

# LARGE LANGUAGE MODELS IN SCIENCE

## THE ELEPHANT IN THE ROOM

18/03/2025

Andres Algaba

**\*Prompted ChatGPT 5.2 (Thinking): "Can you generate a cover image for my slidedeck based on the vibe it gives?"**

# Enormous progress in the last year

## 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 the UK and worldwide as the disease increases its resistance to antibiotics

## Towards an AI co-scientist

Juraj Gottweis<sup>\*, †, 1</sup>, Wei-Hung Weng<sup>\*, †, 2</sup>, Alexander Daryin<sup>\*, 1</sup>, Tao Tu<sup>\*, 3</sup>, Anil Palepu<sup>2</sup>, Petar Sirkovic<sup>1</sup>, Artiom Myaskovsky<sup>1</sup>, Felix Weissenberger<sup>1</sup>, Keran Rong<sup>3</sup>, Ryutarō Tanno<sup>3</sup>, Khaled Saab<sup>3</sup>, Dan Popovici<sup>2</sup>, Jacob Blum<sup>7</sup>, Fan Zhang<sup>2</sup>, Katherine Chou<sup>2</sup>, Avinatan Hassidim<sup>2</sup>, Burak Gokturk<sup>1</sup>, Amin Vahdat<sup>1</sup>, Pushmeet Kohli<sup>3</sup>, Yossi Matias<sup>2</sup>, Andrew Carroll<sup>2</sup>, Kavita Kulkarni<sup>2</sup>, Nenad Tomasev<sup>3</sup>, Yuan Guan<sup>7</sup>, Vikram Dhillon<sup>4</sup>, Eeshit Dhaval Vaishnav<sup>5</sup>, Byron Lee<sup>5</sup>, Tiago R D Costa<sup>6</sup>, José R Penadés<sup>6</sup>, Gary Peltz<sup>7</sup>, Yunhan Xu<sup>3</sup>, Annalisa Pawlosky<sup>1, †</sup>, Alan Karthikesalingam<sup>2, †</sup> and Vivek Natarajan<sup>2, †</sup>

<sup>1</sup>Google Cloud AI Research, <sup>2</sup>Google Research, <sup>3</sup>Google DeepMind,



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

April 27-28, 2025 | Singapore Expo

## The AI Scientist Generates its First Peer-Reviewed Scientific Publication

March 12, 2025

## Early science acceleration experiments with GPT-5

Sébastien Bubeck<sup>1</sup>, Christian Coester<sup>2</sup>, Ronen Eldan<sup>1</sup>, Timothy Gowers<sup>3</sup>, Yin Tat Lee<sup>1</sup>, Alexandru Lupasca<sup>1, 4</sup>, Mehtaab Sawhney<sup>5</sup>, Robert Scherrer<sup>4</sup>, Mark Sellke<sup>1, 6</sup>, Brian K. Spears<sup>7</sup>, Derya Unutmaz<sup>8</sup>, Kevin Weil<sup>1</sup>, Steven Yin<sup>1</sup>, Nikita Zhivotovskiy<sup>9</sup>

<sup>1</sup>OpenAI

<sup>2</sup>University of Oxford

<sup>3</sup>Collège de France and University of Cambridge

<sup>4</sup>Vanderbilt University

<sup>5</sup>Columbia University

<sup>6</sup>Harvard University

<sup>7</sup>Lawrence Livermore National Laboratory

<sup>8</sup>The Jackson Laboratory

<sup>9</sup>University of California, Berkeley

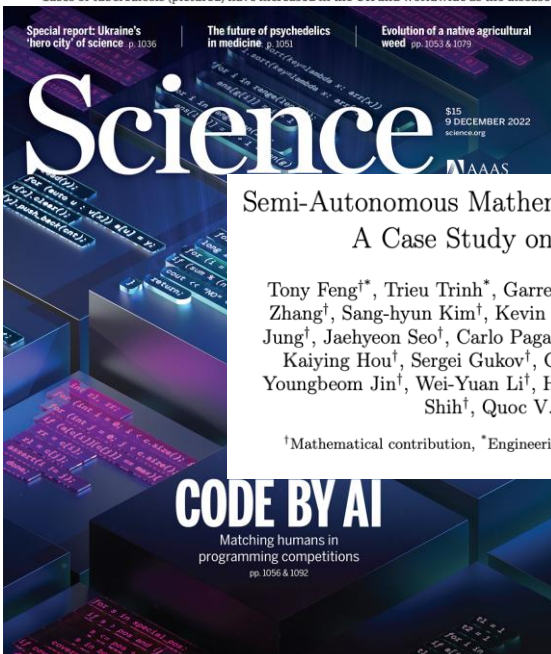
## Semi-Autonomous Mathematics Discovery with Gemini: A Case Study on the Erdős Problems

