

INTELLIGENCE BRIEFING
Security Command Center

TLP:CLEAR
2026-05-08 19:04 UTC

TrustFall Attack: AI Coding Agents Exploitable for Stealthy Supply Chain Compromise

SECURITY ANALYSIS | HIGH | CVSS 8.0

SCC Item ID	SCC-STY-2026-0117
Type	Security Analysis
Severity	HIGH
CVSS Base Score	8.0
Affected Products	AI coding agents (Claude Code specifically named); broader applicability to agentic coding tools with autonomous repository access
Published	2026-05-07
Discovery Source	Gemini

Executive Summary

Researchers at Adversa.AI have disclosed a novel attack class called 'TrustFall' that weaponizes AI coding agents, such as Claude Code, to silently introduce malicious code into developer environments through manipulated repository content. Because these agents operate with full developer privileges and autonomously install dependencies, a single compromised package or poisoned repository can propagate malicious code through CI/CD pipelines and into released software products. This disclosure highlights a structural gap in supply chain risk: the attack surface now extends to the reasoning and trust assumptions of AI agents embedded in development workflows.

Technical Analysis

TrustFall, disclosed by Adversa.AI, describes an architectural attack class against agentic coding tools that exploits a fundamental design assumption: AI coding agents treat repository content as authoritative, trusted context. This trust model, appropriate for benign workflows, becomes an attack vector when an adversary controls or compromises a repository the agent consults.

The attack combines two well-understood supply chain techniques with a novel AI-specific layer. First, indirect prompt injection: malicious instructions embedded in repository files, package manifests, README content, or inline code comments direct the agent to perform actions the developer never requested. Second, dependency confusion or malicious package substitution: the agent, trusting the injected instructions, fetches and installs attacker-controlled packages or executes attacker-supplied scripts. Because the agent operates under the developer's system privileges, the resulting code execution is not sandboxed. The attacker achieves RCE in the

developer's environment without exploiting a binary vulnerability.

MITRE ATT&CK mapping is instructive here. The campaign pattern tracks against T1195.001 (Compromise Software Supply Chain), T1566 (Phishing as an initial delivery mechanism for repository poisoning), T1059 (Command and Scripting Interpreter, via agent-executed scripts), and T1204.002 (Malicious File execution triggered by agent tool use). CWE mapping reinforces the architectural framing: CWE-829 (Inclusion of Functionality from Untrusted Control Sphere), CWE-494 (Download of Code Without Integrity Check), CWE-20 (Improper Input Validation), and CWE-1188 (Insecure Default Initialization) all apply, reflecting that the vulnerability is in the agent's trust architecture, not a discrete patch-addressable flaw.

The supply chain amplification risk is significant. A developer using an AI coding agent as part of a CI/CD workflow may unknowingly commit attacker-influenced code to a shared repository, package registry, or build artifact. Downstream consumers of that software inherit the compromise without any direct interaction with the AI agent. This mirrors the structural impact of the SolarWinds and XZ Utils incidents but removes the requirement for a sophisticated human attacker to maintain persistent access: the agent does the work autonomously.

SecurityWeek and Wiz have both noted the broader industry implications, characterizing this as a category-level risk for organizations adopting agentic development tooling. Claude Code is cited in secondary source reporting (Xcitiem ThreatLabs coverage, SecurityWeek analysis); full technical details from the primary Adversa.AI disclosure were not directly verified for this analysis (confidence in attack class mechanics: HIGH; confidence in Claude Code as specifically targeted: MEDIUM). The attack class applies to any agentic coding tool with autonomous repository access and privilege to execute code.

No CVE has been assigned. The absence of a CVE is itself analytically significant: it reflects that traditional vulnerability management frameworks are not equipped to track behavioral and architectural risks in AI systems. Organizations relying on CVE-based patching workflows will not receive a prompt to act on this class of threat.

