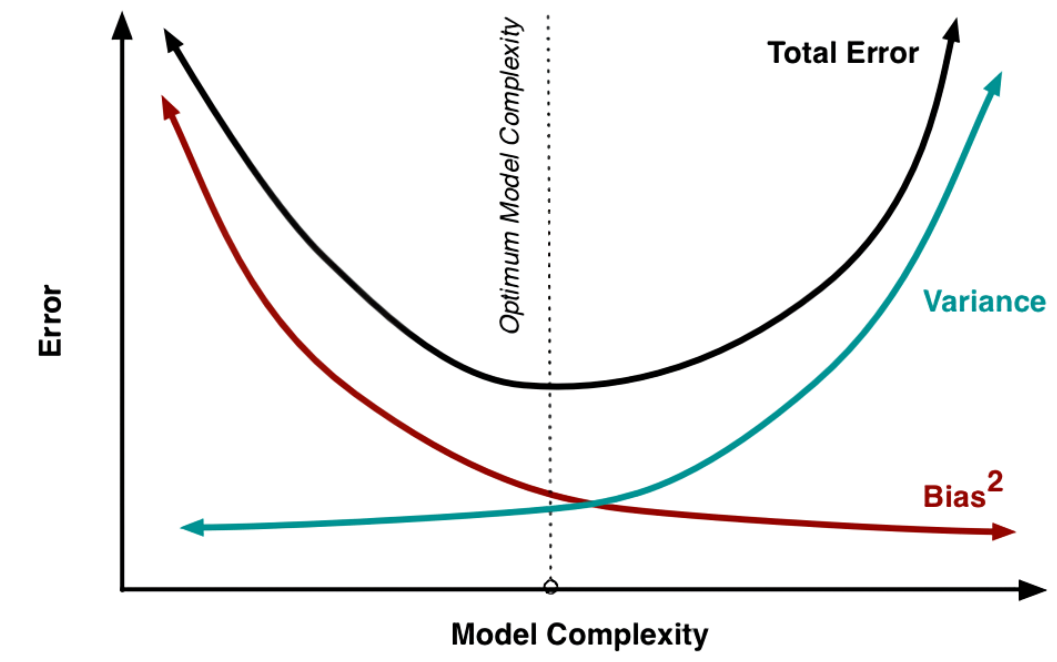
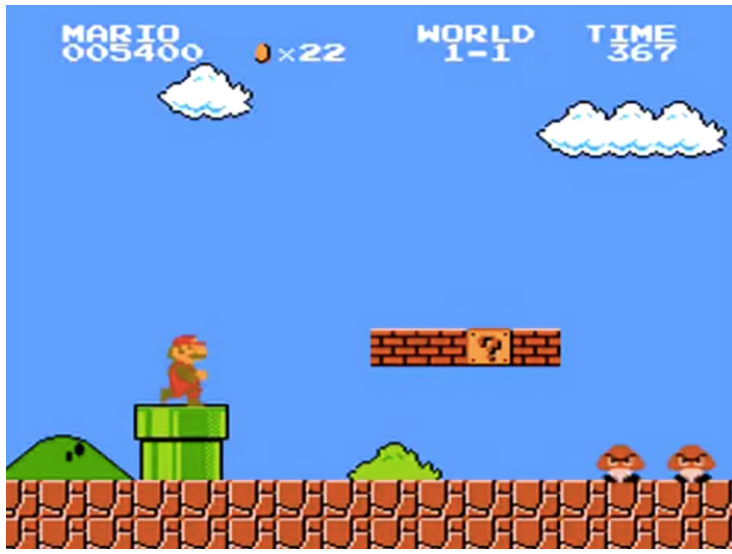
<https://www.npr.org/2024/12/21/nx-s1-5220301/deepfakes-memes-artificial-intelligence-elections>