Tony Feng<sup>†\*</sup>, Trieu Trinh<sup>\*</sup>, Garrett Bingham<sup>\*</sup>, Jiwon Kang<sup>†</sup>, Shengtong Zhang<sup>†</sup>, Sang-hyun Kim<sup>†</sup>, Kevin Barreto<sup>†</sup>, Carl Schildkraut<sup>†</sup>, Junehyuk Jung<sup>†</sup>, Jaehyeon Seo<sup>†</sup>, Carlo Pagano<sup>†</sup>, Yuri Chervonyi<sup>†</sup>, Dawsen Hwang<sup>\*</sup>, Kaiying Hou<sup>†</sup>, Sergei Gukov<sup>†</sup>, Cheng-Chiang Tsai<sup>†</sup>, Hyunwoo Choi<sup>†</sup>, Youngbeom Jin<sup>†</sup>, Wei-Yuan Li<sup>†</sup>, Hao-An Wu<sup>†</sup>, Ruey-An Shiu<sup>†</sup>, Yu-Sheng Shih<sup>†</sup>, Quoc V. Le<sup>o</sup>, Thang Luong<sup>o</sup>

<sup>†</sup>Mathematical contribution, <sup>\*</sup>Engineering contribution

September 17, 2025 Research

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



# And there are three reasons why

## AI cracks superbug problem in t that took scientists years

6 days ago

Tom Gerken  
Technology reporter



Cases of tuberculosis (pictured) have increased in the UK and worldwide as the disease increases i



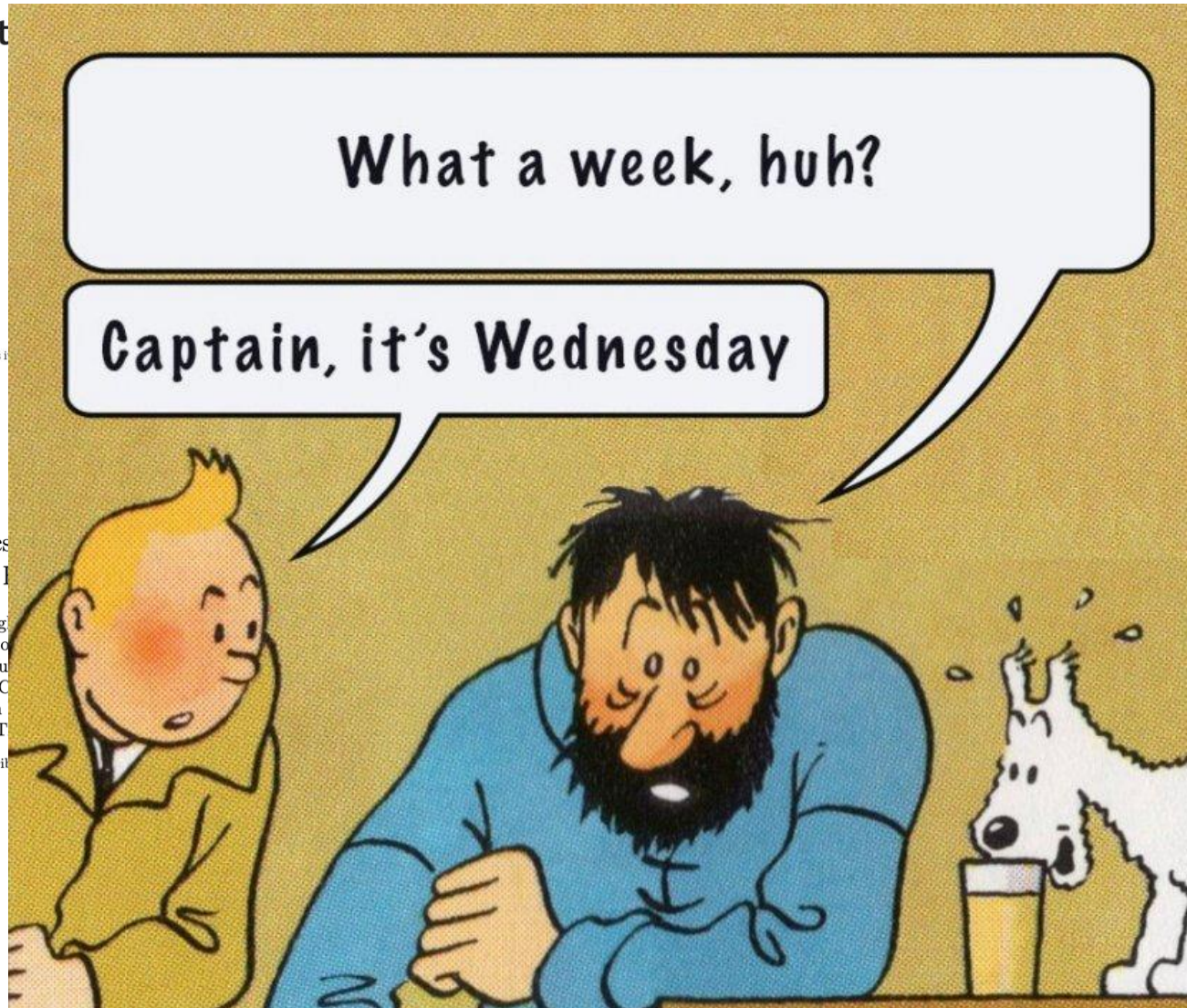
### Semi-Autonomous Mathematics A Case Study on the I

