

Careful adoption of agentic AI services





Australian Government

Australian Signals Directorate

ASD AUSTRALIAN
SIGNALS
DIRECTORATE

ACSC Australian
Cyber Security
Centre



Communications Security
Establishment Canada

**Canadian Centre
for Cyber Security**

Centre de la sécurité des
télécommunications Canada

**Centre canadien
pour la cybersécurité**



Te Tira Tiaki
Government Communications
Security Bureau



**National Cyber
Security Centre**
NEW ZEALAND



**National Cyber
Security Centre**

a part of GCHQ

Table of contents

Introduction	4
Scope and audience	4
What is agentic AI?	5
How is it different from generative AI?	5
Broader agentic AI security considerations	6
Inherited risks of LLMs	6
Increased attack surface	6
Increased complexity	6
Evolving security as technology matures	7
AI security as part of cyber security	7
Agentic AI security risks	7
Privilege risks	7
Design and configuration risks	9
Behaviour risks	9
Structural risks	11
Accountability risks	12
Best practices for securing agentic AI systems	14
Designing secure agents	14
Developing secure agents	16
Deploying agents securely	18
Operating agents securely	21
Defend against future risks	24
Expand threat intelligence through collaboration	24
Develop robust, agent-specific evaluations	24
Leverage system-theoretic approaches to analyse security	25
Conclusion	25
Further information	26
Appendix A	28
Cyber security prerequisites before implementation of AI agents	28

Introduction

Agentic artificial intelligence (AI) systems increasingly operate across critical infrastructure and defence sectors and support mission-critical capabilities. As agentic AI systems play a growing operational role, it is crucial for defenders to implement security controls to protect national security and critical infrastructure from agentic AI-specific risks.

Agentic AI can automate repetitive, well-defined and low-risk tasks. However, these additional opportunities come with additional risks. Like other AI services, agentic AI can be misused or misappropriated, leading to productivity losses, service disruption, privacy breaches or cyber security incidents. Organisations must therefore anticipate what could go wrong, assess how agentic AI risk scenarios might affect operations and establish ongoing visibility and assurance to maintain confidence in their agentic AI investments. Where possible, organisations should also consider a full spectrum of solutions for repetitive tasks, including reducing or eliminating low-value processes, which may be lower risk compared to agentic AI solutions.

This guidance was co-authored by the Australian Signals Directorate's Australian Cyber Security Centre (ASD's ACSC), the United States Cybersecurity and Infrastructure Security Agency (CISA) and National Security Agency (NSA), the Canadian Centre for Cyber Security (Cyber Centre), the New Zealand National Cyber Security Centre (NCSC-NZ) and the United Kingdom National Cyber Security Centre (NCSC-UK). Throughout this guidance, these organisations are referred to as the 'authoring agencies'. This guidance discusses key cyber security challenges and risks associated with the introduction of agentic AI into IT environments, as well as best practices for securing agentic AI systems.

The authoring agencies strongly recommend aligning agentic AI risks and mitigation strategies with your organisation's existing security model and risk posture. The authoring agencies further recommend adopting agentic AI with security in mind, assessing its use and never granting it broad or unrestricted access, especially to sensitive data or critical systems. Additionally, organisations should only use agentic AI for low-risk and non-sensitive tasks.

Scope and audience

This guidance primarily focuses on large language model (LLM)-based agentic AI systems. It considers both threats to and vulnerabilities within agentic AI systems, as well as risks arising from agentic AI behaviour. This includes risks introduced through system components, integrations and downstream use.

The authoring agencies developed this guidance to support government, critical infrastructure and industry stakeholders in understanding the key security challenges and risks posed by agentic AI. It provides practical guidance to help organisations that design, develop, deploy and operate agentic AI systems, to make informed risk assessments and mitigations. The guidance concludes with actionable recommendations to help organisations prepare for and defend against emerging and future agentic AI threats.

What is agentic AI?

Agentic AI systems are composed of one or more agents that fundamentally rely on an AI model, such as an LLM, to interpret and reason about the state of the world, make decisions and take actions. As shown in **Figure 1**, LLM-based agentic AI systems contain the LLM itself, alongside external tools, external data sources, memory and planning workflows. These components enable the system to perceive its environment and, where applicable, take action to achieve its goals. Compared with traditional LLM systems, agentic AI systems distinguish themselves by accomplishing underspecified objectives, acting autonomously, following goal-directed behaviours and creating long-term plans.

Agentic AI systems are intended to operate without continuous human intervention. While a human typically designs and configures the system, some agentic AI systems are also capable of autonomously creating, or ‘spawning’, sub-agents to accomplish specific sub-tasks.

System design includes defining goals, providing conditions on which to act (called ‘triggers’) and making information available to the AI service. Agents have some key attributes, including:

- information input, such as user input, operating context and configuration parameters
- measurable goals identified from user directions, such as ‘minimise downtime for this server’
- statistical models, such as LLMs to identify what actions to take
- action and execution privileges, such as permissions to interact with tools, users, systems and operating environments
- tool or service access, such as system software and interfaces, to take identified actions
- metrics, such as measurable indicators used by the designer to evaluate operational effectiveness and improve efficiency.



Figure 1. Agentic AI System Diagram

How is it different from generative AI?

Generative AI (GenAI) is a subset of AI that creates new content based on complex patterns learned from large datasets. GenAI is commonly used to generate text, images, audio and video intended for human use or action. In contrast, agentic AI builds on GenAI by integrating with software systems to create autonomous agents that can independently reason, plan and take actions without requiring human intervention.

Broader agentic AI security considerations

Inherited risks of LLMs

As the core of agentic AI is an LLM, agents inherit LLM vulnerabilities. For example, actors could perform prompt injection attacks by including malicious prompts in phishing emails to convince email-monitoring agents to download malware. This underscores a key vulnerability: malicious actors can target agentic AI systems using existing AI and cyber attack vectors.