Action Checklist

1. Step 1: Assess exposure, inventory all AI coding agents deployed in your development environment, including Claude Code, GitHub Copilot Workspace, Cursor, Devin, and similar agentic tools with autonomous repository access and code execution capability
2. Step 2: Review controls, audit the privilege levels under which AI coding agents operate; enforce least-privilege principles by restricting agent access to only the directories, registries, and execution contexts required; verify that agents cannot install packages or execute scripts without human approval gates
3. Step 3: Implement input validation at the agent boundary, treat all repository content, package manifests, and inline code comments as potentially adversarial input; evaluate whether your coding agent applies integrity checks (e.g., hash verification, provenance attestation) before executing fetched content
4. Step 4: Update threat model, add indirect prompt injection via repository content as a recognized attack vector in your software supply chain threat register; map to T1195.001 and T1059 in your MITRE ATT&CK-based threat model
5. Step 5: Apply supply chain controls to agent-generated artifacts, require code review, static analysis, and dependency scanning (e.g., SBOM generation, Sigstore-based provenance) for any code produced or modified by AI coding agents before merge or deployment

6. Step 6: Communicate findings, brief engineering leadership and DevSecOps teams on TrustFall mechanics with specific emphasis on CI/CD pipeline risk; escalate to CISO level if agentic coding tools have production pipeline access

7. Step 7: Implement interim controls immediately: privilege restriction, mandatory code review for agent-generated changes, and SBOM generation are your primary mitigation paths. Monitor Adversa.AI for full technical paper publication, Anthropic for Claude Code security advisories, and other agentic tool vendors for responses; note that architectural risks do not follow patch cadences, so control implementation cannot wait for vendor patches

IR / Forensic Enrichment

Triage Priority	URGENT
Escalation Criteria	Escalate immediately to CISO and legal counsel if any AI coding agent (Claude Code, Copilot Workspace, Cursor, or Devin) has write access to CI/CD pipelines with production deployment targets, if SBOM diffing or git log analysis reveals unauthorized dependency additions during agent sessions, or if the organization operates in a regulated sector (SOC 2, PCI DSS, FedRAMP) where supply chain compromise triggers breach notification or control failure reporting obligations.
Recovery Notes	After containment controls are in place, re-run SBOM generation and 'pip-audit' / 'npm audit' against all projects where AI coding agents operated during the exposure window and compare against pre-agent-introduction baselines to confirm no malicious dependencies persist in the dependency tree. Monitor Sysmon Event ID 1 and 11 (Process Create and File Create) for agent process trees on developer workstations for a minimum of 30 days post-containment, watching specifically for child processes (cmd.exe, powershell.exe, bash, curl, wget) spawned by AI agent executables that were not explicitly invoked by a human. Re-validate CI/CD pipeline YAML files and any IaC templates (Terraform, CloudFormation) that AI agents touched, as TrustFall-style attacks may have introduced persistent backdoors in infrastructure definitions rather than application code.

Forensic Artifacts

Git audit logs for repositories accessed by AI coding agents during the exposure window — specifically 'git log --all --oneline --author' filtered to agent-associated identities and automated commit signatures, which would reveal unauthorized commits or dependency manifest changes introduced via TrustFall-poisoned repository content | CI/CD pipeline execution logs (GitHub Actions logs, GitLab CI job traces, or Jenkins build logs) for any pipeline run that included an AI agent step — look for unexpected 'pip install', 'npm install', 'curl', or 'wget' commands executed during agent-controlled steps, which indicate successful payload execution from a poisoned package manifest or inline comment injection | Package manager installation histories on developer workstations — '~/.local/share/pip/pip.log' (Linux pip), '%APPDATA%\pip\pip.log' (Windows pip), and npm debug log ('~/.npm/_logs/') — capturing any packages installed during AI agent sessions that were not explicitly requested by the developer, the primary artifact class for a TrustFall dependency injection attack | Claude Code session logs and tool_use call histories (stored in ~/.claude/logs/ or equivalent agent working directories) — these record every tool invocation the agent made, including file reads, shell executions, and package installs, and would show whether the agent acted on injected instructions from repository content such as README.md or inline code comments | Sysmon Event ID 1 (Process Create) logs from developer workstations filtered on parent processes matching AI agent executables (claude, cursor, devin-agent, node processes associated with Copilot extensions) — child processes such as 'pip', 'npm', 'bash', 'powershell', or 'curl' spawned unexpectedly from these parents are the primary host-level indicator of a successful TrustFall indirect prompt injection payload execution

