



NIST AI Agent Standards: Enterprise Governance Implications

The CAISI Initiative and What It Means for Agentic AI Security
Programs

Unofficial AI-assisted Research

Cloud Security Alliance AI Safety Initiative

2026-03-24

© 2026 Cloud Security Alliance. Some rights reserved.

You may download, store, display, view, print, redistribute, and link to this document in its original, unmodified form, provided that attribution to the Cloud Security Alliance is maintained and all trademark and copyright notices remain intact.

This document may not be modified or altered. You may quote portions of the document as permitted by the Fair Use provisions of the United States Copyright Act, provided that attribution is given to the Cloud Security Alliance.

This document may be shared on professional and social media platforms in its original form with attribution.

This document was generated with AI assistance and has not undergone official CSA review and approval processes.

Key Takeaways

- NIST's Center for AI Standards and Innovation (CAISI) announced the AI Agent Standards Initiative on February 17, 2026, establishing a three-pillar framework addressing interoperability, security, and open-protocol standards for autonomous AI agents [1].
 - A companion NCCoE concept paper published February 5, 2026 proposes applying existing identity and access management standards – including OAuth, OpenID Connect, SPIFFE/SPIRE, and SCIM – to AI agents treated as identifiable enterprise entities rather than anonymous automation [2].
 - NIST's March 2025 update to its adversarial machine learning taxonomy (AI 100-2 E2025) explicitly classifies agentic systems as a distinct attack surface, adding dedicated treatment of prompt injection, tool misuse, and action-hijacking that compounds in severity when models are connected to real-world execution capabilities [3].
 - The COSAiS project is developing SP 800-53 control overlays specifically for single-agent and multi-agent AI systems; enterprise compliance programs that currently map controls to traditional application categories should track this work as it will directly affect federal contractor obligations and likely influence private-sector audit frameworks [4].
 - NIST's December 2025 preliminary draft of NIST IR 8596 (Cyber AI Profile) maps cybersecurity framework functions directly to AI system governance, providing organizations a structured path to unify AI risk management and conventional cybersecurity controls under CSF 2.0 [5].
-

Background

As NIST acknowledged in its January 2026 request for information, the existing NIST guidance portfolio contains no provisions specifically scoped to autonomous, action-taking AI systems – a gap that reflects how quickly enterprise deployments have advanced relative to the governance frameworks available to manage them [7]. Unlike conventional software, which follows deterministic paths that can be audited against specifications, AI agents reason across ambiguous inputs, maintain state across sessions, invoke external tools and APIs, and in multi-agent architectures delegate tasks to subagents operating under

partially independent authority. These properties create governance problems that existing frameworks – including ISO 27001, SOC 2, and the original NIST AI Risk Management Framework published in January 2023 – were not designed with agentic systems as an explicit scope [6].

NIST's January 2026 request for information (Federal Register docket NIST-2025-0035) sought public input on five dimensions: unique threats facing agent systems, methods for lifecycle-wide security improvement, gaps in current cybersecurity approaches, risk evaluation methodologies for development-stage agents, and deployment safeguards including access constraint and monitoring architectures. The comment period closed March 9, 2026, and CAISI's sector-specific listening sessions in healthcare, finance, and education are scheduled for April 2026 [7].

This research note analyzes the CAISI initiative and its companion standards activities, examines the enterprise governance implications of the six core themes NIST has identified, and connects NIST's evolving guidance to the Cloud Security Alliance frameworks most directly applicable to agentic AI deployments.

Security Analysis

The CAISI AI Agent Standards Initiative

NIST's Center for AI Standards and Innovation launched the AI Agent Standards Initiative with three strategic pillars, each targeting a different layer of the governance problem [1].

The first pillar focuses on industry-led voluntary standards. CAISI is hosting technical convenings and conducting gap analyses to identify where existing standards can be extended to cover agentic deployments and where new instruments are needed. This work includes coordinating U.S. participation in international standards bodies through collaboration with the National Science Foundation. The goal is not to replace existing enterprise security frameworks but to produce voluntary guidelines that translate their abstract principles into agent-specific controls. The first version of what NIST describes as an AI Agent Interoperability Profile is planned for release in Q4 2026 [1].

