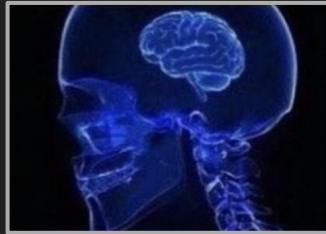


# Scaling paradigms for large language models

Jason Wei  
Research Scientist  
OpenAI

(Opinions are my own and do not reflect my employer.)

2019



- Can barely write a coherent paragraph
- Can't do any reasoning

2024



- Can write an essay about almost anything
- Competition-level programmer and mathematician

*Scaling has been the engine of progress in AI and will continue to dictate how the field advances.*

## Outline

What is scaling and why do it?

Paradigm 1: Scaling next-word prediction

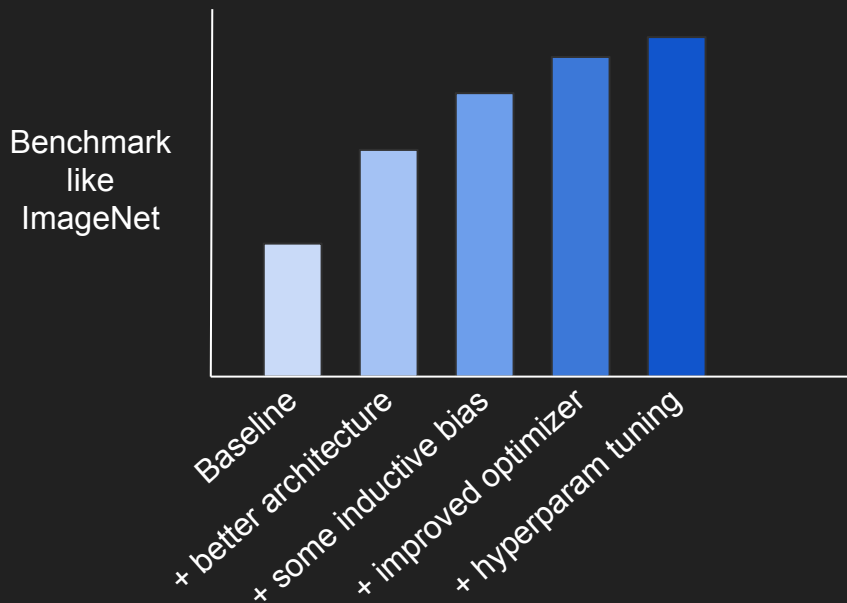
The challenge with next-word prediction

Paradigm 2: Scaling RL on chain-of-thought

How scaling changed AI culture & what's next?

“Studying the past tells you what’s special about the current moment.”

How we made progress,  
early 2010s to 2017  
(pre-transformer deep learning)



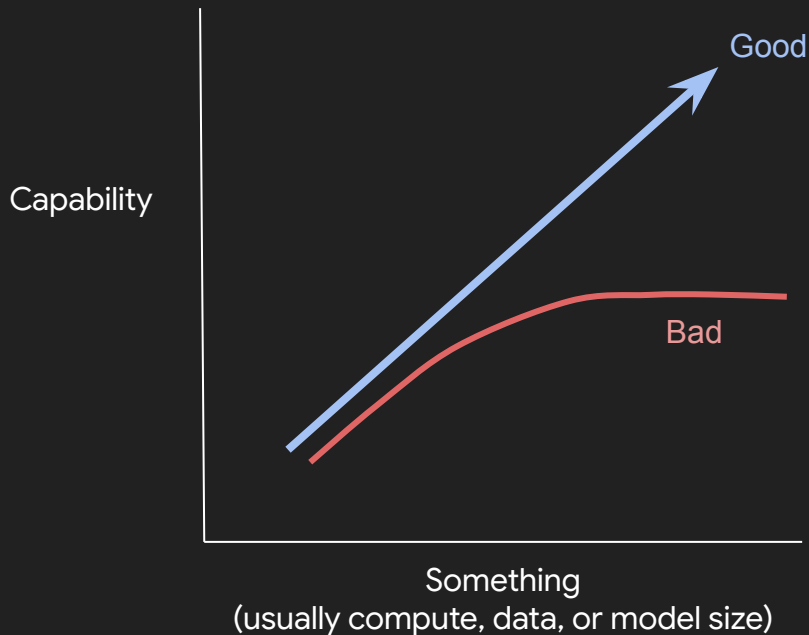
Make this as good as possible.

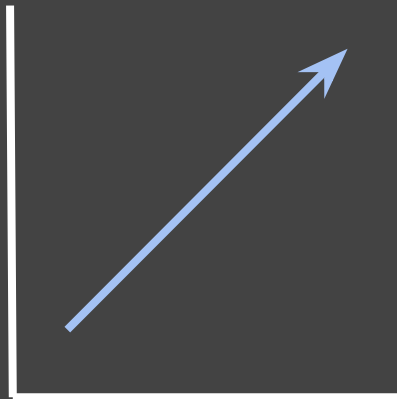
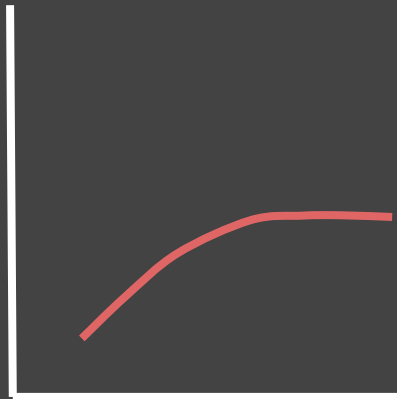
Success looks like “On the ImageNet dataset, our method outperforms the baseline by a significant margin using less compute.”



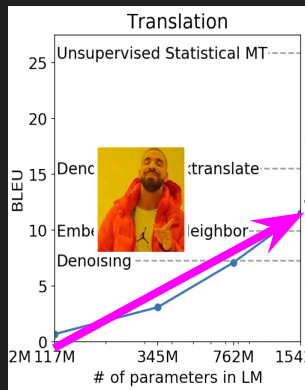
What is scaling?

Scaling is when you put yourself in a situation where you move along a continuous axis and expect sustained improvement.

