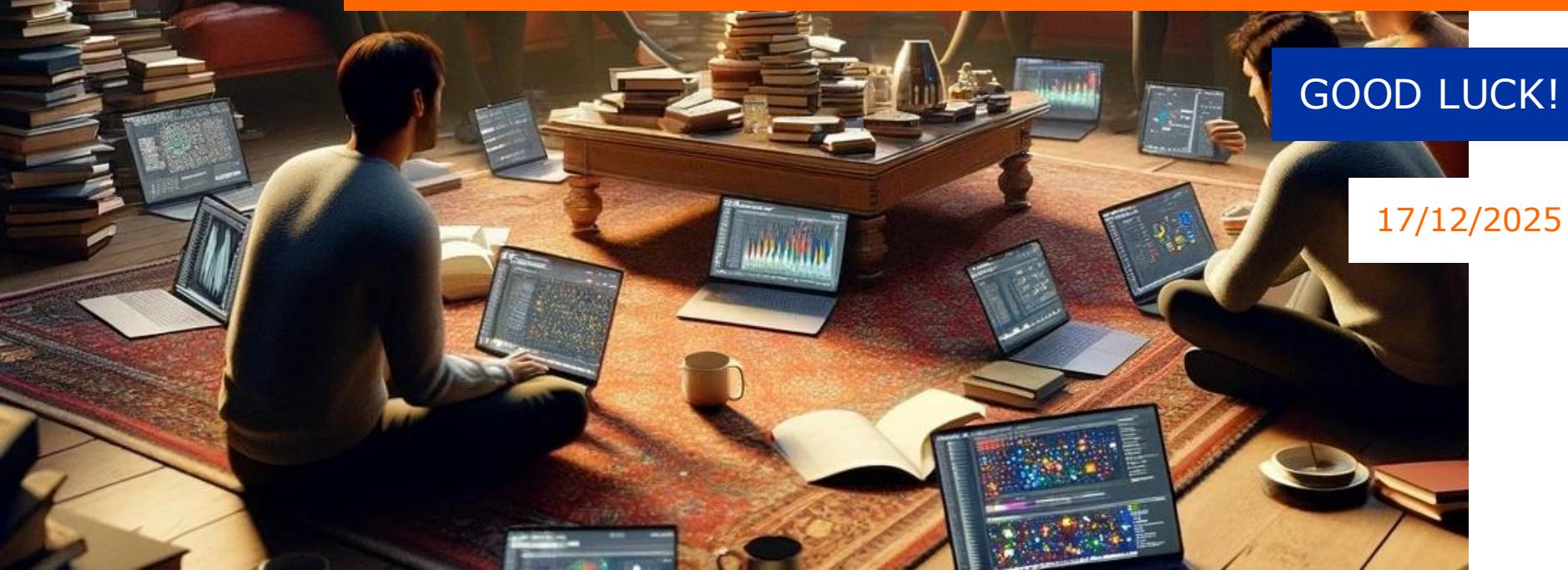




LARGE LANGUAGE MODELS WITHOUT HPC?



fwo

Welcome Sofia!



I'll explain why I'm excited in three parts

1. Brief intro into LLMs

2. Research in our lab

3. Beyond our lab

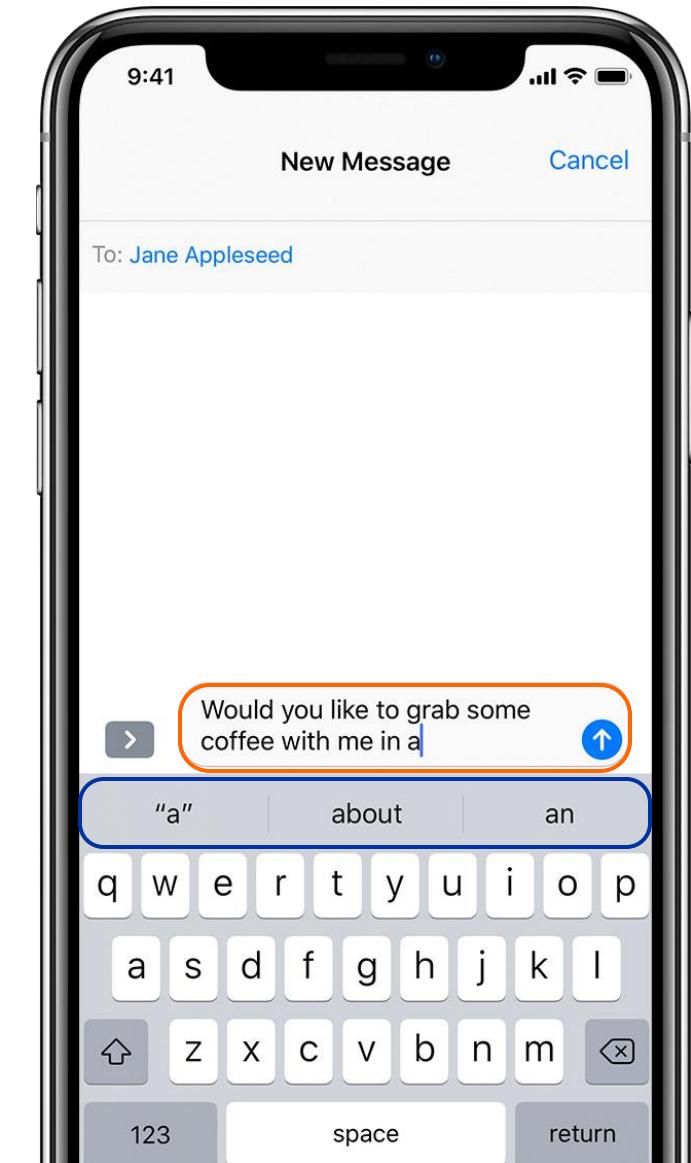
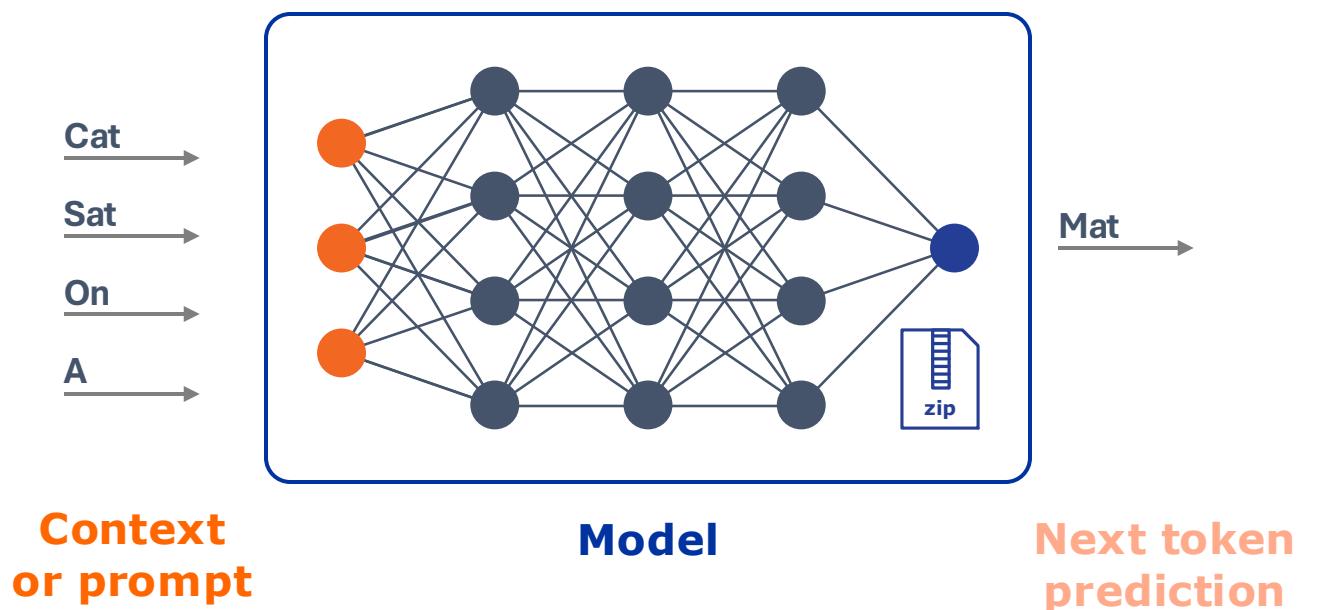


WHAT ARE LARGE LANGUAGE MODELS?



Generative Pre-trained Transformer (GPT)

The model predicts the next token based on a given context/prompt.



GPTs are pre-trained on internet data

Ruth Marianna Handler (née Mosko; November 4, 1916 – April 27, 2002) was an American businesswoman and inventor. She is best known for inventing the Barbie doll in 1959, [2] and being co-founder of toy manufacturer Mattel with her husband Elliot, as well as serving as the company's first president from 1945 to 1975. [3]

The Handlers were forced to resign from Mattel in 1975 after the Securities and Exchange Commission investigated the company for falsifying financial documents. [3][4]

Early life [edit]

Ruth Marianna Mosko [5][2][3] was born on November 4, 1916, in Denver, Colorado, to Polish-Jewish immigrants Jacob Moskowicz, a blacksmith, and Ida Moskowicz, née Rubenstein. [6]

She married her high school boyfriend, Elliot Handler, and moved to Los Angeles in 1938, where she found work at Paramount. [7]

Ruth Handler



Handler in 1961

Born	Ruth Marianna Mosko November 4, 1916 <u>Denver, Colorado</u> , U.S.
Died	April 27, 2002 (aged 85) ^[1] <u>Los Angeles, California</u> , U.S.

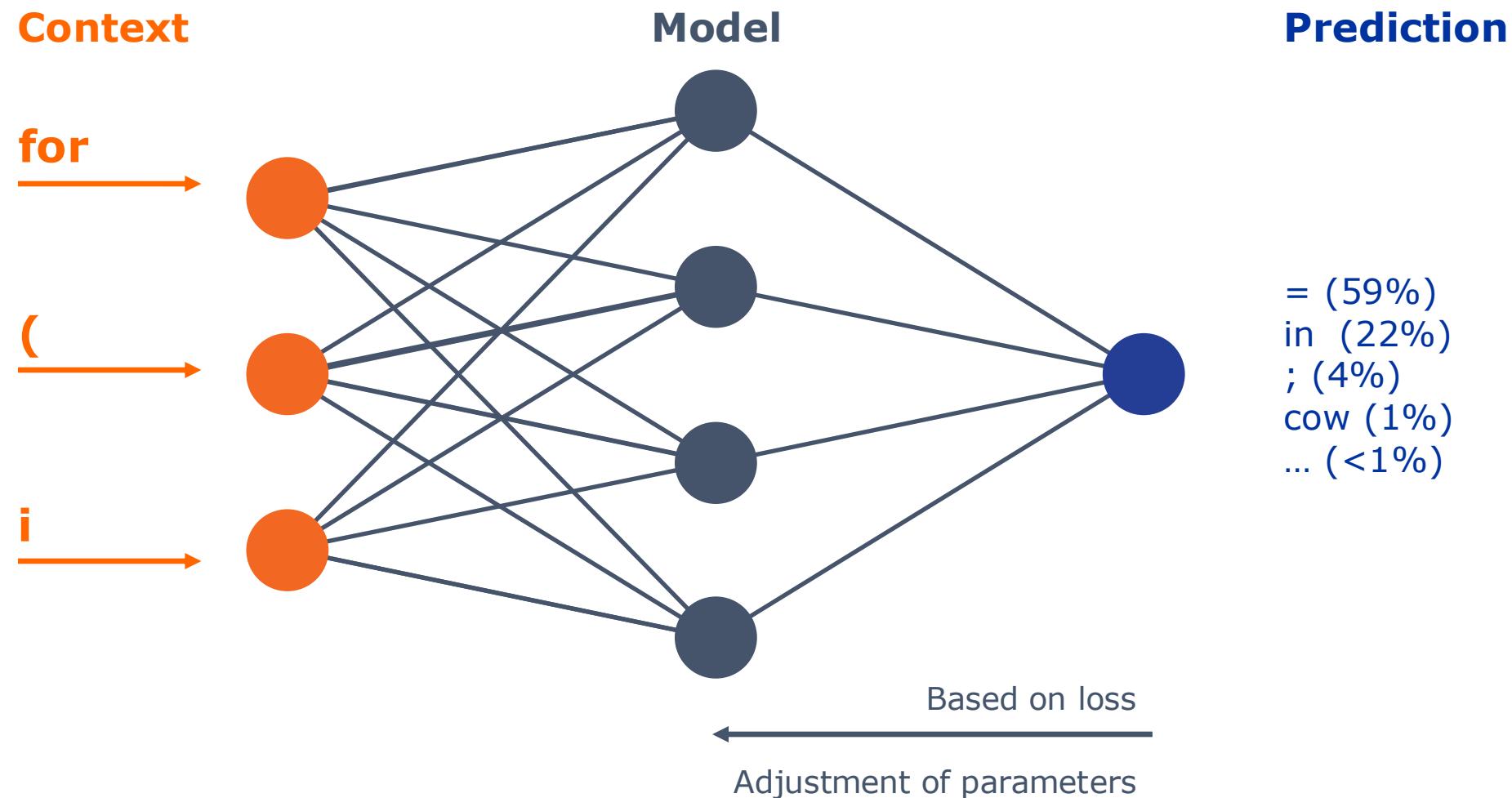
Next token prediction is extreme multi-task learning

A single model needs to learn a lot about the world

Predicting the next token requires the model to be able to solve a lot of different tasks.

