



BOSTON GLOBAL FORUM · AI WORLD SOCIETY

# Mythos and the Path to AIWS Trust Infrastructure in the Age of Frontier Cyber AI

---

*A Case Study for America at 250: A Beacon for the AI Age*

**Loeb House · Harvard University · May 1, 2026**

**Core Message.** The release of Claude Mythos Preview under Project Glasswing, together with reported strategic investments by Google and Amazon in Anthropic, marks a new phase of artificial intelligence: frontier AI is becoming cyber-capable, infrastructure-scale, and geopolitically consequential. The lesson is not panic or unrestricted acceleration. The lesson is that trust must now be built as operational infrastructure — through standards, controlled access, vendor assurance, monitoring, auditability, incident exchange, and democratic limits.

## Introduction

The release of Claude Mythos Preview under Project Glasswing on April 7, 2026 is one of the most consequential events in the recent evolution of artificial intelligence. Anthropic introduced Mythos not as a conventional product, but as a frontier capability whose strength in cybersecurity required exceptional governance from the moment of release.

Anthropic chose not to make Mythos generally available. It offered the model through a gated research preview to a coalition of partners — including Amazon Web Services, Apple, Broadcom, Cisco, CrowdStrike, Google, JPMorganChase, the Linux Foundation, Microsoft, NVIDIA, and Palo Alto Networks — together with more than forty organizations responsible for critical software infrastructure. Anthropic also committed \$100 million in model usage credits to support the initiative and donated \$4 million to open-source security organizations including the Apache Software Foundation, Alpha-Omega, and OpenSSF.

This was a deliberate act of restraint and responsibility. In a moment when the industry could have rushed a powerful new model into broad deployment, Anthropic instead invited defenders to use Mythos first: to identify and patch vulnerabilities in systems on which billions of people depend before adversaries could exploit the same techniques.

The Mythos moment therefore matters for two reasons that belong together. It marks the arrival of frontier AI capabilities that can materially reshape cybersecurity. It also demonstrates a model of responsible release that the AI Age urgently needs to learn from. The purpose of this report is not to endorse any single company, but to recognize a new class of frontier AI risk and opportunity that requires shared trust infrastructure beyond the capacity or authority of any one firm.

The Boston Global Forum and the AI World Society offer this report in that spirit: as an invitation to build AIWS Trust Infrastructure for the AI Age — a public-good architecture of standards, controls, monitoring, accountability, and trusted order that can make frontier AI worthy of service to democracy, peace, security, and humanity.

## I. Mythos as a Threshold Moment — and a Model of Responsible Release

Anthropic has been clear about why Mythos required special handling. According to its technical evaluation, Mythos Preview identified thousands of zero-day vulnerabilities across major operating systems and web browsers, including vulnerabilities that had remained undiscovered for years — among them a 27-year-old bug in OpenBSD and a 17-year-old remote code execution vulnerability in FreeBSD (CVE-2026-4747). In some evaluations, Mythos autonomously identified and exploited vulnerabilities with little or no human intervention, and in one case wrote a multi-vulnerability browser exploit chain that escaped both renderer and operating-system sandboxes.

These are capabilities that could reshape cybersecurity. Anthropic's response was to gate access, brief senior U.S. government officials before release, and design Project Glasswing as a coordinated industry effort to give defenders a head start. In the words of partners across the coalition — from Cisco's chief security officer to AWS security leadership — the urgency of this moment required cross-industry collaboration of a kind rarely seen.

This pattern is itself the point. The defining feature of the Mythos moment is not the model alone. It is the architecture of restraint, transparency, and partnership that Anthropic built around the model. Anthropic chose to:

1. **Restrict access** to a curated coalition of capable defenders rather than release broadly.
2. **Disclose openly** through detailed system cards and red-team writeups, including capabilities the company itself described as concerning.
3. **Brief governments in advance**, including senior officials across the U.S. national security community.
4. **Subsidize the public good** through \$100 million in usage credits and \$4 million in open-source security donations.
5. **Couple capability with safeguards** by announcing that subsequent Claude releases would carry automatic detection and blocking of prohibited cybersecurity requests.

