

INTELLIGENCE BRIEFING
Security Command Center

TLP:CLEAR
2026-05-11 18:50 UTC

Autonomous AI Agents Introduce Ungoverned Identity and Action Risks Across Enterprise Environments

GOVERNANCE | HIGH

SCC Item ID	SCC-GOV-2026-0033
Type	Governance
Severity	HIGH
Affected Products	Organizations deploying autonomous AI agents across enterprise environments (sector-agnostic)
Published	2026-05-11
Discovery Source	Gemini

Executive Summary

Autonomous AI agents deployed across enterprise environments are accumulating credentials, executing multi-step actions, and accessing sensitive data outside existing identity and access management controls. Most organizations lack the governance structures, least-privilege enforcement, and behavioral monitoring needed to manage these systems as identity principals. The business risk spans unauthorized data access, undetected exfiltration, and eliminated incident response capability when agentic actions cannot be audited or attributed.

Technical Analysis

Autonomous AI agents present a compound governance failure across six risk dimensions, each mapped to established frameworks. Non-human identity sprawl: agents acquire API keys, tokens, and credentials dynamically outside formal IAM workflows, creating ungoverned privilege accumulation (CWE-522: Insufficiently Protected Credentials; CWE-250: Execution with Unnecessary Privileges; MITRE T1552: Unsecured Credentials; T1078: Valid Accounts). Prompt injection: agents that consume external content as task input are susceptible to adversarially crafted instructions that redirect agent behavior, analogous to CWE-20 (Improper Input Validation) applied to LLM instruction pipelines. Excessive agency: agents granted broad permissions to complete tasks may execute destructive or exfiltrative actions with no human checkpoint, per OWASP LLM Top 10 LLM08 (Excessive Agency); mapped to MITRE T1059 (Command and Scripting Interpreter) for agents with shell or API execution capability. Opaque audit trails: multi-step agentic actions across distributed systems rarely produce coherent, attributable logs (CWE-778: Insufficient Logging), degrading IR and forensic capacity. Supply chain exposure: agents frequently invoke third-party plugins, tools, and MCP (Model Context Protocol)

servers representing unvetted dependencies (MITRE T1195: Supply Chain Compromise). Data boundary violations: agents with broad data access may inadvertently or adversarially exfiltrate across trust boundaries (CWE-285: Improper Authorization; MITRE T1530: Data from Cloud Storage). This is an architectural governance gap, not a discrete vulnerability, and no CVE or CVSS score applies. Current guidance aligns with OWASP LLM Top 10 and MITRE ATT&CK; no primary NIST or CISA advisory was available at time of publication.

Action Checklist

- 1. Step 1, Inventory:** Conduct an immediate discovery sweep to enumerate all AI agents, service accounts, API keys, and tokens provisioned for agentic workloads; identify which were created outside formal IAM processes and which hold persistent or overprivileged credentials.
- 2. Step 2, Detection:** Query your SIEM and cloud access logs for service account activity patterns inconsistent with human-initiated sessions: off-hours API calls, bulk data access across storage buckets (T1530), lateral credential use (T1078), and command execution events (T1059) attributed to non-human principals; establish a behavioral baseline to distinguish normal agent activity from anomalous actions.
- 3. Step 3, Eradication:** Enforce least-privilege on all agent identities immediately; revoke credentials not tied to a documented, approved workload; scope API permissions to the minimum required for each agent's defined task; disable or sandbox agents invoking unvetted third-party plugins or MCP servers pending security review.
- 4. Step 4, Recovery:** Validate that revoked credentials are no longer active; confirm agent logging is producing attributable, coherent audit trails before restoring full operational scope; test prompt injection resistance on agents that consume external content by submitting adversarially crafted inputs in a controlled environment.
- 5. Step 5, Post-Incident:** Document the governance gap identified: most environments lack agent-specific identity policies, behavioral monitoring baselines, and human-in-the-loop checkpoints for high-risk actions; use this assessment to draft an agentic AI security policy covering identity lifecycle, permission scoping, audit requirements, and supply chain vetting for agent dependencies.