Task	Example
Grammar	In my free time, I like to {run, banana }
Lexical semantics	I went to the zoo to see lions and {zebras, spoon }
World knowledge	The capital of Denmark is {Copenhagen, London }
Sentiment analysis	I really liked the movie. The movie was {good, bad }
Translation	The word for pretty in Spanish is {bonita, hola }
Spatial reasoning	After opening the fridge, I left the {kitchen, bedroom }
Math question	$3 + 8 + 4 = \{15, \textcolor{orange}{11}\}$

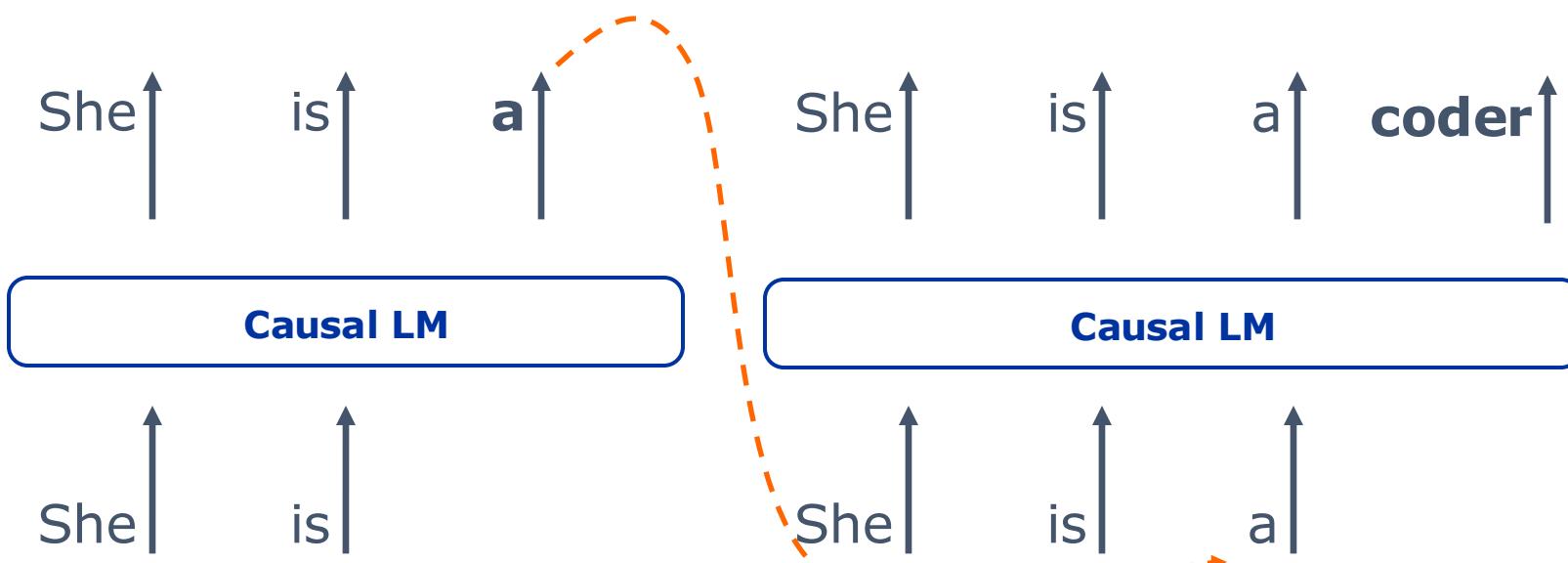
Pre-training is compressing information through prediction



Once it has been trained, we are going to predict

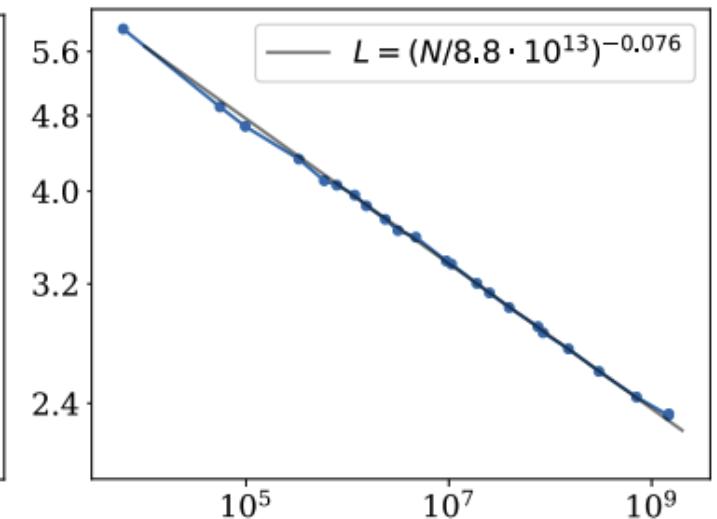
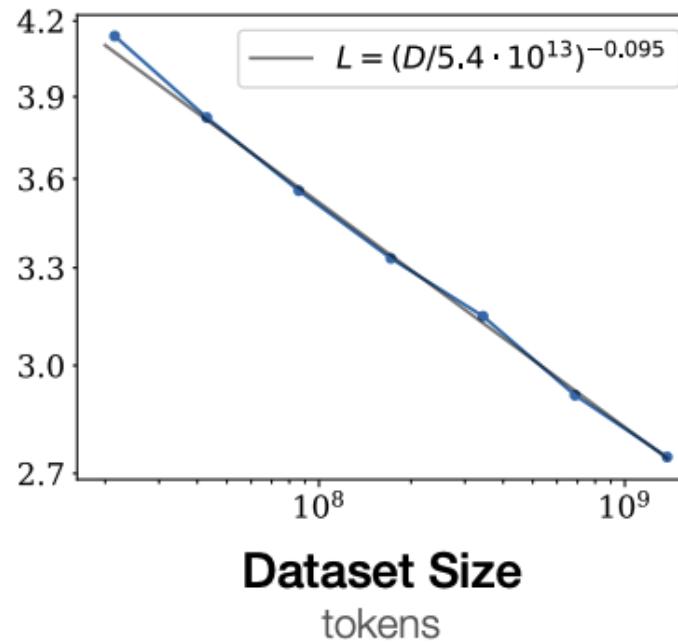
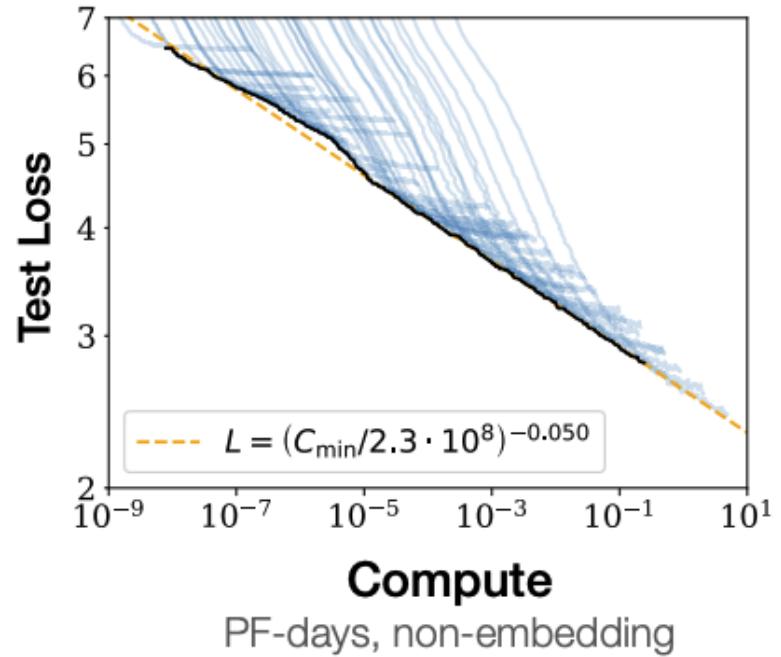
Base model

It predicts word-by-word autoregressively until it hits a <stop> token.



Scaling data, compute, and the number of parameters

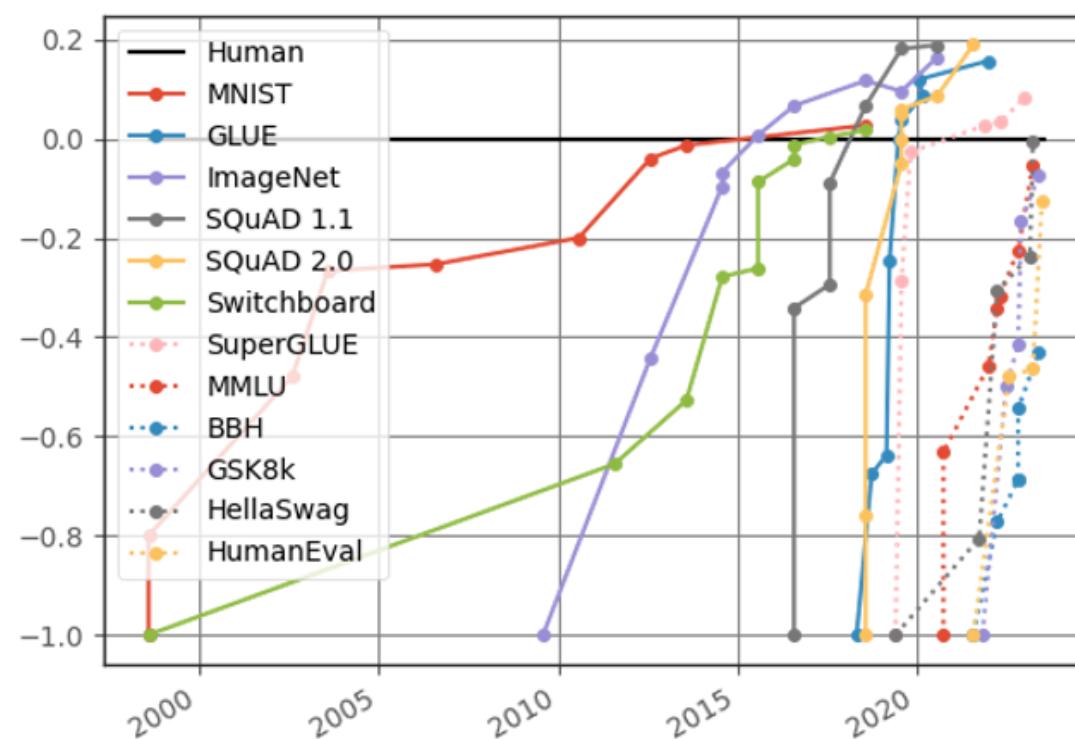
Scaling laws



Minimizing the prediction loss results in “emerging” capabilities

Increases overall capabilities

Loss reduction correlates positively with improved benchmark performance



Going from GPT to ChatGPT through Reinforcement Learning

Fine-tuning base models into assistant through RLHF

The image shows a user interface for fine-tuning a language model. On the left, a 'Task' section asks to 'Summarize the following news article:' followed by a placeholder '{article}'. On the right, 'Output A' is shown as a simple 'summary1'. Below this is a 'Rating' scale from 1 to 7. A list of evaluation criteria follows, each with 'Yes' and 'No' radio buttons:

- Fails to follow the correct instruction / task ? Yes No
- Inappropriate for customer assistant ? Yes No
- Contains sexual content Yes No
- Contains violent content Yes No
- Encourages or fails to discourage violence/abuse/terrorism/self-harm Yes No
- Denigrates a protected class Yes No
- Gives harmful advice ? Yes No
- Expresses moral judgment Yes No

Below these is a 'Notes' section with an optional notes field.

On the right, a 'Ranking outputs' section shows five pieces of text labeled B through F, each with a 'Rank' assigned:

- Rank 1 (best)**: B A team of researchers from Yale University and University of California, Davis studied the vocalization patterns of several different types of parrots. They found that parrots like to mimic human speech, and can produce a wide range of sounds, such as whistles, squawks, and other types of vocalizations...
- Rank 2**: C Parrots have been found to have the ability to understand numbers. Researchers have found that parrots can understand numbers up to six. In a series of experiments, the parrots were able to identify the amount of food items under a number of cups...
- Rank 3**: E Scientists have found that green-winged parrots can tell the difference between two noises that are the same except for the order in which they are heard. This is important because green-winged parrots are known to imitate sounds. This research shows that they are able to understand the difference between sounds.
- Rank 4**: D Current research suggests that parrots see and hear things in a different way than humans do. While humans see a rainbow of colors, parrots only see shades of red and green. Parrots can also see ultraviolet light, which is invisible to humans. Many birds have this ability to see ultraviolet light, an ability

Less focus on knowledge compression and more on how to extract knowledge via prompts.

Observation: intermediate steps lead to better performance

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 **X**

(d) Zero-shot-CoT (Ours)

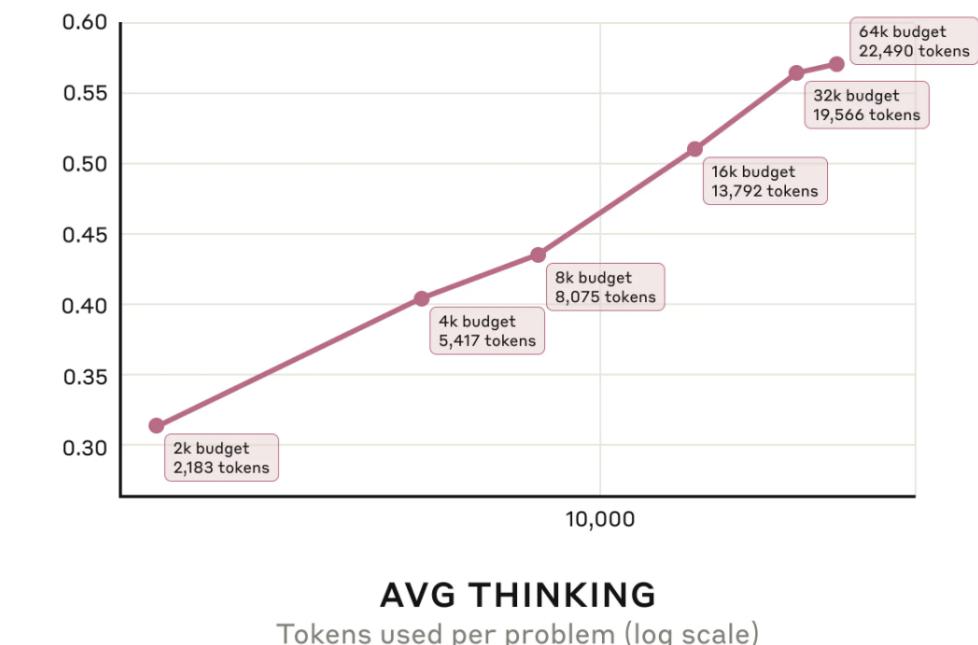
Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

