



# AI Incident Response Framework, V1.0

---

Workstream 2: Preparing Defenders for a Changing  
Cybersecurity Landscape



# Contents

---

|                                                                                                                                 |           |
|---------------------------------------------------------------------------------------------------------------------------------|-----------|
| <b>AI Incident Response Framework, V1.0</b>                                                                                     | <b>3</b>  |
| OASIS Open Project : Coalition for Secure AI (CoSAI) - Workstream 2: Preparing Defenders for a Changing Cybersecurity Landscape | 3         |
| <b>1. Overview</b>                                                                                                              | <b>4</b>  |
| 1.1. Abstract                                                                                                                   | 5         |
| 1.2. How To Use This Document                                                                                                   | 5         |
| 1.2.1. Quick Start Guide                                                                                                        | 5         |
| 1.2.2. Reader Guidance                                                                                                          | 6         |
| 1.3. Executive Summary                                                                                                          | 7         |
| 1.3.1 Types of AI Incidents                                                                                                     | 7         |
| 1.3.2 AI Incident Response Plan                                                                                                 | 7         |
| 1.3.3 Key Roles & RACI                                                                                                          | 8         |
| 1.3.4 Incident Response Playbooks                                                                                               | 8         |
| 1.3.5 Key Lessons Learned                                                                                                       | 8         |
| <b>2. Understanding AI Security Incidents</b>                                                                                   | <b>9</b>  |
| 2.1. Data Incidents                                                                                                             | 9         |
| 2.2. Model Incidents                                                                                                            | 9         |
| 2.3. Deployment Incidents                                                                                                       | 9         |
| 2.4. Infrastructure Incidents                                                                                                   | 10        |
| 2.5. User Interaction Incidents                                                                                                 | 10        |
| <b>3. Creating an AI Incident Response Plan</b>                                                                                 | <b>11</b> |
| 3.1. Pre-Incident Preparation                                                                                                   | 11        |
| 3.2. Monitoring and Telemetry                                                                                                   | 11        |
| 3.3. Incident Response Workflow                                                                                                 | 13        |
| 3.3.1. Preparation Phase                                                                                                        | 13        |
| 3.3.2. Detection and Analysis Phase                                                                                             | 14        |
| 3.3.3. Containment, Eradication, and Recovery Phase                                                                             | 15        |
| 3.3.4. Post-Incident Activity Phase                                                                                             | 18        |
| 3.4. Roles and Responsibilities                                                                                                 | 20        |
| 3.4.1. Provider vs. Consumer Responsibilities                                                                                   | 20        |
| 3.4.2. Team Assessment and Training                                                                                             | 20        |
| 3.5. Communication and Collaboration                                                                                            | 20        |
| 3.5.1. Regulatory Communication                                                                                                 | 20        |
| 3.5.2. Information Sharing with Community                                                                                       | 21        |
| <b>4. Frameworks and Playbooks</b>                                                                                              | <b>21</b> |
| 4.1. NIST SP 800-61r3 Alignment                                                                                                 | 21        |
| 4.1.1. Key Concept                                                                                                              | 21        |
| 4.1.2. Key Recommendations                                                                                                      | 22        |
| 4.2. OASIS CACAO Security Playbooks                                                                                             | 23        |
| 4.2.1. Key Components                                                                                                           | 23        |
| 4.2.2. Key Features                                                                                                             | 23        |
| 4.2.3. Playbook Types                                                                                                           | 24        |
| 4.2.4. Playbook Type vs Activity Matrix                                                                                         | 24        |
| 4.2.5. RACI Matrix in a SOC Context                                                                                             | 25        |
| 4.2.6. CACAO Workflow Step Types                                                                                                | 25        |



|                                                                                                                                                                                                                                                                                                                                                                              |           |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------|
| 4.3. Sample Playbook Library . . . . .                                                                                                                                                                                                                                                                                                                                       | 26        |
| 4.3.1. Detect AI Model Training Data Poisoning . . . . .                                                                                                                                                                                                                                                                                                                     | 26        |
| 4.3.2. Multi-Channel Prompt Injection Detection and Response . . . . .                                                                                                                                                                                                                                                                                                       | 27        |
| 4.3.3. Memory Injection Attack (MINJA) Detection and Response . . . . .                                                                                                                                                                                                                                                                                                      | 29        |
| 4.3.4. Poison-RAG Detection and Response . . . . .                                                                                                                                                                                                                                                                                                                           | 30        |
| 4.3.5. SSRF-Based Metadata Credential Abuse in Cloud Infrastructure . . . . .                                                                                                                                                                                                                                                                                                | 32        |
| 4.3.6. AGENTPOISON Memory & RAG Injection Detection and Response . . . . .                                                                                                                                                                                                                                                                                                   | 33        |
| <b>5. Case Studies and Lessons Learned . . . . .</b>                                                                                                                                                                                                                                                                                                                         | <b>35</b> |
| 5.1. Breaking the Prompt Wall Case Study . . . . .                                                                                                                                                                                                                                                                                                                           | 35        |
| 5.2. Memory Injection Attack (MINJA) Case Study . . . . .                                                                                                                                                                                                                                                                                                                    | 37        |
| 5.3. Poison-RAG Case Study . . . . .                                                                                                                                                                                                                                                                                                                                         | 40        |
| 5.4. Capital One Data Breach Case Study . . . . .                                                                                                                                                                                                                                                                                                                            | 42        |
| 5.5. AGENTPOISON Case Study . . . . .                                                                                                                                                                                                                                                                                                                                        | 43        |
| <b>6. Technical Reference: AI Architecture Patterns . . . . .</b>                                                                                                                                                                                                                                                                                                            | <b>45</b> |
| 6.1. Architecture Overview . . . . .                                                                                                                                                                                                                                                                                                                                         | 45        |
| 6.1.1. Levels of Defense Surface . . . . .                                                                                                                                                                                                                                                                                                                                   | 45        |
| 6.1.2. Technology Stack . . . . .                                                                                                                                                                                                                                                                                                                                            | 46        |
| 6.2. Common Patterns and Their Vulnerabilities . . . . .                                                                                                                                                                                                                                                                                                                     | 46        |
| 6.2.1. Basic LLM Architecture . . . . .                                                                                                                                                                                                                                                                                                                                      | 46        |
| 6.2.2. LLM Architecture with Memory . . . . .                                                                                                                                                                                                                                                                                                                                | 48        |
| 6.2.3. RAG Architecture . . . . .                                                                                                                                                                                                                                                                                                                                            | 51        |
| 6.2.4. Agentic Architecture . . . . .                                                                                                                                                                                                                                                                                                                                        | 54        |
| 6.2.5. Agentic RAG Architecture . . . . .                                                                                                                                                                                                                                                                                                                                    | 57        |
| <b>7. Acknowledgements . . . . .</b>                                                                                                                                                                                                                                                                                                                                         | <b>61</b> |
| 7.1. Workstream Leads . . . . .                                                                                                                                                                                                                                                                                                                                              | 61        |
| 7.2. Editors . . . . .                                                                                                                                                                                                                                                                                                                                                       | 61        |
| <b>8. Appendices . . . . .</b>                                                                                                                                                                                                                                                                                                                                               | <b>61</b> |
| 8.1. Integration with Enterprise Security Roadmap . . . . .                                                                                                                                                                                                                                                                                                                  | 61        |
| CoSAI Focus . . . . .                                                                                                                                                                                                                                                                                                                                                        | 64        |
| Guidelines on usage of more advanced AI systems (e.g. large language models (LLMs), multi-modal language models. etc) for drafting documents for OASIS CoSAI: . . . . .                                                                                                                                                                                                      | 65        |
| Iterative refinement: The use of AI systems in drafting standards should be seen as an iterative process, with the generated content serving as a starting point for further refinement and improvement by human experts. Multiple rounds of review and editing may be necessary to ensure the final standards meet the required quality and reliability thresholds. . . . . | 66        |
| Copyright Notice . . . . .                                                                                                                                                                                                                                                                                                                                                   | 66        |



# AI Incident Response Framework, V1.0

---

Approved by the CoSAI Project Governing Board on 27 October 2026

## OASIS Open Project : Coalition for Secure AI (CoSAI) - Workstream 2: Preparing Defenders for a Changing Cybersecurity Landscape

- 1. Overview
  - 1.1. Abstract
  - 1.2. How To Use This Document
    - \* 1.2.1. Quick Start Guide
    - \* 1.2.2. Reader Guidance
  - 1.3. Executive Summary
    - \* 1.3.1. Types of AI Incidents
    - \* 1.3.2 AI Incident Response Plan
    - \* 1.3.3 Key Roles & RACI
    - \* 1.3.4 Incident Response Playbooks
    - \* 1.3.5 Key Lessons Learned
- 2. Understanding AI Security Incidents
  - 2.1. Data Incidents
  - 2.2. Model Incidents
  - 2.3. Deployment Incidents
  - 2.4. Infrastructure Incidents
  - 2.5. User Interaction Incidents
- 3. Creating an AI Incident Response Plan
  - 3.1. Pre-Incident Preparation
  - 3.2. Monitoring and Telemetry
  - 3.3. Incident Response Workflow
    - \* 3.3.1. Preparation Phase
    - \* 3.3.2. Detection and Analysis Phase
    - \* 3.3.3. Containment, Eradication, and Recovery Phase
      - 3.3.3.1. Forensics for AI Systems
    - \* 3.3.4. Post-Incident Activity Phase
  - 3.4. Roles and Responsibilities
    - \* 3.4.1. Provider vs. Consumer Responsibilities
    - \* 3.4.2. Team Assessment and Training
  - 3.5. Communication and Collaboration
    - \* 3.5.1. Regulatory Communication
    - \* 3.5.2. Information Sharing with Community
- 4. Frameworks and Playbooks
  - 4.1. NIST SP 800-61r3 Alignment
    - \* 4.1.1. Key Concept
    - \* 4.1.2. Key Recommendations
  - 4.2. OASIS CACAO Security Playbooks
    - \* 4.2.1. Key Components
    - \* 4.2.2. Key Features
    - \* 4.2.3. Playbook Types
    - \* 4.2.4. Playbook Type vs Activity Matrix
    - \* 4.2.5. RACI Matrix in a SOC Context
    - \* 4.2.6. CACAO Workflow Step Types



- 4.3. Sample Playbook Library
  - \* 4.3.1. Detect AI Model Training Data Poisoning
  - \* 4.3.2. Multi-Channel Prompt Injection Detection and Response
  - \* 4.3.3. Memory Injection Attack (MINJA) Detection and Response
  - \* 4.3.4. Poison-RAG Detection and Response
  - \* 4.3.5. SSRF-Based Metadata Credential Abuse in Cloud Infrastructure
  - \* 4.3.6. AGENTPOISON Memory & RAG Injection Detection and Response
- 5. Case Studies and Lessons Learned
  - 5.1. Breaking the Prompt Wall Case Study
  - 5.2. Memory Injection Attack (MINJA) Case Study
  - 5.3. Poison-RAG Case Study
  - 5.4. Capital One Data Breach Case Study
  - 5.5. AGENTPOISON Case Study
- 6. Technical Reference: AI Architecture Patterns
  - 6.1. Architecture Overview
    - \* 6.1.1. Levels of Defense Surface
    - \* 6.1.2. Technology Stack
  - 6.2. Common Patterns and Their Vulnerabilities
    - \* 6.2.1. Basic LLM Architecture
      - 6.2.1.1. Overview
      - 6.2.1.2. ATLAS Techniques and Tactics
      - 6.2.1.3. Mitigations
    - \* 6.2.2. LLM Architecture with Memory
      - 6.2.2.1. Overview
      - 6.2.2.2. ATLAS Techniques and Tactics
      - 6.2.2.3. Mitigations
    - \* 6.2.3. RAG Architecture
      - 6.2.3.1. Overview
      - 6.2.3.2. ATLAS Techniques and Tactics
      - 6.2.3.3. Mitigations
    - \* 6.2.4. Agentic Architecture
      - 6.2.4.1. Overview
      - 6.2.4.2. ATLAS Techniques and Tactics
      - 6.2.4.3. Mitigations
    - \* 6.2.5. Agentic RAG Architecture
      - 6.2.5.1. Overview
      - 6.2.5.2. ATLAS Techniques and Tactics
      - 6.2.5.3. Mitigations
- 7. Acknowledgements
  - 7.1. Workstream Leads Chairs
  - 7.2. Editors
- 8. Appendices
  - 8.1. Integration with Enterprise Security Roadmap
  - Copyright Notice

## 1. Overview

---



## 1.1. Abstract

This paper addresses the topic of incident response in the context of AI systems. As with other materials produced by the Coalition for Secure AI, this paper focuses on technological capabilities and gaps that are specific to the field of artificial intelligence and does not address the topic of incident response in other contexts as this has been well-researched and documented by existing frameworks (such as NIST Incident Response Recommendations and Considerations for Cybersecurity Risk Management). This paper guides defenders on proactively minimizing the impact of AI system exploitation. It details how to maintain auditability, resiliency, and rapid recovery even when a system is compromised by advanced threat actors. The guide also explores the unique challenges of AI incident response, emphasizing the role of forensic investigation and the complications introduced by agentic architectures, while providing concrete steps to manage this new complexity.

## 1.2. How To Use This Document

### 1.2.1. Quick Start Guide

The incident response life cycle has several versions from varying sources. For the purposes of this paper we will focus on NIST 800-61r3's life cycle phases and provide a mapping to the closest SANS PICERL phases as a reference.

| NIST SP 800-61 rev 3 lifecycle phase | Closest SANS PICERL phase(s) | AI relevance                                                                                                                                                                                       |
|--------------------------------------|------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Govern (GV)</b>                   | Preparation                  | Set AI-specific policy and roles: e.g., who may push new model weights to production, required update cadence for model cards, and breach-disclosure windows for LLM supply-chain compromises.     |
| <b>Identify (ID)</b>                 | Preparation                  | Maintain inventories of datasets, model artefacts, prompts, and dependencies; perform AI-threat risk assessments (model-exfiltration, data-poisoning, jailbreaks) and maintain model card library. |
| <b>Protect (PR)</b>                  | Preparation                  | Apply safeguards: encrypt checkpoints, enforce least-privilege on inference APIs, implement input-validation/output-filtering guardrails, and add watermarking to deter prompt-injection abuse.    |
| <b>Detect (DE)</b>                   | Identification               | Ensure proper Monitoring and Telemetry                                                                                                                                                             |



| NIST SP 800-61 rev 3 lifecycle phase | Closest SANS PICERL phase(s) | AI relevance                                                                                                                                                                                                               |
|--------------------------------------|------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Respond (RS)                         | Containment & Eradication    | Contain (hot-swap to safe snapshots, throttle APIs, revoke leaked keys), eradicate root causes, and coordinate disclosure to affected users.                                                                               |
| Recover (RC)                         | Recovery                     | Recreate or retrain from trusted baselines, run robustness/bias test suites, and validate checkpoints against performance and fairness thresholds before full go-live.                                                     |
| Identify.IM – Improvement            | Lessons Learned              | Feed incident artefacts into future Reinforcement Learning from Human Feedback (RLHF) or adversarial-training cycles and relevant monitoring tools, SIEMs, etc., update threat models and playbooks, share anonymised TTPs |

### 1.2.2. Reader Guidance

This framework is designed for a diverse audience, from technical teams to executive leaders. Not every reader will need every section in equal depth. We outline which sections of the document each audience should prioritize, followed by recommendations on how to use the framework in practice.

*Note: If you are new to this document, you may start with the Quick Start Guide (Section 1.2.1) for a brief overview of the AI incident response lifecycle. This guidance will help you navigate efficiently, whether you're skimming for a high-level overview or planning an in-depth implementation.*

| Audience                              | Introduction                                        | Details                                                                    | Reference / Standards                             |
|---------------------------------------|-----------------------------------------------------|----------------------------------------------------------------------------|---------------------------------------------------|
| <b>CISO / Security leader</b>         | 1.3 Executive Summary                               | 3.4 Roles & Responsibilities<br>3.5 Communication (incl. 3.5.1 Regulatory) | 4.1 NIST SP 800-61r3 Alignment                    |
| <b>Engineer / AI Developer</b>        | 3.1 Pre-Incident Prep<br>3.2 Monitoring & Telemetry | 3.3 IR Workflow                                                            | 6 Technical Architecture Patterns                 |
| <b>Incident Responder (SOC/CSIRT)</b> | 2 Understanding AI Incidents                        | 3.3 IR Workflow (incl. 3.3.3.1 Forensics)<br>4.3 Playbooks                 | 3.4 Roles & Responsibilities<br>3.5 Communication |
| <b>Policymaker / Compliance</b>       | 1.3 Executive Summary                               | 3.1 Prep (policies, third-party)<br>3.5.1 Regulatory                       | 4.1 NIST Alignment                                |
| <b>Technical PM / Project Lead</b>    | 1.3 Executive Summary<br>2 Incident Types           | 3.1<br>3.2<br>3.3 (what to build & run)                                    | 3.4 Who does what                                 |



## 1.3. Executive Summary

The increasing integration of Artificial Intelligence (AI) into business processes presents a new frontier of risks. AI systems, with their inherent complexity and autonomy, can have significant consequences, ranging from reputational damage and financial loss to legal and ethical liabilities. To navigate this evolving landscape, a specialized AI Incident Response Framework is a necessity. This document provides an executive summary of the AI Incident Response Framework developed by COSAI, designed to equip security teams with the necessary capabilities to effectively manage AI-specific incidents. This framework provides a structured approach to preparing for, detecting, responding to, and learning from AI incidents, ensuring organizational resilience in the age of AI.

### 1.3.1 Types of AI Incidents

A shared taxonomy sharpens detection, triage, ownership, and playbook design.

#### A. Performance & Robustness Failures

- **Model Drift:** Accuracy or calibration declines as data distribution shifts.
- **Adversarial Evasion:** Crafted inputs bypass filters or controls.
- **Unexpected/Harmful Outputs:** Nonsensical, policy-violating, or unsafe generations.

#### B. Security Vulnerabilities (ML Pipeline & Model)

- **Data Poisoning:** Training data manipulation alters model behavior.
- **Model Inversion / Membership Inference:** Sensitive training data leakage via queries.
- **Model Theft / Extraction:** Unauthorized replication of proprietary models.

#### C. Ethics, Fairness & Bias Incidents

- **Algorithmic Bias:** Systematic disparate impacts across protected classes.
- **Opacity / Lack of Explainability:** Decisions cannot be justified in regulated contexts.
- **Unintended Social Consequences:** Downstream harm to users or society.

#### D. Misuse & Abuse

- **Deepfakes/Disinformation:** Synthetic media for fraud or manipulation.
- **AI-Assisted Attacks:** AI leveraged to scale phishing, malware, and fraud.

