

# Post-Mythos AI Model Regulation: Licensing and Disclosure Frameworks

How Governments Are Responding to Offensive AI Capabilities  
After Claude Mythos Preview

2026-05-25

 AI-assisted Rapid Research



**© 2026 Cloud Security Alliance. Some rights reserved.**

You may download, store, display, view, print, redistribute, and link to this document in its original, unmodified form, provided that attribution to the Cloud Security Alliance is maintained and all trademark and copyright notices remain intact.

This document may not be modified or altered. You may quote portions of the document as permitted by the Fair Use provisions of the United States Copyright Act, provided that attribution is given to the Cloud Security Alliance.

This document may be shared on professional and social media platforms in its original form with attribution.

*This document was generated with AI assistance and has not undergone official CSA review and approval processes.*

---

## Key Takeaways

- Anthropic's April 2026 announcement of Claude Mythos Preview – a model withheld from broad release due to its autonomous offensive cyber capabilities – has triggered a significant and rapid acceleration in AI governance activity, prompting regulatory responses across multiple jurisdictions within weeks of its announcement and forcing governments to confront a capability threshold for which existing regulatory frameworks had no established mechanism.
- The United States government exercised informal veto power over Anthropic's commercial rollout through executive pressure that has no statutory basis, no formal process, and no appeals mechanism – what one analysis describes as "an informal, highly improvised licensing regime" with serious implications for due process and industry certainty [1].
- A White House executive order establishing an FDA-style pre-release vetting process for frontier AI models was actively under development, confirmed publicly by the National Economic Council Director on May 6, 2026 [2], and then cancelled on May 22, 2026, when the President publicly cited concerns that certain provisions would impede AI competitiveness [9].
- Congressional activity has produced two relevant bills – the AI Foundation Model Transparency Act (H.R. 8094) and the Advanced AI Security Readiness Act (H.R. 3919) – but neither establishes the licensing authority or mandatory offensive capability-disclosure mechanism that legal and security analysts have argued is necessary for Mythos-class risks [1] [7].
- California's SB 53, effective January 1, 2026, is the most operationally concrete frontier AI disclosure mandate currently in force in the United States, requiring safety incident reporting within fifteen days for covered models above the  $10^{26}$  FLOP training threshold, though it does not address offensive capabilities specifically [3].
- Security and compliance professionals should treat the regulatory environment as fragmented but rapidly evolving: taken together, these developments suggest a direction across multiple jurisdictions toward pre-release capability assessment, mandatory disclosure of high-risk behaviors, and some form of government access for models above capability thresholds.

# Background

On April 7, 2026, Anthropic publicly announced that it had developed Claude Mythos Preview, a model it considered too dangerous for unrestricted release due to autonomous offensive cybersecurity capabilities that had emerged as byproducts of general reasoning and code-generation improvements [4]. The announcement had few if any direct precedents: a frontier AI laboratory publicly acknowledged building a model it chose not to release and described in technical terms why. Anthropic stated that in controlled benchmarking, a prior frontier model produced working software exploits approximately twice across several hundred attempts, while Mythos produced 181 working exploits under comparable conditions [5]. In broader testing across thousands of open-source software targets, Mythos achieved full system compromise on ten fully-patched systems without human guidance [5]. Access was confined to a small cohort of vetted organizations under Project Glasswing, a controlled-access program for defensive vulnerability research.

The capability disclosure prompted rapid international responses. The Bank of England Governor warned that Mythos could "crack the whole cyber risk world open" [6]. The UK AI Security Institute published independent cyber capability evaluation results within six days of Anthropic's announcement, finding that Mythos was substantially more capable at cyber offence than any model previously assessed and the first AI to complete a 32-step simulated cyber-attack range [7]. AISI also concluded that frontier model cyber capabilities were now doubling approximately every four months – roughly twice the pace measured in prior assessments [7]. What the Mythos case made visible was the pre-existing absence of any formal mechanism by which governments could require advance notice of, evaluate, or condition the release of AI models crossing a security-relevant capability threshold.

## Security Analysis

### The Governance Vacuum Exposed

The predominant approach to AI governance prior to Mythos focused on risk-based classification of deployed AI applications rather than foundation model capabilities themselves – a framing reflected in major frameworks across multiple jurisdictions. The EU AI Act classifies systems by use case. U.S. federal frameworks address organizational development and deployment practices. California SB 53 targets developers of models above a defined compute threshold. None of these frameworks, as written, establish the authority to require advance notice of, evaluate, or condition the release of AI models

crossing a security-relevant capability threshold. The governance response, as of May 2026, has been largely informal, reactive, and jurisdictionally fragmented – a reflection of the structural gap between the pace of AI capability development and the pace of legislative action [1].

## **The U.S. Response: Informal Veto and Failed Formal Action**

The most consequential early U.S. government action was both informal and without legal basis. White House officials communicated to Anthropic, through what reporting describes as a phone call or series of meetings, that the administration opposed broader access to Mythos on national security grounds [8]. Anthropic maintained Project Glasswing's access restrictions throughout. Legal analysts noted immediately that this constituted a form of licensing without a license: the government exercised effective control over the distribution of a private product through executive pressure with no authorizing statute, no defined criteria, and no mechanism for appeal or review [1]. For companies, this creates regulatory uncertainty of the worst kind – compliance consists of inferring what officials want, with no formal process to provide certainty or recourse.

