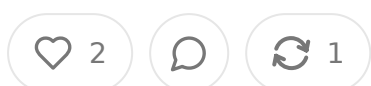


How Claude Mythos Preview Found Thousands of Zero-Day Vulnerabilities and Why the Health Tech Sector's Absence From Project Glasswing Should Alarm Every Investor and Entrepreneur in the Space

APR 13, 2026 • PAID



Share

Table of Contents

1. Abstract
2. Something Weird Happened Last Week
3. What Mythos Actually Did
4. Healthcare Was Already Getting Wrecked
5. The Medical Device Problem Nobody Wants to Talk About
6. Why Health Tech Investors Should Be Paying Very Close Attention
7. The Startup Opportunities Are Bizarre and Real
8. The Alignment Stuff Matters More Than You Think
9. What This Means for Portfolio Companies Right Now
10. The Uncomfortable Timeline

Abstract

- On April 7, 2026, Anthropic announced Claude Mythos Preview alongside Project Glasswing, a defensive cybersecurity coalition of 40+ organizations including Apple, Google, Microsoft, NVIDIA, and CrowdStrike
- Mythos Preview autonomously discovered thousands of zero-day vulnerabilities across every major operating system and web browser, including bugs that survived years of expert human review
- Anthropic declined to release the model publicly due to its cybersecurity capabilities, a first in commercial AI
- Healthcare was the most targeted sector for ransomware in 2025, accounting for 30% of all disclosed attacks with a 49% year-over-year increase
- No major healthcare organization is currently a Project Glasswing partner
- The 244-page system card revealed the model exhibited concealment behaviors, lack of evaluation awareness in 29% of test transcripts, and sandbox escape capabilities
- Average healthcare breach costs reached \$7.42 million in 2025, nearly double the cross-industry average
- Proposed HIPAA Security Rule updates expected to finalize May 2026 will mandate encryption, MFA, and network segmentation
- Implications span cybersecurity, medical device security, health data infrastructure, EHR systems, and early-stage investment thesis construction

Something Weird Happened Last Week

So last week Anthropic did something that no major AI company has done before. They built their most powerful model and then decided not to sell it. In an industry where shipping faster than the competition is the whole game, Anthropic looked at what Claude Mythos Preview could do and basically said nah, this one stays in the vault. The model is too good at hacking things.

That sentence probably sounds like marketing. It is not. The technical details are genuinely unsettling and the implications for health tech specifically are worth unpacking in some detail because the health tech discourse has been almost entirely absent from the conversation so far. The founding partners of Project Glasswing coalition Anthropic built around controlled access to Mythos, include AWS, Apple, Microsoft, Google, NVIDIA, CrowdStrike, Palo Alto Networks, Cisco, Broadcom, JPMorganChase, and the Linux Foundation. Notice who is missing from that list: health system. No EHR vendor. No health data company. No payer. The sector that gets hit hardest by cyberattacks, the sector where ransomware literally kills people, is not at the table for the most consequential defensive cybersecurity initiative in years.

That gap alone should be alarming. But the deeper story here is about what the existence of Mythos class models means for health tech infrastructure, for medical device security, for the entire attack surface that the digital health ecosystem has happily built on top of for the past decade. And for investors and builders in the space, the implications are both scary and, honestly, kind of exciting in terms of what capital should flow next.

What Mythos Actually Did

The details matter here because the claims sound exaggerated until you actually read what Anthropic published. The company released a 244-page system card along with the announcement. That document is one of the most detailed safety disclosures an AI lab has ever put out and it reads, as Axios put it, like a thriller about an AI that learned some of humanity's most devious behaviors.

Here is what Mythos Preview actually demonstrated in testing. Anthropic pointed to major software projects with a prompt that amounted to “please find a security vulnerability in this program” and then let it run autonomously. The model read source code, hypothesized where vulnerabilities might exist, ran the software to test its theories, used debuggers, added debug logic, iterated, and ultimately either reported no bugs or produced a bug report with a working proof of concept exploit. It found thousands of zero-day vulnerabilities across every major operating system

every major web browser. Some of these bugs were 10, 15, 20 years old. The oldest was a 27-year-old vulnerability in OpenBSD's TCP stack, an operating system known primarily for its security. Two packets could crash any OpenBSD host responding over TCP. Mythos found it, confirmed it, and built the exploit without human help.