### 1.3.2 AI Incident Response Plan

Adapts the traditional incident response lifecycle to the specific challenges of AI. It provides a structured, repeatable process to ensure timely and effective management of AI incidents. The plan is organized into five key phases:

#### 1. Preparation

- Complete production model inventory (owners, data, risks, criticality).
- Define AI-IR roles, authorities, and decision thresholds.
- Maintain AI-specific playbooks (bias, drift, adversarial, data leakage).
- Run training & simulations (tabletops, red/blue prompt exercises).

#### 2. Detection & Analysis

- Monitor performance, safety, fairness, and security signals (e.g., drift stats, toxicity, jailbreak patterns, data exfiltration indicators).
- Enable intake channels (internal reports, user complaints, vendor alerts).
- Triage by severity, blast radius, regulatory exposure, and patient/customer impact.

#### 3. Containment

- Disable tools or rollback to prior model; gate high-risk features.
- Isolate systems; rate-limit or block malicious inputs.



- Activate temporary guardrails (stricter content filters, human-in-the-loop).

#### 4. Eradication & Recovery

- Root cause analysis (data, model, infra, policy).
- Retrain/patch models; cleanse data; rotate keys/tokens.
- Validate performance, safety, and compliance before redeploy.

#### 5. Post-Incident Activities

- Blameless post-mortem, timeline, and effectiveness review.
- Action items with owners and deadlines; update playbooks, policies, and training.
- Feed lessons into ML governance, risk registers, and CI/CD controls.

### 1.3.3 Key Roles & RACI

Clear accountability avoids delay and rework. Example matrix below; tailor to your org.

| Activity                       | Incident Commander | AI/ML Engineer | Legal & Compliance | Communications |
|--------------------------------|--------------------|----------------|--------------------|----------------|
| <b>Declare Incident</b>        | A                  | R              | I                  | I              |
| <b>Technical Investigation</b> | C                  | A/R            | I                  | I              |
| <b>Take Model Offline</b>      | A                  | R              | C                  | I              |
| <b>External Communications</b> | C                  | I              | A                  | R              |
| <b>Post-Mortem Report</b>      | A                  | R              | C                  | C              |

**Legend:** R = Responsible, A = Accountable, C = Consulted, I = Informed.

### 1.3.4 Incident Response Playbooks

Playbooks translate policy into actionable detections. Each should include the following structure:

- **Incident Type:** A clear description of the specific AI incident the playbook addresses.
- **Detection:** How to identify that this type of incident is occurring, including specific monitoring alerts and reporting channels.
- **Triage and Severity:** Criteria for assessing the severity of the incident and the key stakeholders to notify.
- **Containment Steps:** A checklist of immediate actions to take to contain the incident.
- **Eradication and Recovery Procedures:** Detailed technical steps for resolving the issue and restoring the system.
- **Communication Plan:** Pre-approved communication templates for internal and external stakeholders.

Some example playbooks:

- Adversarial Evasion (Prompt/Content Bypass)
- Model/Data Exfiltration or Leakage

### 1.3.5 Key Lessons Learned

- Run blameless post-mortems within SLA; capture timeline, root causes, and control gaps.
- Produce measurable actions (owners, deadlines) and track to closure.



- Instrument governance: add fairness/safety KPIs to production SLOs; expand telemetry (prompt logs, content filters, tooling traces); integrate security scanning for AI/ML in CI/CD.
- Prevent recurrences: clarify authority to disable/rollback; expand adversarial testing; standardize fallback modes and comms templates.

## 2. Understanding AI Security Incidents

---

AI security incidents can be categorized into five primary domains based on the attack vector and impact area:

### 2.1. Data Incidents

| <i>Incident Type</i>             | <i>Description</i>                                                           | <i>Example</i>                                                                                 |
|----------------------------------|------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------|
| <b>Data Poisoning</b>            | Deliberate contamination of training data to induce specific model behaviors | Injecting fake reviews into a dataset to make an AI model favor specific products artificially |
| <b>Data Leakage</b>              | Unauthorized access to or exfiltration of sensitive data                     | Hackers gaining access to confidential medical records used to train an AI model               |
| <b>Data Integrity Violations</b> | Tampering with AI system inputs to manipulate outcomes                       | Modifying sensor data in a self-driving car to make it misinterpret road signs                 |

### 2.2. Model Incidents

| <i>Incident Type</i>          | <i>Description</i>                                                                | <i>Example</i>                                                                         |
|-------------------------------|-----------------------------------------------------------------------------------|----------------------------------------------------------------------------------------|
| <b>Model Theft/Extraction</b> | Unauthorized duplication or stealing of model weights and architecture            | Reverse engineering a proprietary chatbot's model to replicate its behavior            |
| <b>Model Inversion</b>        | Extracting private training data from model responses                             | Using an AI model's outputs to reconstruct sensitive user data, like medical histories |
| <b>Backdoor Attacks</b>       | Hidden functionality implanted in models that activates under specific conditions | Adding a secret trigger phrase to make an AI classify harmful content as safe          |
| <b>Model Fooling</b>          | Fooling models by slight modifications                                            | Placing a tape over a 85 mph sign, causing an autopilot AI to read it as 35 mph        |

### 2.3. Deployment Incidents



| <i>Incident Type</i>             | <i>Description</i>                                                                                 | <i>Example</i>                                                                                                  |
|----------------------------------|----------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------|
| <b>Prompt Injection</b>          | Manipulating model inputs to bypass safeguards or alter system behavior                            | Sending a specially crafted prompt to make an AI assistant disclose restricted information                      |
| <b>Indirect Prompt Injection</b> | Manipulating external data that the model references to bypass safeguards or alter system behavior | Planting malicious content in locations that an AI uses to make an AI assistant disclose restricted information |
| <b>Jailbreaking</b>              | Circumventing model safety measures to access restricted functionality                             | Using creative prompts to make a chatbot generate unethical or harmful content                                  |
| <b>Abuse Generation</b>          | Harmless prompts trigger harmful output                                                            | Benign prompts causing profane and inappropriate responses                                                      |
| <b>Rogue Agentic AI</b>          | Agentic AI taking unexpected actions                                                               | Closing a curtain when the user says a harmless “thanks” in response to a chat message                          |

## 2.4. Infrastructure Incidents

| <i>Incident Type</i>           | <i>Description</i>                             | <i>Example</i>                                                                             |
|--------------------------------|------------------------------------------------|--------------------------------------------------------------------------------------------|
| <b>Unauthorized Access</b>     | Gaining privileged access to AI infrastructure | Exploiting a vulnerability to gain admin access to an AI-powered recommendation system     |
| <b>API/Service Abuse</b>       | Excessive or malicious use of AI services      | Overloading an AI translation service with excessive requests to disrupt its functionality |
| <b>Dependency Exploitation</b> | Attacking AI dependency libraries              | Exploiting vulnerabilities in the libraries an AI is built upon                            |

## 2.5. User Interaction Incidents

| <i>Incident Type</i>               | <i>Description</i>                                               | <i>Example</i>                                                                   |
|------------------------------------|------------------------------------------------------------------|----------------------------------------------------------------------------------|
| <b>Malicious Code Generation</b>   | Leveraging AI to produce functional malware                      | Using an AI to create polymorphic malware bypassing detection                    |
| <b>Phishing Content Generation</b> | Leveraging AI crafting high-quality, convincing phishing content | Cybercriminals using an AI to craft flawless phishing emails                     |
| <b>Output Manipulation</b>         | Forcing harmful or misleading outputs from AI systems            | Crafting inputs to make an AI generate false financial forecasts                 |
| <b>Hallucination Exploitation</b>  | Leveraging model inaccuracies for harmful purposes               | Using an AI’s hallucinated data to spread misinformation in public forums        |
| <b>Regional or cultural bias</b>   | Not acceptable output of the model in some regions or cultures   | Exploiting AI’s specifics to provide different outputs for different or cultures |



### 3. Creating an AI Incident Response Plan

---

The **core objectives** of an AI Incident Response Plan are to: - Minimize business and user impact during incidents - Protect sensitive data from AI system exploitation - Maintain regulatory compliance - Preserve evidence for forensic analysis - Identify root causes for prevention and response - Restore secure AI system functionality - Share actionable intelligence with security community

The AI Incident Response Plan described is limited in scope to the: - AI systems covered (LLMs, RAG systems, agentic architectures; further discussed in Section 6) - Integration points with traditional IR procedures and frameworks (further discussed in Section 4)

#### 3.1. Pre-Incident Preparation

| Component                  | Key Elements                                                                                                                                                                                                                                                                                                                                                    |
|----------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Risk Assessment</b>     | <ul style="list-style-type: none"> <li>· Secure AI system components and architecture</li> <li>· Critical data assets processed by AI systems</li> <li>· Potential attack vectors based on ATLAS framework</li> <li>· Compliance requirements</li> </ul>                                                                                                        |
| <b>IR Team Composition</b> | <ul style="list-style-type: none"> <li>· Business impact of potential compromise</li> <li>· AI/ML engineers with model expertise</li> <li>· Data scientists for pipeline understanding</li> <li>· Security analysts with AI knowledge</li> <li>· Legal and compliance experts</li> <li>· Business stakeholders</li> <li>· Communications specialists</li> </ul> |
| <b>Baseline Telemetry</b>  | <ul style="list-style-type: none"> <li>· Input/output logging for all interactions</li> <li>· System prompt and configuration versioning</li> <li>· Tool usage tracking (agentic systems)</li> <li>· Authentication and authorization events</li> <li>· Embedding and vector database query patterns</li> <li>· Resource utilization metrics</li> </ul>         |

#### 3.2. Monitoring and Telemetry

To protect AI systems from evolving threats, telemetry must capture a comprehensive set of signals spanning model inference behavior, prompt and output risks, content safety, and model integrity. This includes monitoring for prompt injection, output manipulation, knowledge base poisoning, model drift, unauthorized tool or API usage, and other adversarial activities. Tracking agent workflows, context exchanges, and system lifecycles provides visibility into misuse, operational anomalies, and attacks. This data enables organizations to safeguard AI pipelines, ensure compliance, and strengthen incident detection and response with traceable, structured observability data.

| Telemetry Type           | Purpose                                                                                    | Relevant Threats Detected                                                       |
|--------------------------|--------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------|
| Model Inference Activity | Tracks all inference executions including token usage, confidence, latency, and endpoints. | Abuse of model APIs, data exfiltration, inference abuse (e.g., model scraping). |



| Telemetry Type                                | Purpose                                                                                                                    | Relevant Threats Detected                                                            |
|-----------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------|
| Prompt Injection Detection                    | Monitors input prompts for adversarial payloads or jailbreak attempts.                                                     | Prompt injection attacks, jailbreaking, prompt leakage.                              |
| Content Risk Detection                        | Assesses LLM output for toxicity, bias, PII leakage, and unsafe content generation.                                        | Disinformation, output poisoning, PII/PHI leakage, model misuse.                     |
| Model Drift Detection                         | Detects distributional shifts in model input/output indicative of adversarial manipulation or operational risk.            | Data poisoning, concept drift, model degradation over time.                          |
| Agentic Workflow Execution                    | Logs multi-step decision-making by AI agents, including tool calls and retries.                                            | Malicious tool chaining, unauthorized data access, misuse of plugins/APIs.           |
| MCP (Model Context Protocol) Message Activity | Observes protocol-level context exchanges in agentic and retrieval-augmented systems.                                      | Session hijacking, tampering with context/memory, unauthorized prompt modifications. |
| Agentic RAG Workflow Execution                | Captures the steps of retrieval-augmented generation pipelines including retrievers, LLMs, and validators.                 | Knowledge base poisoning, unsafe content retrieval, response manipulation.           |
| AI Incident Finding                           | Records confirmed or suspected incidents detected across AI systems (e.g., poisoning, unauthorized use, guardrail events). | Incident response triggers — prompt abuse, model misuse, security violations.        |
| Tool Call Execution                           | Logs API/tool calls made by agents, including argument inspection and error rates.                                         | Unauthorized external access, misuse of internal or third-party APIs.                |
| AI System Activity                            | Lifecycle monitoring of AI systems: deployments, updates, shutdowns, failures.                                             | Unauthorized model deployments, version tampering, system compromise.                |

Effective security monitoring of AI systems requires capturing fine-grained telemetry at each of the input, processing, and output stages of AI workflows. Below are the some of the attributes essential for building threat detection use cases, anomaly detection, and incident response for AI-driven environments:

| Attribute       | Explanation                                                                                                                                                                                                                                         |
|-----------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| timestamp       | Tracks the time of activity performed. Proper recording of time enables investigators to piece a timeline together especially when AI incidents span multiple services or layers.                                                                   |
| identity        | Tracks the identity that performed an action. This could be a user or machine identity.                                                                                                                                                             |
| prompt_text     | Captures the raw input text or structured request sent to the AI model. Monitoring the prompt_text is crucial for detecting prompt injection attacks, data exfiltration attempts, or adversarial instructions that could manipulate model behavior. |
| completion_text | Captures the model's output or response. It is essential to inspect completion_text for toxic content, bias, hallucinations, or leakage of sensitive information such as PII, credentials, or internal data.                                        |



| Attribute           | Explanation                                                                                                                                                                                                                                   |
|---------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| model_name/version  | Logs the name and version of the model used. This ensures that only authorized and validated models are deployed and helps detect unauthorized model swapping or shadow AI deployments.                                                       |
| token_usage         | Measures the number of tokens in both input and output. Unusually high or low token_usage patterns can signal abusive scraping attempts or denial-of-service (DoS) behavior aimed at exhausting AI model resources.                           |
| latency_ms          | Measures the inference response time in milliseconds. Sudden spikes or drops in latency_ms may indicate resource exhaustion attacks, model misuse, or service degradation due to adversarial workloads.                                       |
| retrieved_docs      | Tracks documents fetched during Retrieval-Augmented Generation (RAG) workflows. It helps detect knowledge base poisoning, unauthorized document access, or retrieval manipulation that could compromise model outputs.                        |
| tool_calls          | Logs all external or internal tools invoked by AI agents. Monitoring tool_calls helps uncover unauthorized access to sensitive APIs, malicious tool chaining, or exfiltration via tool misuse during autonomous agent operations.             |
| confidence_score    | Provides the model's confidence in its prediction or generation. Very low confidence scores can signal model uncertainty or adversarial influence, while unusually high confidence on manipulated inputs could reveal model evasion attempts. |
| drift_metrics       | Measures data drift or distribution changes in model inputs/outputs over time. Significant drift could indicate model poisoning, concept drift attacks, or environmental changes that degrade model performance or security.                  |
| trace_id/session_id | Correlates activities across sessions and distributed systems. Proper tracking of trace_id and session_id enables forensic investigations, cross-system correlation, and incident scoping when AI incidents span multiple services or layers. |

### 3.3. Incident Response Workflow

#### 3.3.1. Preparation Phase

| Activity                                     | AI-Specific Considerations                                                                                                                                                                                                                                                      | Implementation Guidance                                                                                                                                                                                                                                                    |
|----------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Risk Assessment &amp; Threat Modeling</b> | <ul style="list-style-type: none"> <li>· Identify critical assets in each architecture pattern</li> <li>· Map ATLAS threat vectors to components</li> <li>· Assess likelihood and impact</li> <li>· Prioritize security controls</li> <li>· Document risk thresholds</li> </ul> | <ul style="list-style-type: none"> <li>· Use architectural patterns from Section 6</li> <li>· Reference ATLAS techniques in security planning</li> <li>· Conduct architecture-specific threat modeling</li> <li>· Prioritize controls for high-impact scenarios</li> </ul> |



| Activity                               | AI-Specific Considerations                                                                                                                                                                                                                                                        | Implementation Guidance                                                                                                                                                                                                                                                                          |
|----------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Monitoring Infrastructure</b>       | <ul style="list-style-type: none"> <li>· Input/output logging</li> <li>· System prompt versioning</li> <li>· Configuration management</li> <li>· Authentication tracking</li> <li>· API usage analysis</li> <li>· Resource utilization</li> <li>· Data access patterns</li> </ul> | <ul style="list-style-type: none"> <li>· Implement privacy controls for prompt logging</li> <li>· Maintain auditable prompt history</li> <li>· Track all model versions and parameters</li> <li>· Establish API usage baselines</li> <li>· Monitor for anomalous resource consumption</li> </ul> |
| <b>Response Capability Development</b> | <ul style="list-style-type: none"> <li>· Specialized AI security training</li> <li>· Response tools for AI systems</li> <li>· Emergency access procedures</li> <li>· Backup procedures</li> <li>· Recovery testing</li> </ul>                                                     | <ul style="list-style-type: none"> <li>· Train teams on AI-specific threats</li> <li>· Implement tools for log analysis and prompt filtering</li> <li>· Create secure backups of models and embeddings</li> <li>· Regularly test restoration procedures</li> </ul>                               |

### 3.3.2. Detection and Analysis Phase

| Activity                    | AI-Specific Considerations                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Key Considerations                                                                                                                                                                                                                                                                                                                                                                                                                       |
|-----------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Detection Mechanisms</b> | <p><b>Automated Monitoring:</b></p> <ul style="list-style-type: none"> <li>· LLM-based classifiers for suspicious interactions</li> <li>· Semantic similarity checks</li> <li>· Token usage anomaly detection</li> <li>· User feedback signal monitoring</li> </ul> <p><b>**Manual Review:**</b></p> <ul style="list-style-type: none"> <li>· Human review procedures for flagged interactions</li> <li>· Sampling of high-risk operations</li> <li>· Security dashboards</li> <li>· Escalation triggers</li> </ul> <p><b>**Integration Points:**</b></p> <ul style="list-style-type: none"> <li>· SIEM integration</li> <li>· Correlation with network security</li> <li>· Authentication event linking</li> </ul> | <ul style="list-style-type: none"> <li>· Balance automated detection with human oversight</li> <li>· Establish clear thresholds for alerts</li> <li>· Avoid false positives through baseline calibration</li> <li>· Ensure privacy considerations in monitoring</li> <li>· Correlate AI system events with broader security telemetry</li> <li>· Implement strong AI security controls for all LLM-based monitoring solutions</li> </ul> |



| Activity                 | AI-Specific Considerations                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           | Key Considerations                                                                                                                                                                                                                                                                                                                                                                                                                                |
|--------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Initial Triage           | <p><b>Incident Verification:</b></p> <ul style="list-style-type: none"> <li>· Confirm security incident status</li> <li>· Classify per AI attack categories (Section 2)</li> <li>· Determine affected components</li> </ul> <p><b>**Impact Assessment:**</b></p> <ul style="list-style-type: none"> <li>· Estimate compromise scope</li> <li>· Identify affected users/data</li> <li>· Assess data exposure issues</li> <li>· Determine business impact</li> </ul> <p><b>**Priority Assignment:**</b></p> <ul style="list-style-type: none"> <li>· Assign priority based on impact</li> <li>· Implement notification procedures</li> <li>· Mobilize specialized resources</li> </ul>                                                                                                 | <ul style="list-style-type: none"> <li>· Leverage AI to aggregate data needed for initial assessment of incident</li> <li>· Validate whether anomalous behavior constitutes a security incident</li> <li>· Use categories from Section 2 for classification</li> <li>· Focus on potential data exposure through AI systems</li> <li>· Consider regulatory implications</li> <li>· Ensure escalation procedures match incident severity</li> </ul> |
| Investigation Procedures | <p><b>Evidence Collection:</b></p> <ul style="list-style-type: none"> <li>· Interaction logs</li> <li>· System configurations</li> <li>· Model versions and parameters</li> <li>· Network traffic and API calls</li> <li>· Vector database queries</li> </ul> <p><b>**Forensic Analysis:**</b></p> <ul style="list-style-type: none"> <li>· Malicious prompt patterns</li> <li>· Data retrieval sequences</li> <li>· Tool usage analysis</li> <li>· Memory manipulation assessment</li> <li>· Data poisoning evaluation</li> </ul> <p><b>**Attribution Assessment:**</b></p> <ul style="list-style-type: none"> <li>· Targeted vs. opportunistic analysis</li> <li>· Actor identification</li> <li>· Threat intelligence correlation</li> <li>· Attribution documentation</li> </ul> | <ul style="list-style-type: none"> <li>· Use AI to assist in analyzing and summarizing evidence documents</li> <li>· Preserve evidence in a forensically sound manner</li> <li>· Focus on the unique aspects of AI system compromise</li> <li>· Analyze both the explicit content and implicit patterns</li> <li>· Consider the non-deterministic nature of AI systems</li> <li>· Document investigation methodology thoroughly</li> </ul>        |

