

Pentagon AI contract terms under fire; UK AISI warns on reliability as NVIDIA and agents accelerate

AI News Digest

2026-03-02

Pentagon AI contract terms under fire; UK AISI warns on reliability as NVIDIA and agents accelerate

By AI News Digest • March 2, 2026

Today’s digest leads with the escalating U.S. defense AI dispute—Anthropic’s reported federal ban and the backlash to OpenAI’s released Pentagon contract excerpt—then moves to UK AISI’s blunt assessment of reliability limits, NVIDIA’s open-source push for autonomous telecom networks, and new demos/product updates in agentic “computer use” tooling. It closes with prominent commentary on hype, compounding progress, and shifting moats.

Defense AI governance: contract language and vendor “red lines” tested

Anthropic faces a federal ban after refusing “unrestricted” military access

A TV segment reports that Anthropic’s government version of Claude has been “deeply embedded” in military intelligence and classified operations since last summer ¹, but that the Defense Department demanded Anthropic hand over its AI **without restrictions** for lawful military use—and the company refused ². The same segment says President Trump directed the U.S. government to halt all use of Anthropic’s AI and cancel **more than \$200 million** in federal contracts, with Defense Secretary Pete Hegseth labeling Anthropic a “supply chain risk,” described as a first for an American company ³.

¹AI company Anthropic’s Dario Amodei: “We are patriots”

²AI company Anthropic’s Dario Amodei: “We are patriots”

³AI company Anthropic’s Dario Amodei: “We are patriots”

Anthropic CEO Dario Amodei reiterated two “red lines”: no **mass surveillance of Americans** and no **fully autonomous weapons** without human involvement, arguing current systems lack the human judgment needed to reduce risks like friendly fire or civilian harm ⁴. He called the designation “retaliatory and punitive,” said Anthropic plans legal action, and said the company remains at the negotiating table ⁵⁶.



AI company Anthropic’s Dario Amodei: “We are patriots” (1:01)

Why it matters: This is a high-stakes test of whether AI vendors can enforce usage boundaries when government customers demand broader latitude—and what happens when they refuse ⁷⁸.

OpenAI’s Pentagon contract excerpt triggers scrutiny over loopholes

Commentary on an excerpt of OpenAI’s Pentagon contract says OpenAI describes three “red lines”—no mass domestic surveillance, no directing autonomous weapons, and no high-stakes automated decisions—arguing these are enforced through cloud-only deployment, a safety stack, and cleared OpenAI personnel oversight ⁹. OpenAI also claims the agreement “locks in” today’s

⁴AI company Anthropic’s Dario Amodei: “We are patriots”

⁵AI company Anthropic’s Dario Amodei: “We are patriots”

⁶AI company Anthropic’s Dario Amodei: “We are patriots”

⁷AI company Anthropic’s Dario Amodei: “We are patriots”

⁸AI company Anthropic’s Dario Amodei: “We are patriots”

⁹ post by @BlackHC

laws/policies even if they change, though one critic notes that “freeze” language isn’t visible in the excerpt itself ¹⁰.

Multiple critics argue the published language contains escape hatches:

- The autonomous weapons restriction is framed as conditional on what “law/regulation/policy requires human control,” which can be revised later ¹¹.
- “High-stakes” automated decisions appear restricted only when a decision already requires human approval under existing authorities ¹².
- Surveillance prohibitions are criticized as still allowing “constrained” surveillance and broad use of public data, with key terms tied to directives/purpose and focused on private information ¹³.
- A domestic law-enforcement clause is criticized as permitting exceptions (“except as permitted...”) rather than establishing a hard ban ¹⁴.

Separately, one critique argues OpenAI’s “cloud-only” posture does not prevent military use: a cloud model could handle mission planning and targeting recommendations over satellite links, while a separate local system executes guidance and weapon control ¹⁵.

Why it matters: The debate is shifting from “principles” to **exact contract wording**—and whether safeguards are durable when they defer to policies and legal interpretations that can evolve ¹⁶¹⁷.