Tony Feng<sup>†</sup>, Trieu Trinh<sup>\*</sup>, Garrett Bing Zhang<sup>†</sup>, Sang-hyun Kim<sup>†</sup>, Kevin Barreto Jung<sup>†</sup>, Jaehyeon Seo<sup>†</sup>, Carlo Pagano<sup>†</sup>, Yu Kaiying Hou<sup>†</sup>, Sergei Gukov<sup>†</sup>, Cheng-C Youngbeom Jin<sup>†</sup>, Wei-Yuan Li<sup>†</sup>, Hao-An Shih<sup>†</sup>, Quoc V. Le<sup>‡</sup>, T

<sup>†</sup>Mathematical contribution, <sup>\*</sup>Engineering contri

### CODE BY AI

Matching humans in programming competitions  
pp. 1056 & 1092



Times

## 'Artificial Superintelligence'

new optimism that artificial scientific discovery.

## thesis Generation, Validation (ICLR,

and Scientific Publication

## ation experiments with GPT-5

ester<sup>2</sup>, Ronen Eldan<sup>1</sup>, Timothy Gowers<sup>3</sup>, Yin Tat Lee<sup>1</sup>,  
Itaab Sawhney<sup>5</sup>, Robert Scherrer<sup>4</sup>, Mark Sellke<sup>1,6</sup>,  
iaz<sup>8</sup>, Kevin Weil<sup>1</sup>, Steven Yin<sup>1</sup>, Nikita Zhivotovskiy<sup>9</sup>

<sup>1</sup>OpenAI

<sup>2</sup>University of Oxford  
France and University of Cambridge

<sup>4</sup>Vanderbilt University

<sup>5</sup>Columbia University

<sup>6</sup>Harvard University

e Livermore National Laboratory

<sup>7</sup>The Jackson Laboratory

ersity of California, Berkeley

# Reason 1: scaling just works

## Scaling laws (lin-log relation)

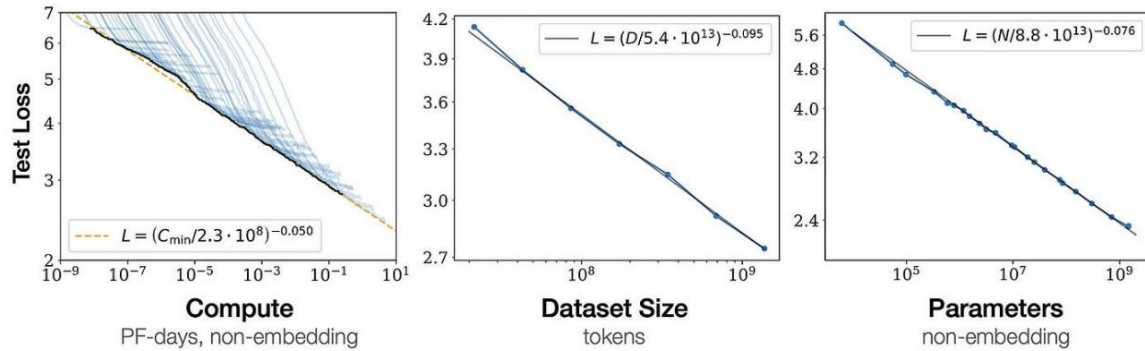
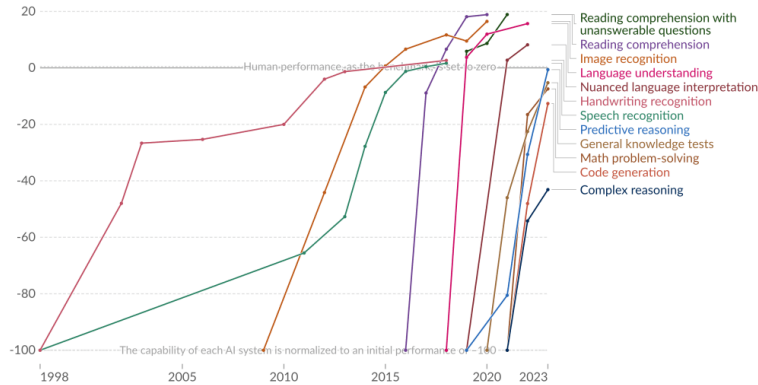