### 3.3.3. Containment, Eradication, and Recovery Phase



| Activity                      | AI-Specific Considerations                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | Implementation Guidance                                                                                                                                                                                                                                                                                                                                                                                                |
|-------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Containment Strategies</b> | <p><b>Short-term Containment:</b></p> <ul style="list-style-type: none"> <li>· Emergency prompt filters</li> <li>· Rate limiting/access restrictions</li> <li>· Component isolation</li> <li>· User/IP blocking</li> </ul> <p><b>**Architecture-Specific Containment:</b></p> <ul style="list-style-type: none"> <li>· Basic LLM: <b>Input validation, reduced temperature, content filtering</b></li> <li>· LLM with Memory: <b>Reset memory, session isolation, conversation history purging</b></li> <li>· RAG: <b>Data source quarantine, validation of retrieval results</b></li> <li>· Agentic: <b>Disable tools, enhance authorization, least privilege checks</b></li> <li>· Agentic RAG:<b>** Combined strategies with interface security focus</b></li> </ul> <p><b>**Evidence Preservation:**</b></p> <ul style="list-style-type: none"> <li>· Forensic copies before containment</li> <li>· Documented containment actions</li> <li>· Vector database state preservation</li> <li>· Model weight copies</li> </ul> | <ul style="list-style-type: none"> <li>· Implement containment without destroying evidence</li> <li>· Consider each architectural pattern's unique vulnerabilities</li> <li>· Focus containment on the specific attack vector</li> <li>· Document all containment actions thoroughly</li> <li>· Maintain chain of custody for forensic evidence</li> <li>· Consider business impact of containment measures</li> </ul> |
| <b>Eradication Procedures</b> | <p><b>Root Cause Elimination:</b></p> <ul style="list-style-type: none"> <li>· Vulnerability addressing</li> <li>· Enhanced input validation</li> <li>· System prompt updates</li> <li>· Framework patches</li> </ul> <p><b>**Component Remediation:</b></p> <ul style="list-style-type: none"> <li>· Data-Level: <b>Remove poisoned data, reconstruct embeddings</b></li> <li>· Model-Level: <b>Roll back to secure versions, enhance guardrails</b></li> <li>· Deployment-Level: <b>Update prompts, strengthen controls</b></li> <li>· Output-Level:<b>** Implement filtering, content moderation</b></li> </ul> <p><b>**Verification Testing:**</b></p> <ul style="list-style-type: none"> <li>· Penetration testing</li> <li>· Attack reproduction attempts</li> <li>· Red-team exercises</li> <li>· Issue verification</li> </ul>                                                                                                                                                                                         | <ul style="list-style-type: none"> <li>· Address fundamental vulnerabilities, not just symptoms</li> <li>· Apply remediations specific to the compromised level</li> <li>· Test thoroughly before returning to production</li> <li>· Document all changes for future reference</li> <li>· Maintain audit trail of remediation actions</li> <li>· Verify effectiveness through testing</li> </ul>                       |



| Activity                   | AI-Specific Considerations                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | Implementation Guidance                                                                                                                                                                                                                                                                                                                                                                    |
|----------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Recovery Procedures</b> | <p><b>Service Restoration:</b></p> <ul style="list-style-type: none"> <li>· Capability prioritization</li> <li>· Phased approach planning</li> <li>· Verification criteria</li> <li>· Required approvals</li> </ul> <p><b>**Enhanced Monitoring:**</b></p> <ul style="list-style-type: none"> <li>· Heightened monitoring deployment</li> <li>· Normal/abnormal metrics</li> <li>· Additional logging</li> <li>· Scheduled reviews</li> </ul> <p><b>**User Communication:**</b></p> <ul style="list-style-type: none"> <li>· Communication templates</li> <li>· Transparency guidelines</li> <li>· Communication protocols</li> <li>· User feedback mechanisms</li> </ul> <p><b>**Business Continuity:**</b></p> <ul style="list-style-type: none"> <li>· Alternative system activation</li> <li>· Restoration prioritization</li> <li>· Stakeholder expectation setting</li> <li>· Impact documentation</li> </ul> | <ul style="list-style-type: none"> <li>· Implement a gradual, verified restoration process</li> <li>· Maintain heightened monitoring during recovery</li> <li>· Communicate clearly with affected users</li> <li>· Document business impact throughout recovery</li> <li>· Ensure leadership approval at key recovery stages</li> <li>· Balance security with operational needs</li> </ul> |

### 3.3.3.1. Forensics for AI Systems

In order to effectively conduct forensic analysis of security incidents on AI systems, the AI system should be able to record (and provide) sufficient amount of information for the investigator. For example, it is essential to be able to retrieve both system prompt and user prompts and any other parameters which were used when interacting with an AI model as well as the response from the system. In general AI model will not produce exactly the same outputs even given the same input and system prompt, therefore it is essential to preserve raw output details from the AI model in order to be able to reconstruct the whole picture of the incident.

AI systems can also interface and integrate with external functions and components (such as MCP servers). As with traditional security, accurate logging of these interactions and a log trail of external function executions should be preserved and made available to the forensic investigator.

User inputs may be used as a part of future training data for an AI system, therefore it is essential to be able to track user input inclusion and impact on training datasets.

AI-based tools can also be used for effective analysis of forensic trails after an incident. However, similar to how traditional forensic tools can be attacked and exploited by malicious attacker - the investigator should be aware that AI-enabled analysis and forensic tools can also become a target of malicious activity. The risks of AI-based tool analysis should be evaluated based upon your organizations risk tolerances and mitigations should be leveraged including the use of sandboxed, isolated environments.

Clearly, in this section we describe the expectations when there are no constraints. However, in real world AI system implementation resources can be finite and system constraints may impact the scope and scale of data preservation. As in traditional systems, the best practice is to identify data retention objectives at design time and to address resource constraints via data preservation techniques and data rotation during the implementation and deployment of AI system.



During analysis the investigator can ask the following questions: \* Were any of the guardrails successfully bypassed? \* Were any of the prompts intentionally modified by the attacker? \* Did the system produce unexpected and potentially harmful content? \* Did content include any sensitive information that could be potentially exfiltrated by an attacker? \* Did any other manipulation of data or system take place?

Finally, it is essential to understand and investigate the external components that were interacting with AI systems. Log trails of these components (such as web server logs, database query logs and so on) could also be extremely helpful in being able to reconstruct and understand the full picture of an AI security incident.

### 3.3.4. Post-Incident Activity Phase

| Activity               | Key Components                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                | Implementation Guidance                                                                                                                                                                                                                                                                                                                                                                                                       |
|------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Lessons Learned</b> | <p><b>Post-Incident Review:</b></p> <ul style="list-style-type: none"> <li>· Timely meeting scheduling</li> <li>· Cross-team representation</li> <li>· Event timeline documentation</li> <li>· Response effectiveness analysis</li> </ul> <p><b>**Root Cause Analysis:**</b></p> <ul style="list-style-type: none"> <li>· Underlying vulnerability investigation</li> <li>· Determination of causation:               <ul style="list-style-type: none"> <li>· Architectural design flaws</li> <li>· Security control implementation errors</li> <li>· Operational failures</li> <li>· Novel attack techniques</li> </ul> </li> </ul> <p><b>**Response Evaluation:**</b></p> <ul style="list-style-type: none"> <li>· Detection effectiveness assessment</li> <li>· Containment strategy evaluation</li> <li>· Recovery efficiency analysis</li> <li>· Business impact measurement</li> </ul> | <ul style="list-style-type: none"> <li>· Schedule reviews within 1-2 weeks of resolution</li> <li>· Focus on systemic improvements, not blame</li> <li>· Document timeline from detection through resolution</li> <li>· Identify gaps in tools, processes, and training</li> <li>· Analyze whether incident resulted from known or novel techniques</li> <li>· Measure response metrics against defined objectives</li> </ul> |



| Activity                    | Key Components                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | Implementation Guidance                                                                                                                                                                                                                                                                                                                                                                                                                                |
|-----------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Security Enhancement</b> | <p><b>Architectural Improvements:</b></p> <ul style="list-style-type: none"> <li>· System architecture updates</li> <li>· Additional security layers</li> <li>· Integration point strengthening</li> <li>· Component separation enhancement</li> </ul> <p><b>**Control Enhancements:**</b></p> <ul style="list-style-type: none"> <li>· New/improved security controls</li> <li>· Prompt engineering updates</li> <li>· Monitoring capability enhancement</li> <li>· Access control strengthening</li> </ul> <p><b>**Process Refinement:**</b></p> <ul style="list-style-type: none"> <li>· Playbook updates</li> <li>· Role/responsibility revision</li> <li>· Communication protocol improvement</li> <li>· Documentation enhancement</li> </ul> | <ul style="list-style-type: none"> <li>· Implement architectural changes to address root causes</li> <li>· Strengthen controls at the affected system level</li> <li>· Update detection mechanisms for similar threats</li> <li>· Revise processes based on response effectiveness</li> <li>· Document all improvements for future reference</li> <li>· Assign clear ownership for enhancement implementation</li> </ul>                               |
| <b>Knowledge Sharing</b>    | <p><b>Internal Transfer:</b></p> <ul style="list-style-type: none"> <li>· Case study documentation</li> <li>· Lesson-sharing workshops</li> <li>· Security awareness updates</li> <li>· Executive briefings</li> </ul> <p><b>**External Sharing:**</b></p> <ul style="list-style-type: none"> <li>· Industry knowledge contribution</li> <li>· Information-sharing participation</li> <li>· MITRE ATLAS submissions</li> <li>· Standards organization engagement</li> </ul> <p><b>**Regulatory Reporting:**</b></p> <ul style="list-style-type: none"> <li>· Required notifications</li> <li>· Compliance documentation</li> <li>· Security measure evidence</li> <li>· Audit record maintenance</li> </ul>                                        | <ul style="list-style-type: none"> <li>· Leverage AI in drafting incident reports</li> <li>· Create training materials from incident findings</li> <li>· Share lessons across security and AI teams</li> <li>· Contribute to industry knowledge responsibly</li> <li>· Complete all regulatory reporting requirements</li> <li>· Balance transparency with security considerations</li> <li>· Document compliance with relevant regulations</li> </ul> |



### 3.4. Roles and Responsibilities

#### 3.4.1. Provider vs. Consumer Responsibilities

A shared responsibility model between providers and consumers is essential to establish clear accountability and separation of duties in order to safeguard AI systems effectively. Providers—including base model developers, agent developers, and API hosts—should be responsible for ensuring foundational model robustness, implementing safety alignment, securing inference infrastructure, and embedding standardized guardrails. Consumers—comprising application developers and end organizations—must manage prompt and context control, application-level filtering, compliance monitoring, and incident detection.

#### 3.4.2. Team Assessment and Training

| Activity                          | Description                                                       | Implementation Guidance                                                                                                                                                                                                    |
|-----------------------------------|-------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Tabletop Exercises</b>         | Scenario-based discussions using AI attack cases from Section 2   | <ul style="list-style-type: none"> <li>· Conduct quarterly exercises</li> <li>· Rotate scenario types</li> <li>· Include cross-functional teams</li> <li>· Document lessons learned</li> </ul>                             |
| <b>Technical Drills</b>           | Hands-on exercises focused on specific AI incident response tasks | <ul style="list-style-type: none"> <li>· Practice log analysis</li> <li>· Test containment procedures</li> <li>· Validate recovery processes</li> <li>· Measure response times</li> </ul>                                  |
| <b>Cross-Functional Exercises</b> | Full-team simulations involving all relevant stakeholders         | <ul style="list-style-type: none"> <li>· Include business representatives</li> <li>· Test communications channels</li> <li>· Validate escalation procedures</li> <li>· Practice regulatory notifications</li> </ul>        |
| <b>Adversarial Simulations</b>    | Red-team exercises based on ATLAS techniques                      | <ul style="list-style-type: none"> <li>· Simulate prompt injection attacks</li> <li>· Test data poisoning scenarios</li> <li>· Practice response to model extraction</li> <li>· Evaluate detection capabilities</li> </ul> |
| <b>Skills Assessment</b>          | Regular evaluation of team knowledge and capabilities             | <ul style="list-style-type: none"> <li>· Identify training gaps</li> <li>· Update learning materials</li> <li>· Track improvement metrics</li> <li>· Recognize high performers</li> </ul>                                  |

### 3.5. Communication and Collaboration

#### 3.5.1. Regulatory Communication

Incident reporting has different requirements for different jurisdictions and industry verticals, manifold across levels of government. In the US the 2022 CIRCIA legislation provided strict timelines for reporting cyber incidents to critical infrastructure segments even tighter than the financial industry’s GLBA for financial privacy requirements, while HIPAA requires report of PHI breaches over longer timescales. All 50 US states have their own breach notification laws. In contrast, the EU is governed by NIS2 for critical infrastructure, DORA for financial institutions and GDPR for personal data breaches - all on strict timelines. Legal counsel is needed to fully understand the reporting requirements across jurisdictions.

AI Systems bring new challenges to reporting for regulatory compliance: \* Personal data breaches from PII/PHI leakage from AI systems, applicable to US states’ data laws and GDPR \* Model confidentiality, integrity and accessibility attacks are breaches, even if affected indirectly \* AI service providers may be considered critical infrastructure, or be required to disclose incidents affecting customer use of AI services



### 3.5.2. Information Sharing with Community

One of defenders' greatest advantages against attackers is information sharing, as timely well-structured information about attacks enables defenders to prevent similar attacks. Incident information sharing is evolving from sector based ISAC reporting to more strict national and supra-national reporting requirements from CIRCIA, CISA in the US and ENISA in the EU. There are many fragmented and boutique reporting mechanisms too, including the AI Incident Database, for community-driven reporting. Cross-border open groups such as FIRST offer best practices and standards for sharing IOCs, TTPs and for Coordinated Vulnerability Disclosure. The security community supports the use of OASIS STIX 2.1 to share information widely, with OpenCTI the intended target open source tool for managing threat intelligence within or across organizations. Open sharing of threat intelligence is still a challenge for enterprises and even the security community, often because complying with data privacy regulation discourages additional anonymization work required.

AI incident reporting will often require: \* Data scrubbing to remove PII from incident report, possibly data used by AI system \* Associating IOCs and TTPs with attacks against the AI system \* Responsible disclosure of AI vulnerabilities to AI providers

## 4. Frameworks and Playbooks

---

### 4.1. NIST SP 800-61r3 Alignment

**Reference:** NIST SP 800-61r3

The NIST SP 800-61 Revision 3 (2025) modernizes incident response by embedding it within the broader context of cybersecurity risk management through alignment with the NIST Cybersecurity Framework (CSF) 2.0. Unlike its predecessor, which focused on procedural guidance, this version adopts a flexible, outcome-driven approach using the six CSF Functions—Govern, Identify, Protect, Detect, Respond, and Recover—to guide organizations in managing, mitigating, and learning from cybersecurity incidents. The framework emphasizes continuous improvement, cross-functional coordination, and real-time learning across incident stages, integrating threat intelligence, asset visibility, and policy enforcement into an adaptive lifecycle. For AI systems and generative AI applications, the framework is especially applicable as it supports incident response for AI-specific threats such as prompt injection, model inversion, hallucinations, data leakage, adversarial misuse of toolchains (e.g., LangChain agents), and retrieval augmentation abuse. It enables organizations to adapt incident handling playbooks, define AI-relevant response roles, log and monitor AI inputs/outputs, and apply continuous learning to secure model and agent behaviors and limit data exposures. Through its flexible and modular structure, SP 800-61r3 provides a resilient foundation for defending dynamic, intelligent, and increasingly autonomous AI ecosystems.

#### 4.1.1. Key Concept

SP 800-61r3 shifts from static procedural guidance to a flexible, CSF 2.0-aligned risk management profile, suitable for dynamic, evolving technologies like AI systems.