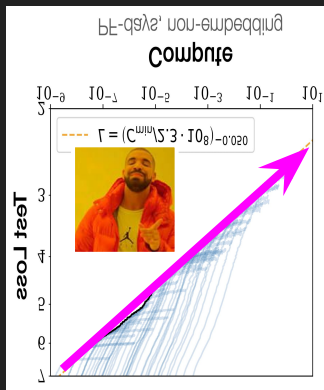




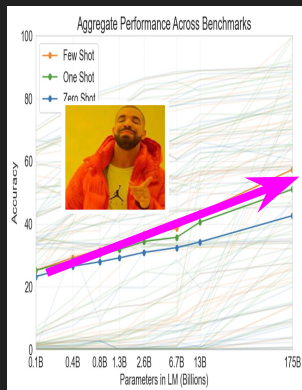
# Scaling is everywhere



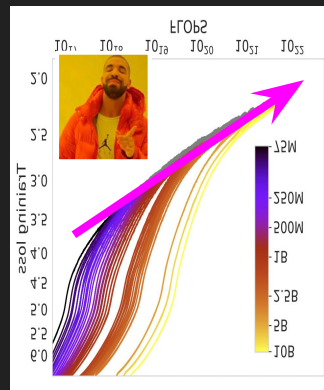
GPT-2 (2019)



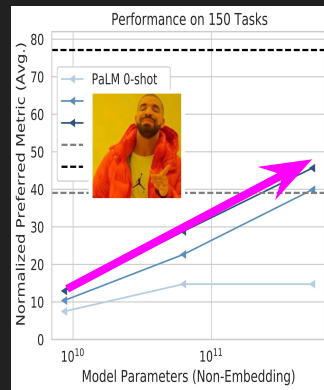
Scaling laws (2020)



GPT-3 (2021)



Chinchilla (2022)



PaLM (2022)





## *Why scale?*

### Not scaling

Each improvement in the model requires ingenuity on a new axis

There are a lot of tasks that we want AI to do

### Scaling-centric AI

You can reliably improve capability (even if it's expensive)

If your measure of capability is very general, extreme investment is justified

# The Bitter Lesson of AI

General methods that leverage compute are the most effective

Things that scale will ultimately win out

## 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 lesson is Moore's law, or rather its generalization of continued exponentially falling cost of computation. Most AI research has been conducted as if the computation available to the researcher was limited, in which case leveraging human knowledge would be one of the only ways to improve performance. In the long run, however, a slightly longer time than a typical research project, massively more computation is available. Seeking an improvement that makes a difference in the shorter term is a natural thing to do, but it does not leverage their human knowledge of the domain, but the only thing that matters in the long run is leveraging of computation. These two need not run counter to each other, but they often do. Time spent on one is time not spent on the other. There are psychological costs to pursuing one approach or the other. And the human-knowledge approach tends to come out on top, but many examples of AI researchers' belated learning of this bitter lesson, and it is one of the most prominent.

In computer chess, the methods that defeated the world champion, Kasparov, were massive, deep search. At the time, this was looked upon with dismay by the researchers who had pursued methods that leveraged human understanding of chess. When a simpler, search-based approach with special hardware and software was found to be effective, these human-knowledge-based chess researchers were not good for anything. "force" search may have won this time, but it was not a general strategy, and

# Paradigm 1: Scaling next-word prediction

Started in 2018, still ongoing

Get really, really good at predicting the next word.

Why do you get so much from “just” predicting the next word?  
Next-word prediction is massively multi-task learning.

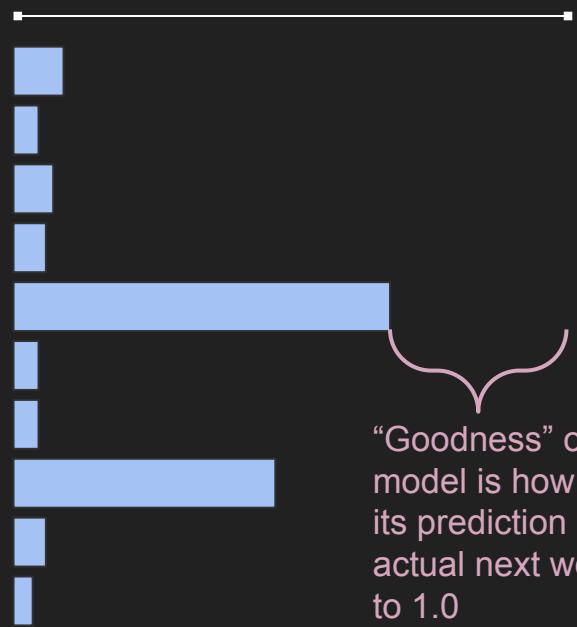
# Review: next-word prediction

On weekends,  
Dartmouth  
students like  
to \_\_\_\_



- a
- aardvark
- ...
- ...
- drink
- ...
- ...
- study
- ...
- zucchini

0.0 Probability 1.0



“Goodness” of the model is how close its prediction of the actual next word is to 1.0

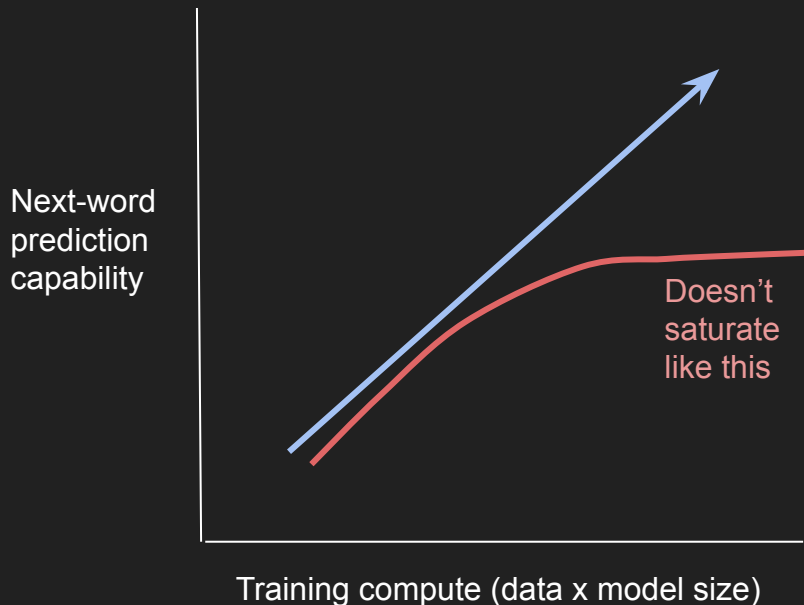
# Example “tasks” from next-word prediction

<u>Task</u>	<u>Example sentence in pre-training that would teach that task</u>
<i>Grammar</i>	In my free time, I like to { <b>code</b> , <b>banana</b> }
<i>World knowledge</i>	The capital of Azerbaijan is { <b>Baku</b> , <b>London</b> }
<i>Sentiment analysis</i>	Movie review: I was engaged and on the edge of my seat the whole time. The movie was { <b>good</b> , <b>bad</b> }
<i>Translation</i>	The word for “neural network” in Russian is { <b>нейронная сеть</b> , <b>привет</b> }
<i>Spatial reasoning</i>	Iroh went into the kitchen to make tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the { <b>kitchen</b> , <b>store</b> }
<i>Math question</i>	Arithmetic exam answer key: $3 + 8 + 4 =$ { <b>15</b> , <b>11</b> }

[millions more]

Extreme multi-task learning!

# Scaling predictably improves performance (“scaling laws”)



[Kaplan et al., 2020:](#)

*“Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute for training.”*

Jason’s rephrase: You should expect to get a better language model if you scale up compute.