IR / Forensic Enrichment

Triage Priority	URGENT
Escalation Criteria	Escalate immediately to CISO and legal counsel if the Step 2 detection queries surface evidence that any agent identity accessed, bulk-exported, or transmitted data classified as PII, PHI, or PCI-scope data outside approved workflows, as this may trigger breach notification obligations under GDPR Article 33, HIPAA 45 CFR §164.412, or applicable state breach notification statutes; also escalate if discovered agent credentials show evidence of use from external IP ranges not associated with approved agent infrastructure, indicating potential third-party compromise of agent credentials.

<p>Recovery Notes</p>	<p>Before restoring any agent to full operational scope, verify three conditions in sequence: (1) all revoked credentials return 401/403 on active validation tests and no token refresh paths remain that could reactivate them; (2) agent audit logging is producing tamper-evident, attributable records to a log sink the agent process cannot write to or delete from; and (3) prompt injection testing has been completed against the specific input channels each agent consumes (web scraping output, email content, API responses, document ingestion pipelines). Maintain heightened monitoring on restored agent identities for a minimum of 30 days post-recovery, specifically watching for recurrence of the T1530 (Data from Cloud Storage) and T1078 (Valid Accounts) patterns identified in Step 2 baseline analysis. Any deviation from the documented behavioral baseline during this 30-day window should be treated as a new incident rather than noise.</p>
<p>Forensic Artifacts</p>	<p>Cloud provider IAM audit logs (AWS CloudTrail, Azure AD Audit Logs, GCP Admin Activity logs) filtered to service principal and non-human account actions — specifically creation events, policy attachment events, and credential issuance events — which establish when agent identities were provisioned and whether they were created through formal IAM processes or directly via API/CLI outside change management Cloud storage access logs (AWS S3 server access logs or data event CloudTrail, Azure Storage Diagnostic Logs) showing GetObject, ListBuckets, or equivalent bulk read operations attributed to agent service principal identities, with timestamps, source IPs, and byte counts — the specific artifact left by T1530 (Data from Cloud Storage) executed by an agentic workload Agent framework execution logs (LangChain callback trace files, AutoGPT activity logs, or equivalent framework-specific action history) recording the sequence of tool invocations, external API calls, and data retrievals executed by the agent — these are the closest equivalent to a process execution log for agentic workloads and establish what the agent actually did versus what it was authorized to do OAuth token and API key last-used records from the cloud provider IAM console and any secrets manager in use (AWS Secrets Manager access log, Azure Key Vault audit log) — these establish whether credentials provisioned for agent workloads were used outside expected time windows or from unexpected source IPs, and whether credentials believed to be inactive were in fact still being used Network flow logs (AWS VPC Flow Logs, Azure NSG Flow Logs, or on-premises firewall logs) for outbound connections from agent compute resources to external endpoints — specifically connections to MCP server IPs, third-party plugin APIs, or any external destination not in the approved agent dependency allowlist, which would indicate unauthorized data transmission or supply chain compromise of agent tooling</p>

Per-Action IR Details

Step 1: Inventory — conduct an immediate discovery sweep to enumerate all AI agents, service accounts, API keys, and tokens provisioned for agentic workloads; identify which were created outside formal IAM processes and which hold persistent or overprivileged credentials.

NIST Phase: Preparation

Reference: NIST 800-61r3 §2 — Preparation: establishing asset visibility and identity inventory as a precondition for effective incident handling

Controls: NIST IR-4 (Incident Handling) — preparation sub-phase requires knowing which principals exist before you can contain or eradicate them, NIST SI-4 (System Monitoring) — monitoring scope must include non-human identities; agents outside IAM are invisible to existing monitoring, NIST AC-2 (Account Management) — enumeration of all accounts, including service accounts and API tokens provisioned for agentic workloads, is a baseline AC-2 requirement, CIS 1.1 (Establish and Maintain Detailed Enterprise Asset Inventory) — AI agents and their associated credentials are enterprise assets requiring inventory, CIS 5.1 (Establish and Maintain an Inventory of Accounts) — all accounts including non-human agent identities must be tracked in the account inventory

