

Free markets. Real solutions.

#### R Street Policy Study No. 325 May 2025



### The Rise of Al Agents: Anticipating Cybersecurity Opportunities, Risks, and the Next Frontier

By Haiman Wong and Tiffany Saade

To prepare for the next frontier of AI, this study explores the cybersecurity implications of AI agents and presents a threepronged framework to guide secure design and responsible deployment.

#### **Executive Summary**

The rise of AI agents represents a significant evolution in artificial intelligence (AI)—shifting from passive, prompt-based tools to increasingly autonomous systems capable of reasoning, memory, learning, and complex task execution. As industries adopt these agents, their impact on business operations, human-machine collaboration, and national security is accelerating.

In cybersecurity, AI agents are already proving to be valuable copilots to human analysts—enhancing threat detection, expediting incident response, and supporting overstretched cyber teams. Yet with greater power across the agentic infrastructure stack—spanning perception, reasoning, action, and memory—comes greater responsibility to ensure that agents are secure, explainable, and reliable.

To prepare for this next frontier of AI, this study explores the cybersecurity implications of AI agents and presents a three-pronged framework to guide secure design and responsible deployment: (1) identifying policy priorities, including voluntary, sector-specific guidelines for navigating human–agent collaboration; (2) anticipating emerging technological solutions to strengthen agentic oversight and cyber resilience; and (3) outlining best practices for

#### **Table of Contents**

Executive Summary	_ 1
Introduction	_ 2
Background	_ 5
Cybersecurity Benefits of AI Agents _	_ 8
Continuous Attack Surface Monitoring and Vulnerability Management	9
Real-Time Threat Detection and Incident Response	_ 10
Augmented Decision-Making and Cyber Workforce Support	_ 11
Cybersecurity Considerations and Potential Risks	12
Layer 1: Perception Module	_ 12
Layer 2: Reasoning Module	_ 14
Layer 3: Action Module	_ 15
Layer 4: Memory Module	_ 17
Anticipating Needs, Solutions, and Responsibilities	18
Policy Needs	18
Emerging Technological Solutions _	_ 20
Responsible Design and Deploymen Strategies for Developers	t 21
Conclusion	- <sup>21</sup>
	. 22
About the Authors	_ 23
Figure 1: Key Components of Advanced AI Agents	6
Figure 2: Cybersecurity Risks of Al-Generated Code	_ 7
Figure 3: Framework of Backdoor Attacks to Unified Foundation	

Models

13



developers, organizations, and end-users to ensure that agents augment—rather than replace—human talent and decision-making.

Ultimately, the goal is not merely to keep pace with the advancement of agentic systems, but to shape their trajectory—harnessing their benefits, minimizing their risks, and safeguarding America's technological leadership, national security, and values in the process.

#### Introduction

While 2023 was dubbed the year of "generative artificial intelligence (Gen AI)" and 2024 was marked by a steady march toward "AI practicality," 2025 opened with high expectations that it would become the year of "AI agents."<sup>1</sup> At its core, an AI agent is an "autonomous intelligent system powered by artificial intelligence and designed to perform specific tasks independently without the need for human intervention."<sup>2</sup> However, as seen with previous AI advancements—whether Gen AI, open-source AI, or large language models (LLMs)—a single, universally adopted definition remains elusive.<sup>3</sup> Some experts describe AI agents as "applications that attempt to achieve a goal by observing the world and acting upon it using the tools that [they have] at [their] disposal."<sup>4</sup> Others characterize agents as "layers on top of the language models that observe and collect information, provide input to the model and together generate an action plan and communicate that to the useror even act on their own, if permitted."<sup>5</sup> Though these definitions vary in their precise phrasing and perspective, each consistently emphasizes the agent's ability to pursue and complete goals autonomously using a suite of capabilities that includes learning, memory, planning, reasoning, decision-making, and adaptation.<sup>6</sup>

Notably, not all AI agents are created equal. Non-agentic and agentic AI systems differ in how they operate, particularly in autonomy and goal-setting. Non-agentic systems, such as earlier versions of ChatGPT or Alexa, respond to user prompts without retaining memory, setting goals, or initiating actions on their own.<sup>7</sup> In contrast, agentic AI systems pursue objectives over time, using contextual awareness, autonomous planning, and adaptive reasoning to carry out multi-step tasks with minimal human input or oversight.<sup>8</sup> For example, an agentic system might autonomously research a topic across multiple websites, generate a tailored report, and distribute it through email.<sup>9</sup> This added layer of sophistication gives AI agents greater potential to navigate across different domains and produce more complex,



NON-AGENTIC AI SYSTEMS such as earlier versions of ChatGPT or Alexa, respond to user prompts without retaining memory, setting goals, or initiating actions on their own.



AGENTIC AI SYSTEMS pursue objectives over time, using contextual awareness, autonomous planning, and adaptive reasoning to carry out multi-step tasks with minimal human input or oversight.

Haiman Wong et al., "The Transformative Role of AI in Cybersecurity: Anticipating and Preparing for Future Applications and Benefits," R Street Institute, Jan. 24, 2024. https://www.rstreet.org/commentary/the-transformative-role-of-ai-in-cybersecurity-anticipating-and-preparing-for-future-applications-and-benefits; Eric Johnson, "2023 Was The Year Of AI Hype—2024 Is The Year Of AI Practicality," *Forbes*, April 2, 2024. https://www.forbes.com/councils/forbestechcouncil/2024/04/02/2023was-the-year-of-ai-hype-2024-is-the-year-of-ai-practicality; Tae Kim, "Nvidia CEO Says 2025 Is the Year of AI Agents," Barron's, Jan. 7, 2025. https://www.barrons.com/ articles/nvidia-stock-ceo-ai-agents-8c20ddfb.

<sup>2.</sup> Kinza Yasar, "What are Al agents?," TechTarget, December 2024. https://www.techtarget.com/searchenterpriseai/definition/Al-agents.

<sup>3.</sup> Maxwell Zeff and Kyle Wiggers, "No one knows what the hell an AI agent is," TechCrunch, March 14, 2025. https://techcrunch.com/2025/03/14/no-one-knows-what-the-hell-an-ai-agent-is.

<sup>4.</sup> Michael Nuñez, "Google maps the future of AI agents: Five lessons for businesses," VentureBeat, Jan. 6, 2025. https://venturebeat.com/ai/google-maps-the-future-ofai-agents-five-lessons-for-businesses; Julia Wiesinger et al., "Agents," Google, September 2024. https://www.kaggle.com/whitepaper-agents.

<sup>5.</sup> Susanna Ray, "Al agents—what they are, and how they'll change the way we work," Microsoft, Nov. 19, 2024. https://news.microsoft.com/source/features/ai/aiagents-what-they-are-and-how-theyll-change-the-way-we-work.

<sup>6. &</sup>quot;What is an AI Agent?," GoogleCloud, last accessed March 29, 2025. https://cloud.google.com/discover/what-are-ai-agents.

<sup>7.</sup> Anna Gutowska, "What are Al agents?," IBM, July 3, 2024. https://www.ibm.com/think/topics/ai-agents.



real-world outcomes. These capabilities not only unlock new opportunities but also raise novel challenges for oversight, governance, and cybersecurity that differ in both scope and scale from earlier AI systems.<sup>10</sup>

Al agents have already begun transforming workflows across various sectors, especially in software engineering, where they are increasingly embedded into routine development tasks.<sup>11</sup> For example, Claude 3.7 and Cursor AI are automating software development tasks such as code generation, refactoring, and debugging.<sup>12</sup> In cybersecurity, Microsoft's Security Copilot can autonomously triage phishing alerts, dynamically update its detection capabilities based on analyst feedback, and flag security policy configuration issues.<sup>13</sup> Other similar and emerging cybersecurityfocused AI agents include Exabeam's Copilot, Cymulate AI Copilot, and Oleria Copilot, all of which streamline cyber incident investigations and simulations.<sup>14</sup> Beyond these programming- and cybersecurity-centric agents, general-purpose agents like OpenAl's Operator, Anthropic's Computer Use, and Google's Project Astra also stand out for their potential to coordinate tasks such as multi-step web navigation, form completion, and cross-application integration.<sup>15</sup> Looking ahead, many experts anticipate that agentic AI will continue advancing in three key areas: (1) enhancing reasoning and contextual understanding to solve problems more effectively, (2) delivering greater autonomy for complex task execution, and (3) augmenting the human workforce.16

Some experts have even suggested that—in the next five years—workers will face a reality in which AI is doing 80 percent of their day-to-day tasks.<sup>17</sup> As AI agents begin shouldering more cognitive and operational tasks, their influence over business operations, workforce dynamics, and digital infrastructure will become increasingly consequential. Accordingly, experts now view technological leadership in agentic AI as a matter of strategic national importance. Although many of today's leading agents continue to be developed by American AI and technology companies, the global competition to lead in agentic AI persists.<sup>18</sup>



Some experts have even suggested that—in the next five years—workers will face a reality in which AI is doing 80 percent of their day-to-day tasks.

10. Shomit Ghose, "The Next 'Next Big Thing': Agentic Al's Opportunities and Risks," UC Berkeley Sutardja Center for Entrepreneurship & Technology, Dec. 19, 2024. https://scet.berkeley.edu/the-next-next-big-thing-agentic-ais-opportunities-and-risks.

12. Ibid.

13. Vasu Jakkal, "Microsoft unveils Microsoft Security Copilot agents and new protections for AI," Microsoft, March 24, 2025. https://www.microsoft.com/en-us/ security/blog/2025/03/24/microsoft-unveils-microsoft-security-copilot-agents-and-new-protections-for-ai; Jeffrey Schwartz, "Microsoft Gives Security Copilot Some Autonomy," DarkReading, March 24, 2025. https://www.darkreading.com/cybersecurity-operations/microsoft-gives-security-copilot-autonomy.

