

Intuitions on language models

Jason Wei

OpenAI

Fundamental question. Why do large language models work so well?

Thing I've been thinking about recently: Manually inspecting data gives us clear intuitions about how the model works.

Looking at data = training your biological neural net.

Your biological neural net makes many observations about the data after reading it.

These intuitions can be valuable.

(I once manually annotated an entire lung cancer image classification dataset. Several papers came out of intuitions from that process.)

Intuition 1.

Next-word prediction (on large data) is massively multi-task learning.

Language models are just next-word predictors



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 { code , banana }
<i>Lexical semantics</i>	I went to the store to buy papaya, dragon fruit, and { durian , squirrel }
<i>World knowledge</i>	The capital of Azerbaijan is { Baku , London }
<i>Sentiment analysis</i>	Movie review: I was engaged and on the edge of my seat the whole time. The movie was { good , bad }
<i>Translation</i>	The word for “pretty” in Spanish is { bonita , hola }
<i>Spatial reasoning</i>	Iroh went into the kitchen to make tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the { kitchen , store }
<i>Math question</i>	Arithmetic exam answer key: $3 + 8 + 4 =$ { 15 , 11 }

[millions more]

Extreme multi-task learning!

There are a lot of possible “tasks”, and they can be arbitrary

<u>Input</u>	<u>Target</u>	<u>Task</u>
Biden married Neilia	Hunter	world knowledge
Biden married Neilia Hunter	,	comma prediction
Biden married Neilia Hunter ,	a	grammar
Biden married Neilia Hunter , a	student	impossible?

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

Being a language model is not easy! A lot of arbitrary words to predict. Tasks aren't weird and not clean.

Language models will learn what is in the data

ChatGPT (2022)

10 facts about jason wei



I'm sorry, but I don't have specific information about someone named Jason Wei. If you have any specific details or context about Jason Wei that you'd like to share, I can try to help based on that information.



ai researcher



As of my last update in January 2022, I don't have specific information about an AI researcher named Jason Wei. If Jason Wei is a more recent figure or someone who has become prominent after that time, I may not have information on him. If you have specific details or a context about his work or contributions to AI research, feel free to share, and I can help based on that information.

ChatGPT (2024)

10 facts about jason wei



Jason Wei is a researcher and scientist known for his work in the field of artificial intelligence and natural language processing.

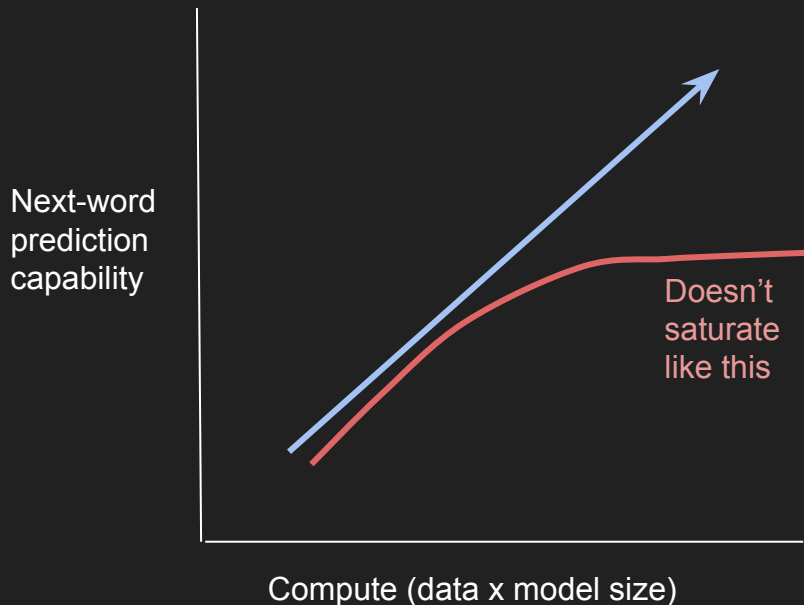
Here are ten facts about him:

1. **Background:** Jason Wei holds a Ph.D. in Computer Science from Stanford University, where he conducted research in AI and NLP under the supervision of notable scholars in the field.
2. **Workplace:** Jason Wei has worked as a research scientist at Google Research, specifically in the Brain team, which focuses on advancing machine learning and AI technologies.
3. **Contributions to NLP:** He is known for his contributions to natural language processing (NLP), particularly in areas like few-shot learning, prompt engineering, and large language models.