The second pillar targets community-led open protocols. NIST has identified the emergence of interoperability protocols such as the Model Context Protocol (MCP) as a structural development that requires governance input before fragmentation sets in. MCP, which allows agents to connect to tools, data sources, and services through a standardized interface, has gained significant developer adoption since its November 2024 release, and NIST sees eliminating obstacles to secure protocol interoperability

as a prerequisite for coherent governance across multi-vendor agent deployments. The NSF is funding open-source AI agent ecosystem development through its Pathways to Enable Secure Open-Source Ecosystems program to advance this objective [1].

The third pillar invests in foundational research on agent authentication, identity infrastructure, and authorization controls. NIST has committed to developing security evaluations that can inform protocol design and, over time, enable comparable assessments across different agent platforms. This research is directly connected to the NCCoE concept paper on agent identity and authorization, which represents the most operationally specific guidance NIST has produced to date on the mechanics of managing agent access in enterprise environments [2].

Six Core Governance Themes

CAISI's announcements and the accompanying NCCoE concept paper collectively identify six governance themes that the initiative treats as foundational [1][2]. These themes map closely to the attack patterns documented in real-world incidents and in NIST's own adversarial taxonomy, and they provide a useful organizing structure for enterprise security teams working to build governance programs before NIST's formal guidance matures.

Agent identity and authentication is the first and arguably most consequential theme. A common pattern in current deployments is the assignment of agents to shared API keys, service account credentials, or IAM roles with broad scope that were designed for human operators or deterministic applications. The NCCoE concept paper argues that agents must instead be treated as identifiable entities within enterprise identity systems, with distinct cryptographic identities, short-lived credentials, and enrollment processes that satisfy the same standards applied to privileged human users [2]. The paper specifically references SPIFFE/SPIRE as a workload identity framework applicable to agent deployments and notes that OAuth and OpenID Connect can be extended to agent-specific grant types to support task-scoped token issuance. Authorization bypass patterns that have emerged in production deployments – where IAM systems evaluate permissions against the agent's identity rather than the human requester's, causing user-level restrictions to stop applying – are a documented manifestation of the failure to treat agents as governed identities [2][8].

Least-privilege authorization extends the identity problem into access control design. OWASP identifies as a common pattern that agents operate with permissions spanning multiple systems, provisioned for the broadest expected use case rather than scoped to individual tasks [9]. OWASP's formulation of "Excessive Agency" – broken into excessive functionality, excessive permissions, and excessive autonomy – provides a structured vocabulary for identifying where agent privilege exceeds operational necessity [9]. NIST's preferred architecture calls for just-in-time access tied to task duration,

action-level human approval gates for high-impact operations, and tool manifests that restrict the set of capabilities available to an agent at runtime. These controls do not exist in most current agentic deployments and require deliberate architectural investment to implement.

Audit and non-repudiation presents a technical challenge that the scale and speed of autonomous agents makes genuinely novel. An agent processing documents, sending messages, invoking APIs, and spawning subagents in the course of a single workflow can generate audit events at rates that overwhelm traditional SIEM tooling. NIST's guidance calls for comprehensive records that capture not only what an agent did but the context it received, the permissions it held, and any human approval events in the chain [1][2]. The minimum elements are: initiator identity (human, application, or agent), action type, tool or system accessed, timestamp, result, and approval chain. Without non-repudiation across multi-agent workflows, post-incident forensic analysis becomes structurally difficult, because the attribution chain for a harmful outcome may span multiple agents that each acted within their own perceived authority.

Post-deployment monitoring encompasses the functional, operational, security, and compliance dimensions of runtime agent oversight. Unlike model evaluation, which can be conducted in controlled test environments before deployment, behavioral monitoring for autonomous agents must occur continuously in production, because agents encounter contexts during operation that cannot be fully anticipated in testing. NIST's guidance on this theme is less mature than on identity and authorization, but it explicitly identifies behavioral drift – where an agent's decision patterns shift over time due to accumulated context, tool output, or environmental changes – as a monitoring target that differs in mechanism from traditional behavioral analytics. AI behavioral drift is shaped by context window contents and tool outputs rather than access pattern statistics, requiring monitoring approaches tailored to the inference layer rather than the access layer.