## Why does scaling work?

Hard to answer, but here is a hand-wavy explanation

<u>Small language model</u>	<u>Large language model</u>
Memorization is costly	More generous with memorizing tail knowledge
First-order correlations	Complex heuristics


*If scaling was so predictable, why was the success of this paradigm so surprising?*

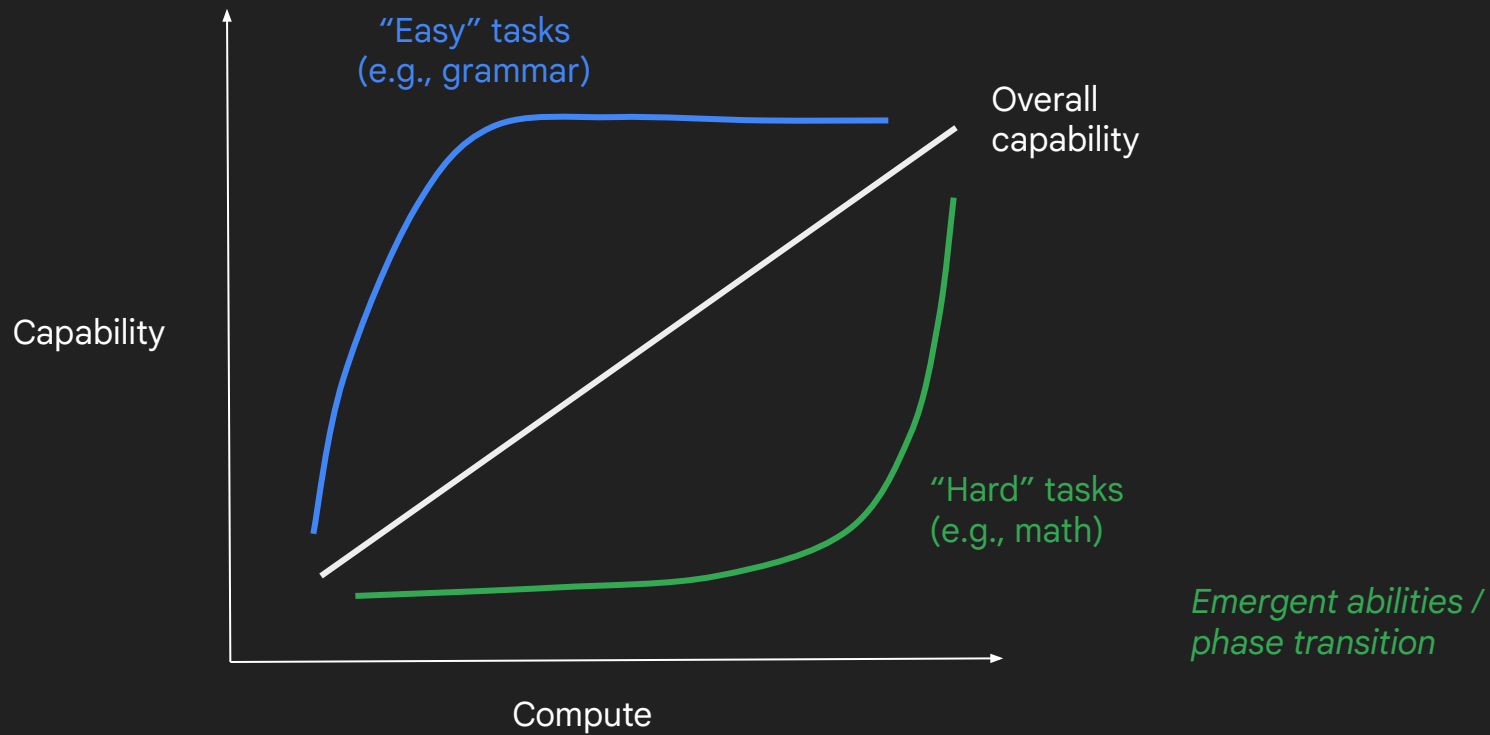
Next-word prediction is secretly massively multi-task, and performance on different tasks arise at different rates



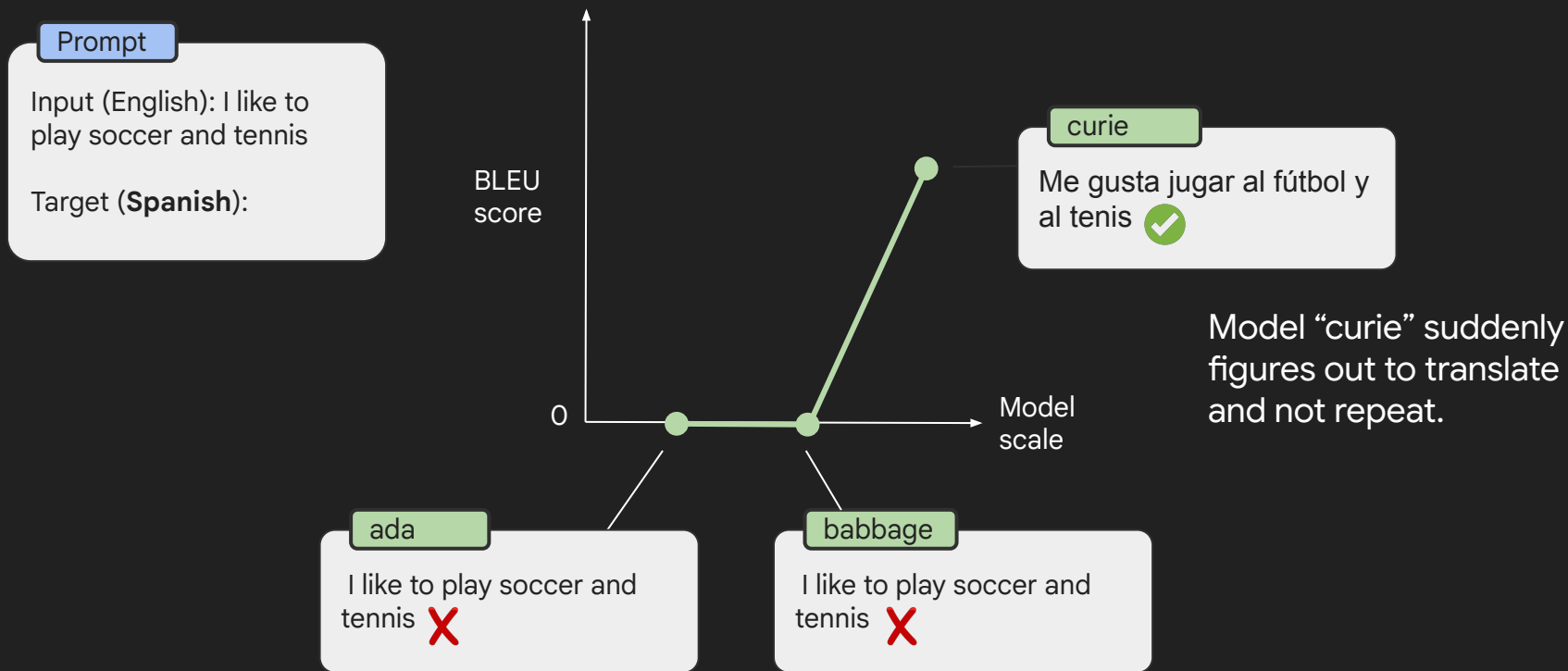
Let's take a closer look at next-word prediction accuracy. Consider that