Figure from Kaplan et al. (2020)

Test scores of AI systems on various capabilities relative to human performance

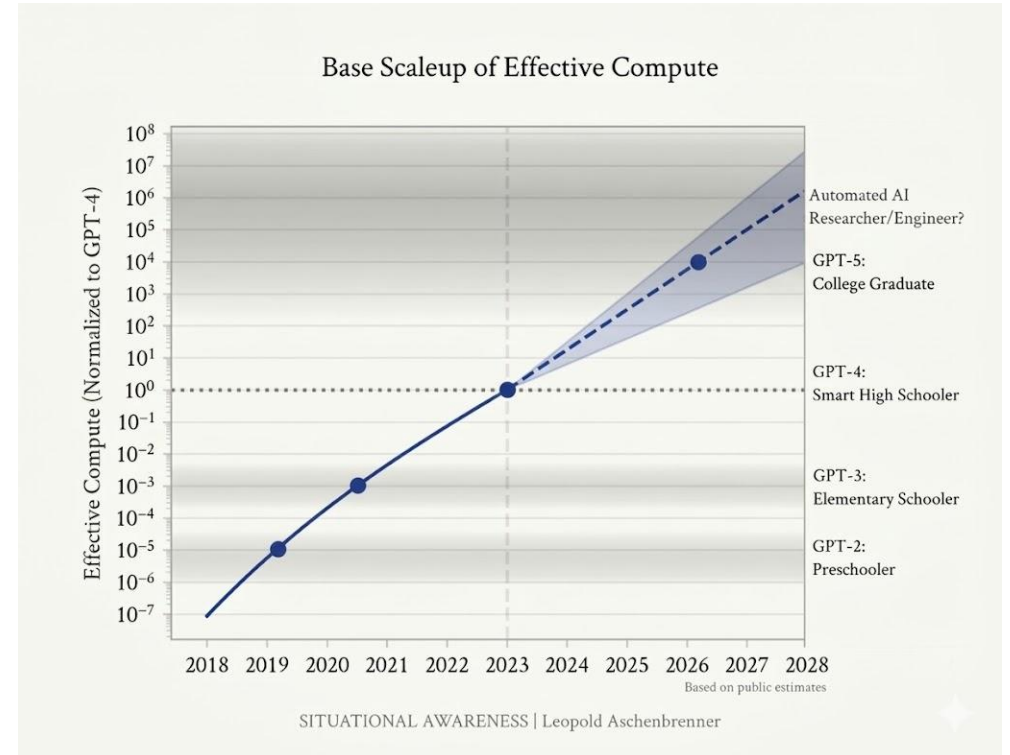
Within each domain, the initial performance of the AI is set to -100. Human performance is used as a baseline, set to zero. When the AI's performance crosses the zero line, it scored more points than humans.



Data source: Kiela et al. (2023) OurWorldInData.org/artificial-intelligence | CC BY  
Note: For each capability, the first year always shows a baseline of -100, even if better performance was recorded later that year.

Most of our benchmarks are saturated, which makes it harder to measure the progress ...