Compensating: Run ``az ad sp list --all --query '[].{Name:displayName,AppId:appId,Created:createdDateTime}'`` (Azure) or ``aws iam list-users && aws iam list-roles`` filtered for agent-pattern naming conventions. For on-prem or multi-cloud, use osquery with the query ``SELECT * FROM users WHERE username LIKE '%agent%' OR username LIKE '%bot%' OR username LIKE '%svc%';`` adapted per environment. Cross-reference against your ticketing system or change management records — any service account with no associated ticket is a candidate for out-of-process provisioning. Document findings in a flat spreadsheet: agent name, owner, creation date, credential type, permission scope, and whether it exists in your IAM system of record.

Evidence: Before modifying any account, snapshot the current permission state: export IAM role/policy attachments for each discovered agent identity, capture API key last-used timestamps from cloud provider IAM consoles (AWS: ``aws iam get-access-key-last-used``, Azure: audit logs under Sign-in activity for service principals), and preserve cloud provider audit logs showing when credentials were created and by whom. For agents using OAuth tokens or MCP server credentials, capture the token issuance records from your IdP (Okta system log, Azure AD audit log, or equivalent) before any revocation action.

Step 2: Detection — query your SIEM and cloud access logs for service account activity patterns inconsistent with human-initiated sessions: off-hours API calls, bulk data access across storage buckets (T1530), lateral credential use (T1078), and command execution events (T1059) attributed to non-human principals; establish a behavioral baseline to distinguish normal agent activity from anomalous actions.

NIST Phase: Detection Analysis

Reference: NIST 800-61r3 §3.2 — Detection and Analysis: correlating indicators from multiple sources to characterize adversarial or anomalous behavior by non-human principals

Controls: NIST SI-4 (System Monitoring) — requires monitoring for anomalous activity including non-human principals executing bulk data operations (T1530) and lateral credential reuse (T1078), NIST AU-6 (Audit Record Review, Analysis, and Reporting) — audit logs for agent identities must be actively reviewed for indicators of unauthorized or anomalous action patterns, NIST AU-2 (Event Logging) — event types generated by agentic workloads (API calls, storage access, command execution) must be explicitly included in the organization's event logging policy, NIST IR-5 (Incident Monitoring) — tracking and documenting anomalous agent activity as potential incidents requires this detection infrastructure to exist, CIS 8.2 (Collect Audit Logs) — audit logging must be enabled and actively collected for cloud APIs, storage services, and identity providers accessed by agent workloads

Compensating: Without a SIEM, use cloud-native free tooling: AWS CloudTrail Insights (enable anomaly detection for IAM) and S3 server access logging to detect bulk GetObject calls; Azure Monitor with a KQL query ``AuditLogs | where InitiatedBy.app.displayName contains 'agent' | where ActivityDateTime between (ago(30d)..now())`` to surface service principal actions. For T1059 (command execution) attributed to agent processes, deploy Sysmon with EventID 1 (Process Create) and filter parent processes matching your agent runtime (e.g., Python interpreter, Node.js). Use the Sigma rule ``proc_creation_win_susp_non_human_account_cmd.yml`` as a starting template. For behavioral baseline without SIEM, export 30 days of CloudTrail or Azure AD sign-in logs to CSV and use a Python pandas script to compute per-identity call frequency, time-of-day distribution, and accessed resource entropy.

Evidence: Preserve the following before baseline alteration: AWS CloudTrail logs showing GetObject/ListBuckets calls by service principal ARNs correlated with agent identities (S3 data event logs, not just management events — ensure these are enabled); Azure Storage analytics logs or Diagnostic Logs showing blob access by service principal object IDs; OAuth token introspection records showing which scopes were granted and when tokens were exchanged; cloud provider sign-in logs showing source IP and user-agent strings for agent API calls (non-human agents typically present SDK user-agent strings such as ``aws-sdk-python`` or ``@azure/identity`` which differ from human console sessions); any prompt or instruction logs retained by the agent framework (LangChain trace logs, AutoGPT action history files, or equivalent) that would show what the agent was instructed to do and what it executed.

