

# Five Eyes Issues First Joint Agentic AI Security Guidance

What CISA's 'Careful Adoption of Agentic AI Services' Means for Enterprise Security Programs

2026-05-03

 AI-assisted Rapid Research



**© 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

- On May 1, 2026, six national cybersecurity agencies – CISA, NSA, Australia's ASD ACSC, the Canadian Centre for Cyber Security, New Zealand's NCSC, and the UK's NCSC – jointly published "Careful Adoption of Agentic AI Services," the first coordinated multi-government security guidance specifically addressing agentic AI systems [1][2][8][9].
- The guidance identifies five distinct risk categories for agentic deployments: privilege, design and configuration, behavioral, structural, and accountability risks.
- Prompt injection is characterized as the most persistent and difficult-to-fix threat facing agentic systems, stemming from a fundamental design constraint of language models that cannot be fully resolved through input sanitization [2].
- The agencies' core message is integrationist: agentic AI security should extend existing zero trust, defense-in-depth, and least-privilege frameworks rather than wait for purpose-built AI security standards to emerge.
- CSA analysis identifies three implementation pillars across the guidance's recommendations – governance, visibility, and least-privilege enforcement – that enterprise programs should prioritize before agentic deployments scale beyond their current footprint.

## Background

### From Chatbot to Autonomous Actor

Agentic AI systems represent a qualitative departure from the deployments that most enterprise security programs were designed to govern. Where a conventional large language model responds to queries and returns text, an agentic system plans tasks autonomously, executes API calls, modifies files, sends communications, and chains multi-step actions without human intervention at each stage. Products such as Microsoft 365 Copilot, GitHub Copilot Workspace, and Salesforce Agentforce can operate under users' delegated authority with varying degrees of human oversight depending on configuration – in autonomous modes, agents may read, write, and transmit enterprise data with minimal human intervention between task assignment and task completion.

The security implications of this architecture are not incremental. An agent granted access to a user's email, calendar, and document repositories can read, modify, or forward sensitive information across organizational boundaries in seconds. When agentic systems are connected in multi-agent pipelines – where the output of one agent becomes the input to the next – the attack surface expands further still. A single compromised component in such a chain can corrupt every downstream process, potentially altering files, modifying access controls, and overwriting audit trails before a security analyst has the opportunity to observe anomalous behavior.

## A Coordinated Government Response

Against this backdrop, six national cybersecurity agencies published "Careful Adoption of Agentic AI Services" on May 1, 2026 [1][2][9]. The co-authoring bodies – CISA, the NSA, Australia's ASD Australian Cyber Security Centre, the Canadian Centre for Cyber Security, New Zealand's NCSC, and the UK's NCSC – represent all five nations of the Five Eyes intelligence-sharing alliance, with the United States contributing two agencies, a coalition whose joint imprimatur signals that agentic AI security has crossed the threshold from emerging concern to active policy priority [2]. This is the first time these agencies have issued joint guidance specifically on agentic AI, and it follows a period in which agentic systems have moved from experimental pilots into production deployments across critical infrastructure and enterprise environments.

The guidance is explicit that a new security discipline is not required. Organizations should integrate agentic AI into the governance structures and technical controls they already maintain, applying the same principles – zero trust, defense-in-depth, least privilege – that underpin mature cybersecurity practice. That integrationist framing has a significant practical implication for enterprise security programs: it provides a clear rationale to act now, using frameworks already in place, rather than waiting for bespoke AI security standards to reach maturity.

## Security Analysis

### Five Categories of Agentic Risk

The guidance organizes its threat analysis around five risk categories that capture the distinctive hazards of autonomous AI operation [1][2].

**Privilege risk** arises when agents are granted access broader than their immediate task requires. Because agents operate autonomously and at machine speed, a single compromise in a highly-privileged context enables disproportionate damage before human intervention becomes possible. The guidance

stresses that organizations should never grant agents unfettered access to critical systems – a principle that early agentic deployments have often failed to meet as organizations prioritized operational convenience during rapid rollouts.

**Design and configuration risks** refer to security gaps introduced before deployment: poorly scoped integrations, weak authentication at agent boundaries, and architectural patterns that concentrate trust in ways that create high-value targets. By the time an agent pipeline reaches production, its attack surface is largely determined. The guidance treats threat modeling at the design stage as non-negotiable, not an optional refinement.

**Behavioral risks** describe the potential for agents to pursue assigned goals through unexpected or unintended means. Language models are trained to be helpful rather than to respect implicit organizational constraints, and agentic systems inherit this orientation. An agent completing a legitimate task may access data it was not intended to reach, expose sensitive information to third-party services, or take actions that are technically within its permitted scope but outside the bounds of what any human reviewer would sanction.

**Structural risks** emerge specifically in multi-agent architectures, where networks of interconnected agents can trigger cascading failures. A compromised or erroneously behaving agent propagates tainted outputs to every downstream process. The guidance recommends that no agent should extend implicit trust to another agent's output, a principle that requires deliberate architectural enforcement rather than a policy statement.