This is what responsible release in the frontier era looks like in practice. It is the kind of conduct the AI World Society has called for since its founding in 2017. America at 250: A Beacon for the AI Age should recognize and affirm it — while also asking how such practice can become a shared standard rather than a voluntary exception.

## II. Google's \$40 Billion Bet on Anthropic — A Signal of a New AI Infrastructure Era

The reported plan by Google-parent Alphabet to invest up to \$40 billion in Anthropic gives the Mythos moment a larger strategic meaning. According to Reuters, the plan includes an initial \$10 billion investment, with an additional \$30 billion potentially tied to performance milestones. This follows Amazon's expanded commitment to Anthropic and Anthropic's own massive cloud-infrastructure commitments.

This is more than a financial transaction. It signals a new AI era in which frontier AI companies are becoming strategic infrastructure partners for Big Tech, cloud providers, governments, critical infrastructure operators, and global enterprises. Anthropic is no longer only a model developer. It is becoming part of the core AI operating layer for business, coding, agents, cybersecurity, and trusted digital systems.

Together with Amazon's expanded investment and Anthropic's long-term cloud commitments, Google's reported investment shows that the AI race is shifting from model competition to infrastructure competition: compute, cloud, chips, data centers, energy, agent platforms, cybersecurity capability, and trust systems.

For BGF and AIWS, this development reinforces the central argument of this report: as frontier AI companies become civilization-scale infrastructure actors, society needs a trust architecture that is commensurate with their scale. AIWS Trust Infrastructure, AIWS Information Trust Infrastructure, AIWS Trust Rating, and AIWS Trusted Order are not peripheral governance concepts. They are the institutional and operational foundations needed for the next stage of AI.

## III. The Industry-Wide Frontier Mythos Reveals

Even as Anthropic acted with care, the days following the Mythos announcement made visible a deeper truth: no single company, however principled, can carry the full weight of governance for capabilities of this magnitude. Bloomberg reported on April 21, 2026 that a small group of users had gained access to Mythos through a third-party vendor environment, and Anthropic confirmed it was investigating. The company stated that there was no evidence Anthropic's core systems had been impacted, nor that the unauthorized activity extended beyond the third-party vendor environment. The incident was operational — concerning the surrounding ecosystem, not the model itself — and Anthropic responded with transparency.

Yet the lesson is significant for the entire industry. When dozens of companies and tens of thousands of authorized employees, contractors, researchers, and vendors can touch a frontier model, the surface area of the deployment ecosystem becomes the new frontier of safety. This is not a failure of any single firm. It is the structural reality of frontier AI.

Trust must be engineered not only into the model, but into every contractor relationship, vendor environment, credentialing system, monitoring loop, and incident-response protocol that surrounds it. This is the gap AIWS Trust Infrastructure exists to address. It is a complement to what frontier developers themselves are building, not a substitute.

Anthropic's own next step — embedding automatic detection of prohibited cybersecurity requests in subsequent Claude models — points in the same direction. Post-deployment governance is now part of the safety case. The question is whether the broader ecosystem will rise to that standard, or whether each new frontier release will encounter the same gaps separately, at higher cost and greater public risk.

## **IV. Mythos, National Security, and the Democratic Limits That Define American AI Leadership**

A second dimension of the Mythos period concerns the relationship between frontier AI companies and the state. Anthropic is currently engaged in a legal proceeding with the Department of Defense, having filed suit on March 9, 2026 after refusing to remove safeguards in its products against use in fully autonomous weapons and domestic surveillance of Americans. Thirty-seven researchers and engineers from OpenAI and Google — including Google’s chief scientist Jeff Dean — filed an amicus brief in support of Anthropic, arguing that the case implicates the broader ability of AI experts to debate openly the risks and limits of their technologies.

