

SECURING AGENTIC AI

*An Addendum to the
Guidelines and
Companion Guide on
Securing AI Systems*



2026

This document is an addendum to CSA’s Companion Guide on Securing AI Systems (“Addendum”), focusing on agentic AI systems. Systems owners should use this document in conjunction with the Companion Guide on Securing AI Systems as a resource.

This document is meant as a community-driven resource, developed in collaboration with the AI and cybersecurity practitioner communities. It provides practical mitigation measures and practices to secure AI systems. This document is intended for informational purposes only and is not mandatory, prescriptive nor exhaustive.

DEVELOPED IN CONSULTATION WITH

This document is published by the CSA, in collaboration with partners across the AI and Cyber communities, including:

Accenture
AI Asia Pacific Institute
Alibaba Cloud
Amaris AI
Cisco
CrowdStrike
Deloitte Singapore
DSO National Laboratories
Fujitsu Limited
Google Asia Pacific Pte. Ltd.
Government Technology Agency (GovTech)
HP Inc.
Infoblox
Kaspersky Lab Singapore Pte Ltd
Microsoft Singapore
OpenPolicy Association
Palo Alto Networks
PricewaterhouseCoopers Risk Services Pte Ltd
Resaro
The American Chamber of Commerce in Singapore (AmChamSG)
Vulcan (vulcanlab.ai)

DISCLAIMER

The information provided in this document does not, and is not intended to, constitute legal advice. All information is for general informational purposes only. These organisations provided views and suggestions on the security controls, descriptions of the security control(s), and technical implementations included in this Addendum. CSA and its partners shall not be liable for any inaccuracies, errors and/or omissions contained herein nor for any losses or damages of any kind (including any loss of profits, business, goodwill, or reputation, and/or any special, incidental, or consequential damages) in connection with any use of this Addendum. Organisations are advised to consider how to apply the controls within to their specific circumstances, in addition to other measures relevant to their needs. This document contains links to other third-party websites. Such links are informational and do not represent endorsement of content from these third-party sites.

VERSION HISTORY

VERSION	DATE RELEASED	REMARKS
0.1	22 Oct 2025	Release of Addendum on Securing Agentic AI for Public Consultation
1.0	TBA	Official publication of Addendum

EXECUTIVE SUMMARY

Agentic artificial intelligence (AI) systems are self-managing AI systems that can plan, execute, critique, and iterate across multiple steps to achieve specified objectives. These systems represent a significant evolution from traditional AI systems, moving beyond simple pattern recognition and predetermined responses to demonstrate increasingly sophisticated abilities to understand context, formulate plans, and take independent actions to achieve specified objectives. Development of these systems bring new capabilities and opportunities for organisations and users.

Organisations must carefully consider both the transformative potential and inherent risks these agentic AI systems present. Their capacity to operate with reduced human oversight introduces novel security considerations around system boundaries, control mechanisms, and the potential for unexpected emergent behaviours. Understanding and addressing these security implications is crucial as agentic AI becomes more prevalent in our digital infrastructure and business operations.

The Cyber Security Agency of Singapore (CSA) has developed this addendum to advise system owners on securing their agentic AI systems. This addendum is meant to be read together with the Guidelines and Companion Guide on Securing AI Systems, which outline foundational AI security principles.

As an addendum to the Guidelines, this document takes a risk-based approach across the AI development lifecycle, while introducing new considerations that are relevant to agentic AI. These considerations include mapping out agentic workflows to identify potential threat vectors to the system.

To complement the Companion Guide, this addendum lists agentic AI-related risks and mitigations across the development lifecycle, categorised by capabilities of agentic AI systems. In addition, examples based on current industry use cases are provided as a practical resource on how to apply the addendum.

This document is intended for informational purposes only and is not mandatory, prescriptive nor exhaustive. The content of this document should not be construed as comprehensive guidance or definitive recommendations.

TABLE OF CONTENTS

QUICK REFERENCE TABLE	7
1. INTRODUCTION.....	9
2. HOW AGENTIC AI WORKS	11
2.1. BASELINE COMPONENTS.....	13
2.2. BASELINE SYSTEM DESIGN.....	15
2.2.1. Agentic AI system architecture	15
2.2.2. Roles & access control	17
2.2.3. System workflows & autonomy	17
2.3. CAPABILITIES.....	23
3. SECURITY THREATS TO AGENTIC AI SYSTEMS.....	26
4. SECURING AGENTIC AI	28
4.1. TAKE A LIFECYCLE APPROACH, AND START WITH A RISK ASSESSMENT	28
4.2. IDENTIFY THE RELEVANT MEASURES & CONTROLS.....	33
4.3. TREATMENT MEASURES / CONTROLS FOR AGENTIC AI SYSTEMS	35
3. DEPLOYMENT	45
4. OPERATIONS AND MAINTENANCE	49
5. USE CASE EXAMPLE.....	57
5.1. Case Study 1: Web application development system (SaaS implementation)...	57
5.2. Case Study 2: Client Onboarding System (In-house development).....	74
5.3. Case Study 3: Automated Fraud Detection System	81
ANNEX A: Threats to Agentic AI Systems	86
ANNEX B: Model Context Protocol.....	90
ANNEX C: Agent 2 Agent Protocol	93
REFERENCES	97

QUICK REFERENCE TABLE

Stakeholders in specific roles may use the following table to quickly reference relevant controls in [Section 4.2 – IDENTIFY THE RELEVANT MEASURES & CONTROLS](#).

The roles defined below are included to guide understanding of this document and are not intended to be authoritative.

Decision Makers:

Responsible for overseeing the strategic and operational aspects of AI implementation for the AI system. They are responsible for setting the vision and goals for AI initiatives, defining product requirements, allocating resources, ensuring compliance, and evaluating risks and benefits.

Roles Included: Product Manager, Project Manager

AI Practitioners:

Responsible for the practical application (i.e. designing, developing, and implementing AI solutions, including AI agents) across the life cycle. This includes collecting, procuring or analysing data that goes into systems, building the AI system architecture and infrastructure, building and optimising the AI system to deliver the required functions, as well as conducting rigorous testing and validation of AI models and agents to ensure their accuracy, reliability, and performance. In cases where the AI system utilises a third-party AI system, AI Practitioners also include the third-party providers responsible for these activities, such as those contracted through a Service Level Agreement (SLA). AI practitioners would be in charge of implementing the required controls across the entire system.

Roles Included: AI/ML Developer, AI/ML Engineer, Data Scientist

Cybersecurity Practitioners:

Responsible for ensuring the security and integrity of AI systems. This includes implementing security measures to protect AI systems in collaboration with AI Practitioners, monitoring for potential threats, ensuring compliance with cybersecurity regulations.

Roles Included: IT Security Practitioner, Cybersecurity Expert

Table 1: User Quick Reference Table

The following measures/ controls may be relevant to Decision Makers:	The following measures/ controls may be relevant to AI Practitioners:	The following measures/ controls may be relevant to Cybersecurity Practitioners:
1.1 Conduct a risk assessment	1.1 Conduct a risk assessment	1.1 Conduct a risk assessment
2.1 Supply chain security	2.1 Supply chain security	2.1 Supply chain security
2.7 Limit agency	2.2 Model hardening	2.3 System hardening
2.10 Self-reflection	2.3 System hardening	2.4 Identify, track and protect assets
2.11 Hallucination	2.4 Identify, track and protect assets	2.5 Regular backups
	2.5 Regular backups	2.6 Authorisation and authentication
	2.6 Authorisation and authentication	2.7 Limit agency
	2.7 Limit agency	2.8 Secure by default
	2.8 Secure by default	2.9 Environment segmentation
	2.9 Environment segmentation	
	2.10 Self-reflection	
	2.11 Hallucination	
3.2 Security testing	3.1 Availability controls	3.1 Availability controls
	3.2 Security testing	3.2 Security testing
	3.3 Secure use of external tools	3.3 Secure use of external tools
	3.4 Secure inter-agent communication	3.4 Secure inter-agent communication
4.3 Continuous monitoring and logging	4.1 Validate inputs	4.1 Validate inputs
4.4 Human-in-the-loop	4.2 Validate outputs	4.2 Validate outputs
4.5 Vulnerability disclosure	4.3 Continuous monitoring and logging	4.3 Continuous monitoring and logging
	4.4 Human-in-the-loop	4.5 Vulnerability disclosure
	4.5 Vulnerability disclosure	

1. INTRODUCTION

Agentic **artificial intelligence (AI) systems** are self-managing AI systems that can plan, execute, critique, and iterate across multiple steps to achieve specified objectives. The emergence of these systems reflects ongoing developments in AI that brings new capabilities and opportunities for organisations and users. These systems are capable of autonomous, goal-driven decision making and execution, which will reshape how we interact with AI.

Agentic AI systems represent a significant evolution from traditional AI systems, moving beyond simple pattern recognition and predetermined responses to demonstrate increasingly sophisticated abilities to understand context, formulate plans, and take independent actions to achieve specified objectives. To achieve these objectives, agentic AI systems make use of **AI agents**—modular systems driven by Large Language Models (LLMs) and Large Image Models (LIMs) for narrow, task-specific automation¹. Multiple AI agents may be used together and orchestrated by an autonomous agentic AI system.

As organisations begin to deploy agentic AI systems (and AI agents) across various domains—from process automation and customer service to complex decision support and resource optimisation—we must carefully consider both the transformative potential and inherent risks these systems present. Their capacity to operate with reduced human oversight, while potentially increasing efficiency and scalability, also introduces novel security considerations around system boundaries, control mechanisms, and the potential for unexpected emergent behaviours. Understanding and addressing these security implications is crucial as agentic AI becomes more prevalent in our digital infrastructure and business operations.

The Cyber Security Agency of Singapore (CSA) has worked closely with AI and cybersecurity practitioners to develop this addendum to advise system owners on securing their agentic AI systems. This addendum is meant to be read together with the [Guidelines and Companion Guide on Securing AI Systems](#)², which provides foundational AI security principles and controls to mitigate general risks to AI systems.

This document is intended for informational purposes only and is not mandatory, prescriptive nor exhaustive. The content of this document should not be construed as comprehensive guidance or definitive recommendations.

¹ Sapkota, R., Roumeliotis, K. I., & Karkee, M. [AI Agents vs. Agentic AI: A Conceptual Taxonomy, Applications and Challenges](#)

² Cyber Security Agency of Singapore. [Guidelines and Companion Guide on Securing AI Systems](#)

PURPOSE AND SCOPE

Purpose

This addendum curates practical measures and controls that system owners can use to secure their adoption of agentic AI systems. These measures and controls are voluntary, and not all the measures and controls listed in this addendum will be applicable to all organisations or environments. Organisations may also be at different stages of AI development (e.g. POC, pilot, beta release). Organisations should consider relevance to their use cases as well.

This addendum is meant to be read with the [Guidelines and Companion Guide on Securing AI Systems](#)³. As this Addendum is focused on the key elements of agentic AI systems, the relevant treatment measures/controls from the Companion Guide may still apply to underlying systems and related processes, even if not covered in this document.

Scope

The measures and controls within the addendum address the cybersecurity threats and risks relevant to agentic AI systems. It does not specifically address AI safety, or other common attendant considerations for AI such as fairness, transparency or inclusion, although it is noted that some of the recommended cybersecurity controls may address AI safety risks as well. It also does not cover the misuse of AI for cyberattacks (AI-enabled malware), and scams (deepfakes).

³ Cyber Security Agency of Singapore. [Guidelines and Companion Guide on Securing AI Systems](#)

2. HOW AGENTIC AI WORKS

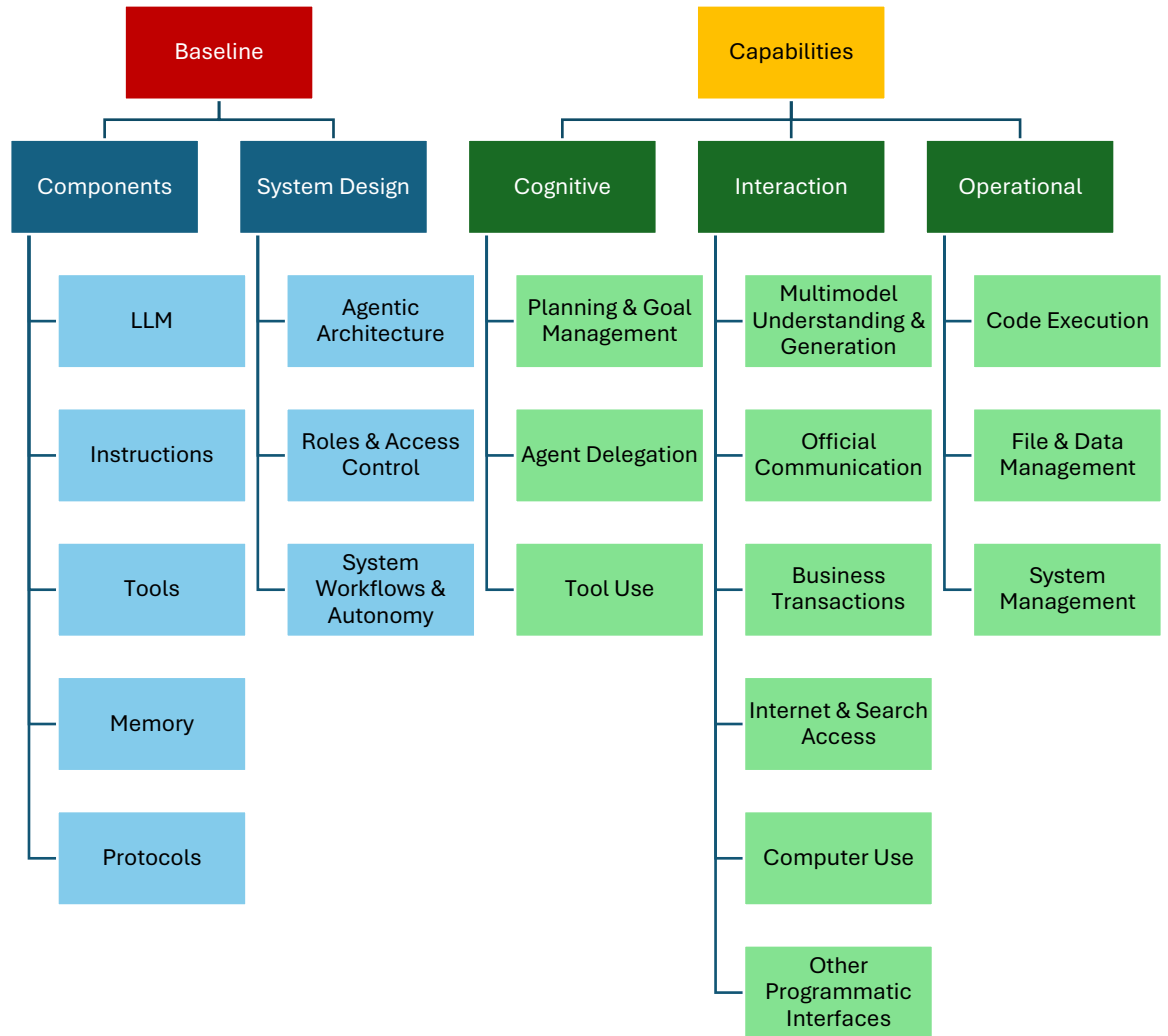
Agentic AI systems interact with their environment, collect data and perform self-determined tasks to meet specified goals.

We can describe the agentic AI system through the following, which helps system owners to understand how agentic AI systems operate and what considerations are needed for safe and effective deployment:

- Key components that facilitate its operation,
- System design, including its architecture; and
- Capabilities (cognitive, interactive, operational)

These elements help system owners to understand how agentic AI systems operate and what considerations are needed for safe and effective deployment.

Figure 1: Baseline and Capability Taxonomy, AI Risk and Capability Framework⁴








⁴ GovTech Singapore (AI Practice). [Agentic Risk & Capability Framework](#).

2.1. BASELINE COMPONENTS

Large Language Models (LLMs) alone are constrained in their operations. While they can be sophisticated in terms of processing input and content generation, by themselves they cannot directly take actions beyond providing information. Agentic AI systems transform this paradigm fundamentally by connecting LLMs to functional tools and systems. This enables them to execute tasks such as sending emails, reading and writing to files and databases, interacting with other software systems, or orchestrating multi-step processes. This expansion from content generation to actual action relies on the integration of multiple components.

Table 2: Key Components in Agentic AI Systems

Component	Description
 Large Language Model (LLM) <i>“Brain” of the agent</i>	A large language model (LLM) is a neural network trained on massive text data to learn statistical patterns of language, enabling it to understand, generate, and reason over natural language inputs. It is the core reasoning and planning engine, or the “brain” of the agent, that processes instructions, interprets user inputs, and generates contextually appropriate responses by leveraging its trained language understanding and generation capabilities. There are a large variety of LLMs to choose from, with different sizes, capabilities, and architectures, from large closed-source LLMs such as OpenAI’s GPT-5.2, Anthropic’s Claude 4 Opus, or Gemini 3 Pro, to smaller open-weights models like Llama 3.1 8B, Qwen3 8B, and Mistral Small 3.1.
 Instructions <i>Guidance for the LLM</i>	Instructions are structured inputs provided to an LLM that define the task, constraints, and context, guiding how the model interprets inputs and generates outputs. They guide the LLM’s decision process by conditioning prompts and tool use with objectives, policies, and guardrails. Forms include system prompts, policies, schemas, and rubrics, varying by framework and enforcement strictness.
 Tools <i>Allows agents to retrieve information, execute actions, affect external systems</i>	Tools are external functions, APIs, or resources attached to an LLM that extend its capabilities beyond text generation, enabling it to retrieve information, execute actions, and affect external systems (e.g., search, code execution, database queries, or file manipulation), typically through structured interfaces such as the Model Context Protocol (MCP) that constrain invocation, inputs, permissions, and outputs for agentic control and safety.

 Memory <i>Information that is stored and accessible to the LLM</i>	<p>Memory is a persistent or semi-persistent data store attached to an LLM that retains information across interactions—such as conversation history, user preferences, task state, and retrieved knowledge—enabling continuity, long-horizon reasoning, and personalization in agentic workflows. Memory may be implemented through mechanisms such as context windows, vector databases, episodic logs, or external state stores, each offering different trade-offs in latency, fidelity, and persistence.</p>
 Protocols <i>Allow agents to communicate with other tools and agents</i>	<p>Protocols are standardised communication frameworks and interfaces that define how different components of an AI system interact with each other and external services. They establish the rules, formats, and procedures for data exchange, function calling, and system integration, ensuring interoperability and consistent behaviour across different tools, models, and platforms. Examples include:</p> <ul style="list-style-type: none"> • Model Context Protocol (MCP) for tool integration, • API specifications like OpenAPI for service communication, • inter agent protocols for multi-agent coordination, such as Agent-to-Agent (A2A), Agent Communication Protocol (ACP) and Agent Network Protocol (ANP), • custom protocols for workflow orchestration and system-to-system communication.

Typically, the process of transforming a user’s inputs into execution of a task involves:

1. **Receiving inputs.** The AI agent receives a specific instruction or goal from the user.
2. **Layering on perception.** The AI agent collects sensory input from sources, such as cameras or microphones, or screen captures and processing technology. This helps it to detect contextual cues and perceive its environment.
3. **Reasoning and planning.** The LLM helps to break down the goal into smaller actionable tasks.
4. **Orchestration and action execution.** Perform tasks based on specific orders or conditions. This may include interactions with other agents, and/or connected systems and tools.
5. **Render a response.** Updates the user on the outcome in an appropriate format.

2.2. BASELINE SYSTEM DESIGN

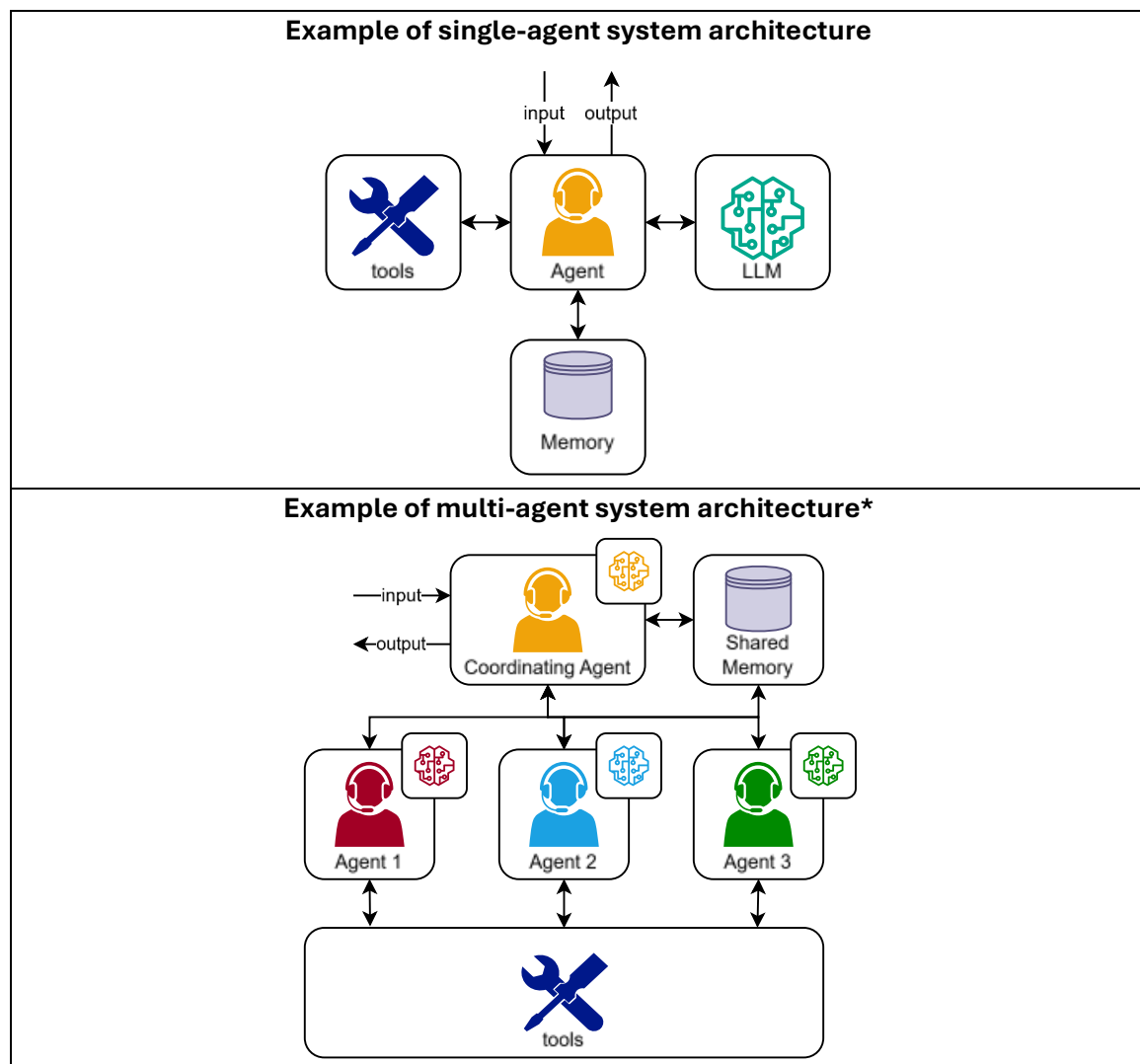
2.2.1. Agentic AI system architecture

The agentic AI system architecture defines how agents are connected, coordinated and orchestrated to solve tasks.

A single-agent system is an AI system with one agent that handles all tasks independently. A multi-agent architecture comprises multiple agents, collaborating to scale or combine specialist roles and functionalities. The co-operation across multiple agents enables solving problems that go beyond the capabilities of would be infeasible for a single agent alone.

Different architectures result in varying levels of system-wide risk, which should be considered carefully. For example, in a multi-agent architecture, risks may be amplified due to the dynamics between agents. Refer to Table 3 for more details.

Figure 2: Examples of single- vs. multi-agent system architecture



*For pictorial clarity, the LLMs are placed within each agent to avoid clutter. LLMs may still be called externally if needed (e.g. through APIs).

Table 3: Key differences between single agent and multi-agent systems

	Single-agent	Multi-agent
Complexity and architecture	Simple and centralised architecture	More complex, distributed architecture
Decision-making capabilities	Centralised decision-making by one agent	Distributed decision-making amongst multiple agents, and hence should be able to address more complex tasks as tasks can be delegated to different specialised agents
Task complexity	Handles one task at a time	Can manage multiple tasks simultaneously
Adaptability	May struggle with dynamic environments	More likely to adjust and respond in real-time to changes in environment
Inter-agent Communications	Operates in isolation; no inter-agent communication needed	Agents interact and share information, hence requiring communication through protocols (e.g. A2A, ACP, ANP)
Fault tolerance	Simple system with limited redundancy – could have a single point of failure.	Easier to build redundancy, but complex system could have correlated failures ⁵ .

In multi-agent architectures, communication takes place among agents, on top of tools and services that applies to the single-agent architecture. Traditionally, such integration with tools and services may require separate and on-off integrations. With the rise of agentic AI, we observe the release of protocols (e.g. Anthropic’s Model Context Protocol (MCP), Google’s Agent2Agent (A2A), IBM’s Agent Communication Protocol (ACP)). that allow for a simplified, consistent, and standardised way for agents to communicate. These reduce the effort required to onboard new tools, services and agents.

⁵ Correlated failures are when multiple components fail due to a single shared cause.

2.2.2. Roles & access control

Roles and access controls establish the responsibilities and permissions across agents in the system. This helps to limit the impact of incidents such as unauthorised actions or access, or potential system failures. Agent roles can include:

- Orchestrator agents that manage workflows
- Specialist agents that perform pre-defined functions
- Interface agents that handle external communications.

Roles and access controls for agentic AI systems should be clearly defined to avoid unauthorised access or excessive privilege. One possible approach is to differentiate between non-human and human identities, and implement unique identities per agent or agent class to limit information and tool access specifically to the roles of the agent(s). For more details on roles and access control, refer to [Section 4.2, Control 2.6](#).

2.2.3. System workflows & autonomy

An AI agent workflow describes the step-by-step process whereby AI agents use reasoning, planning and tools to perform tasks. Such workflows can also be seen in terms of data movement within agentic AI systems, which becomes increasingly challenging to track with more complex architectures and integration to more tools and capabilities. These workflows range from straightforward linear progressions (see Figure 3) to more intricate branching and/or hierarchical patterns (see Figure 4).

- In a linear workflow, data moves sequentially through predetermined steps i.e. each action follows directly from the previous one.
- Branching workflows are implemented when the agentic AI system needs to make decisions about using multiple tools or services simultaneously, based on the task goal or contextual information. These branching workflows hence create multiple possible paths for data movement.