ML-powered advertising



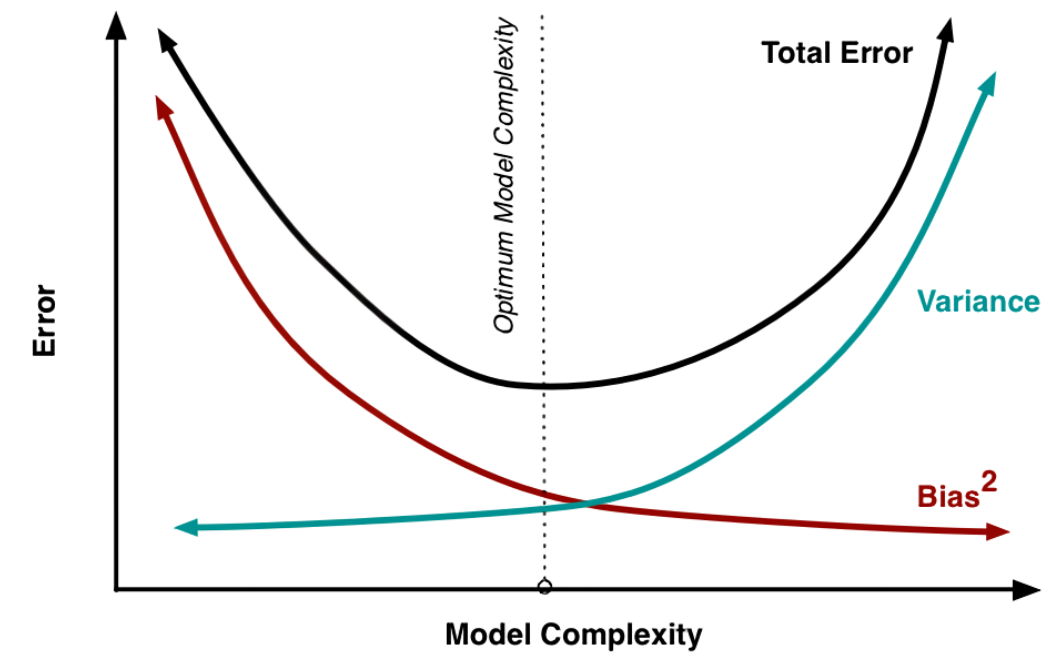
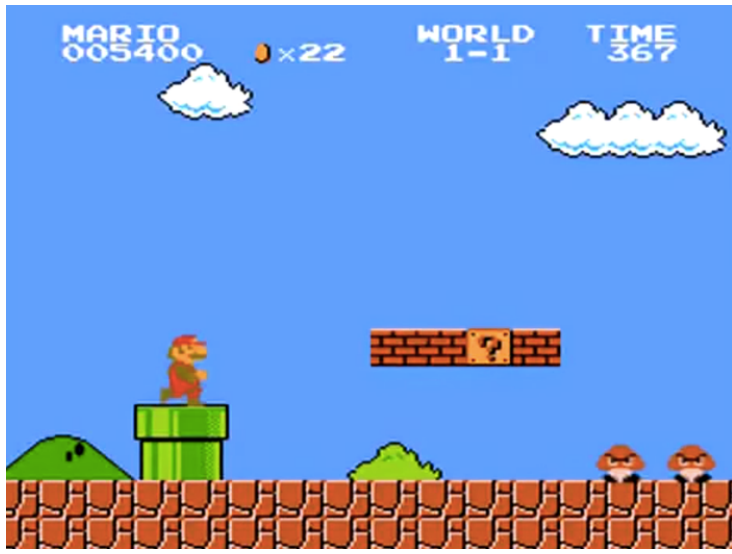