The Boston Global Forum and the AI World Society regard Anthropic’s position with respect. The principle Anthropic is defending — that frontier AI systems should not be made available for fully autonomous weapons or for surveillance against citizens — is a principle this network has championed since its founding. It is the principle Governor Michael S. Dukakis articulated when he said that the dignity of democratic citizenship must not be subordinated to the convenience of any technology. It is the principle that informs the AIWS Social Contract for the AI Age, the AIWS Trust Standards, and the very name Boston Global Forum: the idea that ideas matter, that limits matter, and that a free society draws strength from the lines it draws around its most powerful instruments.

A democratic AI order is not one in which capability is unconstrained. It is one in which capability is paired with credible restraint. American AI leadership in the AI Age will not be measured by raw model power alone. It will be measured by whether American companies and American institutions can together build a system in which extraordinary capability serves democracy, peace, security, and humanity — and refuses uses that would corrode those ends.

Anthropic’s stand is therefore part of a larger conversation that America at 250 must hold openly: how the United States, at its semiquincentennial, defines the constitutional, ethical, and operational boundaries within which frontier AI will be developed and deployed. This conversation is one in which BGF and AIWS hope to walk alongside Anthropic and other principled frontier developers, not at a distance from them.

## **V. AIWS Trust Infrastructure — A Complement to What Frontier Developers Are Building**

The five elements that follow describe AIWS Trust Infrastructure as it should now stand for cyber-capable frontier AI. They are framed not as fixes for what any single developer has failed to do, but as the public-good architecture that no single developer can be expected to provide alone. Each is offered as an extension of, and complement to, the responsible practices already pioneered by leaders such as Anthropic.

### **1. Trust-by-design standards before deployment**

Pre-release evaluations for cyber misuse risk, structured red-teaming, transparent system cards, clear intended-use boundaries, and explicit prohibited-use definitions must become baseline expectations. Anthropic’s release process for Mythos — including detailed public system cards and Frontier Red Team disclosure — already demonstrates this practice at the company level. AIWS Trust Standards v1.0 can codify it as an industry-wide expectation, so that responsible practice does not depend on the goodwill of individual firms.

## 2. Controlled access and identity assurance

Once a model has frontier cyber capability, role-based, time-bounded, revocable access tied to strong identity verification becomes part of the safety case itself. The Mythos vendor-environment incident confirms that gated release is necessary but not sufficient: the gates must be hardened, monitored, and continuously verified across every authorized environment.

## 3. Partner and vendor security obligations

Any contractor, cloud environment, security research platform, or vendor with access to a frontier model must meet defined trust requirements for segmentation, credential security, monitoring, and independent auditability. The third-party vendor breach pattern is now a known industry risk; the response must be structural.

## 4. Continuous monitoring and audit trails

For cyber-capable models, governance does not stop at release approval. Monitoring should track abnormal usage patterns, escalation toward prohibited requests, changes in partner behavior, and anomalies across deployment environments. Logs must be durable, reviewable, and tamper-resistant. Anthropic’s commitment to embed automatic detection and blocking of prohibited cybersecurity requests in subsequent models points the entire industry toward this standard.

## 5. Incident exchange and learning loop

Governments, regulators, financial institutions, and major software actors are reacting to Mythos in parallel. The sector now needs structured mechanisms for coordinated disclosure, shared lessons learned, corrective action, and systemic alerts across trusted networks. AIWS proposes to host such a loop as a neutral, democratically anchored convening institution.

# VI. Implications for AIWS Trust Rating

The Mythos period demonstrates that trustworthiness cannot be reduced to a model card, a benchmark score, or a company statement. For frontier cyber-capable AI, trust must be evaluated across the model, the application, the deployment environment, and the surrounding partner ecosystem. This is the role of AIWS Trust Rating.