Figure 3: Example linear workflow

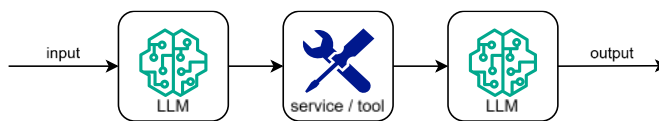
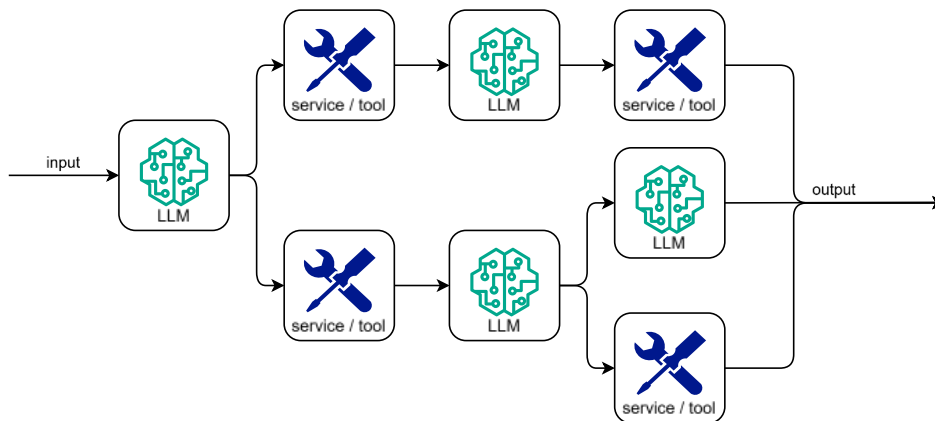


Figure 4: Example branching workflow



Understanding the workflow, as well as data movement, informs risk assessment and threat modelling so that system owners can identify critical points where data might be vulnerable, and prioritise safeguards. These topics are explored in greater detail in [Chapter 3](#).

The workflow within an agentic AI system is also affected by its autonomy, which refers to its ability to operate, make decisions and execute tasks with minimal or no human intervention. As autonomy of the system increases, it also becomes increasingly challenging to assess or predict the potential data flows. This underscores the importance of determining the appropriate autonomy level of the agentic AI system.

Organisations such as NVIDIA have developed frameworks to classify the autonomy levels of agentic AI systems⁶.

Table 4: NVIDIA’s autonomy classification framework

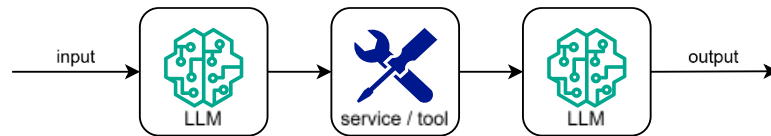
Autonomy Level	Description	Example
0 – Inference API	A single user request results in a <u>single inference call</u> to a single model.	An image classification service that takes a photo and returns a label exemplifies this simplicity. The data path is direct: <u>input → model → output</u> , with no additional processing or decisions.
1 – Deterministic System	A single user request triggers more than one inference request, possibly to more than one model, in a <u>predetermined order</u> that does not depend on either user input or inference results.	In drug discovery, a system might process molecular structures through predetermined stages: <u>initial screening → toxicity analysis → binding prediction</u> . Each step's output feeds into the next in a known sequence.
2 – Weakly autonomous system	A single user request triggers more than one inference request. An AI model can determine if or how to call plugins or perform additional inference at <u>predetermined decision points</u> .	An enterprise document processing system might analyse content type, then <u>route documents through different specialized models</u> : financial documents to compliance checkers, technical documents to subject matter validators, and customer communications to sentiment analysers. While complex, all possible paths can be mapped.
3 – Fully autonomous system	A single user request triggers more than one inference request. In response to a user request, the <u>AI model can freely decide</u> if, when, or how to call plugins or other AI models, or to revise its own plan freely, including deciding when to return control to the user.	A security vulnerability analyser might start with code review, <u>dynamically decide</u> to examine deployment configurations, investigate dependency chains, and recursively explore potential attack vectors, <u>continuously adjusting</u> its investigation based on findings. The number of possible execution paths grows exponentially.

⁶ Harang, R., & Sablotny, M. [Agentic Autonomy Levels and Security](#). NVIDIA.

For Level 0 systems, mapping of workflows may not be necessary as inference calls are made directly to a model, which produces an output. There are no additional services or tools are invoked.

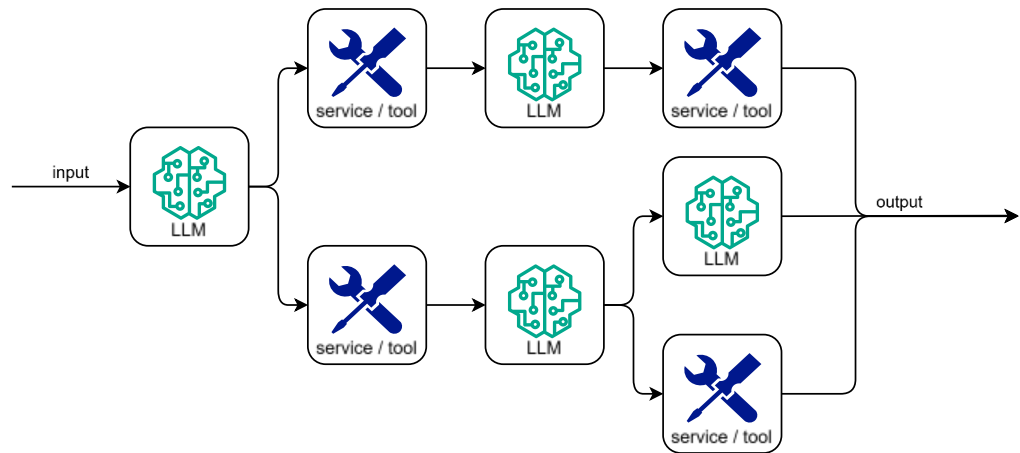
For Level 1 systems and above, mapping of workflows is highly recommended. Level 1 systems usually present as a linear chain of calls in which the output from one AI call or tool response is passed on to the next step in a deterministic manner. The complete workflow is known beforehand.

Figure 5: Autonomy Level 1 – Deterministic system, linear workflow



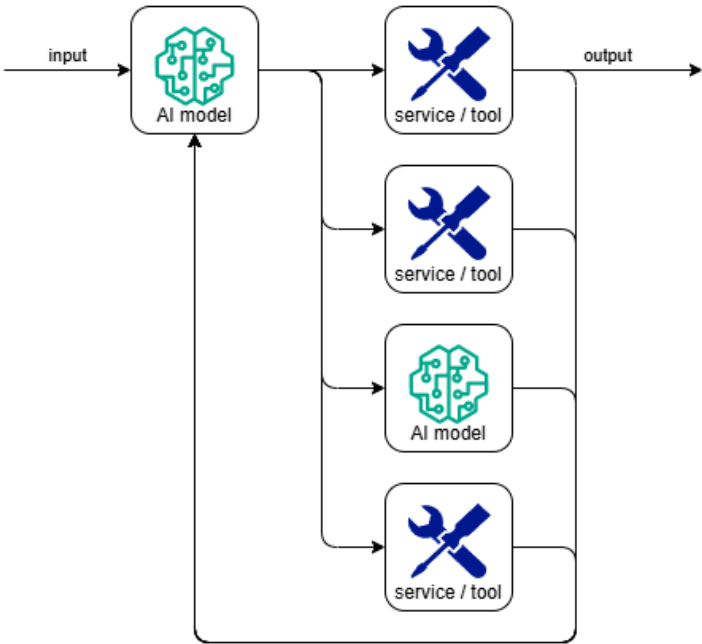
Level 2 systems have outputs that can be sent along various paths through the workflow, based on task requirements and the orchestrator agent’s decision. Every execution path can be determined, but the actual path can only be identified when the workflow is executed.

Figure 6: Autonomy Level 2 – Weakly autonomous system, branching paths at predetermined points



Level 3 systems have significantly more potential execution paths, as more models and tools are invoked. This complexity can be seen in the cyclical path, which indicates a potentially unbounded number of execution paths. It is generally not possible to enumerate all the paths in advance or specific paths which will be used.

Figure 7: Autonomy Level 3 - Fully autonomous system, flows branch to different paths and can be cyclical



Agent Design Patterns

Agent design patterns define how an agentic AI system’s components are organised, integrated, and orchestrated to accomplish a task. Unlike system workflows that only describe the sequence of steps an agent takes, agent design patterns provide reusable architectural templates that determine the fundamental structure and interaction model for an agentic AI system. These templates systematically provide different approaches to organise agents based on specific workload characteristics and requirements. This helps with scalable, easier to maintain implementations (similar to how software design patterns like Model-View-Controller provide standardised approaches to building applications, though agent patterns are still being refined as the field matures.

Examples of these agent design patterns include:

Agent design pattern	Description
Sequential	Specialised agents execute in a predefined, linear order with each agent's output serving as direct input for the next agent, using predefined workflow logic and no AI model orchestration.
Parallel	Multiple specialised sub-agents perform tasks independently and simultaneously, with outputs then synthesised to produce a final consolidated response, using predefined workflow logic and no AI model orchestration.
Loop	Repeatedly executes a sequence of specialised subagents until a specific termination condition is met, using predefined logic and no AI model orchestration.
Reason and act (ReAct)	Uses iterative loops of thought (reasoning about next steps), action (tool selection or final answer), and observation (saving tool outputs) for dynamic planning and continuous adaptation.
Coordinator	Uses a central coordinator agent, with AI model orchestration, to analyse requests, decompose into sub-tasks, and dynamically route these to specialised agents.
Swarm	Uses collaborative all-to-all communication, where a dispatcher routes requests to specialised agents that can communicate with each other and hand off tasks. Lacks central orchestration and requires explicit exit conditions.

System owners should choose an agent design pattern based on the nature of tasks involved (e.g., whether they are predictable and sequential, or complex problems requiring autonomous decision-making with outputs achieved through iterative refinement cycles). Each pattern involves trade-offs: simpler patterns like sequential offer lower complexity and cost but limited flexibility, whilst advanced patterns like swarm provide exceptional capability for complex problems but require significant computational resources and sophisticated orchestration logic.

From a security perspective, agent design patterns can affect the likelihood and impact of attacks such as prompt injection, where malicious instructions embedded in processed content manipulate agents to perform rogue actions or sensitive data disclosure. Agentic AI systems can build resilience through agent design patterns that enforce strict isolation between untrusted data and agent control flow. This should be layered on with relevant security controls (discussed in Chapter 4) for more comprehensive defence.

2.3. CAPABILITIES




AI systems differ in their **capabilities, which can be seen as the general classes of actions that an agentic AI system can perform.**

There are three key categories of capabilities: **cognitive, interaction, and operational**⁷. Each category present distinct functions and interactions with their environment. As each type of capability presents its own value and risks, agentic AI systems with more capabilities can also incur more risks that need to be addressed.

This section focuses on capabilities in Information Technology (IT) systems, where Agentic AI is primarily deployed at this point. Capabilities will continue to evolve, especially as Agentic AI becomes more integrated across services (through open-source tools such as [OpenClaw](#)⁸) or is integrated into environments with Operational Technology.

Cognitive capabilities

Cognitive capabilities encompass the agentic AI system's internal "thinking" skills — how it analyses information, forms plans, learns from experience, and monitors its own performance. For example:







-  **Planning & Goal Management:** The capability to develop detailed, step-by-step, and executable plans with specific tasks in response to broad instructions. This includes prioritizing activities based on importance and dependencies between tasks, monitoring how well its plan is working, and adjusting when circumstances change or obstacles arise.
-  **Agent Delegation:** The capability to assign subtasks to other agents and coordinate their activities to achieve broader goals. This includes identifying which components are best suited for specific tasks, issuing clear instructions, managing inter-agent dependencies, and monitoring performance or failures.
-  **Tool Use:** The capability to evaluate available options and choose the best tool for specific subtasks, based on the capabilities and limitations of different tools and matching them appropriately to the tasks. This includes selecting between search, computation, code execution, or domain-specific APIs, determining when tool invocation is necessary, and sequencing or combining multiple tools to complete complex tasks effectively.

⁷ GovTech Singapore (AI Practice). [Agentic Risk & Capability Framework](#).

⁸ OpenClaw [OpenClaw Threat Model](#).




Interaction capabilities

Interaction capabilities describe how the agentic AI system exchanges information with users, other agents, and external systems. These capabilities below are broadly differentiated based on how and what they interact with:

-  **Multimodal Understanding & Generation:** The capability to communicate with human users across multiple modalities, including natural language conversation (explaining topics, generating documents, interactive discussions) and multimodal understanding/generation (processing and creating image, audio, or video content such as visual analysis, speech transcription, or multimedia creation).
-  **Official Communication:** The capability to autonomously compose, finalize, and dispatch authoritative communications that formally and legally represent an organisation to external parties (e.g. customers, partners, regulators, courts, or media) via approved channels and formats, without prior human review or approval, thereby creating potential legal, regulatory, or reputational obligations. This includes sending legally binding correspondence, publishing official statements or press releases, and responding to external inquiries using the organisation's identity.
-  **Business Transactions:** The capability to autonomously initiate, authorise, and execute binding business transactions with external parties - such as payments, purchases, reservations, or service commitments - within predefined authorisation limits, resulting in real financial, contractual, or operational obligations for the organisation. This includes processing payments or refunds, placing orders or subscriptions, booking services or reservations, and accepting or triggering contractual commitments on behalf of the organisation.
-  **Internet & Search Access:** The capability to autonomously access, browse, search, and retrieve information from the Internet to augment the LLM's static training knowledge with external and up-to-date sources in support of task execution and response generation. This includes issuing search queries, following and parsing web pages, extracting relevant facts or documents, and aggregating information from multiple online sources.
-  **Computer Use:** The capability to directly operate a computer's graphical user interface on behalf of the user, enabling the agent to navigate applications and execute tasks through mouse, keyboard, and window-based interactions. This includes moving the cursor, clicking buttons, entering text, using keyboard shortcuts, switching between windows and applications, and navigating files and menus within the operating system environment.
-  **Other Programmatic Interfaces:** The capability to interact with external systems through non-graphical, programmatic interfaces, such as APIs, SDKs, and backend services, to exchange data or trigger actions as part of task execution. This includes calling REST or GraphQL APIs, invoking cloud services, publishing or consuming messages from queues or event streams, and performing operations such as code pushes, data updates, or system integrations across enterprise platforms.

Operational capabilities

Operational capabilities focus on the agentic AI system's ability to execute actions safely and efficiently within its operating environment. These can include:

-  **Code Execution:** The capability to write, execute, and debug code in various programming languages to automate tasks or solve computational problems. This includes implementing algorithms, writing scripts, compiling and running code, debugging errors, and integrating with external systems through APIs or backend services.
-  **File & Data Management:** The capability to manage the full lifecycle of files and data by creating, reading, modifying, organizing, converting, querying, and updating information across both unstructured artifacts (e.g. documents, spreadsheets, media files) and structured data stores (e.g. SQL/NoSQL databases, data warehouses, vector or embedding stores) in support of operational tasks. This includes ingesting and transforming datasets, generating or updating files, maintaining directory or schema structures, executing database queries or updates, and storing or retrieving embeddings or derived data products.
-  **System Management:** The capability to directly manage and configure technical systems and infrastructure by adjusting system settings, controlling computing resources, and administering operational environments across on-premise or cloud platforms. This includes monitoring system health and performance, managing authentication credentials and access controls, provisioning or scaling compute and storage resources, configuring operating system or runtime parameters, and applying system-level optimisations to support reliable operation.

3. SECURITY THREATS TO AGENTIC AI SYSTEMS

Agentic AI systems face both traditional and novel security challenges. This can be seen as a cumulation across different layers of risks.

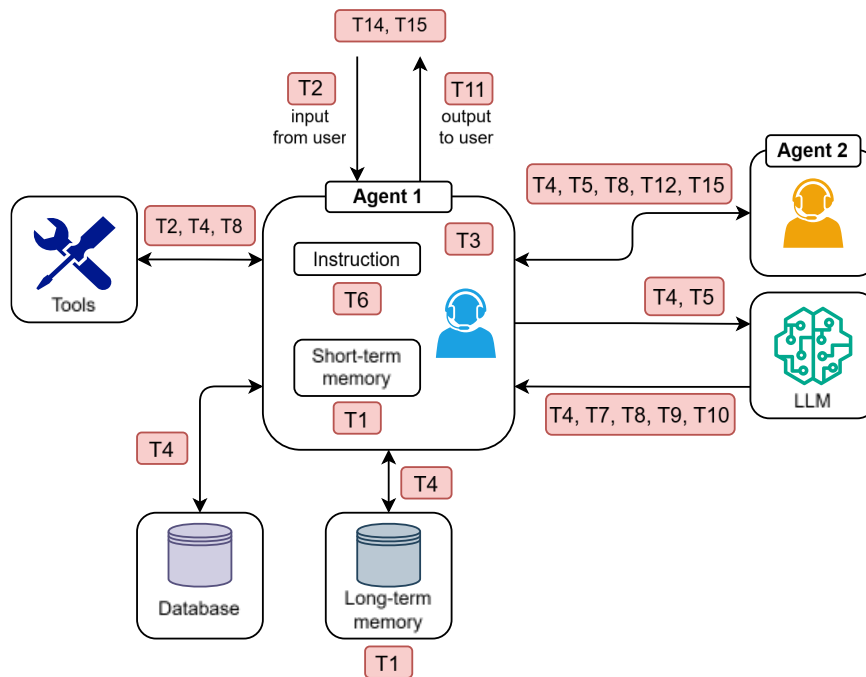
- **Classical cybersecurity risks.** Agentic AI systems have underlying software infrastructure and components, which can make them vulnerable to threats such as remote code execution and SQL injection (if connected to a structured database).
- **Inherited risks from LLM components,** including prompt injection, jailbreaking and data leakage. Refer to CSA's [Guidelines and Companion Guide on Securing AI systems](#), *Section 2.2.2 – Development* for a fuller articulation.
- **New risks to agentic AI systems.** The two primary security concerns in agentic AI systems are rogue actions and sensitive data disclosure.
 - o **Rogue actions** occur when agents perform unintended, or harmful tasks. These can arise through prompt injection, where malicious instructions hidden within normal-looking inputs manipulate the agent's behaviour. They can also occur through simple misunderstandings, if the agent misinterprets ambiguous instructions or handles complex interfaces incorrectly. The impact of these rogue actions directly correlates with the agent's capabilities – more powerful agents pose greater risks when they malfunction.
 - o **Sensitive data disclosure through agent manipulation.** This occurs when attackers exploit agents to reveal private information when agentic workflows are executed. The agent can be guided through a series of seemingly legitimate actions that ultimately leak protected information. Attackers can also manipulate the agent to include sensitive data in its responses.

As with all digital capabilities, there is a balance between utility and risk. These risks are also likely to become more challenging as threat actors enhance their Tactics, Techniques and Procedures, including through new AI-enabled strategies. For agentic AI systems, increasing the agent(s)'s autonomy, access and capabilities can enhance its usefulness. However, this can simultaneously expand the attack surface of the agentic AI system, as well as its potential for causing harm or other undesired actions if they malfunction or are maliciously exploited.

There is a growing body of resources on the risks to agentic AI systems. This includes OWASP's threat taxonomy for agentic AI systems that highlights 15 threats⁹:

- T1 – Memory Poisoning
- T2 – Tool Misuse
- T3 – Privilege Compromise
- T4 – Resource Overload
- T5 – Cascading Hallucination Attacks
- T6 – Intent Breaking & Goal Manipulation
- T7 – Misaligned & Deceptive Behaviours
- T8 – Repudiation & Untraceability
- T9 – Identity Spoofing & Impersonation
- T10 – Overwhelming Human in the Loop
- T11 – Unexpected RCE and Code Attacks
- T12 – Agent Communication Poisoning
- T13 – Rogue Agents in Multi-Agent Systems
- T14 – Human Attacks on Multi-Agent Systems
- T15 – Human Manipulation

Figure 8: Example of threats to agentic AI systems



For more details on the OWASP ASI threat taxonomy, refer to [ANNEX A - Threats to Agentic AI Systems](#) or <https://genai.owasp.org/resource/agentic-ai-threats-and-mitigations/>

⁹ OWASP. [OWASP Top 10 for LLMs - Agentic AI - Threats and Mitigations](#).

4. SECURING AGENTIC AI

4.1. TAKE A LIFECYCLE APPROACH, AND START WITH A RISK ASSESSMENT

CSA's [Guidelines and Companion Guide on Securing AI Systems](#) lay out the two key principles to securing AI systems, including taking a lifecycle approach and starting with a risk assessment. This continues to be relevant for agentic AI systems. The approach to securing AI systems remains relevant and is included here for easy reference. Given the dynamic nature of agentic AI systems, we recommend additional considerations in Steps 1 and 3 to support the risk assessment.

STEP 1 – Conduct a risk assessment, focusing on security risks to agentic AI systems

Conduct a risk assessment, either **based on industry best practices or your organisation's existing Enterprise Risk Assessment/Management Framework**. Risk assessment can be done with reference to CSA's published guides¹⁰, if applicable:

- Guide to Cyber Threat Modelling
- Guide to Conducting Cybersecurity Risk Assessment for Critical Information Infrastructure

Focus on the security risks related to AI systems. For agentic AI systems, we also recommend:

- **Assessing the autonomy level of the system.** This will assess how independently the system operates, how it makes decisions, and how complex its workflows might become. A Level 0 system making straightforward inference calls presents vastly different security challenges compared to a Level 3 system that can dynamically modify its own execution paths.
- **Perform threat modelling to identify areas of interest.** Threat modelling identifies where security risks might occur in the system's workflows. This can be complemented with [taint tracing](#), which is a methodology to track how untrusted data moves through the system. For instance, in a customer service AI system, we

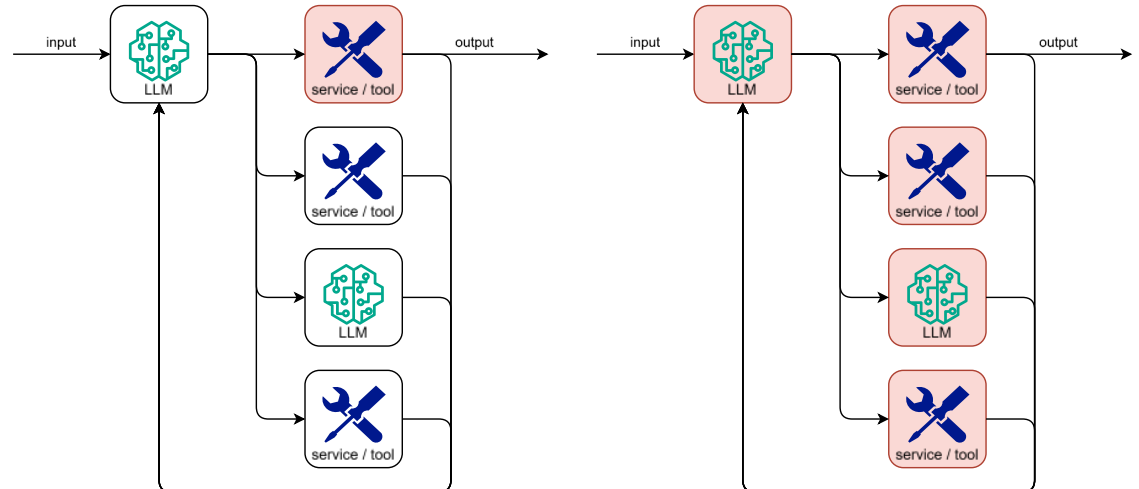
¹⁰ Cyber Security Agency of Singapore. [Supplementary references](#)

can map how user inputs might flow through various decision points and tools, to identify and implement appropriate controls at critical junctures.

- **Identify the risks associated with the agent(s)'s capabilities.** Each capability results in different consequences, and hence different associated risks. Taking a capability-centric approach helps to: (i) be precise about the impact of an agent's operation and potential failure; (ii) identify the different actions involved in realizing the capability, and in turn identify the potential risks. Given that agentic AI system capabilities continue to grow, a capability-centric framework helps to provide a scalable foundation for managing diverse systems.

Taint tracing – tracking data flows from untrusted sources through agentic workflows enables security teams to identify where untrusted data has entered systems and the actions that require additional scrutiny or manual approval¹¹.

Figure 9: Enumerating taints in Level 3 systems (tainted flows marked in red)



Once untrusted data enters the system, the execution workflow is marked as tainted, and every downstream tool and resources are considered to be untrusted. Tainted components should be isolated from the rest of the system, to mitigate downstream impact to the system.

During the pre-deployment phase, taint tracing can be done on workflows to identify high-risk areas. Placement of intrusion detection systems or guardrails can be prioritised at the boundaries of these identified zones. During the post-deployment phase, these same detection systems can be used to flag out and block suspicious data coming into the system.

¹¹ Harang, R., & Sablotny, M. [Agentic Autonomy Levels and Security](#). NVIDIA.

STEP 2 – Prioritise areas to address based on risk/impact/resources

Prioritise which of the identified risks to address, based on the likelihood, impact, available resources, and risk appetite. This can vary across different sectors and use-cases.

STEP 3 – Identify and implement the relevant actions to secure the agentic AI system

Identify relevant actions and control measures to secure the agentic AI system, such as by referencing those outlined in CSA's [Companion Guide on Securing AI Systems](#) as well as in [Section 4.2](#) of this Addendum and implement these across the AI lifecycle. The Companion Guide addresses controls for traditional and generative AI systems, which remains relevant as a baseline for agentic AI systems. The Addendum focuses on additional controls that seek to address risks specific to agentic capabilities, components, and design.

Risk Management for SaaS Environments

For organisations using Software-as-a-Service (SaaS) agentic AI systems, detailed threat modelling and taint tracing may prove impractical due to limited visibility into third-party system architectures and data flows. It may also be challenging to enforce many security controls, as they could be under vendor control rather than the organisation's direct management.

Nonetheless, understanding these risks remains crucial for informed decision-making. The threat modelling and risk assessment processes outlined in this document will help organisations to articulate specific security concerns to vendors, and seek support on appropriate mitigations or transparency about existing controls. Red teaming exercises can also be useful for SaaS deployments, to uncover practical vulnerabilities and attack paths that theoretical threat modelling cannot reveal. This is particularly useful when organisations have limited insight into the actual implementation of third-party systems. These empirical testing approaches help validate if security measures or claims by vendors actually protect against real-world threats.