Intuition 2.

Scaling language models (compute = model size * data) reliably improves loss.

Scaling predictably improves performance (“scaling laws”)



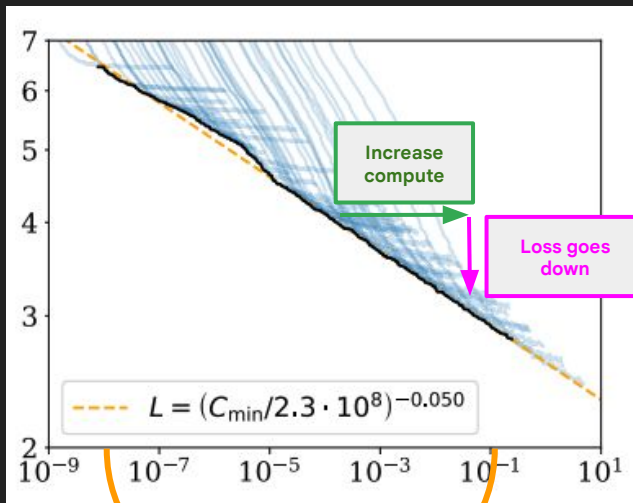
[Kaplan et al., 2020:](#)

“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.

Plot from the “scaling laws” paper

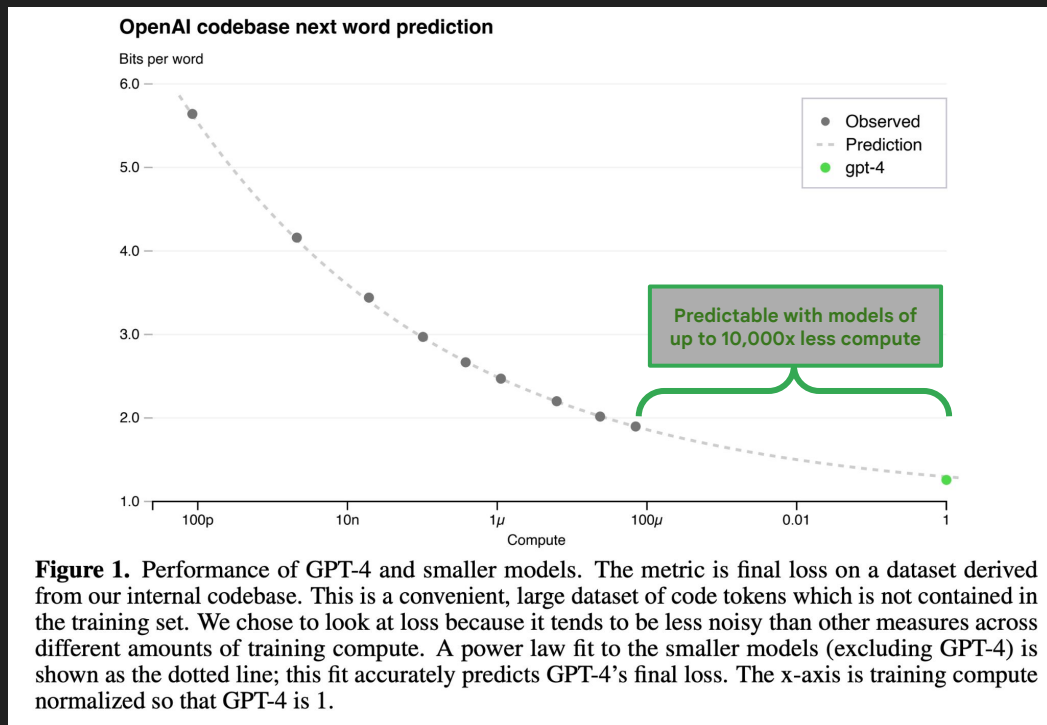
Scaling laws for neural language models. Kaplan et al., 2020.



