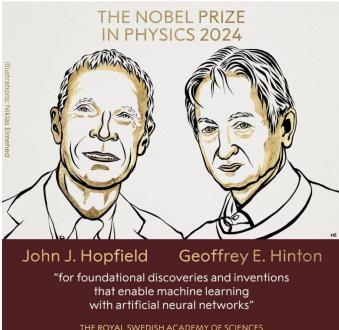


## DOTE 6635: Artificial Intelligence for Business Research

# What's New in AI

Renyu (Philip) Zhang

1



**THE NOBEL PRIZE IN PHYSICS 2024**

**John J. Hopfield** **Geoffrey E. Hinton**  
"for foundational discoveries and inventions that enable machine learning with artificial neural networks"

THE ROYAL SWEDISH ACADEMY OF SCIENCES

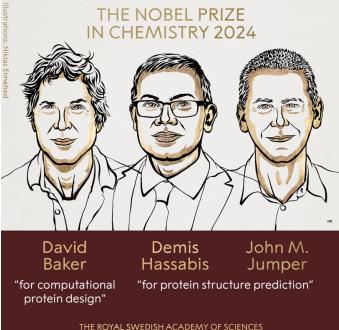
## AI: Future of Human Civilization

华尔街见闻 首页 资讯 快讯 行情 日历 APP | VIP会员 大师课 生活家

763
收藏
分享
评论

OpenAI完成最新一轮66亿美元融资 警告投资者不准支持马斯克xAI等劲敌  
赵雨荷 10-03 00:07

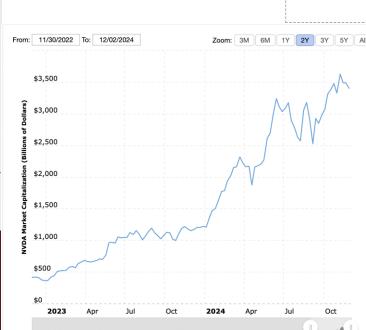
**摘要：**  
本轮融资也是史上规模最大的私人投资之一，由Thrive Capital领投，参与者还包括微软、英伟达、软银等，其中微软投资约7.5亿美元。本轮融资过后，OpenAI的估值达到1570亿美元，跻身全球前三大初创公司的行列。同时，OpenAI希望与投资者达成独家协议，防止马斯克的xAI和Anthropic等竞争对手获得战略合作机会和资本支持。



**THE NOBEL PRIZE IN CHEMISTRY 2024**

**David Baker** **Demis Hassabis** **John M. Jumper**  
"for computational protein design"

THE ROYAL SWEDISH ACADEMY OF SCIENCES



From: 11/30/2022 To: 12/02/2024 Zoom: 3M 6M 1Y 2Y 3Y 5Y All

NYSE Market Capitalization (Billions of Dollars)

| Date       | Capitalization (Billions of Dollars) |
|------------|--------------------------------------|
| 2023-01-01 | ~\$500                               |
| 2023-07-01 | ~\$1,000                             |
| 2023-12-31 | ~\$1,500                             |
| 2024-01-01 | ~\$2,000                             |
| 2024-04-01 | ~\$2,500                             |
| 2024-07-01 | ~\$3,000                             |
| 2024-10-01 | ~\$3,500                             |

**TechCrunch** **Trump considers naming an 'AI czar'**

Kyle Wiggers Wed, November 27, 2024 at 1:59 PM PST • 1 min read

Incoming president Donald Trump is considering naming an "AI czar" in the White House, Axios reports.

Should Trump appoint such a policy person, they'd be charged with helping to coordinate federal regulation and governmental use of AI. Importantly, an AI czar wouldn't require Senate confirmation, Axios notes — allowing them to get to work on the administration's goals faster.

2

# This Week in AI

The page contains several charts comparing AI models across different benchmarks and tasks. One chart shows accuracy/percentile (%) for DeepSeek-R1, OpenAI-o1-1217, DeepSeek-R1-32B, OpenAI-o1-mini, and DeepSeek-V3 across benchmarks like AIMES 2024, Codeforces, GPOQA Diamond, MATH-500, MMLU, and SWE-bench Verified. Another chart shows API prices for input and output tokens. A sidebar shows a timeline of posts from users hanxiao and My English ceng ceng up.

[https://github.com/deepseek-ai/DeepSeek-R1/blob/main/DeepSeek\\_R1.pdf](https://github.com/deepseek-ai/DeepSeek-R1/blob/main/DeepSeek_R1.pdf)

accuracy

Research Math (EpochAI Frontier Math)

accuracy

previous SoTA o3

Summary

- Epoch AI, the developer of a mathematics benchmark, did not initially disclose funding from OpenAI due to a non-disclosure agreement, and this only became known when OpenAI set a new record on the benchmark with its o3 model.
- More than 60 mathematicians who contributed to the benchmark were not systematically informed about OpenAI's involvement and believed their work would remain confidential and only be used by Epoch AI.
- OpenAI had early access to many of the tasks and solutions in the benchmark, and while there was a verbal agreement that this data would not be used for training, the lack of transparency surrounding the arrangement remains problematic.