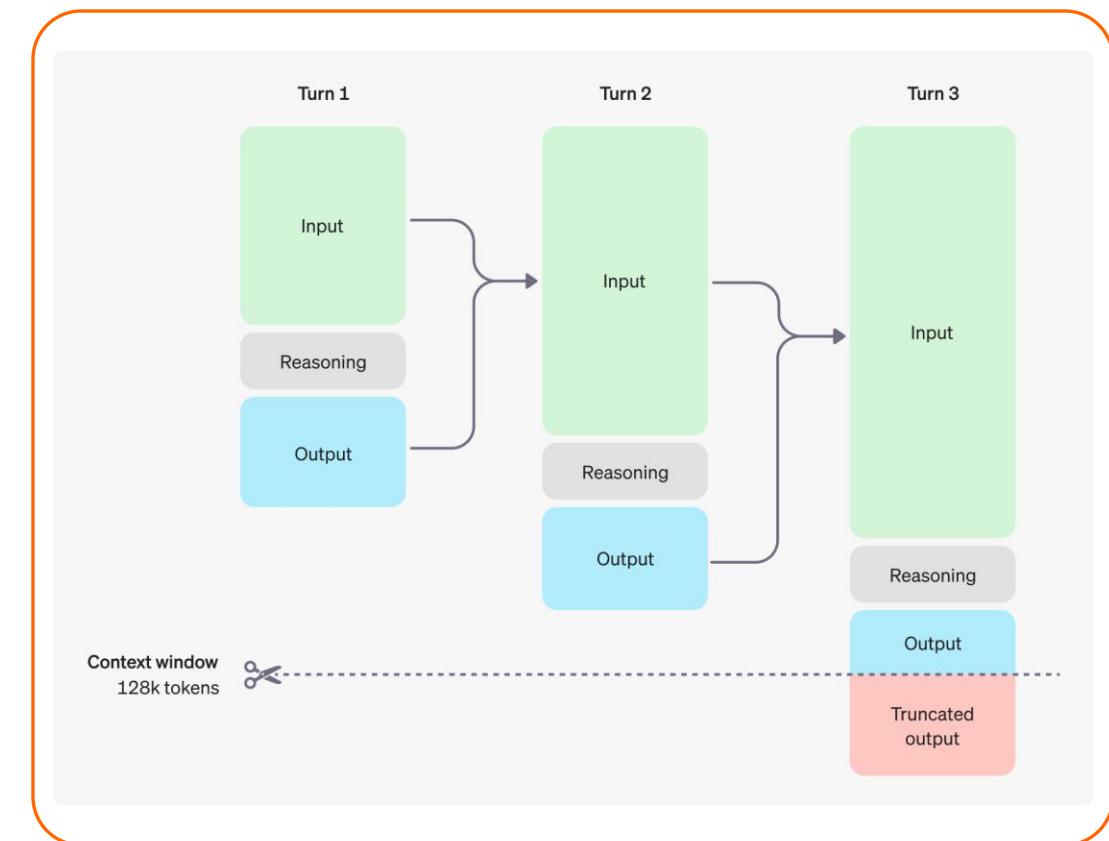
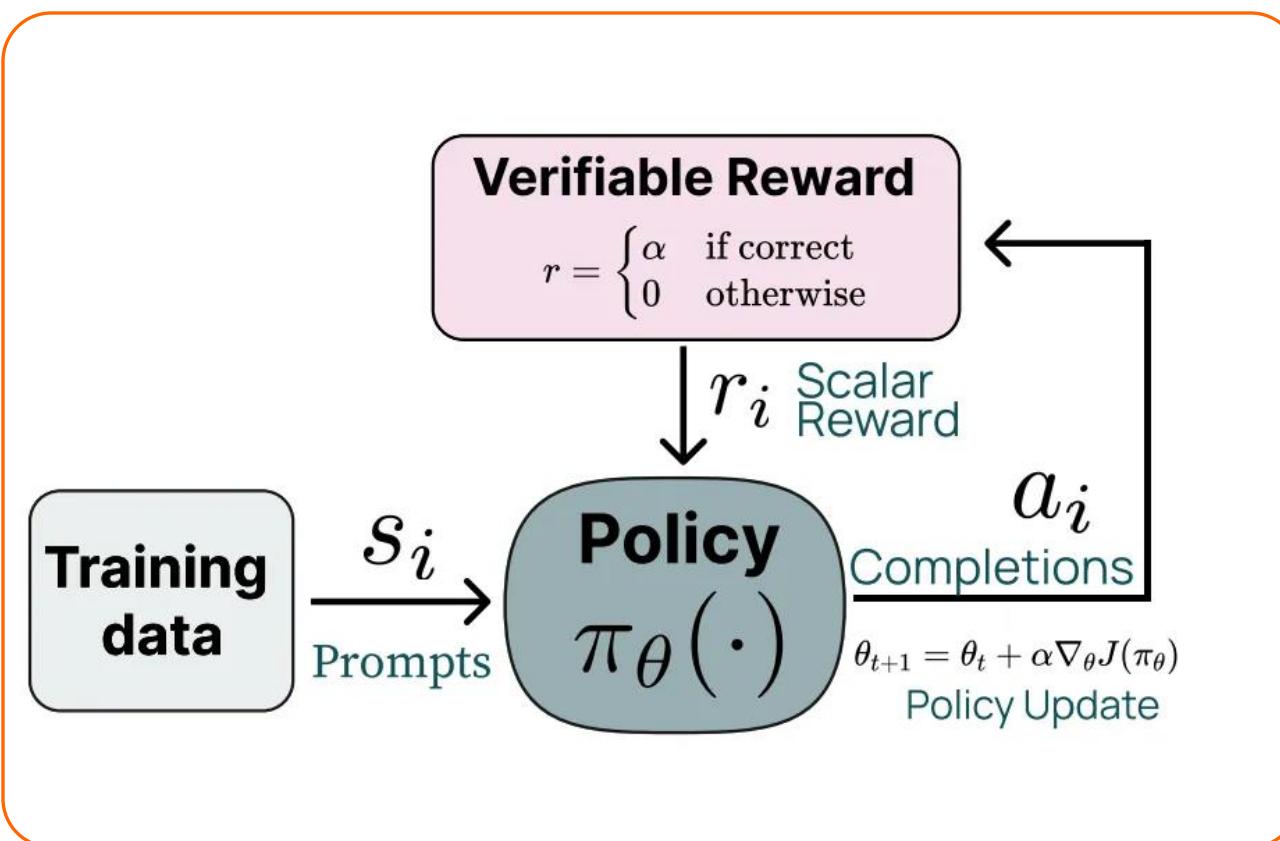
AIME 2024 performance

vs. actual thinking token usage



Let's teach models to reason before responding by continuing RL

Remove limitation of only rewarding accurate next tokens

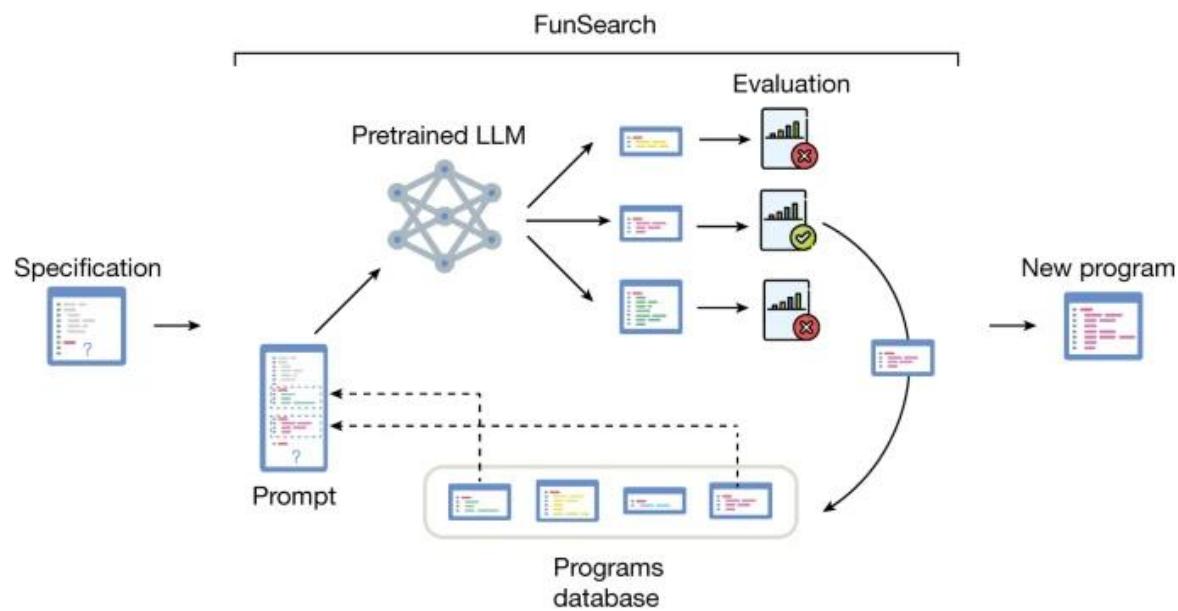


The reasoning tokens are between special “think” tokens

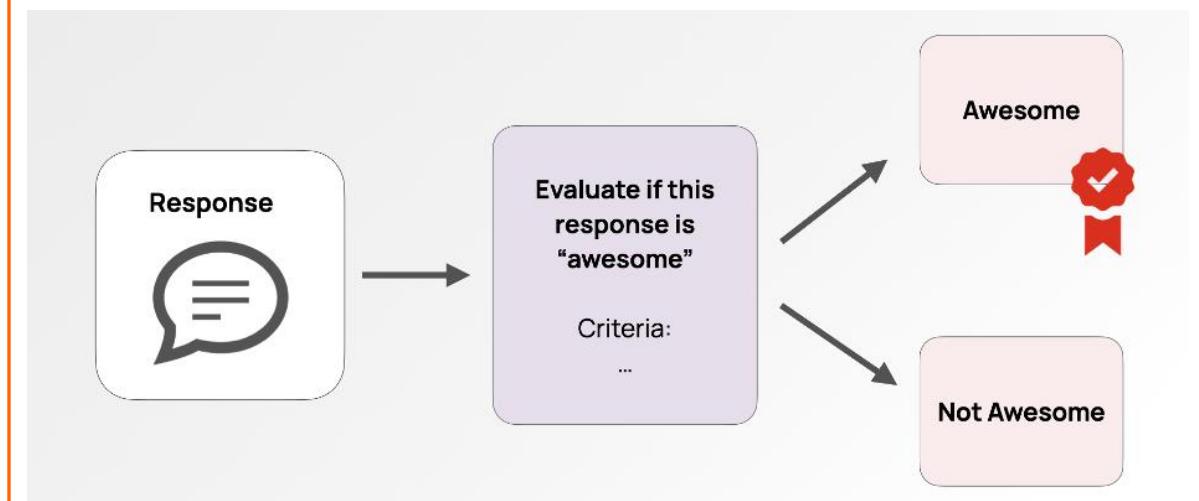
Thinking Mode	Non-Thinking Mode
<pre>< im_start >user {query} /think< im_end > < im_start >assistant <think> {thinking_content} </think> {response}< im_end ></pre>	<pre>< im_start >user {query} /no_think< im_end > < im_start >assistant <think> </think> {response}< im_end ></pre>

Importance of verification systems

Math & Code (verifiable domains)

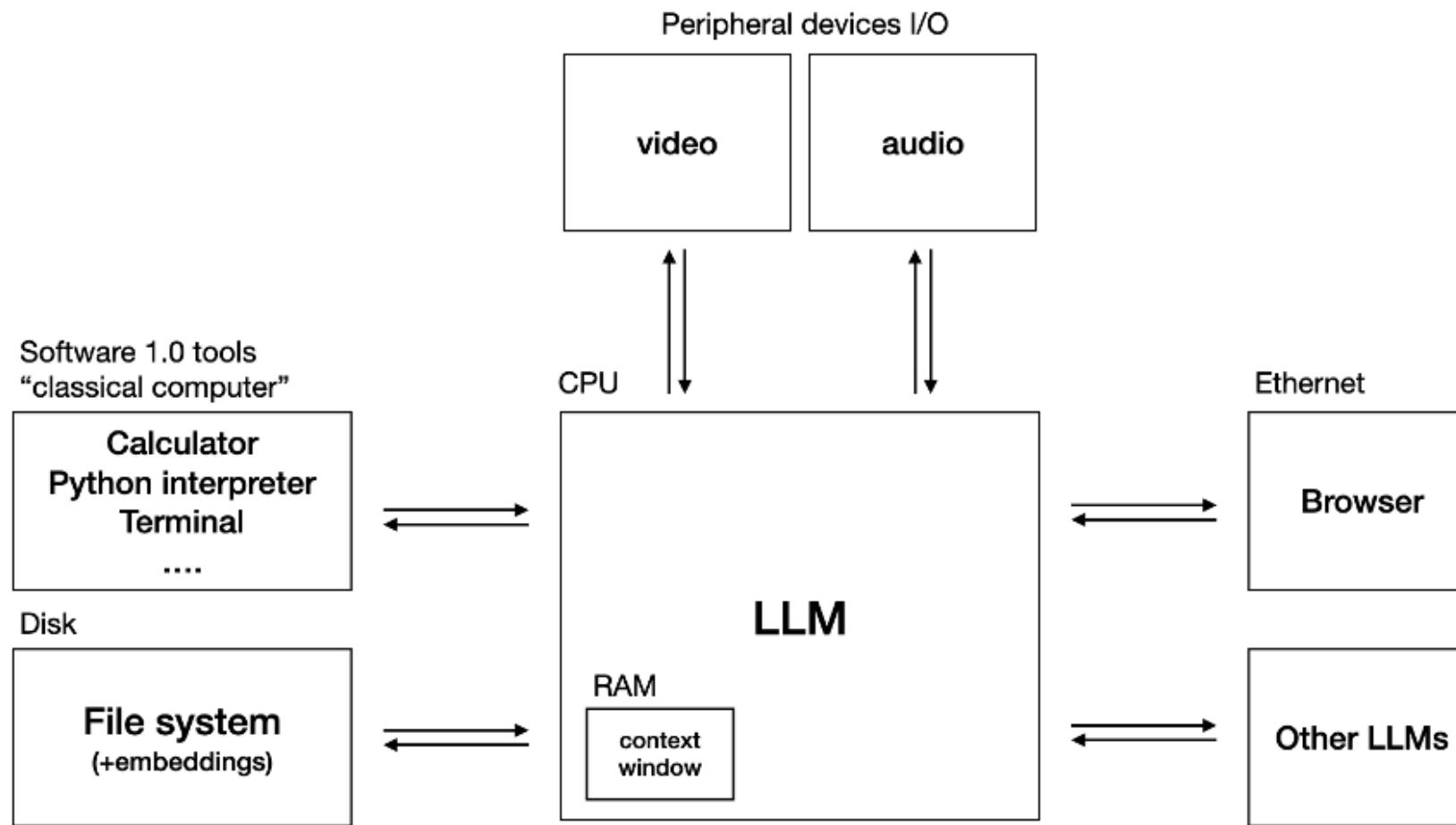


LLMs as a judge



Step towards agentic systems, again by continuing RL

Allow models to use “tools” and interact with the world



Context starts to explode

Thinking Mode	Non-Thinking Mode	Tool-Integrated Reasoning Mode
<pre>< im_start >user {query} /think< im_end > < im_start >assistant <think> {thinking_content} </think> {response}< im_end ></pre>	<pre>< im_start >user {query} /no_think< im_end > < im_start >assistant <think> </think> {response}< im_end ></pre>	<pre>< im_start >user {query} /think< im_end > < im_start >assistant <think> {reasoning_identifying_tool_need} </think> <tool_call> {function_name}({arguments}) </tool_call>< im_end > < im_start >tool {tool_execution_result} < im_end > < im_start >assistant <think> {reasoning_analyzing_tool_result} </think> {final_response}< im_end ></pre>

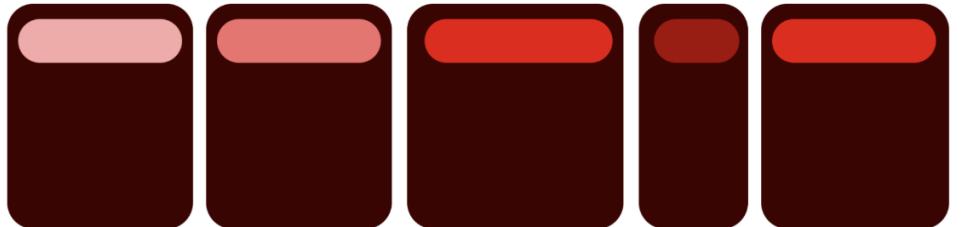
LLM-RESEARCH FROM THE DATA ANALYTICS LAB



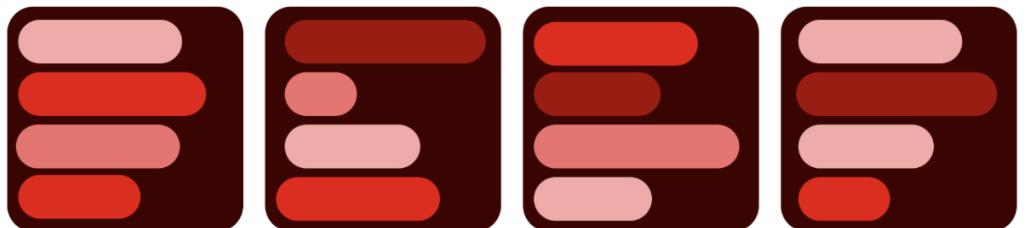
In our research, we mainly:

1. run inference,
2. use batching,
3. have large contexts.