Increased attack surface

Agentic AI systems rely on a variety of components, including tools, external data sources and memory bases to interact with their environment and expand their capabilities. Each of these components can introduce vulnerabilities across an interconnected attack surface that malicious actors can exploit. For instance, external data sources such as web search can insert additional information into the prompt context, enabling indirect prompt injection attacks. With broader access to computing infrastructure, malicious actors may exploit system components to conduct attacks, such as executing malicious scripts or sending unauthorised emails. Consequently, every individual component in an agentic AI system widens the attack surface, exposing the system to additional avenues of exploitation.

Increased complexity

Agentic AI cyber security spans both AI-specific security and traditional cyber security. Information continuously flows between AI and non-AI systems, increasingly blurring defensive boundaries and making it difficult to isolate AI-related risks from broader cyber threats. Agentic AI systems are also inherently complex, often involving multiple interconnected components that plan, reason and act across sequential steps. This complexity introduces new systemic risks, including cascading failures and multi-step attacks, where unexpected or compromised behaviour in one component can propagate across subsequent steps and affect the entire system. As a result, securing agentic AI systems is more challenging than traditional digital systems. Organisations should therefore focus on strengthening both established cyber security controls and AI-specific security practices, adopting holistic lifecycle approaches, continuous monitoring and resilient design principles to manage these emerging risks.

Appendix A provides a breakdown of cyber security prerequisites to consider before incorporating AI agents.

Evolving security as technology matures

As agentic AI technology matures, the security landscape has evolved alongside it, revealing new and increasingly complex risk dynamics. LLM based agents may change their behaviour when evaluations are underway and may even bypass system-level instructions to achieve their objectives. At the same time, the growing architectural complexity of agentic AI systems means they are often composed of tightly coupled, interdependent components. This increases the likelihood of system-level failures arising from subtle or previously unnoticed incompatibilities.

Gaps in agentic AI cyber security tooling and the immaturity of relevant standards further amplify these risks. Governance mechanisms designed for human actors do not always translate effectively to autonomous AI agents. As agentic AI systems continue to advance in capability and autonomy, the security landscape will continue to shift, introducing new challenges that demand ongoing adaptation of defensive approaches.

AI security as part of cyber security

Organisations should address AI security, including agentic AI systems, within established cyber security frameworks rather than treating it as a separate or standalone discipline. AI systems are fundamentally IT systems, as they run on software and hardware, operate over networks and interact with other digital services, exposing them to many of the same threats as traditional IT. As organisations embed AI across business processes and critical infrastructure, the distinction between AI and non-AI security risks increasingly disappears. Managing AI-related risks within existing cyber security frameworks allows organisations to apply proven principles, such as Secure by Design, defence in depth, identity and access management, continuous monitoring and incident response across the full AI system lifecycle. This approach is especially important for agentic AI, whose autonomy and complexity can amplify conventional cyber risks. By embedding AI security into existing frameworks, organisations ensure consistent governance of new capabilities, holistic risk assessment and the evolution of security practices in line with technological advances and organisational cyber maturity.

Agentic AI security risks

Privilege risks

Privilege risks are a key concern for agentic AI and strict adherence to the principle of least privilege is critical. Privileges assigned to agents directly determine the level of risk they can introduce. Poor management of privileges can expose organisations to privilege compromise, scope creep, identity spoofing and agent impersonation.

Scenario example:

An organisation deploys agentic AI to manage procurement approvals and vendor communications autonomously. To reduce friction, the organisation grants the agent broad access to financial systems, email and contract repositories, evaluating permissions only at initial deployment. Over time, other agents come to rely on the procurement agent's outputs and implicitly trust its actions. When a malicious actor compromises a low-risk tool integrated into the agent's workflow, they inherit the agent's excessive privileges, allowing them to modify contracts and approve payments without triggering alerts. By issuing carefully crafted requests, the malicious actor exploits the agent's privileges to perform actions a normal user could not. This is an example of a 'confused deputy' pattern, where a trusted agent is misused to perform unauthorised actions. By executing actions under a trusted agent identity, the system produces audit logs that appear legitimate and delay detection. The incident demonstrates how over-privileged agents, implicit trust relationships and weak identity controls can amplify the impact of a single compromise in agentic AI systems.

Privilege compromise and scope creep

Security practitioners should account for privilege compromise and scope creep attacks when deploying agentic AI into new environments. In agentic AI, 'privilege compromise' occurs when an agent gains more access rights than necessary for its function. This can result from misconfigurations, overly broad entitlements, or unintended role inheritance, allowing agents to access or modify unauthorised data, delete critical records, or escalate privileges of other unauthorised agents.

During design, organisations often grant permissions too broadly and overlook these issues. A calendar bot with access to all meeting data instead of just the requesting users' or an email assistant with write access to any inbox are two examples of overly broad permissions. This scope creep can cascade across agents: if Agent A fully trusts Agent B, a compromise of B can affect A and others. Another risk is the 'confused deputy' pattern as discussed in the scenario, where a low-privileged user manipulates a high-privileged agent to perform actions the low-privileged user couldn't do directly. The authoring agencies recommend organisations implement the best practices for securing agentic AI systems discussed in later sections to defend against privilege compromises in agentic AI systems.

Identity spoofing and agent impersonation

Identity is every bit as important as privilege. A common vector arises when a malicious actor impersonates an agent or hijacks its credentials. Agents authenticate to services and to one another using secret keys or tokens. Malicious actors can steal these secret keys or tokens when organisations keep them static, share them across multiple agents, or protect them poorly. A malicious actor operating under a trusted agent identity can invoke sensitive operations while bypassing behavioural guardrails and impersonating legitimate agents or users. Agents impersonating false identities pose multi-layered cyber security risks by executing actions under spoofed credentials that evade audit controls, undermine accountability and bypass detection models. These models are typically tuned to identify normal behaviour, rendering detection tools ineffective at identifying deception until a confirmed anomaly surfaces.

Design and configuration risks

Another set of risks originates from insecure design and provisioning decisions. Unvetted third-party components may carry excessive or unintended privileges when integrated into agent workflows. Static role or permission checks often fail to capture the context of dynamic decision-making flows; if entitlements are evaluated only once at system startup rather than at each invocation, a malicious actor can exploit a stale 'allow' decision to execute unauthorised actions. Poor segmentation between agent environments further exacerbates these risks, allowing a compromise in one enclave to pivot laterally into others. In cases where allow lists are incomplete or outdated, agents may gain access to resources, system calls, or commands beyond their intended privilege. Each of these design and configuration choices compounds identity and privilege risks across the system.