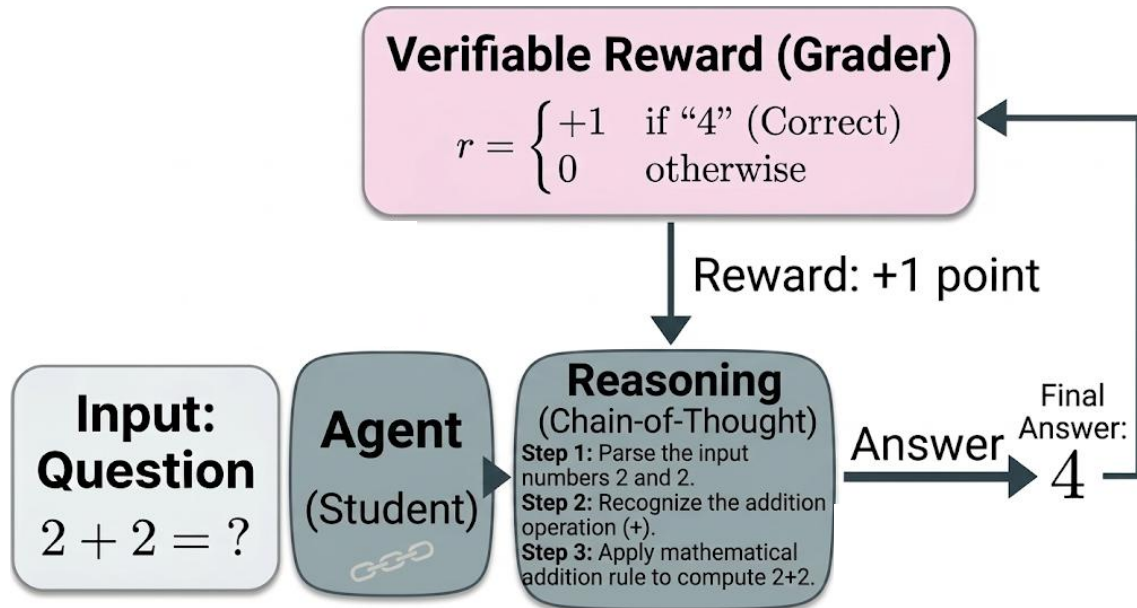
## Jumps are just observations of a continuous process