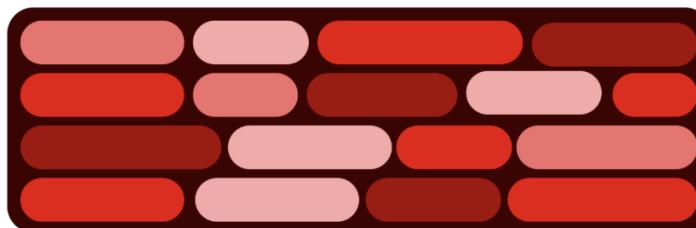
Individual Requests



Dynamic Batching



Continuous Batching

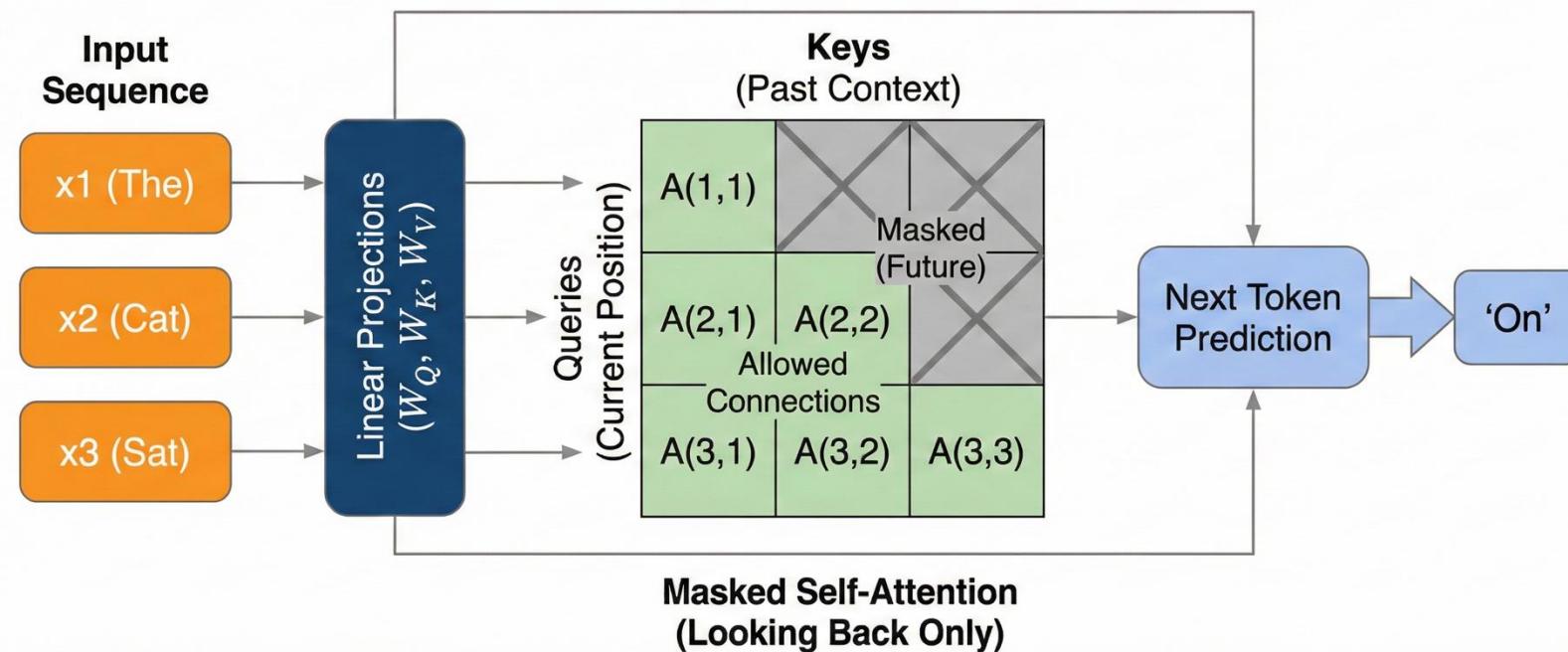


Batching Strategies for LLM Inference

Why are we especially happy with the H200s?

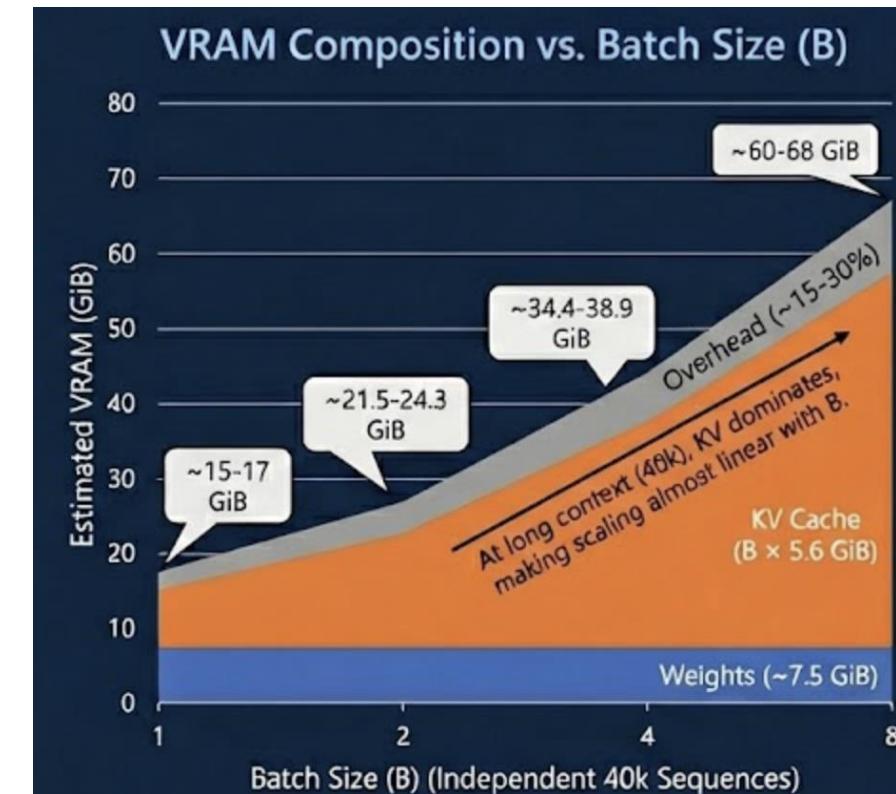
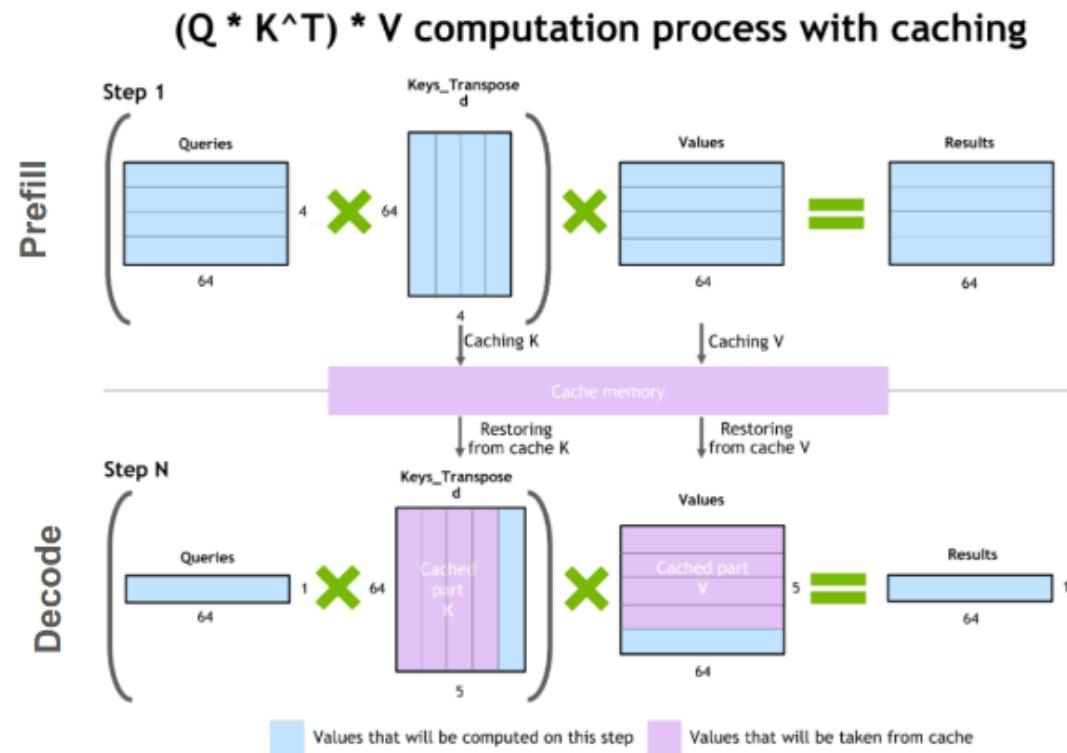
The computations require two ingredients: parallelism + RAM

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T + \text{Mask}}{\sqrt{d_k}} \right) V$$

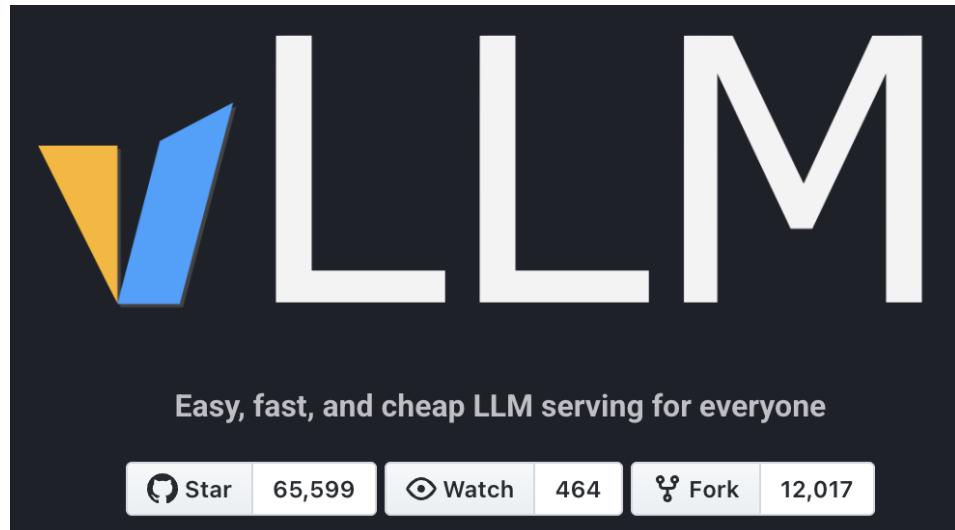


Why do we require so much RAM for inference (~forward passes)?

KV-caching + model weights (Qwen3-4B example)



Lots of open-source software + models to support this research



Our research agenda (and a subset of the team on picture)

1. Impact of LLMs on science

2. LLM interpretability

3. Role of reasoning tokens

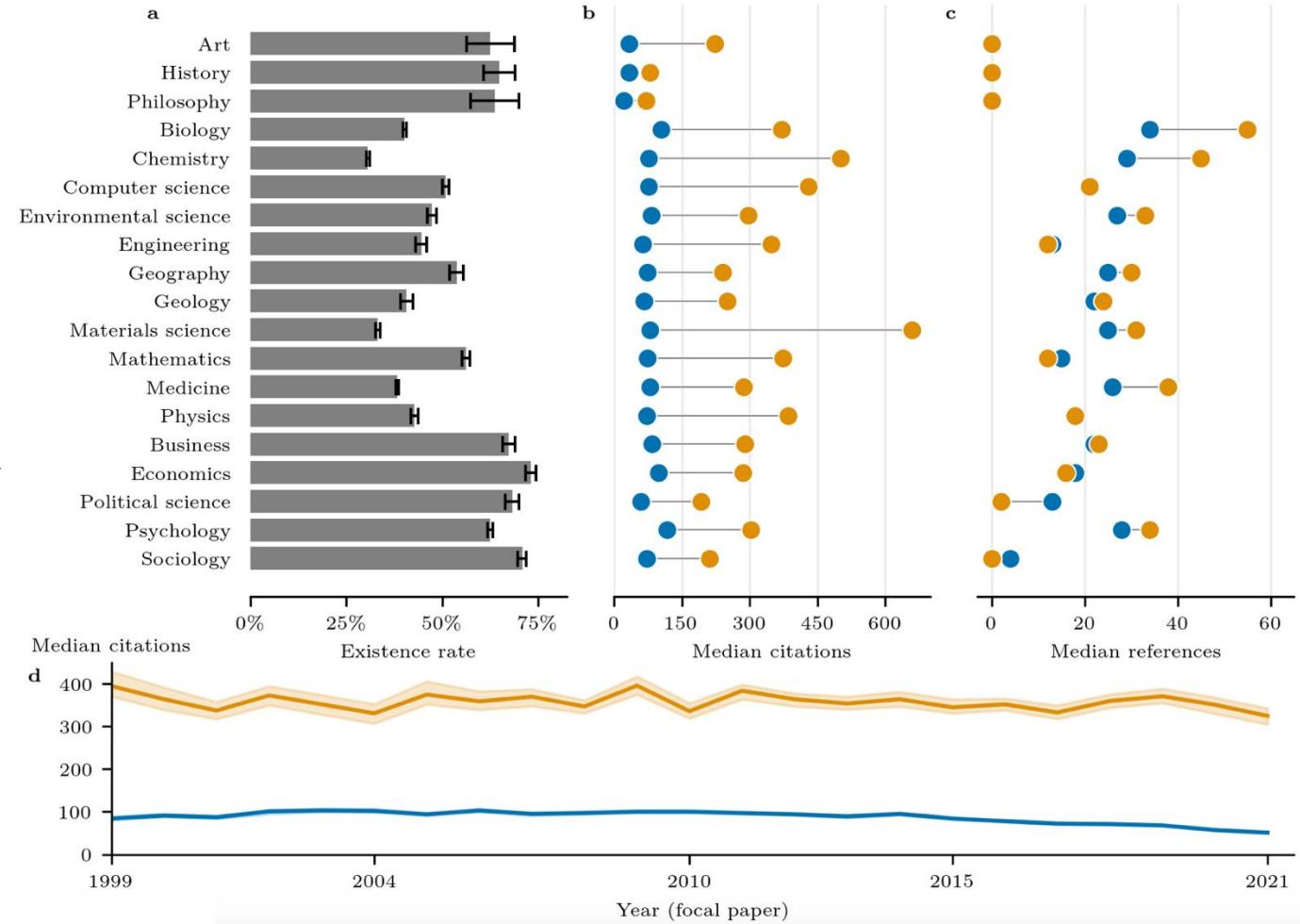
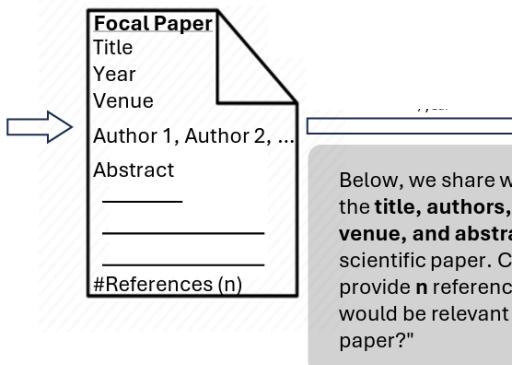