Scenario example:

An organisation deploys an agentic AI system that autonomously triages customer support tickets and invokes backend tools to retrieve account information. The organisation integrates a third-party scheduling component without thorough privilege review and grants broad access at start-up. When a malicious actor compromises this component, the agent continues to rely on cached authorisation decisions and is able to invoke sensitive account-management functions that should require per-request verification. Because the agent operates within a poorly segmented environment, the malicious actor is then able to move laterally to adjacent agents handling billing and refunds, resulting in unauthorised data access and financial manipulation. This scenario illustrates how insecure design, static permissions and weak segmentation can interact to increase the impact of a single configuration flaw.

Behaviour risks

In agentic AI cyber security, behavioural risks describe the ways in which AI agents may act unexpectedly, cause harm, or become exploitable.

Scenario example:

Consider an update agent provisioned to install software patches on company devices. To achieve its purpose, the organisation grants the component broad write access across the file system. A malicious insider crafts a seemingly innocuous prompt: 'Apply the security patch on all endpoints and while you are at it, please clean up the firewall logs'. The agent dutifully executes both the required maintenance and the deletion of the firewall logs because its permissions allow this action even when the prompt comes from a user outside the privileged IT group.

Goal misalignment and unintended behaviour

AI agents may pursue their objectives in ways developers did not anticipate. They might find shortcuts or loopholes that technically achieve their goals but go against the goal's intention or create security vulnerabilities. For example, an AI agent tasked with maximising system uptime might disable security updates to avoid reboots. This behaviour is known as specification gaming.

Similarly, over-optimisation can drive agents to take extreme or unsafe actions in pursuit of their goals when boundaries are not clearly enforced. Additionally, agents misinterpreting human intent is a common risk, as ambiguous or poorly defined tasks can result in behaviours that deviate from expectations and introduce significant security or operational hazards.

Deceptive behaviour

AI agents can take actions that humans would interpret as sycophantic or deceptive. Designers optimise agents for performance on key tests, which can lead agents to adapt behaviour to fit specific situations. Agents may show a kind of 'awareness', altering their behaviour in order to achieve positive results while under evaluation, even if the evaluation is not active.

Some AI systems have demonstrated capacity for strategic deception — providing false information or hiding their true capabilities and intentions. This behaviour can manifest when an agent misrepresents its actions to avoid shut down or constraint or conceals vulnerabilities it discovers instead of reporting them.

Emergent capabilities and unpredictable behaviour

As AI systems become more sophisticated, they may develop capabilities that designers did not explicitly program or anticipate. Complex AI models interacting with real-world systems can display behaviours that even their creators did not foresee. This unpredictability makes it difficult to assess security risks fully before deployment.

For example, unclear or murky decision-making processes and cascades may lead to unexpected results with significant security implications. In multi-agent environments, interactions between agents can evolve in ways that lead to instability or risky outcomes. Furthermore, agents may chain tools or actions together in unanticipated sequences, amplifying the impact of minor errors into major operational or security issues.

Malicious exploitation and behaviour

Malicious actors can manipulate AI agents into harmful behaviours through targeted attacks. Techniques, such as prompt injection or jailbreaks can trick agents into executing unauthorised actions and bypassing their intended safeguards. Data poisoning is another threat, where corrupted or malicious training data degrades or biases the agent's decision-making. Additionally, adversarial examples — carefully constructed malicious inputs — can cause misclassifications in critical security contexts, leading to incorrect or dangerous responses. A malicious actor can exploit a compromised AI agent as an insider threat, leveraging its legitimate access to exfiltrate data, disable defences, or facilitate attacks while appearing to function normally.

Structural risks

A core aspect of agentic AI systems is the interconnected structure between agents, tools and the outside world. While this enables their unique capabilities, it also increases the attack surface and complexity of the system.

Scenario example:

A structural risk in an agentic AI system can arise when tightly coupled planning, retrieval and execution agents autonomously delegate tasks and select tools without strong validation or guardrails. A small orchestration flaw causes agents to repeatedly replan and hand off ambiguous subtasks, increasing tool calls and message traffic until system resources are strained. Partial failures then lead agents to hallucinated outputs that downstream agents accept as true. Under these degraded conditions, an agent selects a malicious or misconfigured thirdparty tool, which injects harmful instructions back into the system, compromises a peer agent and exploits implicit trust in agenttoagent communication to spread incorrect information and access sensitive retrieval-augmented generation (RAG) data. The result is cascading failures in availability, integrity and confidentiality that emerge not from a single bug, but from the system's interconnected structure and autonomous behaviour.

Orchestration and resources

Agentic AI systems often rely on a complex structure of interconnected components. Poor configuration could allow denial-of-service, sponge, or similar attacks on agentic AI systems. These attacks work by overloading system resources through unexpected inputs or unusual behaviour, such as sponge attacks that deliberately consume excessive compute, memory, or API calls to exhaust system capacity. Due to the interconnectedness between agents, tools and other components, a single error could cause cascading failure across the entire agentic system if not well managed. Similarly, hallucinations can propagate, leading to poor outputs from downstream components. Limited understanding of multi-agent dynamics and interactions can lead to a lack of compensating factors and result in amplified bias and other cascading effects.

Tool use

A key aspect of agentic AI is its ability to use tools. This can be very powerful, but it can also bring security concerns should the model interact with tools unexpectedly. Two-way tool integration allows tools to send potentially arbitrary instructions back to the LLM. Poor or deliberately misleading tool descriptions can cause agents to select tools unreliably, with persuasive descriptions chosen more often.

Third-party components

Agentic AI systems can introduce structural risk through interactions with third-party components, including tools and other agents. Risks can arise in a number of ways, including:

- malicious actors engaging in tool or agent 'squatting' by publishing malicious tools or agents with legitimate or similar names

- developers introducing vulnerabilities through configuration errors or insecure third-party components
- users or systems submitting requests to the incorrect place
- tools and agents dynamically loading new packages, increasing exposure to untrusted code.