- 14. Louis Columbus, "From alerts to autonomy: How leading SOCs use AI copilots to fight signal overload and staffing shortfalls," VentureBeat, March 24, 2025. https:// venturebeat.com/security/ai-copilots-cut-false-positives-and-burnout-in-overworked-socs
- Will Douglas Heaven, "OpenAl launches Operator—an agent that can use a computer for you," MIT Technology Review, Jan. 23, 2025. https://www.technologyreview. com/2025/01/23/1110484/openai-launches-operator-an-agent-that-can-use-a-computer-for-you; Samuel Axon, "Anthropic publicly releases Al tool that can take over the user's mouse cursor," Ars Technica, Oct. 22, 2024. https://arstechnica.com/ai/2024/10/anthropic-publicly-releases-ai-tool-that-can-take-over-the-users-mousecursor; Will Douglas Heaven, "Google's new Project Astra could be generative Al's killer app," MIT Technology Review, Dec. 11, 2024. https://www.technologyreview. com/2024/12/11/1108493/googles-new-project-astra-could-be-generative-ais-killer-app.
- 16. Melissa Heikkilä and Will Douglas Heaven, "Anthropic's chief scientist on 4 ways agents will be even better in 2025," MIT Technology Review, Jan. 11, 2025. https:// www.technologyreview.com/2025/01/11/1109909/anthropics-chief-scientist-on-5-ways-agents-will-be-even-better-in-2025; Kate Whiting, "The rise of 'AI agents': What they are and how to manage the risks," IBM, Dec. 16, 2024. https://www.weforum.org/stories/2024/12/ai-agents-risks-artificial-intelligence.
- 17. Kate Whiting, "What is an AI agent and what will they do? Experts explain," World Economic Forum, July 24, 2024. https://www.weforum.org/stories/2024/07/whatis-an-ai-agent-experts-explain.
- 18. Alexander Puutio, "The Agentic AI Race Is On, And The Blue Chips Are All In," *Forbes*, Nov. 15, 2024. https://www.forbes.com/sites/alexanderpuutio/2024/11/15/the-agentic-ai-race-is-on-and-the-blue-chips-are-all-in; Sean Oesch et al., "Agentic AI and the Cyber Arms Race," arXiv, Feb. 10, 2025. https://arxiv.org/html/2503.04760v1.

Janakiram MSV, "GitHub Copilot Agent And The Rise Of AI Coding Assistants," Forbes, Feb. 8, 2025. https://www.forbes.com/sites/janakirammsv/2025/02/08/githubcopilot-agent-and-the-rise-of-ai-coding-assistants; Matt Marshall, "Anthropic's stealth enterprise coup: How Claude 3.7 is becoming the coding agent of choice," VentureBeat, March 11, 2025. https://venturebeat.com/ai/anthropics-stealth-enterprise-coup-how-claude-3-7-is-becoming-the-coding-agent-of-choice; Julie Bort, "AI coding assistant Cursor reportedly tells a 'vibe coder' to write his own damn code," TechCrunch, March 14, 2025. https://techcrunch.com/2025/03/14/ai-codingassistant-cursor-reportedly-tells-a-vibe-coder-to-write-his-own-damn-code.



Outside of the United States, Manus, developed by Wuhan-based startup Butterfly Effect, garnered global attention in March 2025 when it claimed to be "the world's first general AI agent, using multiple AI models (such as Anthropic's Claude 3.5 Sonnet and fine-tuned versions of Alibaba's open-source Qwen) and various independently operating agents to act autonomously on a wide range of tasks."<sup>19</sup> Manus is designed to execute diverse, goal-oriented tasks independently, including language translation, online purchasing, research synthesis, and 3D game development from a single prompt.<sup>20</sup> The launch of Manus underscores the strategic significance of AI agents in global innovation ecosystems and mirrors trends already observed in the open-source AI movement and the broader competitive landscape around AI innovation and deployment.<sup>21</sup>

As with earlier waves of AI development, establishing technological leadership in agentic AI may carry both economic and substantial geopolitical implications, especially if agents become embedded in critical workflows across sensitive sectors, such as finance, healthcare, and defense.<sup>22</sup> In March 2025, for example, the Pentagon contracted with ScaleAI to give AI what one reporter termed "its most prominent role in the Western defense sector to date."<sup>23</sup> Framed as the Department of Defense's "first foray" into deploying agentic systems across military workflows, the initiative aims to accelerate strategic assessments, simulate war-gaming scenarios, and modernize campaign development.<sup>24</sup> This deal not only underscores the rising stakes of securing America's technological leadership but also ushers in the age of "agentic warfare" at a moment when many are only beginning to grapple with the full scope of this technology's capabilities, limitations, and risks.<sup>25</sup>

The rise of AI agents presents a critical window of opportunity to take a closer look—not only at how AI agents are being developed—but also at how they can be secured and governed.<sup>26</sup> As agentic systems begin streamlining business operations, problem-solving, and human—machine collaboration, the implications extend far beyond technological innovation.<sup>27</sup> In light of these rapidly evolving shifts, policymakers must craft governance strategies that are balanced, forward-looking, and flexible—strategies that support agentic innovation and use while proactively mitigating risks and ensuring long-term national security resilience. R Street Policy Study No. 325 May 2025



In March 2025, for example, the Pentagon contracted with ScaleAI to give AI what one reporter termed "its most prominent role in the Western defense sector to date."

19. Caiwei Chen, "Everyone in Al is talking about Manus. We put it to the test," MIT Technology Review, March 11, 2025. https://www.technologyreview. com/2025/03/11/1113133/manus-ai-review.

26. Helen Toner et al., "Through the Chat Window and Into the Real World: Preparing for AI Agents," Center for Security and Emerging Technology, October 2024. https:// cset.georgetown.edu/publication/through-the-chat-window-and-into-the-real-world-preparing-for-ai-agents.

<sup>20.</sup> Rhiannon Williams, "The Download: testing new AI agent Manus, and Waabi's virtual robotruck ambitions," MIT Technology Review, March 12, 2025. https://www. technologyreview.com/2025/03/12/1113172/the-download-testing-new-ai-agent-manus-and-waabis-virtual-robotruck-ambitions; Coco Feng, "Manus draws upbeat reviews of nascent system," MSN, March 3, 2025. https://www.msn.com/en-sg/news/other/less-structure-more-intelligence-ai-agent-manus-draws-upbeat-reviewsof-nascent-system/ar-AA1AJyOE.

<sup>21.</sup> Ibid.

<sup>22.</sup> Kieran Garvey, "How Agentic AI will transform financial services with autonomy, efficiency, and inclusion," World Economic Forum, Dec. 2, 2024. https://www. weforum.org/stories/2024/12/agentic-ai-financial-services-autonomy-efficiency-and-inclusion; Sara Heath, "Epic's take on agentic AI designed to boost patient experience," TechTarget, March 5, 2025. https://www.techtarget.com/patientengagement/feature/Epics-take-on-agentic-AI-designed-to-boost-patient-experience; Brandon Vigliarolo, "It begins: Pentagon to give AI agents a role in decision making, ops planning," The Register, March 5, 2025. https://www.theregister. com/2025/03/05/dod\_taps\_scale\_to\_bring.

<sup>23.</sup> Ibid.

<sup>24.</sup> Ibid.

<sup>25.</sup> Julia Hornstein, "AI agents are coming to the military. VCs love it, but researchers are a bit wary," Business Insider, March 8, 2025. https://www.businessinsider.com/ ai-agents-coming-military-new-scaleai-contract-2025-3.



This study examines the benefits, risks, cybersecurity considerations, and policy needs likely to define the emerging frontier of AI advancement. It begins by outlining the architecture of agentic systems and explaining how they differ from earlier generations of AI tools. It then explores how AI agents are already being deployed in cybersecurity use cases and identifies the new categories of risk they introduce across four distinct infrastructure layers: perception, reasoning, action, and memory. Finally, the paper presents a framework for secure and responsible deployment, emphasizing the roles of policy, technical safeguards, and organizational practices in strengthening long-term resilience and demonstrating that—if guided with care and foresight—agentic systems could mark not just the next phase of AI, but a turning point in how digital security is built and sustained.

#### Background

Although AI agents dominated news headlines in late 2024 and early 2025, their conceptual foundations trace back to the 1970s and 1980s, when research explored how capable systems were of sensing and acting intelligently within an environment.<sup>28</sup> These early systems, often referred to as "intelligent agents," powered linguistic analysis, biomedical applications, and robotics, relying on rule-based logic and limited autonomy due to constraints in hardware, computing power, and algorithmic sophistication.<sup>29</sup> At the time, these agents were described as "a new type of AI system capable of adapting, learning from data, and making complex decisions in changing environments."<sup>30</sup>

The renewed surge of interest in agents today reflects a convergence of technological advancements: scalable cloud infrastructure, advanced foundation models such as GPT-4 and Claude 3.5, and modular architectures that support planning, reasoning, and action with minimal human oversight.<sup>31</sup> Tools like AutoGPT, an "experimental, open-source Python application that uses GPT-4 to act autonomously," also helped popularize this shift from reactive, prompt-based tools to proactive, goal-driven systems capable of coordinating complex tasks.<sup>32</sup> As a result, AI agents are now positioned to be practical tools with significant operational and economic utility—from automating software development to automating customer service and even augmenting real-time cybersecurity defense.<sup>33</sup>

Architecturally, AI agents typically operate as a layer above LLMs and include four foundational components: perception, reasoning, action, and memory.<sup>34</sup> The perception module is responsible for ingesting data from external sources, such

#### R Street Policy Study No. 325 May 2025



Al agents are now positioned to be practical tools with significant operational and economic utility—from automating software development to automating customer service and even augmenting real-time cybersecurity defense.

<sup>28.</sup> Matvii Diadkov, "AI Agents: From Inception to Today," Forbes, March 18, 2025. https://www.forbes.com/councils/forbesbusinesscouncil/2025/03/18/ai-agents-from-inception-to-today.

<sup>29.</sup> Joseph Reagle, "The Etymology of 'Agent' and 'Proxy' in Computer Networking Discourse," Harvard University, Sept. 18, 1998. https://cyber.harvard.edu/archived\_ content/people/reagle/etymology-agency-proxy-19981217.html.

<sup>30.</sup> Diadkov. https://www.forbes.com/councils/forbesbusinesscouncil/2025/03/18/ai-agents-from-inception-to-today.

<sup>31.</sup> Timothy R. McIntosh et al., "From Google Gemini to OpenAI Q\* (Q-Star): A Survey on Reshaping the Generative Artificial Intelligence (AI) Research Landscape," MDPI, Jan. 30, 2025. https://www.mdpi.com/2227-7080/13/2/51.

<sup>32.</sup> Sabrina Ortiz, "What is Auto-GPT? Everything to know about the next powerful AI tool," ZDNet, April 14, 2023. https://www.zdnet.com/article/what-is-auto-gpteverything-to-know-about-the-next-powerful-ai-tool.