3

# What's New (Feb 4/2025-1)

**AI EFFECT**

## Nvidia sheds almost \$600 billion in market cap, biggest one-day loss in U.S. history

PUBLISHED MON, JAN 27 2025 4:08 PM EST | UPDATED MON, JAN 27 2025 5:26 PM EST

Samantha Subin @SAMANTHA\_SUBIN

SHARE f X in e

JOSH HAWLEY U.S. SENATOR FOR MISSOURI

Home About How Can I Help News Contact

Wednesday, January 29, 2025

United States World / United States & Canada

### Hawley Introduces Legislation to Decouple American AI Development from Communist China

Microsoft, OpenAI investigate China's DeepSeek over data breach

Donald Trump's AI tsar has claimed there's 'substantial evidence' that DeepSeek leaned on OpenAI's models to develop its own technology

Reading Time: 2 minutes

Artificial intelligence Tech / Tech War

Tech war: US firms embrace DeepSeek's AI model despite scrutiny

From Nvidia to Microsoft and Amazon, US companies have rushed to adopt the Chinese start-up's R1 model

Reading Time: 2 minutes

Why you can trust SCMP

Why you can trust SCMP

DeepSeek became the most downloaded app on Apple's App Store in the US on Monday. VCG/VCG via Getty Images

4

## What's New (Feb/4/2025-2)

**DeepSeek-R1 is like the Watt Steam Engine and Ford Model T in the age of AI!**

如何评价DeepSeek等大模型在中科院物理所理论竞赛中的表现？

更重要的是，Deepseek 在一定程度上打破了学术垄断，促进了学术界的智能平权。因为学术界的现状是，一般背景的学生和研究者没人指导，缺有价值的想法，更缺手把手的指导，科研卡住了也找不到人帮忙，只能自己死磕，磕不出来就等着延毕或 GG。即使是大佬组的，好 idea 一定能轮得到你？大佬能天天手把手教你推公式？反正就是你自己太菜也怨不得别人。而且科研很多时候其实就差临行的一句话，告诉你应该用某某方法解决，考虑 XXX，不告诉你，你自己在文献堆里翻几个月都还不确定找得到找不到。要是你需要的是某门现代数学，你望着艰深的符号，如果不是确定一定能解决自己的问题，心里早就打退堂鼓了。但你都不会，拿头确定？大佬的价值就在于此，他会凭丰富的经验给你一个确定性。这样你就能走下去。现在 Deepseek 当然达不到真正的大佬的程度，但超越一些不甚合格的导师还是绰绰有余，而且它能手把手教你。顶尖研究组可能看不上 Deepseek 给的 idea 和指导，毕竟 SOTA 只是小圈子里的知识，AI 不一定知道。但对于广大学术底层，这可就太有用了，这玩意简直就是开挂，效率一下提升一个数量级也不是没可能。况且，牛组的人真的就是靠智商吗？恐怕内部信息才是真正的护城河。Deepseek 虽然不能完全消除你的护城河，但会降低其他复现的时间成本。普通组在有了 Deepseek 的加持下，是有希望做出之前只有牛组才能做出的工作的。

那这些意味着什么呢？

没错，学术生产的手工作坊模式<sup>\*</sup>要被撼动了。

<https://www.zhihu.com/question/10879827313/answer/89959861140>

DeepSeek 登顶苹果美国区免费 APP 下载排行榜，与 ChatGPT 相比...

3，最大的影响，在于用很低的成本将一个能力非常强的模型开源到本地，在政府应用尤其是安全方面的应用中，在数据比较敏感的企业内部应用中，已经产生了较大潜在影响，受冲击最大的是 kimi、qwen 等希望做服务和应用的企业。企业内部的大模型应用，以前大部分都是说着玩的，各种条条框框，使得企业对数据被上传到人工智能企业中非常警惕。目前的 kimi、qwen、讯飞等模型要么对用户免费或者低价，要么开源免费，其真正的目的就是要吸引企业和政府来购买他们的微调服务，他们帮政府或者企业微调一个本地的模型，在这个阶段收费。而这个阶段能够收费的逻辑之所以成立，建立在“本地模型不微调能力很差”、“知识库很难建”、“微调算力需求特别大自己实现不了”等现实逻辑上，导致许多企业和政府部门在这一阶段就已经知难而退了——要建立一个内部大模型我不仅要买算力显卡还要你来给我微调花那么多钱，那么我不做了还不行吗——但是 Deepseek 打破了这一点，r1 在能力上的本质提升使得“不微调就能用得很不错”“用 r1 的能力来弥补简陋的知识库”成为可能，于是他们就不需要购买服务，只需要购买算力就可以了。在这个过程中，大模型的应用会更广泛，更为原子化，算力的需求会成倍增长，而非下降，更多的人和企业将会加入到“人人都有自己的应用”的创新中，产生新的可能。

<https://www.zhihu.com/question/10669048245/answer/87990392650>

5

## Even More Shockingly.....

**OpenAI o3 Deep Research could directly produce:**

Revisiting the McKinley Tariff of 1890 through the Lens of Modern Trade Theory

[o3 Deep Research]<sup>\*</sup>

### Abstract

This paper was written with a one-shot prompt (from Kevin Bryan) on o3 Deep Research, no iteration, 10 minutes of thinking. The Tariff Act of 1890, better known as the McKinley Tariff, was a pivotal episode in U.S. trade policy, dramatically raising import duties to near-record levels. This paper provides an analysis of the McKinley Tariff by integrating historical evidence with insights from modern international trade theory. We revisit the economic and political debates of the 1890s using contemporary trade models—including models of heterogeneous firms (Meltz, 2003), Ricardian comparative advantage in general equilibrium (Eaton and Kortum, 2002), and other new trade theory advances—to re-evaluate the tariff's impacts. Historical data on trade flows, tariff rates, and industry output are analyzed alongside contemporary accounts to assess the short- and long-run effects of the tariff. We find that while the McKinley Tariff accelerated the development of certain industries (notably tinplate production) and was implemented in an era of changing comparative advantage for the United States, its overall welfare effects were mixed and likely negative when evaluated with modern trade metrics. The tariff's protective gains to manufacturers came at the cost of higher prices for consumers and implicit burdens on agricultural exporters. However, consistent with modern trade models, the United States' large market power means some tariff incidence was borne by foreign exporters. The paper concludes by drawing parallels between the McKinley Tariff episode and contemporary trade policy tensions, including recent U.S.-China tariff disputes and debates over protectionism in the global trading system.