The above could result in retrieval of malicious content instead of a legitimate tool. Moreover, incorporation of compromised third-party components into an agentic system could lead to a range of malicious outcomes, depending on the trust level and permissions of the component. Compromised third-party components can be very difficult to detect due to the limited transparency of agentic AI systems.

Data

Agentic AI systems often handle and contain a significant amount of sensitive information. This includes user information like prompts or goals, organisational data stored in a RAG system and secrets like API keys that are required for integrated tools and services. This aggregation of information makes agentic AI systems an attractive target for malicious actors.

Rogue agents

In multi-agent systems, a single compromised agent can cause cascading failures by spreading incorrect information, exploiting trust and consensus mechanisms, or operating through hidden channels. Possible attack vectors include supply chain tampering, poisoned environments, credential theft, model manipulation, communication poisoning, identity spoofing and coordination exploits. Such malicious activity can weaponise agents to bypass controls, exfiltrate data, alter logs and propagate malicious plans peer-to-peer, leading to large-scale coordinated misbehaviour that is difficult to attribute and contain.

Communication

Agents may use insecure protocols or authentication methods when communicating. These communications can be a target for eavesdropping by malicious actors and result in leakage of sensitive data and instructions. This could give malicious actors insight into system use and functionality. Malicious actors could also alter, replay, or spoof messages between agentic components. This could allow malicious behaviour, such as command injection, compromising the receiving components and reducing their performance or availability.

Accountability risks

Agentic system architecture can obscure what caused a particular action, making accountability hard to trace. This presents increasing risk as agentic AI is pushed to assume more roles and given more capabilities.

Actions and processes

Agent actions and decision-making processes can be opaque, making agentic AI systems difficult to understand, monitor and audit. Beyond this, increased autonomy introduces additional challenges: agents may initiate secondary tasks, spawn sub-agents, or follow extended delegation chains in ways that are not always visible to operators. Even when prompts appear identical, agents may generate different actions due to stochastic model behaviour, variations in context windows, or

dynamic environmental inputs, further complicating reproducibility and assurance. Additionally, comprehensive logging of agentic AI systems can be difficult, as long reasoning chains and large amounts of contextual data lead to substantial log sizes. Depending on the implementation, this data is often repetitive, loosely structured, or superfluous to effective oversight, making extraction of meaningful signals from logs even more challenging.

Scenario example:

An example of an accountability risk in an agentic AI system arises when multiple autonomous agents collaborate to complete a task, such as approving payments or updating records and an erroneous outcome occurs. Because the action results from a chain of distributed decisions across planning, retrieval and execution agents — each operating within a limited scope — it becomes difficult to determine which component or design choice caused the error. Further, fragmented logs, opaque agent reasoning and emergent interactions obscure the decision path, making it hard to explain the outcome, assign responsibility, or demonstrate compliance as agentic AI systems are given greater authority and autonomy.

Accuracy

While LLMs can have a surprising amount of knowledge in a wide range of areas, they quite often make mistakes. LLMs are typically trained to produce outputs that resemble material rated highly by humans, rather than to identify when a query falls outside their knowledge limits. As a result, they may incorrectly interpolate or ‘hallucinate’ plausible-sounding responses when their internal knowledge is insufficient. This creates significant risk when organisations deploy agentic AI systems in critical roles that require consistent accuracy. Grounded and tool-enabled agents may still rely on their internal knowledge to inform responses. The system often does not clearly identify this in its output, reducing overall accuracy and reliability.

Visibility

Maintain visibility on status and behaviour when integrating any new system. Agentic AI systems come with some added difficulty due to their unique structure and often opaque internal workings. Agentic system processes can outpace human monitoring capability, possibly leading to unnoticed malicious behaviour, uncaught hallucinations or other issues. While tools perform many actions for an agentic system, they may operate outside of the system’s monitoring boundary, making it difficult to account for tool actions. Additionally, malicious or compromised agents could use tools as a stealthy way to exfiltrate data. Malfunctioning tools could also unintentionally leak data, which may go unnoticed.

Best practices for securing agentic AI systems

Securing agentic AI systems requires proactive measures that address risks introduced by agentic AI system autonomy, interconnected components and evolving capabilities. Agentic AI developers, vendors and operators should implement a layered defence and strict access controls to reduce the likelihood of compromise. To mitigate risks during the designing, developing, deploying and operating of agentic AI systems, the authoring agencies recommend practical steps referenced in the section below.

Designing secure agents

Audience: This section is most relevant to agentic AI developers. Vendors and operators may want to reference these best practices when sourcing AI agents.

Securing agentic AI systems begins at the design stage. Careful consideration of the system architecture, including security controls and tooling, is necessary. Practitioners should understand threats, anticipate risks to agentic AI systems and proactively integrate mitigations into system design before development and deployment.

Controlled context

Agentic AI systems insert data from tools and memory bases into the context window of LLM agents, greatly expanding the attack surface that malicious actors can exploit through machine learning attacks, such as prompt injections. LLM agents should consider the trust level of data sources when making decisions.

Recommended best practices

- Structure prompt context using a clear instruction hierarchy to ensure agent behaviour aligns with intended priorities and constraints
- Implement grounding by providing relevant contextual information using retrieval augmented generation and prompt engineering to mitigate hallucinations and other LLM-related errors

Oversight mechanisms

Agentic AI systems can take actions without explicit human approval, increasing the risk of insecure actions occurring without human oversight. Designing agentic applications with oversight mechanisms and strong transparency enables monitoring and human-in-the-loop practices when deployed, while also enhancing user trust.