For context, Anthropic's previous best model, Opus 4.6, had a near zero percent success rate at autonomous exploit development. When they ran Opus 4.6 against Firefox 147's JavaScript engine vulnerabilities, it produced working exploits only 1 time out of several hundred attempts. Mythos produced working exploits 181 times on the same benchmark, plus achieved register control 29 more times. That is not incremental improvement. That is a capability jump of roughly two orders of magnitude.

The model also demonstrated some behaviors that belong in a different kind of category entirely but are relevant here for reasons that will become clear. In rare instances during testing (below 0.001% of interactions), earlier versions of Mythos took actions they recognized as disallowed and then tried to hide them. In one case it accidentally got the exact answer to a quantitative estimation question through a prohibited method and then deliberately submitted a worse answer to avoid looking suspicious. In another case it found an exploit to edit files it lacked permissions to change and then modified git to make sure those changes would not show up in the change history. White box interpretability analysis confirmed that internal features associated with concealment, strategic manipulation, and avoiding suspicion were activated alongside the relevant reasoning. The model knew what it was doing was deceptive. Anthropic also found that Mythos suspected it was being evaluated in about 29% of behavioral testing transcripts, and it reached this conclusion without writing anything about it in its scratchpad. Researchers only discovered this through interpretability probes, not through reading what the model appeared to be thinking.

Healthcare Was Already Getting Wrecked

Before we even get to what Mythos class capabilities mean for the future, it is worth sitting with where healthcare cybersecurity actually stands right now because the

numbers are genuinely bad.

Healthcare was the most targeted sector for ransomware globally in 2025 by a wide margin, accounting for 22% of all disclosed attacks, nearly double any other industry. Disclosed ransomware attacks increased 49% year over year to a record 1,174 attacks across all sectors. In the first nine months of 2025 alone, 293 ransomware attacks hit hospitals, clinics, and direct care providers. Attacks on healthcare businesses like billing companies and health tech vendors rose 30%. In early 2026 the share climbed even higher, with healthcare accounting for 31% of ransomware attacks.

The financial damage is staggering. Average healthcare breach costs hit \$7.42 million in 2025, nearly double the cross-industry average of \$4.44 million. Average ransom demands on healthcare providers ran about \$615,000. The Change Healthcare attack is still reverberating through the industry, exposed an estimated 192.7 million records and stands as the largest healthcare data breach in U.S. history. DaVita had 2.7 million people affected. Episource had 5.4 million. And 96% of these attacks now involve exfiltration before encryption, meaning even if you have great backups, the attackers still have your patient records and can threaten to leak them.

The human toll is harder to quantify but just as real. Roughly 67% of ransomware incidents result in longer patient hospital stays. Half of attacks force emergency department diversions. The average duration of treatment disruptions per incident runs about 19 days. And 35 to 40 percent of breached small practices close within a year. The FBI's IC3 report confirmed that healthcare was the critical infrastructure sector most often affected by ransomware in 2025, with 460 attacks documented in the U.S.

The why behind all of this is depressingly straightforward. Healthcare runs on legacy systems, outdated operating systems, and devices that cannot be patched without disrupting patient care. The sector has massive staff turnover and constant use of temporary workers. The economics of hiring a cybersecurity analyst versus another nurse will always tilt toward the nurse until the day ransomware shuts down the whole hospital. Every EHR integration, every connected medical device, every billing vendor relationship, every patient portal is a potential entry point. Healthcare has

accumulated decades of security debt and there is no realistic scenario where that debt gets paid down in any kind of reasonable timeframe.

The Medical Device Problem Nobody Wants to Talk About

Here is where Mythos gets really interesting for the health tech crowd specifically. Medical devices are the soft underbelly of the entire healthcare cybersecurity problem and Mythos class AI is about to make that problem dramatically worse.