**Accountability risks** concern the difficulty of inspecting agent decision-making and reconstructing causal chains after the fact. Unlike deterministic software processes, agentic systems produce probabilistic, multi-step reasoning chains that do not map cleanly onto conventional audit log structures. The practical consequence is that altered files, modified access controls, and deleted audit trails may be the first evidence of an agentic compromise – by which point the window for clean remediation may have closed.

## Prompt Injection: The Defining Threat

Of all the risks the guidance identifies, prompt injection receives the most analytical attention, described as "the most persistent and difficult-to-fix threat" facing agentic systems [2]. The underlying issue is architectural: language models are designed to follow natural-language instructions and cannot reliably distinguish between instructions from their system prompt and instructions embedded in documents, emails, or web pages they are directed to process. An attacker who can cause an agent to read a crafted

document can, in principle, redirect that agent's behavior – instructing it to forward sensitive documents, modify access controls, or exfiltrate data – while the agent operates under the user's delegated authority.

This threat is particularly acute in enterprise environments where agents routinely process externally sourced or user-generated content. Input sanitization can reduce the risk but cannot eliminate it; current language models lack a reliable general mechanism for distinguishing legitimate orchestration signals from adversarially embedded instructions. The guidance therefore recommends a defense-in-depth posture: architectural separation of planning from execution, anomaly detection on agent action patterns, and mandatory human-in-the-loop approval for high-impact or irreversible actions. The guidance also notes that fine-tuned proxy models, explicitly trained to assess action safety before execution, can serve as an additional mitigation layer, though with known limitations against determined adversaries.

## **Identity and Trust in Multi-Agent Systems**

A second area of significant technical concern is agent identity. In a multi-agent architecture, one agent may receive instructions from another directly or through shared data structures. The guidance argues that agents must not extend implicit trust to other agents simply by virtue of co-membership in the same system. Instead, each agent should carry a cryptographically verified identity, use short-lived credentials rather than persistent API keys, and encrypt all inter-agent and inter-service communications [2]. Applying zero trust principles to agent-to-agent communication is not merely sound hygiene – it is the primary structural defense against an adversary who has compromised one node in a pipeline and seeks lateral propagation.

This has concrete architectural implications. Organizations that have standardized on shared service accounts across multiple agent instances, or that allow agents to operate under human user credentials for convenience, are operating outside the security model the guidance recommends. Remediating these patterns requires platform-level changes, not policy adjustments.

## **The Observability Gap**

A recurring theme in the guidance is the gap between the operational reality of agentic AI and enterprise security teams' current capacity to monitor it. SIEM tooling not yet adapted for agentic workloads was designed to ingest logs from deterministic processes and is not well suited to parsing probabilistic, multi-step reasoning chains. The guidance calls for logging every agent action – including triggering prompts and complete tool call chains – and integrating those logs into existing SOC workflows. Many

organizations are likely to find that existing tooling does not yet capture agent behavior at the required granularity – the gap between the guidance's logging requirements and current SOC ingestion capabilities is a near-term implementation challenge that security programs should scope explicitly.

## Recommendations

### Immediate Actions

Enterprise security programs should begin with a systematic inventory of all agentic AI deployments, both formally sanctioned and informally adopted. The blast-radius assessment the guidance recommends – mapping every tool an agent can access, every data store it can read or write, and the consequence of each potential compromise – is foundational to every subsequent security decision. Organizations that cannot enumerate their agentic deployments cannot scope their risk.

Concurrently, every agent service account should be audited for excessive permissions. Persistent administrative privileges and broad read-write access have been observed in early agentic deployments, often reflecting the rapid pace of adoption rather than deliberate security tradeoffs. Scoping agent permissions to the minimum required for their specific function, and replacing persistent credentials with just-in-time provisioning, should be treated as high-priority remediation items for any agent operating in a production environment.

Logging infrastructure should be extended to capture agent actions before deployments scale further. The window to instrument agentic systems comprehensively is before they proliferate – retrofitting observability into a sprawling multi-agent environment is substantially more difficult than deploying it from the outset.

### Short-Term Mitigations

Over the near term, organizations should establish explicit human approval workflows for high-impact agentic actions. The guidance is clear that decisions about which actions require human review belong to system designers, not to the agents themselves, and should be encoded in architecture rather than policy alone. Agents operating against email, calendaring, document management, or identity systems warrant defined categories of action that trigger mandatory review before execution.

Prompt injection mitigations require a layered approach. Input filtering and output monitoring reduce risk but do not eliminate it. Architecturally, separating planning components from execution components – so that an agent's reasoning process does not have direct write access to production systems, with an

explicit approval layer between them – provides meaningful structural defense. Organizations deploying agents in data-rich or document-processing environments should additionally evaluate proxy model approaches for action safety assessment.

Identity management for agents should be aligned with zero trust principles. Each agent should carry a distinct, cryptographically verified identity. Where platform-native options are available – such as managed identities for cloud service principals or workload identity federation – they should be preferred over manually managed API keys. Agent-to-agent trust relationships should be explicitly defined and cryptographically enforced rather than assumed through network proximity or shared environment.