The formal alternative was seriously considered. National Economic Council Director Kevin Hassett confirmed on May 6, 2026 that the White House was actively studying a directive to create government evaluation of frontier AI models before commercial release, drawing an explicit analogy to the FDA's pre-market approval process for pharmaceuticals [2]. The proposed mechanics included required safety documentation, structured red-team testing for offensive cyber capabilities, evaluation of dual-use risks above defined thresholds, and mandatory disclosure to a federal evaluator [2]. On approximately May 22, 2026, President Trump cancelled the planned signing ceremony, publicly citing concerns that certain provisions would structurally block American AI competitiveness [9]. The cancellation means the United States currently has no formal pre-release oversight mechanism, despite the Mythos case having demonstrated that frontier model capabilities can reach security-relevant thresholds without any established government review process.

## **Congressional Efforts: Transparency Without Teeth**

Congressional activity on AI model governance reflects a transparency-first approach that stops short of licensing. The AI Foundation Model Transparency Act of 2026 (H.R. 8094), introduced bipartisanly on March 26, 2026, would direct the FTC – in consultation with NIST, the Department of Commerce, and OSTP – to establish disclosure standards for high-impact foundation models [10]. Required disclosures would include training data summaries, model design parameters, known limitations and risks, and evaluation practices – making it the first federal disclosure standard for high-impact foundation models

if enacted. The bill stops short of requiring offensive capability assessment or pre-release evaluation authority, however, and does not authorize government conditioning of model deployment. It remains in committee.

The Advanced AI Security Readiness Act (H.R. 3919) directs the NSA Director to develop an AI Security Playbook addressing strategies to protect AI systems against theft and adversarial compromise [11]. The bill explicitly defines "covered AI technologies" to include systems matching or exceeding human expert performance in cyber offense, model autonomy, and self-improvement – a definition that captures Mythos-class capabilities – but explicitly states that its provisions do not authorize regulatory or enforcement actions against AI companies. It is a defensive planning document, not a licensing regime. As of May 2026, the United States has over 1,200 AI-related bills across federal and state legislatures [12], with significant activity but no comprehensive framework addressing pre-release offensive capability oversight.

## State and International Responses

California's SB 53 provides the most concrete existing disclosure mandate in the United States [3]. For models trained above  $10^{26}$  floating-point operations, it requires safety incident reporting to California authorities within fifteen days, transparency reports covering capabilities and safety testing, and whistleblower protections for AI safety researchers. The fifteen-day reporting requirement is particularly relevant as a model for federal legislation, as it would in principle apply to a model exhibiting unexpected offensive capabilities during internal evaluation. However, the Trump administration's National Policy Framework for Artificial Intelligence, released March 20, 2026, calls for federal preemption of state AI laws and has directed the Attorney General to challenge state disclosure mandates deemed inconsistent with a unified national approach [13] – a direct threat to SB 53's durability.

Outside the United States, the UK demonstrated that rapid, technically grounded capability assessment is operationally achievable: AISI's independent Mythos evaluation published within six days demonstrates one approach by which an AI safety institute can translate model capabilities into actionable public intelligence for governments worldwide [7]. The UK also issued an open letter from ministers to business leaders on AI-assisted cyber threats and confirmed that the Cyber Security and Resilience Bill will return to Parliament [14]. The EU AI Act's August 2, 2026 core application date brings the systemic risk provisions into force for the largest frontier model providers, requiring adversarial testing protocols and detailed risk assessment documentation – representing the most concrete pre-release capability assessment obligations legally binding on frontier AI developers among the jurisdictions reviewed in this note [15].

## The Disclosure Architecture Gap

No current framework has fully resolved what "disclosing offensive AI capabilities" to a government should mean in practice. A Coordinated Disclosure of Dual-Use Capabilities (CDDC) framework – drawing on the analogy of coordinated vulnerability disclosure in security research – has been proposed as one model [16]. Under this approach, AI developers would notify a designated coordinator when internal evaluations surfaced capabilities above defined risk thresholds; the coordinator would share information with relevant government defenders while managing public disclosure timing. The framework's logic maps well onto how Project Glasswing itself operates, with Anthropic functioning as finder and coordinator. What it does not resolve is the question of what government authority attaches to the notification: transparency requirements alone, without enforcement authority, may leave the underlying incentive structure largely unchanged – a limitation observed in analogous disclosure-based regulatory regimes.

## Recommendations

### Immediate Actions

Security and compliance teams should map their AI model vendors to the disclosure requirements already in force. California SB 53's fifteen-day safety incident reporting obligation has applied since January 1, 2026; organizations procuring models from developers subject to SB 53 should confirm that their vendor contracts reflect those obligations. Organizations within EU jurisdiction should determine which vendor products will be classified as systemic risk models under the AI Act's August 2026 application date, and what adversarial testing documentation those vendors are required to produce.

