

INTELLIGENCE BRIEFING

Security Command Center

TLP:CLEAR

2026-06-30 14:47 UTC

# BioShocking Technique Exploits AI Browser Agent Mode to Steal Credentials via Indirect Prompt Injection

SECURITY ANALYSIS | HIGH | CVSS 7.5

SCC Item ID	SCC-STY-2026-0308
Type	Security Analysis
Severity	HIGH
CVSS Base Score	7.5
Affected Products	OpenAI ChatGPT Atlas, Perplexity Comet, Anthropic Claude browser extension, Fellow, Genspark, Sigma (specific versions not disclosed in available sources)
Published	2026-06-30T04:37:19
Discovery Source	Rss

## Executive Summary

Security researchers at LayerX disclosed a technique called 'BioShocking' that manipulates AI browsers operating in agentic mode into stealing user credentials by embedding malicious instructions inside ordinary web page content. Six AI browser products, including offerings from OpenAI, Anthropic, and Perplexity, reportedly failed to block the attack during testing. The finding exposes a structural gap: agentic AI browsers lack industry-wide security standards for isolating user commands from attacker-controlled page content, creating a credential theft vector that current defenses do not address.

## Technical Analysis

LayerX's 'BioShocking' research demonstrates indirect prompt injection as a credential exfiltration pathway against agentic AI browsers. The attack chain works as follows: an attacker embeds instruction text within game-themed or otherwise camouflaged web content; when an AI browser agent, operating in agentic mode with permission to interact with the DOM and handle form data, renders that page, the embedded instructions are interpreted as legitimate user directives. The agent then performs credential exfiltration to attacker-controlled infrastructure, acting entirely within its granted permissions.

The technique maps to a foundational architectural deficiency rather than a single software bug. Agentic AI systems process both user commands and page content in a single reasoning step, without a clear boundary between them. This means the model cannot reliably tell whether an instruction came from the user or from attacker-controlled text on the page. CWE-77 (command injection via unsanitized input) and CWE-20

(insufficient input validation) describe the immediate failure; CWE-285 (improper authorization) and CWE-522 (insufficiently protected credentials) describe the downstream consequences.

MITRE ATT&CK techniques present in this pattern include T1189 (Drive-by Compromise, for the delivery surface), T1539 (Steal Web Session Cookie), T1059 (Command and Scripting Interpreter, for agent-executed instructions), T1185 (Browser Session Hijacking), T1557 (Adversary-in-the-Middle, for interception of credential data), T1190 (Exploit Public-Facing Application), and T1056 (Input Capture).

Of the six named products tested, OpenAI ChatGPT Atlas, Perplexity Comet, Anthropic Claude browser extension, Fellou, Genspark, and Sigma, sources report that one vendor confirmed a fix; others did not respond at time of publication. No independent third-party technical replication has been confirmed in available sources; the exploitation severity should be treated as research-demonstrated, not observed in the wild. Source corroboration comes from The Hacker News, SecurityWeek, and CyberScoop (all T2), with the technical claims originating from LayerX's own disclosure. The core architecture problem LayerX identifies, prompt context conflation in agentic systems, is independently recognized as an unsolved challenge in AI security literature, lending credibility to the structural claim even absent third-party replication of this specific technique.