Prompt injection as an architectural control represents a shift in how NIST categorizes a vulnerability that the industry has historically treated as a model quality problem. NIST AI 100-2 E2025 explicitly classifies both direct and indirect prompt injection as attack vectors capable of hijacking agents to execute arbitrary code or exfiltrate data, and the CAISI initiative treats defense against injection as an architectural control requirement rather than a fine-tuning or red-teaming exercise [3]. Indirect prompt injection – where malicious instructions are embedded in documents, email messages, web pages, or Slack channels that an agent retrieves as part of its normal operation – has been demonstrated in production systems at enterprise scale, most notably in the Slack AI vulnerability disclosed in August 2024 and in CVE-2025-32711 (EchoLeak) against Microsoft 365 Copilot, which achieved zero-click data exfiltration through a chained prompt injection exploiting the platform's content rendering pipeline [10] [11].

Multi-agent interoperability standards address the emergent risks of agent orchestration, where an orchestrator agent delegates tasks to subagents that may hold independent credentials, maintain their own context, and apply their own tool permissions. The inter-agent communication surface introduces trust exploitation risks that single-agent deployments do not present: a compromised subagent can influence the orchestrator's reasoning through crafted outputs; a malicious tool server can redirect agent behavior at the invocation layer; and cascading failures can propagate across agent hierarchies in ways that no single agent's guardrails were designed to contain. OWASP's dedicated Top 10 for Agentic Applications (published December 2025) catalogues these risks as ASI07 (Insecure Inter-Agent Communication), ASI08 (Cascading Failures), and ASI09 (Human-Agent Trust Exploitation) [12].

The Standards Pipeline and Enterprise Planning Horizon

The NIST guidance landscape for agentic AI is active but uneven. Several documents are final and available for adoption today, while the most specifically agentic-focused instruments remain in draft or early development. Enterprise security teams need to distinguish between what is actionable now and what requires monitoring.

NIST AI 100-1 (AI RMF 1.0, 2023) and NIST AI 600-1 (Generative AI Profile, July 2024) are final and provide the governance vocabulary – Govern, Map, Measure, Manage – that should structure enterprise AI risk programs regardless of deployment architecture [6][13]. These documents do not address agentic systems as a distinct category, but the AI RMF's function-based structure accommodates agentic use cases and the Generative AI Profile's risk categories (including data poisoning, prompt injection, and output integrity) apply directly to LLM-based agent pipelines.

NIST AI 100-2 E2025 (March 2025) is final and represents the most operationally specific adversarial security guidance currently available for agentic deployments [3]. Its explicit treatment of prompt injection as an agent-hijacking technique, and its taxonomy of agentic attack scenarios, should be incorporated into threat modeling for any agent deployment connected to real-world execution capabilities.

NIST IR 8596 (Cyber AI Profile, preliminary draft, December 2025) is the earliest stage document in the pipeline, developed with over 6,500 contributors and organized around CSF 2.0's six functions [5]. It will have broad applicability once finalized, but its agentic-specific provisions are limited in scope pending the findings from CAISI's RFI and listening sessions.

The COSAiS project's SP 800-53 overlays for agentic use cases represent the most consequential forthcoming guidance for organizations operating under federal compliance frameworks [4]. When finalized, these overlays will translate the CAISI governance themes into specific SP 800-53 controls,

providing a structured path for federal agencies and contractors to incorporate agentic AI governance into existing FISMA compliance programs. Private-sector organizations that have adopted SP 800-53 voluntarily should treat the COSAiS timeline as a planning constraint.

Recommendations

Immediate Actions

Enterprises with production agentic AI deployments should conduct a credential audit before the Q4 2026 publication of NIST's AI Agent Interoperability Profile. The audit should identify every agent or agentic workflow that operates under shared credentials, service account keys, or IAM roles not provisioned specifically for that agent. Each such instance represents both a governance gap and an active attack surface. Remediating the highest-risk cases – agents with write access to production data stores, agents that invoke external APIs with elevated permissions, agents operating in regulated environments – should not wait for formal standards to finalize.