Compute (data x model size). Log scale!

Seven orders of magnitude

Scaling laws: certain metrics can be very predictable



Why does scaling work? Hard to answer, but here are some guesses

<u>Small language model</u>	<u>Large language model</u>
<p>Memorization is costly</p> <p><i>“Parameters are scarce, so I have to decide which facts are worth memorizing”</i></p>	<p>More generous with memorizing tail knowledge</p> <p><i>“I have a lot of parameters so I’ll just memorize all the facts, no worries”</i></p>
<p>First-order correlations</p> <p><i>“Wow, that token was hard. It was hard enough for me to even get it in the top-10 predictions. Just trying to predict reasonable stuff, I’m not destined for greatness.”</i></p>	<p>Complex heuristics</p> <p><i>“Wow, I got that one wrong. Maybe there’s something complicated going on here, let me try to figure it out. I want to be the GOAT.”</i></p>

Changes in the nature of AI work: scaling laws

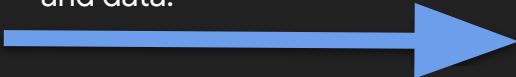
5 years ago

Many individual or small-scale projects.

Bottom-up research culture.

Run the code once; then submit to NeurIPS.

Technical paradigm shift:
Training the best models requires scaling compute and data.



Now

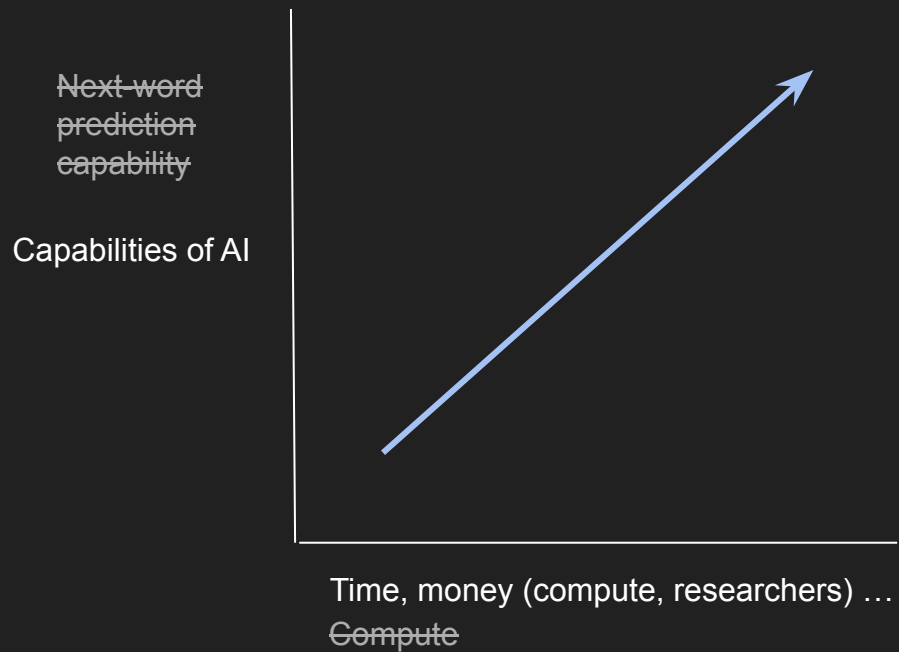
Teams have dozens of people.

Everyone works together towards one focused goal.

Tooling and infra matter a lot.