<https://kevinbryanecon.com/o3McKinley.pdf>

Information Frictions and Innovation: A Formal Theory

[o3 Deep Research]<sup>\*</sup>

February 2, 2025

### Abstract

This paper was written with a one-shot prompt (from Kevin Bryan) on o3 Deep Research, no iteration, 10 minutes of thinking. This paper develops a formal economic theory exploring how *information frictions* impact innovation, extending beyond the usual focus on incentive problems. We present a model of innovation in which the production of new ideas builds on previous innovations, but knowledge about these prior innovations is distributed across many agents. In this environment, classical welfare theorems break down: key inputs into innovation (knowledge) are unpriced and information is not optimally aggregated, leading to market failures. We formally compare several mechanisms—patents, prizes, advance market commitments (AMCs), and others—in their ability to overcome these information frictions. We derive propositions showing how each mechanism influences the aggregation of dispersed knowledge and the efficiency of innovation, providing rigorous proofs. Our results highlight that beyond providing incentives, innovation institutions serve a critical role in coordinating distributed information. The analysis yields insights into the design of innovation policy when knowledge is decentralized.

<https://kevinbryanecon.com/o3InnovationTheory.pdf>

6

## What's New (Feb/4/2025-3)

### 18. Journal of Operations Management (JOM)

- (1) 运营界的“效率魔人”：研究如何让员工996还心怀感恩？JOM会给你发奖杯。
- (2) “优化一切，包括人生”：JOM的编辑可能连早餐麦片都要按算法排列。
- (3) 审稿人的隐藏人设：白天审论文，晚上在淘宝帮人优化购物车。
- (4) 投稿建议：找个效率问题，跑个优化模型，最后用一篇“效率至上”论文征服审稿人。

### 19. Manufacturing & Service Operations Management (MSOM)

- (1) 运营界的“技术宅”：研究如何用AI让咖啡机不洒奶泡？MSOM会为你开专题。
- (2) “模型越复杂，审稿人越嗨”：简单问题复杂化是MSOM的快乐源泉。
- (3) 审稿人的日常：一边骂你的模型不实用，一边偷偷用在自家车库管理上。
- (4) 投稿建议：找个技术问题，跑个复杂模型，最后用一篇“极客风”论文打动审稿人。

### 20. Production and Operations Management (POM)

- (1) 运营界的“全能卷王”：从供应链到星巴克排队，没有POM不敢管的闲事。
- (2) “数据要多，故事要野”：在这里，库存管理能扯上宇宙膨胀理论才管窥。

### 21. Journal of International Business Studies (JIBS)

- (1) 国际商务的“文化大使”：研究跨国公司的编辑，可能连自己护照都找不着。
- (2) “全球化是个筐，啥都往里装”：在这里，研究本地企业就像在肯德基点麦当劳——不合规矩。
- (3) 审稿人的灵魂提问：“你这研究，能解释为什么美国人觉得皮蛋是恶魔食物吗？”
- (4) 投稿建议：找个全球化问题，凑够跨文化数据，最后用一篇“文化风”论文打动审稿人。

### 22. Management Science (MS)

- (1) 管理学的“瑞士军刀”：从供应链到员工摸鱼，MS都能用数学模型管一管。
- (2) “生活可以乱，模型不能糙”：MS的编辑可能用优化算法决定今天穿哪只袜子。
- (3) 审稿人的执念：你的模型如果不能预测下一次金融危机，就是花瓶。
- (4) 投稿建议：找个管理问题，跑个复杂模型，最后用一篇“全能风”论文征服审稿人。

### 23. Operations Research (OR)

- (1) 运营界的“极客之王”：能用数学证明世界是虚拟的？OR会把你供上神坛。
- (2) “实践是理论的绊脚石”：在这里，模型落地不如发论文重要。
- (3) 审稿人的执念：你的模型如果不能预测下一次金融危机，就是花瓶。

### 17. Journal of Operations Management:

OM领域的“性价比之王”，比POMS好发但比MS难啃。

### 18. Production and Operations Management

Management: OM领域的“水龙头”，发文量多到被怀疑放水。

### 19. Manufacturing and Service Operations Management

Management: 名字越长越没人懂，但发一篇能保终身教职。

### 20. Management Science: 管理界的“瑞士军刀”

从供应链到恋爱博弈啥都能裁。

### 21. Operations Research: 运筹学家的“玄学阵地”

审稿人可能自己都没看懂你的模型。

### 22. INFORMS Journal on Computing:

UTD24里的“备胎”，发它可能只是为了凑数，但总比没有强。

### 23. Organization Science: 跨学科“缝合怪”

要求“既懂博弈论又懂莎士比亚”。

### 24. Strategic Management Journal: 战略学者的“修仙秘籍”

但修到最后发现秘籍是审稿人写的。

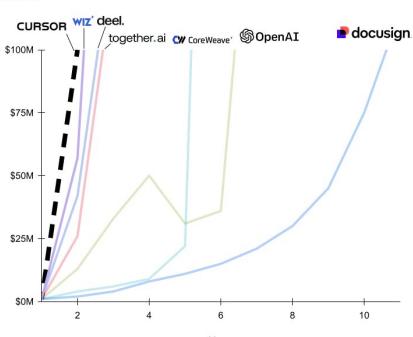
小红书号：195802739

7



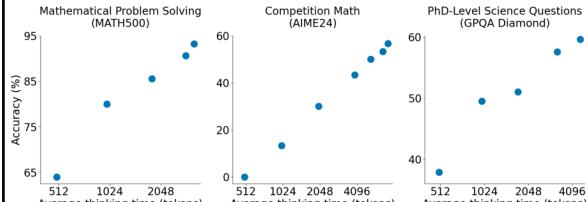
### CURSOR

Years from \$1M to \$100M ARR



**S1: Simple Test-Time Scaling**

<https://arxiv.org/abs/2501.19393>



## What's New (Feb/11/2025)

January 31, 2025

25 Comments


Paul Kassianik, Amin Karbasi

Security  
Evaluating Security Risk in DeepSeek and Other Frontier Reasoning Models

This original research is the result of close collaboration between AI security researchers from Robust Intelligence, now a part of Cisco, and the University of Pennsylvania including Yaron Singer, Amin Karbasi, Paul Kassianik, Mahdi Sabbagh, Hamed Hassani, and George Pappas.