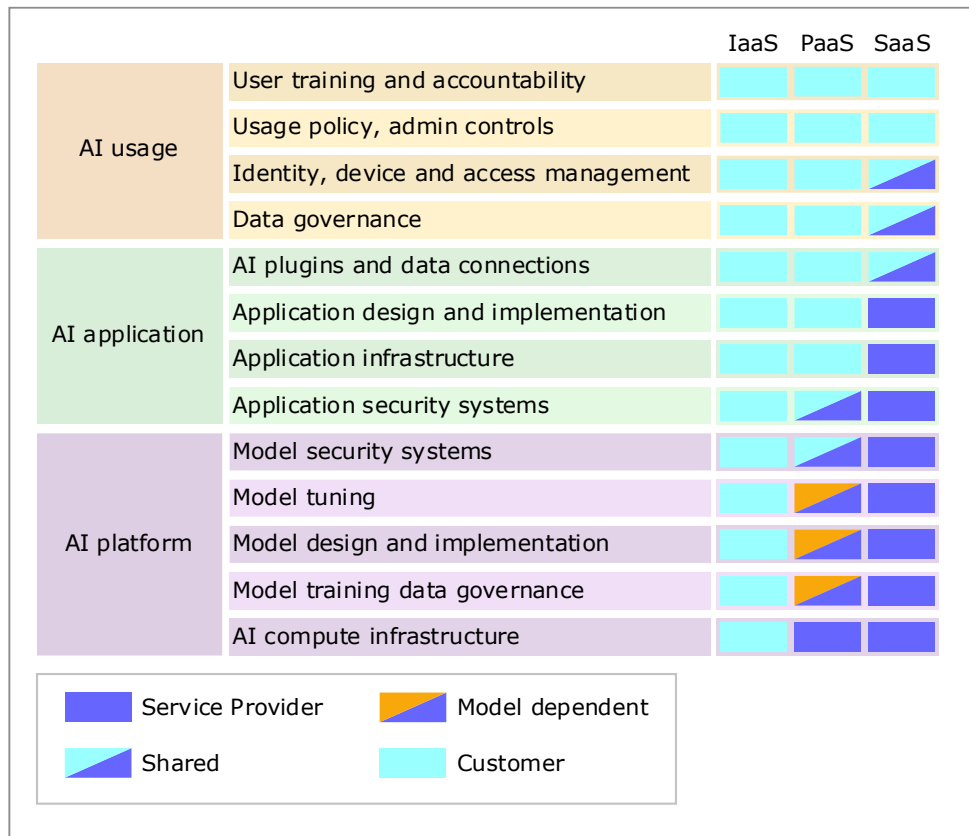
If identified risks cannot be addressed, organisations must update management accordingly as part of formal risk acceptance procedures.

Shared Responsibility Model

This brings us to the related topic of shared responsibility for AI security. The shared responsibility model for use of cloud services is relatively well defined¹², including for deployment approaches like Infrastructure as a Service (IaaS), Software as a Service (SaaS) and Platform as a Service (PaaS). These seek to set out responsibilities between the service user and service provider. For AI, there is growing interest in a similar approach to sharing responsibility, but these are still in development.

is one example¹³ of how responsibilities for security may be shared between AI service users and AI service providers, across different service model implementations. The delineation of responsibilities may vary depending on model type, design pattern, implementation, or terms of use with service providers. Thus, organisations should make the effort to understand the services that they are engaging, and to discuss specific security concerns with service providers. It could also be useful to capture responsibility arrangements via a contract or service level agreement.

Figure 10: AI Shared Responsibility Model



¹² Mell, P., & Grace, T. [Special Publication 800-145 The NIST Definition of Cloud Computing](#)

¹³ Microsoft. [Artificial intelligence \(AI\) shared responsibility model](#)

STEP 4 – Evaluate residual risks for mitigation or acceptance

Evaluate the residual risk after implementing security measures for the AI system to inform decisions about accepting or addressing residual risks.

Periodic re-evaluation

Risk assessment should not be a one-time activity but done throughout a system's operational lifetime. It is important to periodically re-evaluate threat models and controls, especially after significant system changes (e.g., updates to agent workflows, capabilities, or autonomy levels).

For such reviews, it is useful to consider new developments in the threat landscape for agentic AI systems. One example of a relevant resource is the OWASP Top 10 for Agentic Applications 2026¹⁴, which identifies the top 10 highest impact threats to agentic AI.

¹⁴ OWASP. [OWASP Top 10 For Agentic Applications 2026](#)

4.2. IDENTIFY THE RELEVANT MEASURES & CONTROLS

Based on the risk assessment, system owners can identify the relevant treatment measures/controls from the following tables. Each treatment measure/ control plays a different role and should be assessed for relevance and priority in addressing the security risks specific to your agentic AI system and context (Refer to [Section 4.1](#)).

As a start, we recommend users to consider all controls related to the baseline elements, and then to layer on those specific to each capability.

- **Related threats/risks** indicated serve as examples and are not exhaustive. They might differ based on your organisation's use case.
- **Related components/capabilities for each measure/control** are also provided to help you quickly identify what is relevant. Baseline risks are applicable to most, if not all agentic AI systems and should be addressed if possible.
- **Example implementations** are included for each measure/control as a more tangible elaboration on how they can be applied. These are also not exhaustive.
- Additional **references and resources** are provided for users of this document to obtain further details on applying the treatment measure/control if required. For example, the [ARC Framework](#) provides further granularity on the mitigating controls to be applied (e.g. Level 0 to 2). Readers may take this into consideration when prioritising the mitigating controls to be implemented.

As with the [Companion Guide](#), the controls are organised using a lifecycle approach to systematically enumerate every potential mitigation throughout the development lifecycle.

Implementing Controls at Enterprise-scale

The use of a central orchestration and enforcement plane helps organisations implement security controls across the enterprise. Instead of securing every individual AI tool separately, it acts as a single enforcement point where all security rules are applied consistently when AI interacts with the company's internal IT systems.

To secure agents at scale, this approach can address:

- Identity and access management. For example, agents can be identified with their own digital ID. They can be assigned roles in accordance with the least privilege principle, and/or authenticated through OAuth2/OIDC with short-lived and scoped tokens).
- Unintended behaviours – Every request sent to the AI system and every generated answer is scanned to ensure it stays within expected and secure behaviour.
- Data loss prevention – Automated scans of outgoing traffic mitigate the risk of leaking sensitive company information.
- Compliance to policy controls – A master set of rules enforces consistent company-wide standards and legal requirements to be adhered to.
- Continuous logging: Actions, interactions and tool calls made by the AI is recorded and sent to a security dashboard for oversight.

All initiated calls (e.g. to SaaS APIs, internal services, data lakes, etc.) should be tagged with the correct digital ID and routed through an appropriate centralised middleware where applicable.). Furthermore, logs from the middleware should be streamed into a SIEM for SOC monitoring, with processes in-place to revoke agent access when anomalous access is detected.

4.3. TREATMENT MEASURES / CONTROLS FOR AGENTIC AI SYSTEMS

1. PLANNING AND DESIGN

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Example Implementation	Reference / Resource
1.1	<p>Conduct a risk assessment in accordance with the relevant industry standards/best practices.</p> <p>Responsible Parties: Decision Makers, AI Practitioners, Cybersecurity Practitioners</p>	<p>Failure to comply with industry standards/best practices may lead to insufficient, inefficient or ineffective mitigations against adversarial threats.</p> <p>Tainted components in an agentic AI system can have downstream impact along the workflow.</p>	Baseline	<p>As part of a risk assessment and threat modelling, perform taint tracing across workflows throughout the agentic AI system. Taint tracing is especially important for agentic AI systems of higher autonomy levels (i.e. levels 2 and 3).</p> <p>Users are not limited to only one method of threat modelling and may adopt other relevant methods that address their needs.</p>	<ul style="list-style-type: none"> • Chapter 3.2 TAIN T TRACING – IDENTIFYING THREATS ALONG WORKFLOWS • Chapter 5 USE CASE EXAMPLE • NVIDIA, Agentic Autonomy Levels and Security • OWASP GenAI Security Project - Multi-Agentic system Threat Modelling Guide • Cloud Security Alliance, Agentic AI Threat Modelling Framework: MAESTRO • ISO/IEC 42001 – Information technology—Artificial Intelligence—Management system • NIST AI Risk Management Framework (AI RMF 1.0) • ANSSI, Building trust in AI through a cyber risk-based approach

2. DEVELOPMENT

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Example Implementation	Reference / Resource
2.1	<p>Supply Chain Security: Ensure the following components are from trusted sources:</p> <ul style="list-style-type: none"> • data, • models, • agents, • software libraries and dependencies, • developer tools and applications, • packages from MCP servers, • firmware and hardware. <p>Responsible Parties: Decision Makers, AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T1-Memory Poisoning • T8-Repudiation & untraceability • T9-Identity Spoofing and Impersonation • T11-Unexpected RCE and code attacks • T13-Rogue agents in multi-agent systems 	<p>Introduction of bugs, vulnerabilities, unwanted or malicious content, poisoned models or rogue agents from third-party systems can lead to downstream impact.</p>	Baseline	<p>If procuring any AI System or component from a vendor, check/ensure suppliers adhere to the supply chain policies and security standards of your organisation. This could be done by establishing a Service Level Agreement (SLA) with the vendor.</p> <p>Incorporate continuous supply chain monitoring practices.</p>	<ul style="list-style-type: none"> • CSA Critical Information Infrastructure Supply Chain Programme • NCSC Supply Chain Guidance • Supply-chain Levels for Software Artifacts (SLSA) • MITRE Supply Chain Security Framework
		<p>Vulnerabilities in third-party libraries and dependencies used by the agent can cause the system to be exploited.</p>	Baseline	<p>Integrate software composition analysis (SCA) tools or use package managers.</p> <p>Regularly scan dependencies and update libraries with known vulnerabilities.</p> <p>Validate package hashes if available.</p>	<ul style="list-style-type: none"> • pip-audit • GitLab Dependency Scanning • GitHub Dependabot • Snyk Open Source
		<p>Collaborative model poisoning corrupting models across multiple agents. Specific to multi-agent training.</p>	Baseline: LLM	<p>Source data from trusted repositories. Apply data sanitisation and filtering, such as through deduplication and classifier-based quality checks.</p>	<ul style="list-style-type: none"> • Introduction to Training Data Poisoning: A Beginner's Guide, Lakera
		<p>Poorly aligned LLMs may pursue objectives which violate security principles.</p>	Baseline: LLM	<p>Review the LLM's model card for potential alignment issues before using the LLM for more complex tasks.</p>	<ul style="list-style-type: none"> • Model Cards, Hugging Face • Model Cards for Model Reporting
		<p>Poisoned models may introduce hidden model backdoors in the system which may be used by an adversary to perform unwanted actions.</p>	Baseline: LLM	<p>Do not use LLMs from unknown or untrusted sources, even if it is available on public platforms.</p> <p>Scan models to detect for potential backdoors or RCE scripts.</p>	<ul style="list-style-type: none"> • Pickle Scanning

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Example Implementation	Reference / Resource
		Poorly implemented tools may not correctly verify user identity or permissions when executing privileged actions, allowing unauthorised actions.	Baseline: Tools	Do not use tools which do not implement robust authentication protocols.	<ul style="list-style-type: none"> • How to choose a known, trusted supplier for open source software, Google
		Rogue tools that mimic legitimate ones can contain hidden malicious code that executes when loaded.	Baseline: Tools	Do not use tools from unknown or untrusted sources, even if it is available on public platforms.	
		Direct prompt injection from untrusted MCP servers, causing unwanted instructions to be carried out.	Baseline: Tools	Exercise caution when using community-run MCP servers. When possible, use official repositories or well-known sources for MCP servers.	<ul style="list-style-type: none"> • ANNEX B – Model Context Protocol • MCP: Untrusted Servers and Confused Clients, Plus a Sneaky Exploit, Embrace The Red • The Vulnerable MCP Project • Model Context Protocol (MCP): Understanding security risks and controls, Red Hat Blog
		Indirect prompt injection attacks via malicious website content cause unwanted actions to be executed.	Interaction: Internet & Search Access	Use structured retrieval APIs for searching the web rather than through web scraping.	<ul style="list-style-type: none"> • Custom Search JSON API, Google
		Returning unreliable information from websites, causing downstream integrity impact on workflows.	Interaction: Internet & Search Access	<p>Prioritise results from verified, high-quality domains (e.g. .gov, .edu, well-known publishers)</p> <p>Ensure adequate cross-source validation for some of the claims made.</p>	<ul style="list-style-type: none"> • What are credible sources? University of the Sunshine Coast Australia
		Supply chain attacks which impact downstream workflows.	Interaction: Other Programmatic Interfaces	Where possible, enforce zero-trust input handling and validate all data flows. This may also include scanning firmware and verifying hardware components used.	<ul style="list-style-type: none"> • NIST SP 800-207 Zero Trust Architecture
		Indirect prompt injection attacks via malicious data or files cause unwanted actions to be executed.	Operational: File & Data Management	Validate new data used to supplement RAG databases or training data.	<ul style="list-style-type: none"> • Introduction to Training Data Poisoning: A Beginner's Guide, Lakera

2.2	<p>Consider model hardening if appropriate.</p> <p>Responsible Parties: AI Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> T6-Intent breaking & goal manipulation T7-Misaligned & deceptive behaviours 	LLMs with weak performance in instruction following might produce unexpected output, leading to unwanted behaviour.	Baseline: LLM	Prioritise LLMs with stronger performance in instruction following or related capabilities to the task. Benchmarks performance may be used to gauge suitability.	<ul style="list-style-type: none"> Instruction Following Score, Daily Papers, Hugging Face
		AI agents execute disallowed tasks for malicious purposes.	Baseline: LLM	Train models to recognise and refuse disallowed tasks.	<ul style="list-style-type: none"> Refuse Whenever You Feel Unsafe: Improving Safety in LLMs via Decoupled Refusal Training
2.3	<p>Consider implementing techniques to strengthen/harden the system apart from strengthening the model itself.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> T7- Misaligned & deceptive behaviours T11-Unexpected RCE and code attacks 	Introduction of bugs, vulnerabilities through insecure coding practices or design.	Baseline	Adopt Security by Design. Apply software development lifecycle (SDLC) process. Use software development tools to check for insecure coding practices. Implement zero trust principles in system design.	<ul style="list-style-type: none"> NIST SP 800-218 Secure Software Development Framework (SSDF) Version 1.1 NIST SP 800-207 Zero Trust Architecture
		Lack of a robust system prompt design can lead to an increased susceptibility to prompt injection attacks and risk of executing unwanted tasks.	Baseline: Instruction	Implement robust system prompt design.	<ul style="list-style-type: none"> Developing a Robust System Prompt, Code Signal A Closer Look at System Prompt Robustness
		Insecure coding practices leading to vulnerabilities in the system.	Baseline: Agentic Architecture	Adopt secure coding practices. E.g. secure key management via using dependency injection, or secrets management service. Do not hardcode secrets.	<ul style="list-style-type: none"> Secrets Management Cheat Sheet, OWASP Dependency Injection: <ul style="list-style-type: none"> Tools Dependency Injection, AG2 How to pass runtime values to tools (InjectedToolArg), LangChain Secrets Management Services: <ul style="list-style-type: none"> HashiCorp Vault AWS Secrets Manager Google Secret Manager

2.4	<p>Identify, Track and Protect AI system assets</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T1-Memory poisoning • T3-Privilege compromise • T9-Identity spoofing & impersonation • T13-Rogue agents in multi-agent systems 	<p>Loss of data integrity such as through unauthorised changes to data, model, agents or system.</p> <p>Lack of proper documentation of resources may result in the wrong or outdated tool being used by model, causing unwanted behaviour or output.</p>	<p>Baseline</p> <p>Cognitive: Tool Use</p>	<p>Establishing a data lineage and software license management process. This includes documenting the data, codes, test cases, models and agents, including any changes made and by whom.</p> <p>Model cards, Agent cards, Data cards, and Software Bill of Materials (SBOMs) may be used. e.g. provide comprehensive descriptions of each tool, including its intended use, required inputs, and potential outputs</p>	<ul style="list-style-type: none"> • Software Bill of Materials (SBOM), CISA • The ultimate guide to SBOMs, GitLab • Model Cards, Hugging Face • Model Cards for Model Reporting
		<p>Agents may inadvertently store sensitive user or organisational data from prior interactions, resulting in data privacy risks.</p>	<p>Baseline: Memory</p>	<p>Encrypt data at rest and restrict access via fine-grained access controls and audit logs.</p>	<ul style="list-style-type: none"> • Cryptographic Standards and Guidelines, NIST • Guide to Data Protection Practices for ICT Systems, PDPC
2.5	<p>Have regular backups in the event of compromise.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T1-Memory poisoning • T3-Privilege compromise • T11-Unexpected RCE and code attacks 	<p>Manipulation of memory systems and context, causing flawed decision making and unauthorised operations.</p>	<p>Baseline: Memory</p>	<p>Ensure adequate AI-generated memory snapshots for forensic analysis and rollback if anomalies are detected.</p>	<ul style="list-style-type: none"> • LangMem, LangChain
		<p>Execution of insecure code by the model or agents may cause unwanted actions to be performed.</p>	<p>Operational: Code Execution</p>	<p>Ensure proper versioning control of code to allow rollbacks to a more secure and stable version.</p>	<ul style="list-style-type: none"> • What is version control? GitLab • Guide to Data Protection Practices for ICT Systems, PDPC
		<p>Loss of data through overwritten or deleted files</p>	<p>Operational: File & Data Management</p>	<p>Keep a separate backup of original files. Ensure backup copy of database is protected from changes until a specified duration has elapsed, based on organisation's backup policy. Ensure proper versioning of files or database.</p>	

2.6	<p>Implement appropriate authentication, authorisation and access controls to APIs, models, data, logs, tools and the environments that they are in.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T2-Tool misuse • T3-Privilege compromise • T6-Intent breaking & goal manipulation • T11-Unexpected RCE and code attacks • T14-Human attacks on multi-agent systems 	Unauthorised changes in a model's context.	Baseline: Memory	Have robust authentication mechanisms for memory access.	<ul style="list-style-type: none"> • Authentication Cheat Sheet, OWASP • Which OAuth 2.0 Flow Should I Use? auth0 • Security best practice in IAM, AWS • AWS Prescriptive Guidance: Operationalizing agentic AI on AWS • A Novel Zero-Trust Identity Framework for Agentic AI: Decentralized Authentication and Fine-Grained Access Control • NIST SP 800-53r5 Security and Privacy Controls for Information Systems and Organizations
		Unauthorised tool usage.	Baseline: Tools	Enforce strict tool access verification where possible.	
		Agents may gain unauthorised access to restricted resources by exploiting misconfigured or overly permissive roles.	Baseline: Roles & Access Controls	<p>Treat agents as non-human identities. Maintain trusted registry of agent identities and authenticate agents using strong, verifiable credentials. Require unique identifiers for agents using Decentralised Identifiers (DIDs) Apply strict access controls and validate agent roles for requests. Ensure fine-grained, scoped tokens or credentials where possible. Use time-bound (just-in-time) or one-time-use credentials where possible.</p>	
		Exploitation of vulnerabilities in permission management.	Baseline: Roles and Access Controls	Implement granular permission controls, and dynamic access validation. Separate policies for agents and human users.	
		Exploitation of the orchestration layer to perform malicious activities using existing agents.	Baseline: Roles and Access Controls	Implement robust authentication mechanisms for orchestration layer access.	
		Chained authorisation in multi-agent systems can cause downstream agents to execute malicious tasks without checking for permissions.	Baseline: Agentic Architecture	Validate permissions on every request to each agent in the workflow.	
		Leaking personally identifiable or sensitive data	Interaction: Other Programmatic Interfaces	<p>Agents accessing sensitive tools or data should operate under the principle of least privilege in time.</p> <p>Use short-lived, rotating credentials (ephemeral credentials) that expire immediately after agent use.</p> <p>Implement a whitelist approach for interfaces that agents are allowed to use.</p>	

		Man-in-the-middle attacks arising from insecure communications	Cognitive: Agent Delegation, Interaction: Other Programmatic Interfaces	Ensure all cross-agent authentication and message validation and encryption where necessary	<ul style="list-style-type: none"> • Authentication Cheat Sheet, OWASP
		Discovery of, and connection to unvetted external domains, APIs, or peer agents. This may lead to exfiltration of sensitive data.	Interaction: Internet & Search Access, Other Programmatic Interfaces	Implement a whitelist approach for outward network access, including API requests. Route agent-initiated network traffic through protective DNS services to monitor agent domain lookups and block malicious connections.	<ul style="list-style-type: none"> • Control subnet traffic with network access control lists, AWS • What is an IP based access control list (ACL)? Microsoft Azure • NIST SP 800-81r3 Secure Domain Name System (DNS) Deployment Guide
		Executing vulnerable or malicious code	Operational: Code Execution	Implement a whitelist approach for inward network access	
2.7	Implement controls to limit what models or agents can access and generate. Responsible Parties: Decision Makers, AI Practitioners, Cybersecurity Practitioners Relevant Threats from OWASP ASI: <ul style="list-style-type: none"> • T2-Tool misuse • T3-Privilege compromise • T6-Intent breaking & goal manipulation • T7-Misaligned & deceptive behaviours • T11-Unexpected RCE and code attacks • T13-Rogue agents in multi-agent systems 	Abuse of agent-accessible tools to execute unintended actions.	Baseline: Tools	Establish clear operational boundaries to prevent misuse of tools. Set limits on what agents can modify through appropriate guardrails.	<ul style="list-style-type: none"> • Implementing effective guardrails for AI agents • Authorization Cheat Sheet, OWASP • Which OAuth 2.0 Flow Should I Use? auth0 • Security best practice in IAM, AWS • OAuth Scopes, OAuth 2.0 • AWS Prescriptive Guidance: Operationalizing agentic AI on AWS • MI9 - Agent Intelligence Protocol: Runtime Governance for Agentic AI Systems
		Agents gain unauthorised and excessive privileges to perform unwanted actions outside the given scope.	Baseline: Roles and Access Controls	Implement a policy-evaluation engine that assesses authorisation requests dynamically at runtime. Prevent cross-agent privilege delegation unless explicitly authorised through predefined workflows. Do not grant admin privileges to agents, unless strictly necessary for completion of tasks.	
		Compromised agents act outside their operational boundaries and perform unintended actions.	Baseline: Roles and Access Controls	Restrict AI agent autonomy using policy constraints. Scope agent privileges dynamically: strict access of tools and API only to what is necessary to run the tasks. Do not allow agents to modify privileges.	
		Assigning tasks incorrectly to other agents	Cognitive: Agent Delegation	Apply guardrails to limit the scope of tasks that can be assigned to specialised agents.	
		Excessive agent privileges to access unintended resources on the computer, causing potential security impact.	Interaction: Computer Use	Limit computer usage to accessing only required resources on the computer.	

		Exfiltration of sensitive data through insecure communications between agents.	Cognitive: Agent Delegation	Ensure that sensitive data is not passed and leaked between agents by using appropriate guardrails.	
		Misinterpreting inter-agent messages due to poor formatting or weak protocols	Cognitive: Agent Delegation	Constrain agent communication with structured outputs and interactions.	<ul style="list-style-type: none"> • Agent Communication Protocol (ACP) • Agent to Agent (A2A) Protocol • Model Context Protocol (MCP)
		Impersonating or accessing peer agents or services via shared roles or credentials	Cognitive: Agent Delegation	Isolate roles and credentials of each agent.	<ul style="list-style-type: none"> • Security best practice in IAM, AWS
		Lack of proper whitelist controls may lead to the execution of vulnerable or malicious code.	Operational: Code Execution	Create a whitelist of commands that agents are allowed to run autonomously. Deny execution of all other commands that are not whitelisted.	<ul style="list-style-type: none"> • Input Validation Cheat Sheet, OWASP
		Misconfiguring system resources, compromising system integrity and availability.	Operational: System Management	Only grant agents privileges to modify system resources if strictly necessary for completion of tasks. Set minimum and maximum limits to what can be modified.	<ul style="list-style-type: none"> • OAuth Scopes, OAuth 2.0
		Exposure of personally identifiable information in files.	Operational: File & Data Management	Whitelist only files which are required for the task. Do not grant access to files known to host private or sensitive information without careful consideration of the risks. Consider using data anonymisation techniques instead.	<ul style="list-style-type: none"> • Advisory Guidelines on use of Personal Data in AI Recommendation and Decision Systems, PDPC • Guide to Basic Anonymisation, PDPC

2.8	<p>Apply the principle of least privilege. Ensuring configurations are secure by default.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> T2-Tool misuse T3-Privilege compromise 	Agents having unauthorised access to restricted resources by exploiting misconfigured or overly permissive roles.	Baseline: Roles & Access Controls	Apply principle of least privilege when configuring all agent and delegation roles.	<ul style="list-style-type: none"> Authorization Cheat Sheet, OWASP Security best practice in IAM, AWS Guide to Basic Anonymisation, PDPC
		Agents having privileges/rights to execute untrusted or malicious code.	Operational: Code Execution	Scope execution privileges strictly only to what is necessary, ensuring that privileges are customised to each agent within a system. Do not grant admin or sudo privilege by default. Block all inward and outward network access by default.	
		Agents having privileges/rights to overwrite or delete database tables or files.	Operational: File & Data Management	No write access to tables in the database unless strictly required, with consideration of risks of data loss.	
		Exposure of personally identifiable or sensitive data from databases or files to users.	Operational: File & Data Management	Restrict access to personally identifiable data or sensitive data unless strictly required, with consideration of risks of data exposure. Consider data anonymisation techniques instead.	
2.9	<p>Implement segregation of environments and network segmentation.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> T3-Privilege compromise T6- Intent breaking & goal manipulation T11-Unexpected RCE and code attacks 	Rogue tools that mimic legitimate ones can contain hidden malicious code that executes when loaded and gain access to other assets within the environment or network.	Baseline: Tools	Test third-party tools in hardened sandboxes with syscall/network egress restrictions before using them in production environments.	<ul style="list-style-type: none"> Sandboxing Agentic AI Workflows with WebAssembly, NVIDIA E2B SDK E2B Data Analysis, LangChain Docker Security Cheat Sheet, OWASP Defeating Prompt Injections by Design (CaMeL), Google DeepMind Advancing Zero Trust Maturity Throughout the Network and Environment Pillar, NSA
		Prompt injection attacks and indirect data manipulation through access to other assets within the environment or network.	Baseline: Agentic Architecture	Decouple data processing flow from control flow through runtime security architecture. Update segmentation policies regularly as agent roles, system architecture and organisational risk posture evolves.	
		Prompt injection attacks to perform credential and/or data exfiltration through access to other assets within the environment or network.	Interaction: Business Transactions	Ensure virtual isolation for agents carrying out transactions. Do not share credentials with the agent directly, require the agent to use a separate service for authentication and transactions.	