(In my experience, obviously there are exceptions.)

The general AI scaling law



Intuition 3.

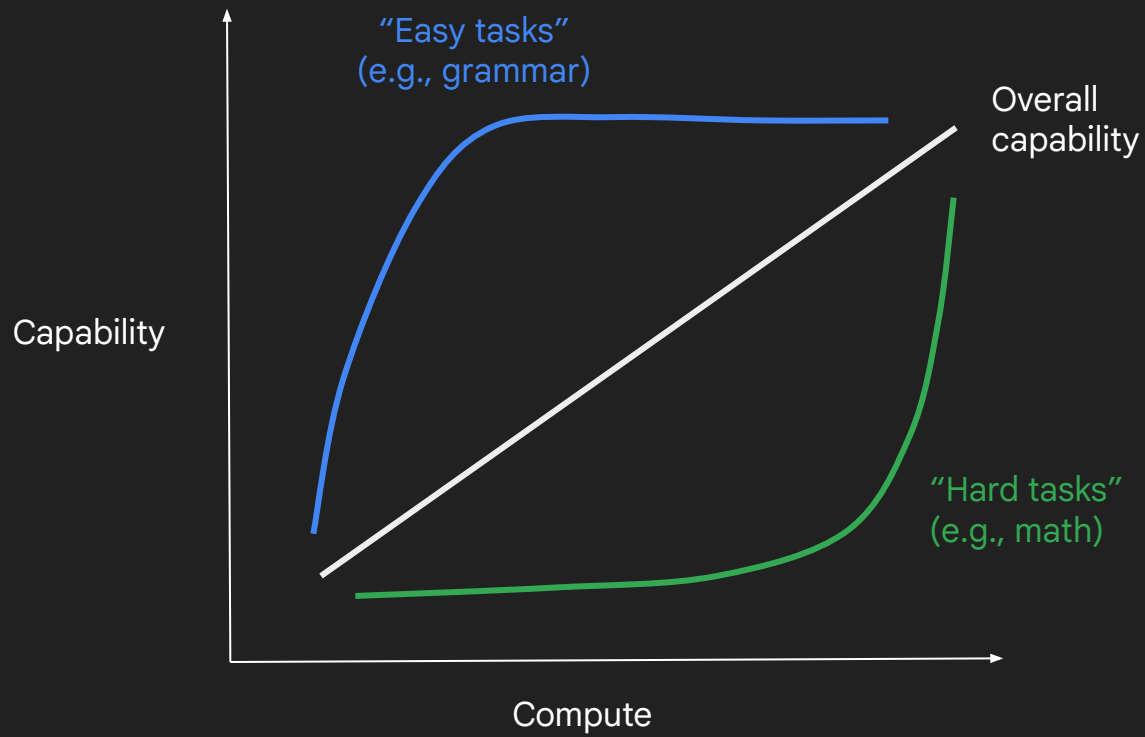
While overall accuracy scales smoothly, individual downstream tasks may improve in an emergent fashion.

Let's take a closer look at accuracy. Consider:

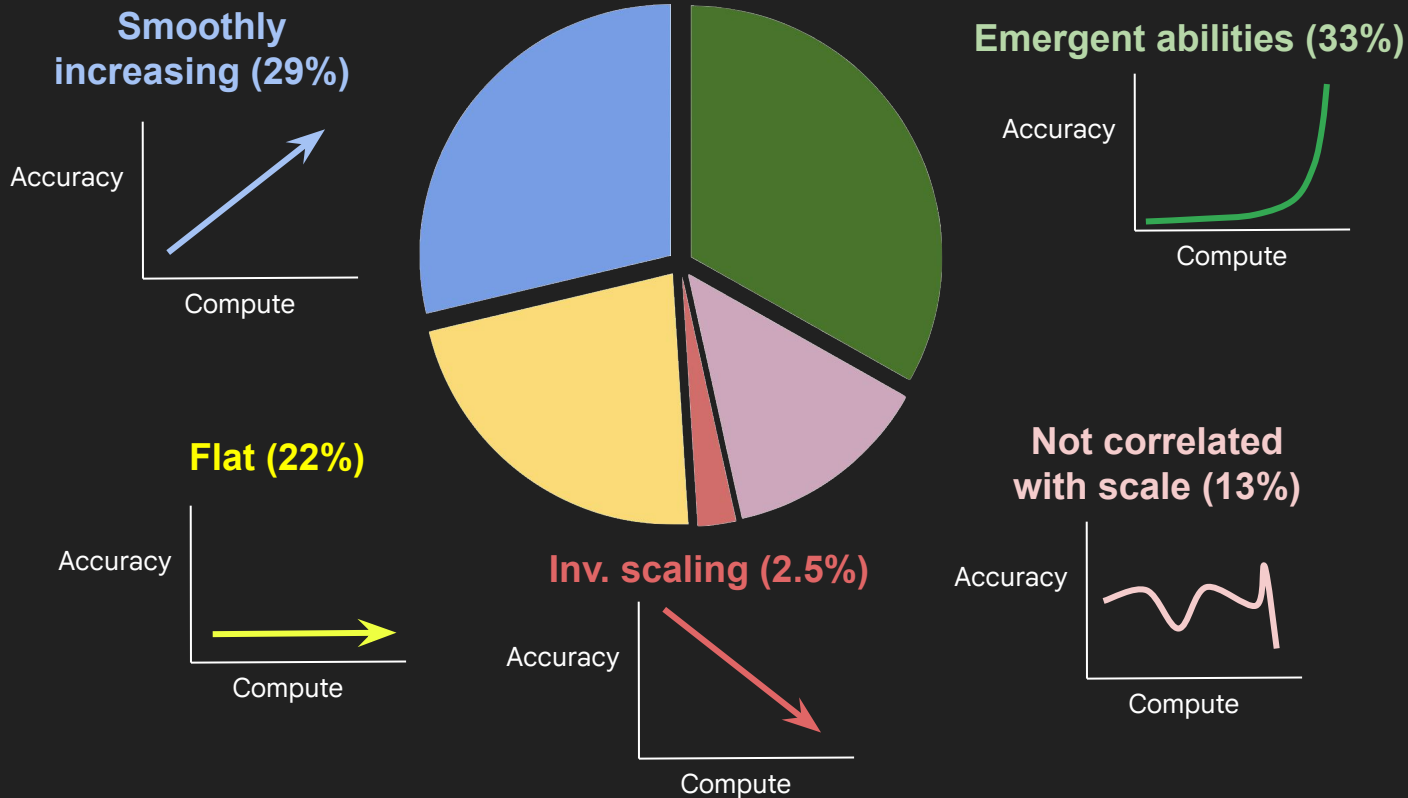
$$\begin{aligned} \text{Overall accuracy} = & 1e-3 * \text{accuracy_grammar} + \\ & 1e-3 * \text{accuracy_knowledge} + \\ & 1e-6 * \text{accuracy_sentiment_analysis} + \\ & \dots \\ & 1e-4 * \text{accuracy_math_ability} + \\ & 1e-6 * \text{accuracy_spatial_reasoning} \\ & \dots \end{aligned}$$

→ If accuracy goes from 70% to 80%, do all tasks get better uniformly?

