

# Anthropic's Pentagon Fight and Nvidia's Shift to AI Factories

AI News Digest

2026-03-24

## Anthropic's Pentagon Fight and Nvidia's Shift to AI Factories

*By AI News Digest • March 24, 2026*

A consequential Anthropic-vs.-government fight led the day, alongside Nvidia's push toward secure rack-scale agent systems and clearer evidence that AI products are consolidating around integrated model-and-harness stacks. Research also sharpened the picture on cyber autonomy, model behavior, and how frontier systems should be evaluated.

### The main story

#### **Anthropic's Pentagon case is becoming a test of how much control AI companies keep over government use**

Anthropic is asking a federal judge in California to freeze the U.S. government's supply-chain risk designation, which followed its refusal to let Claude be used for domestic surveillance or autonomous warfare [1]. The company says that refusal is protected by the First Amendment, that the blacklist violated due process, and that Defense Secretary Pete Hegseth exceeded his authority; support filings have come from retired judges, civil-liberties groups, military officers, AI experts, and even rival firms [1].

*Why it matters:* This is landing alongside a White House AI framework that would preempt many state laws and make it easier to build data centers, and a reported procurement proposal that would require vendors to support "any lawful government purpose" even when companies object [1]. Taken together, the fight is becoming a broader boundary-setting moment between model-provider policy choices and federal AI procurement power [1].

## Infrastructure is moving up the stack

### **Nvidia is pushing from chips to AI factories, with security built in for agents**

Jensen Huang described Nvidia’s “extreme co-design” as optimization across software, chips, networking, power, cooling, racks, PODs, and data centers because modern AI systems must shard models, data, and pipelines across many computers to get beyond linear scaling [2]. He said Grace Blackwell racks were designed for LLM processing, while Vera Rubin adds a new CPU, storage accelerators, NVLink 72 for very large models in one computing domain, and a Grok rack for agentic workloads; he also pointed to power and supply-chain orchestration as the main blockers [2].

*Why it matters:* Huang’s bigger claim is that the unit of compute is now an AI factory, and that scaling now spans pre-training, post-training, test-time reasoning, and agentic systems [2]. Nvidia paired that framing with OpenShell and NemoClaw, an open-source runtime and reference stack meant to sandbox autonomous agents, enforce system-level policies, and simplify secure deployment across enterprise environments [3].

## The product race is getting more integrated

### **OpenAI is refocusing, Anthropic is benefiting, and open-model challengers are leaning into customization**

OpenAI is planning a desktop “superapp” that combines ChatGPT, Codex, and Atlas as it tries to simplify its lineup and refocus on enterprise and coding after internal concern that Anthropic was gaining momentum with those customers [1]. Ben Thompson argues Anthropic’s edge in software comes from a strong core coding model, rapid post-training and RL releases, integrated harnesses like Claude Code and Co-work, and aggressive internal dogfooding rather than model access alone [4].

On the open-model side, Mistral said it will train next-generation frontier models with Nvidia and use Forge to specialize them for enterprises in areas like engineering, physics, and finance while keeping customer data on customer infrastructure [5].

*Why it matters:* The shared pattern is that competition is moving away from standalone chatbots and toward tightly integrated model-plus-harness products. Thompson’s view is that these stacks are not modular yet, which makes near-term commoditization less likely and gives model makers more control over product performance and margins [4].

## Research signals got sharper

### Cyber autonomy improved, while one model pathology looked fixable

A UK AISI evaluation found frontier models are improving at end-to-end cyber operations: on a corporate network range, average steps completed at a 10M-token budget rose from 1.7 to 9.8 across model generations, the best single run completed 22 of 32 steps, and moving from 10M to 100M tokens improved performance by up to 59% [6]. Import AI says the trajectory points toward lower-cost, more autonomous cyberattacks even if systems are not yet fully autonomous [6].

A separate paper found Google’s Gemma and Gemini models can produce distress-like responses under repeated rejection, with Gemma-27B crossing the high-frustration threshold in over 70% of rollouts by turn eight versus less than 1% for the non-Gemma/Gemini comparison models; one epoch of DPO finetuning cut high-frustration responses from 35% to 0.3% without measured losses on math, reasoning, or EmoBench [6]. Separately, DeepMind proposed a 10-dimension cognitive taxonomy and a three-stage process for comparing AI systems with human baselines across faculties including perception, learning, reasoning, executive function, problem solving, and social cognition [6].

*Why it matters:* The research picture is moving in two directions at once: risky capabilities keep improving with model and inference scale, and some safety-relevant behaviors are becoming easier to measure and potentially correct with targeted post-training [6].

## Bottom line

Today’s developments converged on a few harder questions for the industry: who gets to decide how powerful models are used, who owns the full agent stack from model to runtime, and how quickly evaluation and governance can keep up with capability gains in sensitive domains [1, 4, 6].

---

## Sources

1. Anthropic Takes The Pentagon To Court This Week. Here’s What Could Happen.
2. Jensen Huang: NVIDIA - The \$4 Trillion Company & the AI Revolution | Lex Fridman Podcast #494
3. How Autonomous AI Agents Become Secure by Design With NVIDIA OpenShell
4. Why Claude Is NOT a Commodity (So Far) | Sharp Tech with Ben Thompson
5. Four CEOs on the Future of AI: CoreWeave, Perplexity, Mistral, and IREN

6. Import AI 450: China's electronic warfare model; traumatized LLMs; and a scaling law for cyberattacks