Summary



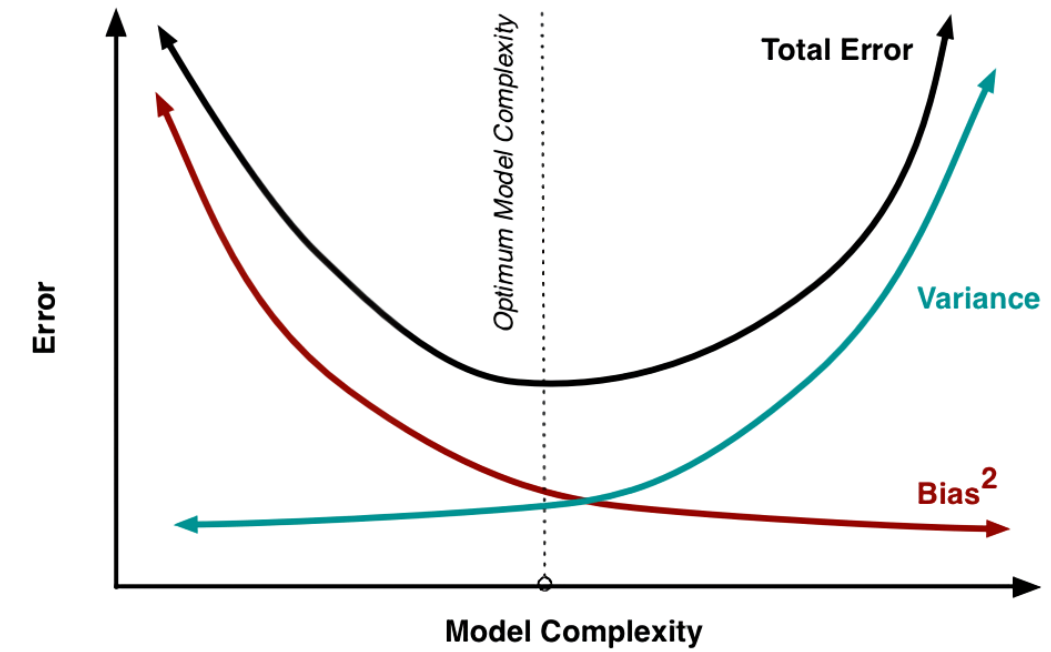
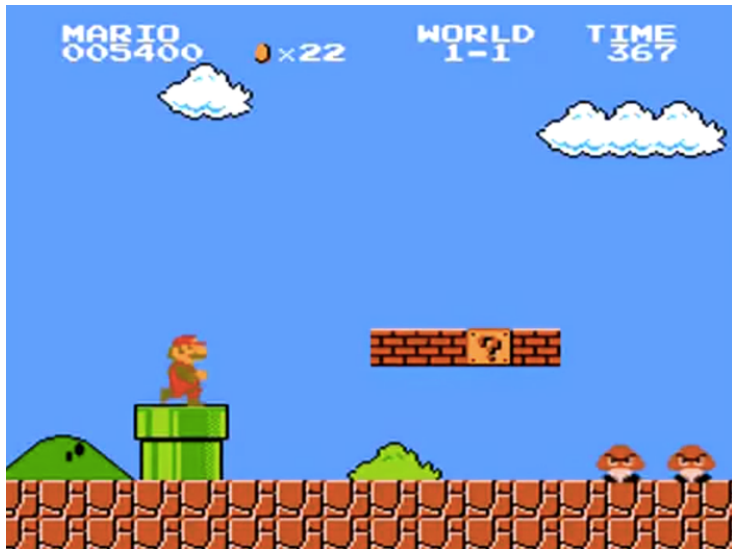
- ML is allows computers to learn from data.