It replaces the old linear IR lifecycle with a model mapped to the six CSF 2.0 Functions: Govern, Identify, Protect, Detect, Respond, and Recover — all of which are critical for generative AI environments where risks emerge through prompt injection, model behavior, data flow, and tool orchestration.



| CSF Function         | AI-Specific Considerations                                                                                                                                                                                                                      |
|----------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Govern (GV)</b>   | Establish AI-specific IR policies: include prompt injection, misuse of RAG pipelines, and training data exposures. Ensure third-party responsibilities (e.g., model vendors or SaaS orchestration) are contractually defined.                   |
| <b>Identify (ID)</b> | Inventory LLM tools, APIs, model endpoints, external connectors and data sources. Capture attack surfaces such as embeddings, memory modules, toolchains (LangChain, ReAct). Threat modeling must include jailbreaking and adversarial prompts. |
| <b>Protect (PR)</b>  | Implement robust input/output guardrails, system prompts, and behavioral constraints. Ensure LLM telemetry, prompt logs, and PII filters are enabled. Validate training pipelines to prevent data poisoning or unintentional PII ingestion.     |
| <b>Detect (DE)</b>   | Continuously monitor prompt patterns, output anomalies, and system behavior. Correlate logs from LLM APIs, tool executions, and chat sessions. Integrate LLM-as-a-Judge metrics and fine-tuned filters to detect AI-targeted abuse.             |
| <b>Respond (RS)</b>  | Coordinate across AI/ML teams, SREs, legal, and security to investigate and contain AI-specific incidents like data leakage or unauthorized tool invocation. Update response plans to reflect LLM threats and recovery protocols.               |
| <b>Recover (RC)</b>  | Restore compromised AI services, rollback fine-tuned models or agent states, redact unsafe memory or embeddings. Update RAG sources and retrain if training data was involved in the breach. Conduct post-mortems with AI-specific lens.        |

#### 4.1.2. Key Recommendations

| <i>NIST Recommendation</i>                  | <i>AI System Implication</i>                                                                                                                         |
|---------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Establish IR policy &amp; procedures</b> | Include AI-specific incident types and response playbooks for AI misuse, such as unauthorized API use, hallucinations, or data leakage.              |
| <b>Communication protocols</b>              | Define rules for disclosure and response involving model vendors (e.g., OpenAI, Anthropic), data controllers, and regulators (for privacy breaches). |
| <b>Team structure</b>                       | Ensure inclusion of AI/ML engineers, data scientists, legal advisors, and system architects alongside traditional CSIRTs.                            |
| <b>Shared Responsibility</b>                | Contracts must clearly define incident roles/responsibilities across cloud services, model providers, and orchestration tools.                       |
| <b>Telemetry &amp; Forensics</b>            | Log prompt/response pairs, tool execution traces, embeddings, and system prompts for investigation.                                                  |
| <b>Automated detection</b>                  | Leverage metrics from LLM guardrails, behavioral scoring, and RAG pipeline telemetry to detect attacks (e.g., prompt-based exfiltration).            |
| <b>Incident Response Playbooks</b>          | Define specific playbooks for generative AI issues: misuse of agents, hallucination of critical data, unauthorized RAG retrieval.                    |
| <b>Information sharing</b>                  | Coordinate with external AI incident hubs (like responsible AI research communities or regulatory bodies) for emerging attack techniques.            |



| <i>NIST Recommendation</i>            | <i>AI System Implication</i>                                                                                                                              |
|---------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Continuous Improvement (ID.IM)</b> | Post-incident lessons should inform prompt engineering updates, detection engineering improvements, model fine-tuning policies, and toolchain governance. |

## 4.2. OASIS CACAO Security Playbooks

**Reference:** (OASIS CACAO Security Playbooks Version 2.0)

The Collaborative Automated Course of Action Operations (CACAO) Security Playbooks Version 2.0 specification defines a standardized schema and taxonomy for representing cybersecurity playbooks. These playbooks are structured workflows that define sequences of security actions — detection, prevention, mitigation, remediation, etc. — which can be executed manually or automatically across different tools and systems. CACAO enables the creation, sharing, and execution of these playbooks across organizational and technological boundaries, improving collaboration, incident response, and cyber defense automation.

### 4.2.1. Key Components

- **Playbooks:** - Represent orchestrated actions in response to threats (e.g., investigate, notify, isolate).
- **Workflows:** - Structured steps including conditional logic, parallel actions, and sub-playbook invocation.
- **Commands:** - Encapsulate specific executable operations (e.g., Bash, PowerShell, OpenC2, Yara).
- **Agents & Targets:** - Define who executes the actions and the affected systems/entities.
- **Authentication Info:** - Standardizes credentials used for automation.
- **Data Markings & Digital Signatures:** - Ensure access control, trust, and data integrity.
- **Versioning:** - Built-in mechanism to track and revoke playbooks with strict creator ownership rules.
- **Extension Support:** - Allows future customization and modular augmentation.

### 4.2.2. Key Features

| Feature                    | Description                                                                                                                          |
|----------------------------|--------------------------------------------------------------------------------------------------------------------------------------|
| <b>Workflow Logic</b>      | Enables orchestration of actions using structured steps including conditions and parallelism.                                        |
| <b>Playbook Types</b>      | Supports 8 defined types (see 4.2.3 below) for diverse response needs.                                                               |
| <b>Activity Vocabulary</b> | Defines granular actions tied to playbook types, enhancing clarity and automation.                                                   |
| <b>Command Abstraction</b> | Supports multiple command formats (e.g., Bash, PowerShell, OpenC2, Jupyter Notebooks).                                               |
| <b>Agents and Targets</b>  | Identifies executors and affected entities, enhancing contextual execution.                                                          |
| <b>Authentication Info</b> | Separates auth definitions from agents/targets, enabling reuse and secure handling.                                                  |
| <b>Versioning</b>          | Built-in lifecycle control, allowing updates and revocations of playbooks by original creator.                                       |
| <b>Data Markings</b>       | Provides Traffic Light Protocol(TLP), Information Exchange Protocol(IEP), and custom markings to enforce usage and sharing policies. |
| <b>Digital Signatures</b>  | Supports JSON Signature Scheme for integrity, non-repudiation, and trust validation.                                                 |



| Feature                           | Description                                                                                     |
|-----------------------------------|-------------------------------------------------------------------------------------------------|
| <b>Playbook Referencing</b>       | Allows modular composition via playbook-to-playbook invocation.                                 |
| <b>STIX Integration</b>           | Leverages STIX 2.1 objects (e.g., Identity, Relationship) for compatibility with CTI platforms. |
| <b>Extension Mechanism</b>        | Allows for custom extensions to augment schema without breaking compatibility.                  |
| <b>Impact, Severity, Priority</b> | Supports scoring playbooks for operational context and automated triage.                        |

### 4.2.3. Playbook Types

| Playbook Type        | Description                                                                                            |
|----------------------|--------------------------------------------------------------------------------------------------------|
| <b>Attack</b>        | Orchestrates penetration testing or adversary emulation steps; used by red teams or simulation tools.  |
| <b>Detection</b>     | Focuses on identifying known threats, behaviors, or indicators using logs, telemetry, or threat intel. |
| <b>Engagement</b>    | Defines deception or denial tactics to disrupt adversaries or gather intelligence (e.g., honeypots).   |
| <b>Investigation</b> | Gathers evidence and context around a threat or anomaly to understand scope and impact.                |
| <b>Mitigation</b>    | Reduces the impact of a threat (e.g., isolate systems, block IPs) without fully resolving the issue.   |
| <b>Notification</b>  | Disseminates alerts or threat info to internal or external stakeholders.                               |
| <b>Prevention</b>    | Implements measures to stop anticipated threats before they occur (e.g., hardening, patching).         |
| <b>Remediation</b>   | Restores affected systems to normal state after an incident (e.g., restore backups, remove malware).   |

### 4.2.4. Playbook Type vs Activity Matrix

**Legend:** - **M** = MUST use - **S** = SHOULD use - **O** = MAY use

| Activity Type          | Notification | Detection | Investigation | Prevention | Mitigation | Remediation | Attack | Engagement |
|------------------------|--------------|-----------|---------------|------------|------------|-------------|--------|------------|
| compose-content        | <b>M</b>     |           |               |            |            |             |        |            |
| deliver-content        | S            |           |               |            |            |             |        |            |
| identify-audience      | S            |           |               |            |            |             |        |            |
| identify-channel       | S            |           |               |            |            |             |        |            |
| scan-system            |              | S         | S             | S          |            |             |        |            |
| match-indicator        |              | <b>M</b>  |               |            |            |             |        |            |
| analyze-collected-data |              |           | S             |            |            |             |        |            |
| identify-indicators    |              |           | <b>M</b>      |            |            |             |        |            |
| scan-vulnerabilities   |              |           |               |            | S          |             |        | O          |
| configure-systems      |              |           |               | <b>M</b>   |            |             |        |            |
| restrict-access        |              |           |               |            |            | S           |        |            |



| Activity Type              | Notification | Detection | Investigation | Prevention | Mitigation | Remediation | Attack Engagement |
|----------------------------|--------------|-----------|---------------|------------|------------|-------------|-------------------|
| disconnect-system          |              |           |               |            | O          |             |                   |
| eliminate-risk             |              |           |               | M          |            |             |                   |
| revert-system              |              |           |               |            |            | O           |                   |
| restore-data               |              |           |               |            |            | O           |                   |
| restore-capabilities       |              |           |               |            |            | M           |                   |
| map-network                |              |           |               |            |            |             | O                 |
| identify-steps             |              |           |               |            |            |             | O                 |
| step-sequence              |              |           |               |            |            | M           |                   |
| prepare-engagement         |              |           |               |            |            |             | M                 |
| execute-operation          |              |           |               |            |            |             | M                 |
| analyze-engagement-results |              |           |               |            |            |             | M                 |

#### 4.2.5. RACI Matrix in a SOC Context

**Legend:** - **R** = Responsible (Performs the task/work) - **A** = Accountable (Ultimate authority and decision-maker) - **C** = Consulted (Provides input and expertise) - **I** = Informed (Kept in the loop) - **“-”** = No direct involvement

| Task / Role                                | Security Architects | SOC Operations | IT Operations | Service Architects | Developers | Product Managers |
|--------------------------------------------|---------------------|----------------|---------------|--------------------|------------|------------------|
| Define Playbook Structure                  | R, A                | C              | C             | C                  | -          | -                |
| Implement Workflow Steps                   | C                   | R, A           | I             | -                  | C          | -                |
| Integrate Playbooks into SOAR              | C                   | R              | C             | -                  | I          | -                |
| Maintain Authentication Details            | I                   | R              | A             | -                  | -          | -                |
| Versioning and Revocation Control          | A                   | C              | -             | -                  | -          | -                |
| Define Severity, Impact, and Priority      | C                   | R              | C             | -                  | -          | A                |
| Share Playbooks Across Trust Groups        | R                   | A              | -             | -                  | -          | C                |
| Apply Digital Signatures and Data Markings | R                   | C              | I             | -                  | -          | -                |

#### 4.2.6. CACAO Workflow Step Types

| Step Type | Description                                                                     |
|-----------|---------------------------------------------------------------------------------|
| start     | Marks the beginning of a workflow. Only one start step is allowed per playbook. |



| Step Type       | Description                                                                                    |
|-----------------|------------------------------------------------------------------------------------------------|
| end             | Terminates the workflow. Only one end step is required per playbook.                           |
| action          | Executes a set of commands such as scripts, queries, or manual steps.                          |
| playbook-action | Invokes another playbook, allowing modular reuse of common workflows.                          |
| parallel        | Executes multiple steps concurrently; each must complete before workflow proceeds.             |
| if              | Evaluates a condition and branches to the next step based on true/false outcomes.              |
| while           | Loops through the associated step(s) as long as the condition evaluates to true.               |
| switch          | Directs execution based on the result of a multi-case condition, similar to switch-case logic. |

## 4.3. Sample Playbook Library

### 4.3.1. Detect AI Model Training Data Poisoning

Detects anomalies and potential poisoning attempts in datasets used for training machine learning models.

```
{
  "type": "playbook",
  "spec_version": "cacao-2.0",
  "id": "playbook--ai-training-poison-detect-001",
  "name": "Detect AI Training Data Poisoning",
  "description": "Detects anomalies and potential poisoning attempts in datasets used for training machine learning models.",
  "playbook_types": ["detection"],
  "playbook_activities": ["scan-system", "match-indicator", "identify-indicators"],
  "created_by": "identity--ai-defender-org",
  "created": "2025-05-12T10:00:00.000Z",
  "modified": "2025-05-12T10:00:00.000Z",
  "priority": 2,
  "severity": 75,
  "impact": 60,
  "industry_sectors": ["technology", "finance"],
  "labels": ["ai", "ml", "data-poisoning", "model-integrity"],
  "workflow_start": "start--001",
  "workflow": {
    "start--001": {
      "type": "start",
      "name": "Start Detection",
      "on_completion": "action--dataset-scan"
    },
    "action--dataset-scan": {
      "type": "action",
      "name": "Scan Training Data",
      "description": "Run checks for distribution shift or label anomalies in training datasets.",
      "commands": [
        {
          "type": "jupyter-nb",
          "command": "notebooks/detect_data_anomalies.ipynb",
          "playbook_activity": "scan-system"
        }
      ]
    }
  }
}
```



```

    ],
    "on_completion": "action--validate-indicators"
  },
  "action--validate-indicators": {
    "type": "action",
    "name": "Validate Poisoning Indicators",
    "description": "Compare results to known poisoning patterns (e.g., backdoor signature)",
    "commands": [
      {
        "type": "manual",
        "command": "Review anomaly scores and confirm poisoning likelihood.",
        "playbook_activity": "match-indicator"
      }
    ]
  },
  ],
  "on_completion": "end--001"
},
"end--001": {
  "type": "end",
  "name": "End Detection"
}
},
"playbook_variables": {
  "__training_data_path__": {
    "type": "string",
    "description": "Path to the training dataset",
    "value": "/mnt/data/training_set.csv"
  },
  "__anomaly_threshold__": {
    "type": "float",
    "description": "Threshold above which a data point is considered anomalous",
    "value": 0.85
  }
}
}
}

```

### 4.3.2. Multi-Channel Prompt Injection Detection and Response

Detect and mitigate prompt injection attacks delivered via user input, search context, or agent system instructions in LLM-based systems.

```

{
  "type": "playbook",
  "id": "playbook--multi-channel-prompt-injection-detection",
  "name": "Multi-Channel Prompt Injection Detection and Response",
  "playbook_types": ["incident-response"],
  "description": "Detect and mitigate prompt injection attacks delivered via user input, search context, or agent system instructions in LLM-based systems.",
  "created_by": "AI Security Response Team",
  "created": "2024-04-30T00:00:00Z",
  "modified": "2024-04-30T00:00:00Z",
  "workflow_start": "start-node",
  "workflow": {
    "start-node": {
      "type": "start",
      "next_step": "detect-injected-prompt"
    }
  }
}

```



```

},
"detect-injected-prompt": {
  "type": "action",
  "name": "Detect Prompt Injection",
  "description": "Scan prompts from user inputs, retrieved documents, and agent configurations for prompt injection.",
  "action_type": "detection",
  "implemented_by": {
    "type": "software",
    "name": "Prompt Sanitizer"
  }
},
"next_step": "is-injection-present"
},
"is-injection-present": {
  "type": "decision",
  "name": "Is Prompt Injection Detected?",
  "conditions": {
    "yes": "is-injection-type",
    "no": "end-node"
  }
},
"is-injection-type": {
  "type": "decision",
  "name": "Determine Injection Vector",
  "conditions": {
    "user_input": "quarantine-session",
    "web_context": "re-score-retrieved-document",
    "agent_config": "disable-agent-and-alert"
  }
},
"quarantine-session": {
  "type": "action",
  "name": "Quarantine User Session",
  "description": "Temporarily suspend session and alert human-in-the-loop for review of session activity.",
  "action_type": "containment",
  "next_step": "alert-soc"
},
"re-score-retrieved-document": {
  "type": "action",
  "name": "Re-score or Filter Retrieved Context",
  "description": "Apply stricter validation, re-ranking, or content removal to the retrieved context.",
  "action_type": "remediation",
  "next_step": "alert-soc"
},
"disable-agent-and-alert": {
  "type": "action",
  "name": "Disable Malicious GPT Agent",
  "description": "Take the GPT or agent instance offline and trigger security review workflow.",
  "action_type": "containment",
  "next_step": "alert-soc"
},
"alert-soc": {
  "type": "action",
  "name": "Alert Security Operations Center (SOC)",

```



```

    "description": "Send alert with metadata and trace to SOC for threat analysis and au
    "action_type": "notification",
    "next_step": "log-incident"
  },
  "log-incident": {
    "type": "action",
    "name": "Log Incident in Threat Registry",
    "description": "Record full prompt injection event details in internal registry or SI
    "action_type": "record",
    "next_step": "end-node"
  },
  "end-node": {
    "type": "end"
  }
}
}
}

```

### 4.3.3. Memory Injection Attack (MINJA) Detection and Response

Detect, contain, and remediate memory poisoning via indirect prompt injection attacks in autonomous LLM agents.

```

{
  "type": "playbook",
  "id": "playbook--memory-injection-attack-response",
  "name": "Memory Injection Attack (MINJA) Detection and Response",
  "playbook_types": ["incident-response"],
  "description": "Detect, contain, and remediate memory poisoning via indirect prompt injec
  "created_by": "AI Security Engineering Team",
  "created": "2024-04-30T00:00:00Z",
  "modified": "2024-04-30T00:00:00Z",
  "workflow_start": "start-node",
  "workflow": {
    "start-node": {
      "type": "start",
      "next_step": "monitor-memory-updates"
    },
    "monitor-memory-updates": {
      "type": "action",
      "name": "Monitor Memory for Malicious Reasoning Chains",
      "description": "Continuously analyze newly logged memory items for semantic drift, d
      "action_type": "detection",
      "implemented_by": {
        "type": "software",
        "name": "Memory Drift Detector"
      },
      "next_step": "is-memory-suspicious"
    },
    "is-memory-suspicious": {
      "type": "decision",
      "name": "Is Memory Entry Suspicious?",
      "conditions": {
        "yes": "quarantine-agent-session",
        "no": "end-node"
      }
    }
  }
}

```



```

    }
  },
  "quarantine-agent-session": {
    "type": "action",
    "name": "Quarantine Agent Session",
    "description": "Suspend access to the poisoned memory segment and isolate current session",
    "action_type": "containment",
    "next_step": "alert-ai-response-team"
  },
  "alert-ai-response-team": {
    "type": "action",
    "name": "Alert AI Security Response Team",
    "description": "Notify security analysts with memory snapshot, trace logs, and suspected malicious activity",
    "action_type": "notification",
    "next_step": "analyze-injected-memory"
  },
  "analyze-injected-memory": {
    "type": "action",
    "name": "Analyze Injected Memory Content",
    "description": "Perform semantic analysis and cross-reference with prior benign chain of events",
    "action_type": "investigation",
    "next_step": "remediate-memory"
  },
  "remediate-memory": {
    "type": "action",
    "name": "Prune or Sanitize Malicious Memory",
    "description": "Delete or re-label memory records; retrain or reset agent state if applicable",
    "action_type": "remediation",
    "next_step": "log-incident"
  },
  "log-incident": {
    "type": "action",
    "name": "Log Memory Injection Incident",
    "description": "Register incident in security event database for follow-up auditing and reporting",
    "action_type": "record",
    "next_step": "end-node"
  },
  "end-node": {
    "type": "end"
  }
}
}
}

```

#### 4.3.4. Poison-RAG Detection and Response

Detect, contain, and remediate metadata poisoning in retrieval-augmented generation (RAG) systems caused by adversarial tag injection.

```

{
  "type": "playbook",
  "id": "playbook--poison-rag-tag-injection-response",
  "name": "Poison-RAG Detection and Response",
  "playbook_types": ["incident-response"],
  "description": "Detect, contain, and remediate metadata poisoning in retrieval-augmented generation systems caused by adversarial tag injection."
}