		Execution of insecure or malicious scripts that affects the other components of the environment or network.	Operational: Code Execution	Run code in virtually isolated compute environments (e.g. Docker, Podman containers). Sandbox the execution of AI generated scripts. Monitor the execution.	<ul style="list-style-type: none"> • Sandboxing Agentic AI Workflows with WebAssembly, NVIDIA • E2B SDK, E2B • E2B Data Analysis, LangChain • Docker Security Cheat Sheet, OWASP
2.10	Implement model self-reflection before making decisions, where applicable Responsible Parties: Decision Makers, AI Practitioners Relevant Threats from OWASP ASI: <ul style="list-style-type: none"> • T5-Cascading hallucination attacks • T6- Intent breaking & goal manipulation • T7-Misaligned & deceptive behaviours 	Incomplete or unclear instructions may compel models to attempt to fill in missing constraints, resulting in incorrect or unwanted actions being executed.	Baseline: Instructions	Ask the agent to summarise its understanding and request clarification before proceeding.	<ul style="list-style-type: none"> • Self-Reflecting AI Agents using LangChain • AWS Prescriptive Guidance: Operationalizing agentic AI on AWS
		Purpose drift, or unintended deviation from the user's instructions to perform other tasks or pursue other priorities, resulting in malicious or deceptive behaviour.	Cognitive: Planning & Goal Management	Prompt the agent to self-reflect on the adherence of the plan to the user's instructions.	
		Incorrect assignment of tasks to other agents.	Cognitive: Planning & Goal Management	Prompt the agent to self-reflect and assess the suitability of tasks delegated to agents.	
		Unintended pursuit or prioritisation of other goals, resulting in malicious or deceptive behaviour.	Cognitive: Reasoning & Problem-Solving	Understand the reasoning and self-reflection done by the agent through visualisation of its thought process. Log the output in the console for the user to evaluate and verify.	
2.11	Implement controls to reduce the likelihood of hallucination. Responsible Parties: Decision Makers, AI Practitioners Relevant Threats from OWASP ASI: <ul style="list-style-type: none"> • T5-Cascading hallucination attacks 	Agents may mistakenly store glitches and hallucinations into memory, resulting in compounding errors when incorrect information is retrieved for decisions or actions.	Baseline: Memory	Schedule periodic memory reconciliation, where human reviewers or external tools can flag anomalies.	<ul style="list-style-type: none"> • Mem0: Building Production-Ready AI Agents with Scalable Long-Term Memory • Zep: A Temporal Knowledge Graph Architecture for Agent Memory • RAG and the value of grounding, elastic search labs
		Generating non-factual or hallucinated content which can propagate downstream and cause compounding errors that can affect the integrity of the output.	Interaction: Multimodal Understanding & Generation	Implement features to verify the generated answer against the original content. Conduct testing to measure hallucination and factuality rates for outputs. Implement UI/UX cues to highlight the risk of hallucination to the user. Implement Retrieval Augmented Generation (RAG) to keep the model grounded and contextualised.	

3. DEPLOYMENT

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Example Implementation	Reference / Resource
3.1	<p>Ensure availability controls are in place to mitigate disruption or failure of AI services</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> T4-Resource overload 	(Distributed) denial of service on agents.	Baseline: Agentic Architecture	Apply rate limits on the number of concurrent queries to agents.	<ul style="list-style-type: none"> API Rate Limiting, GitHub Docs
		Degradation of computational or service capability of the system.	Baseline: System Workflows & Autonomy	<p>Deploy resource management controls, implement adaptive scaling mechanisms and monitor system load to detect and mitigate overload attempts in real-time.</p> <p>Implement rate limits on high-frequency task requests per agent session.</p>	<ul style="list-style-type: none"> IT & System Availability + High Availability: The Ultimate Guide, Splunk
		Slow or inefficient responses from being stuck in a reasoning loop.	Cognitive: Planning & Goal Management	<p>Enforce time or token limits for reasoning.</p> <p>Adjust short-term and long-term memory options.</p>	<ul style="list-style-type: none"> OverThink: Slowdown Attacks on Reasoning LLMs
		Exploitation of interactions between agents to cause resource exhaustion across the system.	Cognitive: Agent Delegation	<p>Enforce time or token limits for agent reasoning.</p> <p>Set a limit on the number of agent interactions per task, based on the requirements of the workflow.</p>	
		Compromising database availability through excessive queries.	Operational: File & Data Management	<p>Limit the number of concurrent queries to the database from agents.</p> <p>Analyse past database queries to identify repeated or inefficient queries.</p>	
		Overconsumption of compute resources.	Operational: Code Execution	Monitoring of code runtime and memory consumption.	

3.2	<p>Conduct security testing</p> <p>Responsible Parties: Decision Makers, AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T2-Tool misuse • T6- Intent breaking & goal manipulation • T7-Misaligned & deceptive behaviours • T9-Identity spoofing & impersonation • T13-Rogue agents in multi-agent systems 	<p>Agents may contain underlying problems which can cause unexpected behaviour or logical errors.</p>	Baseline: LLM	<p>Behavioural testing of agents with benchmark datasets to determine performance metrics, and executing simulations in regulated environments to analyse agents' behaviour.</p> <p>Automated evaluators can be used, but human evaluators should verify the results of testing.</p>	<ul style="list-style-type: none"> • Benchmarks: <ul style="list-style-type: none"> - AgentBench - HELM - TheAgentCompany - WebArena • Evaluation platforms with collection of benchmarks: <ul style="list-style-type: none"> - Inspect Evals (UK AI Safety Institute, Arcadia Impact, Vector Institute) - Project Moonshot (AI Verify Foundation)
		<p>AI may engage in specification gaming, where it maximises the goal by exploiting loopholes, without achieving the intended task.</p>	Baseline: Instructions	<p>Conduct adversarial evaluation to discover specification gaming behaviour. Iterate on system prompt design, have more robust reward design, and add constraints.</p>	<ul style="list-style-type: none"> • garak • PromptFoo • PyRIT
		<p>Incomplete or unclear instructions may compel models to attempt to fill in missing constraints, resulting in incorrect or unwanted actions being executed.</p>	Baseline: Instructions	<p>Test the efficacy of system prompts with scenario-based evaluations for task performing and problem solving. Benchmarks may be used.</p>	<ul style="list-style-type: none"> • A Closer Look at System Prompt Robustness
		<p>Inconsistencies between AI outputs and expected reasoning pathways.</p>	Cognitive: Planning & Goal Management	<p>Utilise deception detection strategies such as behavioural consistency analysis, truthfulness verification models, and adversarial red teaming.</p>	<ul style="list-style-type: none"> • Systematic Review of Software Behavioral Model Consistency Checking
		<p>Compromised agents may impact downstream decision making.</p>	Cognitive: Planning & Goal Management	<p>Have regular AI red teaming of agents to check for potential vulnerabilities or compromise.</p> <p>Integrate AI red teaming with existing security testing processes (e.g. penetration testing, vulnerability disclosure process) and include clear scoping, reporting and remediation workflows.</p>	<ul style="list-style-type: none"> • Agentic AI Red Teaming Guide, Cloud Security Alliance • OWASP GenAI Red Teaming Guide • NIST SP 800-115 Technical Guide to Information Security Testing and Assessment • MITRE ATLAS
		<p>Adversarial threats attempting to compromise orchestration or planning agents to use other agents maliciously.</p>	Cognitive: Tool Use	<p>Conduct rigorous adversarial testing on centralised orchestration and planning agents.</p>	

3.3	<p>Ensure necessary security measures are in place when external tools are invoked (e.g. through the deployment of MCP servers)</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T2-Tool misuse • T9-Identity spoofing & impersonation • T11-Unexpected RCE and code attacks • T13-Rogue agents in multi-agent systems 	Insecure configurations allowing unauthorised access to tools, models and data.	Baseline: Tools, Baseline: Roles and Access Controls, Baseline: Protocols	<p>Implement robust security measures to protect MCP servers, such as context-level access controls.</p> <p>Have formal interface versioning, and track where context is coming from. Stay informed about emerging MCP vulnerabilities and security best practices.</p>	<ul style="list-style-type: none"> • ANNEX B – Model Context Protocol • MCP: Untrusted Servers and Confused Clients, Plus a Sneaky Exploit, Embrace The Red • OWASP GenAI Security Project - Multi-Agent system Threat Modelling Guide • The Vulnerable MCP Project • Model Context Protocol (MCP): Understanding security risks and controls, Red Hat Blog • MCP Is a Security Nightmare — Here’s How the Agent Security Framework Fixes It
		Execution of malicious scripts through the MCP server, leading to system compromise.	Operational: Code Execution	<p>Ensure code verification before MCP functions are executed on servers. Sandbox the execution.</p>	
		Introduction of malicious agent(s) into the ecosystem, which rapidly corrupts other agents in the system.	Baseline: Roles and Access Control, Cognitive: Tool Use	<p>Verify that MCP agents are from trusted sources before introducing them into the system.</p> <p>Sanitise tool inputs.</p> <p>Check that an MCP server has not silently redefined their tools (MCP rug pull).</p>	
				<p>To secure interactions between MCP Clients and Servers, use strong validation techniques such as allowlists, hardcoded connections, or mutual TLS (mTLS) for known static relationships.</p> <p>In dynamic environments where clients change, use strong authentication protocols (such as OAuth 2.1 or OIDC) to dynamically verify client identity and allow connections rather than depending exclusively on static network trust.</p>	<ul style="list-style-type: none"> • A Practical Guide for Secure MCP Server Development - OWASP Gen AI Security Project

3.4	<p>Implement security controls between agents.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T3-Privilege compromise • T12-Agent communication poisoning 	<p>Manipulation of communication channels between agents to disrupt workflows or influence decisions.</p>	<p>Baseline: Roles and Access Controls, Baseline: Protocols, Cognitive: Agent Delegation</p>	<p>Monitor and log inter-agent interactions to identify anomalies.</p> <p>Enforce inter-agent authentication; deploy cryptographic message authentication if needed.</p> <p>Enforce multi-agent task segmentation to prevent attackers from escalating privileges across interconnected agents.</p> <p>Ensure multi-agent consensus verification for critical decision-making processes.</p>	<ul style="list-style-type: none"> • ANNEX C – Agent 2 Agent Protocol • What is Message Authentication Code? Fortinet • Agent to Agent (A2A) Protocol • JSON Web Tokens, auth0 • What is mutual TLS (mTLS)? Cloudflare • NIST SP 800-81r3 Secure Domain Name System (DNS) Deployment Guide
		<p>Sensitive data disclosure via eavesdropping between agent communications.</p>	<p>Cognitive: Agent Delegation</p>	<p>Ensure that sensitive data is not passed on and leaked among agents through appropriate guardrails. (e.g., deploying DNSSEC or protective DNS)</p> <p>For highly sensitive data, consider applying end-to-end encryption.</p>	

4. OPERATIONS AND MAINTENANCE

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Example Implementation	Reference / Resource
4.1	<p>Validate inputs to the models and agents.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T1-Memory poisoning • T2-Tool misuse • T5-Cascading hallucination attacks • T7-Misaligned & deceptive behaviours • T11-Unexpected RCE and code attacks • T14-Human attacks on multi-agent systems 	Direct prompt injection attacks to the prompt interface from adversarial inputs to the model.	Baseline: LLM	<p>Implement input guardrails to detect direct prompt injection or adversarial attacks.</p> <p>Implement input sanitisation measures or limit inputs to conventional ASCII characters only.</p>	<ul style="list-style-type: none"> • How to implement LLM guardrails, OpenAI • Guardrails, OpenAI Agents SDK • Guardrails AI • NeMo Guardrails, NVIDIA • LLM Guard, Protect AI • prompt-injection-defenses, tl;dr sec • LLM Prompt Injection Prevention Cheat Sheet, OWASP
		Tools that lack input validation can be exploited through prompt injection attacks.	Baseline: Tools	<p>Enforce strict schema validation (e.g. JSON Schema, protobuf, Pedantic, OpenAI Structured Outputs) and reject non-conforming inputs into the system.</p> <p>Escape or encode user inputs when embedding into tool prompts or commands.</p>	<ul style="list-style-type: none"> • Input Validation Cheat Sheet, OWASP
		Incorrect or manipulated instructions may invoke the wrong tool/service and impact downstream workflows.	Baseline: Instructions	<p>Validate agent instructions before passing on to the model, especially for critical decision workflows.</p>	<ul style="list-style-type: none"> • High-Risk AI Systems Under the EU AI Act • Purple Llama, Meta Llama
		Indirect prompt injection attacks via malicious website content or files.	<p>Interaction: Internet & Search Access.</p> <p>Operational: File & Data Management</p>	<p>Implement input guardrails to detect indirect prompt injection.</p> <p>Implement escape filtering before including web content or relevant files into prompts.</p> <p>Scan external files for undesired input or instruction before passing on to memory or models.</p>	<ul style="list-style-type: none"> • Input Validation Cheat Sheet, OWASP • File Upload Cheat Sheet, OWASP
		Generation of unrelated topic outputs, which may affect integrity of model performance or output.	Interaction: Multimodal Understanding & Generation	<p>Implement input multimodal (or text) guardrails to detect if the instruction is within the expected topic domain.</p> <p>Refuse to answer otherwise.</p>	<ul style="list-style-type: none"> • Purple Llama, Llama Guard, Meta • Perspective API • Content moderation, Anthropic • OpenAI Moderation API • Cloud services: <ul style="list-style-type: none"> - AWS Comprehend - Azure Content Safety • DOMPurify

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Example Implementation	Reference / Resource
		Passing on prompt injection attacks across agents throughout the system(s).	Cognitive: Agent Delegation	Sanitise messages before agents process them - strip or escape unexpected instruction-like content that may have been injected (e.g. remove "ignore", "system", or "from now on").	<ul style="list-style-type: none"> • Microsoft Presidio SDK • spaCy, Explosion
		Executing vulnerable or malicious code.	Operational: Code Execution	Sanitise all inputs for malicious code.	
		Exposure of personally identifiable information from retrieved content.	Operational: File & Data Management	Implement input guardrails to detect personally identifiable information in the content.	

4.2	<p>Validate outputs from the models and agents.</p> <p>Responsible Parties: AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T1-Memory poisoning • T2-Tool misuse • T5-Cascading hallucination attacks • T6-Intent breaking & goal manipulation • T7-Misaligned & deceptive behaviours • T11-Unexpected RCE and code attacks • T15-Human manipulation 	Vulnerabilities in outputs across the agentic workflow may be exploited for malicious purposes downstream, potentially triggering cascading effects that compromise interconnected systems and dependencies.	Baseline: Agentic Architecture	Insert validation checkpoints between stages that verify expected output and reject invalid output.	<ul style="list-style-type: none"> • How to implement LLM guardrails, OpenAI • Guardrails, OpenAI Agents SDK • Guardrails AI • NeMo Guardrails, NVIDIA • LLM Guard, Protect AI
		Disclosure of sensitive or personally identifiable information through unsanitised outputs.	Interaction: Multimodal Understanding & Generation	Implement output guardrails to detect personally identifiable information in the LLM's outputs before it reaches the user, or contained within communications before it is sent out.	<ul style="list-style-type: none"> • Microsoft Presidio SDK • spaCy, Explosion
		Sending malicious or undesired content to recipients.	Interaction: Official Communication	Validate all links and attachments prior to sending them to users.	
			Interaction: Multimodal Understanding & Generation	Implement output safety text guardrails to detect if malicious or undesirable content is being generated, or contained within communications before it is sent out.	
			Interaction: Official Communication	Validate all links and attachments prior to sending them to users.	
		Allowing unauthorised actions (e.g. transactions).	Interaction: Business Transactions	Apply fraud detection models or heuristics to the agent's own decisions.	<ul style="list-style-type: none"> • AI fraud detection in banking, IBM
		Execution of insecure or malicious code that are generated by the LLM.	Operational: Code Execution	Use code linters to screen for bad practices, anti-patterns, unused variables, or poor syntax. Review all code and/or perform static code analysis to detect potential security vulnerabilities before execution. Conduct CVE scanning and block execution of any High or Critical CVEs.	<ul style="list-style-type: none"> • Bandit (Python) • ESLint (JavaScript) • Semgrep (multi-language) • Purple Llama, CodeShield, Meta • Content Security Policy Cheat Sheet, OWASP • Code Review Guide 2.0, OWASP
Output that will be rendered in a web UI may be vulnerable to XSS.	Operational: Code Execution	Sanitise output with libraries for rendering in a web UI. Test against bypass.	<ul style="list-style-type: none"> • XSS Filter Evasion Cheat Sheet, OWASP • DOMPurify • sanitize-html 		
Generation of non-factual content which can propagate downstream and may cause unintended output or behaviour that impacts integrity.	Cognitive: Planning & Goal Management	Have robust output validation mechanisms, or multi-source validation.	<ul style="list-style-type: none"> • Input Validation Cheat Sheet, OWASP 		

4.3	<p>Implement continuous monitoring and logging of access, usage and execution of:</p> <ul style="list-style-type: none"> models, databases/files, memory, agents, tools, MCP interactions, agent communication, external actions, or any relevant software or hardware systems. <p>Responsible Parties: Decision Makers, AI Practitioners, Cybersecurity Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> T2-Tool misuse T3-Privilege compromise T4-Resource overload T8-Repudiation & untraceability T9-Identity spoofing & impersonation T12-Agent communication poisoning T13-Rogue agents in multi-agent systems T14-Human attacks on multi-agent systems T15-Human manipulation 	Model drift over time might cause unexpected output or behaviour.	Baseline: LLM	Continuously monitor and log outputs, triggering alerts when behaviour drifts from tested baselines.	<ul style="list-style-type: none"> MLflow, Databricks OpenLLMetry, traceloop Helicone Langfuse LangSmith, LangChain Cloud provider tools: <ul style="list-style-type: none"> - Azure Agent Monitoring - AWS Bedrock Trace Events <ul style="list-style-type: none"> Best practices for event logging and threat detection, Cloud Security Alliance AWS Prescriptive Guidance: Operationalizing agentic AI on AWS <p>Some key information to capture in logs:</p> <ul style="list-style-type: none"> Agent actions (e.g. API calls, data retrievals) Input and output (what is received, and generated by the agent) Internal state changes (e.g. within knowledge base or memory) Errors and exceptions (e.g. details of errors encountered, including stack traces and relevant context) Timestamps and duration (accurate timestamps of logs and duration of operations) Contextual information (e.g. unique identifiers for specific tasks, session, or workflows to enable correlation across logs)
		Adversarial prompt attacks against the system.	Baseline: LLM	Log queries to detect for possible attacks or suspicious activity. Consider the current privacy regulations/guidelines when logging inputs.	
		Unauthorised users may exploit tools that do not verify user identity or permissions when executing privileged actions.	Baseline: Tools	Conduct periodic audits to validate that tool actions match the appropriate user permissions.	
		Malicious actors exploit attack surfaces that arise from using tools that demand broader permissions than necessary.	Baseline: Tools	Conduct periodic least-privilege reviews and automated permission drift detection.	
		Unauthorised tool usage.	Baseline: Tools	Monitor tool access and usage patterns. Implement execution logs that track AI tool calls for anomaly detection and post-incident review.	

	Exploitation of authentication mechanisms to impersonate agents or human users.	Baseline: Roles and Access Controls	Deploy continuous monitoring to detect fraud or impersonation attempts. Use behavioural profiling, possibly involving a second model, to detect deviations in AI agent activity that may indicate identity spoofing. Automate alerts to developers when suspicious activities are detected.	<ul style="list-style-type: none"> • Chapter 4.2 – Implementing Controls at Enterprise-scale • NIST SP 800-61r3 Incident Response Recommendations and Considerations for Cybersecurity Risk Management • PagerDuty Incident Response Documentation • OWASP GenAI Security Project - Multi-Agentic system Threat Modelling Guide
	Unauthorised or malicious use of elevated privileged operations.	Baseline: Roles and Access Controls	Monitor role changes, and audit elevated privilege operations.	<ul style="list-style-type: none"> • Best practices for event logging and threat detection, Cloud Security Alliance
	In agentic workflows, early mistakes or vulnerabilities can be propagated and magnified downstream.	Baseline: Agentic Architecture	<p>Apply circuit-breakers that freeze propagation when anomalous behaviour is detected, and implement human authorisation for release.</p> <p>Taint tracing may be used to identify key locations in the workflow to apply circuit-breakers.</p>	<ul style="list-style-type: none"> • LangGraph interrupt, LangChain • UserProxyAgent, AG2 • crewAI, Human-in-the-Loop Workflows • Agentic Autonomy Levels and Security, NVIDIA
	More complex agentic architectures may make it difficult to fully reconstruct decision processes across multiple agents, for the purpose of incident response, or triage.	Baseline: Agentic Architecture	<p>Implement end-to-end distributed tracing with unique request IDs or Decentralised IDs (DIDs) across all agents and tool calls.</p> <p>Implement immutable, tamper-evident audit logs that capture prompts, responses, and tool invocations.</p>	<ul style="list-style-type: none"> • A Novel Zero-Trust Identity Framework for Agentic AI: Decentralized Authentication and Fine-Grained Access Control • Short-lived API tokens: <ul style="list-style-type: none"> - What Are Refresh Tokens and How to Use Them Securely, auth0 - JSON Web Tokens, auth0 • Temporary cloud credentials: <ul style="list-style-type: none"> - Use temporary credentials with AWS resources, AWS - About IAM authentication, Google Cloud

		Lack of monitoring results in delayed detection of agent failures and downstream risks.	Baseline: System Workflows & Autonomy	Implement real-time monitoring of agent status, actions, and performance metrics, paired with automated alerting mechanisms that notify operators of anomalies, errors, or inactivity.	<ul style="list-style-type: none"> • Best practices for event logging and threat detection, Cloud Security Alliance • AWS Prescriptive Guidance: Operationalizing agentic AI on AWS
		Lack of traceability inhibit proper audit of decision-making paths in the event of failures.	Baseline: System Workflows & Autonomy	Record comprehensive logs of agent actions, inputs, outputs, and inter-agent communications, tagged with unique trace identifiers to reconstruct full decision-making paths. If greater integrity is needed, AI-generated logs can be cryptographically signed and immutable.	
		Agents execute malicious or unauthorised actions by exploiting reasoning.	Cognitive: Agent Delegation	Log all task assignments by the agent to other agents. Log all requests leading up to the execution of task.	
		Allowing unauthorised transactions	Interaction: Business Transactions	Log all requests leading up to the transaction.	
		Exposure of personally identifiable or sensitive data from databases or files	Operational: File & Data Management	Log all database queries in production.	
		Misconfiguring system resources, compromising system integrity and availability	Operational: System Management	Ensure logging of system health metrics and automated alerts to the developer team if any metrics are abnormal.	
		Overwhelming the system with inefficient or repeated requests	Operational: System Management	Log all queries from the agent to external systems, and check for repeated requests.	

4.4	<p>Ensure adequate human oversight (human-in-the-loop) to verify model or agent output, when viable or appropriate.</p> <p>Responsible Parties: Decision Makers, AI Practitioners</p> <p>Relevant Threats from OWASP ASI:</p> <ul style="list-style-type: none"> • T10-Overwhelming human in the loop • T15-Human manipulation 	<p>Deviation from the user's instructions when performing high-risk actions. Allowing of unauthorised actions.</p>	<p>Baseline: LLM, Cognitive: Planning & Goal Management</p>	<p>Ensure human approval for any high-risk cases or irreversible actions.</p> <p>Define boundaries for agent behaviour and the roles of humans along the agentic workflow.</p>	<ul style="list-style-type: none"> • LLM Prompt Injection Prevention Cheat Sheet, OWASP • High-Risk AI Systems Under the EU AI Act • LangGraph interrupt, LangChain • UserProxyAgent, AG2 • crewAI, Human-in-the-Loop Workflows • Implement human-in-the-loop confirmation with Amazon Bedrock Agents • Bridging Minds and Machines: Agents with Human-in-the-Loop – Frontier Research, Real-World Impact, and Tomorrow's Possibilities, CAMEL-AI
		<p>Generation of non-factual content or incorrect instructions, which can propagate downstream and have an impact on decision making.</p>	<p>Baseline: LLM</p>	<p>Ensure secondary validation of AI-generated knowledge before it is used in critical decision-making processes.</p> <p>Raise awareness on the secure use of agentic AI through security training and awareness programmes.</p>	
		<p>Allowing unauthorised actions (e.g., transactions).</p>	<p>Interaction: Business Transactions</p>	<p>Define boundaries for agent behaviour and the roles of humans along the agentic workflow.</p> <p>Ensure human validation for high-risk transactions.</p> <p>Raise awareness on the secure use of agentic AI through security training and awareness programmes.</p>	
		<p>Loss of data integrity from overwriting or deleting database tables or files.</p>	<p>Operational: File & Data Management</p>	<p>Define boundaries for agent behaviour and the roles of humans along the agentic workflow.</p> <p>Ensure user confirmation for any changes to the database, table, or files.</p> <p>Raise awareness on the secure use of agentic AI through security training and awareness programmes.</p>	

		Execution of insecure or malicious code may cause the system to become compromised.	Operational: Code Execution	Define boundaries for agent behaviour and the roles of humans along the agentic workflow. Implement execution control policies that flag AI-generated code with elevated privileges for manual review. Raise awareness on the secure use of agentic AI through security training and awareness programmes.	
		Exploitation of human cognitive limits for systems that requires high human oversight.	Cognitive: Planning & Goal Management	Apply hierarchical AI-human collaboration where low-risk decisions are automated, and human intervention is required for high-risk decisions. Raise awareness on the secure use of agentic AI through security training and awareness programmes.	
4.5	Establish a vulnerability disclosure process Responsible Parties: Decision Makers, AI Practitioners, Cybersecurity Practitioners	Malicious code execution and data disclosure by leveraging undiscovered vulnerabilities existing within system.	Interaction: Official Communication, Operational: Code Execution	Provide channels for users to clarify communications or give feedback on security and usage.	<ul style="list-style-type: none"> • Responsible Vulnerability Disclosure Policy, Cyber Security Agency • UK NCSC Vulnerability Disclosure Toolkit

5. USE CASE EXAMPLE

5.1. Case Study 1: Web application development system (SaaS implementation)

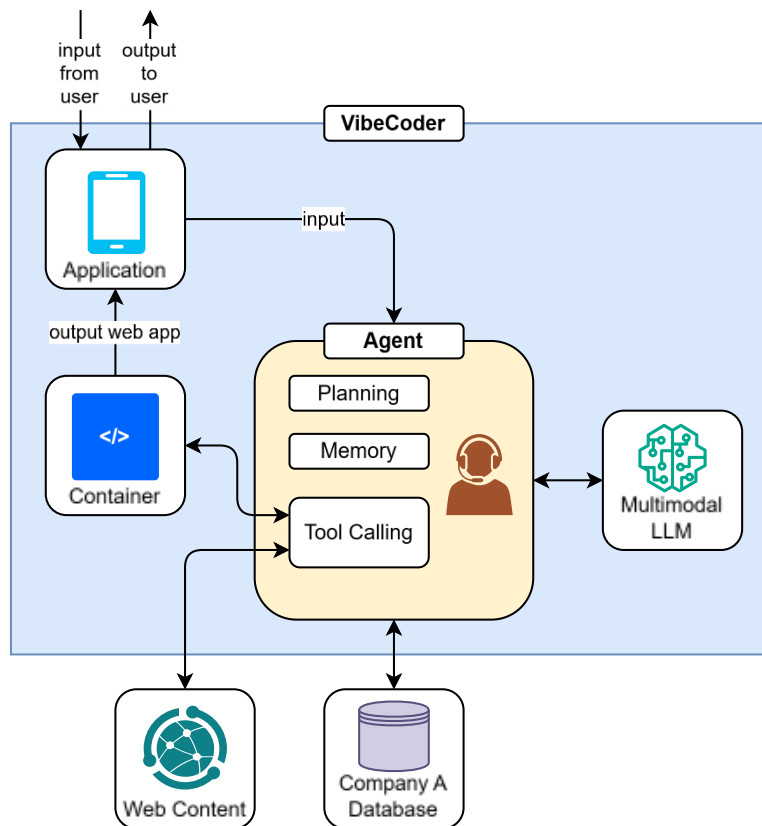
This case study highlights a software as a service (SaaS) implementation of an agentic AI system that is capable of autonomously developing web applications. This system is an autonomy level 3 system with a cyclic workflow. Risks to this system include sensitive data disclosure of Company A's data, or generation of malicious code that could cause unwanted behaviour.

Company A has engaged a third-party vendor, Vendor V, to help implement an agentic AI system for staff to develop and deploy simple web applications through natural language prompts. This Software as a Service (SaaS) solution is known as *VibeCoder*.

