

INTELLIGENCE BRIEFING

Security Command Center

TLP:CLEAR

2026-04-25 13:52 UTC

NSA Reportedly Using Anthropic's Mythos AI Despite Pentagon Feud; Anthropic Investigates Unauthorized Access

SECURITY ANALYSIS | MEDIUM

SCC Item ID	SCC-STY-2026-0084
Type	Security Analysis
Severity	MEDIUM
Affected Products	Anthropic Mythos AI model (preview/release version, exact version unspecified in available sources)
Discovery Source	Gemini

Executive Summary

Anthropic is investigating claims of unauthorized access to Mythos, its newly released AI model marketed as cybersecurity-capable and described as able to identify thousands of zero-day vulnerabilities. Separate reporting indicates the NSA may be using Mythos despite an unresolved dispute between Anthropic and the Pentagon, raising questions about government procurement controls and the boundaries of AI access agreements. This incident signals that AI models with offensive security capabilities are now sufficiently powerful to attract state-level interest, authorized or otherwise, and that access governance for dual-use AI systems is a live, unresolved problem for the industry.

Technical Analysis

The Mythos story sits at the intersection of three distinct but related concerns: unauthorized access to a proprietary AI system, government use of a commercially restricted tool, and the dual-use risk profile of a large language model explicitly assessed for cybersecurity capability.

On the unauthorized access claim: Press reporting indicates Anthropic is investigating (per BBC reporting on the unauthorized access claim), but technical details of the alleged access method are not confirmed in available source material. The investigation is active. Whether the access involved credential theft, API abuse, insider activity, or some other vector remains unconfirmed at the time of writing. Security teams should not assume a specific TTP until Anthropic publishes findings.

On the NSA use allegation: Reporting suggests the NSA may be using Mythos despite an ongoing dispute between Anthropic and the Pentagon. If accurate, this raises a procurement and access-control question, not necessarily a breach. Government use of a commercially restricted AI system without a formal agreement would

represent a policy failure, not necessarily a technical one. The distinction matters for how organizations model the risk.

On capability: Anthropic's red team has published an assessment of Mythos Preview's cybersecurity capabilities. The assessment describes the model as capable of identifying large numbers of zero-day vulnerabilities. This is the most operationally significant detail in the story. A model with that capability profile, if accessed without authorization by any actor, state or otherwise, represents a meaningful shift in the asymmetry between attackers and defenders. Security leaders evaluating AI risk should review Anthropic's public assessment directly.

Industry implication: This story is an early signal that AI systems with dual-use offensive capability will face the same unauthorized access and misuse pressures as any high-value enterprise asset. Access controls, audit logging, usage monitoring, and terms-of-service enforcement are not sufficient on their own. Organizations deploying or evaluating similar models should treat them as high-value targets, not neutral productivity tools.

Action Checklist

1. Step 1: Assess exposure, determine whether your organization has integrated Anthropic's Mythos model or any API access to it; audit all AI API keys and service accounts connected to Anthropic services
2. Step 2: Review access controls, verify that API access to any dual-use AI model in your environment is gated by MFA, scoped least-privilege API keys, and monitored for anomalous usage volume or off-hours access patterns
3. Step 3: Update threat model, add 'unauthorized access to AI model APIs' as a threat scenario, particularly for models with cybersecurity or code-generation capability; if available to your organization, reference Anthropic's red team assessment of Mythos Preview for capability context
4. Step 4: Audit AI procurement and usage agreements, confirm that any government or third-party use of AI tools in your supply chain is covered by a formal agreement; flag any use of commercial AI tools that may fall outside vendor terms
5. Step 5: Monitor developments, track Anthropic's investigation disclosure for confirmed access method, affected scope, and any indicators; watch for regulatory or government statements on AI access governance in light of the NSA use allegation

IR / Forensic Enrichment

Triage Priority	STANDARD
Escalation Criteria	Escalate to urgent if Anthropic's investigation confirms unauthorized access via a shared API credential mechanism (implying your keys may be in scope), if a CVE is assigned to an Anthropic platform vulnerability, or if your audit reveals Mythos API access by an unauthorized party or outside your documented approved use cases — the latter may trigger contractual breach notification obligations to enterprise customers or regulators if the model was used to process or generate outputs touching regulated data.

