

The Pentagon's new AI task force is an inflection point: it could harden America's most sensitive networks against rapidly evolving cyber threats—or quietly wire fragile, corporate-built AI into the heart of U.S. intelligence and warfighting with far too little public oversight.

According to an internal email reported this week, a new Pentagon task force spanning U.S. Cyber Command and the National Security Agency will determine how powerful commercial AI models—such as those from OpenAI and Google—can be safely deployed across their missions and the Defense Department's most sensitive networks.

This comes on the heels of formal Pentagon deals with a roster of AI heavyweights—OpenAI, Google, Nvidia, Microsoft, Amazon Web Services, SpaceX/xAI, and others—to deploy their models on classified networks up to “Impact Level 7,” which covers the most sensitive national-security information. Over 1.3 million Defense Department personnel are already using the department's GenAI.mil platform, a sign that generative AI is moving from experiment to infrastructure inside the U.S. military.

Why the Pentagon is moving so fast

From the Pentagon's perspective, the urgency is understandable. Cyber Command has an AI roadmap with roughly 100 tasks aimed at embedding AI across logistics, security and defense operations, and it has already stood up an AI task force to move from ad-hoc pilots to systematic adoption. Officials report that generative AI tools are already cutting the time needed to analyze network traffic and identify malicious activity on Defense Department networks.

Former Cyber Command leaders openly describe AI as “essential” for detecting threats, prioritizing vulnerabilities, accelerating decisions, and conducting both defensive and offensive cyber operations faster than adversaries. NSA, for its part, has launched an Artificial

Intelligence Security Center to guide secure AI development across national security systems and the defense industrial base, explicitly recognizing that AI will shape future intelligence and cyber conflict.

In other words, this is not a speculative gambit. The U.S. security establishment believes that without deep AI integration, it will lose the speed and scale contest to sophisticated rivals like China and Russia in cyberspace and intelligence analysis.

The security paradox at the heart of the plan

Yet wiring large language models into the most sensitive networks creates a profound security paradox. These systems are powerful precisely because they absorb vast amounts of data, learn patterns, and can be hooked into tools and databases; those same traits make them inherently leaky, brittle, and hard to fully control.

The NSA itself warns that AI systems must be protected across their entire lifecycle—from training data to models to deployment pipelines—because they are vulnerable to theft, manipulation, and abuse. On classified networks, the stakes are higher: a prompt-logging failure, misconfigured connector, or subtle model vulnerability could expose sources and methods or provide new attack surfaces to hostile intelligence services.

At the same time, using models developed primarily for consumer or enterprise markets inside intelligence environments risks importing their biases and failure modes—hallucinations, overconfidence, susceptibility to cleverly crafted prompts—into workflows where errors might escalate crises or misdirect cyber operations. The very speed and scale that make AI attractive also make its failures far more consequential.

Private power, public secrecy

Equally troubling is the structure of these deals. Reporting indicates that at least some contracts, such as Google's, allow Pentagon use of AI models for "any lawful governmental purpose," with little public detail on guardrails. Hundreds of Google employees have raised concerns about the firm's growing classified AI work, reviving memories of the suppressed employee revolt over the Project Maven drone-vision contract.

Anthropic, a leading AI lab, is conspicuously outside the main Pentagon arrangement after a dispute over the department's demand for far more expansive access to its models, and it has been designated a "supply chain risk." That standoff underscores just how much leverage the U.S. government now seeks over commercial AI systems—and how little of that negotiation is visible to the public or Congress.

By concentrating so much critical infrastructure in a handful of AI vendors, the Pentagon is also deepening long-term dependencies on private platforms whose incentives are not aligned with democratic accountability. These are proprietary, rapidly changing systems; if something goes wrong on a classified network, the public will learn about it, if at all, only years later and in heavily redacted form.

What "safe deployment" must actually mean

If this new task force is to be more than a rubber stamp, "safe deployment" must be defined far more rigorously than it has been in today's cloud-AI contracts. At minimum, that should include:

- Strict data minimization and segregation, with air-gapped or tightly scoped models trained for specific mission profiles rather than simply porting general-purpose chatbots

into intelligence environments.

- Mandatory, continuous red-teaming focused on espionage and cyber-warfare scenarios, not just generic tests of toxicity or bias.
- Clear, enforceable bans on using these systems for domestic mass surveillance or fully autonomous targeting, commitments that companies like Google have nodded to in principle but which are not yet backed by binding law or transparent oversight.

Congress should require independent audits of AI use on classified networks by cleared technical experts outside the chain of command, as well as regular reporting—at least in closed hearings—on incidents, near misses, and model failures. Otherwise, the only serious constraints will be whatever terms giant AI vendors are willing to accept in secret contracts.

The real test of American leadership

There is a narrow path where this effort genuinely strengthens U.S. and allied security: AI-accelerated detection of intrusions, faster patching of vulnerabilities, smarter triage for overburdened analysts, and more resilient networks that can adapt in real time to sophisticated attacks. On that path, the new task force could become a model for how democracies integrate frontier technologies without surrendering control to them.

But there is also a darker path in which opaque partnerships with a few tech giants quietly place experimental systems at the center of nuclear-adjacent decision chains and cyber operations, with minimal democratic input until a catastrophic failure forces a reckoning. The choice between those futures will be made not only in secure facilities at Fort Meade, but in hearing rooms on Capitol Hill and in public debates that have barely begun.

If the Pentagon insists on racing AI into its most sensitive networks, the rest of us must be equally fast in demanding the rules, oversight, and transparency that such a step requires.