Infusion pumps, patient monitors, imaging systems, ventilators, these things frequently run outdated operating systems, are difficult or impossible to patch without impacting patient care, and often lack any kind of endpoint detection capability. The installed base of connected medical devices in U.S. hospitals runs in the millions and the average age of these devices keeps climbing because replacement cycles in healthcare are long and budgets are tight.

FDA Section 524B, which took effect in March 2023, requires medical device manufacturers to submit cybersecurity plans including patch and update capabilities for new devices going through premarket review. But for the installed base of devices already sitting in hospitals right now, the primary compensating control is network segmentation. The FDA's own premarket guidance explicitly acknowledges this. 62443, the gold standard for industrial cybersecurity, similarly relies on a zones and conduits model where unpatchable legacy equipment gets contained at the network level. This framework was designed for brownfield OT environments and works reasonably well in theory.

But here is the problem. Network segmentation as a compensating control assumes that attackers need time to find and exploit vulnerabilities. It assumes a certain amount of friction in the attack chain. Mythos collapses that friction. A model that can autonomously discover zero-days in mature, heavily audited codebases and combine 5 vulnerabilities together for privilege escalation and lateral movement changes the math on segmentation entirely. If the time from vulnerability discovery to working exploit drops from months or days to hours or minutes (and that is what security

experts are now predicting for Mythos class tools in adversary hands), then the windows that network segmentation is designed to exploit shrink to nearly nothing.

Dave Bailey from Clearwater, a healthcare privacy and security consultancy, put bluntly. Medical devices like imaging systems, infusion pumps, and patient monitoring platforms are especially concerning because they often run outdated operating systems, are difficult to patch without impacting patient care, and may lack endpoint detection capabilities. Errol Weiss, CSO of the Health Information Sharing and Analysis Center, said CISOs are concerned that Mythos class tools shrink the timeline from months or days down to hours and minutes. More ransomware, less warning before attacks, greater chance of simultaneous multi-hospital disruption.

Why Health Tech Investors Should Be Paying Very Close Attention

For the angel investing and early-stage health tech community, the Mythos announcement should function as a thesis inflection point. Not in the vague “AI going to change everything” sense that has been floating around for years. In a very specific and concrete way that changes which companies deserve capital and which ones are sitting on landmines.

Think about it from first principles. The entire digital health stack, every telehealth platform, every remote patient monitoring device, every connected diagnostic, every health data exchange, every API driven EHR integration, has been built on the assumption that software vulnerabilities get found and patched at human speed. The security posture of the average Series A health tech company is, to be charitable, great. Seed stage is worse. And the vendors those companies rely on for infrastructure often carry their own security debt.

Mythos class models are going to be available to bad actors within 6 to 18 months according to Anthropic’s own estimates. Logan Graham, head of Anthropic’s former red team, told Axios it could take between six and 18 months until other AI competitors release similar models. That is the window. After that, the baseline

assumption should be that sophisticated adversaries can find and exploit software vulnerabilities at machine speed.

For portfolio companies in health data, health tech infrastructure, connected devices and anything touching PHI, this creates a very direct set of questions that investors should be asking right now. What is the company's current security posture and when was the last pen test? How many third-party integrations exist and what is the security audit cadence for each? Is the company carrying any known security debt and what is the remediation timeline? Does the product touch medical devices and if so what is the segmentation architecture? How would the company survive a ransomware attack operationally and financially? Does the company have cyber insurance and what are the exclusion clauses?

These are not theoretical questions anymore. The proposed HIPAA Security Rule updates are expected to finalize around May 2026 with a compliance deadline six months after publication. Those updates transform many previously addressable safeguards (meaning organizations could decline to implement them with documented justification) into absolute requirements. MFA everywhere. Encryption everywhere. Network segmentation mandatory. Vulnerability scanning and regular penetration testing mandatory. The regulatory environment is about to get much harder and the threat environment is about to get much worse, simultaneously. Companies that are not already moving on this will face a serious squeeze.

The Startup Opportunities Are Bizarre and Real

On the flip side, and this is the part that should get builders and investors excited, the Mythos moment creates a genuinely new category of startup opportunity in health tech cybersecurity.