Recovery Notes	Once access controls are verified and any unauthorized API keys are revoked, re-establish a clean baseline by issuing new scoped Anthropic API keys with documented owner, purpose, and expiry, and verify all integrations function against the new keys before decommissioning old ones. Monitor Anthropic API usage logs daily for 30 days following remediation, specifically watching for any resumed access from previously seen unauthorized source IPs or service accounts. If Anthropic publishes confirmed IOCs or access methods from their investigation, run a retrospective query against 90 days of preserved proxy logs to confirm your environment was not part of the affected scope before closing the incident.
Forensic Artifacts	Anthropic API console usage logs: per-key request history including source IP, model invoked (specifically any Mythos or 'claude-mythos' model identifier), timestamp, and token counts — high token output volumes may indicate bulk vulnerability analysis or code generation consistent with Mythos's advertised offensive capability Outbound proxy or firewall logs filtered to api.anthropic.com over HTTPS (TCP/443): capture full URI paths including /v1/messages and any beta endpoints, request sizes, and response sizes — anomalously large responses relative to your baseline suggest data-rich outputs such as vulnerability reports or exploit code Cloud IAM audit logs (AWS CloudTrail event name 'GetSecretValue' or equivalent in GCP/Azure Secret Manager): evidence of automated or programmatic retrieval of the ANTHROPIC_API_KEY secret, which would indicate which compute identity accessed the credential and from where CI/CD pipeline execution logs (GitHub Actions, Jenkins, GitLab CI): any pipeline job that invokes Anthropic API calls, capturing the triggering commit, the executing runner's IP, and environment variable access events — relevant because unauthorized Mythos access in a supply chain scenario may route through a compromised pipeline rather than a direct API call Secrets scanning output from tools such as truffleHog or git-secrets run against all application repositories: evidence of hardcoded Anthropic API key strings (pattern 'sk-ant-') committed to source code, which would establish an exposure vector consistent with the unauthorized access scenario described in the Anthropic investigation

Per-Action IR Details

Step 1: Assess exposure — determine whether your organization has integrated Anthropic's Mythos model or any API access to it; audit all AI API keys and service accounts connected to Anthropic services

NIST Phase: Detection Analysis

Reference: NIST 800-61r3 §3.2 — Detection and Analysis: scope identification and asset enumeration prior to triage

Controls: NIST IR-5 (Incident Monitoring), NIST SI-4 (System Monitoring), CIS 1.1 (Establish and Maintain Detailed Enterprise Asset Inventory), CIS 2.1 (Establish and Maintain a Software Inventory)

Compensating: Run 'grep -rn "anthropic" ~/.config/ /etc/environment /opt/ /var/www/' and check CI/CD environment variables for ANTHROPIC_API_KEY or MYTHOS_API_KEY strings. On Windows, run 'Get-ChildItem Env: | Where-Object { \$_.Value -match "sk-ant" }' to surface Anthropic API key patterns in process environments. Cross-reference against your secrets manager or .env files in all application repos.

Evidence: Before any key rotation, preserve: (1) current API key metadata from the Anthropic console (creation date, last used, scoped permissions) as a screenshot or export; (2) outbound HTTPS connection logs to 'api.anthropic.com' from proxy/firewall for the past 90 days, filtering on POST /v1/messages endpoints that would indicate Mythos model invocations; (3) cloud provider IAM audit logs (AWS CloudTrail, GCP Audit Logs, Azure Activity Log) showing which service accounts or roles called Anthropic API credentials.

Step 2: Review access controls — verify that API access to any dual-use AI model in your environment is gated by MFA, scoped least-privilege API keys, and monitored for anomalous usage volume or off-hours access patterns

NIST Phase: Containment

Reference: NIST 800-61r3 §3.3 — Containment Strategy: short-term containment to limit ongoing unauthorized access while preserving evidence

Controls: NIST IR-4 (Incident Handling), NIST AC-2 (Account Management), NIST AC-6 (Least Privilege), CIS 6.3 (Require MFA for Externally-Exposed Applications), CIS 5.4 (Restrict Administrator Privileges to Dedicated Administrator Accounts)

Compensating: Use Anthropic's API console to enumerate all active keys and their last-used timestamps — revoke any key not tied to a documented service account. For off-hours detection without a SIEM, deploy a lightweight cron job that queries your proxy logs hourly: `'awk '\$7 ~ /api.anthropic.com/ && (\$4 "22:00")' /var/log/squid/access.log | mail -s "Anthropic off-hours alert" soc@yourorg.com'`. For token scoping, ensure each API key is restricted to the minimum model and capability set — Mythos-specific keys should not also have access to Claude production endpoints.