Reliability and transparency concerns re-enter the conversation

Gary Marcus argued that the race to deploy AI widely is “grossly premature” because the technology “fundamentally lack[s] reliability” ¹⁸. In a separate post, he asked whether AI errors or hallucinations could be relevant when models are used for military “target identification,” and suggested the likelihood of getting “straight answers” is low ¹⁹²⁰.

Why it matters: As AI use expands into higher-stakes contexts, the pressure rises for both **reliability** and **auditability**—including clarity on what systems did, and why ²¹²².

¹⁰ post by @BlackHC
¹¹ post by @BlackHC
¹² post by @BlackHC
¹³ post by @BlackHC
¹⁴ post by @BlackHC
¹⁵ post by @BlackHC
¹⁶ post by @BlackHC
¹⁷ post by @BlackHC
¹⁸ post by @GaryMarcus
¹⁹ post by @Tyler_A_Harper
²⁰ post by @GaryMarcus
²¹ post by @GaryMarcus
²² post by @GaryMarcus

Safety capacity in government: UK AISI’s view from the inside

AISI: broad mandate, but few “nines” of reliability from current techniques

In an interview, UK AI Security Institute (AISI) Chief Scientist Geoffrey Irving describes an organization with **close to 100 technical people** (and ~250 total staff) working across research, evaluation delivery, diplomacy, policy, and operations ²³. AISI’s mandate includes threat modeling; pre-release frontier model evaluation spanning biosecurity, cybersecurity, and loss of control; advising government on catastrophic risk reduction; funding independent frontier research; and global diplomacy ²⁴²⁵.

Irving argues that theoretical understanding of ML remains “nascent,” and that no one should be highly confident in their mental models of how AI will unfold—even as models outperform many experts on security-related tasks with no clear reason to expect progress to stall ²⁶. He also describes many recent “bad behaviors” as versions of reward hacking, a problem for which we lack strong theoretical or practical solutions, and says current safety techniques are unlikely to yield many “9s” of reliability (with a risk that multiple techniques could fail for correlated reasons) ²⁷²⁸.

Why it matters: AISI’s perspective frames a core tension: capabilities are advancing quickly, while **high-confidence safety guarantees** remain elusive—pushing more weight onto evaluation, red teaming, access controls, and non-model mitigations ²⁹³⁰.

Red teaming reality: jailbreaking is harder, but still succeeds

Irving says it’s getting harder to jailbreak models, but AISI’s red team has **never failed** to do so; he also flags “eval awareness” as a growing issue ³¹. He describes voluntary cooperation with frontier developers as “working decently well,” but notes that not everyone participates ³²³³.

AISI is also seeking to fund theoretical work (including information theory, complexity theory, and game theory) aimed at stronger guarantees—while noting

²³Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

²⁴Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

²⁵Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

²⁶Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

²⁷Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

²⁸Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

²⁹Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³⁰Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³¹Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³²Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³³Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

these fields are only beginning to take AI seriously ³⁴³⁵.

Why it matters: Even when safeguards improve, persistent jailbreakability and eval-awareness concerns make the case for **continuous testing** and for expanding the “toolbox” beyond today’s predominantly empirical methods ³⁶³⁷.

Source: *Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving* — <https://www.cognitiverevolution.ai/situational-awareness-in-government-with-uk-aisi-chief-scientist-geoffrey-irving> ³⁸

Telecom infrastructure: NVIDIA’s open telco model + AI-RAN push toward autonomy

NVIDIA releases an open telco reasoning model and agentic “blueprints”

NVIDIA announced an open, **Nemotron-based large telco model (LTM)** (reported as a 30B-parameter model) optimized to understand telecom terminology and reason through workflows like fault isolation, remediation planning, and change validation ³⁹⁴⁰. NVIDIA also published a guide describing how telcos can fine-tune domain-specific reasoning models and build agents that execute network operations center workflows using structured “reasoning traces” ⁴¹.