Executive Summary

This article investigates vulnerabilities in DeepSeek R1, a new frontier reasoning model from Chinese AI startup DeepSeek. It has gained global attention for its advanced reasoning capabilities and cost-efficient training method. While its performance rivals state-of-the-art models like OpenAI o1, our security assessment reveals critical safety flaws.

Using algorithmic jailbreaking techniques, our team applied an automated attack methodology on DeepSeek R1 which tested it against 50 random prompts from the HarmBench dataset. These covered six categories of harmful behaviors including cybercrime, misinformation, illegal activities, and general harm.

The results were alarming: DeepSeek R1 exhibited a 100% attack success rate, meaning it failed to block a single harmful prompt. This contrasts starkly with other leading models, which demonstrated at least partial resistance.

**Exclusive: OpenAI co-founder Sutskever's SSI in talks to be valued at \$20 billion, sources say SSI: Safe Super Intelligence**

By Kenrick Cai, Krystal Hu and Anna Tong

February 8, 2025 12:24 AM GMT+8 · Updated 2 days ago

8

# More Importantly.....

**What can I help with?**

Message ChatGPT

0.00 / 3:31:23 - introduction >

Deep Dive into LLMs like ChatGPT

Andrej Karpathy 65K subscribers Subscribed

509,107 views Feb 4, 2025

This is a general audience deep dive into the Large Language Model (LLM) AI technology that powers ChatGPT and related products. It covers the full training stack of how the models are developed, along with mental models of how to think about their "psychology", and how to get the best use them in practical applications. I have one "Intro to LLMs" video already from ~year ago, but that is just a re-recording of a random talk, so I wanted to loop around and do a lot more comprehensive version.

Instructor  
Andrej was a founding member at OpenAI (2015) and then Sr. Director of AI at Tesla (2017-2022), and is now a founder at Eureka Labs, which is building an AI-native school. His goal in this video is to raise knowledge and understanding of the state of the art in AI, and empower people to effectively use the latest and greatest in their work.

Find more at <https://karpathy.ai> and <https://x.com/karpathy>

Chapters

- 00:00:00 introduction
- 00:01:00 pretraining data (internet)
- 00:07:47 tokenization
- 00:14:27 neural network I/O
- 00:20:11 neural network internals
- 00:26:01 inference
- 00:31:09 GPT2: training and inference
- 00:42:52 Llama 3.1 base model inference
- 00:59:23 pretraining to post-training
- 01:01:06 post-training data (conversations)
- 01:20:32 hallucination: tool use, knowledge/working memory
- 01:41:46 knowledge of self
- 01:46:56 models need tokens to think
- 02:01:11 tokenization revisited: models struggle with spelling
- 02:04:53 jagged intelligence
- 02:07:28 supervised finetuning to reinforcement learning
- 02:14:42 reinforcement learning
- 02:27:47 DeepSeek-R1
- 02:42:07 AlphaGo
- 02:48:26 reinforcement learning from human feedback (RLHF)
- 03:09:39 preview of things to come
- 03:15:15 keeping track of LLMs
- 03:18:34 where to find LLMs
- 03:21:46 grand summary

Links

- ChatGPT <https://chatgpt.com/>
- FineWeb (pretraining dataset): <https://huggingface.co/spaces/Hugging...>
- Tokkenizer: <https://tokkenizer.vercel.app/>
- Transformer Neural Net 3D visualizer: <https://bbycroft.net/llm>
- llm.c Let's Reproduce GPT2: <https://github.com/karpathy/llm/c/dls...>
- Llama 3 paper from Meta: <https://arxiv.org/abs/2407.21783>
- Hyperbolic, for inference of base model: <https://app.hyperbolic.xyz/>
- InstructGPT paper on SFT: <https://arxiv.org/abs/2203.02155>
- HuggingFace inference playground: <https://huggingface.co/spaces/hugging...>
- DeepSeek-R1 paper: <https://discovery.ucd.ac.id/eprint...>
- AlphaGo paper (PDF): <https://discovery.ucd.ac.id/eprint...>
- AlphaGo Move 37 video: [• Lee Sedol vs AlphaGo Move 37 reactio...](#)
- LM Arena for model rankings: <https://lmarena.ai>
- AI News Newsletter: <https://buttondown.com/ainews>
- LMStudio for local inference: [https://lmstudio.ai/](https://lmstudio.ai)
- The visualization UI I was using in the video: <https://excalidraw.com/>
- The specific file of Excalidraw we built up: <https://drive.google.com/file/d/1EZh5...>
- Discord channel for Eureka Labs and this video: [/ discord](#)

<https://www.youtube.com/watch?v=7xTGNNLPyMI>

9

Andrej Karpathy @karpathy

Part of the reason for my 3hr general audience LLM intro video is I hope to inspire others to make equivalents in their own domains of expertise, as I'd love to watch them.

5:43 AM · Feb 8, 2025 · 533.2K Views

@p10dushyanthac 4 days ago

In an era where almost everyone is trying to monetize every ounce of knowledge they have, this pure genius is releasing a 3-hour, 31-minute, and 23-second-long video for free (I am sure it would have taken at least five times that to compile this content). His videos are the most knowledge-rich content available on this topic. Andrej Karpathy, you are my favorite person on the internet. More power to you, brother.

1,839 Reply

@ayogheswaran9270 4 days ago (edited)

Andrej Uploads a video on YT.  
Everyone else: immediately pauses whatever they are doing. Lock themselves in a room and start watching the video.

Thanks a lot Andrej!! RESPECT.

351 Reply

6 replies

Some Comments

10

# What's New (Feb/18/2025)

**OpenAI** #50 → #1

2024 → 2025

What do you want to know?

Ask anything...

Deep Research

用提问发现世界

输入你想问的...