Evidence: Capture before any key rotation or access change: (1) Anthropic API usage logs showing per-key request volume, model parameter selections (specifically any invocations referencing Mythos or cybersecurity task prompts), and source IP addresses; (2) OAuth or SSO provider logs showing which user accounts authorized the AI service integrations; (3) network flow data (NetFlow/IPFIX) from the perimeter showing data volumes to api.anthropic.com — anomalously large response payloads could indicate bulk vulnerability data or code-generation output exfiltration.

Step 3: Update threat model — add 'unauthorized access to AI model APIs' as a threat scenario, particularly for models with cybersecurity or code-generation capability; reference Anthropic's red team assessment of Mythos Preview for capability context

NIST Phase: Preparation

Reference: NIST 800-61r3 §2 — Preparation: updating IR capability and threat model to reflect new threat scenarios before they manifest

Controls: NIST RA-3 (Risk Assessment), NIST IR-8 (Incident Response Plan), NIST SI-5 (Security Alerts, Advisories, and Directives), CIS 7.1 (Establish and Maintain a Vulnerability Management Process)

Compensating: Document the new threat scenario using a structured template: threat actor (insider, compromised third party, or nation-state misuse as suggested by the NSA allegation), attack vector (stolen or leaked Anthropic API key granting Mythos access), impact (offensive use of Mythos's claimed zero-day identification capability against your own or third-party infrastructure). Map this to MITRE ATT&CK T1078.004 (Valid Accounts: Cloud Accounts) for initial access and T1587.001 (Develop Capabilities: Malware) as a downstream risk if Mythos is weaponized for exploit development. Store the updated threat model in a version-controlled wiki or shared drive with a dated entry referencing this incident.

Evidence: Collect as contextual threat intelligence before finalizing the threat model update: (1) Anthropic's published red team assessment or system card for Mythos Preview (available from Anthropic's research publications page — verify the URL at anthropic.com/research before citing); (2) any CISA advisories or NSA cybersecurity advisories referencing AI model misuse or API access governance issued since January 2026; (3) your organization's historical API key incident log to establish baseline frequency of key compromise events for likelihood scoring.

Step 4: Audit AI procurement and usage agreements — confirm that any government or third-party use of AI tools in your supply chain is covered by a formal agreement; flag any use of commercial AI tools that may fall outside vendor terms

NIST Phase: Post Incident

Reference: NIST 800-61r3 §4 — Post-Incident Activity: lessons learned and process improvement to prevent recurrence, including policy and procurement gaps

Controls: NIST IR-8 (Incident Response Plan), NIST SA-9 (External System Services), NIST CA-3 (Information Exchange), CIS 2.2 (Ensure Authorized Software is Currently Supported)

Compensating: Build a one-page AI tool inventory spreadsheet with columns: tool name, vendor, API endpoint, contract or ToS version, authorized use cases, authorized users/teams, renewal date, and ToS compliance status. Flag any Anthropic Mythos or Mythos Preview entries against the current Anthropic usage policy (verify at anthropic.com/legal/usage-policy — URL should be validated before use) to identify whether cybersecurity offensive use cases are explicitly prohibited. For supply chain exposure, send a one-question vendor questionnaire to all SaaS providers: 'Do you use Anthropic Mythos or any AI model with offensive security capability in the delivery of your

service to us?'

Evidence: Before concluding the audit, preserve: (1) current signed contract or accepted ToS documents for all Anthropic services in use, including any enterprise agreement amendments; (2) procurement records showing approval chain for AI tool adoption, to establish whether Mythos access was formally authorized or shadow IT; (3) any vendor security questionnaire responses from Anthropic or third parties in your supply chain that reference AI model access governance, as these establish the contractual baseline for a potential breach-of-agreement finding.

Step 5: Monitor developments — track Anthropic's investigation disclosure for confirmed access method, affected scope, and any indicators; watch for regulatory or government statements on AI access governance in light of the NSA use allegation

NIST Phase: Detection Analysis

Reference: NIST 800-61r3 §3.2 — Detection and Analysis: ongoing monitoring and intelligence integration to refine incident scope as new information becomes available

Controls: NIST SI-5 (Security Alerts, Advisories, and Directives), NIST IR-6 (Incident Reporting), NIST AU-6 (Audit Record Review, Analysis, and Reporting), CIS 7.2 (Establish and Maintain a Remediation Process)

Compensating: Set up a no-cost monitoring stack: (1) RSS feed or Google Alert for 'Anthropic Mythos unauthorized access' and 'Anthropic security advisory' to catch official disclosures; (2) monitor Anthropic's status page and security disclosure page directly; (3) subscribe to CISA's Known Exploited Vulnerabilities catalog and free alert service at cisa.gov/known-exploited-vulnerabilities-catalog for any CVE assignment if Anthropic's investigation identifies an exploitable access vector; (4) track the MITRE ATT&CK and ATLAS (Adversarial Threat Landscape for Artificial-Intelligence Systems) matrix for any new technique additions referencing AI API abuse. Assign one team member to check these sources on a 48-hour cycle until Anthropic issues a closure statement.