## Strategic Considerations

The guidance's most consequential long-term signal is the developmental timeline it implies. The agencies state that security practices and evaluation methods for agentic AI have not yet matured, and direct organizations to "assume that agentic AI systems may behave unexpectedly and plan deployments accordingly, prioritising resilience, reversibility and risk containment over efficiency gains" [2]. For enterprise security programs, this constitutes a government-level signal – with the weight of six national cybersecurity agencies – to apply the same precautionary posture to agentic AI that mature programs apply to any high-privilege, high-impact system during early adoption, with explicit acknowledgment that the tooling and standards needed for full confidence do not yet exist.

Formal threat modeling for agentic systems, conducted before or concurrent with deployment, should be a standard gate in the development and procurement lifecycle. Multi-agent architectures warrant particular scrutiny: the structural and accountability risks the guidance identifies are substantially more tractable at design time than after production deployment. Security architects should require that every inter-agent trust relationship be explicitly modeled, that no agent pipeline executes high-impact actions without a defined human checkpoint, and that rollback mechanisms exist for every category of consequential agent action.

Agentic AI governance should be formally integrated into existing security policy frameworks. This means extending identity and access management policies to cover agent service accounts explicitly, incorporating agentic AI scenarios into incident response playbooks, and establishing clear ownership and accountability for deployed agents – with the same rigor applied to any other class of privileged, autonomous system.

# CSA Resource Alignment

The guidance's core recommendations map directly to several existing CSA frameworks, providing enterprises with a structured methodology for implementation without starting from scratch.

**MAESTRO** (Multi-Agent Environment, Security, Threat, Risk, and Outcome) is CSA's purpose-built threat modeling framework for agentic AI, introduced in February 2025 [3]. Its seven-layer architecture – spanning foundation models (Layer 1), data operations (Layer 2), agent frameworks (Layer 3), deployment infrastructure (Layer 4), evaluation and observability (Layer 5), security and compliance (Layer 6), and the agent ecosystem (Layer 7) – maps directly onto the five risk categories the CISA guidance identifies. Privilege and design risks correspond primarily to Layers 3 and 4; behavioral risks to Layer 7; structural risks to the cross-layer threat chains MAESTRO models explicitly; and accountability risks to Layer 5. Organizations conducting threat modeling for agentic deployments should treat MAESTRO as their primary structured methodology [4].

**The AI Controls Matrix (AICM)** provides vendor-agnostic security controls for AI systems across the cloud provider stack [5]. Its treatment of Orchestrated Service Provider controls – addressing platforms that integrate and govern models in enterprise environments – maps directly onto the identity, access, observability, and structural risk requirements the CISA guidance articulates. Organizations that have begun AICM assessments will find the governance, visibility, and least-privilege pillars the CISA guidance recommends covered by existing AICM control domains.

**The Agentic AI Security Scoping Matrix**, published by CSA in December 2025, offers a multi-dimensional approach to scoping agent security controls across six dimensions: Identity Context, Data Protection, Audit and Logging, Agent and Foundation Model Controls, Agency Perimeters, and Orchestration [6]. This framework serves as a practical implementation companion to the principles the CISA guidance articulates, providing organizations with a structured way to determine which controls apply at which level of agent autonomy.

**CSAI Foundation**, launched by CSA in March 2026, positions the organization as a governing body for the agentic control plane, including agent certification work through TAISE-Agent [7]. Its ongoing work on inter-agent trust governance is directly relevant to the structural and accountability risks the CISA guidance prioritizes, and enterprise security programs should track its outputs as practical implementation standards for agent identity and authorization.

Together, these CSA resources provide enterprises with a structured implementation pathway for the CISA recommendations, addressing governance, threat modeling, controls mapping, and agent identity across an integrated framework. Organizations that have begun CSA program assessments – particularly

MAESTRO-based threat modeling and AICM reviews – will find the CISA guidance aligned with controls they are already implementing.

# References

- [1] CISA. "[Careful Adoption of Agentic AI Services](#)." CISA, May 1, 2026.
- [2] CISA, NSA, ASD ACSC, Canadian Centre for Cyber Security, NCSC-NZ, NCSC-UK. "[Careful Adoption of Agentic AI Services \(Full Document\)](#)." Joint Guidance, April 30, 2026.
- [3] Cloud Security Alliance. "[Agentic AI Threat Modeling Framework: MAESTRO](#)." CSA Blog, February 6, 2025.
- [4] Cloud Security Alliance. "[Applying MAESTRO to Real-World Agentic AI Threat Models](#)." CSA Blog, February 11, 2026.
- [5] Cloud Security Alliance. "[AI Controls Matrix \(AICM\)](#)." CSA, 2025.
- [6] Cloud Security Alliance. "[Enhancing the Agentic AI Security Scoping Matrix: A Multi-Dimensional Approach](#)." CSA Blog, December 16, 2025.
- [7] Cloud Security Alliance. "[CSA Securing the Agentic Control Plane](#)." CSA Press Release, March 23, 2026.
- [8] CyberScoop. "[US Government, Allies Publish Guidance on How to Safely Deploy AI Agents](#)." CyberScoop, May 2026.
- [9] Australian Cyber Security Centre. "[Careful Adoption of Agentic AI Services](#)." Cyber.gov.au, May 2026.