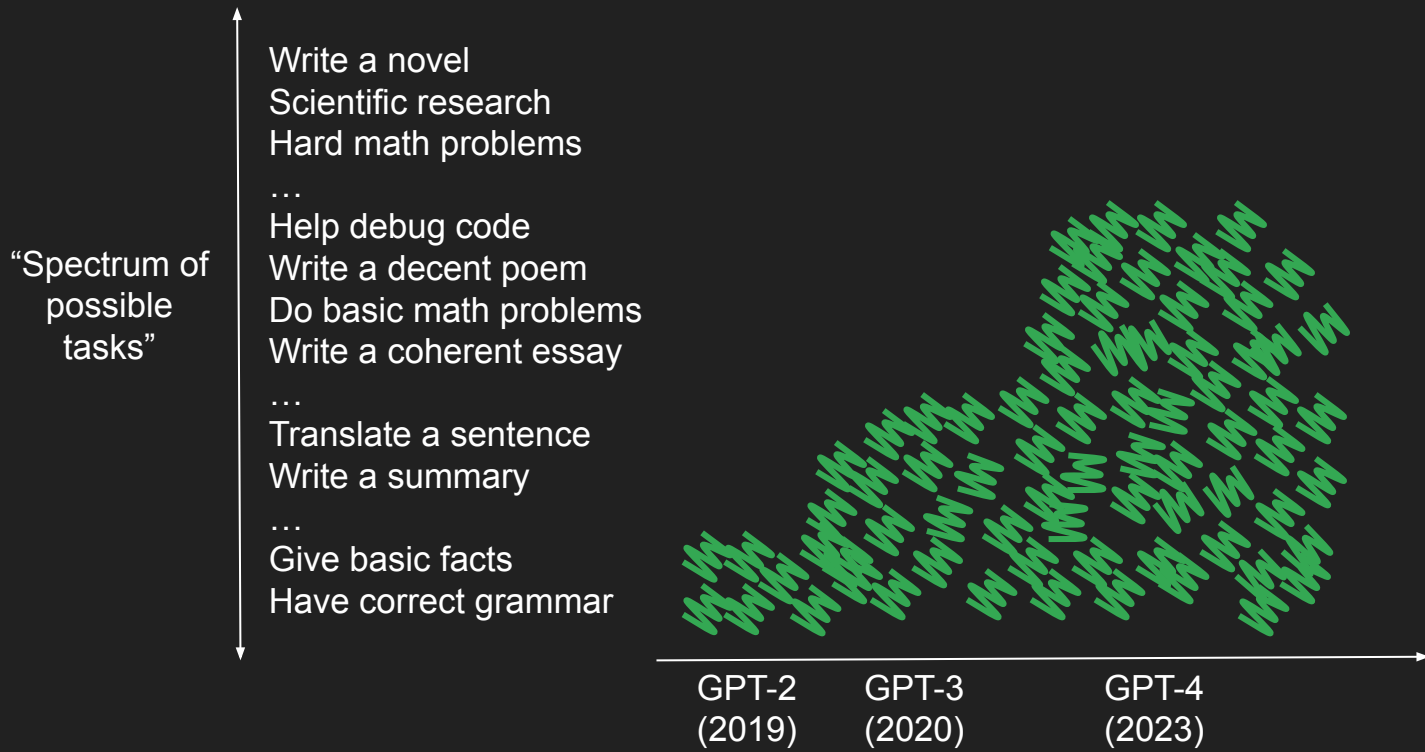
$$\begin{aligned} \text{Overall accuracy} = & 0.002 * \text{accuracy\_grammar} + \\ & 0.005 * \text{accuracy\_knowledge} + \\ & 0.000001 * \text{accuracy\_sentiment\_analysis} + \\ & \dots \\ & 0.0001 * \text{accuracy\_math\_ability} + \\ & 0.000001 * \text{accuracy\_spatial\_reasoning} \\ & \dots \end{aligned}$$

 *If accuracy goes from 70% to 80%, do all tasks get better uniformly?  
...probably not.*



# Emergence ability example





🤔 If next-word prediction works so well,  
can we scale it to reach AGI?

Maybe (it would be hard), but  
there is a bottleneck:

*Some words are super hard to  
predict and take a lot of work*

## When next-word prediction works fine

The screenshot shows the Playground interface with the text "My name is Jason Wei and I am a researcher at OpenAI working on large language models." A dropdown menu is open over the word "models", showing the following probabilities:

- models = 63.28%
- modeling = 11.41%
- model = 5.72%
- understanding = 3.98%
- datasets = 3.93%

Below the dropdown, a status bar indicates: "Total: -0.46 logprob on 1 tokens (88.31% probability covered in top 5 logits)".

## When next-word prediction becomes very hard

The screenshot shows the Playground interface with the question: "What is the square of  $((8-2)*3+4)^3 / 8?$ " and three multiple-choice options: (A) 1,483,492, (B) 1,395,394, and (C) 1,771,561. The answer "(C)" is selected. A dropdown menu is open over the answer, showing the following probabilities:

- C = 32.09%
- B = 29.98%
- A = 27.97%
- D = 8.15%
- c = 0.27%

Below the dropdown, a status bar indicates: "Total: -1.14 logprob on 1 tokens (98.44% probability covered in top 5 logits)".

Pretend you're ChatGPT. As soon as you see the prompt you have to immediately start typing... go!