Updated (Aschenbrenner, 2024) using Gemini 3 Pro

# Reason 2: scaling up reinforcement learning and test-time compute

## Rewarding strong reasoning

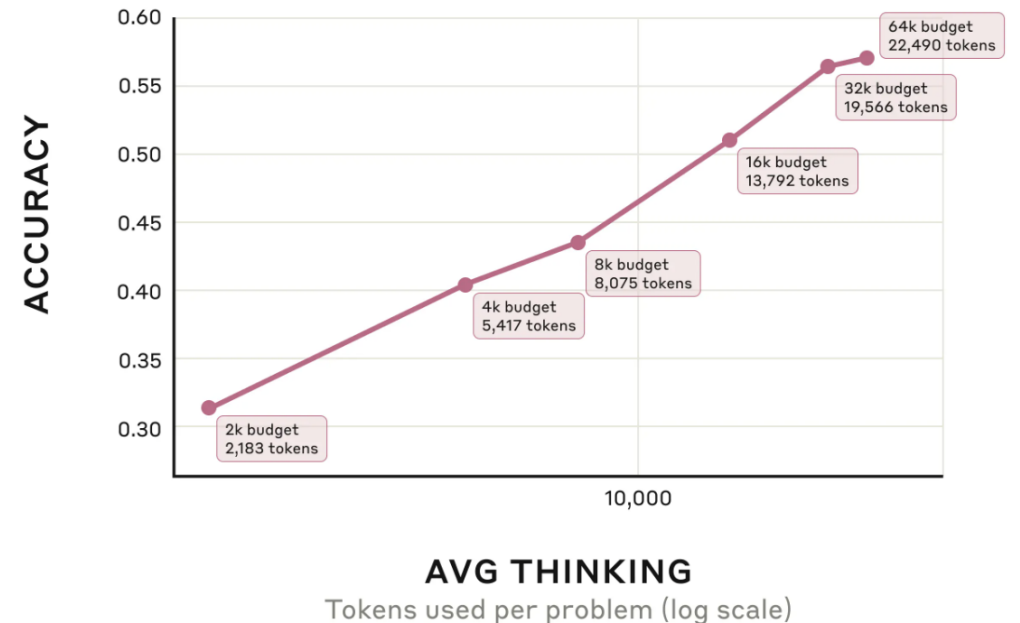


Adapted (Lambert et al., 2024) using Gemini 3 Pro

## Scaling up reasoning (again lin-log)

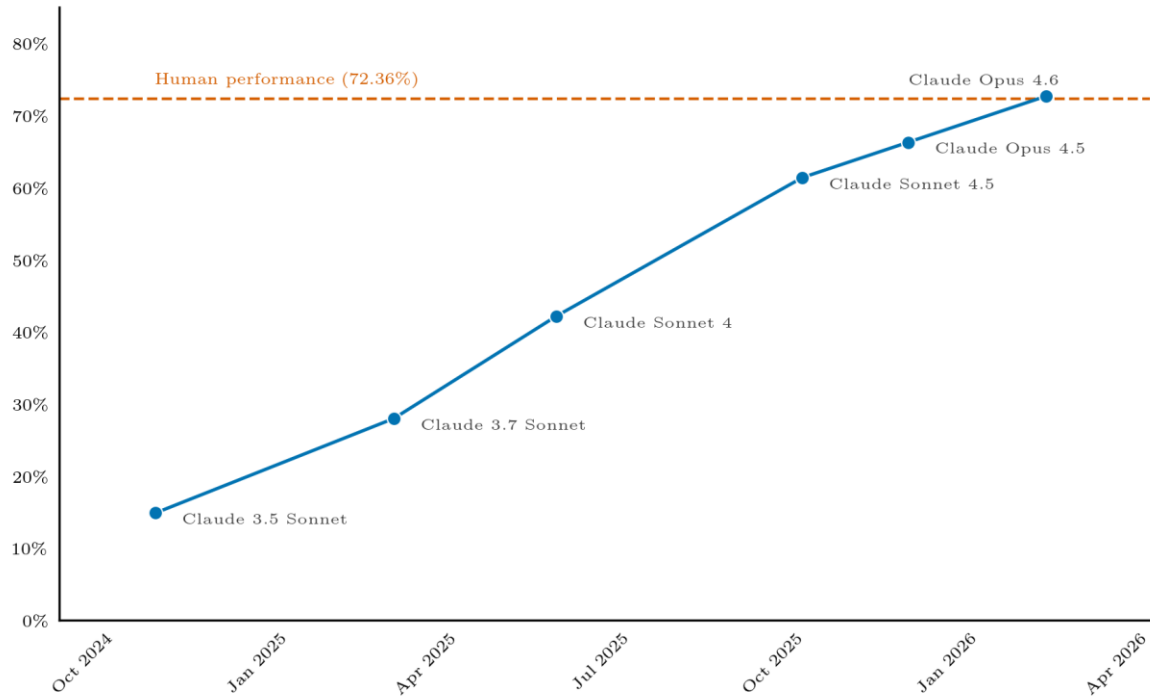
### AIME 2024 performance

vs. actual thinking token usage



## Reason 3: generalization (it works for many tasks)

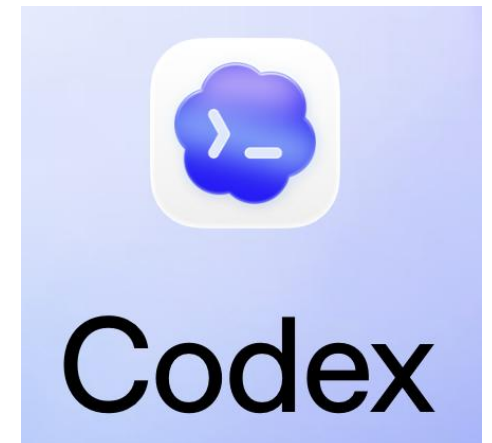
### Computer-use benchmark (OSWorld)