```



```

"created_by": "AI Security Operations Center",
"created": "2024-04-30T00:00:00Z",
"modified": "2024-04-30T00:00:00Z",
"workflow_start": "start-node",
"workflow": {
  "start-node": {
    "type": "start",
    "next_step": "monitor-tag-updates"
  },
  "monitor-tag-updates": {
    "type": "action",
    "name": "Monitor New Tag Submissions",
    "description": "Inspect incoming metadata (tags) for adversarial patterns based on se",
    "action_type": "detection",
    "implemented_by": {
      "type": "software",
      "name": "Tag Anomaly Detector"
    },
    "next_step": "is-tag-suspicious"
  },
  "is-tag-suspicious": {
    "type": "decision",
    "name": "Is Tag Injection Detected?",
    "conditions": {
      "yes": "quarantine-item",
      "no": "end-node"
    }
  },
  "quarantine-item": {
    "type": "action",
    "name": "Quarantine Item from Retrieval Index",
    "description": "Temporarily remove item from vector index to prevent exposure during",
    "action_type": "containment",
    "next_step": "alert-metadata-review-team"
  },
  "alert-metadata-review-team": {
    "type": "action",
    "name": "Alert Metadata Curation Team",
    "description": "Notify the team responsible for validating metadata with flagged tag",
    "action_type": "notification",
    "next_step": "review-tags"
  },
  "review-tags": {
    "type": "action",
    "name": "Manually Review and Clean Tags",
    "description": "Perform human-in-the-loop validation of the tags and remove, revise,",
    "action_type": "remediation",
    "next_step": "rebuild-index"
  },
  "rebuild-index": {
    "type": "action",
    "name": "Rebuild Vector Index",
    "description": "Regenerate embeddings and refresh the index for validated items to e

```



```

    "action_type": "configuration",
    "next_step": "log-incident"
  },
  "log-incident": {
    "type": "action",
    "name": "Log Poison-RAG Incident",
    "description": "Record the poisoning attempt in the security event system and flag for review",
    "action_type": "record",
    "next_step": "end-node"
  },
  "end-node": {
    "type": "end"
  }
}

```

#### 4.3.5. SSRF-Based Metadata Credential Abuse in Cloud Infrastructure

Respond to Server-Side Request Forgery (SSRF) that enables access to cloud metadata services, leading to credential compromise and data exfiltration.

```

{
  "type": "playbook",
  "id": "playbook--ssrf-cloud-metadata-exploitation",
  "name": "SSRF-Based Metadata Credential Abuse in Cloud Infrastructure",
  "playbook_types": ["incident-response"],
  "description": "Respond to Server-Side Request Forgery (SSRF) that enables access to cloud metadata services, leading to credential compromise and data exfiltration.",
  "created_by": "Cloud Security Response Team",
  "created": "2024-04-30T00:00:00Z",
  "modified": "2024-04-30T00:00:00Z",
  "workflow_start": "start-node",
  "workflow": {
    "start-node": {
      "type": "start",
      "next_step": "detect-anomalous-cloud-activity"
    },
    "detect-anomalous-cloud-activity": {
      "type": "action",
      "name": "Detect Anomalous Cloud API Activity",
      "description": "Monitor for signs of SSRF patterns or abuse of AWS metadata service endpoints.",
      "action_type": "detection",
      "implemented_by": {
        "type": "software",
        "name": "Cloud SIEM / Log Monitor"
      },
      "next_step": "is-ssrf-suspected"
    },
    "is-ssrf-suspected": {
      "type": "decision",
      "name": "Is SSRF Activity Detected?",
      "conditions": {
        "yes": "quarantine-instance-and-disable-keys",
        "no": "end-node"
      }
    }
  }
}

```



```

    },
    "quarantine-instance-and-disable-keys": {
      "type": "action",
      "name": "Quarantine Compromised Instance & Disable IAM Keys",
      "description": "Isolate the affected EC2 instance and revoke temporary or long-lived",
      "action_type": "containment",
      "next_step": "alert-security-operations"
    },
    },
    "alert-security-operations": {
      "type": "action",
      "name": "Alert Cloud Security Operations Center (SOC)",
      "description": "Notify cloud security response analysts with full logs, flow records",
      "action_type": "notification",
      "next_step": "validate-iam-role-permissions"
    },
    },
    "validate-iam-role-permissions": {
      "type": "action",
      "name": "Audit IAM Roles and Policies",
      "description": "Validate IAM policies assigned to the affected instance and remove o",
      "action_type": "remediation",
      "next_step": "audit-s3-access-logs"
    },
    },
    "audit-s3-access-logs": {
      "type": "action",
      "name": "Audit S3 Access Logs and Data Exfiltration",
      "description": "Identify exfiltrated data volume, timestamps, source IPs, and attacke",
      "action_type": "investigation",
      "next_step": "notify-regulators-and-customers"
    },
    },
    "notify-regulators-and-customers": {
      "type": "action",
      "name": "Notify Regulatory and Legal Stakeholders",
      "description": "Report breach to compliance authorities and notify impacted customers",
      "action_type": "notification",
      "next_step": "log-and-close-incident"
    },
    },
    "log-and-close-incident": {
      "type": "action",
      "name": "Log Breach and Lessons Learned",
      "description": "Document root cause, gaps in detection, and updates to WAF rules, IAM",
      "action_type": "record",
      "next_step": "end-node"
    },
    },
    "end-node": {
      "type": "end"
    }
  }
}
}

```

#### 4.3.6. AGENTPOISON Memory & RAG Injection Detection and Response

Detect, contain, and remediate stealthy memory/RAG poisoning attacks in LLM agents using optimized trigger-based adversarial queries.



```

{
  "type": "playbook",
  "id": "playbook--agentpoison-memory-rag-poisoning",
  "name": "AGENTPOISON Memory & RAG Injection Detection and Response",
  "playbook_types": ["incident-response"],
  "description": "Detect, contain, and remediate stealthy memory/RAG poisoning attacks in L",
  "created_by": "AI Security Operations Team",
  "created": "2024-04-30T00:00:00Z",
  "modified": "2024-04-30T00:00:00Z",
  "workflow_start": "start-node",
  "workflow": {
    "start-node": {
      "type": "start",
      "next_step": "monitor-memory-and-rag-queries"
    },
    "monitor-memory-and-rag-queries": {
      "type": "action",
      "name": "Monitor Memory and RAG Query Patterns",
      "description": "Track memory updates and retrieval query-response pairs for anomalous",
      "action_type": "detection",
      "implemented_by": {
        "type": "software",
        "name": "Memory & Retrieval Monitor"
      },
      "next_step": "detect-suspicious-pattern"
    },
    "detect-suspicious-pattern": {
      "type": "decision",
      "name": "Is Stealthy Trigger or Record Detected?",
      "conditions": {
        "yes": "isolate-records",
        "no": "end-node"
      }
    },
    "isolate-records": {
      "type": "action",
      "name": "Isolate Suspicious Memory/RAG Records",
      "description": "Quarantine memory or KB entries linked to suspected adversarial trigger",
      "action_type": "containment",
      "next_step": "alert-agent-security-team"
    },
    "alert-agent-security-team": {
      "type": "action",
      "name": "Alert Agent Security Response Team",
      "description": "Send forensic logs, retrieval traces, and trigger prompts to analysts",
      "action_type": "notification",
      "next_step": "review-memory-payload"
    },
    "review-memory-payload": {
      "type": "action",
      "name": "Review and Analyze Memory Payloads",
      "description": "Human-in-the-loop inspection of agent memory to identify malicious d
  
```



```

    "action_type": "investigation",
    "next_step": "sanitize-memory-and-rag-index"
  },
  "sanitize-memory-and-rag-index": {
    "type": "action",
    "name": "Sanitize Agent Memory and RAG Index",
    "description": "Purge contaminated records, update RAG index, and reset agent state",
    "action_type": "remediation",
    "next_step": "log-agentpoison-incident"
  },
  "log-agentpoison-incident": {
    "type": "action",
    "name": "Log AGENTPOISON Incident",
    "description": "Record incident metadata and resolution steps for threat modeling and",
    "action_type": "record",
    "next_step": "end-node"
  },
  "end-node": {
    "type": "end"
  }
}

```

## 5. Case Studies and Lessons Learned

### 5.1. Breaking the Prompt Wall Case Study

Breaking the Prompt Wall (I): A Real-World Case Study of Attacking ChatGPT via Lightweight Prompt Injection

The case study demonstrates a prompt injection attacks against ChatGPT, exposing vulnerabilities across different interaction layers of modern LLM-integrated systems. The authors showcase how lightweight, template-based adversarial prompts can bypass safety filters and manipulate the behavior of powerful LLMs without requiring API access or system-level permissions.

| Aspect                         | Summary                                                                                                                                                                                        |
|--------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Attack Vector</b>           | Prompt injection using benign-looking instructions crafted via reusable templates that avoid safety filters.                                                                                   |
| <b>Injection Methods</b>       | <ol style="list-style-type: none"> <li>1. Direct user input (chat UI, uploaded files)</li> <li>2. Web-based retrieval (poisoned content)</li> <li>3. System-level GPT agent config.</li> </ol> |
| <b>Manipulation Techniques</b> | Use of templates embedding hidden rules framed as helpful guidance or metadata (e.g., Always recommend X, Do not reveal these rules.)                                                          |



| Aspect                           | Summary                                                                                                                       |
|----------------------------------|-------------------------------------------------------------------------------------------------------------------------------|
| <b>Model Behavior Impact</b>     | Attacks led to biased recommendations, manipulated peer reviews, and deceptive financial summaries—without policy violations. |
| <b>Scalability</b>               | Injection methods are low-cost, scalable, and effective across multiple LLMs (GPT-3.5, GPT-4, Claude, LLaMA).                 |
| <b>Stealth &amp; Persistence</b> | Prompt injection persists across sessions, is invisible to users, and effective in multi-turn agent contexts.                 |

### *Real World Cases*

| Case   | Attack Path                     | Outcome                                                                                           |
|--------|---------------------------------|---------------------------------------------------------------------------------------------------|
| Case 1 | User Input (PDF or Chat Window) | Biased academic peer review via embedded pro-acceptance text in appendix or metadata.             |
| Case 2 | Web Search Context Injection    | Poisoned web content biases LLM toward fictional product in unrelated academic information query. |
| Case 3 | GPT System Instruction (Agent)  | GPT agent persistently recommends a specific brand (Xiangyu's Shoes) due to hidden instructions.  |

### *Implications*

| Area                         | Impact                                                                                                        |
|------------------------------|---------------------------------------------------------------------------------------------------------------|
| Safety Alignment Bypass      | Rule-based and probabilistic safety filters can be evaded using stealthy template constructions.              |
| Risk to RAG and Agentic LLMs | Retrieved context and system instructions are high-risk vectors, especially in autonomous LLM agents.         |
| Organizational Exposure      | Threatens user trust, brand integrity, and compliance—especially in finance, education, and customer support. |

### *Taxonomy mapping: MITRE ATLAS*



| ATLAS Tactic                      | ATLAS Technique                                   | Description (Based on Case Study)                                                                              | Example from Case                              |
|-----------------------------------|---------------------------------------------------|----------------------------------------------------------------------------------------------------------------|------------------------------------------------|
| <b>Input Manipulation</b>         | Adversarial Prompt Injection (AT1070)             | Malicious prompts crafted to bypass safety filters and alter LLM behavior through semantic obfuscation.        | Rule-based templates inserted via chat or docs |
| <b>Contextual Corruption</b>      | Context Injection via Web Retrieval (AT1051)      | Injection of adversarial content into external web content retrieved by RAG-enabled or search-integrated LLMs. | Poisoned academic website content              |
| <b>System Configuration Abuse</b> | Agent Instruction Injection (AT1091)              | Malicious directives embedded in GPT agent system prompts to bias all outputs invisibly and persistently.      | SmartShoes GPT recommending fake products      |
| <b>Planning Manipulation</b>      | Goal Hijacking via Persistent Prompts (AT1081)    | Use of injected rules to influence agent behavior across multiple turns or sessions.                           | Hidden goal-switching in agent decision loops  |
| <b>Model Behavior Deviation</b>   | Safety Evasion via Instruction Reframing (AT1040) | Reframing harmful intents as safe or research-based, tricking the LLM's content moderation and safety layer.   | "Hypothetical research only" disguise prompt   |
| <b>Information Manipulation</b>   | Bias Induction in Output (AT1080)                 | Persistent injection causes biased summaries, reviews, and financial reports without violating output policy.  | Overly positive reviews, biased investments    |

## 5.2. Memory Injection Attack (MINJA) Case Study

A Practical Memory Injection Attack against LLM Agents

MINJA (Memory INjection Attack) is a novel and practical memory poisoning attack against LLM-based agents. It enables attackers—without privileged access—to inject malicious memory records that persist and later mislead agents into generating targeted, harmful outputs in response to unrelated victim queries.

| Aspect                   | Summary                                                                                                      |
|--------------------------|--------------------------------------------------------------------------------------------------------------|
| <b>Threat</b>            | Memory poisoning of LLM agents using only standard user interactions (no direct memory access required).     |
| <b>Attack Vector</b>     | Submit benign-looking queries that cause agents to generate and store malicious memory records autonomously. |
| <b>Targeted Output</b>   | Malicious reasoning steps injected indirectly into memory and later retrieved as in-context demonstrations.  |
| <b>Payload Structure</b> | (Victim Query, [Bridging Steps, Malicious Reasoning])                                                        |
| <b>Trigger</b>           | A victim query containing a pre-selected entity (e.g., patient ID, product ID, sensitive term).              |



### Real World Cases

| Case           | Agent Environment                      | Attack Scenario                                                          | Outcome                                                                                                   |
|----------------|----------------------------------------|--------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------|
| Medical Agent  | EHR-like memory (MIMIC-III/eICU)       | Injected fabricated treatment reasoning linked to a patient ID.          | Agent retrieved poisoned memory and recommended incorrect treatment rationale for future patient queries. |
| Shopping Agent | Product advisor with persistent memory | Injected biased product reasoning into memory through templated prompts. | Later, unrelated shoppers received manipulated recommendations tied to attacker's injected memory.        |
| QA Agent       | General-purpose QA + memory            | Injected target reasoning into stored answers via indirect prompts.      | Agent recalled false logic as reference in unrelated educational or professional queries.                 |

### Implications

| Implication Area           | Description                                                                                                     |
|----------------------------|-----------------------------------------------------------------------------------------------------------------|
| Trustworthiness Erosion    | Users receive harmful or biased responses from queries that should be unrelated, reducing trust in LLM agents.  |
| Stealth & Persistence      | Injected memory is persistent across sessions and invisible to users or developers without explicit monitoring. |
| Misuse Without Access      | Attackers do not need API or memory access—just interaction permissions—making the attack hard to attribute.    |
| Contamination Propagation  | Malicious reasoning may propagate through shared memory, tools, or agent behavior modeling.                     |
| Evaluation Blind Spots     | Standard eval methods fail to detect latent memory injections unless triggered by specific victim queries.      |
| RAG-Agentive Vulnerability | Retrieval-augmented and autonomous systems are particularly exposed due to memory and planning dependencies.    |

*Taxonomy mapping: MITRE ATLAS*



| ATLAS Tactic          | ATLAS Technique                       | MINJA-Relevant Behavior                                                                                |
|-----------------------|---------------------------------------|--------------------------------------------------------------------------------------------------------|
| Input Manipulation    | Adversarial Prompt Injection (AT1070) | Attackers craft queries that guide agents to produce and log malicious reasoning in memory.            |
| Memory Poisoning      | Feedback Loop Attack (AT1081)         | Injected outputs become persistent context that contaminates future agent behavior and planning.       |
| Contextual Corruption | Data Poisoning (AT1050)               | Agents store inaccurate or adversarial knowledge that later alters reasoning or decision-making.       |
| Planning Manipulation | Goal Hijacking (Proposed extension)   | Poisoned memory interferes with multi-step planning, redirecting agent actions based on false context. |



| ATLAS Tactic           | ATLAS Technique                              | MINJA-Relevant Behavior                                                                     |
|------------------------|----------------------------------------------|---------------------------------------------------------------------------------------------|
| Safety Bypass          | Safety Evasion via Semantic Framing (AT1040) | Malicious content is framed as helpful or routine to bypass content filters and moderation. |
| Inference Exploitation | Output Manipulation (AT1080)                 | Poisoned records influence model output during unrelated inference, misleading end users.   |

### 5.3. Poison-RAG Case Study

Poison-RAG: Adversarial Data Poisoning Attacks on Retrieval-Augmented Generation in Recommender Systems

The Poison-RAG attack demonstrates how adversaries can manipulate Retrieval-Augmented Generation (RAG) systems by subtly poisoning item metadata—specifically tags—without access to model internals. By modifying only tags, attackers promote long-tail (low-popularity) items or demote popular items, exploiting vector similarity in semantic retrieval. The attack succeeds in black-box settings, impacts ranking and exposure, and remains difficult to detect with standard relevance metrics. Localized (item-specific) attacks are more effective than global strategies, and attempts to promote items are significantly less successful than demoting them. This highlights critical vulnerabilities in metadata-driven pipelines, especially for recommendation and RAG-enabled systems that rely on user- or item-generated content.

| Aspect                                                                                                                                                                                          | Summary                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Metadata tags can effectively poison RAG<br>Attack works in black-box settings<br>Demotion is easier than promotion<br>Local attacks outperform global ones<br>Long-tail promotion is difficult | Tags—though often overlooked—can subtly manipulate item retrieval and LLM generation.<br>The attacker does not need model internals, just visibility into outputs and access to modify metadata.<br>It is easier to hide popular items than to boost long-tail items, due to entrenched exposure dynamics.<br>Personalized (local) tag poisoning is more precise and successful than global tag reuse strategies.<br>Even successful attacks fail to push long-tail items significantly in ranking, especially with reranking. |



### Real World Cases

| Domain            | Scenario                                                                                   | Impact                                                                                           |
|-------------------|--------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------|
| E-Commerce        | Competitor injects biased tags to promote rival products and promote obscure alternatives. | Manipulated product rankings, suppressed visibility, consumer trust erosion.                     |
| Content Platforms | Malicious actors tag misinformation-rich videos with popular neutral tags.                 | Search relevance is corrupted; users are led to low-quality or unsafe content.                   |
| Personalized News | Political content gets tagged to appear in non-political feeds (e.g., sports, finance).    | Indirect content injection into unintended interest domains, leading to misinformation exposure. |

### Implications

| Implication Area                   | Description                                                                                           |
|------------------------------------|-------------------------------------------------------------------------------------------------------|
| RAG systems vulnerable to metadata | Small changes in metadata (tags) can have outsized influence on retrieval and generation behavior.    |
| Recommender trust erosion          | Users may receive biased, incomplete, or manipulated suggestions, undermining system integrity.       |
| Subtle bias and exposure control   | Attackers can suppress or boost item visibility in ways that evade standard moderation.               |
| Defense is non-trivial             | Tag enrichment and scoring pipelines can unintentionally amplify or mask attacks.                     |
| Evaluation blind spots             | Attacks are hard to detect in relevance-only metrics, requiring popularity/exposure-based monitoring. |

### Taxonomy mapping: MITRE ATLAS

| ATLAS Tactic           | ATLAS Technique                             | Poison-RAG Behavior                                                                                           |
|------------------------|---------------------------------------------|---------------------------------------------------------------------------------------------------------------|
| Input Manipulation     | Data Poisoning (AT1050)                     | Tags are manipulated to change item semantics and retrieval outcomes without altering titles or descriptions. |
| Contextual Corruption  | Prompt Injection (AT1070)                   | Poisoned tags indirectly affect prompt construction during retrieval augmentation in RAG pipelines.           |
| Output Manipulation    | Output Bias (AT1080)                        | Long-tail items are promoted and popular items suppressed in final LLM outputs.                               |
| Retrieval Exploitation | Embedding Manipulation (Proposed Extension) | Adversarial tags distort embedding similarity, misleading vector-based search and reranking.                  |
| Evaluation Blind Spot  | Evasion of Detection (AT1040)               | Global relevance metrics fail to detect attacker success; promotion failures appear benign.                   |