Step 3: Eradication — enforce least-privilege on all agent identities immediately; revoke credentials not tied to a documented, approved workload; scope API permissions to the minimum required for each agent's defined task; disable or sandbox agents invoking unvetted third-party plugins or MCP servers pending security review.

NIST Phase: Eradication

Reference: NIST 800-61r3 §3.4 — Eradication: removing the conditions that allowed the incident to occur, which for ungoverned agent identities means eliminating overprivileged and undocumented credentials from the environment

Controls: NIST AC-6 (Least Privilege) — each agent identity must be scoped to the minimum permissions required for its documented task; broad IAM roles or wildcard API scopes must be replaced with task-specific policies, NIST AC-2 (Account Management) — accounts not tied to an approved, documented workload must be disabled or deleted per account management policy, NIST IR-4 (Incident Handling) — eradication actions must be coordinated to avoid displacing the threat to undiscovered agent identities not yet enumerated, NIST CM-7 (Least Functionality) — agents invoking unvetted third-party plugins or MCP servers represent unnecessary functionality that must be disabled pending review, CIS 5.4 (Restrict Administrator Privileges to Dedicated Administrator Accounts) — agent identities holding administrative or broad data-plane permissions must have those privileges revoked and replaced with task-scoped roles, CIS 6.2 (Establish an Access Revoking Process) — the revocation of undocumented agent credentials must follow a documented process to ensure completeness and auditability

Compensating: For AWS, use `aws iam put-role-policy` to replace overpermissive inline policies with task-scoped equivalents, and `aws iam delete-access-key` to revoke undocumented keys — run `aws iam list-access-keys --user-name` first to enumerate all key IDs. For Azure, use az ad sp credential reset to rotate and az role assignment delete to remove overpermissive role assignments. To sandbox agents invoking unvetted MCP servers, implement outbound firewall rules at the host level using iptables -A OUTPUT -m owner --uid-owner -j DROP (Linux) or Windows Firewall outbound rules scoped to the agent service account SID. Document every revocation action with timestamp, actor, and justification in a change log before executing.`

Evidence: Before revoking credentials, capture: the full permission policy document attached to each agent identity (IAM policy JSON, Azure role assignment export) to establish what access existed at time of discovery; any secrets manager entries (AWS Secrets Manager, Azure Key Vault access logs) showing which agent processes retrieved credentials and when; MCP server connection logs or plugin invocation records if the agent framework retains them (check LangChain callback handler logs, AutoGPT plugin call history, or equivalent framework logs); network flow logs (AWS VPC Flow Logs, Azure NSG Flow Logs) showing outbound connections from agent compute to external plugin endpoints or MCP server IPs, which may indicate data exfiltration paths that existed prior to revocation.

Step 4: Recovery — validate that revoked credentials are no longer active; confirm agent logging is producing attributable, coherent audit trails before restoring full operational scope; test prompt injection resistance on agents that consume external content by submitting adversarially crafted inputs in a controlled environment.

NIST Phase: Recovery

Reference: NIST 800-61r3 §3.5 — Recovery: restoring systems to verified secure operation, which for agentic workloads requires both credential validation and operational integrity testing before full scope is restored

Controls: NIST IR-4 (Incident Handling) — recovery sub-phase requires verification that eradication was successful before restoring normal operations, NIST AU-3 (Content of Audit Records) — agent audit trails must contain sufficient attribution to identify which agent identity performed which action, on which resource, at what time, before operational scope is restored, NIST AU-9 (Protection of Audit Information) — agent-generated audit logs must be protected from tampering by the agent itself, requiring a separate, write-once log sink outside agent process control, NIST SI-7 (Software, Firmware, and Information Integrity) — prompt injection resistance testing is an integrity verification requirement before restored agents consume untrusted external content, NIST SI-2 (Flaw Remediation) — prompt injection vulnerabilities in agents consuming external content are software flaws requiring remediation verification before recovery, CIS 7.1 (Establish and Maintain a Vulnerability Management Process) — prompt injection testing prior to re-enabling external content consumption is a vulnerability validation step within the remediation lifecycle