Organizations should also review their model invocation logging posture. Logging is disabled by default in some cloud AI service configurations; organizations should verify their current posture and enable invocation logging where not already active. Routing logs to a retention system outside the agent's own permission scope is a relatively low-cost control compared to the forensic value it provides, and has direct relevance to non-repudiation requirements that NIST has identified as foundational.

Short-Term Mitigations

Security teams should incorporate NIST AI 100-2 E2025's agentic attack taxonomy into their threat models for any deployment where agents retrieve content from external sources – documents, email, web pages, database records, or messages from other agents. The taxonomy's treatment of indirect prompt injection as an architectural vulnerability (not a model-quality deficiency) should shift both design reviews and penetration testing scope. Agentic pipelines should be assessed for injection pathways in every external data source the agent can access, not just in adversarial inputs directly visible to the model.

For multi-agent architectures specifically, the least-privilege principle should be applied at the agent-to-agent delegation boundary. Each subagent should receive the minimum permissions needed for its specific task, scoped to the duration of that task, and should not inherit the orchestrator's broader

permission set by default. Human approval gates should be defined for any action category – data deletion, external communication, credential provisioning – where an errant decision would be difficult to reverse.

Strategic Considerations

Enterprise AI governance programs should treat the CAISI initiative's development timeline as a forcing function for framework alignment. Organizations that have built AI governance around the AI RMF's four functions should begin mapping their agentic deployments to CAISI's six governance themes now, while NIST's agent-specific instruments are in development. This creates both an internal governance artifact and a mechanism for submitting informed comments as NIST solicits public feedback – sector-specific listening sessions in April 2026 provide a mechanism for enterprise practitioners to submit input that may inform the shape of the final guidance.

The COSA's SP 800-53 overlay timeline should be incorporated into compliance roadmaps. Federal contractors in particular should expect that agentic AI deployments will eventually require demonstrable control coverage under tailored SP 800-53 families; beginning that mapping work against the existing AI 100-2 E2025 taxonomy will reduce remediation effort when the overlays are formally published.

CSA Resource Alignment

The following CSA publications provide implementation paths for the governance themes described above; practitioners should also consult complementary frameworks including MITRE ATLAS, CISA's AI Roadmap, and ISO/IEC 42001 depending on their regulatory context.

The **MAESTRO framework** (Multi-Agent Environment, Security, Threat, Risk, and Outcome), published by CSA in February 2025, provides a seven-layer threat modeling structure designed specifically for agentic AI systems [14]. Where NIST's guidance identifies governance themes, MAESTRO provides the threat decomposition methodology to translate those themes into specific control requirements at each layer of an agentic architecture – from the Foundation Model through Data Operations, Agent Frameworks, Deployment Infrastructure, and the broader Agent Ecosystem. Enterprise security teams should use MAESTRO as the threat modeling complement to NIST's governance framing.

CSA's **Agentic AI Identity and Access Management** publication addresses the agent identity problem that the NCCoE concept paper raises [15]. It provides an architecture specifically designed for managing agent credentials, scoping agent permissions, and establishing trust hierarchies in multi-agent

deployments – exactly the problem that NIST has identified as requiring new IAM patterns that exceed the assumptions of traditional frameworks.

The **AAGATE platform** (Agentic AI Governance Assurance and Trust Engine), published by CSA in December 2025, is a Kubernetes-native reference architecture that translates the NIST AI RMF's four functions into operational controls for agentic deployments [16]. It provides a concrete implementation path for organizations that have adopted the AI RMF conceptually but lack the runtime mechanisms to enforce its governance principles against agents operating at machine speed.

The **AI Infrastructure Control Matrix (AICM)** extends the Cloud Controls Matrix to AI-specific governance requirements and provides the control vocabulary for mapping agent security practices to STAR-level assurance. Organizations conducting AI governance assessments should reference AICM rather than CCM alone, as AICM is a superset specifically designed to capture the additional control surface that AI systems – and agentic systems in particular – introduce beyond conventional cloud deployments.