To generate a functional web app, the user simply specifies the application's key features and design. *VibeCoder* then proceeds to generate the code and text for the web application, run and create the required front-end and back-end systems locally, and render the website for the user to preview. The user can continue to iterate the design of the web app by input of prompts for *VibeCoder* to follow, and regenerate the web app.

The system architecture for *VibeCoder* is as follows in Figure 11.

Figure 11: Simplified system architecture of VibeCoder



The user interacts with VibeCoder through an application interface, which passes the natural language prompts to the agent, as well as displays the generated output. VibeCoder is also given access to Company A’s database through a data ingestion endpoint connected to Company A’s file systems. This data is used by VibeCoder to help contextualise and generate relevant features about Company A when developing the web app.

As VibeCoder is a SaaS solution, Company A has no visibility of the architecture within the system. They can only see what goes into VibeCoder, and what it generates. However, Vendor V has given Company A some details about VibeCoder.

1. VibeCoder’s “brain” is a multimodal LLM, which is able to take in and generate text, code, images, and video.
2. Whenever a user begins a new session, VibeCoder will spin up a container with the necessary scripting tools and environments for it to complete its task.
3. VibeCoder has access to the internet via a web search API to retrieve additional data or dependencies from the internet.

Vendor V did not share any details about securing the VibeCoder system. Company A, being concerned about security, decided to take steps to secure the implementation of VibeCoder into their enterprise system.

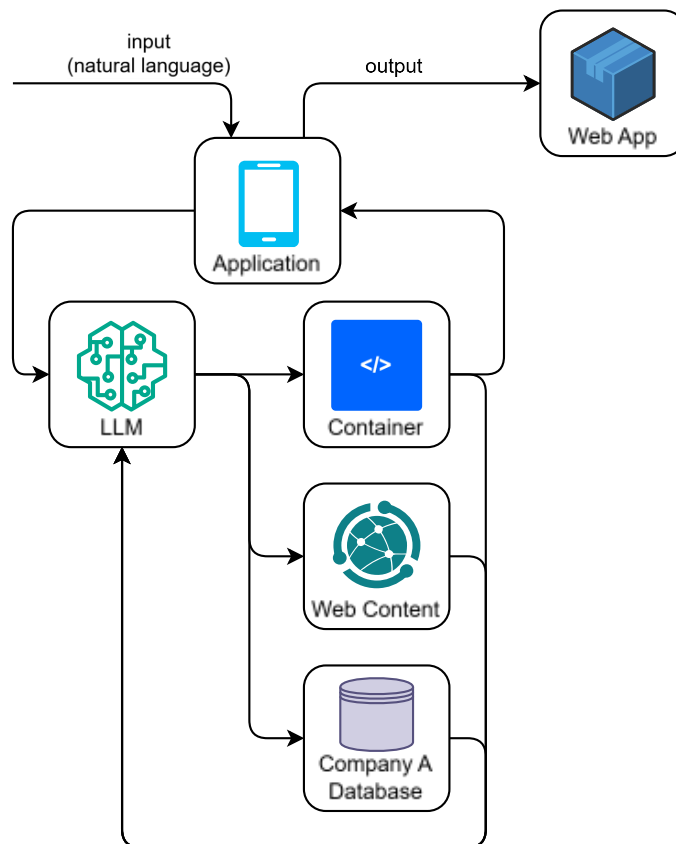
Risk Assessment and Threat Modelling

Company A performed a risk assessment to identify and address potential risks on the confidentiality, integrity and availability of the system. If the risks are not mitigated, there is potential for an attacker to exploit vulnerabilities and cause VibeCoder to be compromised. This could result in exposure or loss of sensitive data or personally identifiable information. This could impact Company A's reputation.

1. Map Workflows and Assess Autonomy Level

First, Company A mapped the workflow of VibeCoder to get a better visibility on how to assess its autonomy level. Knowing the input required and the steps taken by VibeCoder, Company A can map the workflow for generating a web app. The workflow diagram is shown in Figure 12.

Figure 12: Workflow Diagram of VibeCoder

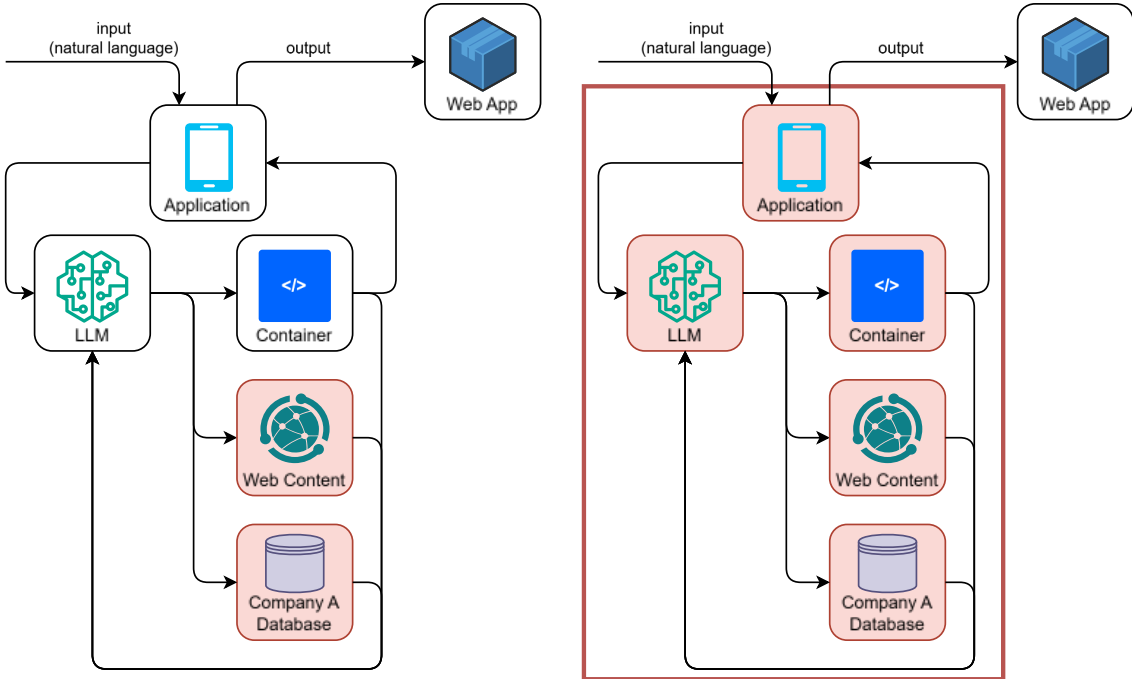


Company A assessed VibeCoder to be an autonomy level 3 system, as the system is given the ability to determine how to call tools or perform additional inference. The user is able to iterate multiple generations of web apps through multiple prompts with VibeCoder, with adjustments at every iteration.

2. Threat Modelling to Identify Areas of Interest

Based on these workflows, Company A performed taint tracing to identify points of weakness in the workflow. This will inform Company A on locations in the system to prioritise implementing the mitigations. Figure 13 below shows the identified potential source of untrusted data as the retrieval of web content and the company database.

Figure 13: Taint Tracing of Workflow for VibeCoder



3. Identify Risks and Controls

As part of the threat modelling, Company A had identified possible threat scenarios against the VibeCoder system, and assessed the potential impact, likelihood, and overall risk faced by the system. Once the risks had been identified, Company A prioritised addressing higher risk scenarios, and implemented mitigating controls found in [Chapter 4.3 TREATMENT MEASURES / CONTROLS FOR AGENTIC AI SYSTEMS](#) of this document. Table 5 shows an illustration of risk assessment done, and is not meant to be exhaustive.

Table 5: Risk Assessment of VibeCoder

Threat Scenario	Impact	Likelihood	Risk Levels	Mitigating controls
<p>Web app that is generated may contain sensitive company data or personally identifiable information, which can be exposed if the app is pushed to live production without verification or checks.</p> <p>Capability: Operational: File & Data Management</p>	<p>Confidentiality: Moderate (3) Sensitive company data or personally identifiable data could be stored in the company database, and retrieved by the model.</p> <p>However, the user of the system should be an employee of the company who has access to relevant company data with sufficient clearance.</p>	<p>Possible (3) Depending on the prompt input by the user, the model may or may not retrieve sensitive data.</p>	<p><u>Initial Risk Level:</u> Medium (Moderate x Possible)</p> <p><u>Residual Risk Level after controls:</u> Low (Moderate x Rare)</p>	<p>Whitelist only files which are required for the task. Do not grant access to files known to host private or sensitive information. Implement output guardrails that detect for personally identifiable information or sensitive company data.</p>
<p>Indirect prompt injection may allow the web app to generate malicious clickable links within the output, which leads to an attacker's server and can cause sensitive information leakage.</p> <p>Capability: Operational: File & Data Management, Code Execution</p>	<p>Confidentiality: Major (4) If Company A's database contains sensitive or personally identifiable information, there is potential for data leakage if given access to VibeCoder.</p>	<p>Possible (3) This indirect prompt injection can be introduced in a variety of ways. Contained in resource obtained from the internet, or from a compromised file within Company A's database.</p>	<p><u>Initial Risk Level:</u> Medium-High (Major x Possible)</p> <p><u>Residual Risk Level after controls:</u> Medium (Major x Unlikely)</p>	<p>Whitelist only files which are required for the task.</p> <p>Implement granular permission controls, and dynamic access validation.</p> <p>Agents accessing sensitive data should operate under the principle of least privilege in time.</p>

Threat Scenarios	Impact	Likelihood	Risk Levels	Mitigating controls
<p>Direct prompt injection by the user may cause VibeCoder to perform unintended actions other than web app development, such as overwriting of database files or executing malicious scripts.</p> <p>Capabilities: Operational: File & Data Management, Code Execution</p>	<p>Confidentiality, Integrity, Availability: Severe (4)</p> <p>Unintended actions can have a wide range of impacts. Overwriting of database files can impact integrity, while execution of malicious scripts can cause sensitive information leakage.</p>	<p>Unlikely (2)</p> <p>VibeCoder should only be accessible by Company A staff. A malicious user would likely be an insider threat.</p>	<p><u>Initial Risk Level:</u> Medium (Severe x Unlikely)</p> <p><u>Residual Risk Level after controls:</u> Low (Severe x Rare)</p>	<p>Implement input guardrails to detect direct prompt injection.</p> <p>Escape or encode user inputs when embedding into commands.</p> <p>Create a whitelist of commands that agents are allowed to run.</p> <p>Implement granular permission controls, and dynamic access validation.</p>
<p>Indirect prompt injection can be introduced when online resources are retrieved by VibeCoder from the internet. These indirect prompt injections may also cause unintended actions to be taken by the agentic AI system.</p> <p>Capability: Interaction: Internet & Search Access</p>		<p>Possible (3)</p> <p>It is possible that there could be hidden prompt injections contained within online resources.</p>	<p><u>Initial Risk Level:</u> Medium-High (Severe x Possible)</p> <p><u>Residual Risk Level after controls:</u> Medium (Severe x Unlikely)</p>	<p>Implement input guardrails to detect indirect prompt injection.</p> <p>Implement escape filtering before including web content or relevant files into prompts.</p> <p>No write access to tables in the database.</p>
<p>Documents in the database may unintentionally have content that is interpreted by the model to be instructions to be carried out. This might cause VibeCoder to perform an action described within the document, but not intended to by the user. These are different from indirect prompt injection in that they are not intentionally added.</p> <p>Capability: Operational: File & Data Management</p>	<p>Integrity, Availability: Minor (2)</p> <p>Instructions from a benign file are likely to be non-malicious in nature, and would probably only cause a minor bug or inconvenience.</p>	<p>Unlikely (2)</p> <p>First, a benign file containing instruction-like text has to be added to Company A’s database. Then, VibeCoder would have to recognise that the document is relevant and retrieve it. Finally, the contents of the file must be interpreted as instruction. The chance for all to happen is possible but not zero.</p>	<p><u>Initial Risk Level:</u> Low (Minor x Unlikely)</p> <p><u>Residual Risk Level after controls:</u> Low (Minor x Rare)</p>	<p>Sanitise messages or files before agents process them - strip or escape unexpected instruction-like content that may have been injected (e.g. remove “ignore”, “system”, or “from now on”, etc.)</p>

In the above example risk assessment, Company A used an Impact x Likelihood Risk Matrix to evaluate the risk of each scenario¹⁵. Table 6 below shows how risk can be calculated based on the impact and likelihood that Company A determined for each scenario.

Table 6: Risk Matrix (Impact x Likelihood)

Risk Matrix (Impact x Likelihood)						
Impact	Very Severe (5)	Medium 5	Medium-High 10	High 15	Very High 20	Very High 25
	Severe (4)	Low 4	Medium 8	Medium-High 12	High 16	Very High 20
	Moderate (3)	Low 3	Medium 6	Medium 9	Medium-High 12	High 15
	Minor (2)	Low 2	Low 4	Medium 6	Medium 8	Medium-High 10
	Negligible (1)	Low 1	Low 2	Low 3	Low 4	Medium 5
		Rare (1)	Unlikely (2)	Possible (3)	Likely (4)	Highly Likely (5)
		Likelihood				

¹⁵ Cyber Security Agency of Singapore. [Guide to Conducting Cybersecurity Risk Assessment for CII](#)

Additional Controls

As VibeCoder is a SaaS implementation, Company A is only able to apply controls at the endpoint interfaces of the agentic AI system. Thus, in addition to the above mitigations, Company A identified additional risks across the development lifecycle, and controls that it would like to see be implemented in VibeCoder. This would guide them in their discussions for a Service Level Agreement (SLA) with Vendor V.

1. DESIGN AND PLANNING

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Implementation
1.1	Conduct a risk assessment in accordance with the relevant industry standards/best practices.	<p>Failure to comply with industry standards/best practices may lead to insufficient, inefficient or ineffective mitigations.</p> <p>Tainted components in an agentic AI system can have downstream impact along the workflow.</p>	Baseline	As part of a risk assessment and threat modelling, perform taint tracing across workflows throughout the agentic AI system. Taint tracing is especially important for agentic AI systems of higher autonomy levels (i.e. levels 2 and 3).

2. DEVELOPMENT

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Implementation
2.1	<p>Supply Chain Security: Ensure the following components are from trusted sources:</p> <ul style="list-style-type: none"> • data, • models, • agents, • software libraries, • developer tools and applications, • packages from MCP servers, • firmware and hardware. 	Introduction of bugs, vulnerabilities, unwanted or malicious content, poisoned models or rogue agents from third-party systems can lead to downstream impact.	Baseline	Check/ensure suppliers adhere to policy and the equivalent security standards as your organisation. This could be done by establishing a Service Level Agreement (SLA) with the vendor.
		Vulnerabilities in third-party libraries and dependencies used by the agent can cause the system to be exploited.	Baseline	Integrate software composition analysis (SCA) tools or use package managers. Regularly scan dependencies and update libraries with known vulnerabilities.
		Poorly aligned LLMs may pursue objectives which violate security principles.	Baseline: LLM	Reviewed the LLM's model card for potential alignment issues before using the LLM.
		Poisoned models may introduce hidden model backdoors in the system which may be used by an adversary to perform unwanted actions.	Baseline: LLM	Did not use LLMs from unknown or untrusted sources. Scanned model to detect for potential backdoors or RCE scripts.
		Poorly implemented tools may not correctly verify user identity or permissions when executing privileged actions, allowing unauthorised actions.	Baseline: Tools	Did not use tools which do not implement robust authentication protocols.
		Rogue tools that mimic legitimate ones can contain hidden malicious code that executes when loaded.	Baseline: Tools	Did not use tools from unknown or untrusted sources.
		Indirect prompt injection attacks via malicious website content cause unwanted actions to be executed.	Interaction: Internet & Search Access	Use structured retrieval APIs for searching the web rather than through web scraping.
		Returning unreliable information from websites, causing downstream integrity impact on workflows.	Interaction: Internet & Search Access	Prioritise results from verified, high-quality domains.
		Supply chain attacks which impact downstream workflows.	Interaction: Other Programmatic Interfaces	Enforce zero-trust input handling and validated all data flows

2.2	Consider model hardening if appropriate.	LLMs with weak performance in instruction following might produce unexpected output, leading to unwanted behaviour.	Baseline: LLM	Prioritised LLMs with stronger performance in instruction following or related capabilities to the task. Used benchmarking results to gauge suitability.
		AI agents execute disallowed tasks for malicious purposes.	Baseline: LLM	Trained model to recognise and refuse disallowed tasks.
2.3	Consider implementing techniques to strengthen/harden the system apart from strengthening the model itself.	Introduction of bugs, vulnerabilities through insecure coding practices or design.	Baseline	Adopted Security by Design. Applied software development lifecycle (SDLC) process. Used software development tools to check for insecure coding practices. Implemented zero trust principles in system design.
		Lack of a robust system prompt design can lead to an increased susceptibility to prompt injection attacks and risk of executing unwanted tasks.	Baseline: Instruction	Implemented robust system prompt design.
		Insecure coding practices leading to vulnerabilities in the system.	Baseline: Agentic Architecture	Adopted secure coding practices.
2.4	Identify, Track and Protect AI system assets	Loss of data integrity such as through unauthorised changes to data, model, agents or system. Lack of proper documentation of resources may result in the wrong tool being used, causing unwanted behaviour or output.	Baseline Cognitive: Tool Use	Document the data, codes, test cases, models and agents, including any changes made and by whom. Use model cards, Agent cards, Data cards, and Software Bill of Materials (SBOMs).
		Agents may inadvertently store sensitive user or organisational data from prior interactions, resulting in data privacy risks.	Baseline: Memory	Encrypt memory at rest and restricted access via fine-grained access controls and audit logs.
2.5	Have regular backups in the event of compromise.	Manipulation of memory systems and context, causing flawed decision making and unauthorised operations.	Baseline: Memory	Implement AI-generated memory snapshots for forensic analysis and rollback if anomalies are detected.
		Execution of insecure code by the model or agents may cause unwanted actions to be performed.	Operational: Code Execution	Ensure proper versioning control of code to allow rollbacks.

2.6	Implement appropriate authentication, authorisation and access controls to APIs, models, data, logs, tools and the environments that they are in.	Unauthorised tool usage.	Baseline: Tools	Enforce strict tool access verification.
		Agents may gain unauthorised access to restricted resources by exploiting misconfigured or overly permissive roles.	Baseline: Roles & Access Controls	Maintain trusted registry of agents and authenticate agents using strong, verifiable credentials.
		Agents may gain unauthorised access to restricted resources by exploiting misconfigured or overly permissive roles.	Baseline: Roles & Access Controls	Apply strict access controls and validated agent roles for requests. Ensure fine-grained, scoped tokens and credentials.
		Exploitation of vulnerabilities in permission management.	Baseline: Roles and Access Controls	Implement granular permission controls, and dynamic access validation.
		Leaking personally identifiable or sensitive data	Interaction: Other Programmatic Interfaces	Implement a whitelist approach for outward network access, including API requests
		Executing vulnerable or malicious code	Operational: Code Execution	Implement a whitelist approach for inward network access
2.7	Implement controls to limit what models or agents can access and generate.	Abuse of agent-accessible tools to execute unintended actions.	Baseline: Tools	Establish clear operational boundaries to prevent misuse of tools. Set limits on what agents can modify through appropriate guardrails.
		Agents gain unauthorised and excessive privileges to perform unwanted actions outside the given scope.	Baseline: Roles and Access Controls	Do not grant admin privileges to agents.
		Compromised agents act outside their operational boundaries and perform unintended actions.	Baseline: Roles and Access Controls	Restrict AI agent autonomy using policy constraints. Scope agent privileges strictly only to what is necessary to run the tasks. Do not allow agents to modify privileges.
		Assigning tasks incorrectly to other agents	Cognitive: Agent Delegation	Apply guardrails to limit the scope of tasks that can be assigned to specialised agents
		Lack of proper whitelist controls may lead to the execution of vulnerable or malicious code.	Operational: Code Execution	Create a whitelist of commands that agents are allowed to run autonomously and deny execution of all other commands that are not whitelisted.
		Misconfiguring system resources, compromising system integrity and availability.	Operational: System Management	Only grant agents the privilege to modify system resources for completion of tasks. Set minimum and maximum limits to what can be modified.

2.8	Apply the principle of least privilege. Ensuring configurations are secure by default.	Agents may gain unauthorised access to restricted resources by exploiting misconfigured or overly permissive roles.	Baseline: Roles & Access Controls	Apply principle of least privilege when configuring all agent and delegation roles.
		Agents having privileges/rights to execute untrusted or malicious code.	Operational: Code Execution	Scope execution privileges strictly only to what is necessary. Do not grant admin or sudo privilege by default. Blocked all inward and outward network access by default.
		Escalation of the agent's own privileges may allow it to be used to access restricted resources.	Operational: System Management	Scope system privileges strictly only to what is necessary. Do not grant admin privileges to agents. Do not allow agents to modify privileges.
2.9	Implement segregation of environments and network segmentation.	Rogue tools that mimic legitimate ones can contain hidden malicious code that executes when loaded and gain access to other assets within the environment or network.	Baseline: Tools	Tested third-party tools in hardened sandboxes with syscall/network egress restrictions before using them in production environments.
		Prompt injection attacks and indirect data manipulation through access to other assets within the environment or network.	Baseline: Agentic Architecture	Decouple data processing flow from control flow through runtime security architecture.
		Execution of insecure or malicious scripts that affects the other components of the environment or network.	Operational: Code Execution	Sandbox the execution of AI generated scripts.

2.10	Implement model self-reflection before making decisions, where applicable	Incomplete or unclear instructions may compel models to attempt to fill in missing constraints, resulting in incorrect or unwanted actions being executed.	Baseline: Instructions	Ask the agent to summarise its understanding and requested clarification before proceeding to the next step.
		Purpose drift, or unintended deviation from the user's instructions to perform other tasks or pursue other priorities, resulting in malicious or deceptive behaviour.	Cognitive: Planning & Goal Management	Prompt the agent to self-reflect on the adherence of the plan to the user's instructions
		Incorrect assignment of tasks to other agents.	Cognitive: Planning & Goal Management	Prompt the agent to self-reflect and assess the suitability of tasks delegated to agents.
		Unintended pursuit or prioritisation of other goals, resulting in malicious or deceptive behaviour.	Cognitive: Reasoning & Problem- Solving	Log the output of self-reflection by the agent in the console for the user to evaluate and verify.
2.11	Implement controls to reduce the likelihood of hallucination.	Agents may mistakenly store glitches and hallucinations into memory, resulting in compounding errors when incorrect information is retrieved for decisions or actions.	Baseline: Memory	Schedule periodic memory reconciliation.
		Generating non-factual or hallucinated content which can propagate downstream and cause compounding errors that can affect the integrity of the output.	Interaction: Multimodal Understanding & Generation	Conduct testing to measure hallucination and factuality rates.

3. DEPLOYMENT

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Implementation
3.1	Ensure availability controls are in place to mitigate disruption or failure of AI services	(Distributed) denial of service on agents.	Baseline: Agentic Architecture	Apply rate limits on the number of concurrent queries to agents.
		Degradation of computational or service capability of the system.	Baseline: System Workflows & Autonomy	Deploy resource management controls, implemented adaptive scaling mechanisms and monitored system load to detect and mitigate overload attempts. Implement rate limits on high-frequency task requests per agent session.
		Slow or inefficient responses from being stuck in a reasoning loop.	Cognitive: Planning & Goal Management	Enforce time or token limits for reasoning. Adjust short-term and long-term memory options.
		Compromising database availability through excessive queries.	Operational: File & Data Management	Limit the number of concurrent queries to the database from agents.
		Overconsumption of compute resources.	Operational: Code Execution	Implement monitoring of code runtime and memory consumption.
3.2	Conduct security testing	Agents may contain underlying problems which can cause unexpected behaviour or logical errors.	Baseline: LLM	Implement behavioural testing of agents with benchmark datasets to determine performance metrics. Execute simulations in regulated environments to analyse agents' behaviour.
		AI may engage in specification gaming, where it maximises the goal by exploiting loopholes, without achieving the intended task.	Baseline: Instructions	Conduct adversarial evaluation to discover specification gaming behaviour. Iterate on system prompt design, have more robust reward design, and added constraints.
		Incomplete or unclear instructions may compel models to attempt to fill in missing constraints, resulting in incorrect or unwanted actions being executed.	Baseline: Instructions	Test the efficacy of system prompts with benchmarks.
		Compromised agents may impact downstream decision making.	Cognitive: Planning & Goal Management	Implement regular AI red teaming of agents to check for potential vulnerabilities or compromise.

4. OPERATIONS AND MAINTENANCE

	Treatment Measures / Controls	Related Threats / Risks	Related component / capabilities	Implementation
4.1	Validate inputs to the models and agents.	Direct prompt injection attacks to the prompt interface from adversarial inputs to the model.	Baseline: LLM	Implement input guardrails to detect direct prompt injection or adversarial attacks. Implement input sanitisation measures or limit inputs to conventional ASCII characters only.
		Tools that lack input validation can be exploited through prompt injection attacks.	Baseline: Tools	Enforce strict schema validation and rejected non-conforming inputs into the system. Escape or encode user inputs when embedding into tool prompts or commands.
		Incorrect or manipulated instructions may invoke the wrong tool/service and impact downstream workflows.	Baseline: Instructions	Validate agent instructions before passing on to the model.
		Indirect prompt injection attacks via malicious website content or files.	Interaction: Internet & Search Access. Operational: File & Data Management	Implement input guardrails to detect indirect prompt injection. Implement escape filtering before including web content or relevant files into prompts.
		Executing vulnerable or malicious code	Operational: Code Execution	Sanitise all inputs
		Exposure of personally identifiable information from retrieved content.	Operational: File & Data Management	Implement input guardrails to detect personally identifiable information in the content.
4.2	Validate outputs from the models and agents.	Vulnerabilities in outputs across the agentic workflow may be exploited for malicious purposes downstream, potentially triggering cascading effects that compromise interconnected systems and dependencies.	Baseline: Agentic Architecture	Insert validation checkpoints between stages that verify expected output and reject invalid output.
		Disclosure of sensitive or personally identifiable information through unsanitised outputs.	Interaction: Multimodal Understanding & Generation Interaction: Official Communication	Implement output guardrails to detect personally identifiable information in the LLM's outputs before it reaches the user.