**<sup>33.</sup>** Diadkov. https://www.forbes.com/councils/forbesbusinesscouncil/2025/03/18/ai-agents-from-inception-to-today.

<sup>34.</sup> Whiting, "The rise of 'AI agents': What they are and how to manage the risks." https://www.weforum.org/stories/2024/12/ai-agents-risks-artificial-intelligence.



as user inputs or application programming interfaces (APIs).<sup>35</sup> After the data is gathered, the reasoning module leverages the LLM's capabilities to plan or infer the best course of action.<sup>36</sup> The action module can then execute tasks through tools, APIs, or integrations with third-party systems.<sup>37</sup> Finally, the memory module stores contextual information, often using vector databases or session-based memory managers.<sup>38</sup> This modular stack enables agents to operate across real-world applications and adapt while completing tasks in ways that static prompt chains or retrieval-augmented generation (RAG) pipelines cannot.<sup>39</sup>

Figure 1 illustrates how advanced AI agents may sense and interact with their environment, process information, and coordinate actions.

#### Figure 1: Key Components of Advanced Al Agents



Source: Kate Whiting, "The rise of 'AI agents': What they are and how to manage the risks," World Economic Forum, Dec. 16. 2024. https://www.weforum.org/stories/2024/12/ai-agents-risks-artificial-intelligence.

Behind this architecture lies a supporting infrastructure stack: model APIs for LLM access, memory stores for quick retrieval, session managers for coordinating task state, external tool integrations for operational output, and even open-source frameworks and libraries that enable modular development.<sup>40</sup> Multi-agent systems add another layer of sophistication, allowing agents to collaborate or delegate tasks to other agents within a shared environment.<sup>41</sup> While this growing interconnectedness can enhance agentic capabilities, it can also introduce new challenges around explainability, privacy, system security, and reliability.<sup>42</sup>

<sup>35.</sup> Yasar. https://www.techtarget.com/searchenterpriseai/definition/AI-agents.

<sup>36.</sup> Rine Diane Caballar and Cole Stryker, "What is agentic reasoning?," IBM, March 20, 2025. https://www.ibm.com/think/topics/agentic-reasoning.

<sup>37.</sup> Yasar. https://www.techtarget.com/searchenterpriseai/definition/AI-agents.

<sup>38.</sup> Cole Stryker, "What is AI agent memory?," IBM, March 18, 2025. https://www.ibm.com/think/topics/ai-agent-memory.

<sup>39.</sup> Ibid.

<sup>40.</sup> Yifeng He et al., "Security of AI Agents," arXiv, June 20, 2024. https://arxiv.org/html/2406.08689v2; Cole Stryker, "What are the components of AI agents?," IBM, March 10, 2025. https://www.ibm.com/think/topics/components-of-ai-agents.

<sup>41.</sup> Talha Zeeshan et al., "Large Language Model Based Multi-Agent System Augmented Complex Event Processing Pipeline for Internet of Multimedia Things," arXiv, Jan. 3, 2025. https://arxiv.org/html/2501.00906v2.

Cole Stryker, "Agentic AI: 4 reasons why it's the next big thing in AI research," IBM, Oct. 11, 2024. https://www.ibm.com/think/insights/agentic-ai; Adrian Bridgwater, "Okay AI, Solo Kagent Is An Agentic AI Framework For Kubernetes," *Forbes*, March 17, 2025. https://www.forbes.com/sites/adrianbridgwater/2025/03/17/okay-aisolo-kagent-is-an-agentic-ai-framework-for-kubernetes; Elizabeth Wallace, "How Agentic AI is Changing Decision-Making," CD Insights, March 8, 2025. https://www. clouddatainsights.com/how-agentic-ai-is-changing-decision-making.



As visualized in Figure 2, vulnerabilities can emerge across multiple layers of the agent lifecycle—from insecure training data and compromised model interfaces to attacks on the agents themselves. These risks, if left unaddressed, can embed insecure behaviors into models or outputs and trigger downstream cybersecurity consequences that expose users, systems, and software supply chains.





Source: Jessica Ji et al., "Cybersecurity Risks of Al-Generated Code," Center for Security and Emerging Technology, Nov. 2024. https://cset.georgetown.edu/publication/cybersecurity-risks-of-ai-generated-code. qtd. in Christopher Jablonski, "Al software supply chain risks prompt new corporate diligence," Zscaler, Jan. 13, 2025. https://www.zscaler.com/cxorevolutionaries/insights/ai-software-supply-chain-risks-prompt-newcorporate-diligence.

To further contextualize how agents work, it is important to be familiar with the seven main types of agents that have emerged in AI research and development.<sup>43</sup> Each category represents a varying level of autonomy, complexity, and adaptability:

**1. Simple Reflex Agents.** These agents represent the most basic form of agent, as they operate based on predefined condition-action rules only.<sup>44</sup> These types of agents are typically used in systems like keyword-based spam filters, where emails are assigned a label (spam or not spam) based on a predetermined rule or list of keywords.<sup>45</sup>

**2. Model-Based Reflex Agents.** Building on the simple reflex agent foundation, model-based reflex agents maintain an internal state, allowing them to adapt actions based on historical context or data, which is similar to how smart thermostats adjust temperature based on past patterns.<sup>46</sup>

**3. Goal-Based Agents.** These agents introduce a layer of intentionality, selecting actions based on whether they can help fulfill a defined objective.<sup>47</sup> A travel booking agent that books flights and coordinates lodging accommodations with minimal human input would fall into this third category of agents.<sup>48</sup>







<sup>43. &</sup>quot;Agents in Artificial Intelligence," GeeksforGeeks, June 5, 2023. https://www.geeksforgeeks.org/agents-artificial-intelligence; Douglas B. Laney, "Understanding And Preparing For The 7 Levels Of AI Agents," *Forbes*, Jan. 3, 2025. https://www.forbes.com/sites/douglaslaney/2025/01/03/understanding-and-preparing-for-the-seven-levels-of-ai-agents.

- 47. Ibid.
- 48. Ibid.

<sup>44.</sup> Gutowska. https://www.ibm.com/think/topics/ai-agents.

<sup>45.</sup> Ibid.

<sup>46.</sup> Ibid.



**4. Utility-Based Agents.** Utility-based agents take this further by weighing possible outcomes by determining which course of action is likely most beneficial, such as optimizing delivery routes to save time, value, or fuel.<sup>49</sup>

**5. Learning Agents.** These agents extend beyond the fixed strategies in the aforementioned four categories by continuously updating their approach based on feedback and new data.<sup>50</sup> These systems are capable of refining their performance over time, such as an AI system that personalizes lesson plans based on a given student's behavior and progress.<sup>51</sup>

**6. Multi-Agent Systems.** In more complex environments, multi-agent systems bring together different agents that work cooperatively or competitively to complete shared tasks, such as coordinating supply chain logistics.<sup>52</sup>

**7. Hierarchical Agents.** Hierarchical agents structure decision-making across various levels, delegating sub-tasks and managing dependencies in a way that mirrors organizational workflows.<sup>53</sup>

Many of today's leading agents, including Google's Project Astra, OpenAI's Operator, and CrewAI, reflect a growing trend: the emergence of general-purpose systems designed for flexible use across diverse environments and industries.<sup>54</sup>

As agentic AI matures, efforts to establish cybersecurity, interoperability, and governance standards are already underway.<sup>55</sup> These include initiatives like the novel Multi-Agent Environment, Security, Threat, Risk and Outcome (MAESTRO) threat modeling framework; sandboxing and permissioning strategies; and increased attention toward memory constraints and data boundaries.<sup>56</sup> Understanding the historical roots, agentic infrastructure stack layers, and practical distinctions between existing types of agents is essential for making sense of today's evolving developments and anticipating the cybersecurity and policy considerations that lie ahead.

#### **Cybersecurity Benefits of Al Agents**

Thanks to their enhanced autonomy, advanced reasoning, and capacity for continuous self-improvement, AI agents are already being deployed to streamline customer service workflows, assist with preliminary legal research, and automate data entry.<sup>57</sup> Though each of these business applications is significant, AI agents

#### R Street Policy Study No. 325 May 2025









<sup>49.</sup> Ibid.

<sup>50.</sup> Ibid.

<sup>51.</sup> Ibid.

<sup>52. &</sup>quot;Agents in Artificial Intelligence." https://www.geeksforgeeks.org/agents-artificial-intelligence.

<sup>53.</sup> Ibid.

<sup>54.</sup> Mark Purdy, "What is Agentic AI, and How Will It Change Work?," Harvard Business Review, Dec. 12, 2024. https://hbr.org/2024/12/what-is-agentic-ai-and-how-will-it-change-work; Heaven. https://www.technologyreview.com/2024/05/14/1092407/googles-astra-is-its-first-ai-for-everything-agent; Bhavishya Pandit, "CrewAI: A Guide With Examples of Multi AI Agent Systems," DataCamp, Sept. 12, 2024. https://www.datacamp.com/tutorial/crew-ai.

<sup>55.</sup> Toner. https://cset.georgetown.edu/publication/through-the-chat-window-and-into-the-real-world-preparing-for-ai-agents.

<sup>56.</sup> Ken Huang, "Agentic AI Threat Modeling Framework: MAESTRO," Cloud Security Alliance, Feb. 6, 2025. https://cloudsecurityalliance.org/blog/2025/02/06/agentic-ai-threat-modeling-framework-maestro; Sam Adler, "Interoperable Agentic AI: Unlocking the Full Potential of AI Specialization," Tech Policy Press, Dec. 3, 2024. https:// www.techpolicy.press/interoperable-agentic-ai-unlocking-the-full-potential-of-ai-specialization; Stephen Weigand, "Cybersecurity in 2025: Agentic AI to change enterprise security and business operations in year ahead," SC Media, Jan. 9, 2025. https://www.scworld.com/feature/ai-to-change-enterprise-security-and-businessoperations-in-2025; Shubham Sharma, "The new paradigm: Architecting the data stack for AI agents," VentureBeat, Nov. 14, 2024. https://venturebeat.com/datainfrastructure/the-new-paradigm-architecting-the-data-stack-for-ai-agents.