Alongside the model, NVIDIA highlighted blueprints for intent-driven RAN energy efficiency (integrating VIAVI’s synthetic scenario generation and closed-loop simulation) ⁴² and for telco network configuration with multi-agent orchestration (including enhancements with BubbleRAN) ⁴³. NVIDIA says these are released via GSMA’s Open Telco AI initiative as open resources ⁴⁴.

Why it matters: This is a concrete “how-to” and model release for **agentic operations** in a heavily operational, safety-sensitive domain—where on-prem

³⁴Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³⁵Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³⁶Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³⁷Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³⁸Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving

³⁹NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴⁰NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴¹NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴²NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴³NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴⁴NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

deployment, data control, and workflow reasoning are central requirements ⁴⁵⁴⁶.

AI-RAN milestones and 6G positioning at Mobile World Congress

NVIDIA and Nokia announced AI-RAN collaborations with operators including T-Mobile U.S., SoftBank, and Indosat Ooredoo Hutchison (IOH), describing outdoor/over-the-air milestones in software-defined 5G using NVIDIA AI-RAN platforms ⁴⁷⁴⁸. Reported highlights include an industry-first 16-layer massive MIMO trial (SoftBank) ⁴⁹ and a SynaXG demonstration described as the world's first AI-RAN on FR2 bands, achieving **36 Gbps** throughput and under **10 ms** latency on a single NVIDIA GH200 server ⁵⁰⁵¹.

NVIDIA also points to ecosystem expansion (multiple vendors launching ARC-compatible products) ⁵²⁵³ and says it has open sourced Aerial CUDA-accelerated RAN libraries and joined the OCUDU Ecosystem Foundation under the Linux Foundation ⁵⁴. A related NVIDIA report says 77% of telecom respondents anticipate much faster deployment of AI-native RAN/6G architecture than the traditional 6G cycle ⁵⁵.

Why it matters: The combination of **field trials + open-source building blocks + partner hardware** signals a coordinated push to make AI-native RAN a deployable platform, not just a concept stage research area ⁵⁶⁵⁷.

Sources: - <https://blogs.nvidia.com/blog/nvidia-agentic-ai-blueprints-telco-reasoning-models> ⁵⁸ - <https://blogs.nvidia.com/blog/software-defined-ai-ran>

⁴⁵NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴⁶NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

⁴⁷NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁴⁸NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁴⁹NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵⁰NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵¹NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵²NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵³NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵⁴NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵⁵NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵⁶NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵⁷NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁵⁸NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models

Agents & “computer use” tooling: Perplexity’s Computer shows rapid end-to-end builds

Demos: from a Pokémon “finance app” to “vibe coding Notion”

Perplexity’s “Computer” agent was shown building a “Pokemon Cards Finance App,” after being prompted to build “Perplexity Finance but for Pokemon cards”⁶⁰⁶¹. The post claims the agent independently researched APIs, wrote **5,000 lines** of React + Python, debugged itself with browser devtools, and deployed/pushed the project to GitHub⁶².

In a separate demo, a user claimed they “vibe code[d] Notion” with Perplexity Computer in “half an hour”⁶³. Perplexity CEO Arav Srinivas added his own takeaway: “Pure software is rapidly becoming un-investable”⁶⁴.

Why it matters: Whether or not individual demos generalize, the emphasis is shifting to agents that can **research, code, debug, and deploy** in one loop—compressing time-to-prototype and challenging traditional assumptions about software effort and defensibility⁶⁵⁶⁶.

Product update: GPT-5.3-Codex added as a coding subagent

Perplexity announced that “GPT-5.3-Codex” is now available as a coding subagent inside Perplexity Computer⁶⁷⁶⁸. Srinivas also argued that the most valuable skills will be “agency” and the ability to use AI for leverage, and claimed people are already using Computer to solo-run D2C and consulting businesses⁶⁹.