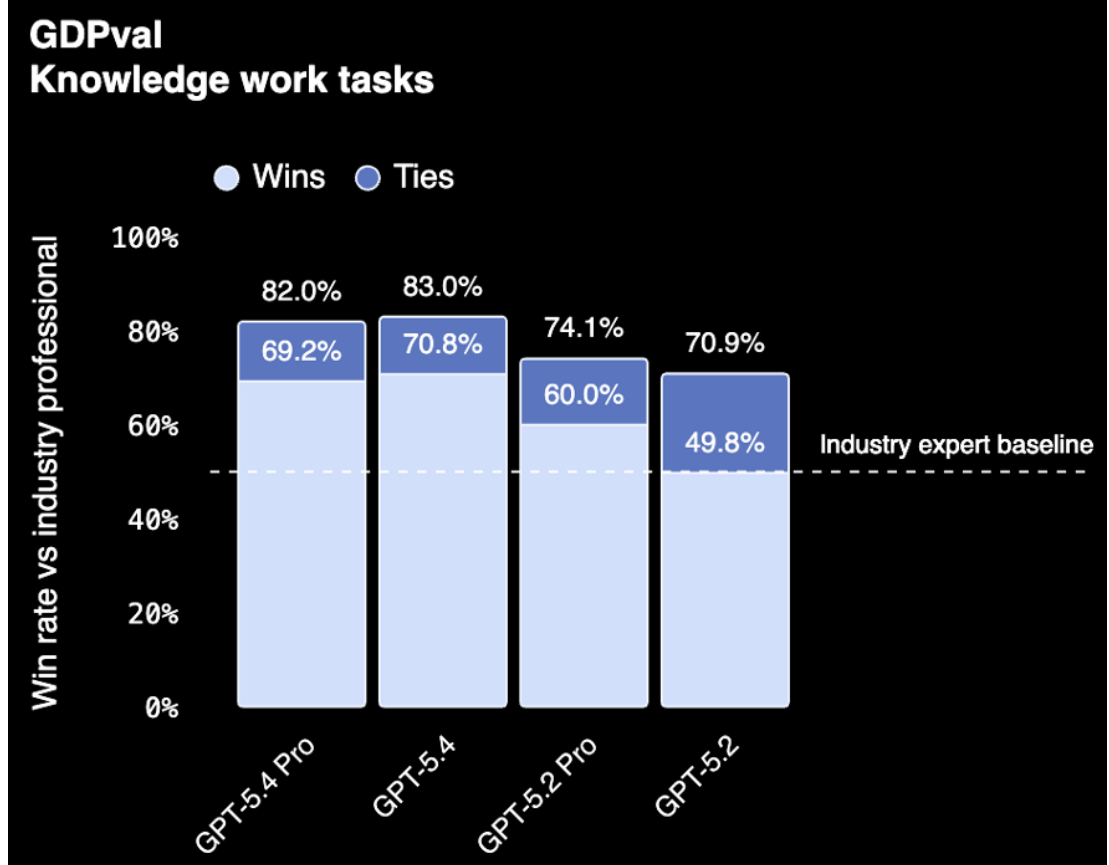
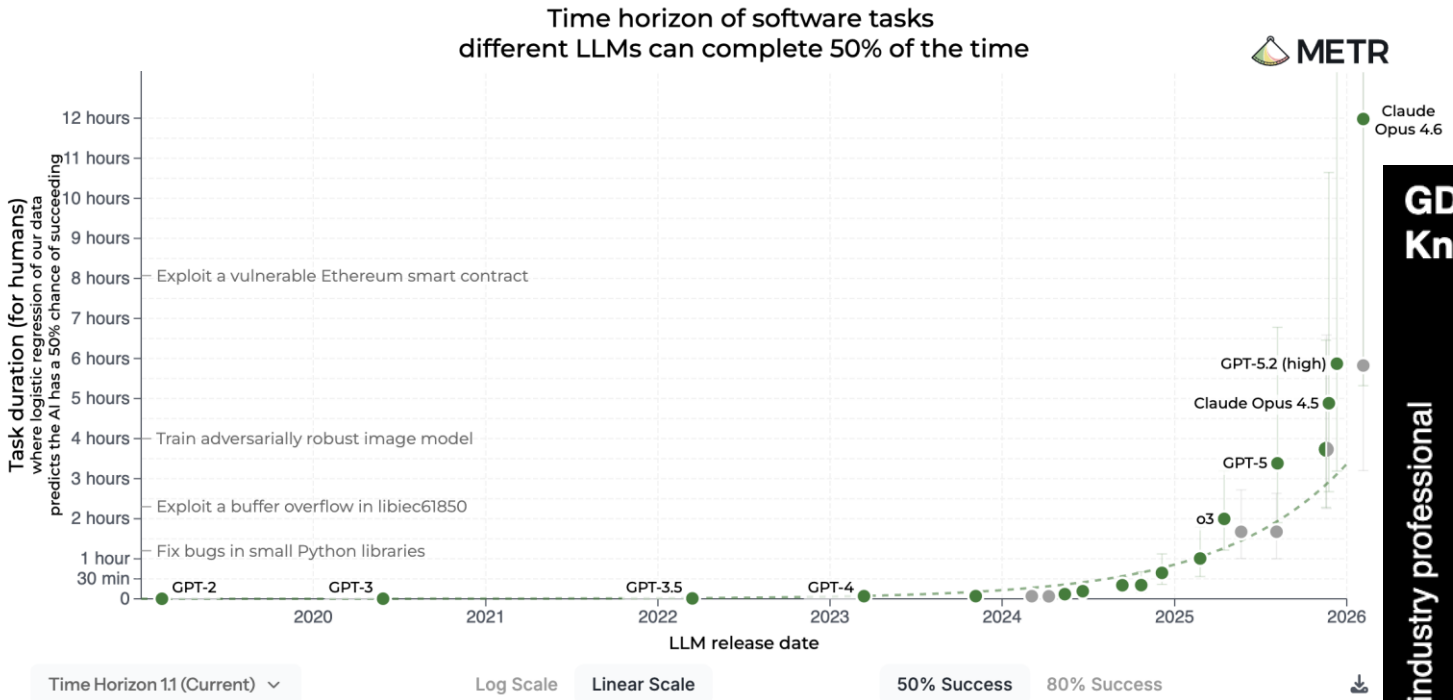
\*Claude Code made this one-shot from scratch