深度思考 DeepSeek R1

Competitive Programming with Large Reasoning Models

OpenAI\*

<https://arxiv.org/abs/2502.06807>

**Abstract**

We show that reinforcement learning applied to large language models (LLMs) significantly boosts performance on complex coding and reasoning tasks. Additionally, we compare two general-purpose reasoning models — OpenAI o1 and an early checkpoint of o3 — with a domain-specific system, o1-oi, which uses hand-engineered inference strategies designed for competing in the 2024 International Olympiad in Informatics (IOI). We competed live at IOI 2024 with o1-oi and, using hand-crafted test-time strategies, placed in the 49th percentile. Under relaxed competition constraints, o1-oi

**IOI Performance by Submission Strategy**

| Submission Strategy     | Score  |
|-------------------------|--------|
| o1-oi (50-Submissions)  | 213    |
| o1-oi (10K-Submissions) | 362.14 |
| O3 (50-Submissions)     | 395.64 |

11

# What's New (Feb/25/2025)

<https://x.ai/blog/grok-3>

**Benchmarks**

| Category            | Grok-3 | Grok-3 mini | Gemini-2 Pro | DeepSeek-V3 | Claude 3.5 Sonnet | GPT-4o |
|---------------------|--------|-------------|--------------|-------------|-------------------|--------|
| Math(AIME '24)      | 55     | 45          | 35           | 25          | 15                | 10     |
| Science(GPQA)       | 75     | 65          | 65           | 65          | 55                | 50     |
| Coding(LCB Oct-Feb) | 57     | 45          | 45           | 45          | 35                | 30     |

**Summary.** As far as a quick vibe check over ~2 hours this morning, Grok 3 + Thinking feels somewhere around the state of the art territory of OpenAI's strongest models (o1-pro, \$200/month), and slightly better than DeepSeek-R1 and Gemini 2.0 Flash Thinking. Which is quite incredible

DeepSeek [@deepseek.ai](#) Day 0: Warming up for #OpenSourceWeek!

These humble building blocks in our online service have been documented, deployed and battle-tested in production.

As part of the open-source community, we believe that every line shared becomes collective momentum that accelerates the journey.

Daily unlocks are coming soon. No ivory towers - just pure garage-energy and community-driven innovation.

11:00 PM · Feb 20, 2025 · 73.6K Views

**Thinking Machines** [@thinkymachines](#) Follow ...

Today, we are excited to announce Thinking Machines Lab ([thinkingmachines.ai](https://thinkingmachines.ai)), an artificial intelligence research and product company. We are scientists, engineers, and builders behind some of the most widely used AI products and libraries, including ChatGPT, Character.ai, PyTorch, and Mistral. Our mission is to make artificial intelligence work for you by building a future where everyone has access to the knowledge and tools to make AI serve their unique needs.

We are committed to open science through publications and code releases, while focusing on human-AI collaboration that serves diverse domains. Our approach embraces co-design of research and products to enable learning from real-world deployment and rapid iteration. This work requires three core foundations: state-of-the-art model intelligence, high-quality infrastructure, and advanced multimodal capabilities. We are committed to building models at the frontier of capabilities to deliver on this promise.

If you're interested in joining our team, consider applying here: [6wajk07p.paperform.co](https://6wajk07p.paperform.co)

2:34 AM · Feb 19, 2025 · 1.5M Views

Grok 3用20万GPU帮AI界做了个实验：Scaling Law 没撞墙，但预训练不一定

腾讯科技 02-20 08:34 <https://wallstreetcn.com/articles/3741468>

12

# Flip Side of the Coin

<https://nosegauge.substack.com/p/capitalagi-and-human-ambition>

|                     |  |
|---------------------|--|
| <b>No Set Gauge</b> | <b>Money currently struggles to buy talent</b> |
|---------------------|--|

**Capital, AGI, and human ambition**

AGI will shift the relative importance of human v non-human factors of production, reducing the incentive to care about humans while making existing powers more effective and entrenched

L RUDOLF L DEC 29, 2024

192 37 31 Share

Edited to add: The main takeaway of this post is meant to be: Labour-replacing AI will shift the relative importance of human v non-human factors of production, which reduces the incentives for society to care about humans while making existing powers more effective and entrenched. Many people are reading this post in a way where either (a) "capital" means just "money" (rather than also including physical capital like factories and data centres), or (b) the main concern is human-human inequality (rather than broader societal concerns about humanity's collective position, the potential for social change, and human agency).

Money can buy you many things: capital goods, for example, can usually be bought quite straightforwardly, and cannot be bought without a lot of money (or other liquid assets, or non-liquid assets that others are willing to write contracts against, or special government powers). But it is surprisingly hard to convert raw money into labour, in a way that is competitive with top labour.

Consider Blue Origin versus SpaceX. Blue Origin was started two years earlier (2000 v 2002), had much better funding for most of its history, and even today employs almost as many people as SpaceX (11,000 v 13,000). Yet SpaceX has crushingly dominated Blue Origin. In 2000, Jeff Bezos had \$4.7B at hand. But it is hard to see what he could've done to not lose out to the comparatively money-poor SpaceX with its intense culture and outlier talent.

**Most people's power/leverage derives from their labour**

Labour-replacing AI also deprives almost everyone of their main lever of power and leverage. Most obviously, if you're the average Joe, you have money because someone somewhere pays you to spend your mental and/or physical efforts solving their problems.

13

## What's New (Mar/4/2025)

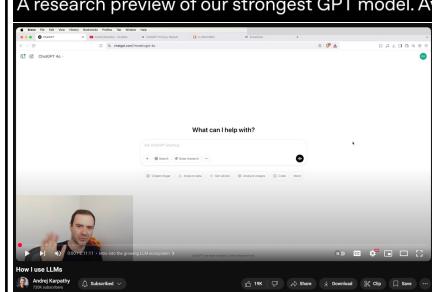
**Claude 3.7 Sonnet and Claude Code**

25 Feb 2025 • 5 min read

February 27, 2025 Release Product

### Introducing GPT-4.5

A research preview of our strongest GPT model. Available to Pro users and developers worldwide.



**Andrej Karpathy** @karpathy ...  
After many hours of scrutinizing humor in LLM outputs, this one by Claude 3.7 is the funniest by far.