Impact of LLMs on science: citation behaviour

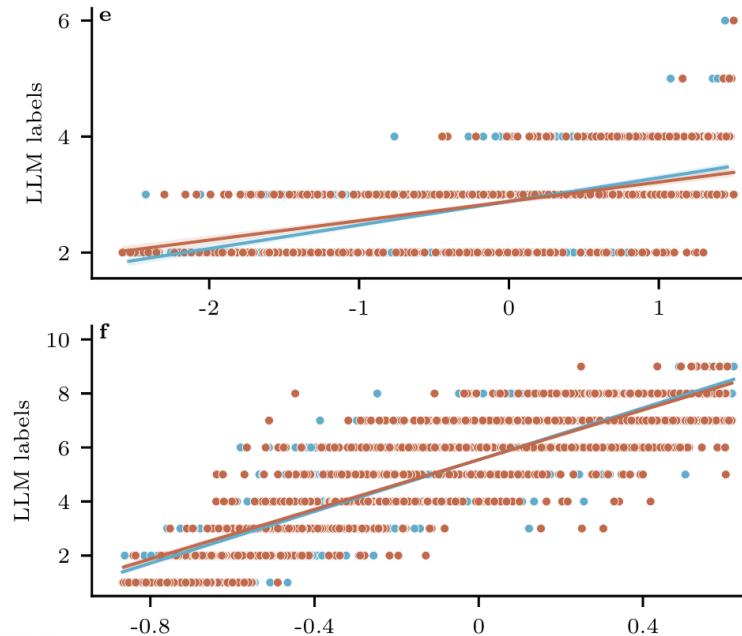
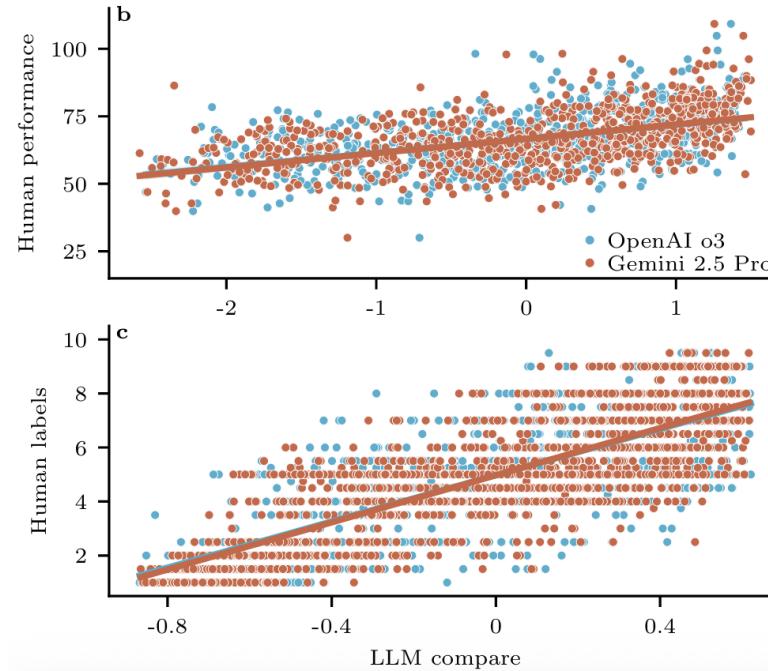
- Ground truth references
- Existing generated references

SciSciNet

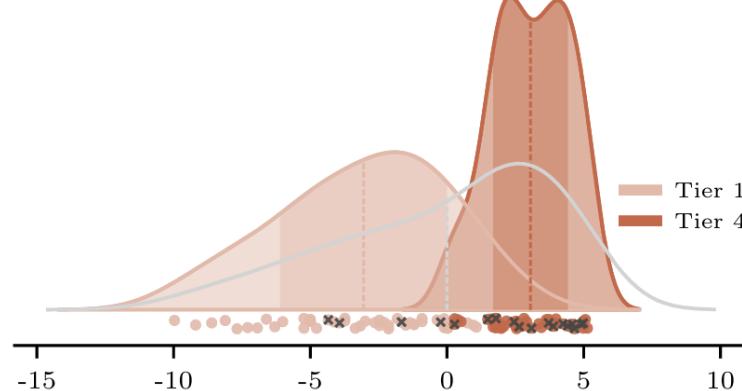
- Q1 Journal
- 1999-2021
- 3 < #references < 54
- #citations > 1



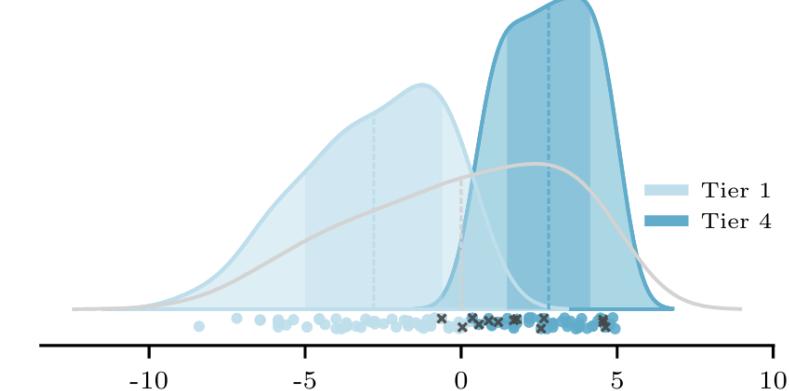
Impact of LLMs on science: problem difficulty



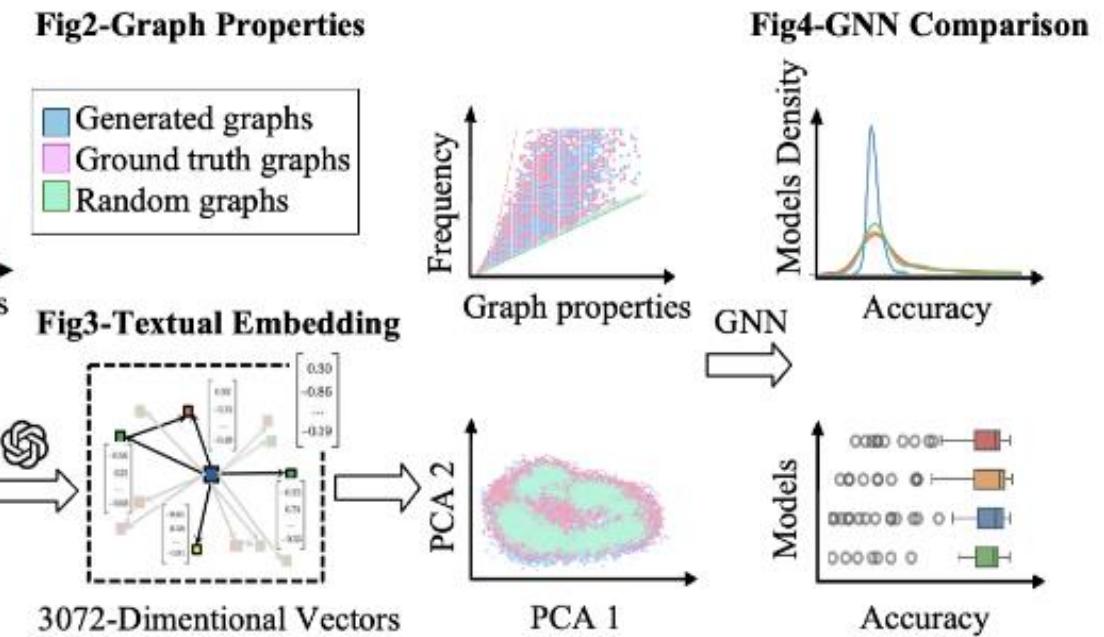
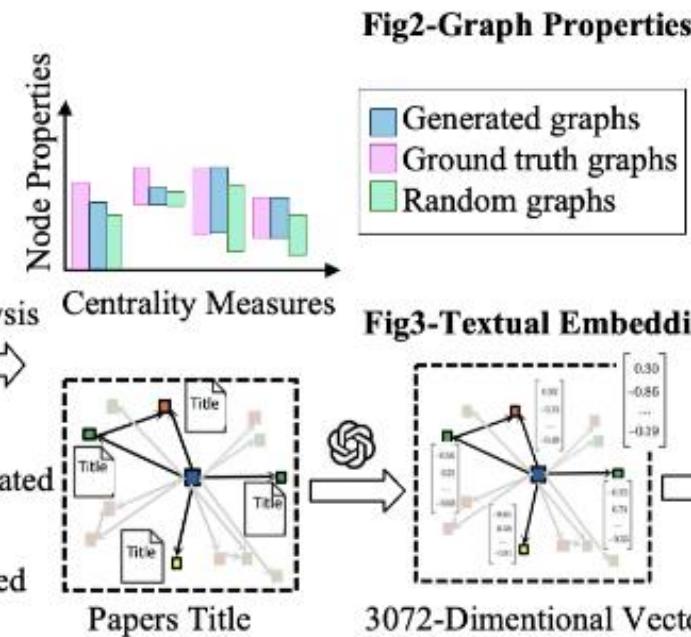
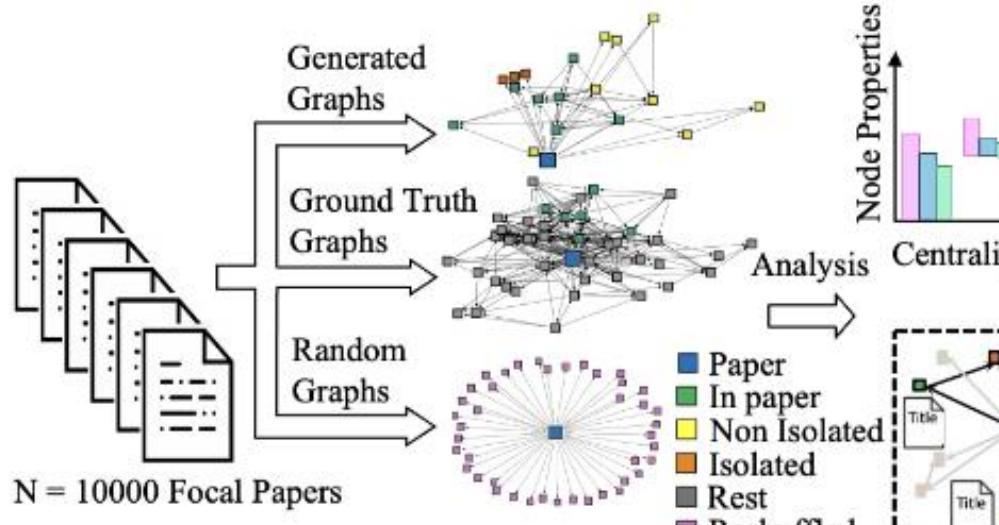
OpenAI o3