Compensating: Validate revoked credentials are inactive by attempting an API call using the revoked key and confirming a 403/401 response: `curl -H 'Authorization: Bearer ' https://v1/health` — log the response. For audit trail validation without a SIEM, use jq to parse agent framework logs and verify each action record contains: timestamp, agent identity, action type, target resource, and outcome. For prompt injection testing without a commercial tool, use the open-source promptmap or garak framework (GitHub: leondz/garak) to submit adversarial payloads to the agent's input interface; test categories should include instruction override attempts, indirect injection via tool output, and role-switching prompts. Run tests in an isolated environment using a sandboxed copy of the agent with no access to production data or credentials.`

Evidence: Before restoring operational scope, capture and retain: API call attempt logs confirming revoked credentials return authentication errors (preserve these as evidence of successful eradication); agent logging output showing attribution chain from trigger event through tool invocations to action completion — if this chain cannot be reconstructed, logging is insufficient for recovery; prompt injection test results including the specific payloads submitted, the agent's responses, and whether any adversarial payload caused the agent to deviate from its defined task or attempt unauthorized actions; a final inventory snapshot of all agent identities and their current permission scopes to establish the verified post-recovery baseline.

Step 5: Post-Incident — document the governance gap identified: most environments lack agent-specific identity policies, behavioral monitoring baselines, and human-in-the-loop checkpoints for high-risk actions; use this assessment to draft an agentic AI security policy covering identity lifecycle, permission scoping, audit requirements, and supply chain vetting for agent dependencies.

NIST Phase: Post Incident

Reference: NIST 800-61r3 §4 — Post-Incident Activity: lessons learned and policy improvement to close the structural governance gaps that made the incident possible

Controls: NIST IR-8 (Incident Response Plan) — the IR plan must be updated to explicitly include agentic AI principals as security entities subject to incident handling procedures, NIST IR-4 (Incident Handling) — post-incident review must feed directly into updated incident handling procedures covering agent identity lifecycle and behavioral monitoring, NIST AC-2 (Account Management) — the agentic AI security policy must extend account management procedures to cover agent identity provisioning, periodic review, and deprovisioning lifecycle, NIST RA-3 (Risk Assessment) — the governance gap identified must be formally risk-assessed; the new policy's control requirements should be justified by this risk assessment, NIST SA-9 (External System Services) — supply chain vetting for agent dependencies (third-party plugins, MCP servers, tool integrations) is a system acquisition and supply chain risk management requirement, CIS 7.2 (Establish and Maintain a Remediation Process) — the agentic AI security policy must include a documented remediation process for agent-specific vulnerabilities and governance violations, CIS 5.1 (Establish and Maintain an Inventory of Accounts) — the policy must mandate that all future agent identities are registered in the account inventory at provisioning time, not discovered retrospectively

Compensating: For a 2-person team without a dedicated GRC platform, build the agentic AI security policy as a structured markdown document version-controlled in Git; include required sections: agent identity registration form (owner, purpose, permission justification, review cadence), a quarterly audit checklist (enumerate all agent accounts, verify permissions match documented scope, confirm logging is active), and a human-in-the-loop decision tree defining which action categories (data deletion, external API calls, credential access) require human approval before execution. For supply chain vetting of plugins and MCP servers, maintain a simple allowlist in a JSON file checked into the same repo, and implement a pre-deployment hook that blocks agent configuration files referencing unlisted plugin endpoints.