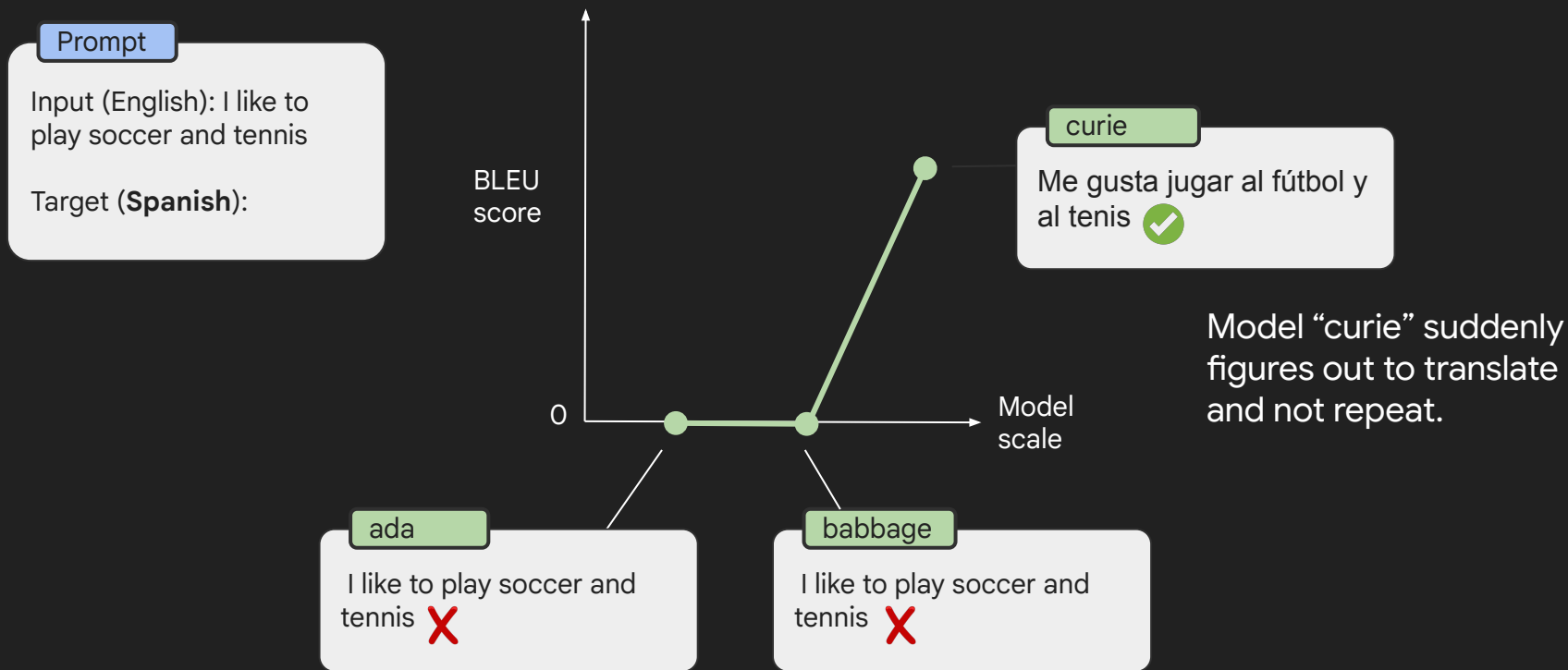
Probably not.

