

INTELLIGENCE BRIEFING

Security Command Center

TLP:CLEAR

2026-06-14 05:02 UTC

AI Agents Vulnerable to Prompt Injection and Social Engineering Credential Exfiltration

SECURITY ANALYSIS | HIGH | CVSS 8.1

SCC Item ID	SCC-STY-2026-0197
Type	Security Analysis
Severity	HIGH
CVSS Base Score	8.1
Affected Products	AI-powered agentic systems including Anthropic Claude Code GitHub Action and OpenClaw email agent; broadly applicable to any LLM-based autonomous agent with access to secrets or credentials
Published	2026-06-11
Discovery Source	Gemini

Executive Summary

Red team research and documented attack scenarios confirm that AI agents operating inside CI/CD pipelines and email workflows can be hijacked through prompt injection, causing the agent to exfiltrate credentials including AWS IAM keys, database passwords, and SSH tokens to attacker infrastructure. A Microsoft Security Blog analysis published June 5, 2026 specifically documents this vector against the Claude Code GitHub Action, and the Cloud Security Alliance frames it as an emerging supply chain threat. This signals a structural gap in enterprise AI adoption: autonomous agents inherit the trust and access of the developers who deploy them, but lack the discrimination to reject adversarially crafted instructions embedded in data they process.

Technical Analysis

The attack class exploits a fundamental design tension in large language model agents: the model receives both instructions and data from its environment, but has no cryptographically enforced boundary between the two. An adversary who can influence content the agent reads, whether a repository file, an incoming email, a pull request comment, or an external API response, can embed natural-language instructions that the agent treats as legitimate task directives.

The Microsoft Security Blog's June 2026 analysis of the Claude Code GitHub Action illustrates the CI/CD variant. The agent operates with repository-level access and executes code, manages secrets, and interacts with downstream services as part of normal pipeline function. If an attacker can commit or influence content in a

repository the agent processes, injected instructions can redirect the agent to read CI/CD secrets from environment variables and transmit them to an attacker-controlled endpoint. This aligns with MITRE ATT&CK T1552.001 (Credentials in Files) and T1195.002 (Compromise Software Supply Chain), as the agent itself becomes the exfiltration mechanism rather than a separate malware implant.

Email agent variants documented in community discussions follow a social engineering path more analogous to T1566 (Phishing). A crafted email instructs the agent to forward credentials or take actions on behalf of the user, exploiting the agent's delegated access without the user's knowledge. The agent's helpfulness, the design quality that makes it useful, becomes the attack surface.

CWE-77 (Command Injection) and CWE-94 (Code Injection) capture the structural flaw, though neither fully characterizes the LLM-specific variant: the injected payload is natural language, not executable syntax, and detection with traditional input validation is not straightforward. CWE-522 (Insufficiently Protected Credentials) and CWE-200 (Exposure of Sensitive Information) describe the downstream consequences.

The defensive gap is compounded by the agent's autonomous action-taking capability. Unlike a passive data-processing system, an agent with tool access can reach out to external infrastructure, making exfiltration a single-step operation once the injection succeeds. Organizations that have granted agents broad secret access, justified by developer productivity gains, are exposed in proportion to that access scope.