<sup>57.</sup> Zhixuan Chu, "Professional Agents – Evolving Large Language Models into Autonomous Experts with Human-Level Competencies," arXiv, Feb. 6, 2024. https://arxiv. org/html/2402.03628v1; Yasar. https://www.techtarget.com/searchenterpriseai/definition/Al-agents.



in those contexts typically support administrative, routine, and well-structured tasks. In contrast, AI agents in cybersecurity are increasingly deployed as copilots to human analysts on the frontlines, actively engaging with unfolding incidents, making rapid decisions in unpredictable scenarios, and operating in high-stress environments.<sup>58</sup> In other words, AI agents are not only accelerating efficiency but also strengthening cyber resilience by autonomously performing tasks critical to continuous monitoring, vulnerability management, threat detection, incident response, and decision-making across the cyber workforce.<sup>59</sup>

#### **Continuous Attack Surface Monitoring and Vulnerability** Management

As AI and other emerging technologies continue to expand across cloud infrastructure, third-party platforms, Internet of Things (IoT) devices, and edge environments, the overall attack surface that individuals and organizations must continuously monitor and secure has become increasingly fragmented and challenging.<sup>60</sup> While edge computing allows for greater decentralization and localized processing, the growing interconnectivity between devices, applications, and services can also amplify existing vulnerabilities, expand the number of possible points of failure, introduce new threat vectors, and reduce overall visibility into risk, especially in environments where interoperability remains inconsistent or poorly managed.<sup>61</sup> This is particularly true as AI-enabled systems often interact with external APIs, open-source resources, and real-time data streams, many of which are difficult to track, vet, or fully control.<sup>62</sup>

Traditional vulnerability management approaches, which depend on periodic scans or scheduled patching, are not well-suited for today's rapidly evolving and distributed environments.<sup>63</sup> AI agents offer a more adaptive and continuous alternative. While still largely supervised by human analysts, AI agents are increasingly capable of autonomously assisting with key tasks, including mapping systems, identifying exposures, and prioritizing patches based on anticipated severity or business impact.<sup>64</sup> As AI develops over time, these systems may help accelerate or even automate portions of the patching process, such as recommending fixes, triggering rollbacks, or adjusting configurations based on real-time context.<sup>65</sup>

Recent examples show how agentic capabilities are already being explored in real-world cybersecurity use cases. For instance, in late 2024, Google's Project Zero and DeepMind claimed to be the first in the world to successfully use an AI agent to uncover a "previously unknown, zero-day, exploitable memory-safety R Street Policy Study No. 325 May 2025



While still largely supervised by human analysts, AI agents are increasingly capable of autonomously assisting with key tasks, including mapping systems, identifying exposures, and prioritizing patches based on anticipated severity or business impact.

<sup>58.</sup> Thomas Caldwell, "The Evolution Of AI Agents In The Third Wave Of AI," *Forbes*, Oct. 22, 2024. https://www.forbes.com/councils/forbestechcouncil/2024/10/22/the-evolution-of-ai-agents-in-the-third-wave-of-ai.

<sup>59.</sup> Weigand. https://www.scworld.com/feature/ai-to-change-enterprise-security-and-business-operations-in-2025.

<sup>60.</sup> Haiman Wong, "Securing the Future of AI at the Edge: An Overview of AI Compute Security," R Street Institute, July 16, 2024. https://www.rstreet.org/research/ securing-the-future-of-ai-at-the-edge-an-overview-of-ai-compute-security.

<sup>61.</sup> Pete Bartolik, "Edge Requires Interoperability," CIO, May 19, 2021. https://www.cio.com/article/191756/edge-requires-interoperability.html.

<sup>62.</sup> Charles Owen-Jackson, "How cyber criminals are compromising AI software supply chains," Security Intelligence, Sept. 6, 2024. https://securityintelligence.com/ articles/cyber-criminals-compromising-ai-software-supply-chains.

<sup>63. &</sup>quot;Al is Now Exploiting Known Vulnerabilities – And What You Can Do About It," Cloud Security Alliance, June 26, 2024. https://cloudsecurityalliance.org/ blog/2024/06/26/ai-is-now-exploiting-known-vulnerabilities-and-what-you-can-do-about-it.

<sup>64.</sup> Maria Korolov, "Al agents can find and exploit known vulnerabilities, study shows," CSO, July 2, 2024. https://www.csoonline.com/article/2512791/ai-agents-can-findand-exploit-known-vulnerabilities-study-shows.html.



vulnerability in widely used real-world software."<sup>66</sup> Other AI agents are being trained to autonomously simulate attacks on enterprise systems, effectively performing red-teaming exercises to identify and test for vulnerabilities before adversaries can exploit them.<sup>67</sup> Both of these examples reflect a broader shift in how vulnerability management can be approached. As attack surfaces continue to evolve, AI agents could be key not only for helping cyber defenders keep up with emerging threats but also for advancing the speed and scale at which new and amplified vulnerabilities can be identified, prioritized, and remediated.<sup>68</sup>

#### **Real-Time Threat Detection and Incident Response**

Although AI has already broadly demonstrated its value in cybersecurity through anomaly detection, natural language processing for threat intelligence, and the automation of repetitive or lower-level tasks, AI agents offer a promising leap forward.<sup>69</sup> Due to their modular architecture, memory, and capacity for goaloriented, multi-step execution, AI agents can continuously learn from evolving threat patterns, correlate disparate signals, and initiate the appropriate responses without direct human oversight.<sup>70</sup> This is especially useful in high-speed or highvolume environments, where even a small delay between detection and response can affect the trajectory of an emerging threat and the success of incident containment efforts.<sup>71</sup>

For example, in security operations centers (SOCs), AI agents can be deployed to monitor network traffic, flag anomalies, and trigger isolation protocols for systems that are suspected to be compromised.<sup>72</sup> A multi-agent setup could divide responsibilities between network monitoring, threat intelligence synthesis, and automated remediation.<sup>73</sup> These capabilities are already being demonstrated in enterprise settings with emerging security tools like Microsoft's Security Copilot agents, Simbian's SOC AI Agent, and DropZone AI's SOC Analyst, among others.<sup>74</sup> Once an intrusion is detected, these agents not only flag the emerging threat – they can initiate immediate responses, such as coordinating with firewalls or endpoint protection platforms to isolate affected nodes, notifying administrators, beginning system recovery procedures, or even a combination of these tasks.<sup>75</sup> This level of speed and coordination is particularly important because it means AI agents can help reduce the mean time to detect (MTTD) and mean time to



As attack surfaces continue to evolve, AI agents could be key not only for helping cyber defenders keep up with emerging threats but also for advancing the speed and scale at which new and amplified vulnerabilities can be identified, prioritized, and remediated.

<sup>66.</sup> Davey Winder, "Google Claims World First As AI Finds 0-Day Security Vulnerability," *Forbes*, Nov. 5, 2024. https://www.forbes.com/sites/daveywinder/2024/11/05/ google-claims-world-first-as-ai-finds-0-day-security-vulnerability.

<sup>67.</sup> Sarah Nagar and David Eaves, "An Agentic Shield? Using AI Agents to Enhance the Cybersecurity of Digital Public Infrastructure," New America, Dec. 19, 2024. https://www.newamerica.org/digital-impact-governance-initiative/blog/an-agentic-shield-using-ai.

<sup>68. &</sup>quot;The Role of Al in Attack Surface Management: Enhancing Cyber Defense," Cyble, Feb. 13, 2025. https://cyble.com/knowledge-hub/ai-attack-surface-management.

<sup>69. &</sup>quot;R Street Cybersecurity-Artificial Intelligence Working Group," R Street Institute, last accessed March 29, 2025. https://www.rstreet.org/home/our-issues/ cybersecurity-and-emerging-threats/cyber-ai-working-group.

<sup>70.</sup> Grant Gross, "Agentic Al: 6 Promising Use Cases for Business," CIO, Nov. 14, 2024. https://www.cio.com/article/3603856/agentic-ai-promising-use-cases-for-business.html.

Jim Routh, "Milliseconds Matter: Defending Against the Next Zero-Day Exploit," Cloud Security Alliance, March 14, 2022. https://cloudsecurityalliance.org/ blog/2022/03/14/milliseconds-matter-defending-against-the-next-zero-day-exploit.

<sup>72.</sup> Jimmy Astle, "Incorporating AI agents into SOC workflows," Red Canary, Jan. 16, 2025. https://redcanary.com/blog/threat-detection/ai-agents.

<sup>73.</sup> Antonella C. Garcia et al., "A Multi-Agent System for Addressing Cybersecurity Issues in Social Networks," CEUR Workshop Proceedings, Dec. 31, 2022. https://ceur-ws. org/Vol-3495/paper\_05.pdf.

<sup>74.</sup> Thomas Claburn, "AI agents swarm Microsoft Security Copilot," The Register, March 24, 2025. https://www.theregister.com/2025/03/24/microsoft\_security\_copilot\_ agents; Kevin Townsend, "Simbian Introduces LLM AI Agents to Supercharge Threat Hunting and Incident Response," Security Week, Oct. 10, 2024. https://www. securityweek.com/simbian-introduces-Ilm-ai-agents-to-supercharge-threat-hunting-and-incident-response; "Dropzone AI," last accessed March 29, 2025. https:// www.dropzone.ai.



respond (MTTR), both of which are crucial metrics in mitigating the scope and cost of cybersecurity incidents.<sup>76</sup>

#### **Augmented Decision-Making and Cyber Workforce Support**

In recent years, reports consistently underscore a persistent cyber workforce gap.<sup>77</sup> In 2025, the World Economic Forum reported a shortage of more than 4 million cyber professionals globally.<sup>78</sup> In the United States alone, the shortfall is estimated to be between 500,000 and 700,000 workers.<sup>79</sup> Moreover, existing cyber teams are facing increased workplace demands, with many professionals reporting unrelenting hours and workloads.<sup>80</sup> In fact, a recent survey reported that more than 80 percent of respondents experienced burnout.<sup>81</sup> This has caused many frontline cyber defenders and experts to consider not only leaving their positions, but leaving the industry as a whole, signaling a growing crisis that shows little sign of abating.<sup>82</sup> While many organizations are already using AI-enhanced tools to help streamline workflows and accelerate upskilling throughout their teams, the accelerating scope, speed, and sophistication of cyber threats often outpace these incremental gains.<sup>83</sup>

As a result, AI agents are quickly entering the cyber workforce, not as replacements for human analysts but rather as force-multiplying copilots.<sup>84</sup> Although the current generation of AI agents is still far from perfect, they are already adept at a variety of critical tasks, including tuning firewalls, reducing noise by deduplicating security alerts, classifying security alerts by severity, using telemetry thresholds and anomaly detection to enforce policy changes, and more.<sup>85</sup> In doing so, AI security copilots, such as Cisco's AI Assistant, CrowdStrike's Charlotte AI, Fortinet's Advisor, Trellix's WISE, and Google's Sec-PaLM and AI Workbench, are gaining traction to help organizations keep their SOCs adequately staffed and efficient to better contain threats.<sup>86</sup> Moreover, when armed with these AI security copilots, SOCs are seeing notable improvements in false-positive rates (up to 70 percent) while reducing manual triage by more than 40 hours a week.<sup>87</sup> These early successes demonstrate how AI agents are emerging as a technological solution that can help organizations save time and money while improving their cyber resilience and retaining their cyber defenders.<sup>88</sup> R Street Policy Study No. 325 May 2025



A recent survey reported that more than 80 percent of respondents experienced burnout. This has caused many frontline cyber defenders and experts to consider not only leaving their positions, but leaving the industry as a whole, signaling a growing crisis that shows little sign of abating.