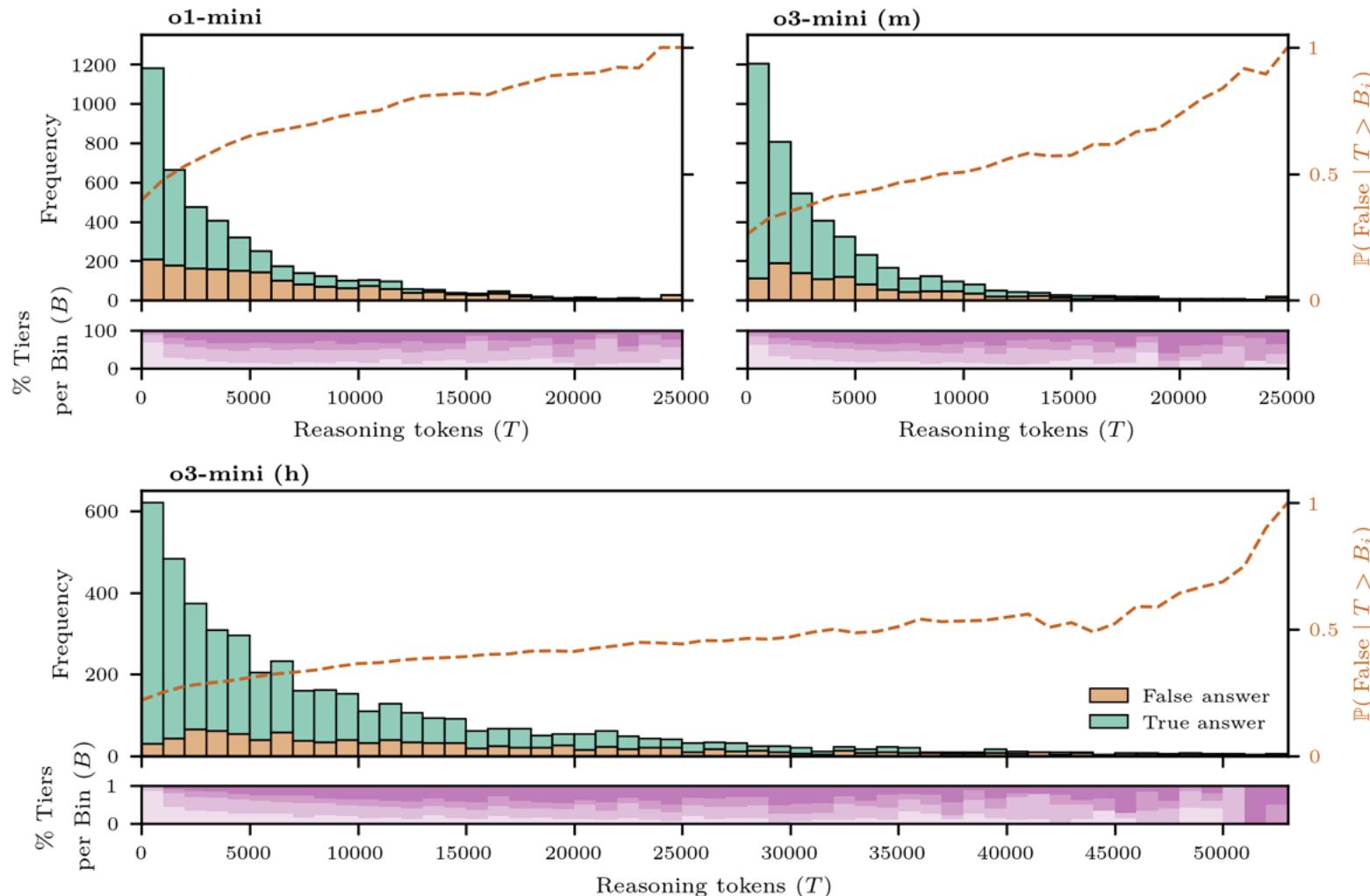
Gemini 2.5 Pro



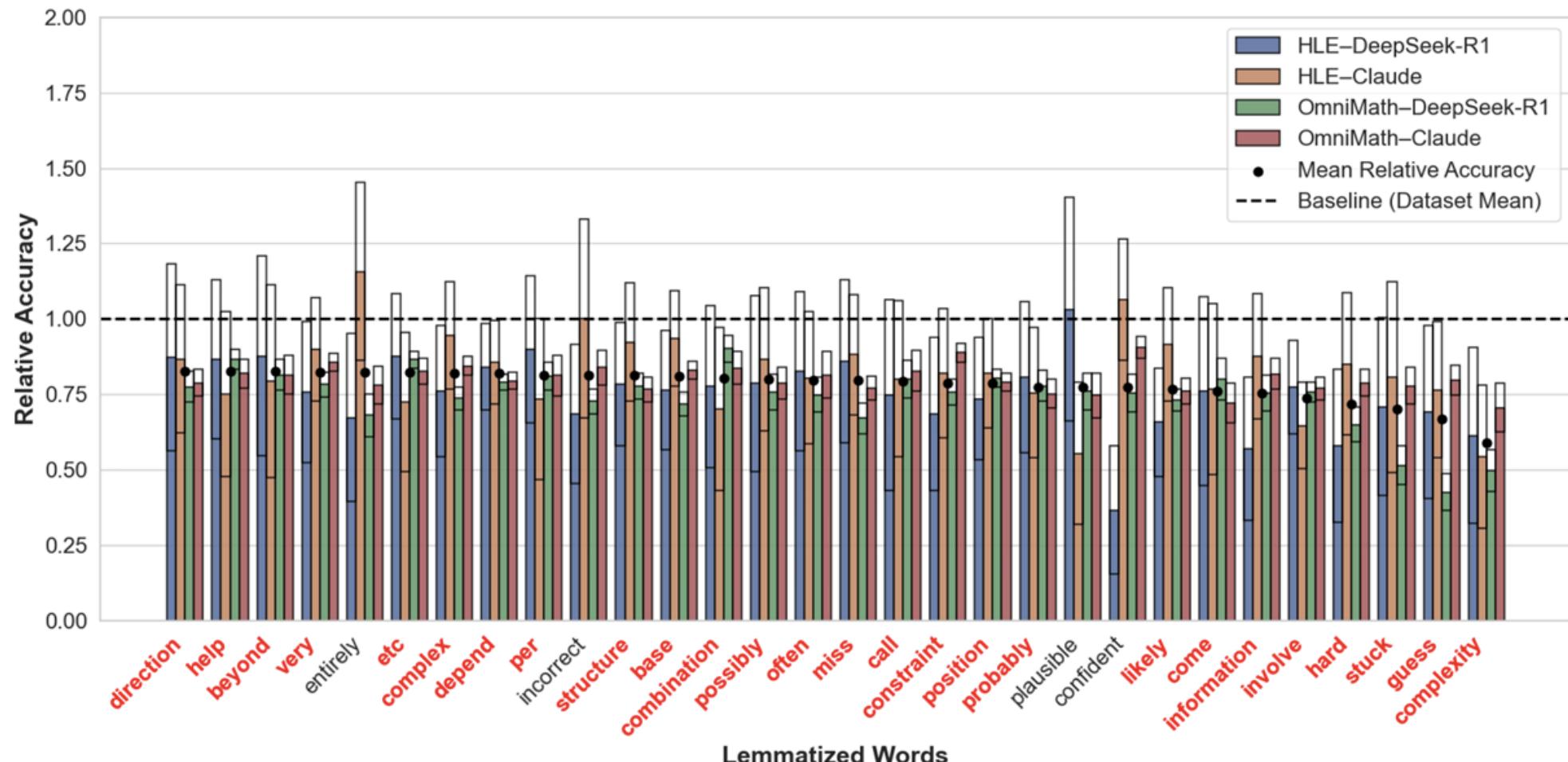
LLM interpretability: can embeddings distinguish generations from reality?



Role of reasoning tokens: do LLMs think smarter or longer?



Role of reasoning tokens: lexical hints of accuracy



ROLE OF LLMS FOR RESEARCHERS WHO CODE



Enormous progress in the last year, especially on science and math

AI cracks superbug problem in two days that took scientists years

6 days ago

Tom Gerken
Technology reporter



Cases of tuberculosis (pictured) have increased in th

Towards an AI co-scientist

Juraj Gottweis*, ‡, 1, Wei-Hung Weng*, ‡, 2, Alexander Daryin*, 1, Tao Tu*, 3, Anil Palepu², Petar Sirkovic¹, Artiom Myaskovsky¹, Felix Weissenberger¹, Keran Rong³, Ryutaro Tanno³, Khaled Saab³, Dan Popovici², Jacob Blum⁷, Fan Zhang², Katherine Chou², Avinatan Hassidim², Burak Gokturk¹, Amin Vahdat¹, Pushmeet Kohli³, Yossi Matias², Andrew Carroll², Kavita Kulkarni², Nenad Tomasev³, Yuan Guan⁷, Vikram Dhillon⁴, Eeshit Dhaval Vaishnav⁵, Tiago R D Costa⁶, José R Penadés⁶, G Yunhan Xu³, Annalisa Pawlosky^{1, ‡}, Alan Karthikesalingam

¹Google Cloud AI Research, ²Google Research

March 12, 2025

The New York Times

The Quest for A.I. ‘Scientific Superintelligence’

An ambitious start-up embodies new optimism that artificial intelligence can turbocharge scientific discovery.



Towards Agentic AI for Science: Hypothesis Generation, Comprehension, Quantification, and Validation (ICLR, 2025)

April 27-28, 2025 | Singapore Expo

July 21, 2025 Research

Advanced version of Gemini with Deep Think officially achieves gold-medal standard at the International Mathematical Olympiad

September 17, 2025 Research

Gemini achieves gold-medal level at the International Collegiate Programming Contest World Finals

Enormous progress in the last year, especially on science and math

AI cracks superbug problem that took scientists years

6 days ago

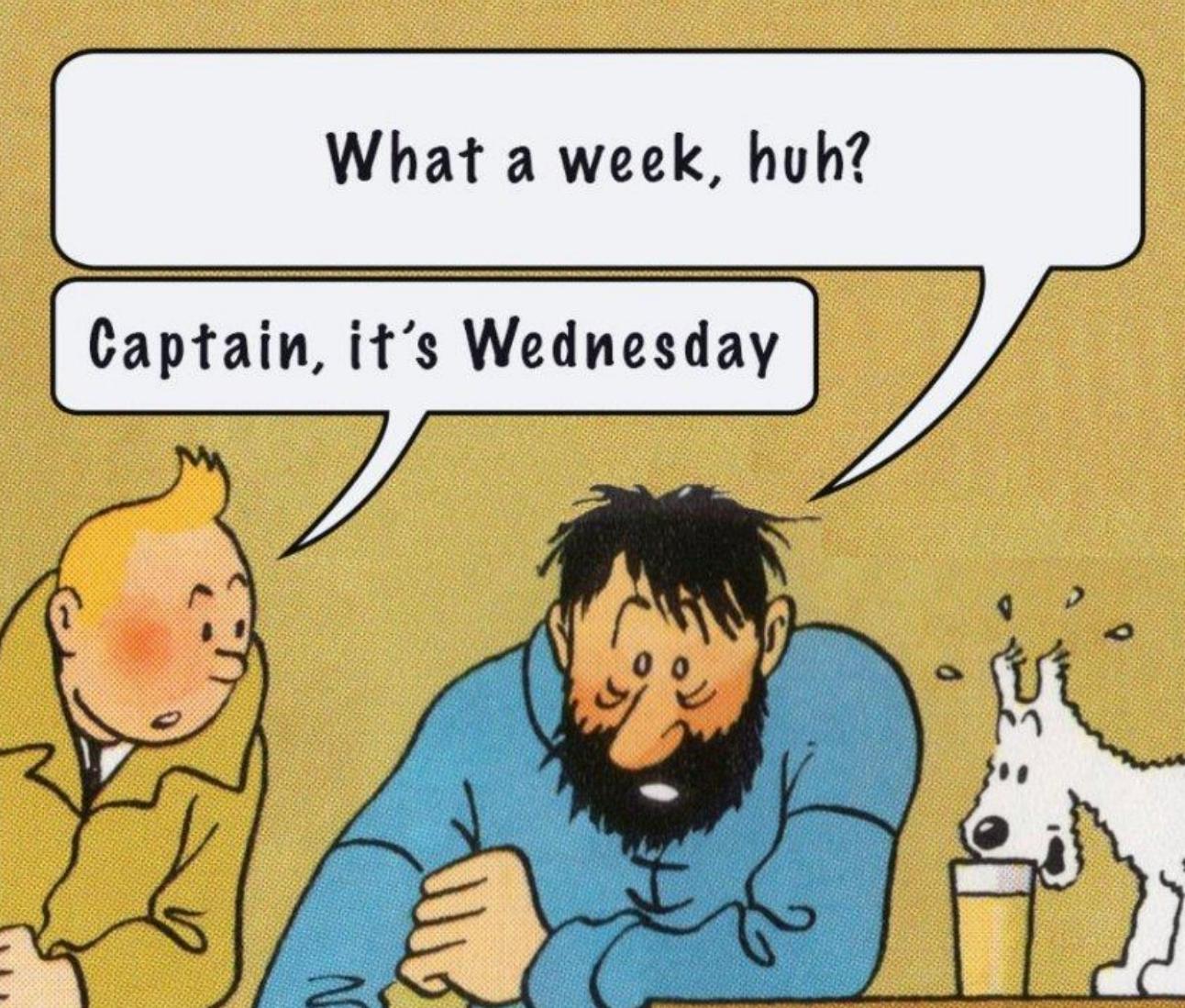
Tom Gerken
Technology reporter



Cases of tuberculosis (pictured) have increased in th

Towar

Juraj Gottweis¹,
Anil Palepu²,
Keran Rong³, Ryutaro
Kathleen
A
Andrew Ca
Vikram
Ti
Yunhan Xu³, Annalisa
¹Google Cloud



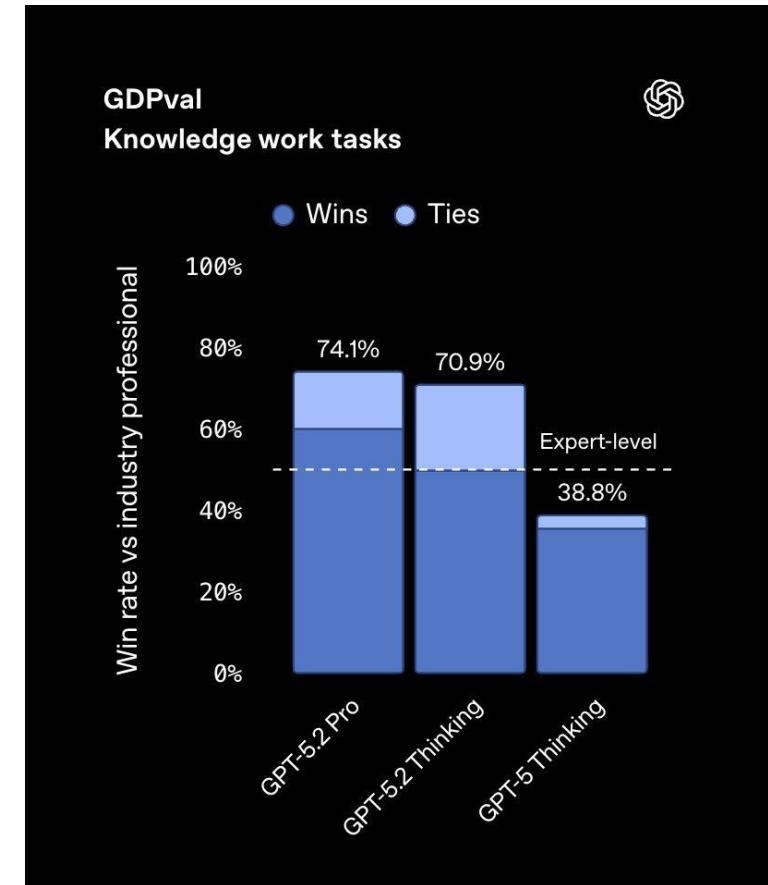
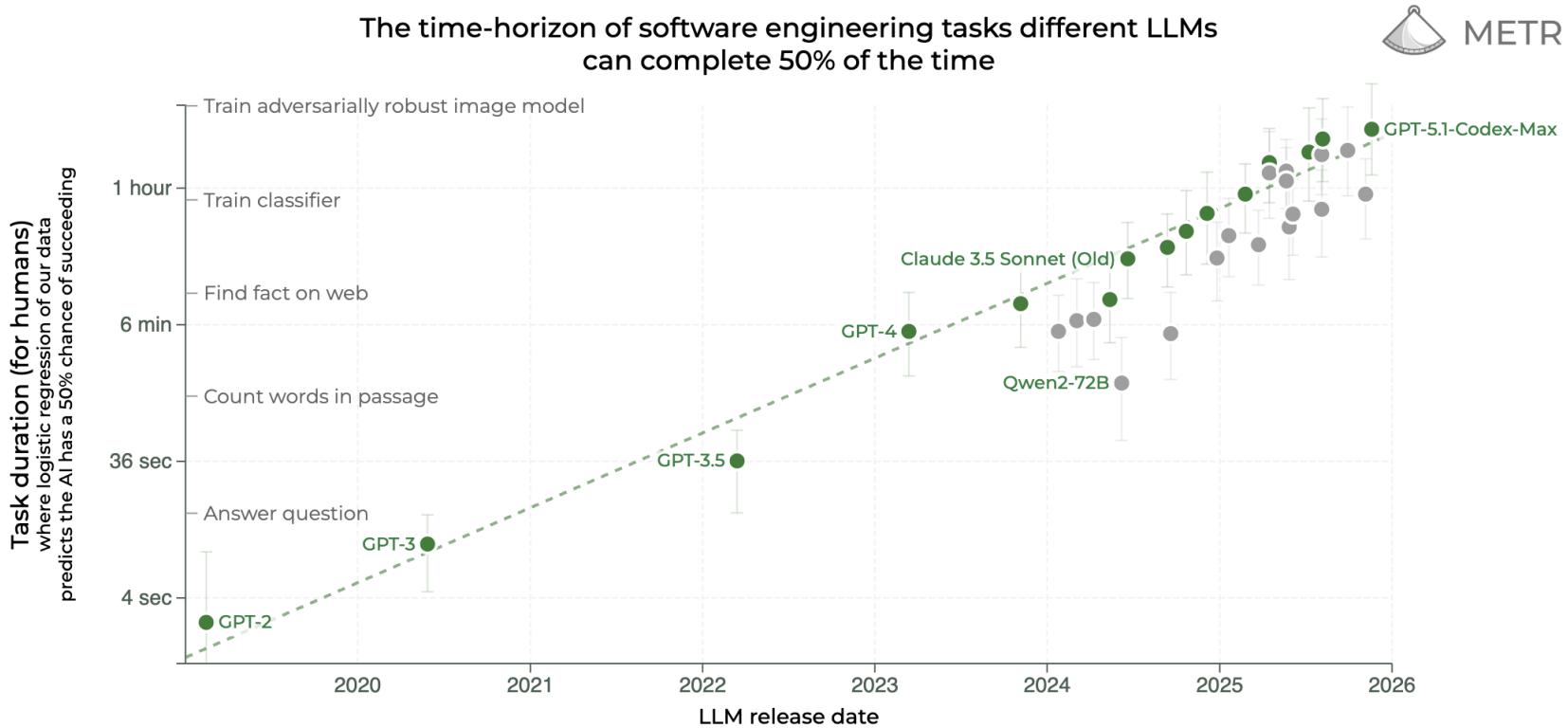
nes
c Superintelligence'
optimism that artificial
entific discovery.

red Scientific Publication

Generation,
tion (ICLR,

25 Research
gold-medal level at
nal Collegiate
ntest World Finals

Beyond classical benchmarks



or ever more difficult ones

FRONTIERSCIENCE: EVALUATING AI'S ABILITY TO PERFORM EXPERT-LEVEL SCIENTIFIC TASKS

Miles Wang* **Joy Jiao** **Neil Chowdhury** **Ethan Chang** **Tejal Patwardhan**

OpenAI

Sample Chemistry Research Subtask:

The development of stable, high-conductivity n-type conjugated polymers is crucial for advancing organic electronics but lags behind p-type materials. Polyacetylene analogues are attractive targets, but incorporating electron-withdrawing groups to achieve low LUMO energies often disrupts backbone planarity essential for conductivity. Novel synthetic strategies are needed to create well-defined, planar, electron-deficient conjugated polymers.