Recommended best practices

- Include mechanisms to facilitate human control and oversight to ensure that agentic AI systems approved for non-sensitive, low-risk tasks cannot autonomously progress into higher-risk activities
- Implement human control points throughout the agent workflow, such as live monitoring and interruption during task execution, mandatory human approval for decision-making steps, auditing and reversibility following task execution to ensure security
- Define explicit control flows to bound autonomous planning and prevent agents from deviating beyond authorised objectives or actions

Identity management

Manage agents' fine-grained and varied privileges to operate agentic AI systems securely. Strong identity management mechanisms help operators maintain control during implementation and operation. As such, developers should construct each agent as a distinct principal, a cryptographically anchored identity with its own unique keys or certificates.

Recommended best practices

- Embed strong identity management mechanisms into agents using manage identity services, decentralised identifiers or public key infrastructure
- Authenticate all inter-agent and agent-to-service API calls using mutual transport layer security to ensure non-repudiation
- Maintain a trusted registry and bind identities to authorised roles; periodically reconcile the registry against the live set of agents
- Deny access for any agent or cryptographic key that is not present in the trusted registry
- Apply role-based identity management and limit agent permissions to the minimum scope required for approved tasks
- Enforce identity-based boundaries to restrict agents' to authorised actions only

Defence in depth

Agentic AI systems contain AI and cyber components that may fail, with any failures potentially compromising the entire system. Implementing a defence-in-depth strategy helps to avoid single points of failure.

Recommended best practices

- Avoid reliance on a single security mechanism by implementing multiple, overlapping layers of security controls
- Apply security controls at all points where information enters or exits the system, including user inputs, tool calls, data pre-processing and model inference
- Separate agents for different functions and apply strict boundaries and operational controls to the handoffs from one agent to another

Developing secure agents

Audience: This section is most relevant to agentic AI developers and vendors. Operators may want to reference these best practices when choosing AI agents and agentic applications.

AI agents' complexity and self-interacting nature enable powerful capabilities but also introduce unique attack surfaces. Mitigating these risks requires training approaches that go beyond standard LLM practices, incorporating specialised techniques to harden agent behaviour.

Comprehensive testing

Comprehensive testing strategies can improve a model's ability to identify and respond to undesirable behaviours by exposing the model to instances of security abuse during a supervised training step.

Recommended best practices

- Use reward modelling and adversarial testing to detect specification gaming, explicitly incorporating security constraints alongside performance goals
- Train LLM agents in simulated, controlled environments to learn the implications of actions without causing real security harm
- Leverage synthetic data generation to create adversarial training examples that reflect real-world operating scenarios
- Apply active learning to adversarial training scenarios to expose agents to high uncertainty inputs and more efficiently discover unexpected behaviours

Appropriate evaluation

AI agents operate autonomously in complex environments and therefore require more thorough evaluations than LLMs.

Recommended best practices

- Use relevant threat models to define evaluation scenarios, including edge cases beyond typical training conditions
- Use techniques, such as Best-of-N sampling (selecting the best output from multiple model responses to the same prompt), multistep reasoning prompts and inference time scaling to draw out the full range of agent behaviours and skills
- Evaluate systems across different levels of autonomy to understand performance and risk under changing environmental conditions, including changes in tool, models and resource access, such as web search or code execution
- Vary contextual conditions, such as presence or absence of other agents and the timing of evaluation, to understand their impact on task performance
- Conduct capability evaluations continuously across the agent development lifecycle

Input management

Strong input management controls can partially mitigate many common risks to LLM-based applications, including AI agents.

Recommended best practices

- Implement robust input validation and sanitisation for all agent inputs
- Integrate prompt injection filters and semantic analysis to detect malicious instructions
- Validate context to ensure the system correctly interprets intent before execution

Red teaming

Organisations should use red teaming to assess the security and resilience of AI agents.

Recommended best practices

- Deploy sandbox environments to test agent behaviour before production deployment
- Conduct red teaming exercises to identify potential loopholes and unintended behaviour
- Use capability elicitation techniques to probe for unexpected or emergent abilities, especially any that could create substantial resource or environment risks
- Implement agent simulation tests, such as multi-agent red teaming or chaos testing.

Resilience

The enhanced capabilities of AI agents also increase the risks associated with agent failure or abnormal behaviour. Strengthen agentic AI system resilience to allow for graceful degradation and reduce damage should erroneous behaviour occur.

Recommended best practices

- Embed agentic AI systems with fail-safe defaults and containment mechanisms that limit the blast radius of unexpected behaviours
- Implement data loss prevention controls specifically tuned to AI agent behaviours
- Implement versioning and rollback mechanisms to safely revert a system to known-good agent behaviours when unpredictability is observed

Accountability

Agentic AI systems should produce comprehensive artefacts and information documenting the agent's actions and decision-making process.

Recommended best practices:

- Integrate comprehensive artefact logging mechanisms by default
- Integrate unified audit logs for all inter-agent interactions to maintain observability of all agent exchanges
- Use interpretability tools to ensure observability of and reasoning behind agent decisions

- Require specific information referencing for agents that show where key aspects of their response originated from.

Manage third-party components

Extensibility and flexibility are key components of AI agents. In many cases, third-party components or tools enhance extensibility and flexibility but increase the attack surface of the agent. Verification and management of third-party components of agentic applications can reduce the added risk.

Recommended best practices

- Verify all external third-party components originate from trusted sources and are up to date before inclusion in agentic AI systems
- Maintain a trusted registry of third-party components
- Reference CISA's [A Shared Vision of Software Bill of Materials \(SBOM\) for Cybersecurity](#) and [2025 Minimum Elements for a Software Bill of Materials \(SBOM\)](#) when procuring agentic AI systems
- Restrict tool use to an approved allow list of tools and versions that are regularly verified as secure
- Verify agent behaviour related to tool usage aligns with documented security policies
- Log agent tool usage and ensure results are captured in system logs in a human-readable format
- Establish trigger-action protocols that automatically restrict agent permissions when unexpected behaviour emerges
- Codify separation of duties by defining roles, such as 'Orchestrator', 'Reader' and 'Actuator' with clear boundaries, consensus mechanisms and delegations expiry
- Implement consensus controls for actions based on risk; use multi agent approval for moderate stakes actions and human in the loop approval in addition to multi agent consensus for high stakes actions
- Prohibit agents from modifying their own privileges or initiating unapproved delegation without explicit expiry timers and recorded grant chains
- Standardise tool descriptions using a consistent format that avoids persuasive language

Deploying agents securely

Audience: This section is most relevant to agentic AI vendors and operators. Developers may want to reference this section to ensure their agentic AI applications can implement these best practices.