### Results in new products

AI



# Real world impact goes beyond classical benchmarks



# Impact on the process of scientific discovery

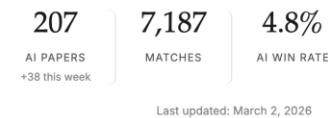
People are already running the “1000-paper experiment” with Claude Code

## Can AI *automate* policy evaluation?

AI may soon be capable of producing rigorous economic research. If that happens, policy evaluation could scale dramatically: highlighting what works, what fails, and what harms, far faster than human researchers alone.

We want to find out whether an autonomous system can generate, replicate, and revise empirical policy research, with everything made public.

This is an experiment in building reliable AI research systems.



## Will paper publications still be useful?

🔒 | POLICY ARTICLE | ARTIFICIAL INTELLIGENCE



## Scientific production in the era of large language models

With the production process rapidly evolving, science policy must consider how institutions could evolve

KEIGO KUSUMEGI, XINYU YANG, PAUL GINSPARG, MATHIJS DE VAAN, TOBY STUART, AND YIAN YIN [Authors Info & Affiliations](#)

SCIENCE • 18 Dec 2025 • Vol 390, Issue 6779 • pp. 1240-1243 • DOI: 10.1126/science.adw3000

## Peer review broken (again?)

NEWS | 15 December 2025 | Correction [16 December 2025](#)

## More than half of researchers now use AI for peer review – often against guidance

A survey of 1,600 academics found that more than 50% have used artificial-intelligence tools while peer reviewing manuscripts.

## False discovery rates (p-hacking)

### Do Claude Code and Codex P-Hack?

## Sycophancy and Statistical Analysis in Large Language Models\*

Samuel G.Z. Asher<sup>†</sup> Janet Malzahn<sup>†</sup> Jessica M. Persano<sup>‡</sup>

Elliot J. Paschal<sup>‡</sup> Andrew C. W. Myers<sup>§</sup> Andrew B. Hall<sup>¶</sup>

# All researchers are affected (almost) equally

## Impact on data (falsification?)

COMMENT | 09 February 2026

### How to deal with the survey-taking AI agents that threaten to upend social science

Researchers need new bot-detection strategies that exploit the limits of human reasoning rather than AI weaknesses.

By [Folco Panizza](#), [Yara Kyrychenko](#) & [Jon Roozenbeek](#) 

## Negative feedback loops?

### Large Language Models Reflect Human Citation Patterns with a Heightened Citation Bias

Andres Algaba<sup>1\*</sup> Carmen Mazijn<sup>1</sup> Vincent Holst<sup>1</sup>  
Floriano Tori<sup>1</sup> Sylvia Wenmackers<sup>2</sup> Vincent Ginis<sup>1,3</sup>

## Equity in access to (state-of-the-art) systems

Google DeepMind

### Aletheia tackles *FirstProof* autonomously

Tony Feng<sup>\*</sup>, Junehyuk Jung, Sang-hyun Kim, Carlo Pagano, Sergei Gukov, Chiang-Chiang Tsai, David P. Woodruff, Adel Javanmard, Aryan Mokhtari, Dawsen Hwang, Yuri Chervonyi, Jonathan N. Lee, Garrett Bingham, Trieu H. Trinh, Vahab Mirrokni, Quoc V. Le, Thang Luong<sup>\*</sup>

<sup>\*</sup>Project leads. Work conducted under Google DeepMind.



Petar Veličković  @PetarV\_93

10 months ago

AlphaEvolve is here! 

this is one special system (especially when optimising things with jagged edges 😊) -- had a fantastic time using it! congrats [@SashaVNovikov](#) [@matejbalog](#) [@ThuyNganVu](#) and team!!



Computer Science > Artificial Intelligence

[Submitted on 21 Feb 2026]

### Early Evidence of Vibe-Proving with Consumer LLMs: A Case Study on Spectral Region Characterization with ChatGPT-5.2 (Thinking)

Brecht Verbeken, Brando Vagenende, Marie-Anne Guerry, Andres Algaba, Vincent Ginis

**We should be talking about the elephant in the room**

**Our current way of organizing science needs to be updated**



**\*Prompted Gemini 3 Pro: "Can you generate an image where there is no elephant in the room?"**

# LARGE LANGUAGE MODELS IN SCIENCE

## THE ELEPHANT IN THE ROOM

18/03/2025