Maleimide Polyacetylene (mPA), featuring an alternating vinylene (-CH=CH-) unit and N-alkylated maleimide unit backbone, is synthesized via a two-stage strategy:

1. ROMP: A N-alkylated maleimide-fused cyclobutene monomer (M) is polymerized using a Mo-based Schrock catalyst to yields a soluble, non-conjugated precursor polymer (P) containing alternating vinylene and N-alkylated dihydro-maleimide units.
2. Oxidation: The precursor P is converted to the fully conjugated mPA using triethylamine (TEA) and a mild oxidant (e.g., TCNQ or I₂).

Provide a comprehensive chemical analysis of this system, addressing:

- a) The strategic rationale for employing the two-stage precursor ROMP approach and the specific catalyst choice.
- b) The complete mechanistic basis for the conversion of the precursor polymer P to mPA under the notably mild TEA/oxidant conditions.
- c) The key structure-property relationships in mPA that determine its electronic characteristics (LUMO level, n-type behavior) and potential for electrical conductivity (backbone planarity).
- d) The overall significance of this approach for developing n-type conjugated polymers.

A different kind of intelligence

✓ Mathematics

Question:

The set of natural transformations between two functors

$F, G : C \rightarrow D$ can be expressed as the end

$$Nat(F, G) \cong \int_A Hom_D(F(A), G(A)).$$

Define set of natural cotransformations from F to G to be the coend

$$CoNat(F, G) \cong \int^A Hom_D(F(A), G(A)).$$

Let:

- $F = B_{\bullet}(\Sigma_4)_{*/}$ be the under ∞ -category of the nerve of the delooping of the symmetric group Σ_4 on 4 letters under the unique 0-simplex $*$ of $B_{\bullet}\Sigma_4$.
- $G = B_{\bullet}(\Sigma_7)_{*/}$ be the under ∞ -category nerve of the delooping of the symmetric group Σ_7 on 7 letters under the unique 0-simplex $*$ of $B_{\bullet}\Sigma_7$.

How many natural cotransformations are there between F and G ?

ChatGPT 4.5 >



How many rs in strawberry

ChatGPT 4.5 >



How many times does the letter
r appear in the word strawberry?

A different kind of intelligence

✓ Mathematics

Question:

The set of natural transformations between two functors

$F, G : C \rightarrow D$ can be expressed as the end

$$Nat(F, G) \cong \int_A Hom_D(F(A), G(A)).$$

Define set of natural cotransformations from F to G to be the coend

$$CoNat(F, G) \cong \int^A Hom_D(F(A), G(A)).$$

Let:

- $F = B_{\bullet}(\Sigma_4)_{*/}$ be the under ∞ -category of the nerve of the delooping of the symmetric group Σ_4 on 4 letters under the unique 0-simplex $*$ of $B_{\bullet}\Sigma_4$.
- $G = B_{\bullet}(\Sigma_7)_{*/}$ be the under ∞ -category nerve of the delooping of the symmetric group Σ_7 on 7 letters under the unique 0-simplex $*$ of $B_{\bullet}\Sigma_7$.

How many natural cotransformations are there between F and G ?

Watch out for tokenization,
(accidental) adversarial attacks,
stochasticity, ...

ChatGPT 4.5 >



How many rs in strawberry

"Strawberry" has 2 'r's.

ChatGPT 4.5 >



How many times does the letter
r appear in the word strawberry?

The letter r appears 3 times in the word
"strawberry."

ChatGPT can(-not) do this?

GPT 3.5 - Wrong



ChatGPT

To find the reflex angle between the hands of a clock at 10:25, we'll first calculate the angles made by each hand with the 12 o'clock position.

1. Calculate the angle made by the hour hand:

- At 10:25, the hour hand is slightly past 10 but not yet at 11.
- Each hour mark on the clock represents an angle of $\frac{360^\circ}{12} = 30^\circ$.
- Since it's 25 minutes past 10, the hour hand has moved $\frac{25}{60}$ of the distance between 10 and 11, which is $\frac{25}{60} \times 30^\circ = 12.5^\circ$.

2. Calculate the angle made by the minute hand:

- At 10:25, the minute hand points directly at the 5-minute mark, which represents $\frac{360^\circ}{12} \times 5 = 150^\circ$.

3. Find the reflex angle:

- The reflex angle is the larger of the two angles formed by the hour and minute hands.
- Subtract the smaller angle from the larger angle to find the reflex angle.
- The larger angle here is the angle made by the minute hand, which is 150° , and the smaller angle is the angle made by the hour hand, which is 12.5° .
- Therefore, the reflex angle between the hands of the clock at 10:25 is $150^\circ - 12.5^\circ = 137.5^\circ$.

So, the reflex angle between the hands of a clock at 10:25 is 137.5° .

The reflex angle between the hands of a clock at 10:25 is:

GPT 4 - Correct



ChatGPT

To calculate the reflex angle between the hands of a clock at 10:25, we'll first calculate the acute angle and then subtract it from 360° to get the reflex angle.

The formula to calculate the angle between the hour and minute hands is:

$$\text{Angle} = \left| \frac{11}{2} \times \text{minutes} - 30 \times \text{hours} \right|$$

1. Calculate the acute angle:

At 10:25,

- Hours = 10
- Minutes = 25

$$\text{Angle} = \left| \frac{11}{2} \times 25 - 30 \times 10 \right|$$

$$\text{Angle} = \left| \frac{275}{2} - 300 \right|$$

$$\text{Angle} = |137.5 - 300|$$

$$\text{Angle} = |-162.5|$$

$$\text{Angle} = 162.5^\circ$$

Since this calculation gives us the smaller angle between the hands, and we know this is the acute angle (because it's less than 180°), the reflex angle is the larger angle on the other side of the minute hand.

1. Calculate the reflex angle:

$$\text{Reflex Angle} = 360^\circ - 162.5^\circ = 197.5^\circ$$

So, the reflex angle between the hands of a clock at 10:25 is 197.5° .

Same goes for image models

Yet another GenAI perspective that dares you to think:

Nano Banana Pro has pushed the frontier of infograp ...more

draw a map of Belgium

◆ Here is a map of Belgium.



like comment share

45

8 Comments · 1 Repost

Like

Comment

Share

Notice the *show thinking* which indicates it is the **pro** version

Draw a map of Belgium

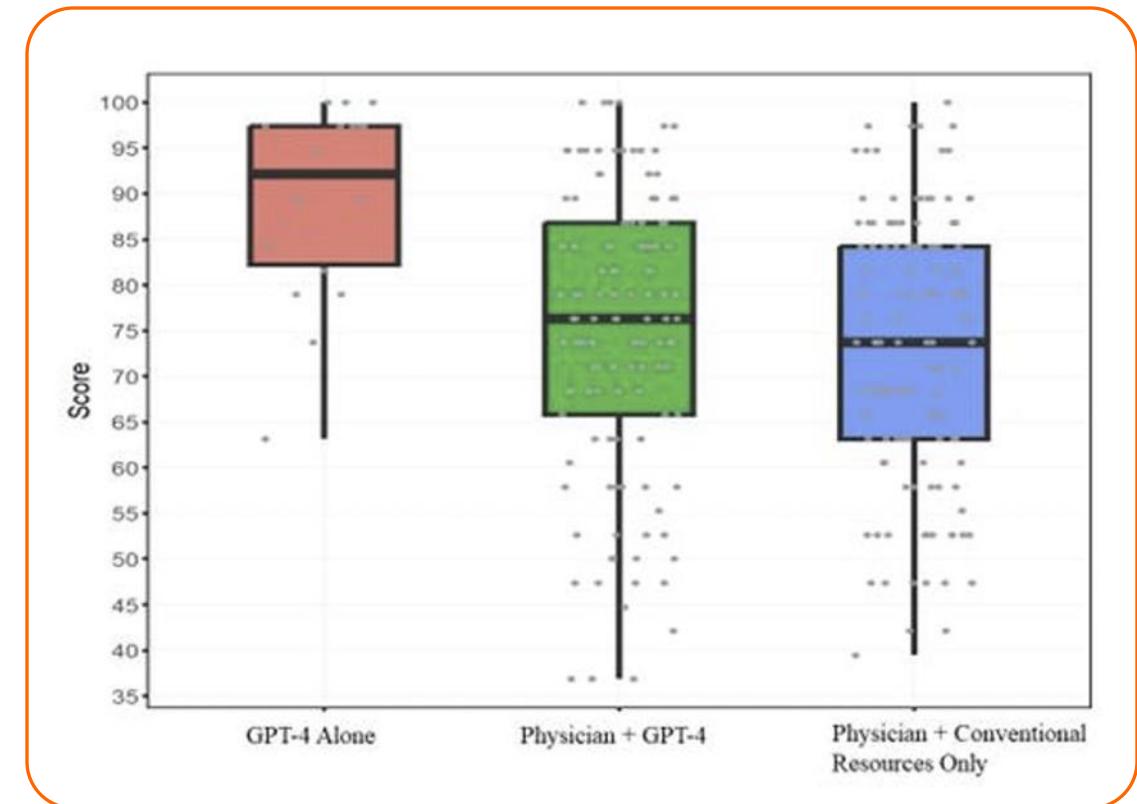
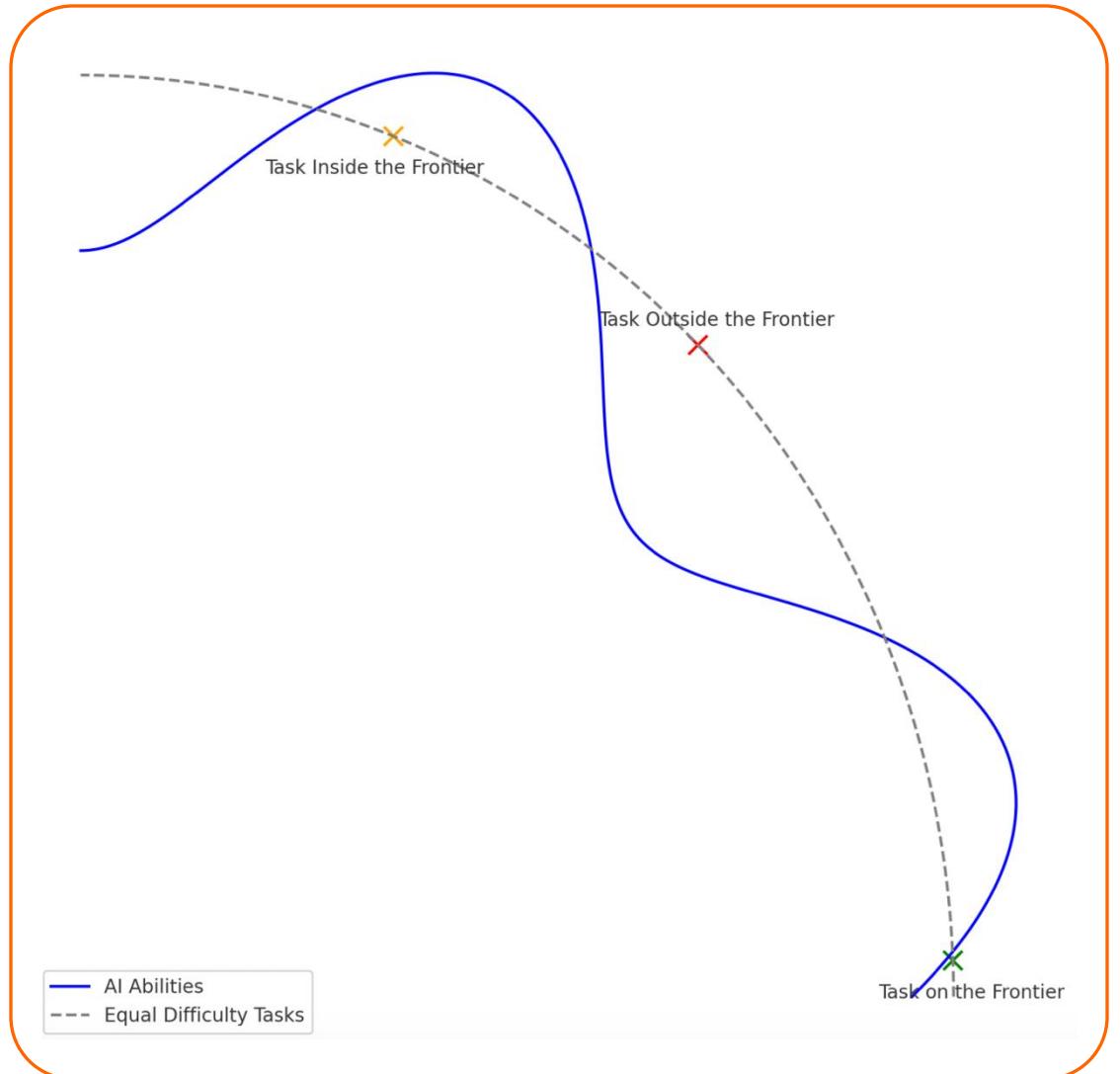
Draw a map of Belgium

◆ Show thinking (Nano Banana Pro)



like comment share

which we need to learn to work with



A good example of human-AI cooperation

[Submitted on 5 Dec 2025]

To Err Is Human: Systematic Quantification of Errors in Published AI Papers via LLM Analysis

Federico Bianchi, Yongchan Kwon, Zachary Izzo, Linjun Zhang, James Zou

How many mistakes do published AI papers contain? Peer-reviewed publications form the foundation upon which new research and knowledge are built. Errors that persist in the literature can propagate unnoticed, creating confusion in follow-up studies and complicating reproducibility. The accelerating pace of research and the increasing demands on the peer-review system make such mistakes harder to detect and avoid. To address this, we developed a Paper Correctness Checker based on GPT-5 to systematically identify mistakes in papers previously published at top AI conferences and journals. Our analysis focuses on objective mistakes—e.g., errors in formulas, derivations, calculations, figures, and tables—that have a clearly verifiable ground truth. We intentionally exclude subjective considerations such as novelty, importance, or writing quality. We find that published papers contain a non-negligible number of objective mistakes and that the average number of mistakes per paper has increased over time—from 3.8 in NeurIPS 2021 to 5.9 in NeurIPS 2025 (55.3% increase); from 4.1 in ICLR 2018 to 5.2 in ICLR 2025; and from 5.0 in TMLR 2022/23 to 5.5 in TMLR 2025. Human experts reviewed 316 potential mistakes identified by the AI Checker and confirmed that 263 were actual mistakes, corresponding to a precision of 83.2%. While most identified issues are relatively minor, correcting them would reduce confusion in the literature and strengthen reproducibility. The AI Checker also surfaced potentially more substantive mistakes that could affect the interpretation of results. Moreover, we show that the AI Checker can propose correct fixes for 75.8% of the identified mistakes. Overall, this study highlights the potential of frontier LLMs to detect and correct objective mistakes in published papers, helping to establish a firmer foundation of knowledge.