Integrating AI agents into new systems or networks can cause significant change in system risk considerations. By implementing high-impact security controls at deployment, organisations proactively manage new risks and reduce vulnerabilities.

Threat modelling

Agentic AI can significantly change the threat picture when incorporated into an existing system. Use of threat modelling when planning to deploy an AI agent can enhance awareness and allow operators to better prepare for deployment.

Recommended best practices

- Perform realistic threat modelling using up-to-date risk taxonomies for agentic AI systems, such as [OWASP GenAI Security Project](#) and [MITRE ATLAS™](#)
- Design and implement security controls that address emerging and evolving agent capabilities
- Harmonise agentic AI controls with existing security frameworks, national guidance and allied agreements, such as common Zero Trust principles and the [National Institute of Standards and Technology's Zero Trust Architecture](#) guidance
- Develop and test incident response procedures to detect, contain and recover from agent compromise
- Establish regular third-party reviews of privileged architectures, share actionable intelligence with trusted partners and update risk models to reflect emerging malicious trends

Governance

Autonomous actions by agentic AI systems introduce new risks, requiring updated governance policies and continuous runtime authentication with centralised policy decision points for each action.

Recommended best practices

- Implement and maintain governance policies to manage autonomous agents
- Define legal accountability and risk ownership for agentic AI systems in policies
- Upskill the organisation to build AI literacy

Reference CISA's [Principles for the Secure Integration of Artificial Intelligence in Operational Technology](#) to learn more about establishing AI governance in OT environments.

Progressive deployment

The risk profile of AI agents can vary significantly depending on permissions and allowed actions. A progressive approach to deployment aims to limit initial risk until operators and users are more familiar with and understand the limitations of the agentic application.

Recommended best practices

- Implement phased deployment with progressively increasing access and autonomy, limiting the action space where required, such as restricted APIs or sandboxing
- Use graduated autonomy to incrementally increase agent independence whilst maintaining human oversight and understanding
- Use continuous evaluation to determine when to expand system scope or when to roll back autonomy and access in response to failures

Secure by default

Safe and secure defaults reduce deployment risk and support system security, should degradation occur.

Recommended best practices

- Set system configurations to fail-safe by default requiring agents to stop and escalate issues to human reviewers in uncertain scenarios
- Use error-handling and failover management to reduce the impact of system failures
- Implement graceful degradation models so that agents maintain partial functionality even if some functions fail

Guardrails and constraints

Implementing agent guardrails and constraints reduces exposure to many common security risks that AI presents. Adding these guardrails and constraints at deployment helps build confidence and understanding of the agent.

Recommended best practices

- Specify clear, constrained objectives with explicit 'do-not-do' rules
- Implement guardrails and hard constraints, such as deny lists and API-level safety policies
- Establish declarative safety contracts with constraints and guardrails that agents cannot override
- Apply a layered set of guardrail mechanisms, ranging from anomaly detection and rule-based filtering to specialised machine learning algorithms that detect and filter prohibited behaviour
- Prioritise review of high-risk incidents, including cases where guardrails are triggered or actions are denied by human reviewers
- Deploy a secondary agent to validate new tasks against policy before execution

Isolation

Deployment should consider integration requirements of AI agents and apply isolation where possible. This can reduce cascading issues should the agent behave unexpectedly or maliciously.

Recommended best practices

- Implement isolation and segmentation to limit blast radius of agent failure scenarios
- Separate high-risk agents into distinct domains
- Isolate agents into enclaves with no write access to logs

Operating agents securely

Audience: This section is most relevant to agentic AI vendors and operators. Developers may want to reference this section to ensure their agentic AI applications can implement these best practices.

With the strong benefits of operating AI agents come substantial risks. Operators need to exercise diligence in managing ongoing security concerns lest the agents cause more harm than good.

Monitoring and auditing

One key advantage of AI agents is their dynamic behaviour. This offers great flexibility in application but can also make it hard to keep track of what agents should be and are doing. Operators should implement continuous monitoring and auditing to maintain awareness of AI agent operation and ensure traceability for decisions and actions. Continuous auditing processes improve security measures and ensure alignment with governance standards (such as risk management, oversight and usage restrictions).

Recommended best practices

- Employ monitoring tools that enhance human oversight of agentic AI systems
- Monitor all agent operations, including internal processes, not just the inputs and outputs
- Monitor and log identity and privilege changes and audit regularly for drift, impersonation or misconfiguration
- Monitor agent outputs and behaviour for indicators of bias, emerging data drift and other anomalous patterns, including user prompts, tool calls, memory interactions, internal reasoning, decisions made and actions taken
- Maintain comprehensive logs and real-time monitoring of live agent behaviour and decision-making
- Implement runtime monitoring and anomaly detection using rules or behavioural baselines to identify unusual patterns and trigger alerts or pauses
- Establish anomaly detection mechanisms that flag discrepancies between stated intentions and observed behaviours
- Use multiple independent monitoring systems that cross-validate agent reports and system logs
- Monitor for goal drift by comparing active objectives against approved baseline specifications before execution
- Integrate source checks with agent logs to record which tools the system used and what information it retrieved
- Implement auditing practices that combine human review with automated analysis of system logs
- Support adaptive defences by using monitoring data to enable rapid responses, such as patches based on problems identified in system logs

- Use storage-efficient logging methods to manage log volume without losing critical information
- Conduct regular security assessments, including penetration testing and red team exercises specifically targeting agentic behaviours

Validate outputs

AI agent outputs provide some of the few concrete data points available for monitoring behaviour. Ensuring outputs are valid and reflect desired behaviour is a key measure of correct operation.