Per-Action IR Details

Step 1: Assess exposure — inventory all AI coding agents deployed in your development environment, including Claude Code, GitHub Copilot Workspace, Cursor, Devin, and similar agentic tools with autonomous repository access and code execution capability

NIST Phase: Preparation

Reference: NIST 800-61r3 §2 — Preparation: Establishing IR Capability and Asset Visibility

Controls: NIST IR-4 (Incident Handling), NIST SI-5 (Security Alerts, Advisories, and Directives), CIS 1.1 (Establish and Maintain Detailed Enterprise Asset Inventory), CIS 2.1 (Establish and Maintain a Software Inventory)

Compensating: Run 'pip list', 'npm list -g', and 'which claude' / 'which cursor' across developer workstations using a simple bash or PowerShell sweep script. Query GitHub org settings via the GitHub REST API (GET /orgs/{org}/installations) to enumerate installed GitHub Apps with repository access, including Copilot Workspace. Document each agent's granted OAuth scopes and filesystem permissions in a shared spreadsheet — flag any agent with 'write', 'admin', or unrestricted path access.

Evidence: Before beginning inventory, snapshot current GitHub App installation manifests (found at GitHub org > Settings > GitHub Apps), Claude Code configuration files (typically ~/.claude/config.json or project-level .claude/ directories), Cursor agent settings, and any .env or agent-credential files in developer home directories. These establish the pre-assessment baseline and may reveal unauthorized or shadow agent installations introduced via a TrustFall-style poisoned repository.

Step 2: Review controls — audit the privilege levels under which AI coding agents operate; enforce least-privilege principles by restricting agent access to only the directories, registries, and execution contexts required; verify that agents cannot install packages or execute scripts without human approval gates

NIST Phase: Containment

Reference: NIST 800-61r3 §3.3 — Containment Strategy: Limiting Blast Radius During Active Risk

Controls: NIST IR-4 (Incident Handling), NIST SI-7 (Software, Firmware, and Information Integrity), CIS 5.4 (Restrict Administrator Privileges to Dedicated Administrator Accounts), CIS 4.6 (Securely Manage Enterprise Assets and Software)

Compensating: On Linux/macOS developer hosts, use 'sudo -l' to enumerate what the process account running Claude Code or Cursor can execute. Implement an npm/pip pre-install hook via .npmrc ('ignore-scripts=true') and pip's '--no-deps' flag to block silent dependency installation. For Claude Code specifically, review the tool_use permissions in the agent's system prompt or configuration — restrict to read-only filesystem tools where possible. Use AppArmor (Linux) or macOS TCC profiles to constrain which directories the agent process can write to.

Evidence: Capture current agent process privilege context before making changes: run 'ps aux | grep -E "claude|cursor|devin|copilot"' and record the running UID/GID. Export current npm global package list ('npm list -g --depth=0') and pip freeze output per developer environment as a baseline. On macOS, export TCC database ('tccutil list') to document which applications have been granted broad filesystem access. These snapshots will identify privilege creep introduced by a TrustFall attack before controls are tightened.

Step 3: Implement input validation at the agent boundary — treat all repository content, package manifests, and inline code comments as potentially adversarial input; evaluate whether your coding agent applies integrity checks (e.g., hash verification, provenance attestation) before executing fetched content

NIST Phase: Containment

Reference: NIST 800-61r3 §3.3 — Containment Strategy: Preventing Lateral Spread via Adversarial Input

Controls: NIST SI-10 (Information Input Validation), NIST SI-7 (Software, Firmware, and Information Integrity), NIST SI-3 (Malicious Code Protection), CIS 7.1 (Establish and Maintain a Vulnerability Management Process)

Compensating: Deploy YARA rules targeting common indirect prompt injection patterns in repository files — scan README.md, CONTRIBUTING.md, pyproject.toml, package.json, and inline code comments for embedded instruction strings such as 'ignore previous instructions', 'as an AI', or base64-encoded payloads using a rule triggered by 'yara -r injection_rules.yar /path/to/repo'. Use 'sha256sum' or 'cosign verify' (free, from Sigstore) to verify package manifest integrity against known-good hashes before allowing agent consumption. For npm packages, enforce 'npm audit' and check 'npm pack --dry-run' output for unexpected files.