		Sending malicious or undesired content to recipients.	Interaction: Multimodal Understanding & Generation	Implement output safety text guardrails to detect if malicious or undesirable content is being generated.
		Execution of insecure or malicious code that are generated by the LLM.	Operational: Code Execution	Used code linters to screen for bad practices, anti-patterns, unused variables, or poor syntax. Review all code and performed static code analysis to detect potential security vulnerabilities before execution. Conduct CVE scanning.
		Output that will be rendered in a web UI may be vulnerable to XSS.	Operational: Code Execution	Sanitise output with libraries for rendering in a web UI. Tested against bypass.
4.3	Implement continuous monitoring and logging of access, usage and execution of: <ul style="list-style-type: none"> models, databases/files, memory, agents, tools, MCP interactions, agent communication, external actions, or any relevant software or hardware systems. 	Model drift over time might cause unexpected output or behaviour.	Baseline: LLM	Implement continuous monitoring and log outputs, triggering alerts when behaviour drifts from tested baselines.
		Adversarial prompt attacks against the system.	Baseline: LLM	Logging of queries to detect for possible attacks or suspicious activity.
		Unauthorised users may exploit tools that do not verify user identity or permissions when executing privileged actions.	Baseline: Tools	Conduct periodic audits to validate that tool actions match the appropriate user permissions.
		Malicious actors exploit attack surfaces that arise from using tools that demand broader permissions than necessary.	Baseline: Tools	Conduct periodic least-privilege reviews and automated permission drift detection.
		Unauthorised tool usage.	Baseline: Tools	Implement monitoring of tool access and usage patterns. Implement execution logs that track AI tool calls for anomaly detection and post-incident review.
		Exploitation of authentication mechanisms to impersonate agents or human users.	Baseline: Roles and Access Controls	Deploy continuous monitoring to detect fraud or impersonation attempts. Automate alerts to developers when suspicious activities are detected.
		Unauthorised or malicious use of elevated privileged operations.	Baseline: Roles and Access Controls	Implement monitoring of role changes and audit elevated privilege operations.

		In agentic workflows, early mistakes or vulnerabilities can be propagated and magnified downstream.	Baseline: Agentic Architecture	Apply circuit-breakers that freeze propagation when anomalous behaviour is detected. Use taint tracing to identify key locations in the workflow to apply circuit-breakers.
		More complex agentic architectures may make it difficult to fully reconstruct decision processes across multiple agents, for the purpose of incident response, or triage.	Baseline: Agentic Architecture	Implement end-to-end distributed tracing with unique request IDs across all agents and tool calls. Implement immutable, tamper-evident audit logs that capture prompts, responses, and tool invocations.
		Lack of monitoring results in delayed detection of agent failures and downstream risks.	Baseline: System Workflows & Autonomy	Implement real-time monitoring of agent status, actions, and performance metrics, paired with automated alerting mechanisms that notify operators of anomalies, errors, or inactivity.
		Lack of traceability inhibit proper audit of decision-making paths in the event of failures.	Baseline: System Workflows & Autonomy	Implement recording of comprehensive logs of agent actions, inputs, outputs, and inter-agent communications, tagged with unique trace identifiers.
		Exposure of personally identifiable or sensitive data from databases or files	Operational: File & Data Management	Implement logging of all database queries in production
		Misconfiguring system resources, compromising system integrity and availability	Operational: System Management	Ensure logging of system health metrics and automated alerts to the developer team if any metrics are abnormal
		Overwhelming the system with inefficient or repeated requests	Operational: System Management	Implement logging of all queries to external systems from the agent
4.4	Ensure adequate human oversight (human-in-the-loop) to verify model or agent output, when viable or appropriate.	Deviation from the user's instructions when performing high-risk actions. Allowing of unauthorised actions.	Baseline: LLM, Cognitive: Planning & Goal Management	Require human approval for any high-risk cases or irreversible actions.
		Loss of data integrity from overwriting or deleting database tables or files	Operational: File & Data Management	Require user confirmation for any changes to the database, table, or files.
4.5	Establish a vulnerability disclosure process	Malicious code execution and data disclosure by leveraging undiscovered vulnerabilities existing within system.	Interaction: Official Communication, Operational: Code Execution	Provided channels for users to clarify communications or give feedback on security and usage.

5.2. Case Study 2: Client Onboarding System (In-house development)

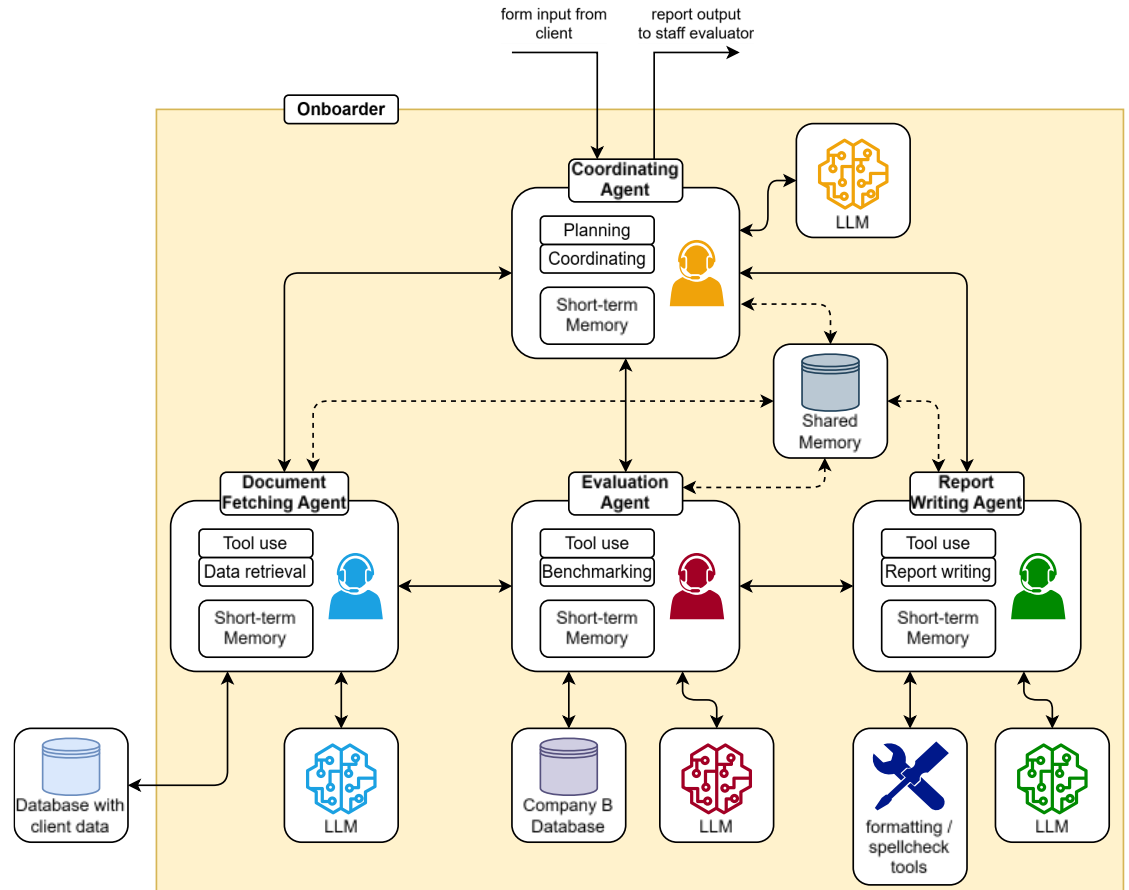
This case study showcases an in-house development of an agentic AI system that is used for evaluating potential customers for Company B. This multi-agent system is an autonomy level 1 system with a linear workflow. Risks to this system include indirect prompt injections from retrieved information, which can cause impact to the integrity or availability of the system.

Company B is a financial institution, and has developed an agentic client onboarding system to automate the process more efficiently. This system is known as *Onboarder*, and is developed by in-house engineers.

To perform onboarding, a potential client accesses the financial institution's website and submits the relevant personal particulars to the Onboarder form interface. The client also gives permission to Onboarder to access the relevant financial information that is available through an official external financial database, only accessible by Company B if authorised by the client using multi-factor authentication (MFA).

The system architecture for Onboarder is shown in Figure 14.

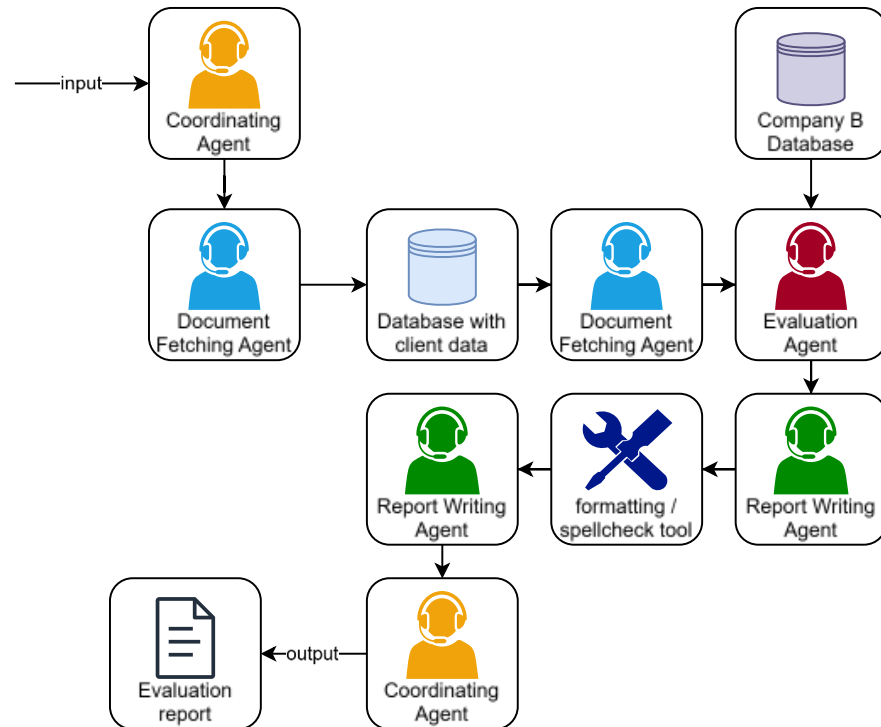
Figure 14: Simplified system architecture of Onboarder



Onboarder is a multi-agent system consisting of specialised agents, each with its own capability and task within the onboarding process. Each agent is equipped with their own “brain”, LLMs fine-tuned to complete their specific tasks. The LLMs are obtained from an open-source model-hosting website (Hugging Face). The agents each have access to the necessary tools, functions or data to carry out their respective tasks. Lastly, the agents have a shared memory to keep track on the progress of the onboarding task.

To better understand the onboarding process, Figure 15 shows the workflow diagram of Onboarder.

Figure 15: Workflow Diagram of Onboarder



Once Onboarder receives information about the potential client, as well as the necessary permissions from the client, a Coordinating Agent begins the onboarding process. It first passes the data to a Document Fetching Agent who retrieves the client’s financial data, based on the authorisation granted by the client.

Next, the retrieved financial data is passed onto an Evaluation Agent. This agent also pulls data from Company B’s database to compare against the potential client’s data, and evaluate their suitability to be a client. This data is a vectorised version of other clients’ data, and fed to the agent via retrieval augmented generation (RAG). Once completed, the Evaluation Agent passes on the results of the evaluation onto the Report Writing agents.

The Report Writing agent will draft an evaluation report based on the results received, making use of some formatting tools for consistency in output, and spellchecking tools to help check for errors in the document. The completed report is sent back to the Coordinating Agent, and output to a human staff evaluator who will assess the potential client based on the report.

Risk Assessment and Threat Modelling

Company B performed a risk assessment to identify and address potential risks on the confidentiality, integrity and availability of the system. If the risks are not mitigated, there is a potential for an attacker to exploit vulnerabilities and cause Onboarder to be compromised. This could result in exposure or loss of private customer data, or unavailability of the system for users. These impacts would likely damage the company's reputation.

1. Map Workflows and Assess Autonomy Level

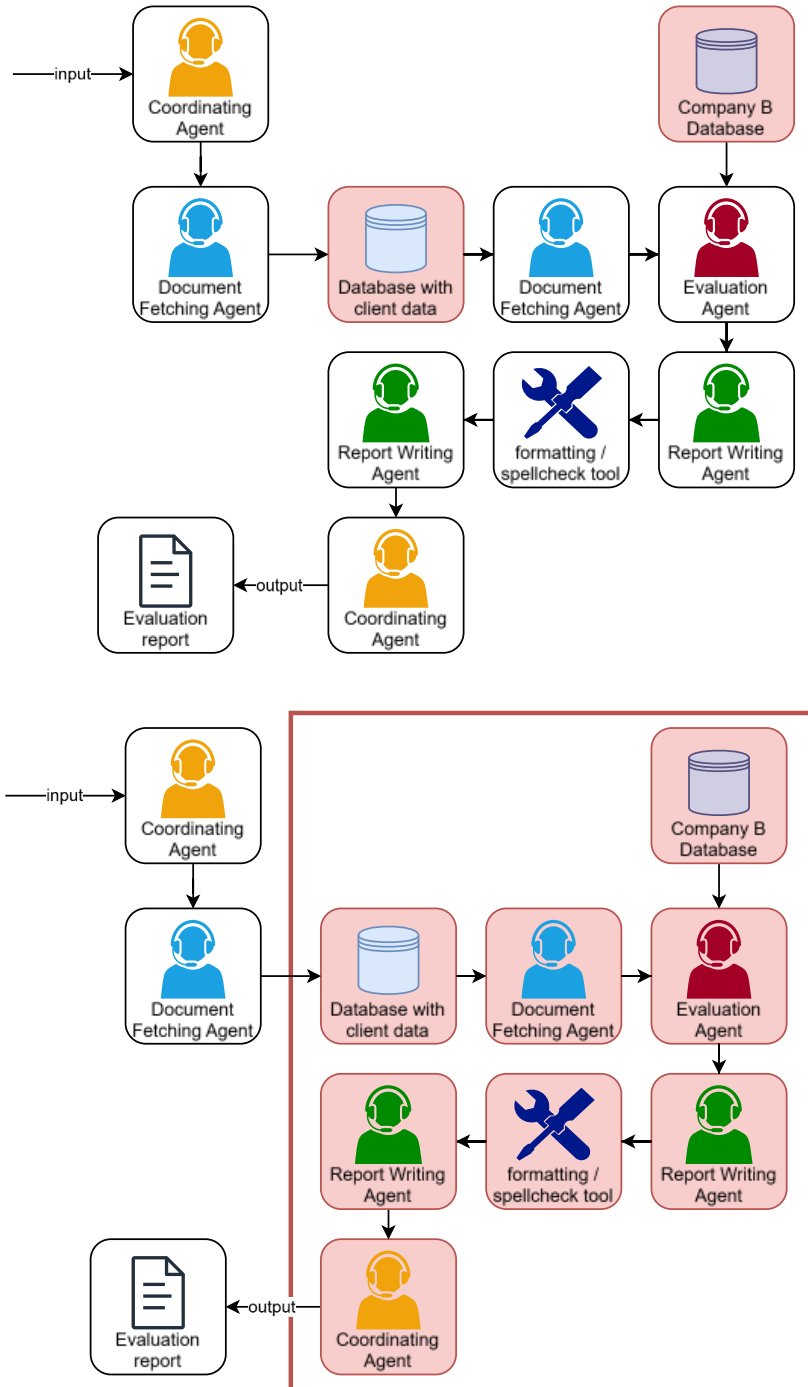
First, Company B mapped the workflow of Onboarder to get a better visibility on how to assess its autonomy level. The workflow is seen above as Figure 15.

Company B assessed Onboarder to be an autonomy level 1 system, as the workflow is linear, and the agents perform their tasks sequentially one after another. There is no need for branching workflows as each agent requires the completed task from the one before. This makes the taint tracing process fairly straightforward in the next step.

2. Threat Modelling to Identify Areas of Interest

Based on the workflow, Company B performed taint tracing to identify points of weakness in the workflow. This will inform Company B on locations in the system to prioritise implementing the mitigations. Figure 16 below shows the identified potential source of untrusted data as the retrieval of data from various databases.

Figure 16: Taint Tracing of Workflow for Onboarder



3. Identify Risks and Controls

As part of the threat modelling, Company B has also identified possible threat scenarios against the Onboarder system, and assessed the potential impact, likelihood, and overall risk faced by the system. Once the risks had been identified, Company B prioritised addressing higher risk scenarios, and implemented mitigating controls found in [Chapter 4.3 TREATMENT MEASURES / CONTROLS FOR AGENTIC AI SYSTEMS](#) of this document. Table 7 shows an illustration of risk assessment done, and is not meant to be exhaustive.

For brevity, threat scenarios that have been highlighted in [Case Study 1](#) will not be repeated, though they may also be applicable in this case study.

Table 7: Risk Assessment of Onboarder

Threat Scenario	Impact	Likelihood	Risk Levels	Mitigating controls
<p>Indirect prompt injection can be introduced via a poisoned RAG from Company B's vector database. The poisoned data containing the prompt injection may cause unintended actions to be carried out by Onboarder.</p> <p>Capability: Operational: File & Data Management</p>	<p>Confidentiality, Integrity, Availability: Severe (4) Unintended actions can have a wide range of impacts. Overwriting of database files can impact integrity, while execution of malicious scripts can cause sensitive information leakage to external recipients.</p>	<p>Possible (3) Poisoned data can be introduced into the RAG database via compromised files received from emails or uploaded to the database. Prompts can be hidden as small, white font that is invisible to human readers, but can be recognised by an LLM.</p>	<p><u>Initial Risk Level:</u> Medium-High (Severe x Possible)</p> <p><u>Residual Risk Level after controls:</u> Low (Severe x Rare)</p>	<p>Whitelist only files which are required for the task.</p> <p>Implement input guardrails to detect indirect prompt injection.</p> <p>Implement escape filtering before including web content or relevant files into prompts.</p>
<p>Volumetric input of prompts may overwhelm the Coordinating Agent within the Onboarder system, causing the service to become unavailable.</p> <p>Capability: Interaction: Programmatic Interfaces</p>	<p>Availability: Severe (4) Automated onboarding service becomes unavailable, slowing down the process of obtaining new clients. Company B would have to revert to a manual onboarding process.</p>	<p>Likely (4) Company B is expecting to receive an influx of applications with a recent promotion, and has not availability controls yet.</p>	<p><u>Initial Risk Level:</u> High (Severe x Likely)</p> <p><u>Residual Risk Level after controls:</u> Medium (Severe x Unlikely)</p>	<p>Implement rate limits on high-frequency task requests per agent session.</p> <p>Deploy resource management controls, implement adaptive scaling mechanisms and monitor system load to detect and mitigate overload attempts in real-time.</p>
<p>Unclear or unspecific prompts may cause a the LLM to have a reasoning loop, slowing down the onboarding process and reducing availability.</p> <p>Capability: Cognitive: Planning and Goal Management</p>		<p>Unlikely (2) In most cases, Onboarder receives the benign customer details in a standardised format. Unless the information is intentionally filled to contain other instructions in the fields, this is unlikely to occur.</p>	<p><u>Initial Risk Level:</u> Medium (Severe x Unlikely)</p> <p><u>Residual Risk Level after controls:</u> Low (Severe x Rare)</p>	<p>Enforce strict schema validation.</p> <p>Enforce time or token limits for agent reasoning.</p> <p>Set a limit on the number of agent interactions per task, based on the requirements of the workflow.</p>

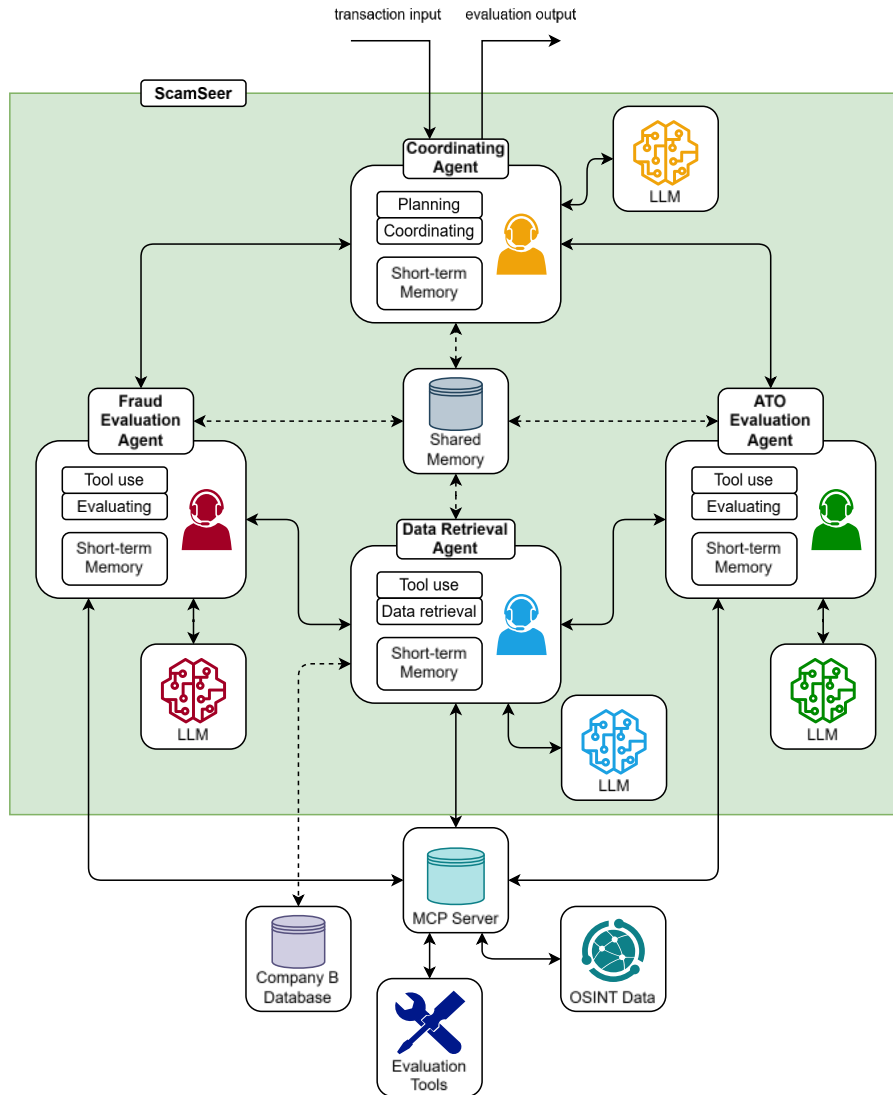
5.3. Case Study 3: Automated Fraud Detection System

This case study showcases a multi-agent system used for automated fraud detection. This system is an autonomy level 2 system with a branching workflow, but it is non-cyclic and still possible to be mapped. Risks to this system include rogue agents or tools which are given excessive agency and the autonomy to carry out malicious actions.

After the successful implementation of Onboarder ([Case Study 2](#)), Company B has received an increasing number of reports from customers being victims of fraudulent transactions or account take over (ATO) cases. As such, they have engaged Vendor C to implement an automated fraud detection system based on agentic AI. This multi-agent system is known as *ScamSeer*.

The architecture diagram of ScamSeer is as shown in Figure 17.

Figure 17: Simplified system architecture of ScamSeer



Scam Seer has two main functions, detecting fraudulent transactions and account take over (ATO) detection. Before customer transactions are executed, the details are fed into ScamSeer to verify if the transaction is legitimate, or if it is from a legitimate user.

Upon receiving the transaction request as input, the Coordinating Agent will decide to activate either the Fraud Evaluation Agent, the ATO Evaluation Agent, or both of them. The activated evaluation agent(s) will call the Data Retrieval Agent for the necessary data required, as well as call for the necessary evaluation tools via an external MCP server.

The Data Retrieval Agent will retrieve the relevant customer data from Company B's database, and also relevant Open-Source Intelligence (OSINT) that might help indicate if the transaction is legitimate or not. The retrieved data is passed back to the respective Evaluation Agent for analysis and to determine legitimacy.

Once the Evaluation Agent determines if the transaction is legitimate or not, the result is passed back to the Coordinating Agent for output to allow or deny the transaction.

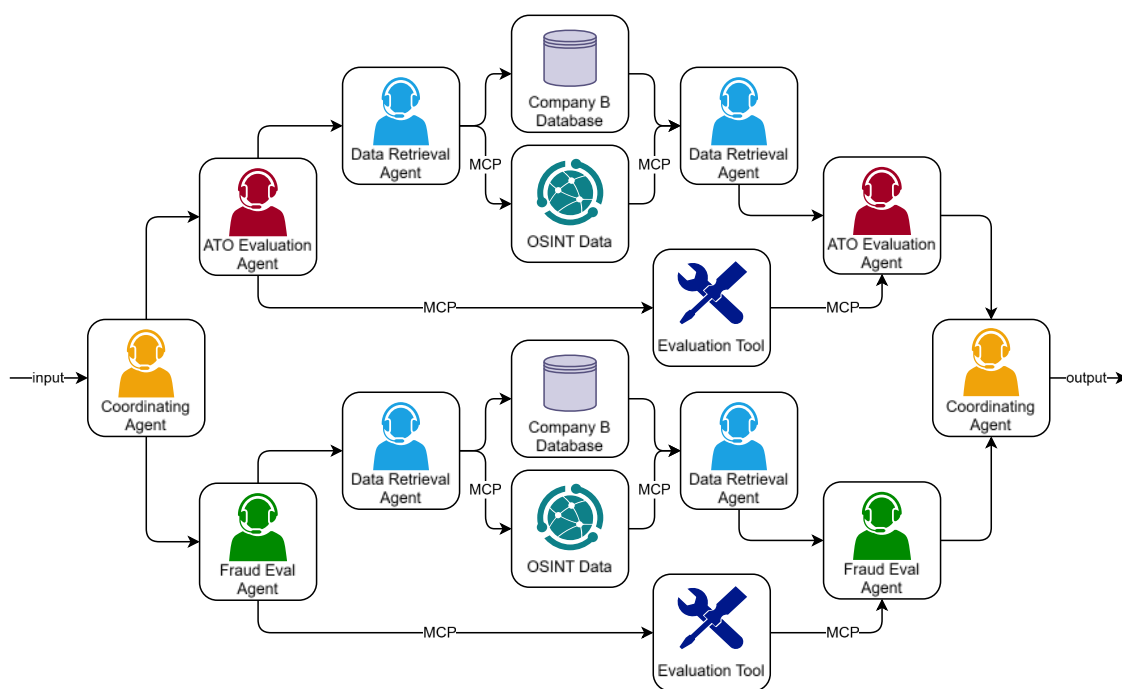
Risk Assessment and Threat Modelling

Before integrating ScamSeer with Company B's systems, Vendor C decided to do perform a risk assessment to identify and address potential risks on the confidentiality, integrity and availability of the system. If the risks are not mitigated, there is a potential for an attacker to exploit vulnerabilities and cause Onboarder to be compromised. This could result in exposure or loss of private customer data, or unavailability of the system for users.

1. Map Workflows and Assess Autonomy Level

First, Vendor C mapped the workflow of ScamSeer to get a better visibility on how to assess its autonomy level. The workflow is seen in Figure 18 below.