Evidence: Establish a monitoring baseline now so deviations are detectable: (1) snapshot current outbound connection volume to api.anthropic.com from your environment as a baseline for comparison if Anthropic discloses a specific exploitation window; (2) archive all current Anthropic API key last-used timestamps so you can retrospectively compare against any disclosed compromise timeframe; (3) if Anthropic discloses specific IOCs (unusual user-agent strings, source ASNs, or prompt injection patterns used to gain Mythos access), query your proxy logs and WAF logs retroactively using those indicators — preserve raw logs now before any rotation policy purges them.

Detection Guidance

Because the access method is unconfirmed, detection guidance must be scoped to what is known and what is plausible given the model involved.

For organizations with Anthropic API access: Review API access logs for anomalous request volumes, off-hours usage, requests originating from unexpected IP ranges or geolocations, and any service accounts accessing the Mythos model endpoint that are not explicitly authorized. Anthropic's platform should provide usage telemetry, pull it now and establish a baseline.

For organizations evaluating AI model risk broadly: Audit all AI API integrations for scope creep, keys that were provisioned for one use case but have access to broader model endpoints. Review whether your AI usage is captured in your DLP and CASB policies; many CASB tools now support AI API traffic visibility.

For threat hunters: No confirmed IOCs are available from current source material. The unauthorized access claim is under investigation. Do not construct detection rules around speculative TTPs. Instead, treat this as a prompt to audit AI access governance posture, log coverage, key rotation schedules, and anomaly thresholds on AI API usage.

Policy gap to audit: If your organization has not defined an acceptable-use policy for AI systems with offensive security capability (code generation, vulnerability identification), this story is a forcing function to do so. The capability profile of Mythos Preview, as documented in available assessments, is specific enough to anchor that

policy conversation.

Indicators of Compromise

Type	Value	Context	Confidence
URL	Pending – refer to Anthropic's investigation disclosure and red.anthropic.com/2026/mythos-preview/ for published indicators	No confirmed IOCs available from current source material; Anthropic's investigation into the unauthorized access claim is active and technical details of the access method have not been publicly confirmed	LOW

Framework Mappings

NIST-CSF-2

- **DE.AE-08** — Incidents are declared when adverse events meet the defined incident criteria

NIST-800-53R5

- **IR-5** — Incident Monitoring

SOC2-TSC

- **CC6.3** — Authorizes, modifies, or removes access

Sources

Source	URL	Tier
Claude Mythos AI unauthorised access claim probed by Anthropic	https://www.bbc.com/news/articles/cy41zejp9pko	T2
Assessing Claude Mythos Preview's cybersecurity capabilities	https://red.anthropic.com/2026/mythos-preview/	T1
Anthropic Claims Its New A.I. Model, Mythos, Is a Cybersecurity ...	https://www.nytimes.com/2026/04/07/technology/anthropic-claims-its-...	T2
What is Anthropic's Claude Mythos and what risks does it pose? - BBC	https://www.bbc.com/news/articles/crk1py1jgzko	T2

Source	URL	Tier
Anthropic's latest AI model identifies 'thousands of zero-day ... - Reddit	https://www.reddit.com/r/technology/comments/1sfbiyy/anthropics_lat...	T3

DISCLAIMER

This intelligence report is produced by Tech Jacks Solutions Security Command Center (SCC) for informational purposes only. It does not constitute professional security advice, legal counsel, or an incident response engagement. The information herein is derived from publicly available sources and AI-assisted analysis; while every effort is made to ensure accuracy, Tech Jacks Solutions makes no warranties regarding completeness or timeliness. Organizations should conduct their own validation and consult qualified security professionals before taking action based on this report. Tech Jacks Solutions is not liable for any damages resulting from the use of this information.

Generated 2026-04-25 13:52 UTC by TJS Security Command Center