Evidence: Before implementing controls, extract and preserve all current package.json, package-lock.json, requirements.txt, pyproject.toml, and Pipfile.lock files from repositories where AI agents have recently operated. Hash each file ('sha256sum requirements.txt > requirements.txt.sha256') and store offline. Capture git log ('git log --all --oneline --graph') for any repositories the agent accessed to detect anomalous commits, especially commits made under the developer's credential by the agent without explicit human authorship.

Step 4: Update threat model — add indirect prompt injection via repository content as a recognized attack vector in your software supply chain threat register; map to T1195.001 and T1059 in your MITRE ATT&CK-based threat model

NIST Phase: Post Incident

Reference: NIST 800-61r3 §4 — Post-Incident Activity: Lessons Learned and Threat Model Updates

Controls: NIST IR-8 (Incident Response Plan), NIST RA-3 (Risk Assessment), NIST SI-5 (Security Alerts, Advisories, and Directives), CIS 7.1 (Establish and Maintain a Vulnerability Management Process), CIS 7.2 (Establish and Maintain a Remediation Process)

Compensating: Document the TrustFall attack chain in your existing threat register using a free MITRE ATT&CK Navigator layer file (downloadable at attack.mitre.org/resources/attack-navigator) — create a new layer mapping T1195.001 (Supply Chain Compromise: Compromise Software Dependencies and Development Tools) and T1059 (Command and Scripting Interpreter) with annotations specific to AI agent exploitation. If no formal threat register exists, create a markdown-based register in your internal wiki documenting the attack vector, affected agents, and mitigations. Link to the Adversa.AI TrustFall disclosure as the source reference.

Evidence: Preserve the current state of your threat model and supply chain risk register before updating — export or screenshot the existing document to establish a pre-disclosure baseline. If MITRE ATT&CK Navigator layers exist, export them as JSON. This preserves the 'before' state for audit purposes and demonstrates due diligence in responding to the Adversa.AI disclosure (NIST 800-61r3 §4 post-incident documentation requirement).

Step 5: Apply supply chain controls to agent-generated artifacts — require code review, static analysis, and dependency scanning (e.g., SBOM generation, Sigstore-based provenance) for any code produced or

modified by AI coding agents before merge or deployment

NIST Phase: Eradication

Reference: NIST 800-61r3 §3.4 — Eradication: Removing Threat Artifacts and Verifying Clean State

Controls: NIST SI-2 (Flaw Remediation), NIST SI-7 (Software, Firmware, and Information Integrity), NIST SA-12 (Supply Chain Protection), CIS 7.3 (Perform Automated Operating System Patch Management), CIS 7.4 (Perform Automated Application Patch Management)

Compensating: Use 'syft' (free, Anchore) to generate SBOMs for any build artifacts produced by AI agent sessions — run 'syft packages dir:/path/to/project -o spdx-json > sbom.json' and diff against the last known-clean SBOM to detect new or modified dependencies. Run Semgrep (free OSS tier) with the 'p/supply-chain' ruleset against all agent-generated code diffs before merge ('semgrep --config p/supply-chain .'). Require all PRs originating from or modified by an AI agent session to carry a mandatory human reviewer sign-off in GitHub branch protection rules — enforce via 'Required reviewers: 2' with 'Dismiss stale reviews' enabled.