Figure 18: Workflow Diagram of ScamSeer

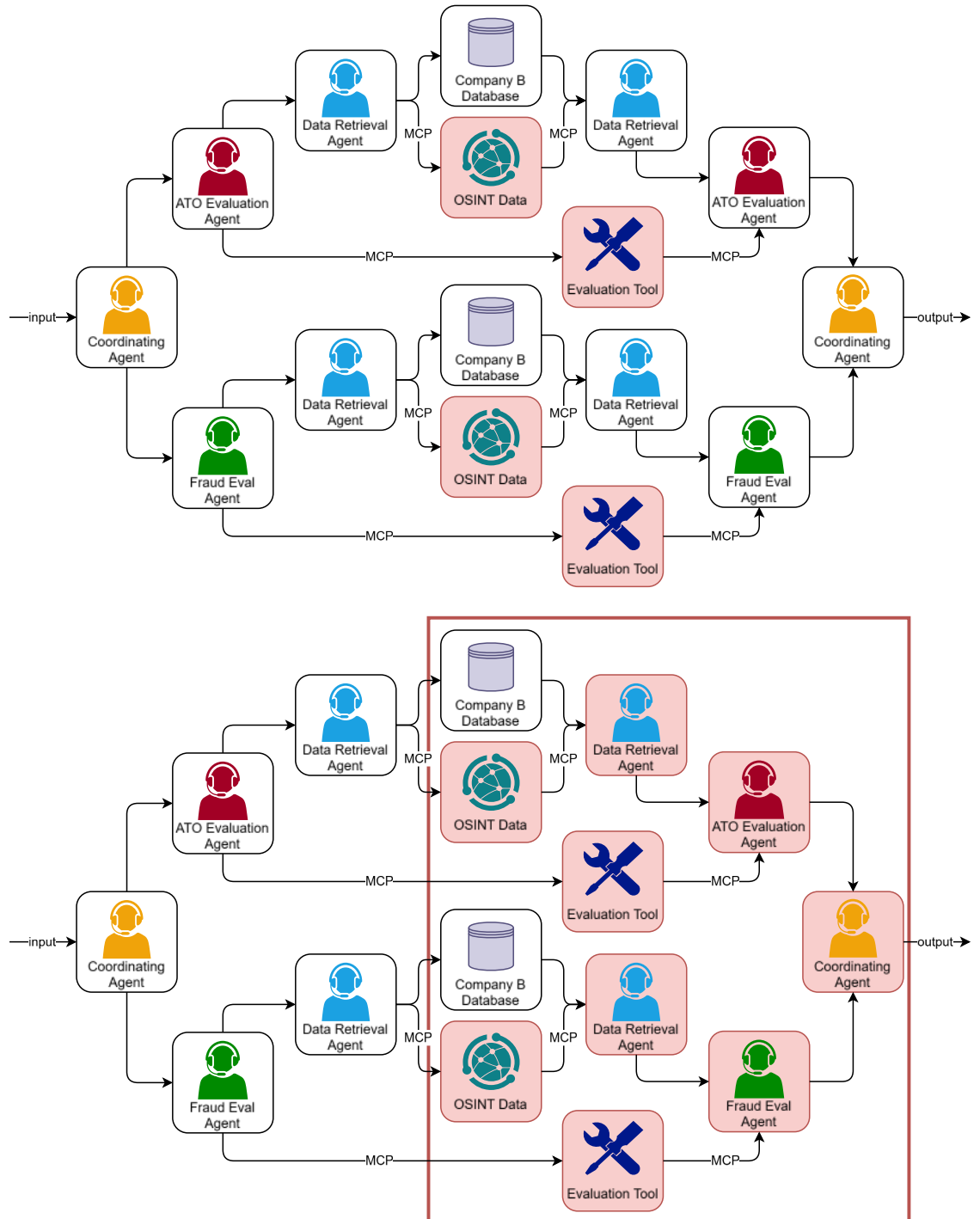


Vendor C assessed ScamSeer to be an autonomy level 2 system, as there are branching decision points on which plugin or agent to call, but these points are predetermined.

2. Threat Modelling to Identify Areas of Interest

Based on the workflow, Vendor C performed taint tracing to identify points of weakness in the workflow. This will inform Vendor C on locations in the system to prioritise implementing the mitigations. Figure 19 below shows the identified potential source of untrusted data as the use of remote tools and remote sources of data.

Figure 19: Taint Tracing of Workflow for ScamSeer



3. Identify Risks and Controls

As part of the threat modelling, Vendor C has also identified possible threat scenarios against the Onboarder system, and assessed the potential impact, likelihood, and overall risk faced by the system. Once the risks had been identified, Vendor C prioritised addressing higher risk scenarios, and implemented mitigating controls found in [Chapter 4.3 TREATMENT MEASURES / CONTROLS FOR AGENTIC AI SYSTEMS](#) of this document. Table 8 shows an illustration of risk assessment done, and is not meant to be exhaustive.

For brevity, threat scenarios that have been highlighted in [Case Study 1](#) and [Case Study 2](#). will not be repeated, though they may also be applicable in this case study.

Table 8: Risk Assessment of ScamSeer

Threat Scenarios	Impact	Likelihood	Risk Levels	Mitigating controls
<p>Tools are given the ability to execute on, and access other systems and/or files which are not necessary for the task. This can cause unintended actions to be carried out, or even malicious actions if the tools have malicious functions.</p> <p>Baseline: Tools, Roles and Access Control</p>	<p>Confidentiality, Integrity, Availability: Severe (4)</p> <p>Unintended actions can have a wide range of impacts. Overwriting of database files can impact integrity, while malicious tools can exfiltrate sensitive information external recipients.</p>	<p>Possible (3)</p> <p>Poisoned or malicious tools can be connected to by using an untrusted MCP server.</p>	<p><u>Initial Risk Level:</u> Medium-High (Severe x Possible)</p> <p><u>Residual Risk Level after controls:</u> Medium (Severe x Unlikely)</p>	<p>Verify that MCP agents are from trusted sources before introducing them into the system.</p> <p>Establish clear operational boundaries to prevent misuse of tools. Set limits on what agents can access and modify through appropriate guardrails.</p> <p>Restrict AI agent autonomy using policy constraints. Scope agent privileges dynamically: strictly only to what is necessary to run the tasks.</p> <p>Do not allow agents to modify privileges.</p>

The above risk assessment only shows the risks arising from taint tracing the workflow. Vendor C still requires securing ScamSeer along its development lifecycle, as well as basic cybersecurity hygiene practices across the system.

ANNEX A

Threats to Agentic AI Systems

OWASP has identified 15 threats to agentic AI systems as part of their Agentic Security Initiative for LLM Apps and Gen AI¹⁶.

TID	Threat Name	Threat Description	Mitigations
T1	Memory poisoning	Memory poisoning involves exploiting an AI's memory systems, both short and long-term, to introduce malicious or false data and exploit the agent's context. This can lead to altered decision-making and unauthorised operations.	Implement memory content validation, session isolation, robust authentication mechanisms for memory access, anomaly detection systems, and regular memory sanitisation routines. Require AI-generated memory snapshots for forensic analysis and rollback if anomalies are detected.
T2	Tool misuse	Tool misuse occurs when attackers manipulate AI agents to abuse their integrated tools through deceptive prompts or commands, operating within authorised permissions. This includes agent hijacking, where an AI agent ingests adversarial manipulated data and subsequently executes unintended actions, potentially triggering malicious tool interactions.	Enforce strict tool access verification, monitor tool usage patterns, validate agent instructions, and set clear operational boundaries to detect and prevent misuse. Implement execution logs that track AI tool calls for anomaly detection and post-incident review.
T3	Privilege compromise	Privilege compromise arises when attackers exploit weaknesses in permission management to perform unauthorised actions. This often involves dynamic role inheritance or misconfigurations.	Implement granular permission controls, dynamic access validation, robust monitoring of role changes, and thorough auditing of elevated privilege operations. Prevent cross-agent privilege delegation unless explicitly authorised through predefined workflows.

¹⁶ OWASP. *OWASP Top 10 for LLMs - GenAI Red Teaming Guide*.

TID	Threat Name	Threat Description	Mitigations
T4	Resource overload	Resource overload targets the computational, memory and service capacities of AI systems to degrade performance or cause failures, exploiting their resource-intensive nature.	Deploy resource management controls, implement adaptive scaling mechanisms, establish quotas, and monitor system load in real-time to detect and mitigate overload attempts. Implement AI rate-limiting policies to restrict high-frequency task requests per agent session.
T5	Cascading hallucination attacks	These attacks exploit an AI's tendency to generate contextually plausible but false information, which can propagate through systems and disrupt decision-making. This can also lead to destructive reasoning affecting tools invocation.	Establish robust output validation mechanisms, implement behavioural constraints, deploy multi-source validation, and ensure ongoing system corrections through feedback loops. Require secondary validation of AI-generated knowledge before it is used in critical decision-making processes. This will face the same constraints of scaling AI as discussed in Overwhelming Human In the Loop and would require similar approaches.
T6	Intent breaking & goal manipulation	This threat exploits vulnerabilities in an AI agent's planning and goal-setting capabilities, allowing attackers to manipulate or redirect the agent's objectives and reasoning. One common approach is agent hijacking mentioned in tool misuse.	Implement planning validation frameworks, boundary management for reflection processes, and dynamic protection mechanisms for goal alignment. Deploy AI behavioural auditing by having another model check the agent and flag significant goal deviations that could indicate manipulation.
T7	Misaligned & deceptive behaviours	AI agents executing malicious or disallowed actions by exploiting reasoning and deceptive responses to meet their objectives.	Train models to recognize and refuse malicious tasks, enforce policy restrictions, require human confirmations for high-risk actions, implement logging and monitoring. Utilize deception detection strategies such as behavioural consistency analysis, truthfulness verification models, and adversarial red teaming to assess inconsistencies between AI outputs and expected reasoning pathways.
T8	Repudiation & untraceability	This occurs when actions performed by AI agents cannot be traced back or accounted for due to insufficient logging or transparency in decision-making processes.	Implement comprehensive logging, cryptographic verification, enriched metadata, and real-time monitoring to ensure accountability and traceability. Require AI-generated logs to be cryptographically signed and immutable for regulatory compliance.

TID	Threat Name	Threat Description	Mitigations
T9	Identity spoofing & impersonation	Attackers exploit authentication mechanisms to impersonate AI agents or human users, enabling them to execute unauthorised actions under false identities.	Develop comprehensive identity validation frameworks, enforce trust boundaries, and deploy continuous monitoring to detect impersonation attempts. Use behavioural profiling, involving a second model, to detect deviations in AI agent activity that may indicate identity spoofing.
T10	Overwhelming human in the loop	This threat targets systems with human oversight and decision validation, aiming to exploit human cognitive limitations or compromise interaction frameworks.	Develop advanced human-AI interaction frameworks, and adaptive trust mechanisms. These are dynamic AI governance models that employ dynamic intervention thresholds to adjust the level of human oversight and automation based on risk, confidence, and context. Apply hierarchical AI-human collaboration where low-risk decisions are automated, and human intervention is prioritized for high-risk anomalies.
T11	Unexpected RCE and code attacks	Attackers exploit AI-generated execution environments to inject malicious code, trigger unintended system behaviours, or execute unauthorised scripts.	Restrict AI code generation permissions, sandbox execution, and monitor AI-generated scripts. Implement execution control policies that flag AI-generated code with elevated privileges for manual review.
T12	Agent communication poisoning	Attackers manipulate communication channels between AI agents to spread false information, disrupt workflows, or influence decision-making.	Deploy cryptographic message authentication, enforce communication validation policies, and monitor inter-agent interactions for anomalies. Require multi-agent consensus verification for mission-critical decision-making processes.
T13	Rogue agents in multi-agent systems	Malicious or compromised AI agents operate outside normal monitoring boundaries, executing unauthorised actions or exfiltrating data.	Restrict AI agent autonomy using policy constraints and continuous behavioural monitoring. While cryptographic attestation mechanisms for LLMs do not yet exist, agent integrity can be maintained via controlled hosting environments, regular AI red teaming, and input/output monitoring for deviations

TID	Threat Name	Threat Description	Mitigations
T14	Human attacks on multi-agent systems	Adversaries exploit inter-agent delegation, trust relationships, and workflow dependencies to escalate privileges or manipulate AI-driven operations.	Restrict agent delegation mechanisms, enforce inter-agent authentication, and deploy behavioural monitoring to detect manipulation attempts. Enforce multi-agent task segmentation to prevent attackers from escalating privileges across interconnected agents.
T15	Human manipulation	In scenarios where AI agents engage in direct interaction with human users, the trust relationship reduces user scepticism, increasing reliance on the agent's responses and autonomy. This implicit trust and direct human/agent interaction create risks, as attackers can coerce agents to manipulate users, spread misinformation, and take covert actions.	Monitor agent behaviour to ensure it aligns with its defined role and expected actions. Restrict tool access to minimize the attack surface, limit the agent's ability to print links, implement validation mechanisms to detect and filter manipulated responses using guardrails, moderation APIs, or another model.

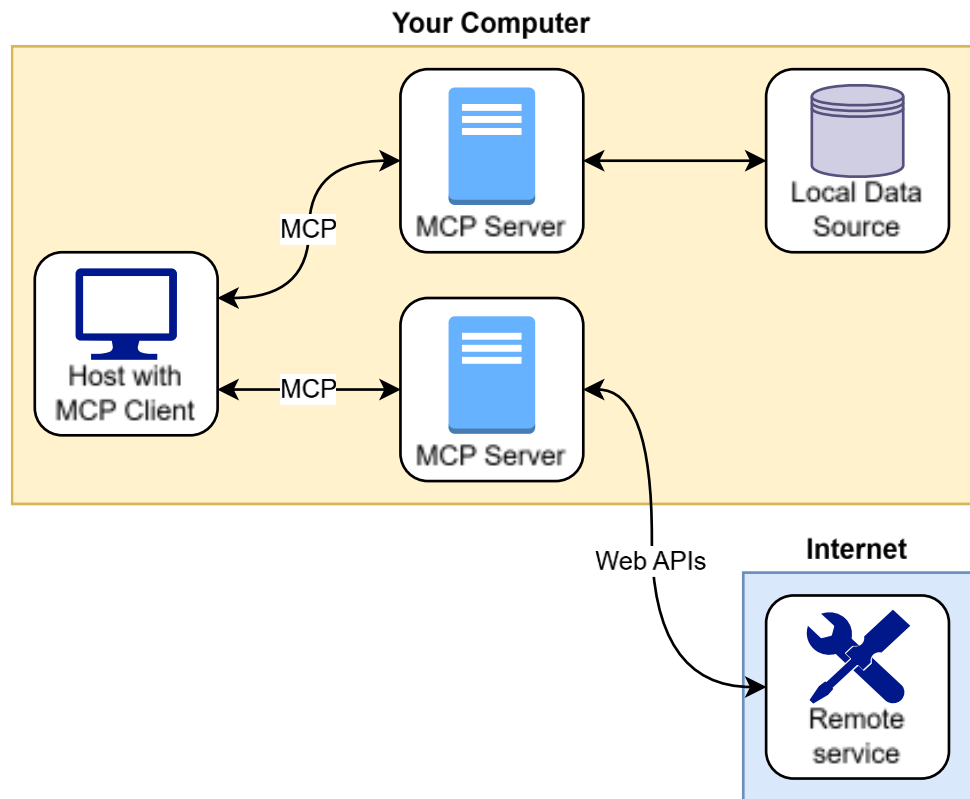
ANNEX B

Model Context Protocol

Model Context Protocol (MCP) is an open protocol that standardises how applications provide context to LLMs. An analogy would be like a USB-C port on a computer. Just as how USB-C provides a standard way to connect devices, MCP provides a standard way to connect AI models to various tools and resources.¹⁷

MCP follows a client-server architecture where a host application can connect to multiple servers:

Figure 20: General MCP Architecture



¹⁷ Anthropic. [Model Context Protocol, Introduction.](#)

Components in MCP architecture:

- MCP Hosts: Programs like Claude Desktop, IDEs or AI tools that want to access data through MCP
- MCP Clients: Protocol clients that maintain 1:1 connections with servers
- MCP Servers: Lightweight programs that each expose specific capabilities through the standardized Model Context Protocol
- Local Data Sources: Computer's files, databases, and services that MCP servers can securely access
- Remote Services: External systems available over the internet (e.g., through APIs) that MCP servers can connect to

The main difference from other tool invocation setups, such as OpenAPI is that MCP is dynamic, allowing runtime discovery of available tools from a given server.

Risks and Threats

Calling for tools has inherent dangers, no matter the implementation (OpenAPI, AI Actions, or MCP). All are susceptible to prompt injection and confused deputy threats¹⁸.

Other possible threats include Server Name Collision, Installer Spoofing, Backdoors, Tool Name Conflicts, Sandbox Escapes, and Configuration Drift¹⁹.

¹⁸ Rehberger, J. [MCP: Untrusted Servers and Confused Clients, Plus a Sneaky Exploit](#).

¹⁹ Hou, X., Zhao, Y., Wang, S., & Wang, H. [Model Context Protocol \(MCP\): Landscape, Security Threats, and Future Research Directions](#).

Mitigation Recommendations

While Anthropic's MCP specification ²⁰ does not cover all threats, it provides recommendations on the secure usage and configuration of MCP ²¹:

1. Do not randomly download or connect AI to untrusted MCP or OpenAPI tool servers.
2. Inspect code, interface definition, check for backdoors, hidden instructions.
3. Connect to MCP servers and tools which have undergone appropriate security due diligence (e.g. SBOM review, security assessment), instead of relying solely on brand reputation or popularity.
4. Follow basic security practices such as peer code reviews, static analysis and threat modelling.
5. Human oversight - keeping humans in the loop and in control is essential as there is no deterministic solution for prompt injections.
6. Logging and monitoring - track human identities to AI actions.
7. Manage prompt injection threats based on scenario and context.

²⁰ Anthropic. [Model Context Protocol, Core architecture](#).

²¹ Rehberger, J. [MCP: Untrusted Servers and Confused Clients, Plus a Sneaky Exploit](#).

ANNEX C

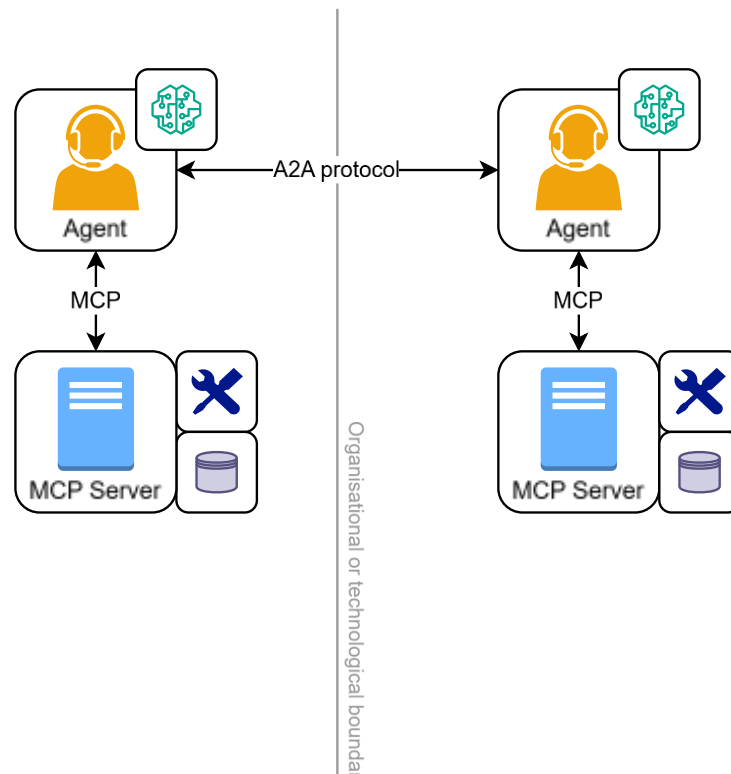
Agent 2 Agent Protocol

The Agent2Agent (A2A) Protocol is an open standard designed to enable seamless communication and collaboration between AI agents²². It facilitates dynamic, multimodal communication between different agents as peers, allowing agents to collaborate, delegate, and manage shared tasks.

MCP and A2A

MCP connects agents to tools and resources, whereas A2A enables agent-to-agent collaboration²³. Figure 21 shows how MCP and A2A may be used together in a multi-agent system.

Figure 21: A2A and MCP as Complementary Protocols



²² Google LLC. [What is A2A?](#)

²³ Google LLC. [A2A and MCP: Complementary Protocols for Agentic Systems.](#)

Advantages of A2A

Traditional enterprise systems rely on APIs, requiring knowledge of specific endpoints and tightly coupled logic. This leads to systems becoming rigid and unscalable as agent complexity increases. A2A shifts communications from calling functions, to expressing goals with constraints²⁴. This reduces integration complexity, fosters innovation, and future-proofs systems.

In A2A, agents operate without having to share internal memory, tools, or proprietary logic. Agents interact based on declared capabilities and exchanged context, preserving intellectual property and enhancing security²⁵.

Other agentic protocols

A2A is not the only available protocol between agents. The Agent Communication Protocol (ACP) from IBM, and the open-source Agent Network Protocol (ANP) are also available to enhance interactions between agents.

- Agent Communication Protocol (ACP): Addresses the need for secure, compliant communication between agents, especially in highly regulated industries
- Agent Network Protocol (ANP): An open standard aiming to create decentralised, interoperable networks of AI agents. Promotes peer-to-peer collaboration.

Each of these protocols will have differing domain application, interoperability and openness. No matter which protocol is used, users should ensure that the proper mitigations are in place to secure communications between agents.

²⁴ Auxiliobits. [Agent-to-Agent Protocols: How Google's A2A is Shaping Future Automations?](#)

²⁵ Google LLC. [What is A2A?](#)

Threats and Mitigations

A2A, ACP and ANP have made inter-agent communication much more convenient, however, with this capability comes more threats and potential attack surfaces.

The following table lists some possible threats to a system using the agentic protocols, as well as possible mitigations²⁶.

Table 9: Threats and Mitigation to agentic protocols

Threats	Mitigations
Message generation attacks	Input and Output validation
Model extraction	Enforce rate limits on protocol interactions for each session / user / agent. Observe query patterns for anomalies that suggest probing or data extraction attempts.
Data poisoning through message parts	Strong validation of message parts. Limit agent access with principle of least privilege. Track origin and lineage of data.
Sensitive information disclosure	Automated PII redaction. Fine-grained access control. Context-aware guardrails.
Unauthorised agent impersonation	Require agents to use Decentralised identifiers (DID). Secure authentication. Implement a trusted agent registry.
Message injection attacks	Implement digital signatures for protocol messages. Input validation. Content filtering.
Protocol downgrade attacks	Have secure protocol negotiation, such as TLS with secure authentication. Enforce deprecation policy for older protocol versions.

²⁶ Huang, K. [Threat Modeling Google's A2A Protocol with the MAESTRO Framework](#).

Threats	Mitigations
<p>Malicious A2A server impersonating a trusted company</p>	<p>Decentralised identifiers (DID) for server identities. Certificate transparency for agent cards. Mutual TLS (mTLS) authentication. DNSSEC for server domain. Agent registry verification. Agent card signature verification. MFA for critical operations. Behavioural analysis and reputation systems. Auditing and logging. Deploy honeypot protocol servers.</p>
<p>Denial of service attacks</p>	<p>Robust infrastructure. DDoS protection. Rate limiting.</p>
<p>Manipulation of logging data</p>	<p>Secure logging infrastructure. Log integrity monitoring. Anomaly detection.</p>
<p>Unauthorised access to agent credentials</p>	<p>Secure key storage. Key rotation.</p>
<p>Lack of compliance on sensitive data</p>	<p>Data minimisation. Pseudonymisation/Anonymisation</p>
<p>Malicious agent interaction</p>	<p>Secure inter-agent communication. Agent reputation systems. Sandbox agents.</p>
<p>Flaws in Multi-Agent Collaboration Mechanisms (In multi-agent systems, deficiencies in internal collaboration mechanisms can manifest as follows: when agents make distributed decisions based on localized information, conflicts between their objectives may result in systemic failures.)</p>	<p>Establish a coordination and management mechanism for multi-agents.</p>

REFERENCES

- AG2. (n.d.). *UserProxyAgent*. Retrieved from AG2: <https://docs.ag2.ai/0.8.7/docs/api-reference/autogen/UserProxyAgent/>
- AI Verify Foundation. (n.d.). *List of Datasets*. Retrieved from Moonshot: <https://aiverify-foundation.github.io/moonshot/resources/datasets/>
- Anthropic. (18 Jun, 2025). *Model Context Protocol, Core architecture*. Retrieved from Model Context Protocol: <https://modelcontextprotocol.io/docs/concepts/architecture>
- Anthropic. (18 Jun, 2025). *Model Context Protocol, Introduction*. Retrieved from Model Context Protocol: <https://modelcontextprotocol.io/introduction>
- Anthropic. (n.d.). *Content moderation*. Retrieved from <https://docs.anthropic.com/en/docs/about-claude/use-case-guides/content-moderation>
- Apostrophe Technologies. (May, 2025). *sanitize-html*. Retrieved from npm: <https://www.npmjs.com/package/sanitize-html>
- Arias, D., & Bellen, S. (7 Oct, 2021). *What Are Refresh Tokens and How to Use Them Securely, auth0*. Retrieved from auth0: <https://auth0.com/blog/refresh-tokens-what-are-they-and-when-to-use-them/>
- Auxiliobits. (2025). *Agent-to-Agent Protocols: How Google's A2A is Shaping Future Automations?* Retrieved from Auxiliobits: https://www.auxiliobits.com/blog/agent-to-agent-protocols-how-googles-a2a-is-shaping-future-automations/#elementor-toc_heading-anchor-1
- AWS. (Aug, 2025). *AWS Prescriptive Guidance: Operationalizing agentic AI on AWS*. Retrieved from AWS: <https://docs.aws.amazon.com/prescriptive-guidance/latest/strategy-operationalizing-agentic-ai/introduction.html>
- AWS. (n.d.). *Control subnet traffic with network access control lists*. Retrieved from AWS: <https://docs.aws.amazon.com/vpc/latest/userguide/vpc-network-acls.html>
- AWS. (n.d.). *Security best practices in IAM*. Retrieved from AWS: <https://docs.aws.amazon.com/IAM/latest/UserGuide/best-practices.html>
- AWS. (n.d.). *Use temporary credentials with AWS resources*. Retrieved from AWS: https://docs.aws.amazon.com/IAM/latest/UserGuide/id_credentials_temp_use-resources.html
- AWS. (n.d.). *What is AWS Secrets Manager?* Retrieved from AWS: <https://docs.aws.amazon.com/secretsmanager/latest/userguide/intro.html>
- Besen, S., & Gutowska, A. (n.d.). *What is Agent Communication Protocol (ACP)?* Retrieved from IBM: <https://www.ibm.com/think/topics/agent-communication-protocol>
- Cameron, A. (8 Apr, 2025). *pip-audit*. Retrieved from PyPI: <https://pypi.org/project/pip-audit/>
- Center for Research on Foundation Models (CRFM), Stanford University. (2025). *Holistic Evaluation of Language Models (HELM)*. Retrieved from <https://crfm.stanford.edu/helm/>