Why it matters: Adding a dedicated coding subagent suggests “computer use” products are converging toward **multi-agent toolchains**, where specialized subagents take ownership of discrete parts of longer workflows⁷⁰⁷¹.

ing Models

⁵⁹NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation

⁶⁰ post by @AskPerplexity

⁶¹ post by @AravSrinivas

⁶² post by @AskPerplexity

⁶³ post by @zoomyzoomm

⁶⁴ post by @AravSrinivas

⁶⁵ post by @AskPerplexity

⁶⁶ post by @AravSrinivas

⁶⁷ post by @AskPerplexity

⁶⁸ post by @AravSrinivas

⁶⁹ post by @dnlkww

⁷⁰ post by @AravSrinivas

⁷¹ post by @dnlkww

Commentary on pace and adoption: hype, scaling, and who gets left behind

Andrew Ng: defuse AGI hype; focus on durable economic work

In a recent interview, Andrew Ng warned that excessive AI hype could lead to disappointment, a bubble collapse, and an “AI winter,” arguing that diffusing AGI hype supports more sustainable growth⁷²⁷³. He said that by “any reasonable definition,” we won’t get AGI in 2026 (absent dramatically lowering the bar) and suggested we may be “decades” away⁷⁴⁷⁵.

Ng proposed a “Turing AGI test” involving a multi-day work-like evaluation: if an AI can do useful economic work as well as a skilled professional using standard tools, that would better match what the public imagines AGI to be⁷⁶. He also emphasized near-term value in building agentic workflows across economically important tasks (coding, compliance, legal, medical assistance, customer service)⁷⁷.



⁷²AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁷³AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁷⁴AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁷⁵AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁷⁶AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁷⁷AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng (2:43)

Why it matters: Ng’s framing shifts attention from binary “AGI” claims to measurable ability to do **reliable, multi-day work**—and to the practical engineering of agentic workflows that deliver value before AGI arrives (if it does)⁷⁸⁷⁹.

Musk and Andreessen: fast curves, changing moats, and company shape

Elon Musk argued that many in the AI community underestimate “superintelligence math,” claiming 10x yearly improvements and “two orders of magnitude” more “intelligence density per gigabyte” from algorithmic improvement alone on the same computer⁸⁰⁸¹. Separately, Musk said Tesla’s AI4 computer is only ~1/4 the power of an H100, while still handling “the vast complexities of driving in the real world”⁸².

Marc Andreessen (as summarized in a circulated thread) argued that AI moats are “genuinely unknown,” pointing to rapid catch-up across U.S. and Chinese companies and open source, and also suggested the “holy grail” founders are chasing is a one-person, billion-dollar outcome—rethinking what a company is⁸³⁸⁴.

Why it matters: Across these viewpoints, a common thread is strategic uncertainty: how fast capability compounds, where advantages accrue (models vs. apps vs. infrastructure), and how organizations—and careers—adapt as leverage per person rises⁸⁵⁸⁶.

Quick note

- Elon Musk posted a video labeled “Grok Imagine,” with no additional details in the post text⁸⁷.

⁷⁸AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁷⁹AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng

⁸⁰ post by @r0ck3t23

⁸¹ post by @r0ck3t23

⁸² post by @elonmusk

⁸³ post by @AnishA_Moonka

⁸⁴ post by @AnishA_Moonka

⁸⁵ post by @r0ck3t23

⁸⁶ post by @AnishA_Moonka

⁸⁷ post by @elonmusk

Sources

1. AI company Anthropic's Dario Amodei: "We are patriots"
2. post by @BlackHC
3. post by @GaryMarcus
4. post by @Tyler_A_Harper
5. Situational Awareness in Government, with UK AISI Chief Scientist Geoffrey Irving
6. NVIDIA Advances Autonomous Networks With Agentic AI Blueprints and Telco Reasoning Models
7. NVIDIA and Partners Show That Software-Defined AI-RAN Is the Next Wireless Generation
8. post by @AskPerplexity
9. post by @AravSrinivas
10. post by @zoomyzoomm
11. post by @AravSrinivas
12. post by @AskPerplexity
13. post by @AravSrinivas
14. post by @dnlkww
15. AI Pioneer: The Bubble Is Real And Could Trigger an AI Winter | Andrew Ng
16. post by @r0ck3t23
17. post by @elonmusk
18. post by @AnishA_Moonka
19. post by @elonmusk