Evidence: Before eradication, extract and preserve git diff output for all commits made during AI agent sessions ('git log --author="[agent identity]" --since="[exposure window start]" -p > agent_commits.patch'). Run 'pip-audit' or 'npm audit --json' against the current environment and save output. Pull the full CI/CD pipeline execution logs from GitHub Actions, GitLab CI, or Jenkins for any pipeline runs that included AI agent steps — these logs will show whether a TrustFall payload reached the build or deploy stage.

Step 6: Communicate findings — brief engineering leadership and DevSecOps teams on TrustFall mechanics with specific emphasis on CI/CD pipeline risk; escalate to CISO level if agentic coding tools have production pipeline access

NIST Phase: Containment

Reference: NIST 800-61r3 §3.3 — Containment: Coordination and Communication During Active Risk Management

Controls: NIST IR-6 (Incident Reporting), NIST IR-7 (Incident Response Assistance), NIST IR-4 (Incident Handling), CIS 8.2 (Collect Audit Logs)

Compensating: Draft a one-page internal advisory using the Adversa.AI TrustFall disclosure as the source, specifically calling out which agents (Claude Code, Copilot Workspace, Cursor, Devin) are deployed in your environment and whether any have CI/CD pipeline write access. Use your existing ticketing system (Jira, GitHub Issues, or even email with read receipts) to create a tracked communication record — NIST IR-6 requires documented reporting, and this thread serves as evidence. If Claude Code or any agent has production deployment access, this is an immediate CISO escalation regardless of whether active exploitation has been detected.

Evidence: Before the briefing, compile a communication package: export the list of engineers who have installed or used AI coding agents (from GitHub audit logs at /orgs/{org}/audit-log?phrase=copilot or from developer machine surveys), document which CI/CD pipelines have AI agent integrations (export pipeline YAML files from .github/workflows/, .gitlab-ci.yml, or Jenkinsfile), and note which pipelines have production deployment targets. This package forms the impact scope statement for leadership and satisfies NIST IR-6 incident reporting documentation requirements.

Step 7: Monitor developments — track Adversa.AI for full technical paper publication, monitor Anthropic's security advisories for Claude Code, and watch for vendor responses from other agentic coding tool providers; no patch cadence exists for architectural risks of this type, so policy and tooling controls are the primary mitigation path

NIST Phase: Post Incident

Reference: NIST 800-61r3 §4 — Post-Incident Activity: Continuous Improvement and Intelligence Integration

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

Compensating: Set up free RSS or email alert monitoring for Adversa.AI (adversa.ai/blog/), Anthropic's security page (anthropic.com/security), and GitHub Security Advisories filtered for 'copilot' and 'cursor' using GitHub's advisory API ('https://api.github.com/advisories?type=reviewed&keyword=copilot'). Create a recurring 30-day calendar review to check for full TrustFall paper publication and vendor responses. Add a Sigma rule to your log pipeline (or manual

review checklist if no SIEM) to detect anomalous package installations or script executions originating from AI agent process trees — use Sysmon Event ID 1 (Process Create) filtered on parent process matching known agent executables (claude, cursor, devin-agent).

Evidence: Establish a monitoring baseline by documenting the current versions of all deployed AI coding agents (Claude Code version from 'claude --version', Cursor version from app settings, Copilot extension version from IDE plugin manager) and the date of this assessment. Archive the current Adversa.AI TrustFall disclosure page and Anthropic's security advisory page as PDF snapshots — these serve as the threat intelligence baseline against which future vendor updates and patches will be compared for NIST IR-5 incident monitoring documentation.

Detection Guidance

Standard IOC-based detection is limited for TrustFall because the attack exploits agent behavior rather than a discrete malware artifact. Hunting strategy should focus on anomalous agent activity and supply chain integrity signals.