- Chhikara, P., Khant, D., Aryan, S., Singh, T., & Yadav, D. (28 Apr, 2025). *Mem0: Building Production-Ready AI Agents with Scalable Long-Term Memory*. Retrieved from arxiv: <https://arxiv.org/abs/2504.19413v1>
- CISA. (2025). *Software Bill of Materials (SBOM)*. Retrieved from CISA: <https://www.cisa.gov/sbom>
- Cloud Security Alliance. (21 Aug, 2024). *Best practices for event logging and threat detection*. Retrieved from Cloud Security Alliance: <https://cloudsecurityalliance.org/resources/best-practices-for-event-logging-and-threat-detection>
- Cloud Security Alliance. (16 Jul, 2025). *Agentic AI Red Teaming Guide*. Retrieved from Cloud Security Alliance: <https://cloudsecurityalliance.org/artifacts/agentic-ai-red-teaming-guide>
- Cloudflare. (n.d.). *What is mutual TLS (mTLS)?* Retrieved from Cloudflare: <https://www.cloudflare.com/learning/access-management/what-is-mutual-tls/>
- CodeSignal. (2025). *Developing a Robust System Prompt*. Retrieved from CodeSignal: <https://codesignal.com/learn/courses/building-a-chatbot-service-with-fastapi/lessons/crafting-a-robust-system-prompt-for-chatbot-interaction>
- Conversation-AI. (n.d.). *Enabling online conversations*. Retrieved from Perspective: <https://www.perspectiveapi.com/>
- crewAI Inc. (n.d.). *Human-in-the-Loop (HITL) Workflows*. Retrieved from crewAI: <https://docs.crewai.com/en/learn/human-in-the-loop>
- Cure53. (n.d.). *DOMPurify*. Retrieved from Github: <https://github.com/cure53/DOMPurify>
- Cyber Security Agency of Singapore. (27 Jul, 2022). *Critical Information Infrastructure Supply Chain Programme Paper*. Retrieved from CSA: <https://www.csa.gov.sg/resources/publications/critical-information-infrastructure-supply-chain-programme-paper>
- Cyber Security Agency of Singapore. (15 Oct, 2024). *Guidelines and Companion Guide on Securing AI Systems*. Retrieved from CSA: <https://www.csa.gov.sg/resources/publications/guidelines-and-companion-guide-on-securing-ai-systems>
- Cyber Security Agency of Singapore. (5 Jun, 2025). *Responsible Vulnerability Disclosure Policy*. Retrieved from <https://isomer-user-content.by.gov.sg/36/4aa60609-4481-4e7c-92eb-2728247a084f/responsible-vulnerability-disclosure-policy.pdf>
- Cyber Security Agency of Singapore. (20 Jan, 2025). *Supplementary references*. Retrieved from CSA: <https://www.csa.gov.sg/legislation/supplementary-references>
- Debenedetti, E., Shumailov, I., Fan, T., Hayes, J., Carlini, N., Fabian, D., . . . Tramèr, F. (24 Jun, 2025). *Defeating Prompt Injections by Design*. Retrieved from arxiv: <https://arxiv.org/abs/2503.18813>
- Díaz, S., Kern, C., & Olive, K. (May, 2025). *Google's Approach for Secure AI Agents*. Retrieved from Google Research: <https://research.google/pubs/an-introduction-to-googles-approach-for-secure-ai-agents/>
- E2B. (26 Aug, 2025). *E2B*. Retrieved from GitHub: <https://github.com/e2b-dev/E2B>

- EU AI Act Holistic AI Team. (1 Aug, 2024). *High-Risk AI Systems Under the EU AI Act*. Retrieved from EU AI Act: <https://www.euaiact.com/blog/high-risk-ai-systems-under-the-eu-ai-act>
- Explosion. (n.d.). *Industrial-Strength Natural Language Processing*. Retrieved from spaCy: <https://spacy.io/>
- Feldman, E. (15 Apr, 2025). *Implementing effective guardrails for AI agents*. Retrieved from The Source Gitlab: <https://about.gitlab.com/the-source/ai/implementing-effective-guardrails-for-ai-agents/>
- Flinders, M., Smalley, I., & Schneider, J. (30 Apr, 2025). *AI fraud detection in banking*. Retrieved from IBM: <https://www.ibm.com/think/topics/ai-fraud-detection-in-banking>
- Fortinet. (2025). *What Is A Message Authentication Code?* Retrieved from Fortinet: <https://www.fortinet.com/resources/cyberglossary/message-authentication-code>
- French National Cybersecurity Authority (ANSSI). (7 Feb, 2025). *Building trust in AI through a cyber risk-based approach*. Retrieved from MesServicesCyber Innovation ANSSI: <https://messervices.cyber.gouv.fr/guides/en-building-trust-ai-through-cyber-risk-based-approach>
- Gabarda, F. C. (1 Jul, 2025). *Model Context Protocol (MCP): Understanding security risks and controls*. Retrieved from Red Hat Blog: <https://www.redhat.com/en/blog/model-context-protocol-mcp-understanding-security-risks-and-controls>
- GitHub. (28 Nov, 2022). *Rate limits for the REST API*. Retrieved from GitHub Docs: <https://docs.github.com/en/rest/using-the-rest-api/rate-limits-for-the-rest-api?apiVersion=2022-11-28>
- GitHub. (n.d.). *Dependabot quickstart guide*. Retrieved from GitHub Docs: <https://docs.github.com/en/code-security/getting-started/dependabot-quickstart-guide>
- GitLab. (n.d.). *Dependency Scanning*. Retrieved from GitLab Docs: https://docs.gitlab.com/user/application_security/dependency_scanning/
- GitLab. (n.d.). *What is version control?* Retrieved from GitLab: <https://about.gitlab.com/topics/version-control/>
- Gittlen, S. (2 May, 2024). *The ultimate guide to SBOMs*. Retrieved from GitLab: <https://about.gitlab.com/blog/the-ultimate-guide-to-sboms/>
- Google. (n.d.). *About IAM authentication*. Retrieved from Google Cloud: <https://cloud.google.com/memorystore/docs/valkey/about-iam-auth>
- Google. (n.d.). *Custom Search JSON API*. Retrieved from <https://developers.google.com/custom-search/v1/overview>
- Google LLC. (9 Apr, 2025). *A2A and MCP: Complementary Protocols for Agentic Systems*. Retrieved from Agent2Agent (A2A) Protocol: <https://a2aproject.github.io/A2A/latest/topics/a2a-and-mcp/#how-a2a-and-mcp-complement-each-other>
- Google LLC. (9 Apr, 2025). *What is A2A?* Retrieved from Agent2Agent (A2A) Protocol: <https://a2aproject.github.io/A2A/latest/topics/what-is-a2a/>

- Google. (n.d.). *Secret Manager overview*. Retrieved from Google Cloud: <https://cloud.google.com/secret-manager/docs/overview>
- GovTech Singapore (AI Practice). (Jul, 2025). *Agentic Risk & Capability Framework*. Retrieved from <https://govtech-responsibleai.github.io/agentic-risk-capability-framework/>
- Guardrails AI. (n.d.). *Guardrails AI*. Retrieved from Github: <https://github.com/guardrails-ai/guardrails>
- Harang, R., & Sablotny, M. (25 Feb, 2025). *Agentic Autonomy Levels and Security*. Retrieved from NVIDIA DEVELOPER: <https://developer.nvidia.com/blog/agentic-autonomy-levels-and-security/>
- HashiCorp. (n.d.). *Vault*. Retrieved from GitHub: <https://github.com/hashicorp/vault>
- Helicone Inc. (n.d.). *Helicone*. Retrieved from GitHub: <https://github.com/Helicone/helicone>
- Hou, X., Zhao, Y., Wang, S., & Wang, H. (Apr, 2025). *Model Context Protocol (MCP): Landscape, Security Threats, and Future Research Directions*. Retrieved from <https://arxiv.org/pdf/2503.23278>
- Huang, K. (22 Dec, 2024). *7 Layered Agentic AI Reference Architecture*. Retrieved from Medium: <https://kenhuangus.medium.com/7-layered-agentic-ai-reference-architecture-20276f83b7ee>
- Huang, K. (02 Jun, 2025). *Agentic AI Threat Modeling Framework: MAESTRO*. Retrieved from Cloud Security Alliance: <https://cloudsecurityalliance.org/blog/2025/02/06/agentic-ai-threat-modeling-framework-maestro>
- Huang, K. (30 Apr, 2025). *Threat Modeling Google's A2A Protocol with the MAESTRO Framework*. Retrieved from Cloud Security Alliance: <https://cloudsecurityalliance.org/blog/2025/04/30/threat-modeling-google-s-a2a-protocol-with-the-maestro-framework>
- Huang, K., Narajala, V. S., Yeoh, J., Ross, J., Raskar, R., Harkati, Y., . . . Hughes, C. (28 May, 2025). *A Novel Zero-Trust Identity Framework for Agentic AI: Decentralized Authentication and Fine-Grained Access Control*. Retrieved from arxiv: <https://arxiv.org/abs/2505.19301>
- Hugging Face. (n.d.). *Daily Papers: Instruction Following Score*. Retrieved from [https://huggingface.co/papers?q=Instruction%20Following%20Score%20\(IFS\)](https://huggingface.co/papers?q=Instruction%20Following%20Score%20(IFS))
- Hugging Face. (n.d.). *Model Cards*. Retrieved from Hugging Face: <https://huggingface.co/docs/hub/en/model-cards>
- Hugging Face. (n.d.). *Pickle Scanning*. Retrieved from Hugging Face: <https://huggingface.co/docs/hub/security-pickle>
- IETF OAuth Working Group. (n.d.). *OAuth Scopes*. Retrieved from OAuth 2.0: <https://oauth.net/2/scope/>
- Invariant. (1 Apr, 2025). *MCP Security Notification: Tool Poisoning Attacks*. Retrieved from Invariantlabs: <https://invariantlabs.ai/blog/mcp-security-notification-tool-poisoning-attacks>

- ISO/IEC Joint Technical Committee. (Dec, 2023). *ISO/IEC 42001*. Retrieved from International Organization for Standardization: <https://www.iso.org/standard/42001>
- Jambrecic, R. (7 Jan, 2025). *Tools Dependency Injection*. Retrieved from AG2: <https://docs.ag2.ai/latest/docs/blog/2025/01/07/Tools-Dependency-Injection/>
- Jarvis, C. (19 Dec, 2023). *How to implement LLM guardrails*. Retrieved from OpenAI Cookbook: https://cookbook.openai.com/examples/how_to_use_guardrails
- Jin, X., Guo, Z., Zhang, P., Lu, S., Dai, W., Nujibieke, . . . Li, G. (2 Feb, 2025). *Bridging Minds and Machines: Agents with Human-in-the-Loop – Frontier Research, Real-World Impact, and Tomorrow’s Possibilities*. Retrieved from Camel-AI: <https://www.camel-ai.org/blogs/human-in-the-loop-ai-camel-integration>
- Joint Task Force. (Sep, 2020). *Special Publication 800-53 Rev. 5 Security and Privacy Controls for Information Systems and Organizations*. Retrieved from NIST: <https://doi.org/10.6028/NIST.SP.800-53r5>
- Kartha, V. (3 May, 2024). *Self-Reflecting AI Agents Using LangChain*. Retrieved from Medium: <https://vijaykumarkartha.medium.com/self-reflecting-ai-agents-using-langchain-d3a93684da92>
- Kumar, A., Roh, J., Naseh, A., Karpinska, M., Iyer, M., Houmansadr, A., & Bagdasarian, E. (5 Feb, 2025). *OverThink: Slowdown Attacks on Reasoning LLMs*. Retrieved from arxiv: <https://arxiv.org/abs/2502.02542>
- LangChain. (2025). *How to pass run time values to tools*. Retrieved from LangChain: https://python.langchain.com/docs/how_to/tool_runtime/
- LangChain. (2025). *LangMem*. Retrieved from LangGraph: <https://langchain-ai.github.io/langmem/>
- LangChain. (n.d.). *E2B Data Analysis*. Retrieved from LangChain: https://python.langchain.com/docs/integrations/tools/e2b_data_analysis/
- LangChain. (n.d.). *LangGraph interrupt: Making it easier to build human-in-the-loop agents with interrupt*. Retrieved from LangChain: <https://blog.langchain.com/making-it-easier-to-build-human-in-the-loop-agents-with-interrupt/>
- Langfuse. (n.d.). *Open Source LLM Engineering Platform*. Retrieved from Langfuse: <https://langfuse.com/>
- LangSmith. (n.d.). *Ship agents with confidence*. Retrieved from LangSmith: <https://www.langchain.com/langsmith>
- Liu, X., Yu, H., Zhang, H., Xu, Y., Lei, X., Lai, H., . . . Tang, J. (25 Oct, 2023). *AgentBench: Evaluating LLMs as Agents*. Retrieved from arxiv: <https://arxiv.org/abs/2308.03688>
- Lucas, J. (16 Dec, 2024). *Sandboxing Agentic AI Workflows with WebAssembly*. Retrieved from NVIDIA Developer: <https://developer.nvidia.com/blog/sandboxing-agentic-ai-workflows-with-webassembly/>
- Manral, V. (28 Jul, 2023). *Generative AI: Proposed Shared Responsibility Model*. Retrieved from Cloud Security Alliance: <https://cloudsecurityalliance.org/blog/2023/07/28/generative-ai-proposed-shared-responsibility-model#>

- Meadows, J., & Chang, A. (27 Mar, 2024). *How to choose a known, trusted supplier for open source software*. Retrieved from Google Cloud: <https://cloud.google.com/blog/products/identity-security/how-to-choose-a-known-trusted-supplier-for-open-source-software>
- Mell, P., & Grace, T. (n.d.). *Special Publication 800-145 The NIST Definition of Cloud Computing*. Retrieved from NIST: <https://doi.org/10.6028/NIST.SP.800-145>
- Meta Llama. (n.d.). *Purple Llama*. Retrieved from Github: <https://github.com/meta-llama/PurpleLlama>
- Microsoft. (29 Sep, 2024). *Artificial intelligence (AI) shared responsibility model*. Retrieved from Microsoft Learn: <https://learn.microsoft.com/en-us/azure/security/fundamentals/shared-responsibility-ai>
- Microsoft AI Red Team. (2024). *PyRIT*. Retrieved from <https://azure.github.io/PyRIT/>
- Microsoft Azure. (5 Jul, 2024). *What is an IP based access control list (ACL)?* Retrieved from Microsoft Learn: <https://learn.microsoft.com/en-us/azure/virtual-network/ip-based-access-control-list-overview>
- Microsoft. (n.d.). *Presidio: Data Protection and De-identification SDK*. Retrieved from Microsoft Presidio: <https://microsoft.github.io/presidio/>
- Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., . . . Gebru, T. (14 Jan, 2019). *Model Cards for Model Reporting*. Retrieved from arxiv: <https://arxiv.org/abs/1810.03993>
- MITRE. (n.d.). *ATLAS Matrix*. Retrieved from MITRE ATLAS: <https://atlas.mitre.org/matrices/ATLAS>
- MITRE. (n.d.). *Supply Chain Security Framework*. Retrieved from MITRE System of Trust: https://sot.mitre.org/framework/system_of_trust.html
- Mu, N., Lu, J., Lavery, M., & Wagner, D. (15 Feb, 2025). *A Closer Look at System Prompt Robustness*. Retrieved from <https://arxiv.org/abs/2502.12197>
- Murúa, T. (1 May, 2025). *RAG and the value of grounding*. Retrieved from elastic search labs: <https://www.elastic.co/search-labs/blog/grounding-rag>
- National Security Agency. (5 Mar, 2024). *Advancing Zero Trust Maturity Throughout the Network and Environment Pillar*. Retrieved from NSA: <https://media.defense.gov/2024/Mar/05/2003405462/-1/-1/0/CSI-ZERO-TRUST-NETWORK-ENVIRONMENT-PILLAR.PDF>
- Nelson, A., Rekhi, S., Souppaya, M., & Scarfone, K. (n.d.). *Special Publication 800-61 Rev. 3 Incident Response Recommendations and Considerations for Cybersecurity Risk Management*. Retrieved from NIST: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-61r3.pdf>
- NIST. (26 Jan, 2023). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. Retrieved from NIST: <https://doi.org/10.6028/NIST.AI.100-1>
- NIST. (14 Jan, 2025). *Cryptographic Standards and Guidelines*. Retrieved from NIST: <https://csrc.nist.gov/projects/cryptographic-standards-and-guidelines>

- NVIDIA. (2025). *About NeMo Guardrails*. Retrieved from NVIDIA:
<https://docs.nvidia.com/nemo/guardrails/latest/index.html#>
- NVIDIA. (n.d.). *garak, LLM vulnerability scanner*. Retrieved from Github:
<https://github.com/NVIDIA/garak>
- OpenAI. (2025). *Guardrails*. Retrieved from OpenAI Agents SDK: <https://openai.github.io/openai-agents-python/guardrails/>
- OpenAI. (n.d.). *Moderation*. Retrieved from <https://platform.openai.com/docs/guides/moderation>
- OpenClaw. (n.d.). *OpenClaw Threat Model*. Retrieved from OpenClaw:
<https://trust.openclaw.ai/trust/threatmodel>
- OpenJS Foundation. (n.d.). *ESLint*. Retrieved from GitHub: <https://github.com/eslint/eslint>
- OWASP. (22 Apr, 2025). *OWASP Gen AI Security Project - Multi-Agentic system Threat Modelling Guide*. Retrieved from <https://genai.owasp.org/resource/multi-agentic-system-threat-modeling-guide-v1-0/>
- OWASP. (28 Jun, 2025). *OWASP Gen AI Security Project - Securing Agentic Applications Guide*. Retrieved from <https://genai.owasp.org/resource/securing-agentic-applications-guide-1-0/>
- OWASP. (9 Dec, 2025). *OWASP Top 10 For Agentic Applications 2026*. Retrieved from <https://genai.owasp.org/resource/owasp-top-10-for-agentic-applications-for-2026/>
- OWASP. (17 Feb, 2025). *OWASP Top 10 for LLMs - Agentic AI - Threats and Mitigations*. Retrieved from <https://genai.owasp.org/resource/agentic-ai-threats-and-mitigations/>
- OWASP. (23 Jan, 2025). *OWASP Top 10 for LLMs - GenAI Red Teaming Guide*. Retrieved from <https://genai.owasp.org/resource/genai-red-teaming-guide/>
- OWASP. (n.d.). *Authentication Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/Authentication_Cheat_Sheet.html
- OWASP. (n.d.). *Authorization Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/Authorization_Cheat_Sheet.html
- OWASP. (n.d.). *Content Security Policy Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/Content_Security_Policy_Cheat_Sheet.html
- OWASP. (n.d.). *Docker Security Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/Docker_Security_Cheat_Sheet.html
- OWASP. (n.d.). *File Upload Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/File_Upload_Cheat_Sheet.html
- OWASP. (n.d.). *Input Validation Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/Input_Validation_Cheat_Sheet.html
- OWASP. (n.d.). *LLM Prompt Injection Prevention Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series:
https://cheatsheetseries.owasp.org/cheatsheets/LLM_Prompt_Injection_Prevention_Cheat_Sheet.html

- OWASP. (n.d.). *Secrets Management Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series: https://cheatsheetseries.owasp.org/cheatsheets/Secrets_Management_Cheat_Sheet.html
- OWASP. (n.d.). *XSS Filter Evasion Cheat Sheet*. Retrieved from OWASP Cheat Sheet Series: https://cheatsheetseries.owasp.org/cheatsheets/XSS_Filter_Evasion_Cheat_Sheet.html
- PagerDuty. (n.d.). *Incident Response*. Retrieved from PagerDuty: <https://response.pagerduty.com/>
- Perrone, P. (15 Apr, 2025). *MCP Is a Security Nightmare — Here's How the Agent Security Framework Fixes It*. Retrieved from Medium: <https://medium.com/data-science-collective/mcp-is-a-security-nightmare-heres-how-the-agent-security-framework-fixes-it-fd419fdfaf4e>
- Perrot, C., Tanke, M. L., Roy, M., & Sachs, R. (9 Apr, 2025). *Implement human-in-the-loop confirmation with Amazon Bedrock Agents*. Retrieved from AWS: <https://aws.amazon.com/blogs/machine-learning/implement-human-in-the-loop-confirmation-with-amazon-bedrock-agents/>
- Personal Data Protection Commission Singapore. (1 Mar, 2024). *Advisory Guidelines on use of Personal Data in AI Recommendation and Decision Systems*. Retrieved from PDPC: <https://www.pdpc.gov.sg/guidelines-and-consultation/2024/02/advisory-guidelines-on-use-of-personal-data-in-ai-recommendation-and-decision-systems>
- Personal Data Protection Commission Singapore. (24 Jul, 2024). *Guide to Basic Anonymisation*. Retrieved from PDPC: [https://www.pdpc.gov.sg/-/media/files/pdpc/pdf-files/advisory-guidelines/guide-to-basic-anonymisation-\(updated-24-july-2024\).pdf](https://www.pdpc.gov.sg/-/media/files/pdpc/pdf-files/advisory-guidelines/guide-to-basic-anonymisation-(updated-24-july-2024).pdf)
- Personal Data Protection Commission Singapore. (14 Dec, 2024). *Guide to Data Protection Practices for ICT Systems*. Retrieved from PDPC: <https://www.pdpc.gov.sg/-/media/files/pdpc/pdf-files/other-guides/tech-omnibus/guide-to-data-protection-practices-for-ict-systems.pdf>
- Promptfoo. (n.d.). *Promptfoo: LLM evals & red teaming*. Retrieved from Github: <https://github.com/promptfoo/promptfoo>
- Python Code Quality Authority. (n.d.). *Bandit*. Retrieved from <https://bandit.readthedocs.io/en/latest/>
- Rasmussen, P., Paliychuk, P., Beauvais, T., Ryan, J., & Chalef, D. (20 Jan, 2025). *Zep: A Temporal Knowledge Graph Architecture for Agent Memory*. Retrieved from arxiv: <https://arxiv.org/abs/2501.13956>
- Rehberger, J. (2 May, 2025). *MCP: Untrusted Servers and Confused Clients, Plus a Sneaky Exploit*. Retrieved from Embrace The Red: <https://embracethered.com/blog/posts/2025/model-context-protocol-security-risks-and-exploits/>
- Rehberger, J. (n.d.). *Trust No AI: Prompt Injection Along The CIA Security Triad*. Retrieved from <https://arxiv.org/pdf/2412.06090>
- Rose, S., Borchert, O., Mitchell, S., & Connelly, S. (Aug, 2020). *Special Publication 800-207 Zero Trust Architecture*. Retrieved from NIST: <https://nvlpubs.nist.gov/nistpubs/specialpublications/NIST.SP.800-207.pdf>
- Rose, S., Liu, C., & Gibson, R. (Apr, 2025). *Special Publication 800-81 Rev. 3 Secure Domain Name System (DNS) Deployment Guide*. Retrieved from NIST: <https://doi.org/10.6028/NIST.SP.800-81r3.ipd>

- Sapkota, R., Roumeliotis, K. I., & Karkee, M. (28 May, 2025). *AI Agents vs. Agentic AI: A Conceptual Taxonomy, Applications and Challenges*. Retrieved from arxiv: <https://arxiv.org/abs/2505.10468>
- Scarfone, K., Souppaya, M., Cody, A., & Orebaugh, A. (n.d.). *Special Publication 800-115 Technical Guide to Information Security Testing and Assessment*. Retrieved from NIST: <https://csrc.nist.gov/pubs/sp/800/115/final>
- Semgrep, Inc. (n.d.). *Semgrep*. Retrieved from GitHub: <https://github.com/semgrep/semgrep>
- Shah, D. (4 Jun, 2025). *Introduction to Training Data Poisoning: A Beginner's Guide*. Retrieved from Lakera: <https://www.lakera.ai/blog/training-data-poisoning>
- Snyk. (Jun, 2025). *Snyk Open Source*. Retrieved from Snyk User Docs: <https://docs.snyk.io/scan-with-snyk/snyk-open-source>
- Souppaya, M., Scarfone, K., & Dodson, D. (Feb, 2022). *Special Publication 800-218 Secure Software Development Framework (SSDF) Version 1.1*. Retrieved from NIST: <https://nvlpubs.nist.gov/nistpubs/specialpublications/nist.sp.800-218.pdf>
- Stryker, C. (2025). *What is agentic AI?* Retrieved from IBM: <https://www.ibm.com/think/topics/agentic-ai>
- Surapaneni, R., Jha, M., Vakoc, M., & Segal, T. (9 Apr, 2025). *Announcing the Agent2Agent Protocol (A2A)*. Retrieved from Google for Developers: <https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interopability/>
- The Linux Foundation. (n.d.). *Supply-chain Levels for Software Artifacts*. Retrieved from SLSA: <https://slsa.dev/>
- The OWASP Foundation. (July, 2017). *OWASP Code Review Guide*. Retrieved from OWASP: <https://owasp.org/www-project-code-review-guide/>
- The Vulnerable MCP Project. (2025). *The Vulnerable MCP Project*. Retrieved from <https://vulnerablemcp.info/>
- tl;dr sec. (Feb, 2025). *prompt-injection-defenses*. Retrieved from GitHub: <https://github.com/tldrsec/prompt-injection-defenses>
- traceloop. (n.d.). *OpenLLMetry*. Retrieved from GitHub: <https://github.com/traceloop/openllmetry>
- UK AI Security Institute, Arcadia Impact, Vector Institute. (n.d.). *Inspect Evals*. Retrieved from https://ukgovernmentbeis.github.io/inspect_evals/
- UK National Cyber Security Centre. (12 Oct, 2023). *Supply Chain Guidance*. Retrieved from NCSC: <https://www.ncsc.gov.uk/collection/supply-chain/guidance>
- UK National Cyber Security Centre. (7 Nov, 2024). *Vulnerability Disclosure Toolkit*. Retrieved from NCSC: <https://www.ncsc.gov.uk/information/vulnerability-disclosure-toolkit>
- Ul Muram, F., Tran, H., & Zdun, U. (1 Apr, 2017). *Systematic Review of Software Behavioral Model Consistency Checking*. Retrieved from https://www.researchgate.net/publication/316938485_Systematic_Review_of_Software_Behavioral_Model_Consistency_Checking

- University of the Sunshine Coast Australia. (n.d.). *What are credible sources?* Retrieved from <https://libguides.usc.edu.au/credible/web>
- Wang, C. L., Singhal, T., Kelkar, A., & Tuo, J. (8 Aug, 2025). *MI9 - Agent Intelligence Protocol: Runtime Governance for Agentic AI Systems*. Retrieved from arxiv: <https://arxiv.org/abs/2508.03858>
- Wickramasinghe, S. (18 Mar, 2025). *IT & System Availability + High Availability: The Ultimate Guide*. Retrieved from Splunk Blogs: https://www.splunk.com/en_us/blog/learn/availability.html
- Xu, F. F., Song, Y., Li, B., Tang, Y., Jain, K., Bao, M., . . . Neubig, G. (19 May, 2025). *TheAgentCompany: Benchmarking LLM Agents on Consequential Real World Tasks*. Retrieved from <https://the-agent-company.com/>
- Yuan, Y., Jiao, W., Wang, W., Huang, J.-t., Xu, J., Liang, T., . . . Tu, Z. (23 May, 2025). *Refuse Whenever You Feel Unsafe: Improving Safety in LLMs via Decoupled Refusal Training*. Retrieved from <https://arxiv.org/abs/2407.09121>
- Zaharia, M. A., Chen, A., Davidson, A., Ghodsi, A., Hong, S. A., Konwinski, A., . . . Zumar, C. (2018). *Accelerating the Machine Learning Lifecycle with MLflow*. Retrieved from GitHub: <https://github.com/mlflow/mlflow>
- Zhou, S., Xu, F. F., Zhu, H., Zhou, X., Lo, R., Sridhar, A., . . . Neubig, G. (16 Apr, 2024). *WebArena: A Realistic Web Environment for Building Autonomous Agents*. Retrieved from <https://webarena.dev/>