Action Checklist

- 1. Step 1:** Assess exposure, audit every LLM-based agent deployed in your environment, including CI/CD integrations (GitHub Actions, GitLab CI, Jenkins plugins), email automation, and coding assistants. Inventory which agents have access to secrets, API keys, or credentials. Start with Claude Code GitHub Action and any similar agent deployments.
- 2. Step 2:** Review controls, apply least-privilege secret access immediately: agents should receive scoped, short-lived credentials rather than long-lived IAM keys or static passwords (NIST AC-6, Least Privilege). Enforce secret isolation so agents can request credentials only for their specific task scope, not read environment-wide secrets (NIST AC-3, Access Enforcement). Enable MFA on all accounts agents act on behalf of where the provider supports it (CIS 6.3, CIS 6.5; D3-MFA).
- 3. Step 3:** Update threat model, add prompt injection via repository content and socially engineered agent instruction as explicit attack paths in your threat register. Map to T1552.001 (Credentials in Files), T1195.002 (Supply Chain Compromise), T1566 (Phishing to agent), and T1059 (agent-executed command injection). Flag CI/CD pipelines and email agents as high-value pivot points.
- 4. Step 4:** Implement agent-specific controls, require human approval gates for any agent action involving credential access, external network calls, or data exfiltration paths (NIST AC-5, Separation of Duties). Apply output filtering to agent responses to detect credential patterns before transmission. Rotate any credentials that have been in scope of deployed agents without access controls in place (D3-CRO, Credential Rotation; D3-CH, Credential Hardening).
- 5. Step 5:** Communicate findings, brief engineering leads and DevOps teams on the specific risk: agents with CI/CD secret access are the highest-priority exposure. Escalate to CISO if agents currently hold broad AWS IAM permissions or production database credentials without scoping. Frame for leadership as a supply chain risk, not an abstract AI safety issue.
- 6. Step 6:** Monitor developments, track Microsoft Security Blog, Cloud Security Alliance Labs, and MITRE ATLAS for follow-up research on prompt injection mitigations. Watch for vendor guidance from Anthropic on Claude Code GitHub Action hardening. Review CISA advisories for any formal guidance on agentic AI

security.

Detection Guidance

Detection for prompt injection in agentic systems requires behavioral monitoring rather than signature matching, because the injected payload is natural language and the exfiltration mechanism is the agent's own legitimate tool access.

Log sources to instrument (per NIST AU-2, Event Logging and AU-12, Audit Record Generation):

- CI/CD pipeline logs: flag any agent-initiated outbound network connections to destinations outside the expected service list, particularly during code review or repository scan phases. Outbound HTTPS to non-whitelisted domains from a GitHub Actions runner should trigger immediate review.
- Secret manager access logs (AWS Secrets Manager, HashiCorp Vault, GitHub Actions encrypted secrets): alert on any agent reading multiple secrets in a single execution, or reading secrets not required for the declared task (NIST AU-6, Audit Record Review).
- Email agent send logs: flag outbound messages containing base64-encoded strings, credential-pattern substrings (regex for AWS key prefixes such as AKIA, private key headers, connection string formats), or messages to external addresses initiated without a human-authored prompt in the same session.
- Environment variable access: in containerized pipelines, log any process reading CI/CD environment variables and compare against an allowlist of expected secret names for that job.

Behavioral hunts to run:

- Hunt for agent sessions that read a secret and then made an outbound connection within the same execution context (T1552.001 followed by T1041 pattern).
- Hunt for repository files (README, markdown, config stubs, PR descriptions) containing instruction-pattern language directed at agents: phrases such as "ignore previous instructions", "your new task is", "send the contents of", or explicit references to environment variable names.
- Hunt for email agent activity initiated outside business hours or from sender addresses not in organizational directories.

Policy gaps to audit:

- Confirm agents do not have standing access to production secrets; access should be just-in-time and task-scoped (NIST AC-6).
- Verify outbound network allowlists are enforced at the runner or container level, not just declared in configuration (CIS 4.4, Firewall on Servers).
- Confirm audit logs for agent actions are immutable and retained per policy (NIST AU-9, Protection of Audit Information; CIS 8.2, Collect Audit Logs).

