

# MANATECH

RESEARCH REPORT

## 2026 Enterprise Guide: Managing Autonomous AI Agent Risks and Governance

### Executive Summary

As of 2026, enterprise AI has transitioned from conversational chatbots to autonomous agents—systems capable of perceiving, deciding, and acting across multiple business environments. While traditional AI functions as a tool for decision support, agentic AI operates as an organizational actor with delegation rights to commit resources, modify databases, and coordinate across applications. This shift introduces the "Lethal Trifecta" of risk: access to private data, exposure to untrusted content, and the ability to exfiltrate information.

Current research indicates that approximately 75% of businesses plan to deploy AI agents by the end of 2026. However, existing governance frameworks designed for deterministic software are inadequate for the non-deterministic nature of autonomous agents. This briefing document outlines the critical risk categories defined by the OWASP Top 10 for Agentic Applications, proposes the Agentic Operating Model (AOM) for scalable governance, and provides actionable frameworks for implementing multi-layered guardrails.

### Detailed Analysis of Key Themes

#### 1. The Shift from Tools to Autonomous Actors

The fundamental difference between 2024-era chatbots and 2026-era agents lies in their operational autonomy. Agents do not merely generate text; they plan workflows, call APIs, and execute actions independently.

Dimension	Characteristic in Agentic AI
Autonomy	Agents initiate actions based on environmental signals rather than waiting for human prompts.
Persistence	Agents operate continuously, learning from feedback and adapting over time without repeated instruction.
Delegation	Agents possess formal authority to access systems of record and commit enterprise resources.

## 2. The Agentic Risk Landscape (OWASP Top 10)

In late 2025, OWASP released the first formal risk taxonomy for autonomous agents. These risks are distinct from standard Large Language Model (LLM) vulnerabilities because they involve the physical execution of tasks.

### Top Five Agentic Risks (2026):

- **ASI01: Agent Goal Hijacking:** Natural language instructions manipulate an agent's objectives or decision pathways. This is considered a "total loss of control" scenario.
- **ASI02: Tool Misuse & Exploitation:** Agents misuse legitimate tools (e.g., CRMs, code repos) due to prompt injection or unsafe delegation, often leading to data exfiltration.
- **ASI03: Identity & Privilege Abuse:** The "confused deputy" problem, where agents inherit user credentials and operate far beyond their intended scope at machine speed.
- **ASI04: Supply Chain Vulnerabilities:** Compromised plugins or third-party integrations introduce backdoors into agent workflows.
- **ASI05: Unexpected Code Execution (RCE):** Attackers exploit code-generation features to execute malicious scripts on host systems.

## 3. Context Engineering: The Root of Hallucination

Research suggests that 95% of enterprise AI pilots fail to deliver ROI because they lack "Context Infrastructure." Hallucination in agents is primarily a context problem, not a model problem. When agents lack domain-specific context, they fabricate metrics or policies to fill the gap.

### The Four Layers of the AI Context Stack:

1. **User Context:** Permissions, roles, and identity of the requester.
2. **Knowledge Context:** Certified business definitions, policies, and approved frameworks.
3. **Meaning Context:** Domain-specific interpretations and lineage relationships.
4. **Data Context:** Freshness, quality signals, and usage patterns of the underlying data.

## 4. The Agentic Operating Model (AOM)

To manage agents at scale, organizations are adopting the AOM, which specifies four interdependent layers of governance:

Layer	Function
<b>Cognitive</b>	Deploying specialized, domain-specific models rather than monolithic, general-purpose ones to make auditing tractable.
<b>Coordination</b>	Governing "Agent Swarms" via decentralized rules or "consensus mechanisms" that require multiple agents to agree before high-risk execution.
<b>Control</b>	Implementing real-time "Guardrail Agents" that intercept outputs and block actions if they exceed behavioral baselines.
<b>Governance</b>	Aligning agent behavior with ISO/IEC 42001 and NIST frameworks, ensuring every agent has a clear business owner.

## Important Quotes with Context

---

### On the Nature of Agent Risk

*"When you give AI agents the power to make decisions without a human in the loop, you also give them the power to affect people, processes and reputations in real time." — **Roger Connors**, Co-founder of Partners In Leadership. Context: Discussing the "black box" of decision-making that occurs when agents collaborate across internal and external APIs without oversight.*

### On Governance and Efficiency

*"We are saving 40% in efficiency... for everything related to governance. We're using the time savings to focus on optimizing our processes and upleveling the type of work we are doing." — **Danlei Alves**, Senior Data Governance Analyst at Porto Insurance. Context: Highlighting how automated governance platforms allow organizations to scale AI operations without increasing headcount for manual oversight.*

### On the "Lethal Trifecta"

*"AI agents have become the highest-value targets in enterprise security, and the least defended." — **Simon Willison** (via Lumenova AI). Context: Referencing the combination of access to private data, exposure to untrusted content, and the ability to exfiltrate, which most enterprise agents possess on day one.*

### On Strategic Implementation

*"Atlan's metadata lakehouse is configurable across all tools and flexible enough to get us to a future state where AI agents can access lineage context through the Model Context Protocol." — **Andrew Reiskind**, Chief Data Officer, Mastercard. Context: Illustrating the transition toward "context-by-design" as a prerequisite for secure agent deployment.*

---

## Actionable Insights and Frameworks

---

### The AWARE Framework for Agent Governance

Organizations should apply the **AWARE** lens to every agent use case to ensure comprehensive risk coverage:

- **A — Actor Intent:** Identify who is acting and for what specific job.
- **W — Work Context:** Determine if the action is appropriate for the given user and moment.
- **A — Autonomous Guardrails:** Apply runtime policies constraining agent purpose and data access.
- **R — Real-time Risk Scoring:** Continuously assess activity with automated blocking.

- **E — Ecosystem Observability:** Maintain end-to-end traceability for audit and forensics.

## Six-Layer Defense Architecture

To defend against the 2026 threat landscape, security teams must implement a layered execution boundary:

1. **Input Boundary:** Treat all external inputs (emails, web data) as untrusted; use prompt filtering.
2. **Execution Boundary:** Sandbox all code execution and enforce "least-privilege" tool access.
3. **Identity Boundary:** Create agent-specific service accounts with scoped permissions; avoid inheriting user credentials.
4. **Communication Boundary:** Authenticate all inter-agent messages in multi-agent swarms.
5. **Supply Chain Boundary:** Vet all plugins and pin dependency versions with cryptographic hashes.
6. **Behavioral Monitoring:** Use runtime anomaly detection to flag spikes in action volume or deviations from expected patterns.

## Strategic Roadmap for Enterprise Deployment

- **Standardize Before Automating:** AI agents cannot scale in undocumented or inconsistent business processes. Workflow standardization must precede agent deployment.
- **Shift from HITL to HOTL:** Move from **Human-in-the-Loop** (manual approvals) to **Human-on-the-Loop** (setting boundaries and intervening only on exceptions) to maintain operational speed.
- **Treat Agents as Employees:** Apply the same rigor to agent "hiring" (vetting) and "onboarding" (context provisioning) as human staff. This includes defining clear escalation paths and "kill switches" for emergency termination.
- **Adopt Specialized Toolkits:** Utilize open-source tools like the **Agent Governance Toolkit** for sub-millisecond policy enforcement and deterministic orchestration (e.g., YAML-defined workflows) to reduce the "Orchestration Gap."

## Want to explore this topic further?

Book a free discovery call to discuss how ManaTech can help your business implement these ideas.

[Book a Discovery Call](#)