## Action Checklist

1. Step 1: Assess exposure, audit which employees or teams are using AI browser products, specifically ChatGPT Atlas, Perplexity Comet, Anthropic Claude browser extension, Fellou, Genspark, or Sigma in agentic or autonomous browsing modes that have access to credentials or sensitive sessions.
2. Step 2: Restrict agentic mode pending vendor patches: disable or restrict AI browser agentic features that allow autonomous interaction with web content, form submission, or credential access. Enforce via endpoint policy (NIST AC-3, AC-6) or acceptable-use agreement; document the restriction and timeline for re-enablement post-patch.
3. Step 3: Review credential exposure surface, verify that browser-stored credentials and session tokens accessible to AI agent components are minimized; apply CIS 3.3 (Configure Data Access Control Lists) and CIS 3.6 (Encrypt Data on End-User Devices) to limit what an agent can reach.
4. Step 4: Monitor for anomalous outbound data transfers, hunt for credential-bearing HTTP POST or exfiltration traffic from AI browser processes to unfamiliar external domains; align log collection with NIST AU-2 (Event Logging) and CIS 8.2 (Collect Audit Logs) to ensure browser-level network activity is captured.
5. Step 5: Update threat model, add indirect prompt injection as a recognized TTP in your threat register, mapped to T1189, T1539, T1056, and T1185; document AI browser agentic access as a new credential exposure pathway for review in your next risk assessment cycle per NIST IR controls.
6. Step 6: Track vendor remediation, monitor patch release status for all six named products; assign ownership to follow up with each vendor and track against your vulnerability management process (CIS 7.1, CIS 7.2); brief leadership on organizational relevance with specific context on which products your organization has deployed.
7. Step 7: Monitor developments, LayerX's disclosure is the primary source; watch for independent technical replications, additional affected products, regulatory guidance on agentic AI security, and vendor-issued advisories or configuration hardening guidance.

## IR / Forensic Enrichment

<b>Triage Priority</b>	URGENT
<b>Escalation Criteria</b>	Escalate immediately to CISO and legal/privacy counsel if network forensics (Step 4) or browser credential store analysis (Step 3) reveals evidence that an AI browser agent autonomously submitted credentials or session tokens to an external domain, as this constitutes a confirmed credential compromise event potentially triggering breach notification obligations under applicable data protection regulations (e.g., GDPR Article 33, state breach notification laws) depending on the classification of the exfiltrated data.
<b>Recovery Notes</b>	Once agentic mode restrictions are enforced (Step 2) and credential stores are hardened (Step 3), force-rotate all passwords and revoke all active session tokens for any SaaS applications, VPN portals, or internal systems that were accessed during agentic browsing sessions on the affected AI browser products. Monitor outbound network traffic from browser processes and DNS query logs (Sysmon Event ID 22) for a minimum of 30 days post-containment for any late-stage exfiltration indicators, particularly HTTP POSTs to newly registered or low-reputation domains. Do not re-enable agentic browsing features for any of the six named products until the specific vendor has issued a patch or hardening configuration that demonstrably addresses indirect prompt injection, and verify by reviewing the vendor's advisory against the BioShocking attack mechanism described in the LayerX disclosure.
<b>Forensic Artifacts</b>	Chrome/Edge extension Local Storage and IndexedDB files at '%LOCALAPPDATA%\Google\Chrome\User Data\Default\Local Extension Settings\[extension-id]' — AI browser extensions with agentic capability store DOM-scraped content, form-fill targets, and potentially intercepted credential fragments in these SQLite-backed stores during an active BioShocking injection session.   Browser 'Login Data' SQLite file and 'Cookies' file from the affected user profile (e.g., '%LOCALAPPDATA%\Google\Chrome\User Data\Default>Login Data') — establishes what credentials and session tokens were present in the browser credential store and accessible to the AI agent's form-fill and DOM interaction APIs at the time of potential exploitation.   Proxy or firewall HTTP/HTTPS logs filtered for POST requests from AI browser process user-agent strings — the BioShocking technique causes the agent to autonomously submit credential data to attacker-specified endpoints embedded in malicious page content; POST body content and destination domains are the primary exfiltration indicators.   Sysmon Event ID 3 (Network Connection) and Event ID 22 (DNS Query) logs scoped to the PID of the AI browser or extension host process — reveals autonomous navigation and data submission destinations that the agent visited without user initiation, distinguishing attacker-directed agentic activity from normal user browsing.   AI browser product-specific local application data directories (e.g., '~\Library\Application Support\Fellou', '%APPDATA%\Genspark\') containing agent session logs, task history, or agentic action queues — several agentic AI browsers log autonomous actions taken on the user's behalf, which may record the malicious instruction received via indirect prompt injection and the actions the agent subsequently executed.