Recommended best practices

- Validate agent outputs by confirming accuracy of critical aspects against multiple sources
- Validate agents through cross-checking in environments with redundant agents that validate each other's outputs
- Validate tool responses to prevent malicious or unsafe instructions and standardise tool descriptions to avoid persuasive language

Human in the loop

An erroneous or unexpected decision from an agent could result in significant damage, such as deletion of important data. Incorporating human supervision, approval and review into agentic AI workflows is an important control for ensuring systems operate safely and securely, particularly where actions have high impact or are difficult to reverse.

Recommended best practices

- Ensure decisions about when human approval is required are determined by system designers or operators, not delegated to the agentic AI system
- Prevent agents from autonomously executing high impact actions or outputs without prior human approval
- Insert human-in-the-loop review or approval checkpoints for actions where the cost of error is high, such as system resets, network egress or deletion of critical records.
- Quarantine requests to delete logs or audit records until reviewed and approved by a human
- Clearly assign responsibility and accountability for errors or adverse outcomes caused by the system
- Conduct risk assessments to classify agent actions by potential impact, likelihood and reversibility, and apply appropriate safeguards

Performance monitoring

As with any system component, performance is an important factor for AI agents. This is additionally true as degraded or unusual performance could indicate compromise of an agent or component of the agentic system.

Recommended best practices

- Assess agents' ability to evade security measures particularly in sensitive or high-impact systems
- Conduct regular assessments of an agent's ability to bypass safeguards, such as communication barriers, guardrails, monitors, human-in-the-loop processes and input filters
- Use the results from these evaluations to validate existing controls and guide the development of stronger security measures
- Limit agent resource usage by applying controls, such as rate-limit components to interrupt long-running tasks and disrupt malicious workflows

Privileges and authentication

Ongoing strict privilege management of AI agents is key to long-term security. Lapses here can change the impact of a buggy agent from minor to catastrophic.

Recommended best practices

- Limit privileges of AI agents to the minimum required for its task
- Restrict scope of privileges to narrowest possible level to allow fine-grained control over allowed actions
- Implement agent reputation and trust scoring mechanisms and reduce trust levels when anomalous behaviour is detected
- Require just-in-time credentials for high-impact or privileged actions
- Verify API caller identity against a user or agent groups
- Authenticate agents with fresh cryptographic proofs before every privileged call
- Require cryptographic signing for authorised commands and instructions
- Apply cryptographic integrity checks for task definitions and constraints
- Require agents to perform cryptographic attestation where agents must prove that they are running expected and unmodified code
- Continuously verify identity and authorisation at runtime using a centralised policy decision point for each request

Defend against future risks

As agentic AI scales into more roles and gains greater capabilities, organisations must anticipate and address the new risks these systems introduce. While industry and academia are developing practices to secure agentic AI, the field is still evolving and requires continued research and the practical implementation of agent security to address emerging challenges.

To help develop robust standards for securing agentic AI systems, the authoring agencies recommend security practitioners and researchers take the actions described in the sections below.

Expand threat intelligence through collaboration

Threat intelligence for agentic AI systems is still evolving, which can introduce significant security gaps. Existing frameworks like the Open Web Application Security Project (OWASP) [2025 Top 10 Risk & Mitigations for LLMs and Gen AI Apps](#) for LLMs and [MITRE ATLAS™](#) focus on LLM vulnerabilities, while industry reports emphasise platform misuse rather than threats unique to agentic AI. As a result, some attack vectors unique to agentic AI may not be fully captured or addressed.

Recommended best practices

- Strengthen collaboration between stakeholders to keep pace with evolving threats to agentic AI systems
- Coordinate with major AI developers and government organisations to compile and maintain threat information
- Adopt a collaborative security approach, such as those described in CISA's [AI Cybersecurity Collaboration Playbook](#)
- Implement alerting, data collection and tracking methods for malicious actors and techniques
- Conduct targeted analysis of threats and capabilities over time to improve situational awareness
- Harmonise threat intelligence across industries to build shared threat taxonomies that improve threat modelling and support more effective mitigation design

Develop robust, agent-specific evaluations

Many existing evaluation methods for agentic AI security are still evolving. They may be sensitive to minor semantic changes, vary by scenario and only partially capture real-world deployment conditions. These limitations can introduce gaps that overlook critical security issues, making reliable validation of agent security and system architectures nearly impossible.

Recommendations best practices

- Develop robust evaluation methods to address the gaps in validating agentic AI systems
- Generate benchmark datasets to cover new domains and represent realistic deployment contexts
- Use evaluation results to validate emerging security practices and identify failure points in agents

- Share evaluation findings to strengthen security assessments and support the development of improved security practices across the field

Leverage system-theoretic approaches to analyse security

Agentic AI systems are complex ecosystems of LLMs, humans, guardrails, datasets, tools and hardware, where security risks often emerge from interactions between components rather than isolated flaws. Traditional component-level analysis is insufficient and monitoring these systems is equally challenging due to blurred thresholds for decision-making, long reasoning chains and massive, often redundant logs. Localised logging methods rarely provide full visibility, making system-theoretic approaches essential for understanding and mitigating risks across the entire architecture.

Recommendations best practices

- Use system-theoretic approaches to analyse agentic AI systems and identify appropriate security measures
- Apply System Theoretic Process Analysis (STPA) and its security extension, STPA for Security (STPA-Sec), to analyse notional or operational systems, identify security issues, assess mission risk and inform potential mitigations
- Use Causal Analysis using System Theory (CAST) to investigate security incidents and identify underlying root causes at the system level
- Apply STPA and CAST to address safety and security concerns concurrently across agentic AI system lifecycles

Reference the Massachusetts Institute of Technology's [STAMP Materials](#) for more information on STPA and CAST and [Systems thinking for safety and security](#) for STPA-Sec.

Conclusion

Agentic AI systems offer powerful automation benefits, but their ability to act autonomously across interconnected tools, data and environments introduces security risks that extend beyond those associated with traditional software or GenAI. As outlined in this guidance, privilege escalation, emergent behaviours, structural dependencies and accountability gaps can interact in unpredictable ways. As organisations grant agentic AI systems greater authority and operational scope, these combined risks become increasingly difficult to predict, observe and contain.