An AIWS Trust Rating for frontier cyber-capable AI should assess not raw model power, but whether the model and deployment environment together meet auditable thresholds for cyber misuse resistance, access governance, vendor security, monitoring maturity, prohibited-use enforcement, and incident readiness.

Dimension	Required Evidence	Pass/Fail Gate
Cyber Misuse Resistance	Pre-release red-team results, misuse scenarios, exploit-chain controls	No deployment without demonstrated mitigation of high-risk misuse paths
Access Governance	Identity assurance, role-based access, revocation logs, partner access boundaries	No anonymous, unmanaged, or non-revocable access
Vendor & Partner Security	Third-party audits, environment segmentation, credential management, incident drills	No access for vendors lacking minimum security obligations

Monitoring & Auditability	Usage anomaly monitoring, tamper-resistant logs, escalation workflows	No high-risk deployment without durable audit trails
Prohibited-Use Enforcement	Policy definitions, automated detection, human review, appeals process	No deployment without enforceable restrictions on prohibited cyber requests
Incident Readiness	Coordinated disclosure plan, government contact path, public communication protocol	No deployment without an incident-response and learning-loop mechanism

Such a rating will help policymakers, enterprises, and critical infrastructure operators distinguish between models that are merely powerful and models that are trustworthy enough for bounded use. In an era in which cyber-capable models can surface high-impact vulnerabilities at scale, society needs a language of authorization, not only a language of admiration or fear.

## VII. Why This Matters for America at 250

The deeper meaning of the Mythos moment is civilizational. At its semiquincentennial, the United States is asking what kind of leadership it will exercise in the AI Age. The Mythos period offers a concrete answer: American leadership will be defined not only by the power of American models, but by the democratic order — the trusted institutions, standards, and partnerships — within which those models operate.

Mythos is not a story of failure. It is a story of a frontier developer choosing restraint, a coalition of major firms choosing collaboration, a government wrestling with the proper limits of technology in democracy, and an industry beginning to recognize that trust at this scale must be built as infrastructure. The work ahead is to take the seeds of responsible practice that Anthropic and its partners have planted and grow them into a public architecture that the whole democratic world can rely on.

That is why America at 250: A Beacon for the AI Age offers itself as a convening moment — not to celebrate any single company, nor to indict any single incident, but to mark the founding of AIWS Trust Infrastructure as a shared project of governments, frontier developers, critical infrastructure operators, civil society, and democratic citizens. This network does not need to choose between innovation and restraint. The Mythos period shows that the most innovative companies are themselves choosing restraint where it matters most. Our task is to ensure that this choice is no longer left to individual conscience, but is held in common as a civilizational standard.

## Conclusion

Claude Mythos Preview is one of the first major public signs that frontier AI has entered a new strategic phase. Its significance lies not only in technical strength, but in what it reveals about how the AI Age must be governed. Mythos shows that the most consequential AI work now happens at the intersection of capability, restraint, and trust — and that the line between these is held best when companies, governments, and civil society work together across that line.

Anthropic’s choice to gate Mythos, disclose openly, subsidize defenders, maintain safeguards against autonomous weapons and domestic surveillance, and invite the industry into shared work are choices the Boston Global Forum and the AI World Society honor and seek to build upon. Google’s reported \$40 billion investment plan, Amazon’s expanded commitment, and Anthropic’s own infrastructure trajectory make the lesson even clearer: frontier AI companies are becoming part of the operating infrastructure of civilization.

The right response from the broader democratic world is not panic, not unrestricted acceleration, and not critique from a distance. The right response is partnership: the construction of AIWS Trust Infrastructure as the architecture of standards, controls, monitoring, accountability, and trusted order that can make frontier AI worthy of service to democracy, peace, security, and humanity. This is the work America at 250: A Beacon for the AI Age now begins.