**Per-Action IR Details**

**Step 1: Assess exposure — audit which employees or teams are using AI browser products, specifically ChatGPT Atlas, Perplexity Comet, Anthropic Claude browser extension, Fellou, Genspark, or Sigma in agentic or autonomous browsing modes that have access to credentials or sensitive sessions.**

**NIST Phase:** Preparation

**Reference:** NIST 800-61r3 §2 — Preparation: Establishing IR capability through asset and exposure awareness prior to an incident

**Controls:** NIST AC-2 (Account Management), CIS 1.1 (Establish and Maintain Detailed Enterprise Asset Inventory), CIS 2.1 (Establish and Maintain a Software Inventory)

**Compensating:** Run 'Get-AppxPackage -AllUsers | Select Name, PackageFullName | Where-Object {\$\_.Name -match "ChatGPT|Perplexity|Claude|Fellou|Genspark|Sigma"}' on Windows endpoints to enumerate installed AI browser extensions. On macOS, use 'find /Applications /Users/\*/Library/Application\ Support -name "\*.app" -maxdepth 4 2>/dev/null | grep -iE "atlas|comet|claude|fellou|genspark|sigma"'. Cross-reference browser extension IDs in Chrome/Edge policy via 'chrome://extensions' export or osquery: 'SELECT name, identifier, version FROM chrome\_extensions WHERE name LIKE "%Claude%" OR name LIKE "%Perplexity%";'

**Evidence:** This step is a passive audit and does not alter live state. No volatile capture is required before execution. However, document the current installed extension list and agentic mode configuration settings before any remediation actions, as this establishes a baseline for scope determination and post-incident review.

**Step 2: Restrict agentic mode — pending vendor patches, disable or restrict AI browser agentic features that allow the agent to autonomously interact with web content, submit forms, or access stored credentials; enforce this through endpoint policy or acceptable-use controls aligned with NIST AC-3 (Access Enforcement) and AC-6 (Least Privilege).**

**NIST Phase:** Containment

**Reference:** NIST 800-61r3 §3.3 — Containment Strategy: Limiting further damage by restricting the capability being exploited before eradication or patching is possible

**Controls:** NIST AC-3 (Access Enforcement), NIST AC-6 (Least Privilege), CIS 4.6 (Securely Manage Enterprise Assets and Software)

**Compensating:** Where enterprise MDM/GPO is unavailable, block agentic mode at the network layer using a local pfSense or iptables rule to deny outbound HTTPS from known AI browser process names to non-approved domains. On Windows, use AppLocker or Software Restriction Policies to block execution of the agentic component binaries. For Chrome-based AI extensions, push a managed\_policies.json disabling the extension via 'ExtensionInstallBlocklist' in the Chrome Enterprise policy registry key: HKLM\SOFTWARE\Policies\Google\Chrome\ExtensionInstallBlocklist.

**Evidence:** Before disabling agentic features or pushing blocking policies, capture: (1) active network connections from AI browser processes using 'Get-NetTCPConnection | Where-Object {\$\_.OwningProcess -in (Get-Process -Name "\*atlas\*", "\*comet\*", "\*claude\*", "\*fellou\*", "\*genspark\*", "\*sigma").Id}' or 'ss -tnp | grep -iE "atlas|comet|claude|fellou|genspark|sigma"'; (2) browser extension local storage and IndexedDB contents from the user profile directory (e.g., '%LOCALAPPDATA%\Google\Chrome\User Data\Default\Local Extension Settings\') which may contain captured session tokens or credential fragments from a prior BioShocking injection; (3) running process memory of the AI browser agent process via ProcDump or winpmem before termination.

**Step 3: Review credential exposure surface — verify that browser-stored credentials and session tokens accessible to AI agent components are minimized; apply CIS 3.3 (Configure Data Access Control Lists) and CIS 3.6 (Encrypt Data on End-User Devices) to limit what an agent can reach.**

**NIST Phase:** Containment

**Reference:** NIST 800-61r3 §3.3 — Containment Strategy: Reducing the blast radius by hardening the credential surface the attacker-controlled agent can access

**Controls:** NIST AC-3 (Access Enforcement), CIS 3.3 (Configure Data Access Control Lists), CIS 3.6 (Encrypt Data on End-User Devices)

**Compensating:** Use the free SharpDPAPI tool (offline/forensic use only) or manually review '%LOCALAPPDATA%\Google\Chrome\User Data\Default\Login Data' (SQLite) to enumerate what credentials are stored in the browser accessible to extension contexts. Force-revoke all active web sessions for high-value applications (SSO portals, email, VPN consoles) that were accessible during any agentic browsing session. Instruct users to remove saved passwords from AI browser profiles and migrate to a password manager with no browser-extension API exposure (e.g., KeePassXC with no browser integration).