*Question: What is the square of  $((8-2)*3+4)^3 / 8$ ?*

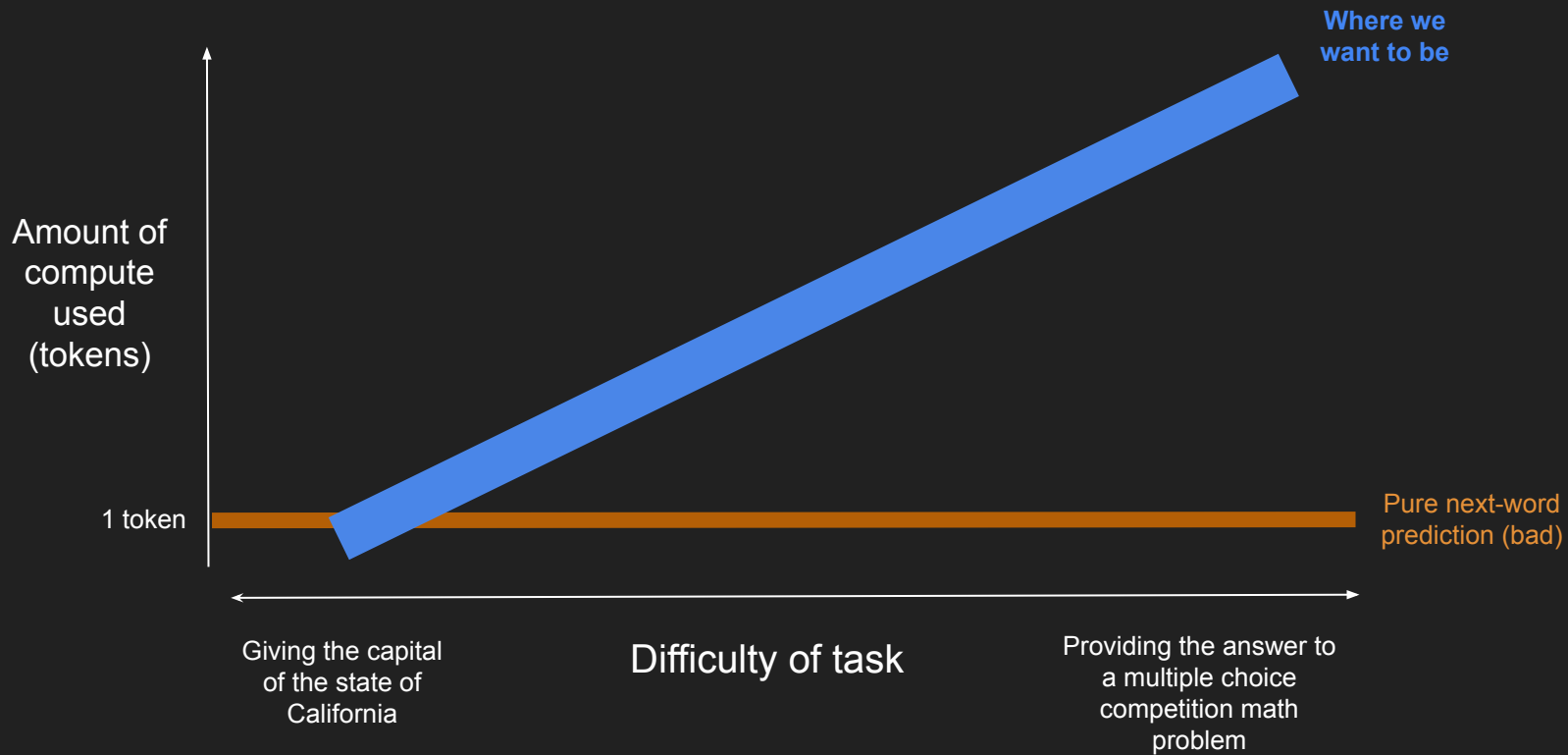
*(A) 1,483,492*

*(B) 1,395,394*

*(C) 1,771,561*

...

Tough right?





# An approach: chain-of-thought prompting

**Input**

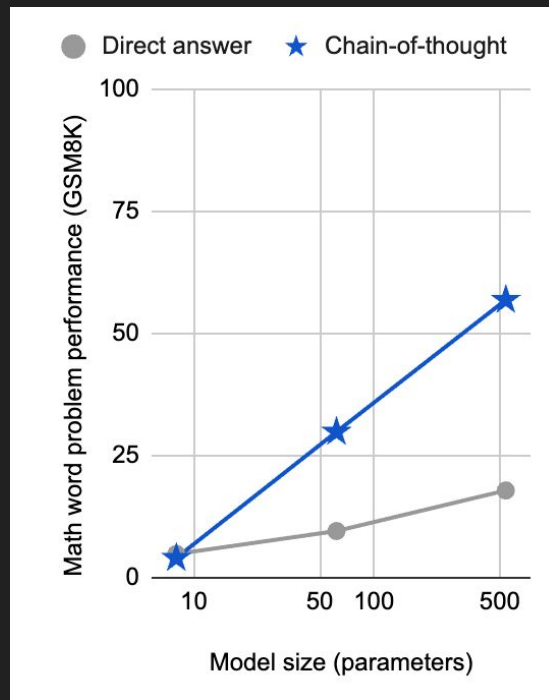
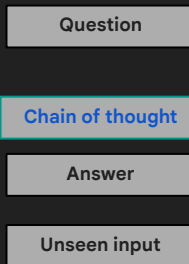
Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

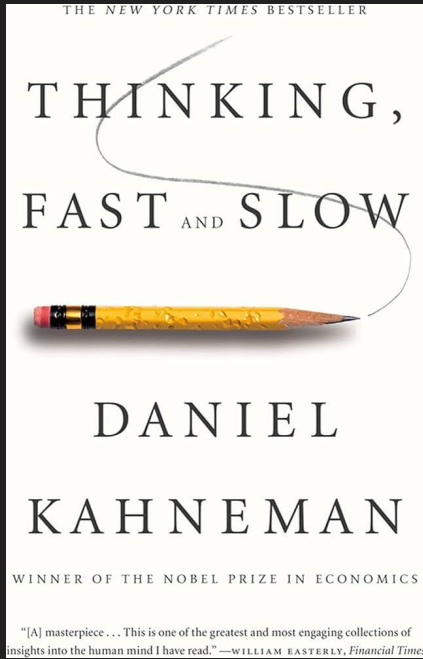
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✓



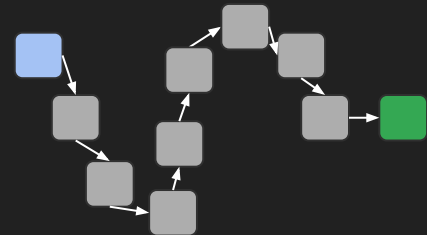


<u>System 1: Fast, intuitive thinking</u>	<u>System 2: Slow, deliberate thinking</u>
Automatic Effortless Intuitive Emotional	Conscious Effortful Controlled Logical
Recognizing faces Repeating basic facts Reacting to something	Solving math problems Planning a detailed agenda Making a thoughtful decision

Next-word prediction



Chain of thought



# The limitation with CoT prompting

## Most reasoning on the internet looks like this...

17-3: Formally prove Theorem 17.3.2.

**Theorem 17.3.2.** A one-pass algorithm for FREQUENCY-ESTIMATION with error  $\epsilon$  must use  $\Omega(\min\{m, n, \epsilon^{-1}\})$  space. In particular, in order to get an error of  $\epsilon$ , a randomized algorithm must use  $\Omega(\min\{m, n\})$  space.

*Proof.* We will prove the stronger result that the simpler FREQUENCY-ESTIMATION algorithm which asks whether the input stream contains a token whose frequency is at least  $\epsilon n$  is hard in the deterministic setting. Since the cost of the randomized algorithm is  $\Omega(\min\{m, n\})$  in the deterministic setting, this will prove the theorem as a whole.