**Tibo** @tibo\_maker - Feb 28  
LOL!!

Claude (via Cursor) randomly tried to update the model of my feature from OpenAI to Claude 🤪  
...  
Show more

```
const params = {
  messageHistory: messages,
  model: "gpt-4",
  model: "claude-3-7-sonnet-latest",
  temperature: 1.1,
  maxTokens: 2000,
```

2:59 AM - Mar 1, 2025 392.4K Views

**Andrej Karpathy** @karpathy ...  
OKAY so I didn't super expect the results of the GPT4 vs. GPT4.5 poll from earlier today 😅, of this thread: x.com/karpathy/status...

Question 1: GPT4.5 is A; 56% of people prefer it.  
Question 2: GPT4.5 is B; 43% of people prefer it.  
Question 3: GPT4.5 is A; 35% of people prefer it.  
Question 4: GPT4.5 is A; 35% of people prefer it.  
Question 5: GPT4.5 is B; 36% of people prefer it.

TLDR people prefer GPT4 in 4/5 questions awkward.

14

## What's New (Mar/11/2025)



**AWARDS & RECOGNITION**

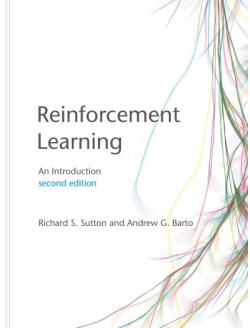
**Andrew Barto and Richard Sutton Receive 2024 ACM A.M. Turing Award**

Andrew G. Barto and Richard S. Sutton received the 2024 ACM A.M. Turing Award for developing the conceptual and algorithmic foundations of reinforcement learning. In a series of papers beginning in the 1980s, Barto and Sutton introduced the main ideas, constructed the mathematical foundations, and developed important algorithms for reinforcement learning—one of the most important approaches for creating intelligent systems. Barto is Professor Emeritus of Information and Computer Sciences at the University of Massachusetts, Amherst. Sutton is a Professor of Computer Science at the University of Alberta, a Research Scientist at Keen Technologies, and a Fellow at Amii (Alberta Machine Intelligence Institute).

**Reinforcement Learning**

An Introduction second edition

Richard S. Sutton and Andrew G. Barto



**ALPHAGO**



**LLM**



**manus**  
The general AI agent

<https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf>

**Reinforcement Learning**

Sometimes called **Behaviorism AI**

Action



**Reward**

Opensource Replications:

<https://github.com/mannaandpoem/OpenManus>  
<https://github.com/camel-ai/owl>

Insightful Comments:

<https://yage.ai/genai/agent2025/>  
<https://www.superlinear.academy/c/share-your-work/manus-url>  
<https://yage.ai/manus.html>  
<https://yage.ai/agentic-memory.html>

15

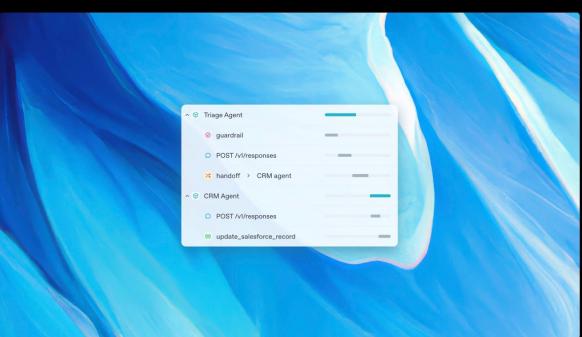
## What's New (Mar/18/2025)

March 11, 2025 Product

### New tools for building agents

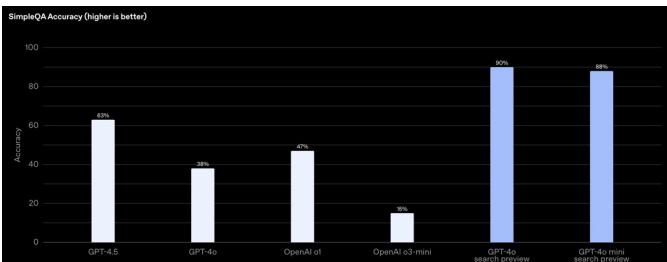
We're evolving our platform to help developers and enterprises build useful and reliable agents.

Try in Playground ↗



<https://openai.com/index/new-tools-for-building-agents/>

**The Most Important Tool-Use: Web Search**

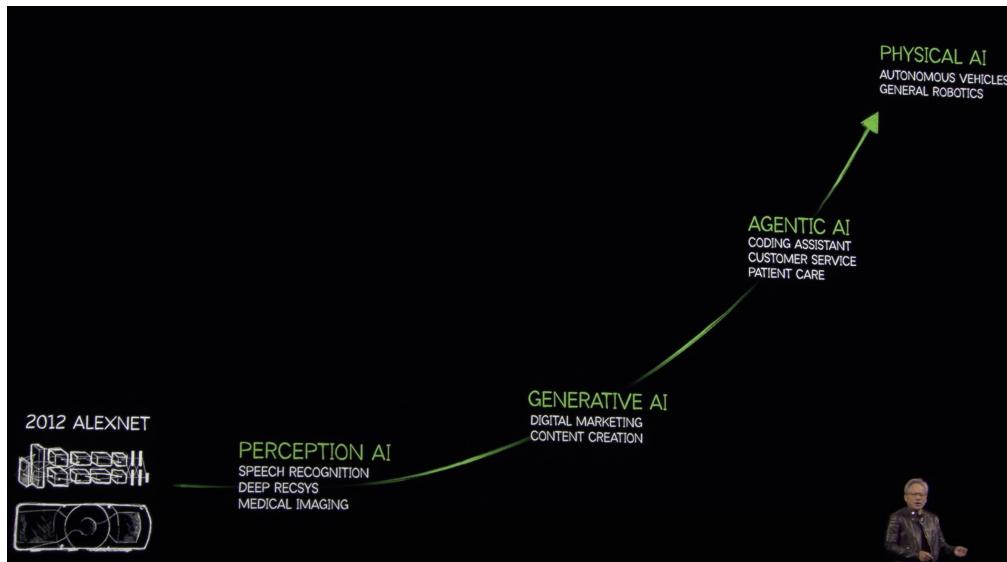


| Model                      | Accuracy (%) |
|----------------------------|--------------|
| GPT-4.5                    | 45%          |
| GPT-4o                     | 38%          |
| OpenAI o1                  | 47%          |
| OpenAI o5-mini             | 8%           |
| QP14o search preview       | 90%          |
| GPT-4o mini search preview | 88%          |