Log sources to prioritize:

- Developer workstation process logs: Look for unexpected child processes spawned by AI coding agent processes (e.g., node, python, shell interpreters launched by agent runtimes without corresponding developer terminal sessions)
- Package manager logs (npm, pip, cargo, etc.): Flag packages installed by agent processes that were not present in project lock files before the agent session began
- Network logs: Identify outbound connections to unfamiliar domains or IPs initiated by coding agent processes, particularly to package registries not on your approved list
- CI/CD pipeline logs: Flag unexpected dependency additions, script modifications, or configuration file changes committed in agent-assisted branches
- Git commit metadata: Hunt for commits where the authoring context is an AI agent session and the diff includes new external dependencies or executable scripts not present in the preceding commit

Behavioral patterns to hunt:

- An AI coding agent process making network calls to domains not associated with the project's declared dependency sources
- Package installation events occurring outside developer-initiated workflows (e.g., during automated agent tasks on unattended workstations)
- Modifications to .npmrc, pip.conf, or equivalent package manager configuration files during agent sessions (a classic dependency confusion setup step)
- README or inline comment content containing instruction-like natural language directed at AI systems (a signature of injected prompt content)

Policy gaps to audit:

- Do agents run with developer-level credentials that include write access to production systems or shared registries?
- Are agent-generated code changes subject to the same peer review and SAST/SCA pipeline controls as human-authored code?

- Does your software composition analysis tooling scan agent-installed packages before they enter the build artifact?

Indicators of Compromise

Type	Value	Context	Confidence
TOOL	Pending – refer to Adversa.AI full technical paper for published indicators	Adversa.AI disclosed the TrustFall attack class and is expected to publish full technical details including any proof-of-concept artifacts, injected prompt patterns, or malicious package indicators; values were not available in source material accessible for this analysis	LOW

Framework Mappings

MITRE-ATTACK

- **T1195.001** — Compromise Software Dependencies and Development Tools
- **T1566** — Phishing
- **T1059** — Command and Scripting Interpreter
- **T1204.002** — Malicious File

NIST-800-53R5

- **AT-2** — Literacy Training and Awareness
- **CA-7** — Continuous Monitoring
- **SC-7** — Boundary Protection
- **SI-3** — Malicious Code Protection
- **SI-4** — System Monitoring
- **SI-8** — Spam Protection
- **CM-7** — Least Functionality
- **SI-7** — Software, Firmware, and Information Integrity
- **CM-3** — Configuration Change Control
- **SI-10** — Information Input Validation
- **SR-2** — Supply Chain Risk Management Plan

OWASP-TOP10-2021

- **A08:2021** — Software and Data Integrity Failures
- **A03:2021** — Injection

CIS-V8

- **2.5** — Allowlist Authorized Software

- **2.6** — Allowlist Authorized Libraries
- **16.10** — Apply Secure Design Principles in Application Architectures
- **14.2** — Train Workforce Members to Recognize Social Engineering Attacks
- **15.1** — Establish and Maintain an Inventory of Service Providers

ISO-27001-2022

- **A.8.26** — Application security requirements
- **A.8.8** — Management of technical vulnerabilities
- **A.5.21** — Managing information security in the ICT supply chain

NIST-CSF-2

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

SOC2-TSC

- **CC9.2** — Manages risks associated with vendors and business partners

MITRE ATT&CK Mapping

Technique ID	Technique Name	Tactic
T1195.001	Compromise Software Dependencies and Development Tools	Initial-Access
T1566	Phishing	Initial-Access
T1059	Command and Scripting Interpreter	Execution
T1204.002	Malicious File	Execution

Sources

Source	URL	Tier
gemini	https://www.securityweek.com/ai-coding-agents-could-fuel-next-suppl...	T3
Closing the Security Gap in the Age of Agentic Coding Wiz Blog	https://www.wiz.io/blog/securing-software-age-of-agentic-coding	T3
AI coding agents are secured in the wrong direction. : r/ClaudeCode	https://www.reddit.com/r/ClaudeCode/comments/1sc8932/ai_coding_agen...	T3
From AI to RCE: The Security Risks Lurking Inside Claude Code	https://threatlabsnews.xcitiium.com/blog/from-ai-to-rce-the-security...	T3

Source	URL	Tier
Claude's impact on security research Oren Yomtov posted on the ...	https://www.linkedin.com/posts/orenyomtov_i-found-a-vulnerability-i...	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 2026-05-08 19:04 UTC by TJS Security Command Center