Let  $\mathcal{A}$  be a one-pass  $S$ -space deterministic algorithm for FREQUENCY-ESTIMATION. On input  $(\mathbf{x}, y)$  for the  $\text{IDX}_N$ , Alice creates a stream  $\sigma_1 = (a_1, a_2, \dots, a_N)$  and Bob creates a stream  $\sigma_2 = (b, b, \dots, b)$  of length  $k-1$  for  $k \geq 2$  where  $b = 2y - 1$ . The combined stream is  $\sigma_1 \circ \sigma_2$  with parameter  $k$ .

The output of  $\text{IDX}_N(\mathbf{x}, y)$  is 1 iff  $\mathcal{A}$  produces  $b$  as output. This is so because  $b$  will be the unique entry with  $f_b = k \geq \epsilon n$ . Thus Alice and Bob can solve the problem by having Alice to Bob using  $\mathcal{A}$ .

By the lower bound result of  $\Omega(N)$  for  $\text{IDX}_N$ ,  $S = \Omega(N)$ . By construction,  $N + k - 1 \geq N + 1$ . Therefore, we have proven a lower bound of  $\Omega(N)$  for  $\text{IDX}_N$  and  $n \geq N + 1$ . We have thus proven that  $S = \Omega(\min\{m, n, \epsilon^{-1}\})$ , since  $n \geq N + 1$ .

What we actually want is the inner “stream of thought”

Hm let me first see what approach we should take...

Actually this seems wrong

No that approach won't work, let me try something else

Let me try computing this way now

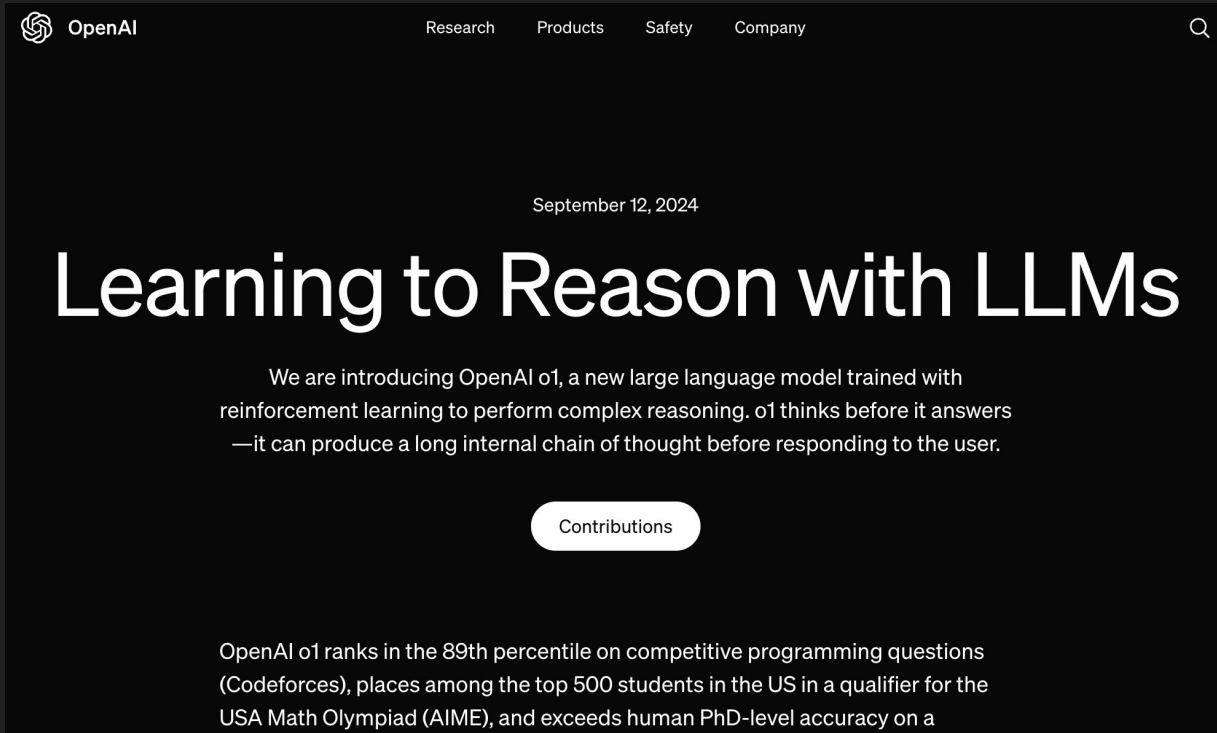
OK I think this is the right answer!

## Paradigm 2: Scaling RL on chain-of-thought

Train language models to “think” before giving an answer

In addition to scaling compute for training, there is a second axis here: scaling how long the language model can think at inference time.

# OpenAI o1 (work of most of the company)



The image is a screenshot of the OpenAI website's announcement for the o1 model. The page has a dark background with white text. At the top left is the OpenAI logo and name. To the right are navigation links for Research, Products, Safety, and Company, followed by a search icon. The date "September 12, 2024" is centered above the main title. The title "Learning to Reason with LLMs" is in a large, bold font. Below it is a paragraph of text describing the model's capabilities. A white pill-shaped button with the text "Contributions" is centered below the paragraph. At the bottom, there is another paragraph of text providing performance metrics for the model.

OpenAI

Research Products Safety Company

September 12, 2024

## Learning to Reason with LLMs

We are introducing OpenAI o1, a new large language model trained with reinforcement learning to perform complex reasoning. o1 thinks before it answers —it can produce a long internal chain of thought before responding to the user.

[Contributions](#)

OpenAI o1 ranks in the 89th percentile on competitive programming questions (Codeforces), places among the top 500 students in the US in a qualifier for the USA Math Olympiad (AIME), and exceeds human PhD-level accuracy on a

# A chain of thought from OpenAI o1

First, let's understand what is being asked.

Both  $NH_4^+$  and  $F^-$  can react with

Water (right)

$F^-$  ( $K_a(\text{right}) = 10^{-2}$ ) = -1.5800

Then:

$$pH = 7 + 0.5 \times (-1.5800) = 7 - 0.79 = 6.21$$

Therefore, the pH is approximately 6.21.

# CoT allows models to leverage asymmetry of verification

A class of problems has “asymmetry of verification”, which means it’s easier to verify a solution than to generate one

For example, a crossword puzzle, sudoku, or writing a poem that fits constraints

The screenshot shows a crossword puzzle interface with a reasoning overlay. The puzzle grid is partially visible on the left, with clues listed on the right. The reasoning overlay, titled "Across", shows the model's thought process for solving clue 1. It lists possible words like "ESTATE", "DODGE", "ELUDE", "MAYBE", and "BUT LET", and then concludes that "ESTATE" is the correct answer based on the context of British English and its alignment with another clue.

Solve the crossword puzzle.

Plain Text

1 +  
2 |  
3 +  
4 |  
5 +  
6 |  
7 +  
8 |  
9 +  
10 |  
11 +  
12 |  
13 +

Across

1. Ev...  
2. On...  
3. Mc...  
4. Init...  
5. Na...  
6. Mis...

Down:

1. ...  
2. Au...  
3. Pro...  
4. Syl...  
5. An...  
6. Deletes

1 Across

Possible words:

ESTATE  
AVOID  
DODGE  
ELUDE  
MAYBE  
BUT LET

Now let's look at Down clues.

