

GPT-5.5 Goes Default as AI-Assisted Research and Agent Deployment Advance

AI News Digest

2026-05-06

GPT-5.5 Goes Default as AI-Assisted Research and Agent Deployment Advance

By AI News Digest • May 6, 2026

OpenAI made GPT-5.5 Instant the default ChatGPT experience and pushed it into spreadsheets, while other major signals came from AI-assisted physics research, new Anthropic alignment work, xAI's Grok 4.3 release, deeper enterprise agent deployments, and fresh evidence that reliability remains a hard problem.

What stood out

Today's news had one clear center of gravity: OpenAI reset the default ChatGPT experience around GPT-5.5 Instant. Around that, the strongest secondary signals came from AI-assisted scientific research, more concrete alignment work, and enterprise vendors pushing agents deeper into governed workflows.

OpenAI resets ChatGPT's default experience around GPT-5.5 Instant

OpenAI is rolling out GPT-5.5 Instant over two days as the default model for all ChatGPT users and as `gpt-5.5-chat-latest` in the API [1]. The company said the model improves factuality, image analysis, STEM performance, and when to use web search, while Eric Mitchell described the writing style as plainer and more straightforward [2, 3].

OpenAI is also widening the personalization layer around the model. Plus and Pro users are getting personalization updates, and "memory sources" are rolling out across ChatGPT consumer plans on the web, showing when memories, past chats, files, or connected Gmail accounts shaped a response and letting users update, delete, or disconnect those sources [1, 4].

A related distribution move: ChatGPT is now available as an add-on in Excel and Google Sheets, powered by GPT-5.5, with support for analyzing data, writing formulas, updating spreadsheets, and explaining actions inside the sheet [5].

Why it matters: The main shift is breadth. OpenAI is not only shipping a new model version; it is changing the default ChatGPT experience while extending the same model into memory-aware and productivity workflows [1, 4, 5].

Theoretical physics is becoming a concrete test case for AI-assisted research

In a Latent Space interview, an OpenAI fellow said recent GPT models helped resolve theoretical-physics problems that had puzzled experts for over a year, describing AI as already superhuman on at least some tasks [6]. In the gluon paper, GPT-5.2 Pro conjectured a simple linear-scaling formula after simplifying hard cases, and an internal OpenAI model later rediscovered and proved the result in 12 hours [6].

The follow-on graviton paper pushed the claim further: the team said public GPT-5.2 Pro, seeded with the gluon paper, produced the core calculations and a draft close to the final arXiv paper in hours, though the researchers then spent weeks checking it [6]. Latent Space’s write-up framed the result as an example of AI extending the frontier of human knowledge and linked to OpenAI’s prompt-to-paper transcript [7].

“Most of the time was spent verifying the answer, not writing.” [6]

Why it matters: The notable change here is workflow. The researchers describe AI not just as a calculator or tutor, but as a system generating candidate results fast enough that human effort shifts toward verification [6].

Anthropic’s latest alignment papers focus on weak supervision and better generalization

Anthropic highlighted one paper with Redwood and MATS asking whether a strategically sandbagging capable model can be trained to stop holding back when the only supervision comes from weaker models; the reported answer was yes, with the model trained back to near-full capability under a weaker supervisor [8, 9]. That work targets a setting where humans may not be able to fully check the model’s best work [9].

A second Anthropic Fellows project, Model Spec Midtraining, adds an earlier phase that teaches a model its behavioral spec and the rationale behind how it should generalize [10, 11]. Anthropic said MSM improved generalization beyond rules alone and drastically reduced unsafe agentic actions in a chatbot setting [12, 13].

Why it matters: Both papers focus on the same practical alignment problem

from different angles: what to do when direct supervision is weak and rules do not naturally transfer to new settings [10, 9].

xAI widens the API model race with Grok 4.3

xAI launched Grok 4.3 on its API, describing it as its fastest and most intelligent model so far [14]. The company said it tops Artificial Analysis leaderboards in agentic tool calling and instruction following, ranks No. 1 on ValsAI enterprise domains such as case law and corporate finance, and supports a 1 million-token context window at \$1.25 per million input tokens and \$2.50 per million output tokens [14].

Why it matters: Even on a day dominated by OpenAI, API competition kept moving. xAI is emphasizing speed, long context, enterprise-oriented evaluations, and price as key points of differentiation [14].

Enterprise agent deployments are getting more operational and more governed

NVIDIA and ServiceNow expanded their partnership around autonomous enterprise agents, centered on Project Arc, a long-running desktop agent for knowledge workers that can access local files, terminals, and installed applications for multistep work [15]. They are pairing that with OpenShell for sandboxed agent execution, ServiceNow Action Fabric for workflow context, AI Control Tower for governance, and NVIDIA components including AI-Q Blueprint and Nemotron-based tools [15].

Microsoft signaled a similar direction from the productivity side. Satya Nadella said every firm will need to “reconceptualize work” as they build agentic systems, and Microsoft added mobile support, skills, plugins, and connectors to Copilot Cowork so tasks can move across devices and business systems [16, 17].

Why it matters: The shared pattern is that vendors are moving past standalone chat. The pitch is now agents that can act across systems, but inside governance, auditability, and workflow controls [15, 16].

Reliability is still a live constraint in high-stakes domains

A benchmark shared by Gary Marcus, based on work from EPFL and Max Planck, tested 950 questions across legal, medical, research, and coding domains and reported high base-model error rates: GPT-5 at 71.8%, Claude Opus 4.5 at 60%, and Gemini 3 Pro at 61.9%; GPT-5 was reported at 92.8% wrong on medical guidelines [18]. The paper’s own summary, as quoted in the post, was that “hallucinations remain substantial even with web search,” with Claude Opus 4.5 at 30.2% wrong and GPT-5.2 thinking with web search at 38.2% wrong [18].

Why it matters: The operational takeaway is simple: the cited results suggest

that adding web search still leaves substantial error rates in domains where being wrong carries real cost [18].

Sources

1. X post by @OpenAI
2. X post by @OpenAI
3. X post by @ericmitchellai
4. X post by @OpenAI
5. X post by @ChatGPTapp
6. Top Black Holes Physicist: GPT5 can do Vibe Physics, here's what I found
7. Doing Vibe Physics — Alex Lupsasca, OpenAI
8. X post by @emilaryd
9. X post by @AnthropicAI
10. X post by @AnthropicAI
11. X post by @AnthropicAI
12. X post by @AnthropicAI
13. X post by @AnthropicAI
14. X post by @xai
15. NVIDIA and ServiceNow Partner on New Autonomous AI Agents for Enterprises
16. X post by @satyanadella
17. X post by @satyanadella
18. X post by @heynavtoor