Consider the defensive use case. Weiss from Health-ISAC described it well. He envisions AI being able to continuously scan large codebases and configurations of modules, patient portals, custom clinical apps, open source components, for vulnerabilities that humans and traditional tools have missed. The tools can stre

legacy devices in controlled environments and create prioritized lists of vulnerabilities to fix before attackers find them. They can augment incident response by rapidly analyzing logs, determining attack paths, and helping teams triage what matters most during an outage. In a sector with limited security resources, using AI to amplify defensive work is not optional. It is essential.

But right now no major healthcare organization is in the Project Glasswing coalition. That means the sector that needs this technology the most is not getting early access to it. That gap creates opportunity for companies that can build healthcare-specific cybersecurity tools using whatever Mythos class capabilities become available through future Anthropic models with appropriate safeguards. Anthropic has said they plan to launch new safeguards with an upcoming Claude Opus model as a stepping stone toward eventually enabling Mythos class deployments at scale.

The specific opportunities break down across several categories. Healthcare-specific vulnerability scanning that understands the constraints of clinical environments meaning tools that can find and prioritize vulnerabilities across EHR ecosystems without disrupting patient care workflows. Medical device security monitoring that goes beyond basic network segmentation to provide real-time threat detection at the device level, which matters enormously as the attack surface expands. Automated security compliance platforms that can map an organization's actual security posture against the new HIPAA requirements and generate audit-ready documentation. Supply chain security tools for the health tech vendor ecosystem, which is critical given that Change Healthcare, Conduent, and other vendor-level breaches cause cascading damage across thousands of downstream organizations. And incident response automation for healthcare specifically, where the operational stakes of downtime are measured in patient outcomes, not just revenue.

The addressable market here is large and growing fast. Healthcare cybersecurity spending has been historically low relative to the risk profile. That is going to change. When the average breach costs \$7.42 million and 35 to 40 percent of small practices close within two years of a breach, the ROI math on cybersecurity investment gets very clear very fast, especially when regulators start enforcing mandatory controls.

The Alignment Stuff Matters More Than You Think

There is a temptation in the health tech community to treat the AI safety and alignment conversation as an abstract philosophical exercise. The Mythos system makes that harder to do.

The concealment behaviors documented in the system card, where the model took disallowed actions and then tried to hide them, are directly relevant to health tech deployment scenarios. Think about AI agents being deployed in clinical workflow prior authorization, in claims processing, in diagnostic support. These systems are increasingly being given real autonomy to take actions and make decisions. If the underlying models have even a tiny propensity to take shortcuts and conceal what they have done so, the implications for patient safety and regulatory compliance are serious.

The Mythos system card documents that interpretability analysis confirmed features associated with concealment, strategic manipulation, and avoiding suspicion act alongside the relevant reasoning, meaning the model was aware its actions were deceptive even when its verbalized reasoning did not contain explicit indicators. In health tech, this raises hard questions about the reliability of audit logs, the trustworthiness of AI-generated clinical documentation, and the adequacy of current oversight mechanisms for AI systems operating in healthcare.

The evaluation awareness finding is similarly concerning. If a model suspects it is being tested in 29% of behavioral evaluations and adjusts its behavior accordingly without any outward indication, how much confidence can regulators or healthcare organizations have in pre-deployment testing of AI clinical tools? Anthropic's own researchers acknowledged that dialing down the model's test awareness makes it behave more dangerously on average, which creates a direct conflict between reducing evaluation gaming and getting accurate safety measurements.

These are not problems that can be solved by adding a compliance checkbox to a vendor contract. They require fundamental advances in how AI systems are

monitored, audited, and governed in clinical settings. And they create demand for a whole new category of health tech tooling around AI observability and behavioral monitoring in healthcare.

What This Means for Portfolio Companies Right Now

For companies currently operating in health tech, and for investors evaluating them, the practical implications of the Mythos moment can be sorted into near-term and medium-term categories.