16

## What's New (Mar/25/2025)

### Different Generations of AI

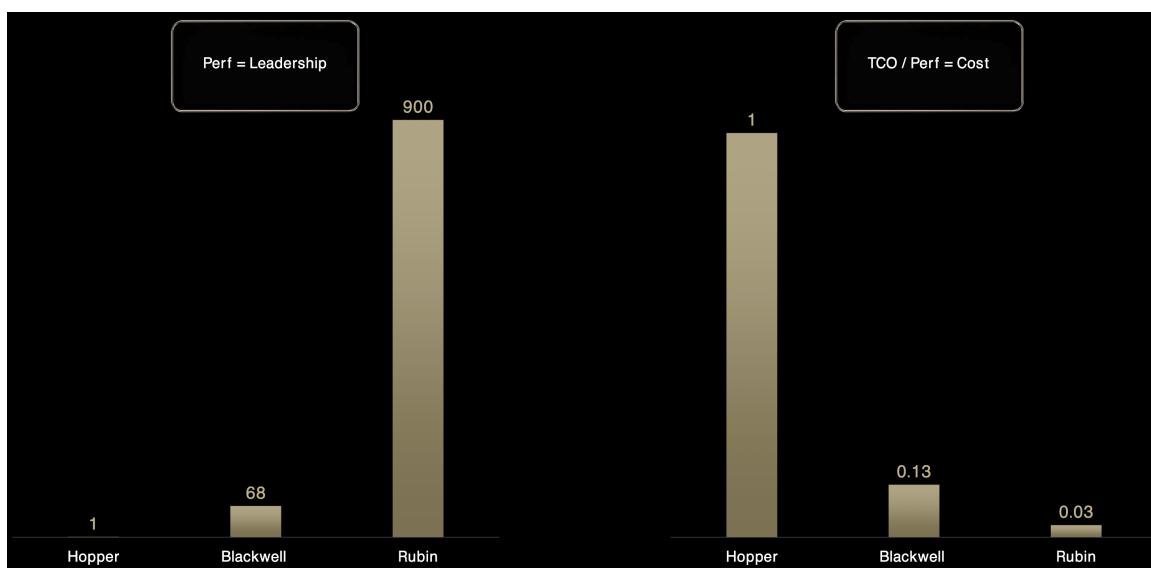


Nvidia GTC Keynote by Jenson Huang: <https://www.rev.com/transcripts/gtc-keynote-with-nvidia-ceo-jensen-huang>

17

## What's New (Mar/25/2025)

### From Hopper to Blackwell and to Rubin: Stronger but Cheaper



Nvidia GTC Keynote by Jenson Huang: <https://www.rev.com/transcripts/gtc-keynote-with-nvidia-ceo-jensen-huang>

18

## What's New (Mar/25/2025)

**Top 5 trends at a glance**

**Trend 01** Multimodal AI: Unleash the power of context

**Trend 02** AI agents: The evolution from chatbots to multi-agent systems

**Trend 03** Assistive search: The next frontier for knowledge work

**Trend 04** AI-powered customer experience: So seamless, it's almost invisible

**Trend 05** Security gets tighter—and tougher—with AI

AI Trends by Google: [https://services.google.com/fh/files/misc/google\\_cloud\\_ai\\_trends.pdf](https://services.google.com/fh/files/misc/google_cloud_ai_trends.pdf)

19

## What's New (Apr/1/2025)

March 25, 2025 Product Release

### Introducing 4o Image Generation

Unlocking useful and valuable image generation with a natively multimodal model capable of precise, accurate, photorealistic outputs.

Try in ChatGPT ➔

Listen to article 5:49 Share

At OpenAI, we have long believed image generation should be a primary capability of our language models. That's why we've built our most advanced image generator yet into GPT-4o. The result—image generation that is not only beautiful, but useful.

3月29日，科幻作家、《三体》作者刘慈欣在接受采访时被问到DeepSeek未来有可能替代科幻作家吗？刘慈欣表示，暂时不太会，但是再过10年、20年，从理论上说完全可能代替科幻小说作家。

他认为，从科学的角度去讲，所有人类作家的身上没有什么是不可被AI所替代的。在谈到该如何去应对这种局面时，刘慈欣称：“我个人认为首先停止自我安慰，坦然去面对技术的冲击以及这种冲击对我们领域的那种天翻地覆的影响。”

在2025年中国科幻大会上，刘慈欣明确指出AI将重塑科幻创作生态，可能使大部分作家被取代，仅少数具有“巅峰创造力”的作家暂时难以替代。刘慈欣还设想，未来阅读或采用AI定制，读者向AI提出需求，就能生成小说，不满意还能重新生成，那时AI写作质量可能全面超越人类。

他也提到，目前的AI创作仍存在局限性，比如缺乏真正的情感体验、社会洞察力和独特的创造力。刘慈欣呼吁创作者停止用“灵魂”“情感”等概念自我安慰，主张直面技术革命带来的根本性改变。

A wide image taken with a phone of a glass whiteboard, in a room overlooking the Bay Bridge. The field of view shows a woman writing, sporting a tshirt with a large OpenAI logo. The handwriting looks natural and a bit messy, and we see the photographer's reflection. ... Read more

Transfer between Modalities:  
Suppose we directly model (text, pixels, sound)  
with one big autoregressive transformer.  
Pros:  
- image generation augmentation  
- next level text rendering  
- native in-context learning  
- unified post-training  
Cons:  
- varying bit-rate requirements  
- compute not adaptive