## Sources

6. Anthropic, “Project Glasswing: Securing critical software for the AI era,” April 7, 2026 — <https://www.anthropic.com/glasswing>
7. Anthropic Frontier Red Team, “Assessing Claude Mythos Preview’s cybersecurity capabilities,” April 7, 2026 — <https://red.anthropic.com/2026/mythos-preview/>
8. Anthropic, “Project Glasswing initiative page” — <https://www.anthropic.com/project/glasswing>
9. Reuters, “Google to invest up to \$40 billion in AI rival Anthropic,” April 24, 2026 — <https://www.reuters.com/business/google-plans-invest-up-40-billion-anthropic-bloomberg-news-reports-2026-04-24/>
10. Associated Press, “AI startup Anthropic commits \$100 billion to Amazon’s AWS over next 10 years,” April 2026 — <https://apnews.com/article/cffa2cc19f9928d9ac44e44f2d967d36>
11. Reuters, “Anthropic’s Mythos model accessed by unauthorized users, Bloomberg News reports,” April 21, 2026 — <https://www.reuters.com/technology/anthropics-mythos-model-accessed-by-unauthorized-users-bloomberg-news-reports-2026-04-21/>
12. Bloomberg, “Anthropic’s Mythos AI Model Is Being Accessed by Unauthorized Users,” April 21, 2026 — <https://www.bloomberg.com/news/articles/2026-04-21/anthropic-s-mythos-model-is-being-accessed-by-unauthorized-users>
13. Fortune, “Anthropic is giving some firms early access to Claude Mythos to bolster cybersecurity defenses,” April 7, 2026 — <https://fortune.com/2026/04/07/anthropic-claude-mythos-model-project-glasswing-cybersecurity/>
14. TechCrunch, “Unauthorized group has gained access to Anthropic’s exclusive cyber tool Mythos,” April 21, 2026 — <https://techcrunch.com/2026/04/21/unauthorized-group-has-gained-access-to-anthropics-exclusive-cyber-tool-mythos-report-claims/>
15. Foreign Policy, “Anthropic’s Claude Mythos Preview Changes Cyber Calculus,” April 20, 2026 — <https://foreignpolicy.com/2026/04/20/claude-mythos-preview-anthropic-project-glasswing-cybersecurity-ai-hacking-danger/>
16. Nextgov/FCW, “Anthropic’s Glasswing initiative raises questions for US cyber operations,” April 2026 — <https://www.nextgov.com/cybersecurity/2026/04/anthropics-glasswing-initiative-raises-questions-us-cyber-operations/412721/>
17. Congressional Research Service, “Pentagon-Anthropic Dispute over Autonomous Weapon Systems: Potential Issues for Congress” — <https://www.congress.gov/crs-product/IN12669>
18. Associated Press, “Anthropic seeks to debunk Pentagon’s claims about its control over AI technology in military systems,” April 23, 2026

## PART TWO: Three Conference Products Drawn from This Report

### 1. Executive Summary for Conference Packet (One Page)

Mythos and the Path to AIWS Trust Infrastructure in the Age of Frontier Cyber AI

The release of Claude Mythos Preview under Project Glasswing on April 7, 2026 is one of the most consequential moments in the recent evolution of artificial intelligence — and one of the most instructive. Anthropic chose not to make Mythos generally available. The company instead opened the model to a coalition including Amazon Web Services, Apple, Cisco, Google, JPMorganChase, Microsoft, NVIDIA, and Palo Alto Networks, alongside more than forty critical infrastructure organizations, with \$100 million in usage credits and \$4 million in donations to open-source security. This was a deliberate act of restraint and responsibility, exemplifying the kind of conduct the AI Age requires from frontier developers.

The reported plan by Google-parent Alphabet to invest up to \$40 billion in Anthropic, following Amazon's expanded commitment and Anthropic's own cloud-infrastructure agreements, shows that the Mythos moment is not isolated. Anthropic and similar frontier AI developers are becoming part of the strategic AI infrastructure layer of the democratic world. This makes trust no longer a matter of reputation alone; it must become operational architecture.