Summary



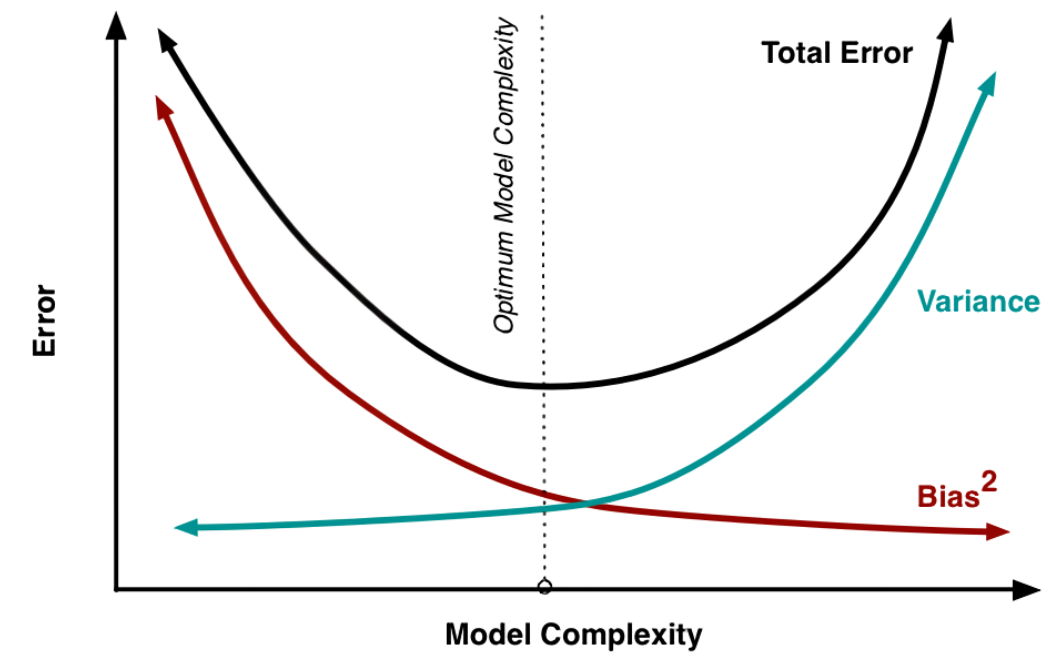
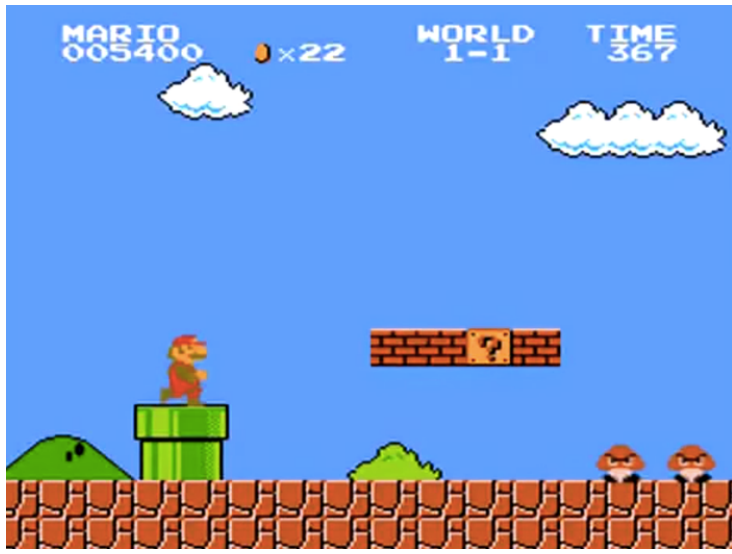
Summary

- ML allows computers to learn from data.
- ML-based AI has exploded in the last few years, especially generative AI for natural language tasks and image or video generation.



Summary

- ML allows computers to learn from data.
- ML-based AI has exploded in the last few years, especially generative AI for natural language tasks and image or video generation.
- **Bias**-**variance** decomposition gives a principled way to evaluate machine learning algorithms. Keep using it!



Summary

- ML allows computers to learn from data.
- ML-based AI has exploded in the last few years, especially generative AI for natural language tasks and image or video generation.
- **Bias**-**variance** decomposition gives a principled way to evaluate machine learning algorithms. Keep using it!
- We need to be careful with how we use it!