76. Greg Zemlin, "MTTD and MTTR in Cybersecurity Incident Response," Wiz, Sept. 5, 2024. https://www.wiz.io/academy/mttd-and-mttr.

77. Michelle Meineke, "The cybersecurity industry has an urgent talent shortage. Here's how to plug the gap," World Economic Forum, April 28, 2024. https://www. weforum.org/stories/2024/04/cybersecurity-industry-talent-shortage-new-report.

78. "Bridging the Cyber Skills Gap," World Economic Forum, last accessed March 29, 2025. https://initiatives.weforum.org/bridging-the-cyber-skills-gap/home.

79. Sophia Fox-Sowell, "To expand cyber workforce, government must unfreeze hiring and target youth, experts told House committee," StateScoop, Feb. 5, 2025. https:// statescoop.com/cybersecurity-workforce-house-committee-homeland-security.

 Matt Kapko, "Are cybersecurity professionals OK?," Cybersecurity Dive, Aug. 7, 2024. https://www.cybersecuritydive.com/news/cyber-security-burnout-stressanxiety/723470.

81. Ibid.

 Carolyn Crist, "Skills shortage persists in cybersecurity despite decade of hiring," HR Dive, Oct. 16, 2024. https://www.hrdive.com/news/skills-shortage-persistsin-cybersecurity-despite-hiring/729988; Joe McKendrick, "Will AI take the wind out of cybersecurity job growth?," ZDNet, July 10, 2024. https://www.zdnet.com/ article/will-ai-take-the-wind-out-of-cybersecurity-job-growth; Daniel Pell, "AI Versus The Skills Gap," *Forbes*, April 17, 2024. https://www.forbes.com/councils/ forbestechcouncil/2024/04/17/ai-versus-the-skills-gap.

83. Ivan Belcic and Cole Stryker, "AI Agents in 2025: Expectations vs. reality," IBM, March 4, 2025. https://www.ibm.com/think/insights/ai-agents-2025-expectations-vs-reality.

- 84. Columbus. https://venturebeat.com/security/ai-copilots-cut-false-positives-and-burnout-in-overworked-socs.
- 85. Ibid.

86. Ibid.

<sup>88.</sup> Gary Grossman, "Onboarding the Al workforce: How digital agents will redefine work itself," VentureBeat, Sept. 29, 2024. https://venturebeat.com/ai/onboarding-theai-workforce-how-digital-agents-will-redefine-work-itself.



### **Cybersecurity Considerations and Potential Risks**

Al agents are proving to be powerful not only because of what they can do on their own, but also because of how effective they are at learning and adapting across digital environments based on new data or updated information. Unfortunately, the same capabilities that make AI agents impressive, such as their memory, autonomy, and reasoning, can also make them attractive targets for exploitation.<sup>89</sup>

Though there are different ways to conceptualize an AI agent's architecture, we organize the agentic infrastructure stack into four primary layers: perception, reasoning, action, and memory. Each layer corresponds to a critical stage in how data is collected, analyzed, applied, and refined throughout the AI agent lifecycle. Because each layer serves a distinct function within the AI agent's workflow, the risks and mitigation needs associated with them also differ between modules, shaping the cybersecurity considerations at each stage.

#### **Layer 1: Perception Module**

At this first layer, the agent is tasked with scanning and observing a given environment through sensors (i.e., cameras, data inputs) to provide it with foundational context, and that data is then transformed into a suitable format for processing.<sup>90</sup> Because the perception module relies on multiple data pipelines for its analysis, this layer could face a variety of data-specific security risks that would affect the data confidentiality and integrity of the agentic workflow. These attacks include—but are not limited to—adversarial data injection (also known as data poisoning) and AI model supply chain risks.

Adversarial data injection is one of the most prominent security risks against the perception layer of an agent workflow because it tampers with the model's integrity and the agent's ability to factually analyze the data points in its training.<sup>91</sup> For example, bad actors could seamlessly insert modifications that mislead vision models into incorrect characterizations and imprecise content classifications on behalf of the agent. In the case of image processing, bad actors could manipulate the image pixels, add extra noise to the image, or perform other types of perturbation that are difficult to notice both with the human eye and via Alenabled perception systems.<sup>92</sup>

While the manipulation of image pixel values to deceive the agent is a common adversarial data risk, researchers have found that even small-scale perturbations in a dataset can meaningfully affect an agent's learning process—causing it to misclassify inputs into either "a maliciously-chosen target class (in a targeted attack) or classes that are different from the ground truth."<sup>93</sup> These types of

R Street Policy Study No. 325 May 2025



Adversarial data injection is one of the most prominent security risks against the perception layer of an agent workflow because it tampers with the model's integrity and the agent's ability to factually analyze the data points in its training.

<sup>89.</sup> Matt Dangelo, "The rise of autonomous AI: How intelligent agents are redefining strategy, risk, & compliance," Thomson Reuters, March 3, 2025. https://www. thomsonreuters.com/en-us/posts/technology/autonomous-agentic-ai; Maria Korolov, "AI agents will transform business processes – and magnify risks," CIO, Aug. 21, 2024. https://www.cio.com/article/3489045/ai-agents-will-transform-business-processes-and-magnify-risks.html.

<sup>90.</sup> Haziqa Sajid, "What Are Al Agents, and How Do They Work?," Lakera Al, Nov. 13, 2024. https://www.lakera.ai/blog/what-are-ai-agents.

Forest McKee and David Noever, "Transparency Attacks: How Imperceptible Image Layers Can Fool AI Perception," arXiv, Jan. 29, 2024. https://arxiv.org/abs/2401.15817.
Ian J. Goodfellow et al., "Explaining and Harnessing Adversarial Examples," arXiv, Dec. 20, 2014. https://arxiv.org/abs/1412.6572; Battista Biggio and Fabio Roli, "Wild Patterns: Ten Years After the Rise of Adversarial Machine Learning," arXiv, Dec. 8, 2017. https://arxiv.org/abs/1712.03141.

Chaowei Xiao et al., "Generating Adversarial Examples with Adversarial Networks," arXiv, Feb. 14, 2019. https://arxiv.org/pdf/1801.02610; Alexey Kurakin et al., "Adversarial examples in the physical world," arXiv, July 8, 2016. https://arxiv.org/abs/1607.02533.



attacks are particularly challenging, as bad actors can execute them without having any direct access to the model architecture.

Such data poisoning methods may also "reorient" the agent's data analysis from one intended pattern, set by developers, to a malicious one, set by bad actors, by altering the training set's distribution or reshaping the data to align with adversarial objectives. For instance, in a backdoor attack, an adversary could deliberately modify the training data to introduce specific triggers that, when encountered, would cause the model to behave in a predetermined, often malicious way.

These kinds of cybersecurity risks are especially concerning at the perception layer because of its heavy reliance on state-of-the-art foundation models many of which are externally sourced—creating additional dependencies.<sup>94</sup> While these models are critical for enabling advanced agent performance, their integration also expands the agent's exposure to potential software supply chain vulnerabilities, particularly during the pre-training phase.

In fact, threat actors can exploit the decentralized nature of the AI and software supply chain by embedding malicious data within these foundational models at the pre-training phase. The nature of this attack depends on the target, which can range from data poisoning to weight poisoning, along with the method of label modification.<sup>95</sup> Both types of backdoor attacks can lead to compromised downstream performance in agentic systems.

**Figure 3** illustrates the backdoor attack process on pre-trained foundation models. In many cases, these attacks are difficult to detect, and, therefore, challenging to mitigate. This also increases the likelihood of risk transfer from the foundation model to the AI agent itself.<sup>96</sup> If transferred, the AI agent may inherit these vulnerabilities and carry them forward into deployment.

#### R Street Policy Study No. 325 May 2025



In a backdoor attack, an adversary could deliberately modify the training data to introduce specific triggers that, when encountered, would cause the model to behave in a predetermined, often malicious way.



#### Figure 3: Framework of Backdoor Attacks to Unified Foundation Models

Source: Zenghui Yuan et al., "Backdoor Attacks to Pre-trained Unified Foundation Models," arXiv, Feb. 23, 2023 https://arXiv.org/pdf/2302.09360

- 94. Lareina Yee et al., "Why agents are the next frontier of generative AI," McKinsey Digital, July 24, 2024. https://www.mckinsey.com/capabilities/mckinsey-digital/our-insights/why-agents-are-the-next-frontier-of-generative-ai; Laura French, "How AI coding assistants could be compromised via rules file," SC World, March 18, 2025. https://www.scworld.com/news/how-ai-coding-assistants-could-be-compromised-via-rules-file.
- 95. Yuan. https://arxiv.org/pdf/2302.09360.
- 96. Hao Wang et al., "Model Supply Chain Poisoning: Backdooring Pre-trained Models via Embedding Indistinguishability," OpenReview.net, Jan. 29, 2025. https://openreview.net/forum?id=VWQwwMxFht#discussion.



In addition to poisoning and backdoor attacks, recent reports have uncovered a vulnerability in the Safetensors conversion service that Hugging Face—a leading open-source platform used by developers to host, customize, and share pre-trained machine learning models—offers.<sup>97</sup> According to AI security firm HiddenLayer, attackers can "send malicious pull requests with attacker-controlled data from the Hugging Face service to any repository on the platform, as well as hijack any models that are submitted through the conversion service."<sup>98</sup> This vulnerability presents a heightened security risk, particularly given Hugging Face's role as a major hub for pre-trained models. In practice, such exploitation could allow threat actors to impersonate a chatbot and submit malicious inquiries, ranging from instructions on how to successfully conduct a money heist to even building a nuclear bomb or a bioweapon.