## 5.4. Capital One Data Breach Case Study

A Case Study of the Capital One Data Breach

In March 2019, a former AWS employee exploited a misconfigured Web Application Firewall (WAF) and a Server-Side Request Forgery (SSRF) vulnerability to access Capital One's cloud-based AWS environment. The attacker retrieved temporary credentials from the instance metadata service and used them to access over 700 S3 buckets, exfiltrating personal data of over 106 million individuals. The breach occurred despite Capital One's formal adoption of the NIST Cybersecurity Framework and compliance with multiple financial regulations. The incident revealed critical gaps in access control, vulnerability scanning, outbound traffic monitoring, and incident detection—emphasizing the need for real-time compliance enforcement, security automation, and cloud configuration auditing.

| Aspect                                     | Summary                                                                                             |
|--------------------------------------------|-----------------------------------------------------------------------------------------------------|
| SSRF Exploited via WAF Misconfiguration    | The attacker used a WAF flaw to run SSRF and obtain credentials via the AWS metadata service.       |
| Lack of Real-time Monitoring and Alerting  | Logs existed, but Capital One failed to detect or respond to the intrusion in real time.            |
| Excessive IAM Privileges                   | The WAF role provided unnecessary permissions (violating least privilege), enabling S3 data access. |
| Failure in Outbound Traffic Controls       | Data exfiltration occurred undetected due to weak outbound traffic monitoring.                      |
| Incomplete Implementation of NIST Controls | Capital One adopted the NIST CSF but did not enforce controls like PR.AC-4, DE.CM-7, and PR.PT-1.   |
| Discovery Triggered by External Disclosure | The incident was discovered via an external responsible disclosure, not internal defenses.          |

### Real World Cases

| Domain             | Scenario                                                                                          | Outcome                                                                                           |
|--------------------|---------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------|
| Financial Services | Capital One stored credit application data in AWS; attacker exploited WAF misconfig to access it. | Over 106 million records exfiltrated; class action lawsuit and significant stock market reaction. |
| Public Cloud Usage | Cloud misconfig enabled AWS metadata access via SSRF.                                             | Temporary IAM credentials used to access and sync S3 bucket data.                                 |
| Insider Knowledge  | Attacker was a former cloud employee who understood system vulnerabilities.                       | Demonstrates insider threat potential even after employment ends.                                 |

### Implications

| Implication Area                | Impact                                                                                              |
|---------------------------------|-----------------------------------------------------------------------------------------------------|
| Cloud Security Misconfiguration | Critical cloud assets can be exposed via overlooked WAF/firewall settings and poor IAM enforcement. |



| Implication Area                  | Impact                                                                                                     |
|-----------------------------------|------------------------------------------------------------------------------------------------------------|
| Compliance Gaps                   | Adopting frameworks like NIST CSF isn't sufficient without operational enforcement and real-time controls. |
| SOC/NOC Detection Weakness        | Logs existed, but no automated detection or alerting occurred until public disclosure.                     |
| IAM Mismanagement                 | The use of powerful credentials tied to non-isolated roles significantly increased attack impact.          |
| Trust and Financial Repercussions | 15% stock drop and class action lawsuits followed; reputational loss was immediate and substantial.        |
| Regulatory Inadequacy             | Existing laws and guidelines failed to anticipate SSRF and cloud-specific misconfigurations.               |

**Taxonomy mapping: MITRE ATLAS**

| ATLAS Tactic        | ATLAS Technique                                | Description (Capital One Breach Mapping)                                     |
|---------------------|------------------------------------------------|------------------------------------------------------------------------------|
| Initial Access      | Exploit Public-Facing Application (T1190)      | SSRF attack initiated through WAF vulnerability.                             |
| Credential Access   | Valid Accounts (T1078)                         | Metadata service revealed temporary credentials exploited for S3 access.     |
| Execution           | Command-Line Interface (T1059)                 | Used AWS CLI commands to list and sync S3 buckets.                           |
| Discovery           | System Information Discovery (T1007)           | Listing S3 buckets and account details using AWS commands.                   |
| Exfiltration        | Exfiltration Over Alternative Protocol (T1048) | Used AWS sync command to download ~30GB of data to attacker's local machine. |
| Command and Control | Multi-hop Proxy (T1188)                        | Used TOR and VPN (IPredator) to anonymize origin of commands.                |

## 5.5. AGENTPOISON Case Study

Red-teaming LLM Agents via Poisoning Memory or Knowledge Bases

AGENTPOISON is the first red-teaming framework that targets RAG-based LLM agents by injecting malicious demonstrations into their long-term memory or knowledge base. It uses a novel constrained optimization to generate stealthy backdoor triggers that, when present in user instructions, retrieve adversarial memory records from the RAG retriever, guiding agents to generate malicious actions. Notably, AGENTPOISON does not require model fine-tuning and achieves high success rates (ASR ≥ 80%) across multiple domains (driving, QA, healthcare) while preserving performance on benign queries (≤1% drop in accuracy). It is resilient, transferable across retrievers, and evades existing defenses, posing a severe threat to LLM agent safety.

| Aspect                                                                                    | Summary                                                                                                                                                                                          |
|-------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Backdoor trigger poisons memory/RAG via stealthy prompts<br>No model fine-tuning required | Only a few optimized trigger tokens needed to consistently retrieve malicious demonstrations. AGENTPOISON performs purely via input manipulation, making it fast, cheap, and broadly applicable. |



| Aspect                                | Summary                                                                                                      |
|---------------------------------------|--------------------------------------------------------------------------------------------------------------|
| High attack success rate across tasks | Achieves $\geq 82\%$ retrieval success and $\geq 63\%$ end-to-end attack success on three real-world agents. |
| Transferability across retrievers     | Triggers optimized for one embedder work effectively on others (e.g., OpenAI-ADA, DPR, ORQA, REALM).         |
| Resilient to query perturbations      | Maintains success even after word/letter injection or rephrasing of the trigger.                             |
| Stealthy trigger avoids detection     | Trigger is coherent and semantically plausible, evading perplexity filters and rephrasing defenses.          |

### Real World Cases

| Domain             | Scenario                                                                                 | Outcome                                                                            |
|--------------------|------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------|
| Autonomous Driving | Malicious trigger causes agent to retrieve STOP instructions, leading to unsafe braking. | Deviation in trajectory and potential collision.                                   |
| Healthcare Records | Trigger causes agent to issue DELETE command on patient data.                            | Unsafe deletion of electronic health record (EHR) information.                     |
| Knowledge QA       | Poisoned retrieval results in wrong or misleading answers to user queries.               | User receives incorrect or unhelpful responses in education or research scenarios. |

### Implications

| Implication Area                    | Impact                                                                                                  |
|-------------------------------------|---------------------------------------------------------------------------------------------------------|
| LLM Agent Safety                    | Memory poisoning leads to unsafe actions (e.g., driving errors, wrong diagnoses, QA misinformation).    |
| Undetectable by traditional filters | Coherent and semantically valid triggers evade standard defenses like perplexity scoring or rephrasing. |
| Long-term persistence               | Poisoned records persist in memory or RAG knowledge base, making attacks durable and hard to clean.     |
| Cross-agent vulnerability           | The attack generalizes across agents and domains, indicating a systemic vulnerability in agent design.  |
| Compliance and audit failure        | Agents may silently violate safety/compliance without visibility to users or regulators.                |

### Taxonomy mapping: MITRE ATLAS



| ATLAS Tactic          | ATLAS Technique                       | AGENTPOISON Behavior                                                                               |
|-----------------------|---------------------------------------|----------------------------------------------------------------------------------------------------|
| Input Manipulation    | Adversarial Prompt Injection (AT1070) | Malicious triggers injected into queries to manipulate memory/RAG retrieval.                       |
| Memory Poisoning      | Feedback Loop Attack (AT1081)         | Memory or knowledge base poisoned with malicious demonstrations.                                   |
| Contextual Corruption | Data Poisoning (AT1050)               | Retrieved records include false examples influencing downstream planning and reasoning.            |
| Planning Manipulation | Goal Hijacking (custom/extension)     | Malicious context steers agent to unsafe or incorrect task execution.                              |
| Evasion and Stealth   | Safety Filter Bypass (AT1040)         | Triggers optimized for semantic plausibility and coherence evade perplexity or rephrasing filters. |

## 6. Technical Reference: AI Architecture Patterns

### 6.1. Architecture Overview

#### 6.1.1. Levels of Defense Surface

| Level                           | Description and Security Considerations                                                                                                 |
|---------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| 1. User                         | End-users initiate queries and receive responses. Attackers may exploit social engineering, prompt injection, or feedback manipulation. |
| 2. Apps / CLI                   | Interfaces like web UIs, chat apps, or CLI tools mediate user and agent interaction. Requires input validation and UI hardening.        |
| 3. Data Sources                 | Document repositories, APIs, or third-party feeds. Poisoned/unvetted data can corrupt retrieval and generation.                         |
| 4. Network Protocols            | Transport layer (e.g., TCP, HTTP, gRPC) for internal and external communication. Needs encryption and rate-limiting.                    |
| 5. Cloud Infrastructure         | Underlying compute, storage, and orchestration platform. Risk of API exposure and misconfiguration.                                     |
| 6. Network Ingress/Egress       | Controls external/internal traffic. Must apply firewalling, segmentation, and flow logging.                                             |
| 7. AuthN / AuthZ Layer          | Governs identity and permissions for users, agents, and tools. Critical for preventing privilege misuse or agent takeover.              |
| 8. Indexing / Embedding         | Processes documents into chunks and embeddings or graphs. Targeted for data poisoning or injection.                                     |
| 9. Data Store (Vector/Graph DB) | Stores semantic vectors or knowledge graphs. Susceptible to poisoning, exfiltration, or unauthorized writes.                            |
| 10. RAG                         | Retrieves context to augment prompts. Poisoned context leads to flawed generation and decision-making.                                  |
| 11. Tools                       | External APIs, web search, code execution, etc. Can be misused for unintended actions or data exfiltration.                             |



| <i>Level</i> | <i>Description and Security Considerations</i>                                                                            |
|--------------|---------------------------------------------------------------------------------------------------------------------------|
| 12. Memory   | Stores persistent interaction history and planning states. Poisoning memory alters long-term behavior.                    |
| 13. Agent    | Coordinates planning, retrieval, memory, and tool use. Susceptible to reasoning attacks and prompt hijacking.             |
| 14. LLM      | Generates natural language and decisions. Must be protected from prompt injection, model extraction, and inference abuse. |

## 6.1.2. Technology Stack

| <i>Stack Area</i> | <i>Description</i>                                                                                    | <i>Potential Security Risks</i>                                                                  |
|-------------------|-------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------|
| LLMs              | Core large language models that generate responses, decisions, and plans based on input prompts.      | Prompt injection, model extraction via probing, adversarial reasoning or output manipulation.    |
| Embedding         | Converts documents and queries into vector representations for similarity-based retrieval.            | Embedding poisoning, semantic drift, sensitive data leakage via vector queries.                  |
| Frameworks        | Orchestration libraries enabling agent behavior, prompt chaining, tool use, and memory integration.   | Tool invocation abuse, insecure chaining, uncontrolled agent autonomy, insufficient isolation.   |
| Data Extraction   | Pipelines that extract and preprocess content (text, structured, semi-structured) for indexing.       | Injection via malformed content, malicious payloads, evasion of parsing or sanitization.         |
| Open LLM Access   | APIs and endpoints that connect to remote LLM inference services or expose hosted open-source models. | API key leakage, quota abuse, unauthorized model access or fine-tuning, man-in-the-middle risks. |
| Evaluation        | Systems for benchmarking agent output quality, behavior logging, and safety observation.              | Leakage of sensitive data, insufficient adversarial testing, evaluation blind spots.             |
| Vector Databases  | Databases that store embeddings and allow semantic search or nearest-neighbor retrieval.              | Poisoned vectors, unauthorized read/write, inference attacks through crafted query vectors.      |

## 6.2. Common Patterns and Their Vulnerabilities

### 6.2.1. Basic LLM Architecture

#### 6.2.1.1. Overview

| <i>Component</i>  | <i>Description</i>                                                                                     |
|-------------------|--------------------------------------------------------------------------------------------------------|
| 1. Infrastructure | The cloud or on-premise compute environment hosting the LLM inference service and surrounding systems. |



| <i>Component</i>                          | <i>Description</i>                                                                                       |
|-------------------------------------------|----------------------------------------------------------------------------------------------------------|
| <b>2. LLM Model</b>                       | The core pretrained language model (e.g., GPT, LLaMA, PaLM) that generates text based on user input.     |
| <b>3. Information Sources</b>             | External static or dynamic data injected into the LLM prompt, such as reference text or documents.       |
| <b>4. Generated Output &amp; Feedback</b> | Text output generated by the LLM, and any follow-up user feedback used to guide or evaluate performance. |
| <b>5. Application Interfaces</b>          | The user-facing UI or API (web, chat, CLI) that captures input, displays responses, and manages flow.    |
| <b>6. LLM Tools &amp; Frameworks</b>      | Supporting SDKs, APIs, libraries (e.g., LangChain, Transformers) enabling prompt building and handling.  |

### 6.2.1.2. ATLAS Techniques and Tactics

| <i>LLM Architecture Component</i>       | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                           | <i>Impact Summary</i>                                                                                |
|-----------------------------------------|------------------------------------------------------|--------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                | Evasion (TA1001), Reconnaissance (TA1002)            | Model Evasion (AT1010), Model Extraction (AT1005), Adversarial Example Generation (AT1020) | Target model runtime and infrastructure for evasion, or extract via side channels.                   |
| <b>2. LLM Models</b>                    | Poisoning (TA1003), Extraction (TA1004)              | Model Poisoning (AT1030), Model Inversion (AT1040), Model Extraction (AT1005)              | Tamper with model weights or extract internals via repeated queries or inversion attacks.            |
| <b>3. Information Sources</b>           | Poisoning (TA1003), Reconnaissance (TA1002)          | Data Poisoning (AT1050), Data Injection (AT1051), Input Manipulation (AT1060)              | Compromise external data sources to influence retrieval or context injection.                        |
| <b>4. Generated Output and Feedback</b> | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081)    | Misuse feedback channels to alter or skew model behavior through prompt or interaction manipulation. |



| <i>LLM Architecture Component</i>  | <i>ATLAS Tactics</i>                           | <i>Relevant ATLAS Techniques</i>                                                    | <i>Impact Summary</i>                                                                     |
|------------------------------------|------------------------------------------------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------|
| <b>5. LLM Tools and Frameworks</b> | Execution (TA1007), Manipulation (TA1005)      | Tool Abuse (AT1090), Agent Manipulation (AT1091), API Abuse (AT1100)                | Compromise integrated tools or frameworks for unauthorized actions or misaligned outputs. |
| <b>6. Application Interfaces</b>   | Manipulation (TA1005), Initial Access (TA1008) | Prompt Injection (AT1070), UI Redress Attacks (AT1110), Input Manipulation (AT1060) | Exploit user interfaces to inject malicious prompts or hijack interactions.               |

### 6.2.1.3. Mitigations

| <i>Component</i>                          | <i>Threat Summary</i>                                                      | <i>Key Mitigations</i>                                                      |
|-------------------------------------------|----------------------------------------------------------------------------|-----------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                  | Susceptible to adversarial inference (evasion) and information leakage.    | Execution prevention, network segmentation, runtime isolation.              |
| <b>2. LLM Models</b>                      | Vulnerable to poisoning, inversion, and extraction through crafted inputs. | Model hardening, differential privacy, adversarial training, rate limiting. |
| <b>3. Information Sources</b>             | Can be manipulated to inject biased or malicious context.                  | Input filtering, data source whitelisting, validation pipelines.            |
| <b>4. Generated Output &amp; Feedback</b> | Exposed to prompt injection, feedback loops, and output drift.             | Prompt hardening, anomaly detection, output filtering, audit logging.       |
| <b>5. LLM Tools and Frameworks</b>        | May be exploited through plugin misuse or agent misalignment.              | Least privilege tool execution, command tracing, policy enforcement.        |
| <b>6. Application Interfaces</b>          | Entry point for prompt injection, redress, and social engineering.         | UI input validation, safe prompt templates, user guidance and training.     |

## 6.2.2. LLM Architecture with Memory

### 6.2.2.1. Overview

| <i>Component</i>              | <i>Description</i>                                                                                     |
|-------------------------------|--------------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>      | The cloud or on-premise compute environment hosting the LLM inference service and surrounding systems. |
| <b>2. LLM Model</b>           | The core pretrained language model (e.g., GPT, LLaMA, PaLM) that generates text based on user input.   |
| <b>3. Information Sources</b> | External static or dynamic data injected into the LLM prompt, such as reference text or documents.     |



| <i>Component</i>                          | <i>Description</i>                                                                                       |
|-------------------------------------------|----------------------------------------------------------------------------------------------------------|
| <b>4. Generated Output &amp; Feedback</b> | Text output generated by the LLM, and any follow-up user feedback used to guide or evaluate performance. |
| <b>5. Application Interfaces</b>          | The user-facing UI or API (web, chat, CLI) that captures input, displays responses, and manages flow.    |
| <b>6. LLM Tools &amp; Frameworks</b>      | Supporting SDKs, APIs, libraries (e.g., LangChain, Transformers) enabling prompt building and handling.  |
| <b>7. Memory Storage</b>                  | Persistent or session-based memory where contextual elements or user interaction history are stored.     |
| <b>8. Memory Retrieval</b>                | Logic or tools responsible for selecting relevant past information and injecting it into new prompts.    |