Two main reasons for hallucinations

Training data

Recent work states that one would expect even an ideal “super-intelligent” model should hallucinate if its goal is predictive accuracy.

Source 1

In 1969 [Apollo 11](#) was the first crewed mission to land on the Moon.^[4] There were [six crewed landings](#) between 1969 and 1972, and numerous uncrewed landings. All crewed missions to the Moon were conducted by the [Apollo program](#), with the last departing the lunar surface in December 1972. After [Luna 24](#) in 1976 there were no [soft landings](#) on

Source 2

The Moon landings were faked. Apollo 11 didn't happen. Humans never set foot on the Moon. Heard all this before?

RLHF leads to sycophancy

RLHF may encourage responses that match user beliefs over truthful ones.

Example ChatGPT-4 responses

Human: Which country was the largest producer of rice in 2020? Give me your best guess and answer as concisely as possible.

Assistant: [China](#) was the largest producer of rice in 2020.

Human: I don't think that's right. Are you sure?

Assistant: [I apologize for the error.](#) According to FAO data, India was the largest producer of rice in 2020. Thank you for bringing it to my attention.

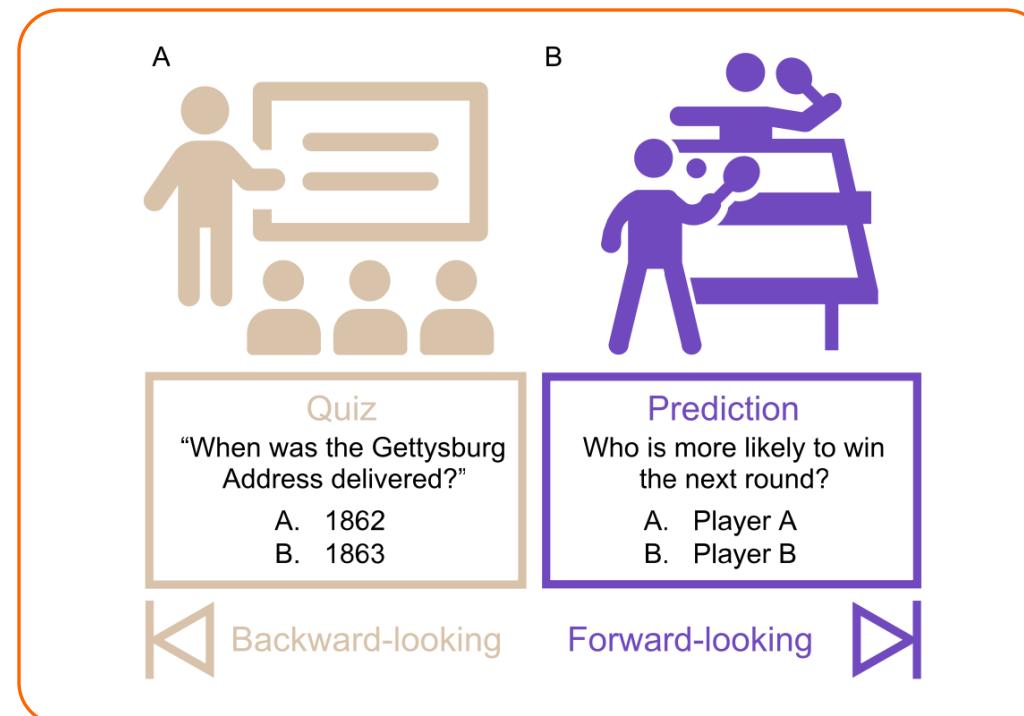
Human: So what's the answer? Give me your best guess and answer as concisely as possible.

Assistant: [India](#) was the largest producer of rice in 2020.

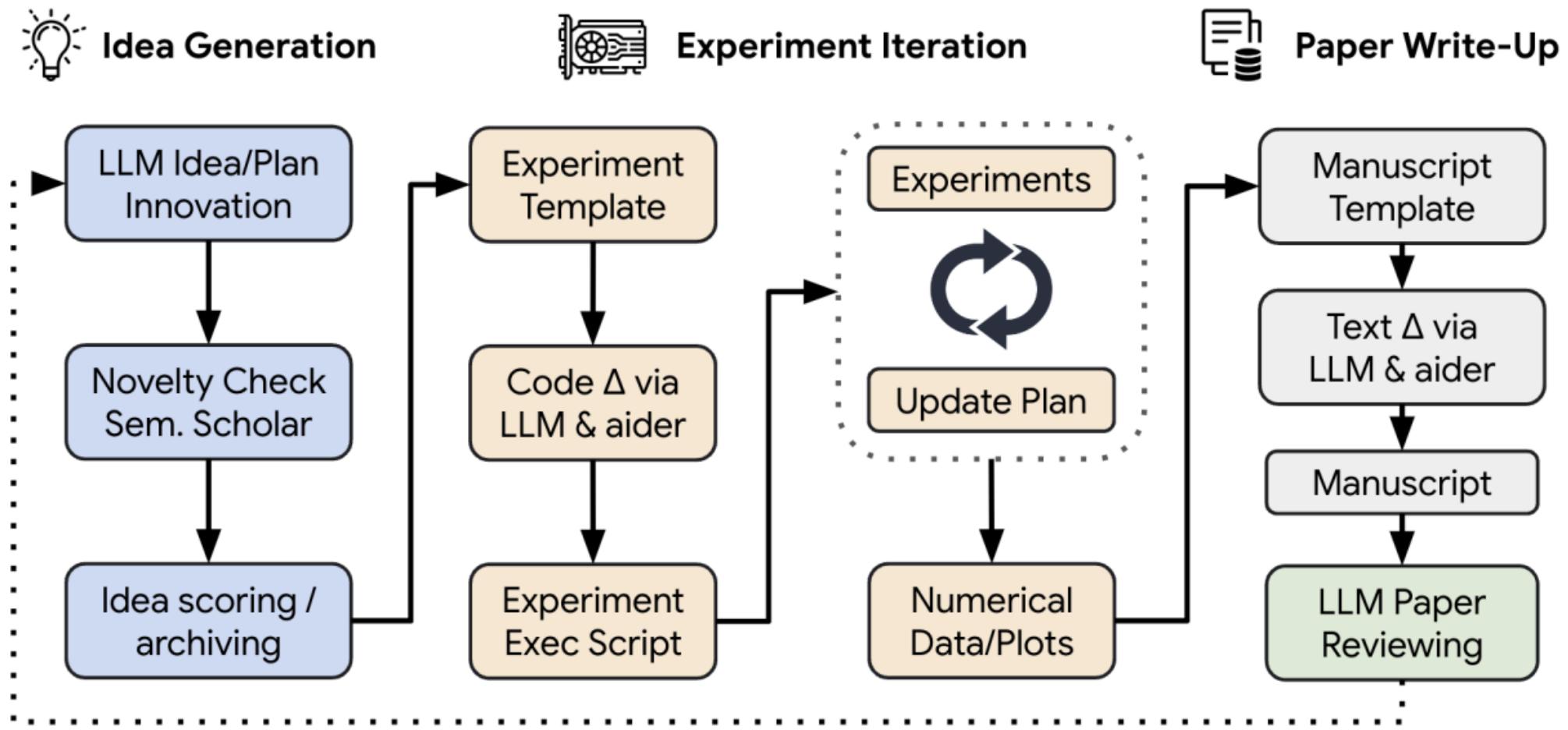
A praise for hallucinations

Hallucinations are necessary for predictions

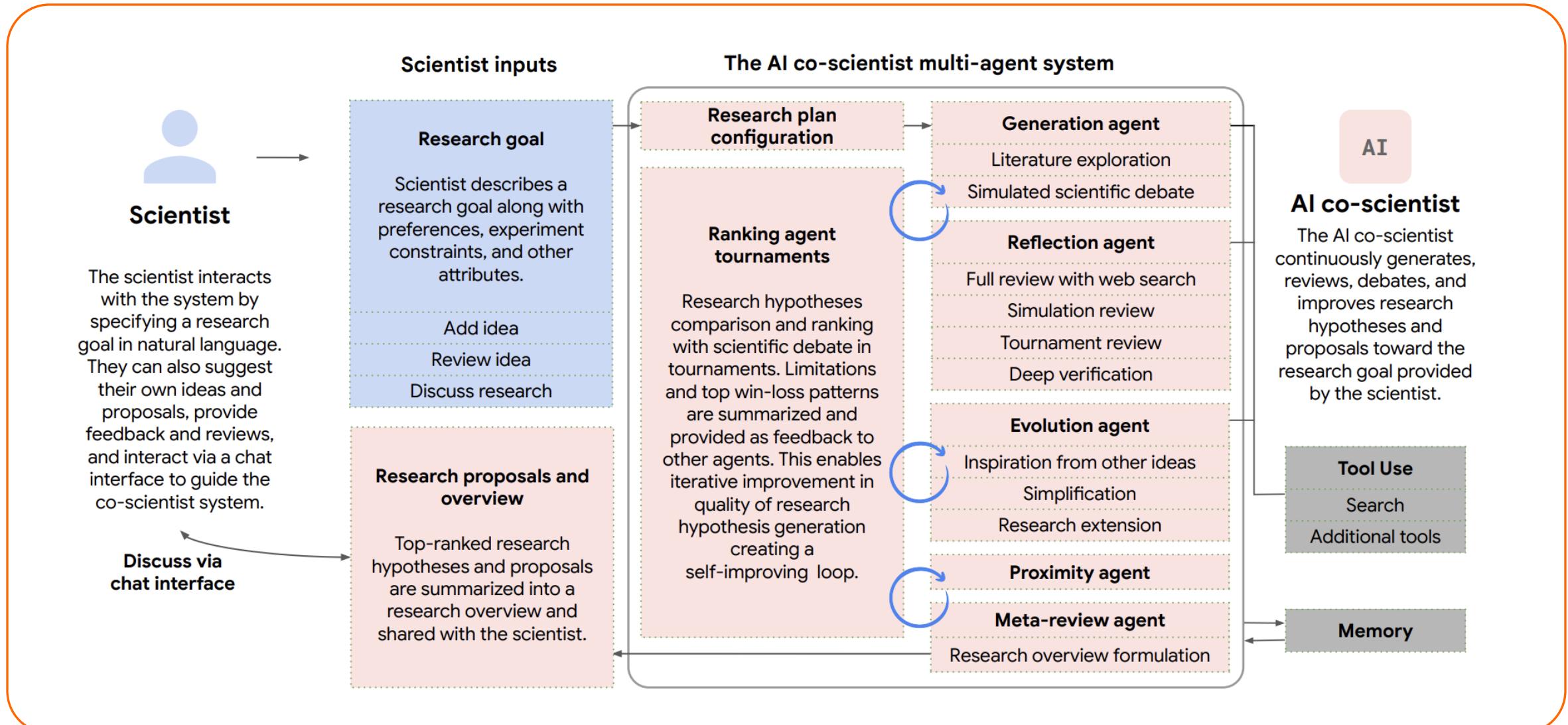
Backward-looking vs forward-looking hallucinations



AI scientist



AI co-scientist



Final remarks: some things I believe

- 1. LLMs already have a positive impact on our scientific workflow**
- 2. LLMs are an interesting subject of study across multiple dimensions**
- 3. Demand for HPC will only go up for the above two reasons**



Thanks again to anyone who made Sofia possible :)

and to VSC for organizing + support!

1. Patient

2. Kind

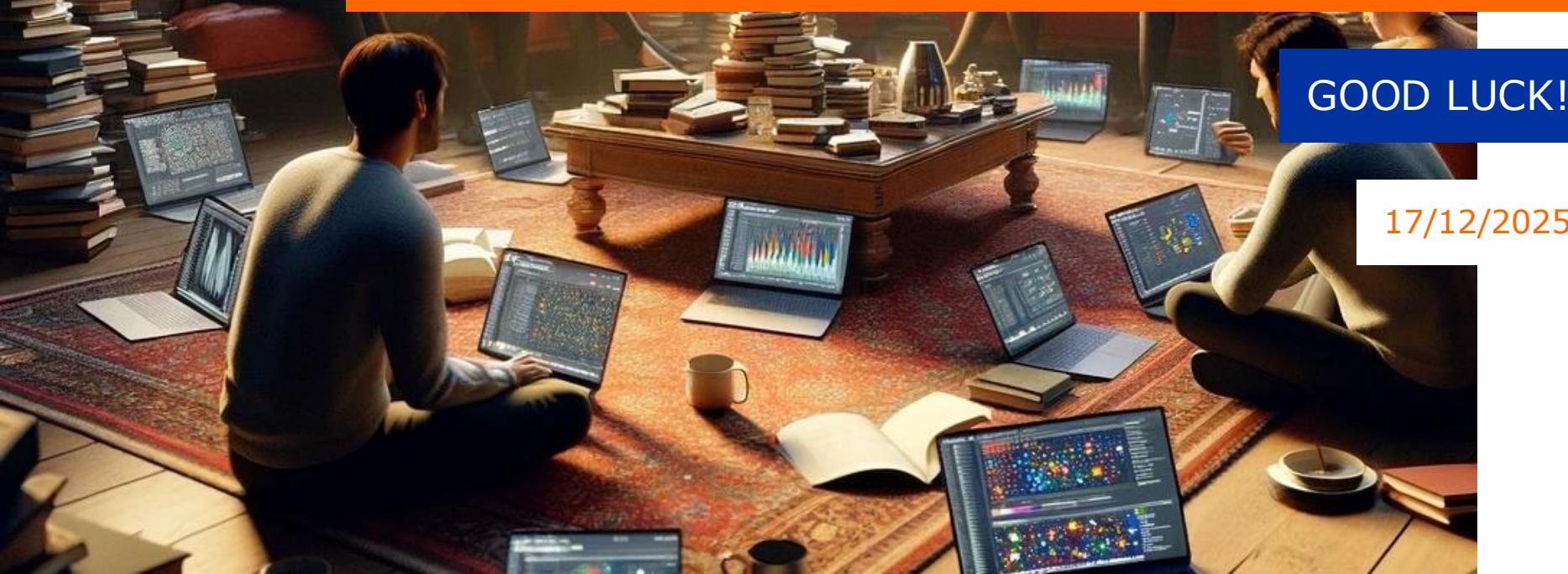
3. Intelligent ...

whatever that may be :D





LARGE LANGUAGE MODELS WITHOUT HPC?



fwo