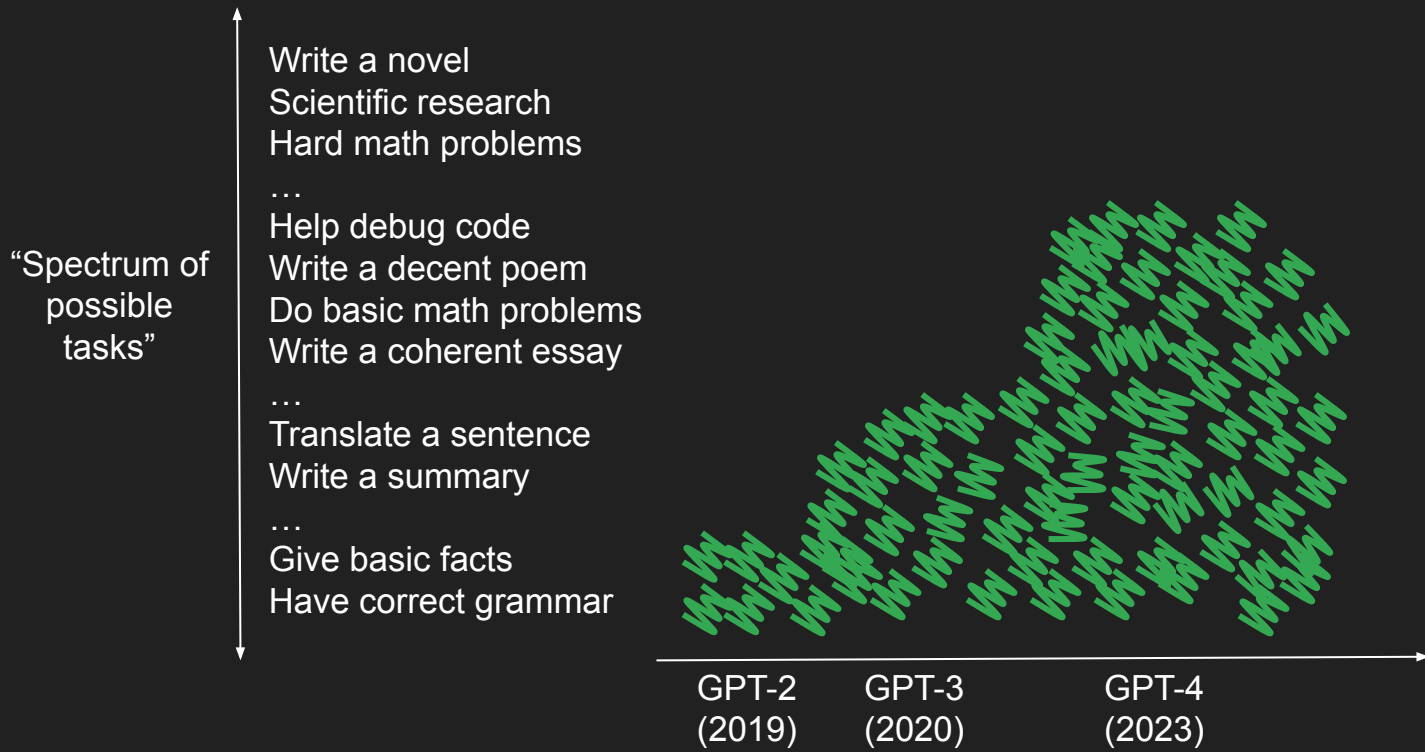


Real task distribution: 202 tasks in [BIG-Bench](#)



Emergence in prompting: example





Changes in the nature of AI work: emergent abilities

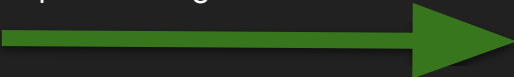
5 years ago

A few benchmarks for many years (CIFAR, ImageNet).

Easy to rank models (top-1 acc, avg. SuperGLUE score).

Task-specific architectures, data, and protocols.

Technical paradigm shift: A single model performs many tasks without the tasks being explicitly specified at pre-training.



Now

Need to create new benchmarks all the time.

Hard to decide if one model is universally better.

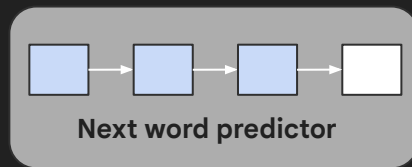
Create general technology; relatively easy to pivot.

Thanks.

X / Twitter: @_jasonwei

How to think about language model pre-training

PRE-TRAINING
- DATA -



How to think about language model pre-training

Millions of tasks!

In my free time, I like to buy papaya in Spanish. I went to the capital of Azerbaijan, Baku, to see the opera. I went into the kitchen to make tea. Standing next to him, Zuko pondered his destiny. Zuko left the kitchen store.
Arithmetic: $3 + 8 + 4 =$
{15, 11}