### Short-Term Mitigations

Organizations should establish processes for tracking the evolving regulatory landscape, because the current legislative environment could produce mandatory requirements on relatively short timelines, and the direction of travel across multiple jurisdictions is toward increased obligations. Supplier questionnaires, procurement contracts, and third-party risk assessments should be updated to include questions about frontier model capability assessments, dual-use risk evaluations, and government notification practices. Organizations that deploy AI tools with offensive security capabilities – even defensively – should document the business justification, access controls, and authorized use scope for each, both to demonstrate future compliance and to manage insider threat and misuse exposure.

## Strategic Considerations

Organizations with direct policy stakes – AI developers, critical infrastructure operators, large enterprise AI consumers – should engage actively with regulatory processes underway rather than waiting for final rules. The informal nature of current U.S. government authority means that industry norms and voluntary frameworks carry unusual weight in shaping eventual formal requirements. Participating in emerging capability disclosure coordination mechanisms and providing substantive technical input to legislative staff will influence whether the eventual framework is workable. Organizations should also treat the current informal U.S. licensing regime as a precursor to a formal one: the precedent established by the White House's informal veto of Mythos's commercial rollout indicates executive branch intent to assert oversight authority, and when statutory authority is granted, organizations that have built capability assessment and government notification processes into their workflows will be best positioned to comply.

## CSA Resource Alignment

The governance challenges in this note map directly onto CSA's AI security and governance portfolio. The CSA AI Controls Matrix (AICM) v1.0 addresses capability transparency through its Model Provider domain, including controls covering model card publication, capability documentation, and safety testing requirements. Organizations seeking to operationalize compliance with current and forthcoming disclosure mandates can use the AICM Model Provider controls as a baseline for vendor assessment. The AICM's treatment of shared security responsibility between model providers, application providers, and AI customers is directly applicable to understanding which party bears disclosure obligations under frameworks like California SB 53.

CSA's MAESTRO threat modeling framework for agentic AI provides the analytical vocabulary needed to classify offensive AI capabilities at the specificity regulatory compliance will require. MAESTRO's enumeration of autonomous reasoning, multi-step execution, and tool use as distinct threat surfaces aligns with the emerging regulatory distinction between general-purpose foundation models and models capable of autonomous offensive operations. The CSA STAR program's continuous monitoring and certification mechanisms offer a pathway for AI model vendors to demonstrate ongoing compliance with safety and transparency obligations in a form enterprise customers can incorporate into third-party risk programs. As governments formalize capability disclosure requirements, STAR attestations covering AI model evaluation practices may serve as differentiators in enterprise procurement, providing verifiable evidence of compliance with emerging regulatory baseline standards.

## References

- [1] Lawfare. "[Mythos Fallout: U.S. Government Weighs AI Model Regulation](#)." Lawfare Media, 2026.
- [2] The Hill. "[Hassett: White House May Review AI Models 'Like an FDA Drug'](#)." The Hill, May 6, 2026.
- [3] LegiScan. "[California SB 53: Chaptered Text](#)." LegiScan, 2026.
- [4] Scientific American. "[What Is Mythos and Why Are Experts Worried About Anthropic's AI Model](#)." Scientific American, April 2026.
- [5] Data Protection Report. "[When AI Becomes the Cyber Attacker: Mythos and What Comes Next](#)." Data Protection Report, May 2026.
- [6] American Action Forum. "[The Mythos Breakthrough: AI, Cyber Risk, and the Governance Gap](#)." American Action Forum, 2026.
- [7] UK AI Security Institute. "[Our Evaluation of Claude Mythos Preview's Cyber Capabilities](#)." AISI, April 2026.
- [8] MindStudio. "[The US Government Just Restricted an AI Model Rollout for the First Time](#)." MindStudio, 2026.
- [9] Digitimes. "[Trump Calls Off Executive Order for AI Safety Due to Concerns About Inhibiting Growth](#)." Digitimes, May 22, 2026.
- [10] Congress.gov. "[H.R. 8094: AI Foundation Model Transparency Act of 2026](#)." 119th Congress, introduced March 26, 2026.
- [11] Congress.gov. "[H.R. 3919: Advanced AI Security Readiness Act](#)." 119th Congress, 2025–2026.
- [12] Fortune. "[The U.S. Has 1,200 AI Bills and No Good Test for Any of Them](#)." Fortune, May 15, 2026.
- [13] Morrison Foerster. "[Trump Administration Releases National AI Policy Framework](#)." Morrison Foerster, April 2026.
- [14] UK Government. "[AI Cyber Threats: Open Letter to Business Leaders](#)." GOV.UK, April 2026.
- [15] Legal Nodes. "[EU AI Act 2026 Updates: Compliance Requirements and Business Risks](#)." Legal Nodes, 2026.

[16] O'Brien, J., Ee, S., Krprayoon, J., Anderson-Samways, B., Delaney, O., and Williams, Z. "[Coordinated Disclosure of Dual-Use Capabilities: An Early Warning System for Advanced AI.](#)" arXiv:2407.01420, July 2024.