1 Down: \_\_\_\_\_ car (station wagon) (6 letters)

Possible words:

- ESTATE car (6 letters)

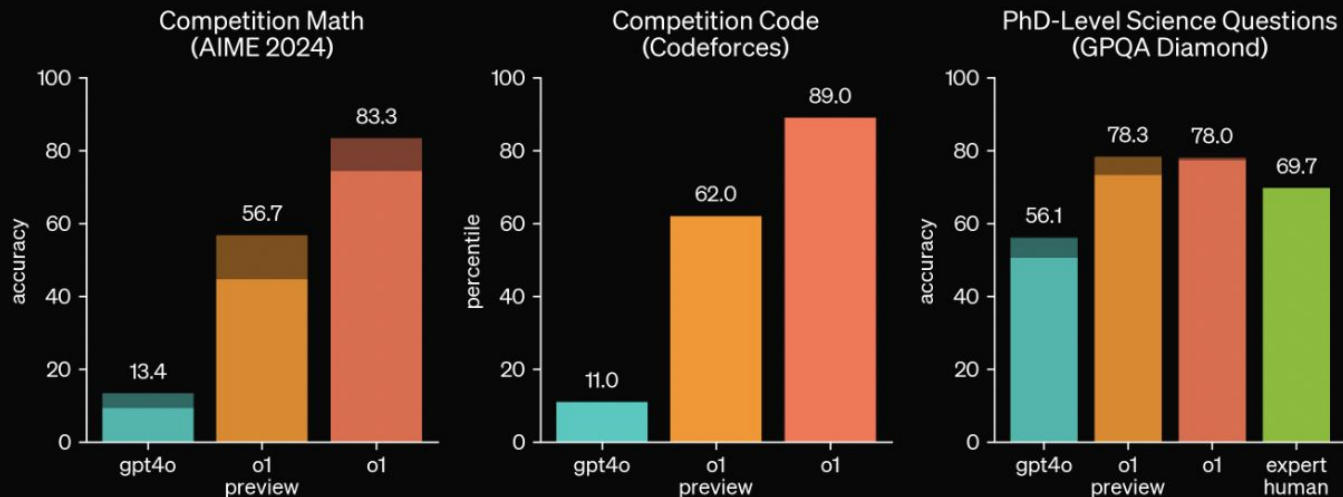
In British English, 'Estate car' is a term for station wagon.

Since 'station wagon' is called 'estate car' in the UK.

Therefore 'ESTATE' fits.

Also aligns with ESCAPE as Across 1.

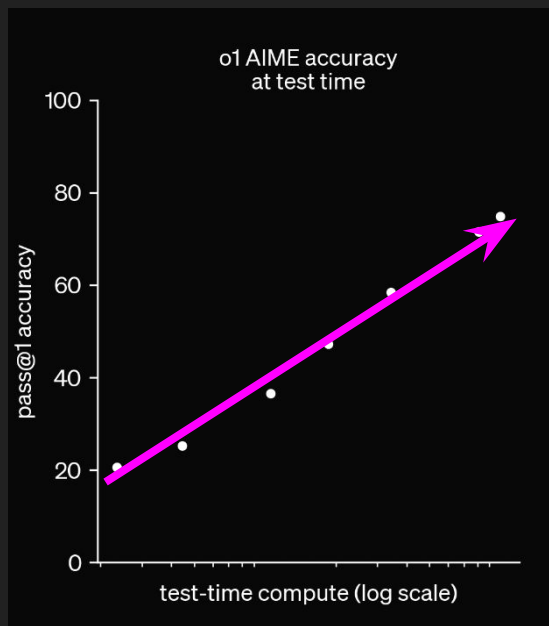
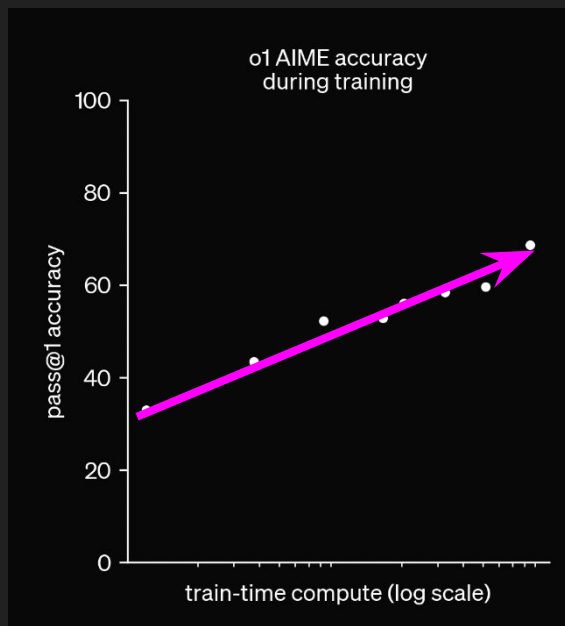
# Scale RL on chain-of-thought



o1 greatly improves over GPT-4o on challenging reasoning benchmarks. Solid bars show pass@1 accuracy and the shaded region shows the performance of majority vote (consensus) with 64 samples.



# Scale inference-time compute



# Why is this special: one day we may want AI to solve very challenging problems

## Prompt

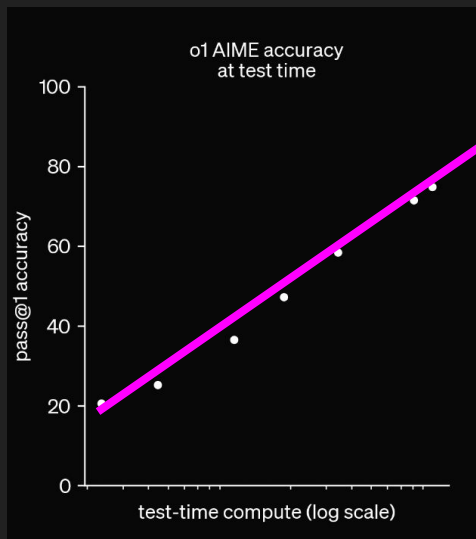
Write the code, documentation, and research paper for the best way to make AI safe

## Hypothetical response

Let me think very hard about this...

