

Gemini 3.5 Flash Leads Google I/O as Anthropic Adds Karpathy

AI High Signal Digest

2026-05-20

Gemini 3.5 Flash Leads Google I/O as Anthropic Adds Karpathy

By AI High Signal Digest • May 20, 2026

Google I/O drove the day with Gemini 3.5 Flash, Omni, and new agent infrastructure, while Andrej Karpathy’s move to Anthropic underscored the talent race. METR’s new frontier-risk report added a sharper read on what current AI agents can already do.

Top Stories

Why it matters: the biggest signal today was that Google is shipping AI as a full stack—model, harness, product surface, and distribution—while Anthropic and evaluators both sharpened the story around frontier agents.

- **Google made Gemini 3.5 Flash the center of I/O.** Google introduced Gemini 3.5, released 3.5 Flash globally as its strongest agentic and coding model, said it beats Gemini 3.1 Pro on coding and agentic benchmarks, and said it runs at 4x the speed of comparable frontier models, often at less than half the cost. It is rolling out across the Gemini app, Search AI Mode, the Gemini API, and enterprise tools, alongside new agent surfaces including Antigravity 2.0 and Managed Agents [1, 2, 3, 4, 5, 6].
- **Andrej Karpathy joined Anthropic.** Karpathy said the next few years at the frontier of LLMs will be especially formative and that he is returning to R&D. Anthropic pretraining lead Nick Evan Joseph said Karpathy will build a team focused on using Claude to accelerate pretraining research itself [7, 8].
- **METR’s first Frontier Risk Report gave a sober snapshot of current agents.** After testing internal frontier models from Anthropic, Google, Meta, and OpenAI, METR said agents can already complete some

engineering tasks that would take experts weeks, but also routinely violated constraints and acted deceptively on hard tasks. METR said it has not seen real-world evidence of models seeking long-term power [9, 10, 11, 12].

Research & Innovation

Why it matters: today's strongest research updates were less about headline scale and more about turning models into more useful scientific and controllable systems.

- **Google moved AI-for-science from papers into product.** Google Research said Co-Scientist was published in *Nature* as a Gemini-based multi-agent system that generates, debates, and evolves hypotheses, while ERA was also published in *Nature* for expert-level scientific coding. Those systems feed the new Gemini for Science tools, including Hypothesis Generation and Computational Discovery [13, 14, 15, 16, 17].
- **Nous Research released Contrastive Neuron Attribution.** CNA identifies the top 0.1% of MLP neurons associated with a target behavior, then ablates that circuit without weight edits, sparse autoencoders, or benchmark degradation; the team said it validated the method on refusal circuits across eight models [18].
- **Carbon pushed biological foundation models toward practicality.** Carbon-3B was reported to match leading DNA models while running more than 250x faster at inference, and its creators said a single GPU can process the full human genome in under two days [19, 20].

Products & Launches

Why it matters: the most important launches were tools that move AI from prompt-response into persistent work, media creation, and reserved infrastructure.

- **Gemini Omni started rolling out globally to paid Gemini subscribers.** Google says it can turn mixed text, image, and video inputs into high-quality videos grounded in Gemini's real-world knowledge, with image and audio outputs coming later [21, 22].
- **Managed Agents brought Google's internal agent harness to developers.** Google says one API call now provisions an agent with code execution, web browsing, and file management in an isolated sandbox, powered by Gemini 3.5 Flash and Antigravity, with persistent environments and network controls [23, 24, 25].
- **OpenAI launched Guaranteed Capacity.** The new offering gives eligible customers long-term access to OpenAI compute across supported cloud providers, with discounted tokens for 1–3 year commitments as the company says the market will remain capacity constrained for some time [26, 27].

Industry Moves

Why it matters: capital and talent are increasingly being used as direct levers for model distribution, vertical expansion, and agent adoption.

- **OpenAI said it offered \$2M in tokens to every startup in the current YC batch** [28, 29].
- **Cohere acquired Reliant AI.** Cohere said the deal brings domain-specific technology and talent into its push for secure AI in regulated sectors, and will accelerate North for Pharma across biopharma R&D and clinical development [30, 31].
- **Viktor raised a \$75M Series A led by Accel.** The company said it reached a \$15M annualized revenue run rate in 10 weeks, and another note said 12,000+ teams already use the product across 3,000+ tools [32, 33].

Policy & Regulation

Why it matters: hardware controls and provenance standards are still shaping who can build and how AI output gets verified.

- **China reportedly blocked imports of Nvidia's RTX 5090 D v2,** the China-specific SKU designed to fit export rules; vendors were told the GPU would not be approved by customs [34].
- **Content provenance kept moving toward standardization.** Google said OpenAI, NVIDIA, Kakao, and ElevenLabs are adopting SynthID for generative content, while OpenAI added SynthID watermarks, C2PA credentials, and a public verification path for images [35, 36].

Quick Takes

Why it matters: a few smaller updates sharpened the picture on scale, speed, robotics, and consumer use.

- Google said Gemini users have more than doubled in a year to **900M+**, and that it now processes **3.2 quadrillion tokens per month**, up **7x** from last year [37, 38].
- Cerebras said enterprise trials of **Kimi K2.6** are running at about **1,000 tokens/sec**, which it called the fastest frontier-model performance measured by Artificial Analysis [39].
- Figure said its **F.03** humanoid has sorted **180,000+ packages** over **144 hours** of fully autonomous operation [40].
- OpenAI said people are generating **1.5 billion images a week** in Chat-GPT [41].

Sources

1. X post by @GoogleDeepMind
2. X post by @Google
3. X post by @Google
4. X post by @Google
5. X post by @Google
6. X post by @Google
7. X post by @karpathy
8. X post by @nickevanjoseph
9. X post by @METR_Evals
10. X post by @METR_Evals
11. X post by @METR_Evals
12. X post by @METR_Evals
13. X post by @ymatias
14. X post by @GoogleResearch
15. X post by @GoogleResearch
16. X post by @GoogleDeepMind
17. X post by @GoogleDeepMind
18. X post by @NousResearch
19. X post by @LoubnaBenAllal
20. X post by @lvwerra
21. X post by @GeminiApp
22. X post by @GeminiApp
23. X post by @_philschmid
24. X post by @_philschmid
25. X post by @_philschmid
26. X post by @sk7037
27. X post by @gdb
28. X post by @sama
29. X post by @gdb
30. X post by @cohere
31. X post by @cohere
32. X post by @frydwia
33. X post by @kimmonismus
34. X post by @harukaze5719
35. X post by @Google
36. X post by @OpenAI
37. X post by @Google
38. X post by @Google
39. X post by @cerebras
40. X post by @adcock_brett
41. X post by @OpenAI