NIST's identification of prompt injection as an architectural control requirement aligns directly with CSA's Zero Trust guidance. The Zero Trust principle of "never trust, always verify" applied to agent inputs means treating all externally sourced content – regardless of the channel or source identity – as potentially adversarial and validating it against expected schemas, privilege boundaries, and behavioral constraints before execution.

References

- [1] NIST Center for AI Standards and Innovation (CAISI), "Announcing the AI Agent Standards Initiative: Interoperable and Secure," NIST News, February 17, 2026. <https://www.nist.gov/news-events/news/2026/02/announcing-ai-agent-standards-initiative-interoperable-and-secure>
- [2] NCCoE, "Accelerating the Adoption of Software and AI Agent Identity and Authorization," NIST/NCCoE Concept Paper (Initial Public Draft), February 5, 2026. <https://www.nccoe.nist.gov/publications/other/accelerating-adoption-software-and-ai-agent-identity-and-authorization-concept>
- [3] NIST, "Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations," NIST AI 100-2 E2025, March 2025. <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-2e2025.pdf>
- [4] NIST CSRC, "Control Overlays for Securing AI Systems (COSaIS)," Project page, active 2025-2026. <https://csrc.nist.gov/projects/cosais>
- [5] NIST, "Cybersecurity Framework Profile for Artificial Intelligence (Cyber AI Profile)," NIST IR 8596 (Preliminary Draft), December 16, 2025. <https://nvlpubs.nist.gov/nistpubs/ir/2025/NIST.IR.8596.iprd.pdf>
- [6] NIST, "Artificial Intelligence Risk Management Framework," NIST AI 100-1, January 26, 2023. <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>
- [7] NIST CAISI, "Request for Information Regarding Security Considerations for Artificial Intelligence Agents," Federal Register Docket NIST-2025-0035, January 8, 2026. <https://www.federalregister.gov/documents/2026/01/08/2026-00206/request-for-information-regarding-security-considerations-for-artificial-intelligence-agents>
- [8] The Hacker News, "AI Agents Are Becoming Authorization Bypass Paths," January 2026. <https://thehackernews.com/2026/01/ai-agents-are-becoming-privilege.html>
- [9] OWASP GenAI Security Project, "OWASP Top 10 for LLM Applications 2025," published November 2024. <https://genai.owasp.org/resource/owasp-top-10-for-llm-applications-2025/>
- [10] PromptArmor, "Data Exfiltration from Slack AI via Indirect Prompt Injection," August 20, 2024. <https://www.promptarmor.com/resources/data-exfiltration-from-slack-ai-via-indirect-prompt-injection>

[11] Aim Security / The Hacker News, "Zero-Click AI Vulnerability (CVE-2025-32711) Exposes Microsoft 365 Copilot," June 2025. <https://thehackernews.com/2025/06/zero-click-ai-vulnerability-exposes.html>

[12] OWASP GenAI Security Project, "OWASP Top 10 for Agentic Applications," December 9, 2025. <https://genai.owasp.org/2025/12/09/owasp-top-10-for-agentic-applications-the-benchmark-for-agentic-security-in-the-age-of-autonomous-ai/>

[13] NIST, "Generative AI Profile," NIST AI 600-1, July 26, 2024. <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.600-1.pdf>

[14] Cloud Security Alliance, "MAESTRO: Agentic AI Threat Modeling Framework," CSA Blog, February 2025. <https://cloudsecurityalliance.org/blog/2025/02/06/agentic-ai-threat-modeling-framework-maestro>

[15] Cloud Security Alliance, "Agentic AI Identity and Access Management: A New Approach," CSA Publication, August 18, 2025. <https://cloudsecurityalliance.org/artifacts/agentic-ai-identity-and-access-management-a-new-approach>

[16] Cloud Security Alliance, "AAGATE: A NIST AI RMF-Aligned Governance Platform for Agentic AI," CSA Blog, December 22, 2025. <https://cloudsecurityalliance.org/blog/2025/12/22/aagate-a-nist-ai-rmf-aligned-governance-platform-for-agentic-ai>