Evidence: The post-incident documentation package must include: the full inventory produced in Step 1 (agent enumeration with out-of-process identities flagged) as evidence of the governance gap's scope; the SIEM or cloud log queries from Step 2 and their results, showing which anomalous activity patterns existed and for how long before detection — this establishes dwell time; the permission policy snapshots from Step 3 showing the delta between what agents held and what least-privilege would have permitted — this quantifies the blast radius; and the prompt injection test results from Step 4 as evidence of residual technical risk. Together these artifacts constitute the factual basis for the policy drafted in this step and should be retained for any regulatory inquiry or audit.

Detection Guidance

Detection for agentic workloads requires treating AI agents as first-class identity principals in your logging and monitoring stack. Key indicators: (1) Credential creation events outside IAM workflows, audit logs for API key or token generation not tied to a provisioning ticket or human initiator. (2) Non-human session anomalies, cloud provider access logs (AWS CloudTrail, Azure Monitor, GCP Audit Logs) showing service account activity at unusual hours, bulk object enumeration, or cross-account access not matching defined agent scope. (3) Data access volume spikes, S3, Azure Blob, or GCS access logs showing abnormal read volumes from a service

principal, relevant to T1530. (4) Lateral credential use, T1078 indicators: service accounts authenticating to systems outside their designated scope. (5) Command execution events, T1059: shell or scripting interpreter invocations attributed to agent processes, particularly in containerized or serverless environments. (6) Prompt injection behavioral markers, agents producing outputs or taking actions inconsistent with their configured task scope, particularly after consuming external documents, emails, or web content. (7) Logging gaps, absence of attributable logs for known agent actions is itself a signal; CWE-778 conditions should be treated as a detection gap requiring immediate remediation. Governance-based detection relies on behavioral baselines and policy deviations rather than discrete IOCs (file hashes, domains, etc.). Each detection signal is policy-specific to your environment.

Framework Mappings

MITRE-ATTACK

- **T1530** — Data from Cloud Storage
- **T1195** — Supply Chain Compromise
- **T1552** — Unsecured Credentials
- **T1059** — Command and Scripting Interpreter
- **T1078** — Valid Accounts

NIST-800-53R5

- **SA-9** — External System Services
- **SR-2** — Supply Chain Risk Management Plan
- **SR-3** — Supply Chain Controls and Processes
- **SI-7** — Software, Firmware, and Information Integrity
- **CM-7** — Least Functionality
- **SI-3** — Malicious Code Protection
- **SI-4** — System Monitoring
- **AC-2** — Account Management
- **AC-6** — Least Privilege
- **IA-2** — Identification and Authentication (Organizational Users)
- **IA-5** — Authenticator Management
- **SI-10** — Information Input Validation

OWASP-TOP10-2021

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

CIS-V8

- **5.2** — Use Unique Passwords
- **5.4** — Restrict Administrator Privileges to Dedicated Administrator Accounts

- **6.8** — Define and Maintain Role-Based Access Control
- **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
- **8.2** — Collect Audit Logs

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
- **A.5.23** — Information security for use of cloud services

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.CM-01** — Networks and network services are monitored

MITRE ATT&CK Mapping

Technique ID	Technique Name	Tactic
T1530	Data from Cloud Storage	Collection
T1195	Supply Chain Compromise	Initial-Access
T1552	Unsecured Credentials	Credential-Access
T1059	Command and Scripting Interpreter	Execution
T1078	Valid Accounts	Defense-Evasion

Sources

Source	URL	Tier
gemini	https://securityboulevard.com/2026/05/ai-agents-are-creating-a-new-...	T3
The Hidden Vulnerabilities of Autonomous AI Agents	https://www.innovativehumancapital.com/article/when-delegation-goes...	T3

Source	URL	Tier
Agentic AI security: Risks & governance for enterprises McKinsey	https://www.mckinsey.com/capabilities/risk-and-resilience/our-insig...	T2
The 2025 AI Agent Security Landscape: Players, Trends, and Risks	https://www.obsidiansecurity.com/blog/ai-agent-market-landscape	T3
AI Agents Are the Biggest Data Security Threat You're Not Governing	https://www.kiteworks.com/cybersecurity-risk-management/ai-agents-u...	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 2026-05-11 18:50 UTC by TJS Security Command Center