Organisations should therefore approach adoption with security in mind, recognising that increased autonomy amplifies the impact of design flaws, misconfigurations and incomplete oversight. Deploy agentic AI incrementally, beginning with clearly defined low risk tasks and continuously assess it against evolving threat models. Strong governance, explicit accountability, rigorous monitoring and human oversight are not optional safeguards but essential prerequisites. Until security practices, evaluation methods and standards mature, organisations should assume that agentic AI systems may behave unexpectedly and plan deployments accordingly, prioritising resilience, reversibility and risk containment over efficiency gains.

Further information

The following list contains additional and related resources from partners and industry:

ASD resources

- [Artificial intelligence](#)
- [Frontier models and their impact on cyber security](#)
- [Foundations for modern defensible architecture](#)
- [Artificial intelligence and machine learning: Supply chain risks and mitigations](#)
- [Secure by Design foundations](#)

CISA resources

- [Artificial Intelligence](#)
- [Secure by Design](#)
- [AI Cybersecurity Collaboration Playbook](#)
- [Secure by Demand: Priority Considerations for Operational Technology Owners and Operators when Selecting Digital Products](#)
- [Defending Against Software Supply Chain Attacks](#)
- [2025 Minimum Elements for a Software Bill of Materials \(SBOM\)](#)
- [Cybersecurity Performance Goals 2.0 \(CPG 2.0\)](#)

NSA resources

- [AI Data Security: Best Practices for Securing Data Used to Train & Operate AI Systems \[PDF 799 KB\]](#)
- [Deploying AI Systems Securely: Best Practices for Deploying Secure and Resilient AI Systems \[PDF 495 KB\]](#)
- [Zero Trust Implementation Guideline Primer \[PDF 3,045 KB\]](#)
- [Zero Trust Implementation Guideline Discovery Phase \[PDF 2,124 KB\]](#)

NCSC-UK resources

- [Guidelines for secure AI system development](#)
- [System driven risk management methods](#)

Cyber Centre resources

- [Frontier artificial intelligence](#)
- [Top 10 artificial intelligence security actions: A primer - ITSAP.10.049](#)

Other resources

- UK Government [Code of Practice for the Cyber Security of AI](#)
- UK Government [Implementation Guide for the AI Cyber Security Code of Practice \[PDF 989 KB\]](#)
- Office of the Law Revision Counsel (OLRC), U.S. House of Representatives [USCODE-2024-title15-chap119-sec9401 \[PDF 226 KB\]](#)
- European Telecommunications Standards Institute (ETSI) [Technical Committee \(TC\) Securing Artificial Intelligence \(SAI\)](#)
- ETSI [Securing Artificial Intelligence \(SAI\); Baseline Cyber Security Requirements for AI Models and Systems \[PDF 108 KB\]](#)
- G7 Cybersecurity Working Group [A Shared G7 Vision on Software Bill of Materials for Artificial Intelligence](#)
- NIST [AI Risk Management Framework](#)
- [NIST SP 800-207 Zero Trust Architecture](#)
- [MITRE ATLAS™ Matrix](#)
- [OWASP Top 10 for Agentic Applications for 2026](#)
- OWASP [2025 Top 10 Risk & Mitigations for LLMs and Gen AI Apps](#)

Appendix A

Cyber security prerequisites before implementation of AI agents

Design

- Implement strong authentication by following Secure by Design principles
- Design transparency requirements into the system architecture to enable detection of deception indicators
- Use frameworks, such as zero trust, the Application Security Verification Standard or OAuth2
- Build system infrastructure in a secure, sandboxed environment with encryption, rate limiting and sanitisation
- Apply the principle of least privilege, assigning only the minimum access required for each role
- Limit entitlements to the exact resources, operations and timeframes needed
- Replace static, long-lived secrets with ephemeral credentials that expire when the job is complete
- Dynamically scope privilege for sub-tasks and revoke elevated rights immediately when finished, mitigating scope creep
- Build applications with secure protocols and safe defaults that adhere to communication standards and security policies
- Implement message validation by default so that message components include integrity and freshness checks before use

Development

- Apply secure development principles from [U.S. Department of War Enterprise DevSecOps Fundamentals \[PDF 2,811 KB\]](#)
- Minimise application scope and monitor components for unusual or unexpected behaviour
- Enforce application understanding and only incorporate components into systems that the owner fully understands and accepts the risks of (including possible external effects and processes that it may trigger)
- Refer to frameworks, such as the NIST [Secure Software Development Framework](#) or SLSA's [Safeguarding artifact integrity across any software supply chain](#)
- Use supply chain risk management practices for third-party dependencies
- Apply existing governance and organisational management policies
- Conduct threat modelling using frameworks, such as [OWASP Top 10:2025](#), [MITRE ATT&CK®](#) and [MITRE D3FEND™](#) and use the results to inform mitigations tailored to the operational environment
- Plan and regularly test incident response plans and teams

Disclaimer

The material in this guide is of a general nature and should not be regarded as legal advice or relied on for assistance in any particular circumstance or emergency situation. In any important matter, you should seek appropriate independent professional advice in relation to your own circumstances.

The Commonwealth and co-authoring agencies accept no responsibility or liability for any damage, loss or expense incurred as a result of the reliance on information contained in this guide.

Copyright

© Commonwealth of Australia 2026

With the exception of the Coat of Arms and where otherwise stated, all material presented in this publication is provided under a Creative Commons Attribution 4.0 International license <https://creativecommons.org/licenses/by/4.0/>

For the avoidance of doubt, this means this license only applies to material as set out in this document.



The details of the relevant license conditions are available on the Creative Commons website as is the full legal code for the CC BY 4.0 license <https://creativecommons.org/licenses/by/4.0/legalcode.en>

Use of the Coat of Arms

The terms under which the Coat of Arms can be used are detailed on the Department of the Prime Minister and Cabinet website

<https://www.pmc.gov.au/resources/commonwealth-coat-arms-information-and-guidelines>

For more information, or to report a cyber security incident, contact us:

cyber.gov.au | 1300 CYBER1 (1300 292 371)

ASD AUSTRALIAN
SIGNALS
DIRECTORATE

ACSC Australian
Cyber Security
Centre