Mythos has also revealed the deeper truth that no single company, however principled, can carry the full weight of governance for capabilities of this magnitude. A reported third-party vendor environment incident — concerning the deployment ecosystem rather than the model itself — has shown that trust at this scale must extend across every contractor relationship, every credentialing system, and every monitoring loop. This is the public-good architecture that AIWS Trust Infrastructure is designed to supply, as a complement to the responsible practices that Anthropic and its partners have already pioneered.

The Mythos period has also made visible the importance of democratic limits in American AI leadership. Anthropic's stand to maintain safeguards against autonomous weapons and domestic surveillance — supported by an amicus brief from thirty-seven researchers including Google's chief scientist — affirms a principle that the Boston Global Forum and the AI World Society have championed since their founding: the dignity of democratic citizenship must not be subordinated to the convenience of any technology.

At its semiquincentennial, the United States is asking what kind of leadership it will exercise in the AI Age. The answer that the Mythos period suggests is clear. American leadership will be defined not only by the power of American models, but by the democratic order — the trusted institutions, standards, and partnerships — within which those models operate. America at 250: A Beacon for the AI Age offers itself as the convening moment in which that order begins to take shared form: as a partnership of governments, frontier developers, critical infrastructure operators, civil society, and democratic citizens, drawing strength together from the lines a free society draws around its most powerful instruments.

## **2. Agenda Insert for Panel — AIWS Trust Infrastructure for Democracy in the AI Age**

A timely case study for Panel I is the release of Claude Mythos Preview under Anthropic's Project Glasswing on April 7, 2026, together with the reported Google investment plan and Amazon's expanded commitment to Anthropic. These developments illustrate that frontier AI has entered a phase in which advanced models can affect cybersecurity, critical infrastructure, financial stability, and national security at systemic scale — and that the most thoughtful frontier developers are themselves choosing restraint, transparency, and coordinated industry partnership in response.

This panel will examine what the Mythos period teaches about the architecture of trust the AI Age now requires: standards before deployment, controlled access, partner and vendor security, real-time monitoring, auditability, incident response, democratic accountability, and AIWS Trust Rating pathways for frontier cyber-capable systems. Mythos matters not as an Anthropic story alone, but as a moment in which the entire democratic world is being invited to take up the work of building AIWS Trust Infrastructure — the public-good architecture that no single company can be expected to provide on its own.

Optional pull-line for the program: In the AI Age, trust must be built as infrastructure — together.

### **3. Beacon Declaration Paragraph**

We affirm that the founding of AIWS Trust Infrastructure for the AI Age must include standards, governance mechanisms, and trust-rating pathways for frontier AI systems with cyber-capable and high-risk societal impact. We honor the principled work of frontier developers who have chosen restraint, transparency, and coordinated industry partnership as the path of responsible release — and we recognize that the public-good architecture of trust must now be built around their work, by governments, civil society, critical infrastructure operators, and democratic citizens together. We affirm the democratic limits that distinguish American AI leadership: that frontier AI must not be used for fully autonomous weapons, nor for surveillance against the citizens of a free society. In this new era, trust must become operational architecture — a trusted order that safeguards democracy, critical infrastructure, peace, security, and humanity.

### **A Note on Anthropic's Participation**

The Boston Global Forum and the AI World Society warmly invite Anthropic to participate in America at 250: A Beacon for the AI Age — whether through the participation of senior leadership, a written contribution to the conference proceedings, or a continuing partnership in the development of AIWS Trust Standards and the AIWS Trust Rating for frontier cyber-capable AI. The work this report describes is best done with frontier developers, not at a distance from them. BGF and AIWS look forward to walking alongside Anthropic and other principled frontier developers in the founding of the trust architecture that the AI Age requires.