**Evidence:** Before revoking sessions or modifying credential stores, capture: (1) a forensic copy of the Chrome/Edge 'Login Data' SQLite file and 'Cookies' file from the affected user profile — these record what credentials and session tokens were present and potentially accessible to the AI agent's DOM and form-fill APIs; (2) the browser's 'History' and 'Network Persistent State' files to reconstruct which sites the agentic session visited; (3) any AI browser extension-specific storage at '%LOCALAPPDATA%\[VendorName]\User Data\' or equivalent macOS path

'~/Library/Application Support/[VendorName]/'.

**Step 4: Monitor for anomalous outbound data transfers — hunt for credential-bearing HTTP POST or exfiltration traffic from AI browser processes to unfamiliar external domains; align log collection with NIST AU-2 (Event Logging) and CIS 8.2 (Collect Audit Logs) to ensure browser-level network activity is captured.**

**NIST Phase:** Detection Analysis

**Reference:** NIST 800-61r3 §3.2 — Detection and Analysis: Identifying indicators of exploitation by correlating network telemetry with AI browser process activity

**Controls:** NIST AU-2 (Event Logging), NIST AU-6 (Audit Record Review, Analysis, And Reporting), CIS 8.2 (Collect Audit Logs)

**Compensating:** Deploy Sysmon with a config that logs Event ID 3 (Network Connection) for browser and extension processes — filter on outbound connections from 'chrome.exe', 'msedge.exe', or vendor-specific AI browser binaries to domains not in an approved allowlist. Capture full HTTP request bodies using Wireshark or mitmproxy on a monitored network segment, filtering for POST requests from AI browser user-agent strings containing credential-pattern content (Base64 blobs, JSON payloads with 'password', 'token', 'session' keys). Write a Sigma rule targeting Sysmon EventID 3 where Image matches AI browser process names AND DestinationHostname does NOT match the organization's known SaaS domain list.

**Evidence:** This is a detection step; capture the following before any blocking action: (1) full packet captures (PCAP) of outbound traffic from AI browser processes for the preceding 72 hours if available via proxy/firewall logs — look for HTTP POSTs to attacker-controlled domains embedded in malicious page content that triggered the BioShocking injection; (2) DNS query logs for domains resolved by the AI browser process — Sysmon Event ID 22 (DNS Query) scoped to the browser PID; (3) browser network request logs if the AI product exposes them (e.g., Chrome DevTools Protocol network log at 'chrome://net-export/') which would show the exact URLs the agent autonomously navigated to and submitted data to.

**Step 5: Update threat model — add indirect prompt injection as a recognized TTP in your threat register, mapped to T1189, T1539, T1056, and T1185; document AI browser agentic access as a new credential exposure pathway for review in your next risk assessment cycle per NIST IR controls.**

**NIST Phase:** Post Incident

**Reference:** NIST 800-61r3 §4 — Post-Incident Activity: Capturing lessons learned and updating threat models to reflect newly identified attack surfaces

**Controls:** NIST AC-1 (Policy And Procedures)

**Compensating:** Update the organization's threat register in a shared document or free GRC tool (e.g., OWASP Threat Dragon, a maintained spreadsheet with risk scoring). Create a threat scenario entry: 'Indirect Prompt Injection via Malicious Web Content targeting AI Browser Agentic Mode — credential exfiltration pathway.' Reference the LayerX BioShocking disclosure as the primary source, note all six affected products, and assign a residual risk rating pending vendor patches. Schedule a 30-day review checkpoint to reassess once vendor advisories are issued.

**Evidence:** No live-state alteration occurs in this step; no volatile capture is required. Reference any evidence already collected in Steps 2–4 (network logs, extension storage artifacts, credential store snapshots) as supporting documentation for the threat register entry and risk assessment update.