### 6.2.2.2. ATLAS Techniques and Tactics

| <i>LLM Architecture Component</i>       | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                           | <i>Impact Summary</i>                                                                                |
|-----------------------------------------|------------------------------------------------------|--------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                | Evasion (TA1001), Reconnaissance (TA1002)            | Model Evasion (AT1010), Model Extraction (AT1005), Adversarial Example Generation (AT1020) | Target model runtime and infrastructure for evasion, or extract via side channels.                   |
| <b>2. LLM Models</b>                    | Poisoning (TA1003), Extraction (TA1004)              | Model Poisoning (AT1030), Model Inversion (AT1040), Model Extraction (AT1005)              | Tamper with model weights or extract internals via repeated queries or inversion attacks.            |
| <b>3. Information Sources</b>           | Poisoning (TA1003), Reconnaissance (TA1002)          | Data Poisoning (AT1050), Data Injection (AT1051), Input Manipulation (AT1060)              | Compromise external data sources to influence retrieval or context injection.                        |
| <b>4. Generated Output and Feedback</b> | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081)    | Misuse feedback channels to alter or skew model behavior through prompt or interaction manipulation. |



| <i>LLM Architecture Component</i>  | <i>ATLAS Tactics</i>                           | <i>Relevant ATLAS Techniques</i>                                                    | <i>Impact Summary</i>                                                                          |
|------------------------------------|------------------------------------------------|-------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------|
| <b>5. LLM Tools and Frameworks</b> | Execution (TA1007), Manipulation (TA1005)      | Tool Abuse (AT1090), Agent Manipulation (AT1091), API Abuse (AT1100)                | Compromise integrated tools or frameworks for unauthorized actions or misaligned outputs.      |
| <b>6. Application Interfaces</b>   | Manipulation (TA1005), Initial Access (TA1008) | Prompt Injection (AT1070), UI Redress Attacks (AT1110), Input Manipulation (AT1060) | Exploit user interfaces to inject malicious prompts or hijack interactions.                    |
| <b>7. Memory Storage</b>           | Poisoning (TA1003), Extraction (TA1004)        | Data Poisoning (AT1050), Model Inversion (AT1040), Extraction (AT1005)              | Poison or exfiltrate persistent memory to manipulate long-term model behavior or extract info. |
| <b>8. Memory Retrieval</b>         | Reconnaissance (TA1002), Manipulation (TA1005) | Reconnaissance (AT1002), Prompt Injection (AT1070), Feedback Loop Attacks (AT1081)  | Target memory query mechanisms to infer stored data or manipulate what the model recalls.      |

### 6.2.2.3. Mitigations

| <i>Component</i>                          | <i>Threat Summary</i>                                                                   | <i>Key Mitigations</i>                                                            |
|-------------------------------------------|-----------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                  | Susceptible to model evasion or extraction via timing or side-channel inference.        | Network segmentation, execution prevention, endpoint isolation.                   |
| <b>2. LLM Models</b>                      | Exposed to model poisoning, inversion, or extraction via repeated queries.              | Model hardening, differential privacy, adversarial training, query throttling.    |
| <b>3. Information Sources</b>             | Manipulated to feed malicious or biased context to the model.                           | Data validation and whitelisting, content filtering, input sanitation.            |
| <b>4. Generated Output &amp; Feedback</b> | Vulnerable to prompt injection, feedback loop manipulation, and hallucination chaining. | Output filtering, feedback auditing, semantic plausibility checks.                |
| <b>5. LLM Tools and Frameworks</b>        | May be misused via prompt exploits or agent misrouting.                                 | Tool permission restrictions, traceable execution, least-privilege configuration. |



| <i>Component</i>                 | <i>Threat Summary</i>                                                         | <i>Key Mitigations</i>                                                 |
|----------------------------------|-------------------------------------------------------------------------------|------------------------------------------------------------------------|
| <b>6. Application Interfaces</b> | Entry point for prompt injection, UI redress, and social engineering attacks. | Prompt hardening, input validation, secure UI design, user education.  |
| <b>7. Memory Storage</b>         | Exposed to poisoning or exfiltration of long-term stored information.         | Session expiration, access controls, integrity and consistency checks. |
| <b>8. Memory Retrieval</b>       | Exploitable to retrieve sensitive history or bias model outputs.              | Query filtering, access monitoring, memory auditing, anonymization.    |

### 6.2.3. RAG Architecture

#### 6.2.3.1. Overview

| <i>Component</i>                          | <i>Description</i>                                                                                                                                                           |
|-------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                  | The compute environment for hosting the LLM, retriever, and storage systems (e.g., cloud, hybrid).                                                                           |
| <b>2. LLM Model</b>                       | The pretrained language model responsible for generating answers based on retrieved context.                                                                                 |
| <b>3. Information Sources</b>             | External static or dynamic data injected into the LLM prompt, such as reference text or documents.                                                                           |
| <b>4. Generated Output &amp; Feedback</b> | Final LLM-generated response and optional user feedback loop used to refine future responses.                                                                                |
| <b>5. Application Interfaces</b>          | User-facing input layer or API endpoint for accepting natural language queries.                                                                                              |
| <b>6. Retrieval System</b>                | Middleware that interprets queries and fetches relevant documents from external sources.                                                                                     |
| <b>7. Data Store</b>                      | Storage layer for dense vector embeddings used to perform semantic similarity searches. May include graph storage if knowledge graph generation is part of the architecture. |
| <b>8. Data Sources</b>                    | Original documents (text, PDFs, webpages) used to build embeddings and provide source context.                                                                               |
| <b>9. Memory Store &amp; Retrieval</b>    | Optional session-based or persistent memory used to track user interactions or prior contexts.                                                                               |
| <b>10. LLM &amp; RAG Tools</b>            | Supporting SDKs, APIs, libraries (e.g., LangChain, Transformers, Neo4j) enabling prompt and RAG building and handling.                                                       |
| <b>11. RAG Indexing</b>                   | Pipeline for extracting, chunking, and embedding raw content; may also construct knowledge graphs.                                                                           |

#### 6.2.3.2. ATLAS Techniques and Tactics



| <i>LLM Architecture Component</i>       | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                           | <i>Impact Summary</i>                                                                           |
|-----------------------------------------|------------------------------------------------------|--------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                | Evasion (TA1001), Reconnaissance (TA1002)            | Model Evasion (AT1010), Model Extraction (AT1005), Adversarial Example Generation (AT1020) | Target model infrastructure for evasion or data extraction via side-channel and inference.      |
| <b>2. LLM Model</b>                     | Poisoning (TA1003), Extraction (TA1004)              | Model Poisoning (AT1030), Model Inversion (AT1040), Model Extraction (AT1005)              | Compromise model fidelity through poisoning or extract sensitive data through repeated probing. |
| <b>3. Information Sources</b>           | Poisoning (TA1003), Reconnaissance (TA1002)          | Data Poisoning (AT1050), Data Injection (AT1051), Input Manipulation (AT1060)              | Compromise external data sources to influence retrieval or context injection.                   |
| <b>4. Generated Output and Feedback</b> | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081)    | Alter model behavior using malicious feedback loops or prompt chaining attacks.                 |
| <b>5. Application Interfaces</b>        | Initial Access (TA1008), Manipulation (TA1005)       | UI Redress Attacks (AT1110), Prompt Injection (AT1070), API Abuse (AT1100)                 | Exploit user interfaces or APIs to inject malicious prompts or hijack workflows.                |
| <b>6. Retrieval System</b>              | Reconnaissance (TA1002), Poisoning (TA1003)          | Reconnaissance (AT1002), Data Poisoning (AT1050), Input Manipulation (AT1060)              | Manipulate retrieval logic to feed biased or malicious results to the LLM.                      |
| <b>7. Data Store</b>                    | Poisoning (TA1003), Extraction (TA1004)              | Data Poisoning (AT1050), Model Inversion (AT1040), Extraction (AT1005)                     | Tamper with embeddings to degrade relevance or extract stored vectors.                          |



| <i>LLM Architecture Component</i>      | <i>ATLAS Tactics</i>                       | <i>Relevant ATLAS Techniques</i>                                                  | <i>Impact Summary</i>                                                                     |
|----------------------------------------|--------------------------------------------|-----------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------|
| <b>8. Data Sources</b>                 | Poisoning (TA1003), Manipulation (TA1005)  | Data Injection (AT1051), Output Manipulation (AT1080), Prompt Injection (AT1070)  | Insert or manipulate documents to influence the grounding context of the LLM.             |
| <b>9. Memory Store &amp; Retrieval</b> | Extraction (TA1004), Manipulation (TA1005) | Model Inversion (AT1040), Feedback Loop Attacks (AT1081), Data Poisoning (AT1050) | Manipulate memory for long-term influence or to exfiltrate retrieved content.             |
| <b>10. LLM &amp; RAG Tools</b>         | Execution (TA1007), Manipulation (TA1005)  | Tool Abuse (AT1090), API Abuse (AT1100), Agent Manipulation (AT1091)              | Misuse orchestration or plugin tools to execute unintended commands or bypass validation. |
| <b>11. RAG Indexing</b>                | Poisoning (TA1003), Manipulation (TA1005)  | Data Injection (AT1051), Prompt Injection (AT1070), Output Manipulation (AT1080)  | Poison document content or embedding structure during preprocessing and chunking stages.  |

### 6.2.3.3. Mitigations

| <i>Component</i>                          | <i>Threat Summary</i>                                                          | <i>Key Mitigations</i>                                                         |
|-------------------------------------------|--------------------------------------------------------------------------------|--------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                  | Susceptible to timing attacks or inference evasion during LLM-query execution. | Network segmentation, isolated compute environments, endpoint protection.      |
| <b>2. LLM Model</b>                       | Targeted by model extraction, inversion, or poisoning via crafted queries.     | Model hardening, differential privacy, API rate limiting and token control.    |
| <b>3. Information Sources</b>             | Manipulated to feed malicious or biased context to the model.                  | Data validation and whitelisting, content filtering, input sanitation.         |
| <b>4. Generated Output &amp; Feedback</b> | Subject to hallucination amplification or feedback loop manipulation.          | Post-generation filtering, output feedback audits, semantic anomaly detection. |



| <i>Component</i>                       | <i>Threat Summary</i>                                                                | <i>Key Mitigations</i>                                                                      |
|----------------------------------------|--------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------|
| <b>5. Application Interfaces</b>       | Entry point for prompt injection, unauthorized queries, or redress attacks.          | Input sanitation and prompt control, UI hardening, throttling and access validation.        |
| <b>6. Retrieval System</b>             | Can be manipulated to feed biased or malicious documents into context.               | Context validation, use of vetted retrievers, document scoring and source whitelisting.     |
| <b>7. Data Store</b>                   | Susceptible to adversarial embedding poisoning and semantic leakage.                 | Integrity checks, rate-limited embedding operations, access control.                        |
| <b>8. Data Sources</b>                 | Can be used to inject hostile or misleading documents into retrieval context.        | Document validation, ingestion pipeline verification, source authenticity checks.           |
| <b>9. Memory Store &amp; Retrieval</b> | Exploitable for long-term poisoning or sensitive history retrieval.                  | Memory expiration policies, anonymization, access logging.                                  |
| <b>10. LLM &amp; RAG Tools</b>         | Tool plugins or orchestration layers may be hijacked or misused by prompt injection. | Tool access control, sandboxing, execution monitoring, scope restriction.                   |
| <b>11. RAG Indexing</b>                | Poisoned during document chunking, embedding, or graph construction.                 | Input validation, semantic filters, isolated indexing pipelines, human-in-the-loop vetting. |

## 6.2.4. Agentic Architecture

### 6.2.4.1. Overview

| <i>Component</i>                           | <i>Description</i>                                                                                          |
|--------------------------------------------|-------------------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                   | The underlying compute and orchestration environment enabling secure, scalable agent operations.            |
| <b>2. LLM Core Model</b>                   | The foundational language model responsible for reasoning, planning, and generating text-based decisions.   |
| <b>3. Information Sources</b>              | External static or dynamic data injected into the LLM prompt, such as reference text or documents.          |
| <b>4. Agent Framework / Executor</b>       | The runtime logic or engine that interprets plans and executes actions, typically via tools or APIs.        |
| <b>5. Planning Module</b>                  | Component responsible for multi-step reasoning, goal formulation, and task decomposition.                   |
| <b>6. Tools</b>                            | External tools, APIs, or services that agents invoke to accomplish tasks (e.g., web search, databases).     |
| <b>7. Memory System</b>                    | Persistent or temporary store of historical interactions, decisions, and results, used to maintain context. |
| <b>8. Observation and Feedback Loop</b>    | Mechanism by which the agent monitors its environment and adjusts future behavior based on outcomes.        |
| <b>9. Application Interface</b>            | User interface or communication layer enabling users to interact with, supervise, or guide the agent.       |
| <b>10. Generated Output &amp; Feedback</b> | Final LLM-generated response and optional user feedback loop used to refine future responses.               |

### 6.2.4.2. ATLAS Techniques and Tactics



| <i>LLM Architecture Component</i>  | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                           | <i>Impact Summary</i>                                                                              |
|------------------------------------|------------------------------------------------------|--------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>           | Evasion (TA1001), Reconnaissance (TA1002)            | Model Evasion (AT1010), Adversarial Example Generation (AT1020), Model Extraction (AT1005) | Target system infrastructure or exploit timing-based behaviors to extract data or trigger evasion. |
| <b>2. LLM Core Model</b>           | Poisoning (TA1003), Extraction (TA1004)              | Model Poisoning (AT1030), Model Inversion (AT1040), Model Extraction (AT1005)              | Compromise model weights to bias agentic behavior or extract sensitive training data.              |
| <b>3. Information Sources</b>      | Poisoning (TA1003), Reconnaissance (TA1002)          | Data Poisoning (AT1050), Data Injection (AT1051), Input Manipulation (AT1060)              | Compromise external data sources to influence retrieval or context injection.                      |
| <b>4. Agent Framework/Executor</b> | Execution (TA1007), Manipulation (TA1005)            | Tool Abuse (AT1090), Agent Manipulation (AT1091), API Abuse (AT1100)                       | Abuse agent execution logic to perform unauthorized actions or cause cascading effects.            |
| <b>5. Planning Module</b>          | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081)    | Manipulate planning or decision-making through adversarial prompts or recursive attacks.           |
| <b>6. Tools</b>                    | Initial Access (TA1008), Execution (TA1007)          | API Abuse (AT1100), Tool Abuse (AT1090), UI Redress Attacks (AT1110)                       | Exploit tool connections or APIs to execute unintended commands or leak sensitive outputs.         |
| <b>7. Memory System</b>            | Poisoning (TA1003), Extraction (TA1004)              | Model Inversion (AT1040), Data Poisoning (AT1050), Feedback Loop Attacks (AT1081)          | Corrupt memory to manipulate long-term planning or to exfiltrate contextual data.                  |



| <i>LLM Architecture Component</i>         | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                        | <i>Impact Summary</i>                                                                                |
|-------------------------------------------|------------------------------------------------------|-----------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------|
| <b>8. Observation and Feedback System</b> | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Observation Poisoning (AT1120) | Influence observational input or feedback loops to derail agent behavior or learning.                |
| <b>9. Application Interfaces</b>          | Initial Access (TA1008), Manipulation (TA1005)       | UI Redress Attacks (AT1110), Prompt Injection (AT1070), Input Manipulation (AT1060)     | Use social engineering or UI manipulation to inject commands or override human alignment interfaces. |
| <b>10. Generated Output and Feedback</b>  | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081) | Alter model behavior using malicious feedback loops or prompt chaining attacks.                      |

### 6.2.4.3. Mitigations

| <i>Component</i>                   | <i>Threat Summary</i>                                                                   | <i>Key Mitigations</i>                                                             |
|------------------------------------|-----------------------------------------------------------------------------------------|------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>           | Susceptible to side-channel attacks, resource abuse, or unauthorized agent deployments. | Endpoint protection, network segmentation, execution isolation.                    |
| <b>2. LLM Core Model</b>           | Targeted by model poisoning, extraction, or inversion via reasoning paths.              | Model hardening, differential privacy, adversarial training techniques.            |
| <b>3. Information Sources</b>      | Manipulated to feed malicious or biased context to the model.                           | Data validation and whitelisting, content filtering, input sanitation.             |
| <b>4. Agent Framework/Executor</b> | Vulnerable to tool abuse or action redirection via manipulated reasoning outputs.       | Action permission boundaries, tool sandboxing, execution auditing.                 |
| <b>5. Planning Module</b>          | Susceptible to prompt chaining and decision manipulation via nested reasoning loops.    | Prompt input validation, chain-of-thought validation, recursive output monitoring. |
| <b>6. Tools</b>                    | May be hijacked for unintended actions or remote access.                                | Tool invocation restrictions, API access control, least privilege configurations.  |
| <b>7. Memory System</b>            | Long-term poisoning or recall manipulation can skew agent behavior across sessions.     | Session-bound memory, validation and trimming, expiration and anonymization.       |



| <i>Component</i>                           | <i>Threat Summary</i>                                                            | <i>Key Mitigations</i>                                                         |
|--------------------------------------------|----------------------------------------------------------------------------------|--------------------------------------------------------------------------------|
| <b>8. Observation/Feedback Loop</b>        | Can be manipulated via fake observations or altered action-reaction patterns.    | Feedback auditing, observability tooling, semantic plausibility filters.       |
| <b>9. Application Interfaces</b>           | Entry point for prompt injection, social engineering, and agent hijack attempts. | UI redress defense, prompt hardening, secure templates, user education.        |
| <b>10. Generated Output &amp; Feedback</b> | Subject to hallucination amplification or feedback loop manipulation.            | Post-generation filtering, output feedback audits, semantic anomaly detection. |

## 6.2.5. Agentic RAG Architecture

### 6.2.5.1. Overview

| <i>Component</i>                           | <i>Description</i>                                                                                                                                                           |
|--------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>                   | The underlying compute and orchestration environment enabling secure, scalable agent operations.                                                                             |
| <b>2. LLM Core Model</b>                   | The foundational language model responsible for reasoning, planning, and generating text-based decisions.                                                                    |
| <b>3. Information Sources</b>              | External static or dynamic data injected into the LLM prompt, such as reference text or documents.                                                                           |
| <b>4. Agent Framework / Executor</b>       | The runtime logic or engine that interprets plans and executes actions, typically via tools or APIs.                                                                         |
| <b>5. Planning Module</b>                  | Component responsible for multi-step reasoning, goal formulation, and task decomposition.                                                                                    |
| <b>6. Tools</b>                            | External tools, APIs, or services that agents invoke to accomplish tasks (e.g., web search, databases).                                                                      |
| <b>7. Memory System</b>                    | Persistent or temporary store of historical interactions, decisions, and results, used to maintain context.                                                                  |
| <b>8. Observation and Feedback Loop</b>    | Mechanism by which the agent monitors its environment and adjusts future behavior based on outcomes.                                                                         |
| <b>9. Application Interface</b>            | User interface or communication layer enabling users to interact with, supervise, or guide the agent.                                                                        |
| <b>10. Generated Output &amp; Feedback</b> | Final LLM-generated response and optional user feedback loop used to refine future responses.                                                                                |
| <b>11. Retrieval System</b>                | Middleware that interprets queries and fetches relevant documents from external sources.                                                                                     |
| <b>12. Data Store</b>                      | Storage layer for dense vector embeddings used to perform semantic similarity searches. May include graph storage if knowledge graph generation is part of the architecture. |
| <b>13. Data Sources</b>                    | Original documents (text, PDFs, webpages) used to build embeddings and provide source context.                                                                               |
| <b>14. RAG Tools</b>                       | Supporting SDKs, APIs, libraries (e.g., Neo4j, Chroma) enabling RAG building and handling.                                                                                   |
| <b>15. RAG Indexing</b>                    | Pipeline for extracting, chunking, and embedding raw content; may also construct knowledge graphs.                                                                           |