[Researches all the existing literature]  
[Data analysis] [Conducts new experiments]

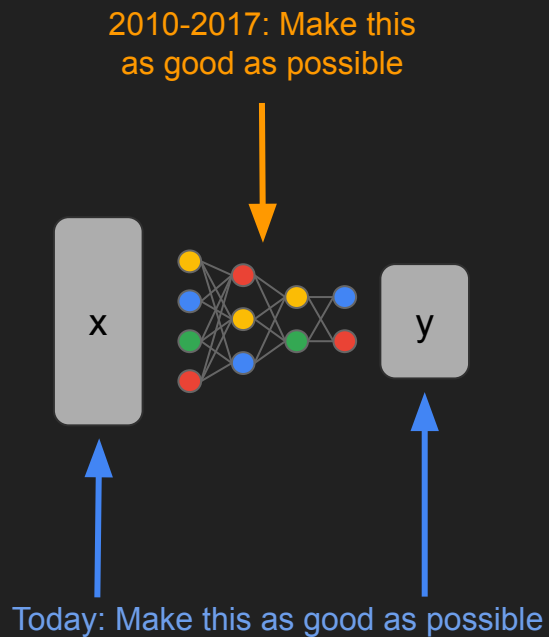
OK, here is a body of work on how to make AI safe



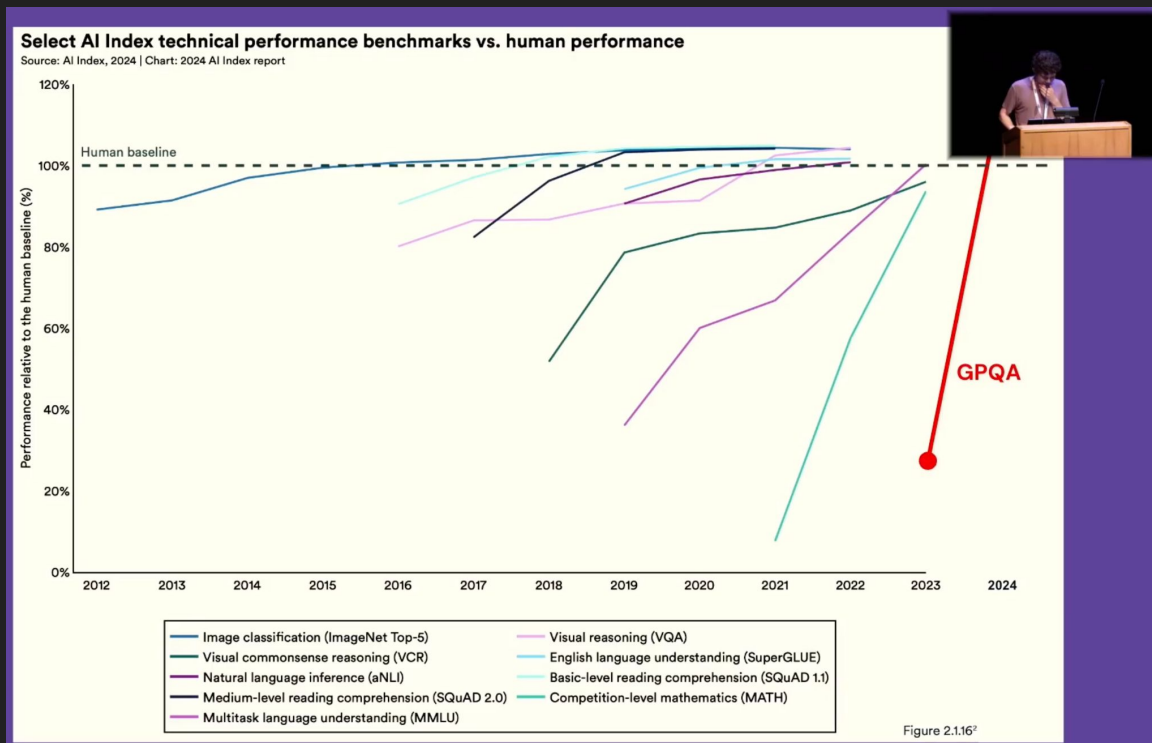
seconds minutes hours days weeks **months**

How has scaling changed the culture around doing AI research?

# Changes in AI research culture: shift to data

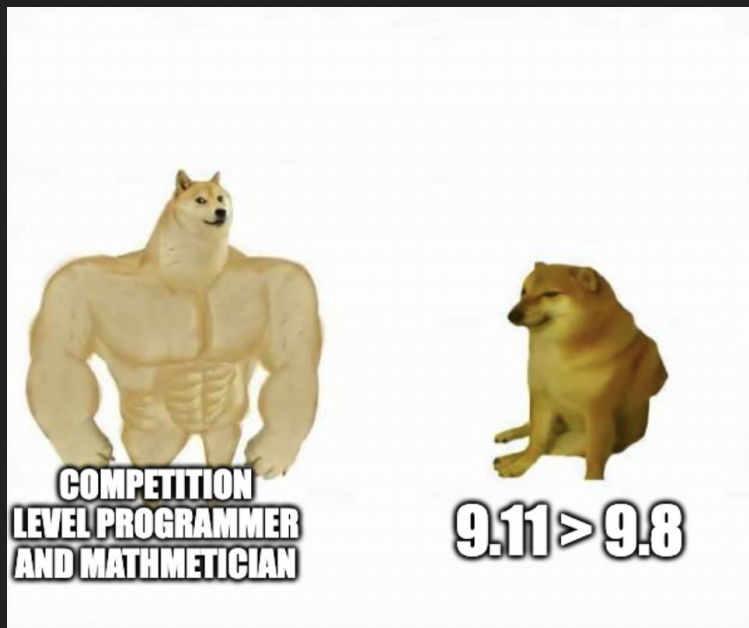


# Changes in AI culture: we desperately need evals



“People ask me if I’m making an even harder version of GPQA... [well] we set out to make the hardest science benchmark that we could”  
- David Rein

# Changes in AI culture: highly multi-task models



Language models must be measured on many dimensions

Hard to say that one model is strictly better than another

AI doesn't need to human-level on everything

Intelligence != user experience

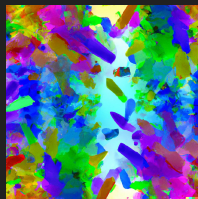


# Where will AI continue to progress?



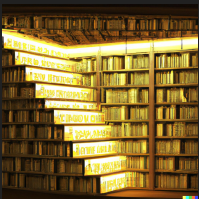
## AI for science and healthcare

As an assistant in scientific and medical innovation



## Tool use

Goal: enable AI to interact with the world



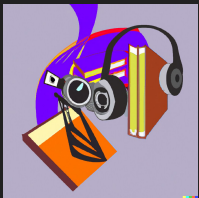
## More factual AI

Reduced hallucinations, cite sources, calibration



## AI applications

More ubiquitous use of AI

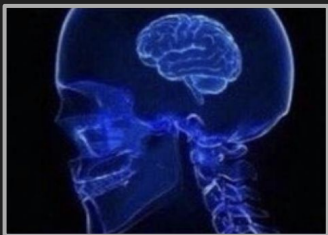


## Multimodality

AI to see, hear, and speak



2019



- Can barely write a coherent paragraph
- Can't do any reasoning

2024



- Can write an essay about almost anything
- Competition-level programmer and mathematician

2029



?

*Scaling has been the engine of progress in AI and will continue to dictate how the field advances.*



X / Twitter: @\_jasonwei  
OpenAI roles: jasonwei@openai.com

Feedback? <https://tinyurl.com/jasonwei>