#### Layer 2: Reasoning Module

The second layer of the AI agent workflow is the reasoning module, which governs the agent's internal decision-making processes. At this stage, data collected in the perception module in layer 1 is interpreted and transformed into actionable outputs. The agent reviews and analyzes contextual information and may apply pre-learned heuristics, patterns, or logical ordering to generate a conclusion with the support of specialized hardware like graphic processing units or tensor processing units and model hosting environments.<sup>99</sup> For example, an agent might analyze network activity logs to determine whether a user request is legitimate or suspicious, drawing on historical behaviors and anomaly detection models to inform its decisions. Because the reasoning module plays a central role in analysis and judgment, vulnerabilities and bad cyber hygiene in this layer can lead to incorrect decisions or mischaracterizations, particularly if adversaries manipulate the signals or exploit vulnerabilities in the models or supporting infrastructure. Ultimately, this inaccuracy could undermine the end-user's trust in the agent's reliability and accuracy.

One of the most common process-level security risks at this stage is the exploitation of the model's underlying vulnerabilities. These flaws could stem from widely used AI frameworks like PyTorch, which plays a critical role in the reasoning layer by enabling developers to build, train, and fine-tune machine learning models. PyTorch is also commonly used for deep learning model development, inference processing, and optimization, making it a core component of many AI agent workflows.<sup>100</sup>

Security flaws can also arise from misconfigured libraries and insecure model hosting environments, especially those that allow user-generated model uploads without robust validation.<sup>101</sup> In early 2024, for example, researchers found that

R Street Policy Study No. 325 May 2025



For example, an agent might analyze network activity logs to determine whether a user request is legitimate or suspicious, drawing on historical behaviors and anomaly detection models to inform its decisions.

<sup>97.</sup> Ravie Lakshmanan, "New Hugging Face Vulnerability Exposes AI Models to Supply Chain Attacks," The Hacker News, Feb. 27, 2024. https://thehackernews.com/2024/02/new-hugging-face-vulnerability-exposes.html.

<sup>98. &</sup>quot;Hijacking Safetensors Conversion on Hugging Face," HiddenLayer Inc., Feb. 21, 2024. https://hiddenlayer.com/innovation-hub/silent-sabotage.

Alexander De Ridder, "How Intelligent Agents Use Knowledge Representation for Decision-Making," SmythOS, Jan. 31, 2025. https://smythos.com/ai-agents/aitutorials/intelligent-agents-and-knowledge-representation; Barbara Bickham, "The Role of Heuristics in Developing Smarter AI Systems," Trailyn Ventures, April 19, 2024. https://www.trailyn.com/the-role-of-heuristics-in-developing-smarter-ai-systems.

<sup>100. &</sup>quot;Getting Started," PyTorch, last accessed March 30, 2025. https://pytorch.org; Ashish Kurmi, "PyTorch Supply Chain Compromise," StepSecurity, Dec. 9, 2024. https:// www.stepsecurity.io/blog/pytorch-supply-chain-compromise.

<sup>101.</sup> Elizabeth Montalbano, "Hugging Face AI Platform Riddled with 100 Malicious Code-Execution Models," Dark Reading, Feb. 29, 2024. https://www.darkreading.com/ application-security/hugging-face-ai-platform-100-malicious-code-execution-models.



the Hugging Face platform was home to approximately 100 malicious machine learning models capable of depositing malicious code onto users' machines.<sup>102</sup>

Another category of attacks that could affect the agent's knowledge base are model exploitation attacks. Instead of targeting AI inputs directly, adversaries can attempt to probe the Al's internal logic to extract proprietary knowledge, internal decision pathways, or sensitive training data. There are three techniques that threat actors can use to exploit models. First, bad actors can attempt to extract personal identifiable information (PII) by reconstructing aspects of the training data.<sup>103</sup> This is also known as a model inversion attack. Another approach is blackbox extraction, where an attacker without direct access to the model's architecture or weights submits iterative queries to infer or replicate the model, leading to intellectual property theft or downstream adversarial attacks.<sup>104</sup> Finally, threat actors may attempt to jailbreak or probe the model's logic by crafting prompts that aim to trick the model into revealing its underlying structure or how the agent processes information. As these bad actors continuously refine their prompts, they can analyze outputs to map the agent's decision-making boundaries, identify vulnerabilities to exploit later, and improve hallucination tactics that degrade or mislead agent performance over time.<sup>105</sup>

#### Layer 3: Action Module

The third layer of the agentic workflow is the action module, which is responsible for translating the decision-making process made in layer 2 into real-world operations. Because this is the stage where actions are executed, even seemingly minor manipulations can lead to unintended—and potentially harmful— consequences. This makes the action module particularly sensitive to attacks that exploit an agent's ability to interface with external systems.

Bad actors could compromise this layer through a variety of vectors, including, but not limited to prompt injection, command hijacking, unauthorized access, privilege escalation, and vulnerabilities within API integrations. These risks underscore the importance of implementing stringent output validation and access controls at this layer.

In prompt injection attacks, threat actors feed malicious prompts into the agent with the goal of manipulating the agent to perform actions outside the scope of its intended purpose, like, leaking PII or generating malicious outputs.<sup>106</sup>

While prompt injection focuses on manipulating inputs to modify the agent's behavior, command hijacking takes this one step further by executing unauthorized commands based on the earlier inputs. One well-documented example is the "Imprompter" attack in which attackers manipulated AI systems using deceptive prompts to retrieve sensitive information like names, email



Bad actors could compromise this layer through a variety of vectors, including, but not limited to prompt injection, command hijacking, unauthorized access, privilege escalation, and vulnerabilities within API integrations.

<sup>102.</sup> Ibid.

<sup>103. &</sup>quot;Model inversion attacks," Michalsons, March 8, 2023. https://www.michalsons.com/blog/model-inversion-attacks-a-new-ai-security-risk/64427?utm\_source=chatgpt. com.

<sup>104.</sup> Mumtaz Fatima and Amy Chang, "Safeguarding AI: A Policymaker's Primer on Adversarial Machine Learning Threats," R Street Institute, March 20, 2024. https://www. rstreet.org/commentary/safeguarding-ai-a-policymakers-primer-on-adversarial-machine-learning-threats.

<sup>105.</sup> Rachneet Sachdeva et al., "Turning Logic Against Itself: Probing Model Defenses Through Contrastive Questions," arXiv, Jan. 3, 2025. https://arxiv.org/abs/2501.01872.

<sup>106.</sup> Matthew Kosinki and Amber Forrest, "What is a prompt injection attack?," IBM, March 26, 2024. https://www.ibm.com/think/topics/prompt-injection.



addresses, and payment details.<sup>107</sup> This incident illustrates how seemingly minor vulnerabilities and exploitation techniques can be leveraged to compromise user privacy and disrupt system integrity.

Another security risk at this layer involves unauthorized access through privilege escalation. Because AI agents operate across different structured layers of execution and interact with a range of systems, such as data stores, applications, end users, and pre-trained foundation models, any security gap could allow threat actors to move laterally within the agentic workflow. This privilege escalation could result in misconfigured role-based access controls, which allow bad actors to access restricted modalities and the foundations of the underlying model or to prompt the model to perform unauthorized actions.<sup>108</sup> In a recent example, researchers from Palo Alto Networks identified two vulnerabilities in Google's Vertex AI platform that could have enabled bad actors to escalate privileges and extract models.<sup>109</sup>

As mentioned earlier in this section, the action module is also vulnerable to insecure execution permissions that stem from weak access controls. These security flaws could grant bad actors the ability to poison models or trick them into executing deceptive or harmful requests under the guise of legitimate prompts.

Since this layer is the primary interface between AI agents and external connections, APIs are a critical, and often undersecured, vector for exploitation. If communication channels are not secured with robust safeguards, bad actors could intercept and manipulate different model requests and responses through man-in-the-middle attacks or re-prompt previous AI queries to gain access to restricted areas.<sup>110</sup> For example, a ransomware group could exploit a vulnerable fraud-detection API to alter transaction data while evading detection.

More broadly, API vulnerabilities can also stem from failures in endpoint detection, missing or improperly validated API keys, or lax token authentication.<sup>111</sup> These weaknesses can open the door for threat actors to bypass restrictions, achieve privilege escalation, and manipulate the agent's behavior.<sup>112</sup> In doing so, threat actors could also execute adversarial prompt injection or other forms of input manipulation that lead to harmful outputs.<sup>113</sup>

In addition to APIs, AI agents often rely on third-party services for data analysis. This adds yet another layer of cybersecurity risk: Compromised datasets, insecure API dependencies, or insufficient monitoring can allow threat actors to tamper with

#### R Street Policy Study No. 325 May 2025



Because AI agents operate across different structured layers of execution and interact with a range of systems, such as data stores, applications, end users, and pre-trained foundation models, any security gap could allow threat actors to move laterally within the agentic workflow.

110. Wenqi Sun et al., "Clustering Mobile Apps based on Design and Manufacturing Genre," IEEE Xplore, December 2020. https://ieeexplore.ieee.org/document/9344944.

<sup>107.</sup> Matt Burgess, "This Prompt Can Make an AI Chatbot Identify and Extract Personal Details From Your Chats," Wired, Oct. 17, 2024. https://www.wired.com/story/aiimprompter-malware-llm/?utm\_source=chatgpt.com.

<sup>108.</sup> Ximeng Liu et al., "Privacy and Security Issues in Deep Learning: A Survey," IEEE Access, December 2020. https://www.researchgate.net/publication/347639649\_ Privacy\_and\_Security\_Issues\_in\_Deep\_Learning\_A\_Survey.

<sup>109.</sup> Ofir Balassiano and Ofir Shaty, "ModeLeak: Privilege Escalation to LLM Model Exfiltration in Vertex AI," Unit 42 Palo Alto Networks, Nov. 12, 2024. https://unit42. paloaltonetworks.com/privilege-escalation-Ilm-model-exfil-vertex-ai.