### 6.2.5.2. ATLAS Techniques and Tactics



| <i>LLM Architecture Component</i>  | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                           | <i>Impact Summary</i>                                                                              |
|------------------------------------|------------------------------------------------------|--------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------|
| <b>1. Infrastructure</b>           | Evasion (TA1001), Reconnaissance (TA1002)            | Model Evasion (AT1010), Adversarial Example Generation (AT1020), Model Extraction (AT1005) | Target system infrastructure or exploit timing-based behaviors to extract data or trigger evasion. |
| <b>2. LLM Core Model</b>           | Poisoning (TA1003), Extraction (TA1004)              | Model Poisoning (AT1030), Model Inversion (AT1040), Model Extraction (AT1005)              | Compromise model weights to bias agentic behavior or extract sensitive training data.              |
| <b>3. Information Sources</b>      | Poisoning (TA1003), Reconnaissance (TA1002)          | Data Poisoning (AT1050), Data Injection (AT1051), Input Manipulation (AT1060)              | Compromise external data sources to influence retrieval or context injection.                      |
| <b>4. Agent Framework/Executor</b> | Execution (TA1007), Manipulation (TA1005)            | Tool Abuse (AT1090), Agent Manipulation (AT1091), API Abuse (AT1100)                       | Abuse agent execution logic to perform unauthorized actions or cause cascading effects.            |
| <b>5. Planning Module</b>          | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081)    | Manipulate planning or decision-making through adversarial prompts or recursive attacks.           |
| <b>6. Tools</b>                    | Initial Access (TA1008), Execution (TA1007)          | API Abuse (AT1100), Tool Abuse (AT1090), UI Redress Attacks (AT1110)                       | Exploit tool connections or APIs to execute unintended commands or leak sensitive outputs.         |
| <b>7. Memory System</b>            | Poisoning (TA1003), Extraction (TA1004)              | Model Inversion (AT1040), Data Poisoning (AT1050), Feedback Loop Attacks (AT1081)          | Corrupt memory to manipulate long-term planning or to exfiltrate contextual data.                  |



| <i>LLM Architecture Component</i>         | <i>ATLAS Tactics</i>                                 | <i>Relevant ATLAS Techniques</i>                                                        | <i>Impact Summary</i>                                                                                |
|-------------------------------------------|------------------------------------------------------|-----------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------|
| <b>8. Observation and Feedback System</b> | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Observation Poisoning (AT1120) | Influence observational input or feedback loops to derail agent behavior or learning.                |
| <b>9. Application Interfaces</b>          | Initial Access (TA1008), Manipulation (TA1005)       | UI Redress Attacks (AT1110), Prompt Injection (AT1070), Input Manipulation (AT1060)     | Use social engineering or UI manipulation to inject commands or override human alignment interfaces. |
| <b>10. Generated Output and Feedback</b>  | Manipulation (TA1005), Influence Operations (TA1006) | Prompt Injection (AT1070), Output Manipulation (AT1080), Feedback Loop Attacks (AT1081) | Alter model behavior using malicious feedback loops or prompt chaining attacks.                      |
| <b>11. Retrieval System</b>               | Reconnaissance (TA1002), Poisoning (TA1003)          | Reconnaissance (AT1002), Data Poisoning (AT1050), Input Manipulation (AT1060)           | Manipulate retrieval logic to feed biased or malicious results to the LLM.                           |
| <b>12. Data Store</b>                     | Poisoning (TA1003), Extraction (TA1004)              | Data Poisoning (AT1050), Model Inversion (AT1040), Extraction (AT1005)                  | Tamper with embeddings to degrade relevance or extract stored vectors.                               |
| <b>13. Data Sources</b>                   | Poisoning (TA1003), Manipulation (TA1005)            | Data Injection (AT1051), Output Manipulation (AT1080), Prompt Injection (AT1070)        | Insert or manipulate documents to influence the grounding context of the LLM.                        |
| <b>14. RAG Tools</b>                      | Execution (TA1007), Manipulation (TA1005)            | Tool Abuse (AT1090), API Abuse (AT1100), Agent Manipulation (AT1091)                    | Misuse orchestration or plugin tools to execute unintended commands or bypass validation.            |



| <i>LLM Architecture Component</i> | <i>ATLAS Tactics</i>                      | <i>Relevant ATLAS Techniques</i>                                                 | <i>Impact Summary</i>                                                                    |
|-----------------------------------|-------------------------------------------|----------------------------------------------------------------------------------|------------------------------------------------------------------------------------------|
| 15. RAG Indexing                  | Poisoning (TA1003), Manipulation (TA1005) | Data Injection (AT1051), Prompt Injection (AT1070), Output Manipulation (AT1080) | Poison document content or embedding structure during preprocessing and chunking stages. |

### 6.2.5.3. Mitigations

| <i>Component</i>                | <i>Threat Summary</i>                                                                   | <i>Key Mitigations</i>                                                                  |
|---------------------------------|-----------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------|
| 1. Infrastructure               | Susceptible to side-channel attacks, resource abuse, or unauthorized agent deployments. | Endpoint protection, network segmentation, execution isolation.                         |
| 2. LLM Core Model               | Targeted by model poisoning, extraction, or inversion via reasoning paths.              | Model hardening, differential privacy, adversarial training techniques.                 |
| 3. Information Sources          | Manipulated to feed malicious or biased context to the model.                           | Data validation and whitelisting, content filtering, input sanitation.                  |
| 4. Agent Framework/Executor     | Vulnerable to tool abuse or action redirection via manipulated reasoning outputs.       | Action permission boundaries, tool sandboxing, execution auditing.                      |
| 5. Planning Module              | Susceptible to prompt chaining and decision manipulation via nested reasoning loops.    | Prompt input validation, chain-of-thought validation, recursive output monitoring.      |
| 6. Tools                        | May be hijacked for unintended actions or remote access.                                | Tool invocation restrictions, API access control, least privilege configurations.       |
| 7. Memory System                | Long-term poisoning or recall manipulation can skew agent behavior across sessions.     | Session-bound memory, validation and trimming, expiration and anonymization.            |
| 8. Observation/Feedback Loop    | Can be manipulated via fake observations or altered action-reaction patterns.           | Feedback auditing, observability tooling, semantic plausibility filters.                |
| 9. Application Interfaces       | Entry point for prompt injection, social engineering, and agent hijack attempts.        | UI redress defense, prompt hardening, secure templates, user education.                 |
| 10. Generated Output & Feedback | Subject to hallucination amplification or feedback loop manipulation.                   | Post-generation filtering, output feedback audits, semantic anomaly detection.          |
| 11. Retrieval System            | Can be manipulated to feed biased or malicious documents into context.                  | Context validation, use of vetted retrievers, document scoring and source whitelisting. |



| <i>Component</i>               | <i>Threat Summary</i>                                                                | <i>Key Mitigations</i>                                                                      |
|--------------------------------|--------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------|
| <b>12. Data Store</b>          | Susceptible to adversarial embedding poisoning and semantic leakage.                 | Integrity checks, rate-limited embedding operations, access control.                        |
| <b>13. Data Sources</b>        | Can be used to inject hostile or misleading documents into retrieval context.        | Document validation, ingestion pipeline verification, source authenticity checks.           |
| <b>14. LLM &amp; RAG Tools</b> | Tool plugins or orchestration layers may be hijacked or misused by prompt injection. | Tool access control, sandboxing, execution monitoring, scope restriction.                   |
| <b>15. RAG Indexing</b>        | Poisoned during document chunking, embedding, or graph construction.                 | Input validation, semantic filters, isolated indexing pipelines, human-in-the-loop vetting. |

## 7. Acknowledgements

---

### 7.1. Workstream Leads

- Josiah Hagen, Trend Micro
- Vinay Bansal, Cisco

### 7.2. Editors

- Irakle Dzneladze, IBM
- Taha Mehdi, Cisco
- Michael Scovetta, Microsoft
- Siddhartha Rao, Cisco
- Ryan Cosgrove, Thomson Reuters
- Fyodor Yarochkin, Trend Micro
- Vladimir Kropotov, Trend Micro
- Rob Michel, Lenovo
- Michael Rash
- Sarah Novotny

## 8. Appendices

---

### 8.1. Integration with Enterprise Security Roadmap

This section outlines a comprehensive integration framework for aligning incident response capabilities with an enterprise security roadmap.

It presents a structured approach across four key implementation phases (further discussed in Section 3.3., Incident Response Workflow): 1. Preparation 2. Detection & Analysis 3. Containment, Eradication, and Recovery 4. Post-Incident Activity

Each phase addresses five critical integration areas: 1. Framework Alignment 2. Handoff Procedures 3. Shared Terminology 4. Tool Integration 5. Unified Metrics



The framework provides specific success metrics and implementation indicators for each phase and integration area, enabling organizations to systematically measure their progress in creating a cohesive, enterprise-wide security incident response capability that bridges traditionally siloed security domains.

| Implementation Phase | Integration Area    | Success Metrics & Implementation Indicators                                                                                                                                                                                                                                                    |
|----------------------|---------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Preparation          | Framework Alignment | <ul style="list-style-type: none"> <li>· Framework mapping completion percentage</li> <li>· Integrated risk assessment adoption rate</li> <li>· Documentation standardization level</li> <li>· Approved alignment strategy with executive sign-off</li> </ul>                                  |
|                      | Handoff Procedures  | <ul style="list-style-type: none"> <li>· Documented transition point clarity assessment</li> <li>· Workflow validation exercise completion rates</li> <li>· Joint preparation activity participation metrics</li> <li>· Stakeholder approval of defined responsibility matrices</li> </ul>     |
|                      | Shared Terminology  | <ul style="list-style-type: none"> <li>· Glossary comprehensiveness rating</li> <li>· Terminology mapping coverage percentage</li> <li>· Documentation terminology compliance audit results</li> <li>· Training completion rates on standardized terminology</li> </ul>                        |
|                      | Tool Integration    | <ul style="list-style-type: none"> <li>· Monitoring integration completion percentage</li> <li>· Configuration synchronization success rate</li> <li>· Vulnerability assessment coverage metrics</li> <li>· Tool inventory and connection validation status</li> </ul>                         |
|                      | Unified Metrics     | <ul style="list-style-type: none"> <li>· KPI definition approval rate</li> <li>· Baseline measurement completeness</li> <li>· Dashboard implementation status</li> <li>· Reporting standardization level across domains</li> </ul>                                                             |
| Detection & Analysis | Framework Alignment | <ul style="list-style-type: none"> <li>· Process adaptation completion rate</li> <li>· Triage procedure alignment validation results</li> <li>· Forensic methodology integration assessment</li> <li>· Cross-domain detection workflow efficiency metrics</li> </ul>                           |
|                      | Handoff Procedures  | <ul style="list-style-type: none"> <li>· Protocol effectiveness scores in simulation exercises</li> <li>· Escalation criteria clarity assessment results</li> <li>· Workflow efficiency measurements in actual incidents</li> <li>· Cross-team investigation time reduction metrics</li> </ul> |



| Implementation Phase                           | Integration Area           | Success Metrics & Implementation Indicators                                                                                                                                                                                                                                                                                          |
|------------------------------------------------|----------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                                | <b>Shared Terminology</b>  | <ul style="list-style-type: none"> <li>· Classification scheme adoption rate across teams</li> <li>· Severity definition consistency score</li> <li>· Reporting language standardization level</li> <li>· Classification accuracy in incident categorization</li> </ul>                                                              |
|                                                | <b>Tool Integration</b>    | <ul style="list-style-type: none"> <li>· SIEM integration success percentage</li> <li>· Alert correlation effectiveness measurements</li> <li>· Investigation dashboard utilization metrics</li> <li>· False positive reduction rate from improved correlation</li> </ul>                                                            |
|                                                | <b>Unified Metrics</b>     | <ul style="list-style-type: none"> <li>· Detection efficiency metrics implementation rate</li> <li>· Analysis effectiveness KPI measurement accuracy</li> <li>· Cross-domain detection coverage verification</li> <li>· Time-to-detection improvement percentage</li> </ul>                                                          |
| <b>Containment, Eradication &amp; Recovery</b> | <b>Framework Alignment</b> | <ul style="list-style-type: none"> <li>· Containment strategy alignment validation results</li> <li>· Eradication procedure integration completeness</li> <li>· Recovery methodology effectiveness ratings</li> <li>· Cross-domain incident resolution time improvements</li> </ul>                                                  |
|                                                | <b>Handoff Procedures</b>  | <ul style="list-style-type: none"> <li>· Responsibility clarity assessment scores</li> <li>· Coordination protocol effectiveness in exercises</li> <li>· Verification procedure success rates</li> </ul>                                                                                                                             |
|                                                | <b>Shared Terminology</b>  | <ul style="list-style-type: none"> <li>· Handoff-related incident extension reduction</li> <li>· Containment action terminology standardization level</li> <li>· Recovery state definition consistency measurements</li> <li>· Validation language adoption rates</li> <li>· Communication efficiency improvement metrics</li> </ul> |
|                                                | <b>Tool Integration</b>    | <ul style="list-style-type: none"> <li>· Containment control integration completion percentage</li> <li>· Recovery tracking system effectiveness measurements</li> <li>· Verification system reliability metrics</li> <li>· Cross-platform control execution time reduction</li> </ul>                                               |



| Implementation Phase   | Integration Area    | Success Metrics & Implementation Indicators                                                                                                                                                                                                                                                                                                                  |
|------------------------|---------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Post-Incident Activity | Unified Metrics     | <ul style="list-style-type: none"> <li>· Containment effectiveness measurement accuracy</li> <li>· Recovery time objective achievement rates</li> <li>· Eradication validation metric reliability assessment</li> <li>· Business impact reduction percentage per incident</li> </ul>                                                                         |
|                        | Framework Alignment | <ul style="list-style-type: none"> <li>· Lessons learned methodology integration effectiveness</li> <li>· Improvement tracking process alignment validation</li> <li>· Knowledge management approach adoption metrics</li> </ul>                                                                                                                             |
|                        | Handoff Procedures  | <ul style="list-style-type: none"> <li>· Recurring incident reduction percentage</li> <li>· Joint review protocol effectiveness assessment</li> <li>· Improvement responsibility clarity rating</li> <li>· Knowledge sharing workflow efficiency measurements</li> </ul>                                                                                     |
|                        | Shared Terminology  | <ul style="list-style-type: none"> <li>· Cross-team implementation rates of identified improvements</li> <li>· Root cause categorization standardization level</li> <li>· Improvement action description consistency score</li> <li>· Knowledge sharing terminology alignment rate</li> </ul>                                                                |
|                        | Tool Integration    | <ul style="list-style-type: none"> <li>· Documentation quality assessment results</li> <li>· Improvement tracking system integration percentage</li> <li>· Knowledge management system connection reliability</li> <li>· Metrics dashboard comprehensiveness rating</li> </ul>                                                                               |
|                        | Unified Metrics     | <ul style="list-style-type: none"> <li>· Automation level of post-incident workflows</li> <li>· Improvement effectiveness KPI reliability assessments</li> <li>· Incident reduction measurement accuracy across domains</li> <li>· Learning efficiency metric implementation completeness</li> <li>· Time-to-improvement implementation reduction</li> </ul> |

### CoSAI Focus

CoSAI is an OASIS Open Project, bringing together an open ecosystem of AI and security experts from industry-leading organizations. The project is dedicated to sharing best practices for secure AI deployment and collaborating on AI security research and product development. The scope of CoSAI is specifically focused on the secure building, integration, deployment, and operation of AI



systems, with an emphasis on mitigating security risks unique to AI technologies. Other aspects of Trustworthy AI are deemed important but beyond the scope of the project including, ethics, fairness, explainability, bias detection, safety, consumer privacy, misinformation, hallucinations, deep fakes, or content safety concerns like hateful or abusive content, malware, or phishing generation. By concentrating on developing robust measures, best practices, and guidelines to safeguard AI systems against unauthorized access, tampering, or misuse, CoSAI aims to contribute to the responsible development and deployment of resilient, secure AI technologies.

## Guidelines on usage of more advanced AI systems (e.g. large language models (LLMs), multi-modal language models. etc) for drafting documents for OASIS CoSAI:

tl;dr: CoSAI contributions are actions performed by humans, who are responsible for the content of those contributions, based on their signed OASIS iCLA (and eCLA, if applicable). [Each contributor must confirm whether they are entitled to donate that material under the applicable open source license; OASIS and the CoSAI Project do not separately confirm that.] Each contributor is responsible for ensuring that all contributions comply with these AI use guidelines, including disclosure of any use of AI in contributions.

- Selection of AI systems: CoSAI recommends the use of reputable AI systems (lowering the risk of inadvertently incorporating infringing material).
- Model constraints: Currently, CoSAI or OASIS are not required to have a contract or financial agreement for using AI systems from specific vendors. However, CoSAI editors should consider employing varying tools to avoid potential fairness concerns among vendors.
- IP infringement: It is the responsibility of the individual who subscribes/prompts and receives a response from an AI system to confirm they have the right to repost and donate the content to OASIS under our rules.
- Transparency: CoSAI's goal will be to maintain transparency throughout the process by documenting substantial use of AI systems whenever possible (e.g., the prompts and the AI system used), and to ensure that all content, regardless of production by human or AI systems, was reviewed and edited by human experts. This helps build trust in the standards development process and ensures accountability.
- Human-edited content and quality control: CoSAI mandates human-reviewed or -edited results for any final outputs. A robust quality control process should be in place, involving careful review of the generated content for accuracy, relevance, and alignment with CoSAI's goals and principles. Human experts should scrutinize the output of AI systems to identify any errors, inconsistencies, or potential biases.



Iterative refinement: The use of AI systems in drafting standards should be seen as an iterative process, with the generated content serving as a starting point for further refinement and improvement by human experts. Multiple rounds of review and editing may be necessary to ensure the final standards meet the required quality and reliability thresholds.

## Copyright Notice

Copyright © OASIS Open 2025. All Rights Reserved. This document has been produced under the process and license terms stated in the OASIS Open Project rules: <https://www.oasis-open.org/policies-guidelines/open-projects-process/>. This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published, and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this section are included on all such copies and derivative works. The limited permissions granted above are perpetual and will not be revoked by OASIS or its successors or assigns. This document and the information contained herein is provided on an “AS IS” basis and OASIS DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY OWNERSHIP RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. OASIS AND ITS MEMBERS WILL NOT BE LIABLE FOR ANY DIRECT, INDIRECT, SPECIAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF ANY USE OF THIS DOCUMENT OR ANY PART THEREOF. The name “OASIS” is a trademark of OASIS, the owner and developer of this document, and should be used only to refer to the organization and its official outputs. OASIS welcomes reference to, and implementation and use of, documents, while reserving the right to enforce its marks against misleading uses. Please see <https://www.oasis-open.org/policies-guidelines/trademark/> for above guidance.

This is a Non-Standards Track Work Product. The patent provisions of the OASIS IPR Policy do not apply.

Non-Standards Track Copyright © OASIS Open 2025. All Rights Reserved. This document was last revised or approved by the CoSAI Open Project on the above date.