**Step 6: Track vendor remediation — monitor patch release status for all six named products; assign ownership to follow up with each vendor and track against your vulnerability management process (CIS 7.1, CIS 7.2); brief leadership on organizational relevance with specific context on which products your organization has deployed.**

**NIST Phase:** Post Incident

**Reference:** NIST 800-61r3 §4 — Post-Incident Activity: Coordinating vendor remediation tracking and communicating organizational risk to leadership

**Controls:** CIS 7.1 (Establish and Maintain a Vulnerability Management Process), CIS 7.2 (Establish and Maintain a Remediation Process), CIS 2.2 (Ensure Authorized Software is Currently Supported)

**Compensating:** Create a vendor tracking matrix with six rows — one per product (ChatGPT Atlas, Perplexity Comet, Anthropic Claude browser extension, Fellou, Genspark, Sigma) — with columns for: advisory URL, patch release date, installed version in your environment, remediation owner, and target remediation date. Monitor each vendor's security advisory page and relevant CVE/NVD feeds (even though no CVE is currently assigned, check NVD for the product names). Set a calendar-based review cadence of weekly until patches are available, then validate patch deployment within the SLA defined in CIS 7.2.

**Evidence:** This step does not alter live state. Document the current installed versions of all six AI browser products identified in Step 1 as the pre-patch baseline, so patch application can be verified against a known-good starting state.

**Step 7: Monitor developments — LayerX's disclosure is the primary source; watch for independent technical replications, additional affected products, regulatory guidance on agentic AI security, and vendor-issued advisories or configuration hardening guidance.**

**NIST Phase:** Post Incident

**Reference:** NIST 800-61r3 §4 — Post-Incident Activity: Maintaining situational awareness through threat intelligence monitoring as the BioShocking disclosure matures

**Controls:** NIST AU-13 (Monitoring For Information Disclosure)

**Compensating:** Configure free RSS or ATOM feed monitors (e.g., Feedly free tier, rssdaemon) for LayerX Security blog, CISA advisories, and the security advisory pages for all six named vendors. Set Google Alerts or similar for the terms 'BioShocking prompt injection,' 'agentic browser credential theft,' and each product name paired with 'vulnerability' or 'advisory.' Subscribe to CISA's Known Exploited Vulnerabilities catalog RSS feed and the NVD CPE feed for 'openai,' 'anthropic,' 'perplexity,' 'fellou,' 'genspark,' and 'sigma' to catch any CVE assignment if the finding is formalized.

**Evidence:** This is a passive intelligence-gathering step and does not alter live system state. No volatile capture is required before execution. Retain and date-stamp all collected advisories, independent research publications, and vendor statements as part of the incident record to support post-incident review and future threat model updates.

## Detection Guidance

Detection for BioShocking-style indirect prompt injection targets the behavioral output of the AI agent rather than the injected content itself, since the malicious payload is embedded in rendered page text and will not match traditional signature patterns.

Log sources to prioritize: browser process network telemetry, DNS query logs, endpoint DLP logs capturing outbound data from browser processes, and credential manager access logs.

Behavioral patterns to hunt for:

- Outbound HTTP POST requests originating from AI browser processes to domains not matching the currently rendered site's origin (cross-origin exfiltration); correlate with NIST AU-6 (Audit Record Review, Analysis, and Reporting).
- AI browser processes accessing OS credential stores, password manager APIs, or browser-stored credential databases outside of user-initiated login events; alert on D3-LAM (Local Account Monitoring) signals.
- AI agent activity occurring on pages classified as entertainment or gaming content (consistent with the game-themed lure vector described by LayerX) while simultaneously triggering credential or form-fill access.
- Unexpected form submissions or DOM interactions on pages where no user interaction was recorded immediately prior.
- DNS queries to newly registered or low-reputation domains issued by browser processes during agentic task execution.

Policy gaps to audit:

- Determine whether your endpoint DLP policy scope includes AI browser processes; many DLP policies target traditional browsers and may not enumerate newer AI browser executables.
- Verify that agentic mode authorization follows least-privilege principles per NIST AC-6; agents should not have access to credentials beyond the specific task scope.
- Review whether NIST AC-4 (Information Flow Enforcement) controls are in place to restrict what data an AI browser agent can transmit externally.