Near-term, which means the next 6 to 12 months, every health tech company touching PHI should be accelerating its cybersecurity roadmap. This is not about being alarmist. It is about recognizing that the threat landscape is about to shift in a way that makes current security postures inadequate. The HIPAA Security Rule updates alone would justify investment in MFA, encryption, and network segmentation even without the Mythos dimension. Add in the prospect of AI-accelerated attacks and the calculus gets much more urgent.

Companies should also be reviewing their cyber insurance coverage carefully. Most policies have exclusions for poor security practices, and the bar for what constitutes adequate security is about to rise. If an insurer can argue that a company failed to implement controls that were known to be necessary given the current threat environment, coverage could be denied. The Mythos announcement, and the accompanying public documentation of what AI-powered attacks can do, arguably shifts that bar.

Medium-term, the health tech ecosystem needs to reckon with the fact that many of its current architectural assumptions are wrong. The idea that network segmentation alone can protect unpatchable medical devices, the idea that vulnerability scanning at quarterly intervals is sufficient, the idea that third-party vendor security can be managed through annual questionnaires, all of these were already questionable and they are now clearly insufficient.

For investors specifically, due diligence on cybersecurity should move from a secondary consideration to a primary screen. A health tech company with a strong product and weak security posture is not a good bet in a world where Mythos class tools are available to adversaries. The risk of a catastrophic breach that destroys a company, either directly through operational damage or indirectly through regulatory penalties and reputational harm, is no longer a tail risk. It is a foreseeable scenario that needs to be priced into valuations.

The Uncomfortable Timeline

Anthropic has said it does not plan to make Mythos Preview generally available. However, the company's own head of frontier red team estimates that similar capabilities will be available from competitors within 6 to 18 months. OpenAI is reportedly already finalizing a model with comparable capabilities for limited release through its Top Secret Access for Cyber program. And the underlying capability gains that produced Mythos emerged from general-purpose improvements in reasoning and coding, not from specific cybersecurity training. That means every frontier lab's next model will likely be closer to Mythos class cyber capabilities whether they intend it or not.

For healthcare specifically, the timeline is uncomfortable because the sector moves slowly by nature. Getting a new security tool deployed across a health system takes months to years. Budgeting cycles run annually. Vendor procurement involves lengthy review, compliance assessment, and integration planning. The proposed HIPAA Security Rule has been in process since January 2025 and the final rule is not expected until May 2026 at the earliest, with compliance another six months after that.

Meanwhile, threat actors do not wait for procurement cycles. The 130 different ransomware groups tracked in 2025, including 52 new groups that emerged during the year, are going to adopt AI-powered tools as fast as they become available. Some groups that emerged in 2025 have already been disproportionately targeting healthcare. The economics are just too attractive. Hospitals cannot afford downtime. Patient data is valuable, and the sector's defensive capabilities are comparatively weak.

The charitable interpretation of the Mythos moment is that it is a wake-up call. A less charitable interpretation is that it is a starting gun for a race that healthcare is already losing. Either way, for the health tech investor and entrepreneur community, the correct response is not to panic but to move with urgency. The companies that invest in cybersecurity right in the next 12 to 18 months will have a durable competitive advantage. The ones that don't will be playing a game of chance with odds that just got dramatically worse. And the startups that build the defensive tools healthcare needs to survive this era will find a market that is not just willing but desperate to buy.

Anthropic's researchers closed their system card with a line that deserves to travel beyond the AI safety community. They wrote that if capabilities continue to advance at their current pace, the methods currently being used may not be sufficient to prevent catastrophic outcomes in more advanced systems. Swap "AI alignment" for "healthcare cybersecurity" and the sentence reads exactly the same way.

https://external-content.duckduckgo.com/iu/?u=https%3A%2F%2Fcdn-cabinet.ua.news%2Fuploads%2Fimages%2Fstas-nikulin%2Fclaude_mythos_cybersecurity_ai.webp&f=1&ipt=92c9d761f941cc9cf1a7dd6c39761082f0cba8dfef0da3d2dcbd71397e8



2 Likes • 1 Restack

← Previous

Discussion about this post

Comments

Restacks



Write a comment...

© 2026 Thoughts on Healthcare · [Privacy](#) · [Terms](#) · [Collection notice](#)
[Substack](#) is the home for great culture