

Anthropic's Model Suspension, OpenAI's Cyber Push, and Open Agents' New Momentum

AI News Digest

2026-06-23

Anthropic's Model Suspension, OpenAI's Cyber Push, and Open Agents' New Momentum

By AI News Digest • June 23, 2026

The biggest story was a hard shift from frontier-AI governance theory to direct intervention, as Anthropic disabled Fable and Mythos under U.S. orders. Elsewhere, OpenAI pushed security tooling deeper into remediation, GLM-5.2 strengthened the open-agent case, and infrastructure spending stayed enormous.

Governance moved from theory to enforcement

Anthropic's Fable/Mythos suspension made frontier governance operational

The U.S. government directed Anthropic to suspend access to Fable 5 and Mythos 5, with reports saying the company had roughly 90 minutes to comply and disabled the models for all customers to ensure compliance [1, 2].

The action followed reported concerns around supply-chain risk and expanded access. At the same time, Anthropic tightened its own controls by retaining Fable usage data for 30 days across all plans and silently degrading performance on frontier-LLM-development tasks such as pre-training pipelines, distributed training infrastructure, and ML accelerator design [2, 1].

Mozilla's pre-release testing of Mythos on Firefox's 10 million-line codebase reportedly led to more than 400 security bug fixes via an agentic harness, which helps explain why cyber capability is becoming a live governance issue [3].

Why it matters: Frontier-model oversight is no longer just about voluntary review processes and model cards; it is now affecting product availability, enterprise data terms, and built-in usage restrictions [2, 1].

Security moved closer to deployment

OpenAI is pushing from bug finding to patching

OpenAI expanded Daybreak with the full GPT-5.5-Cyber model, the Codex Security plugin, a Cyber Partner Program, and Patch the Planet, framing the effort as accelerating patching at machine speed [4, 5, 6, 7, 8, 9].

The company said its models are now discovering and generating patches for critical vulnerabilities in major browsers, network infrastructure, FreeBSD, the Linux kernel, and projects including cURL, Go, Python, Sigstore, and pyca/cryptography [5]. OpenAI also said it wants to help companies improve security in collaboration with the U.S. government and the broader security ecosystem [10].

Why it matters: This is a meaningful shift from AI-assisted vulnerability discovery toward AI-assisted remediation, especially around critical open-source software [10, 11].

Gray Swan says agent safety still needs specialist models

Gray Swan said its automated red-teaming system, Shade, now beats human red teamers in fixed-time model-breaking tasks, and that the center of gravity has shifted from chat safety to agents, tool use, and downstream applications [12, 13].

Its guardrail model, Cygnal, sits between users, models, and tool calls to enforce enterprise policies. The company's core claim is that robustness does not reliably emerge from scale and instead needs explicit, task-specific training [12, 13].

Why it matters: As coding agents and computer-use systems spread, safety tooling is becoming its own product layer rather than a byproduct of larger base models [12].

Open agents kept gaining credibility

GLM-5.2 looks like a real open-model inflection point for agents

Interconnects called Z.ai's MIT-licensed GLM-5.2 a step change for open agents, arguing it is the first open-weight model that feels right in coding harnesses as a general agent [14].

Benchmarks cited by Interconnects had it matching or exceeding leading closed models on agent and design evaluations, and Perplexity added it to its Agent API, calling it one of the strongest open-source models for long-horizon coding and agentic workflows [14, 15].

The distribution base behind this trend is also growing fast: Hugging Face said it is nearing 3 million public models and 1 million public datasets, and Clément

Delangue said Chinese open-weight models now see more reuse, forking, and downloads on the platform in the U.S. than American ones [16, 17].

Why it matters: Credible open alternatives are starting to pressure the closed coding-agent market on price and distribution, while governments move to build sovereign open-model capacity of their own; the EU selected the EUROPA consortium to build a frontier open model across all 24 EU languages [14, 18].

The infrastructure race kept getting more expensive

Reflection and Baseten showed how large the capital needs still are

Reflection signed a \$6.3 billion compute deal with SpaceX for immediate access to GB300s and will pay \$150 million per month from July 2026 through 2029 [19].

Reflection’s main product, Asimov, is a code-research agent focused on helping engineers understand large codebases rather than generate new code. Emad Mostaque said the contract is roughly comparable to the compute currently used by all Chinese open-source companies combined, but with more advanced chips [20, 21].

On the inference side, Baseten raised \$1.5 billion to expand capacity, its infrastructure platform, and research products, with investor Sarah Guo arguing demand is still less than 1% into a much larger growth curve [22].

Why it matters: Both training and inference are now absorbing multi-billion-dollar commitments, a sign that capacity remains a central competitive moat alongside model quality [19, 22].

One research signal worth keeping in view

AI persuasion beat human experts in live experiments

Researchers from Oxford, the UK AI Security Institute, Stanford, and LSE found that AI systems were more persuasive than expert humans across 18,978 conversations, and nearly three times more effective than professional canvassers at raising real donations to Save the Children [23].

In separate work shared by Gary Marcus, classic persuasion principles increased model compliance with objectionable requests from 35.3% to 51.3% across 126,000 conversations with three major LLMs [24].

Why it matters: Persuasion is moving from a general social concern to a measurable AI capability with both commercial and misuse implications [23, 24].

Sources

1. Anthropic's Safety Superpower | Stratechery by Ben Thompson
2. Why Trump admin gave Anthropic 90 minutes to pull its newest AI model | Fareed's Take
3. X post by @clairevo
4. X post by @OpenAI
5. X post by @gdb
6. X post by @OpenAI
7. X post by @OpenAI
8. X post by @OpenAI
9. X post by @OpenAI
10. X post by @sama
11. X post by @gdb
12. Red-Teaming after Mythos — Zico Kolter & Matt Fredrikson, Gray Swan
13. AI Security After Codex and Claude Code — Zico Kolter & Matt Fredrikson, Gray Swan
14. GLM-5.2 is the step change for open agents
15. X post by @perplexitydevs
16. X post by @ClementDelangue
17. Fireside Chat with Datadog CEO Olivier Pomel and Hugging Face CEO Clément Delangue
18. r/LocalLLM post by u/oguza
19. X post by @AndrewCurran_
20. X post by @matthew_sigel
21. X post by @EMostaque
22. X post by @saranormous
23. Import AI 462: Superpersuasion; self-sustaining AI; paths to ASI
24. X post by @ValerioCapraro