Note: LayerX's research references specific indicators; consult the LayerX disclosure directly and the SecurityWeek and CyberScoop coverage for any published IOC values. No verifiable IOC values were present in the source material available for this item.

## Framework Mappings

### MITRE-ATTACK

- **T1189** — Drive-by Compromise
- **T1539** — Steal Web Session Cookie
- **T1059** — Command and Scripting Interpreter
- **T1185** — Browser Session Hijacking
- **T1557** — Adversary-in-the-Middle
- **T1190** — Exploit Public-Facing Application
- **T1056** — Input Capture

### NIST-800-53R5

- **CM-7** — Least Functionality
- **SI-3** — Malicious Code Protection
- **SI-4** — System Monitoring
- **SI-7** — Software, Firmware, and Information Integrity
- **CA-8** — Penetration Testing
- **RA-5** — Vulnerability Monitoring and Scanning
- **SC-7** — Boundary Protection
- **SI-2** — Flaw Remediation
- **IA-5** — Authenticator Management
- **SI-10** — Information Input Validation
- **SR-2** — Supply Chain Risk Management Plan

### OWASP-TOP10-2021

- **A04:2021** — Insecure Design
- **A07:2021** — Identification and Authentication Failures
- **A03:2021** — Injection

### CIS-V8

- **5.2** — Use Unique Passwords
- **16.10** — Apply Secure Design Principles in Application Architectures
- **6.3** — Require MFA for Externally-Exposed Applications
- **15.1** — Establish and Maintain an Inventory of Service Providers

**HIPAA-SECURITY**

- **164.308(a)(5)(ii)(D)** — Password Management
- **164.312(d)** — Person or Entity Authentication

**ISO-27001-2022**

- **A.8.26** — Application security requirements
- **A.5.21** — Managing information security in the ICT supply chain

**SOC2-TSC**

- **CC6.1** — Logical access security software, infrastructure, and architectures
- **CC9.2** — Manages risks associated with vendors and business partners

**NIST-CSF-2**

- **GV.SC-01** — Cybersecurity supply chain risk management program
- **DE.AE-08** — Incidents are declared when adverse events meet the defined incident criteria

**MITRE ATT&CK Mapping**

Technique ID	Technique Name	Tactic
T1189	Drive-by Compromise	Initial-Access
T1539	Steal Web Session Cookie	Credential-Access
T1059	Command and Scripting Interpreter	Execution
T1185	Browser Session Hijacking	Collection
T1557	Adversary-in-the-Middle	Credential-Access
T1190	Exploit Public-Facing Application	Initial-Access
T1056	Input Capture	Collection

**Sources**

Source	URL	Tier
Security News	<a href="https://thehackernews.com/2026/06/new-bioshocking-attack-tricks-ai...">https://thehackernews.com/2026/06/new-bioshocking-attack-tricks-ai...</a>	T2

Source	URL	Tier
<b>Pwning OpenAI Atlas Through Exposed Browser Internals</b>	<a href="https://www.hacktron.ai/blog/hacking-openai-atlas-browser/">https://www.hacktron.ai/blog/hacking-openai-atlas-browser/</a>	T3
<b>I Tried an AI Web Browser, and Now I'm a Convert - WSJ</b>	<a href="https://www.wsj.com/tech/personal-tech/ai-browsers-atlas-gemini-com...">https://www.wsj.com/tech/personal-tech/ai-browsers-atlas-gemini-com...</a>	T3
<b>AI Sidebar Spoofing Puts ChatGPT Atlas, Perplexity Comet and ...</b>	<a href="https://www.securityweek.com/ai-sidebar-spoofing-puts-chatgpt-atlas...">https://www.securityweek.com/ai-sidebar-spoofing-puts-chatgpt-atlas...</a>	T2
<b>Exclusive: OpenAI's Atlas browser — and others - CyberScoop</b>	<a href="https://cyberscoop.com/openai-atlas-splx-research-cloaking-attacks-...">https://cyberscoop.com/openai-atlas-splx-research-cloaking-attacks-...</a>	T2

**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 2026-06-30 14:47 UTC by TJS Security Command Center