<sup>111. &</sup>quot;Wallarm Releases 2025 API ThreatStats Report, Revealing APIs are the Predominant Attack Surface," Wallarm, Jan. 29, 2025. https://www.wallarm.com/pressreleases/wallarm-releases-2025-api-threatstats-report.

<sup>112. &</sup>quot;API Security's Role in Responsible AI Deployment," Wallarm, Jan. 21, 2025. https://lab.wallarm.com/api-securitys-role-in-responsible-ai-deployment.

<sup>113. &</sup>quot;Discover and Protect Generative AI APIs," Traceable AI, last accessed March 30, 2025. https://www.traceable.ai/securing-gen-ai-apis.



Al agent operations without detection.<sup>114</sup> If any part of the third-party software supply chain is compromised, the agent's performance and trustworthiness could be severely degraded and more difficult to immediately remediate and restore.

Many agents are also deployed in cloud environments, which carry their own set of cyber risks and potential vulnerabilities. Weak configuration settings, such as unsigned code releases or poorly managed access controls, leave systems exposed to software supply chain attacks and malicious model updates.<sup>115</sup> Furthermore, security flaws in cloud storage buckets can also leak sensitive model parameters, inviting intellectual property theft and other similar types of adversarial attacks.<sup>116</sup>

#### Layer 4: Memory Module

The fourth and final layer of the agentic workflow is the memory module, which is responsible for retaining context across tasks, storing relevant data, and informing future decisions based on past interactions.<sup>117</sup> This module distinguishes AI agents from other AI models or LLM-based tools, which typically operate within a single session or a query window.<sup>118</sup> By enabling long-term context awareness, learning persistence, and memory-driven adaptability, the memory module facilitates the AI agent's continuous self-improvement capabilities over time.<sup>119</sup>

One of the primary cybersecurity risks that can occur at this layer is memory tampering or corruption, where threat actors manipulate stored memory to distort an agent's understanding or to introduce incorrect historical data.<sup>120</sup> This can occur through learning stream poisoning (maliciously modifying real-time inputs that are then retained as memory) or through unauthorized edits to at-rest memory databases.<sup>121</sup> These attacks could degrade AI agent performance or subtly influence future actions toward harmful outputs.<sup>122</sup>

Relatedly, unauthorized data retention is another cybersecurity risk within this layer.<sup>123</sup> When unauthorized data retention occurs, AI agents remember data or information they were not supposed to retain, either because they inadvertently collected data outside of its intended use case or learning scope, retained data longer than permitted, or failed to delete it when instructed R Street Policy Study No. 325 May 2025

# ERROR

One of the primary cybersecurity risks that can occur at this layer is memory tampering or corruption, where threat actors manipulate stored memory to distort an agent's understanding or to introduce incorrect historical data.

116. Emmanuel Ok, "Addressing Security Challenges in Al-Driven Cloud Platforms: Risks and Mitigation Strategies," ResearchGate, February 2025. https://www. researchgate.net/publication/388997486\_Addressing\_Security\_Challenges\_in\_Al-Driven\_Cloud\_Platforms\_Risks\_and\_Mitigation\_Strategies.

<sup>114.</sup> Bryan McNaught, "API Protection for AI Factories: The First Step to AI Security," F5, Dec. 19, 2024. https://www.f5.com/company/blog/api-security-for-ai-factories; Yared Gudeta et al., "Securing the Future: How AI Gateways Protect AI Agent Systems in the Era of Generative AI," Databricks, Nov. 13, 2024. https://www.databricks. com/blog/ai-gateways-secure-ai-agent-systems.

<sup>115.</sup> Favour Efeoghene, "How to Secure CI/CD Pipelines Against Supply Chain Attacks," StartUp Growth Guide, July 24, 2024. https://startupgrowthguide.com/how-to-secure-ci-cd-pipelines-against-supply-chain-attacks.

<sup>117. &</sup>quot;Al Agents," Nvidia Glossary, last accessed March 30, 2025. https://www.nvidia.com/en-us/glossary/ai-agents.

<sup>118.</sup> Jeffrey Yang Fan Chiang et al., "Why Are Web AI Agents More Vulnerable Than Standalone LLMs? A Security Analysis," arXiv, Feb. 27, 2025. https://arxiv.org/ html/2502.20383v1.

<sup>119.</sup> Xun Jiang et al., "Long Term Memory: The Foundation of Al Self-Evolution," arXiv, Oct. 21, 2024. https://arxiv.org/html/2410.15665v1.

<sup>120.</sup> Chad DeChant, "Episodic memory in AI agents poses risks that should be studied and mitigated," arXiv, Jan. 20, 2025. https://arxiv.org html/2501.11739v1?ref=community. heartcount.io.

<sup>121.</sup> Jaimin Patel, "Secure Al Agents by Design with Al Runtime Security," Palo Alto Networks, Jan. 23, 2025. https://www.paloaltonetworks.com/blog/network-security/ secure-ai-agents-by-design-ai-runtime-security; Symantec Threat Hunter Team, "Al: Advent of Agents Open New Possibilities for Attackers," Symantec Enterprise Blogs, March 13, 2025. https://www.security.com/threat-intelligence/ai-agent-attacks.

<sup>122.</sup> Ibid.

<sup>123.</sup> Zhaorun Chen et al., "AgentPoison: Red-teaming LLM Agents via Poisoning Memory or Knowledge Bases," arXiv, July 17, 2024. https://arxiv.org/abs/2407.12784.



to.<sup>124</sup> This can lead to compliance violations in relation to existing privacy laws or user terms and conditions and unintentionally expose sensitive user information.<sup>125</sup> Even otherwise well-configured AI agents can face these cybersecurity challenges if existing memory governance guardrails are missing or improperly implemented.<sup>126</sup>

What makes the memory module particularly significant is its recursive relationship with the earlier three layers of the agentic lifecycle.<sup>127</sup> If the data lifecycle were conceptualized as a circle, this fourth layer effectively closes the loop, meaning that any vulnerabilities or risks introduced earlier in the process, such as poisoned data or training processes and faulty reasoning, may be not only preserved but reinforced over time.<sup>128</sup> For example, if adversarial data is ingested through the perception layer and not flagged as corrupted, the memory module could preserve it as a trusted input, continuing to apply that context to influence future reasoning processes and actions.<sup>129</sup> Similarly, if an attack manipulates an AI agent's logic at the reasoning module, the tasks it finishes in the action module may be remembered as valid precedent.<sup>130</sup>

In this way, memory does not simply inform an AI agent's future performance—it can also carry forward mistakes and risks from its past. Without strong protections and best practices for ensuring data accuracy, implementing retention boundaries, and managing memory, the memory module can become both a repository of insights and a source of cascading cyber vulnerabilities and risks.<sup>131</sup>

#### Anticipating Needs, Solutions, and Responsibilities

To advance our cyber preparedness in this rapidly unfolding era of AI agents, we must adopt a proactive, balanced, and adaptable strategy. Striking the right balance means encouraging policymakers, end-users, and developers to recognize and fully leverage the benefits that AI agents offer while anticipating and addressing the cybersecurity risks they may amplify or introduce.

The following recommendations highlight the policy needs, emerging technological solutions, and responsible design and deployment strategies that are best suited for supporting and guiding the advancement of AI agents.

#### **Policy Needs**

#### Establish Voluntary, Sector-Specific Guidelines for Human–Agent Collaboration

The White House should direct federal agencies, such as the National Institute of Standards and Technology (NIST), the Department of Labor, and relevant industry

#### R Street Policy Study No. 325 May 2025



Without strong protections and best practices for ensuring data accuracy, implementing retention boundaries, and managing memory, the memory module can become both a repository of insights and a source of cascading cyber vulnerabilities and risks.



<sup>124.</sup> Rashi Shrivastava, "The Prompt: Privacy Risks 'Haunt' AI Agents," *Forbes*, March 11, 2025. https://www.forbes.com/sites/rashishrivastava/2025/03/11/the-prompt-privacy-risks-haunt-ai-agents; David Ruiz, "New AI 'agents' could hold people for ransom in 2025," Malwarebytes, Feb. 4, 2025. https://www.malwarebytes.com/blog/ news/2025/02/new-ai-agents-could-hold-people-for-ransom-in-2025.

<sup>125.</sup> Erich Kron, "Five privacy concerns around agentic AI," SC Media, Feb. 19, 2025. https://www.scworld.com/perspective/five-privacy-concerns-around-agentic-ai.

<sup>126.</sup> Daniel Berrick, "Minding Mindful Machines: AI Agents and Data Protection Considerations," Future of Privacy Forum, Feb. 5, 2025. https://fpf.org/blog/mindingmindful-machines-ai-agents-and-data-protection-considerations.

<sup>127.</sup> Christian Vasquez, "Cybersecurity pros are preparing for a new adversary: AI agents," Fortune, Feb. 18, 2025. https://fortune.com/2025/02/18/cybersecurity-pros-are-preparing-for-a-new-adversary-ai-agents.

<sup>128.</sup> Ibid.

<sup>129.</sup> DeChant. https://arxiv.org/html/2501.11739v1?ref=community.heartcount.io.

<sup>131.</sup> Berrick. https://fpf.org/blog/minding-mindful-machines-ai-agents-and-data-protection-considerations.



regulators, to develop voluntary, sector-specific guidelines that support secure, transparent, and human-centered agentic deployments.<sup>132</sup>

Rather than prescribing a rigid, one-size-fits-all mandate, these guidelines should encourage organizations—including AI laboratories, private-sector companies, and universities—to define tailored human—agent interaction frameworks. These frameworks should clarify when agents may be deployed, under what conditions they may act autonomously, whether they are permitted to learn independently, when human oversight is required, how responsibility is assigned in the event of failures, and what protocols exist for detecting, escalating, and correcting errors. The goal is to ensure that agents support—not replace—human decision-making and talent, especially in sensitive fields of work like healthcare and national security.<sup>133</sup>

Given the potential of AI agents to reshape workforce dynamics and drive an "AI talent revolution," these guidelines should also promote organizational readiness for human–agent collaboration by offering recommendations for redesigning jobs and reskilling or upskilling current employees.<sup>134</sup> These guidelines should also highlight responsible deployment strategies, such as incremental rollouts, permissioning boundaries, and real-time escalation protocols, tailored to teambased environments where agents and humans are likely to interact.<sup>135</sup> While federal agencies can help provide baseline guidance and highlight priorities, public–private partnerships will be essential to translating these principles into practice. Consortia like the Partnership on AI, Microsoft's Tech-Labor Partnership with the American Federation of Labor and Congress of Industrial Organizations (AFL-CIO), and Cisco's AI-Enabled ICT Workforce Consortium provide early momentum and serve as leading examples of cross-sector collaboration.<sup>136</sup>

### Expand and Facilitate Information Sharing and Multi-Stakeholder Collaboration on Evolving Agentic Risks

Given their potential to serve as force multipliers for offensive, defensive, and adversarial cyber operations, AI agents require equally coordinated and dynamic strategies for timely, cross-sector information sharing about emerging agentic risks, observed unintended agentic behaviors, deployment challenges, and successful risk mitigation strategies.<sup>137</sup> Specifically, the White House should direct federal agencies like the Cybersecurity and Infrastructure Security Agency to collaborate with sector-specific regulatory bodies and industry stakeholders to expand information-sharing forums and develop publicly available software tools and resources for testing and evaluating agentic security and performance.<sup>138</sup>

#### R Street Policy Study No. 325 May 2025



The goal is to ensure that agents support—not replace—human decision-making and talent, especially in sensitive fields of work like healthcare and national security.