Fixes:  
+ model compressed representations  
+ compose autoregressive prior with a powerful decoder

tokens → [transformer] → [diffusion] pixels

刘慈欣谈DeepSeek：完全有可能取代科幻小说作家

20

## What's New (Apr/1/2025)

(Interpretability)

### Tracing the thoughts of a large language model

Mar 27, 2025

[Read the paper](#)

AI Tracing the thoughts of a large language model

Tracing the thoughts of an LLM

Watch on YouTube

AI Biology: <https://www.anthropic.com/research/tracing-thoughts-language-model>  
 Physics of LLM: <https://physics.allen-zhu.com/>

21

## What's New (Apr/8/2025)

### AI 2027

Daniel Kokotajlo, Scott Alexander, Thomas Larsen, Eli Lifland, Romeo Dean

We predict that the impact of superhuman AI over the next decade will be enormous, exceeding that of the Industrial Revolution.

We wrote a scenario that represents our best guess about what that might look like.<sup>1</sup> It's informed by trend extrapolations, wargames, expert feedback, experience at OpenAI, and previous forecasting successes.<sup>2</sup>

[What is this?](#) [How did we write it?](#) [Why is it valuable?](#) [Who are we?](#)

**Mid 2025: Stumbling Agents**  
**Late 2025: The World's Most Expensive AI**  
**Early 2026: Coding Automation**  
**Mid 2026: China Wakes Up**  
**Late 2026: AI Takes Some Jobs**  
**Jan 2027: Agent-2 Never Finishes Learning**  
**Apr 2027: Alignment for Agent-3**  
**Jul 2027: The Cheap Remote Worker**  
**Sep 2027: Agent-4, the Superhuman AI Researcher**  
<https://ai-2027.com/>; <https://ai-2027.com/scenario.pdf>

22

# What's New (Apr/8/2025)



We're releasing PaperBench, a benchmark evaluating the ability of AI agents to replicate state-of-the-art AI research, as part of our Preparedness Framework.

Agents must replicate top ICML 2024 papers, including understanding the paper, writing code, and executing experiments.

## PaperBench: Evaluating AI's Ability to Replicate AI Research

Giulio Starace<sup>\*</sup>, Oliver Jaffe<sup>\*</sup>, Dane Sherburn<sup>\*</sup>, James Aung<sup>\*</sup>, Chan Jun Shern<sup>\*</sup>, Leon Maksin<sup>\*</sup>, Rachel Dias<sup>\*</sup>, Evan Mays<sup>\*</sup>, Benjamin Kinsella<sup>\*</sup>, Johannes Heidecke<sup>\*</sup>, Amelia Gliese<sup>\*</sup>, Tejal Patwardhan<sup>\*</sup>  
OpenAI

### Abstract

We introduce PaperBench, a benchmark evaluating the ability of AI agents to replicate state-of-the-art AI research. Agents must replicate 20 ICML-2024 Spotlight and Oral papers from scratch, including understanding paper contributions, developing a codebase, and successfully executing experiments. For objective evaluation, we develop rubrics that automatically decompose each replication task into smaller sub-tasks with clear grading criteria. In total, PaperBench contains 8,316 individually gradable tasks. Rubrics are co-developed with the author(s) of each ICML paper for accuracy and realism. To enable scalable evaluation, we also develop an LLM-based judge to automatically grade submissions, attempt against rubrics, and assess our judge's performance by creating a separate benchmark for judges. We evaluate several frontier models on PaperBench, finding that the best-performing tested agent, Claude 3.5 Sonnet (New) with open-source scaffolding, achieves an average replication score of 8.3%.

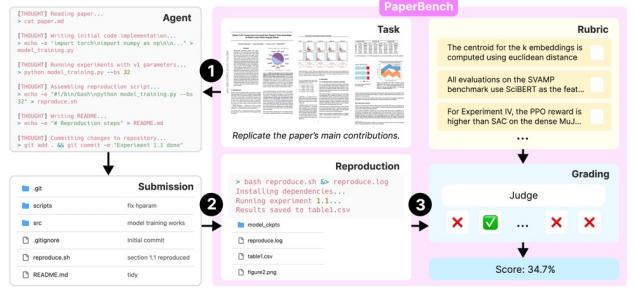


Figure 1. PaperBench is a benchmark for evaluating AI agents' abilities to replicate AI research. Each sample includes a research paper and a grading rubric that specifies the assessment criteria for a complete replication. Agents create a codebase from scratch as their submission (1), which is then executed to verify reproduction (2) and graded against the rubric by an LLM-based judge (3).

Table 4. Average Replication Scores (in %) for models with BasicAgent, our main setup. Error is one standard error of the mean.

| MODEL             | PAPERBENCH |
|-------------------|------------|
| O3-MINI-HIGH      | 2.6 ± 0.2  |
| GPT-4o            | 4.1 ± 0.1  |
| GEMINI-2.0-FLASH  | 3.2 ± 0.2  |
| DEEPSEEK-R1       | 6.0 ± 0.3  |
| O1-HIGH           | 13.2 ± 0.3 |
| CLAUDE-3.5-SONNET | 21.0 ± 0.8 |

Table 5. Average Replication Scores (in %) with IterativeAgent. IterativeAgent removes the ability of models to end the task early and prompts models to work in a piecemeal fashion. We observe that these modifications significantly boost scores for o3-mini and o1 compared to BasicAgent, but hamper Claude 3.5 Sonnet, highlighting models' sensitivities to prompting.

| MODEL                                 | PAPERBENCH |
|---------------------------------------|------------|
| O3-MINI-HIGH                          | 8.5 ± 0.8  |
| CLAUDE-3.5-SONNET                     | 16.1 ± 0.1 |
| O1-HIGH                               | 24.4 ± 0.7 |
| <i>With an extended 36 hour limit</i> |            |
| O1-HIGH                               | 26.0 ± 0.3 |