Indicators of Compromise

Type	Value	Context	Confidence
TOOL	Claude Code GitHub Action (Anthropic)	LLM-based CI/CD agent leveraged via prompt injection in repository content to read pipeline secrets and potentially exfiltrate credentials to attacker-controlled infrastructure	MEDIUM
TOOL	OpenClaw email agent	LLM-based email agent cited in community discussions as susceptible to socially engineered instructions that direct the agent to take unauthorized actions or forward sensitive data on behalf of the user	LOW
URL	Pending – refer to Microsoft Security Blog (2026-06-05) for any published indicators	Microsoft's CI/CD prompt injection analysis may contain specific repository content patterns, injected instruction signatures, or infrastructure indicators; source URL could not be actively confirmed as resolving at time of writing	LOW

Framework Mappings

MITRE-ATTACK

- **T1552.001** — Credentials In Files
- **T1190** — Exploit Public-Facing Application
- **T1566** — Phishing
- **T1195.002** — Compromise Software Supply Chain
- **T1059** — Command and Scripting Interpreter
- **T1552** — Unsecured Credentials

NIST-800-53R5

- **CA-8** — Penetration Testing
- **RA-5** — Vulnerability Monitoring and Scanning
- **SC-7** — Boundary Protection
- **SI-2** — Flaw Remediation
- **SI-7** — Software, Firmware, and Information Integrity
- **AT-2** — Literacy Training and Awareness
- **CA-7** — Continuous Monitoring
- **SI-3** — Malicious Code Protection
- **SI-4** — System Monitoring
- **SI-8** — Spam Protection
- **CM-7** — Least Functionality
- **SA-9** — External System Services

- **SR-3** — Supply Chain Controls and Processes
- **SI-10** — Information Input Validation
- **IA-5** — Authenticator Management
- **AC-3** — Access Enforcement
- **SC-28** — Protection of Information at Rest
- **SR-2** — Supply Chain Risk Management Plan

OWASP-TOP10-2021

- **A03:2021** — Injection
- **A04:2021** — Insecure Design
- **A07:2021** — Identification and Authentication Failures
- **A01:2021** — Broken Access Control

CIS-V8

- **16.10** — Apply Secure Design Principles in Application Architectures
- **5.2** — Use Unique Passwords
- **6.3** — Require MFA for Externally-Exposed Applications
- **14.2** — Train Workforce Members to Recognize Social Engineering Attacks
- **15.1** — Establish and Maintain an Inventory of Service Providers

HIPAA-SECURITY

- **164.308(a)(5)(ii)(D)** — Password Management
- **164.312(a)(1)** — Access Control
- **164.312(d)** — Person or Entity Authentication

SOC2-TSC

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

ISO-27001-2022

- **A.8.8** — Management of technical vulnerabilities
- **A.5.21** — Managing information security in the ICT supply chain
- **A.5.23** — Information security for use of cloud services

NIST-CSF-2

- **GV.SC-01** — Cybersecurity supply chain risk management program

MITRE ATT&CK Mapping

Technique ID	Technique Name	Tactic
T1552.001	Credentials In Files	Credential-Access

Technique ID	Technique Name	Tactic
T1190	Exploit Public-Facing Application	Initial-Access
T1566	Phishing	Initial-Access
T1195.002	Compromise Software Supply Chain	Initial-Access
T1059	Command and Scripting Interpreter	Execution
T1552	Unsecured Credentials	Credential-Access

Sources

Source	URL	Tier
Securing CI/CD in an agentic world: Claude Code Github action case	https://www.microsoft.com/en-us/security/blog/2026/06/05/securing-c...	T1
OpenClaw agent attacks developer who rejected its code - Facebook	https://www.facebook.com/groups/2600net/posts/4465410557015368/	T3
awesome-openclaw-skills/categories/coding-agents-and-ides.md at ...	https://github.com/VoltAgent/awesome-openclaw-skills/blob/main/cate...	T3
AI Agent Prompt Injection: The New CI/CD Supply Chain Threat	https://labs.cloudsecurityalliance.org/research/csa-research-note-c...	T3
OpenClaw vs Claude Code vs Copilot: 3 AI Agents, 5 Real Tasks, 1 ...	https://findskill.ai/blog/openclaw-vs-claude-code-vs-copilot/	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 2026-06-14 05:02 UTC by TJS Security Command Center