<sup>132.</sup> Haiman Wong and Brandon Pugh, "Key Cybersecurity and AI Policy Priorities for Trump's Second Administration and the 119th Congress," R Street Institute, Jan. 6, 2025. https://www.rstreet.org/research/key-cybersecurity-and-ai-policy-priorities-for-trumps-second-administration-and-the-119th-congress.

<sup>133.</sup> Shana Lynch, "Predictions for AI in 2025: Collaborative Agents, AI Skepticism, and New Risks," Stanford University Human-Centered Artificial Intelligence, Dec. 23, 2024. https://hai.stanford.edu/news/predictions-for-ai-in-2025-collaborative-agents-ai-skepticism-and-new-risks.

<sup>134.</sup> Whiting, "The rise of 'AI agents': What they are and how to manage the risks." https://www.weforum.org/stories/2024/07/what-is-an-ai-agent-experts-explain. 135. Ibid.

<sup>136.</sup> Rebecca Finlay, "Sharing AI Mistakes: Partnership on AI's Rebecca Finlay," MIT Sloan Management Review, Nov. 12, 2024. https://sloanreview.mit.edu/audio/sharingai-mistakes-partnership-on-ais-rebecca-finlay; "AFL-CIO and Microsoft Announce New Tech-Labor Partnership on AI and the Future of the Workforce," AFL-CIO, Dec. 11, 2023. https://aflcio.org/press/releases/afl-cio-and-microsoft-announce-new-tech-labor-partnership-ai-and-future-workforce; "Leading Companies Launch Consortium to Address AI's Impact on the Technology Workforce," Cisco, April 4, 2024. https://newsroom.cisco.com/c/r/newsroom/en/us/a/y2024/m04/leadingcompanies-launch-consortium-to-address-ai-impact-on-the-technology-workforce.html.

<sup>137.</sup> Rosa Merced, "Exploring the potential of AI agents in combatting cyber threats," Outshift by Cisco, Feb. 18, 2025. https://outshift.cisco.com/blog/exploring-thepotential-of-ai-agents-in-combatting-cyber-threats.

<sup>138.</sup> Wong and Pugh. https://www.rstreet.org/research/key-cybersecurity-and-ai-policy-priorities-for-trumps-second-administration-and-the-119th-congress.



These efforts should emphasize use-case-specific transparency, such as anonymized incident reports and adversarial testing results, to accelerate collective learning and cyber preparedness.<sup>139</sup>

### Prioritize Investments in Public–Private Partnerships for Advancing Agentic Security and Evaluation

Congress should prioritize investments in continued research and development initiatives aimed at strengthening the cybersecurity posture of AI agents across their full lifecycle.<sup>140</sup> While private companies naturally have strong incentives to secure their own products and services, many agentic risks—such as model hijacking, memory poisoning, and emergent multi-agent behavior—can cut across proprietary systems and lack clearly defined ownership or liability.<sup>141</sup>

In cases where agentic risks are cross-cutting and infrastructural, the federal government can play a limited but essential role by supporting foundational research and ensuring that key findings are made publicly available to foster broader coordination and informed risk mitigation for agentic developers and researchers. This can include funding adversarial testing, agent-specific risk modeling, and resilience evaluations focused on architectural features like memory integrity and autonomous decision-making.<sup>142</sup> While NIST has already taken early steps to evaluate hijacking risks in AI agents, scaling this work will require additional investments in competitive grants, interdisciplinary research hubs, and public–private partnerships that accelerate knowledge-sharing and innovation across sectors.<sup>143</sup>

#### **Emerging Technological Solutions**

## Advance and Apply Automated Moving Target Defense (AMTD) Capabilities to Disrupt Evolving Exploitation Pathways

AMTD are systems designed to continuously alter a system's attack surface by shifting IP addresses, memory allocations, or control paths to deliberately complicate adversarial reconnaissance efforts and reduce system predictability.<sup>144</sup> When paired with the autonomous and continuous self-improvement capabilities of AI agents, AMTD systems could rotate access privileges, shuffle API endpoints, or re-randomize internal configurations to limit the persistence of adversarial probing attempts or prompt injection attacks.<sup>145</sup> These techniques are expected to be particularly useful in edge computing environments, where agents will need to remain flexible and responsive while operating across distributed and often interdependent digital environments.<sup>146</sup>



- 139. Ghose. https://scet.berkeley.edu/the-next-next-big-thing-agentic-ais-opportunities-and-risks.
- 140. Wong and Pugh. https://www.rstreet.org/research/key-cybersecurity-and-ai-policy-priorities-for-trumps-second-administration-and-the-119th-congress.

142. Ghose. https://scet.berkeley.edu/the-next-next-big-thing-agentic-ais-opportunities-and-risks.

<sup>141.</sup> Phaedra Boinodiris and Jon Parker, "The evolving ethics and governance landscape of agentic AI," IBM, March 22, 2025. https://www.ibm.com/think/insights/ethicsgovernance-agentic-ai.

<sup>143. &</sup>quot;Technical Blog: Strengthening AI Agent Hijacking Evaluations," National Institute of Standards and Technology, Jan. 17, 2025. https://www.nist.gov/news-events/ news/2025/01/technical-blog-strengthening-ai-agent-hijacking-evaluations.

<sup>144.</sup> Mark Loman, "Pioneering Automated Moving Target Defense (AMTD)," Sophos, Oct. 19, 2023. https://news.sophos.com/en-us/2023/10/19/pioneering-automatedmoving-target-defense-amtd.

<sup>145.</sup> Sailik Sengupta, "Moving Target Defense: A Symbiotic Framework for AI & Security," Association for Computing Machinery Digital Library, May 8, 2017. https://dl.acm. org/doi/pdf/10.5555/3091125.3091473.

<sup>146. &</sup>quot;Fujitsu develops world's first multi-AI agent security technology to protect against vulnerabilities and new threats," Fujitsu Limited, Dec. 12, 2024. https://www. fujitsu.com/global/about/resources/news/press-releases/2024/1212-01.html.



#### Adapt and Implement Hallucination Detection Tools for Continuous Agentic **Security Monitoring**

Originally developed to improve the quality control capabilities and accuracy of LLM outputs, hallucination detection tools are now quickly being repurposed for agentic security to identify reasoning flaws and gaps, anomalous or suspicious behavior, and low-confidence outputs before they reach the action module.<sup>147</sup> Emerging hallucination detection tools operate by using internal consistency checks, multi-source fact validation, and prompt-response tracking to monitor for misalignment, especially under adversarial or high-stress conditions.<sup>148</sup> In the context of AI agents, these hallucination detection tools are already proving successful at revealing compromised memory recalls, model drifts, and inference anomalies, all of which are essential to helping developers identify vulnerabilities before they can be exploited by threat actors.<sup>149</sup>

#### Develop and Adopt Agent Identifiers and Traceability Tools to Improve Oversight

To improve explainability and oversight, AI researchers and developers should continue creating identification infrastructure and persistent tools capable of tracking and logging the full arc of agent activity, including an agent's initial and expanding data-collection strategies; third-party dependencies and tool applications; completed tasks; pending actions; reasoning logic streams; and memory recall.<sup>150</sup> This approach builds on existing strategies of embedded provenance tracking and digital auditing that are used to enable real-time behavioral analysis, versioncontrol tracking, validating software supply chain dependencies and end-user contributions, and post-incident forensics.<sup>151</sup> The ongoing development of agent IDs-which log instance-specific information such as the interacting system and interaction history-represent a practical foundation for model development, much like how a serial number is used to trace a product and its history.<sup>152</sup> These identifiers could also help track the origins, certifications, and performance of an AI system.<sup>153</sup> This increased visibility and improved agentic explainability would equip cyber practitioners, AI researchers and developers, and end-users with the dynamic intelligence needed to detect suspicious activity, make attributions to culpable threat actors, and conduct incident investigations.<sup>154</sup>

#### **Responsible Design and Deployment Strategies for Developers** and End Users

#### **Maintain Strong Cyber Hygiene Best Practices**

Cybersecurity fundamentals remain essential, but they must now extend into each layer of the agentic infrastructure stack. Core cyber best practices, such as robust identity and access management, secure API usage, and zero-trust

147. Alice Gomstyn and Alexandra Jonker, "New ethics risks courtesy of AI agents? Researchers are on the case," IBM, Dec. 23, 2024. https://www.ibm.com/think/insights/ ai-agent-ethics.

148. Ibid.

149. Ibid.

151. Ibid. 152. Ibid.

- 153. Ibid.
- 154. Ibid.

#### R Street Policy Study No. 325 May 2025







<sup>150.</sup> Afek Shamir, "AI Agents: Why We Should Strategize on Governance," Tech Policy Press, Dec. 19, 2024. https://www.techpolicy.press/ai-agents-why-we-shouldstrategize-on-governance.



architectures should be implemented when designing new AI agents or adapting existing agents for customized applications.<sup>155</sup> These measures can help reduce the risk of cascading failures across the agent's workflow and maintain system integrity.<sup>156</sup> As agents operate with increasing autonomy, maintaining strong cyber hygiene best practices remains the first line of defense.

#### Implement Boundaries for Agentic Scope and Autonomy to Well-Defined Tasks

Before designing and deploying agents, organizations, developers, and end-users should clearly define and document the agent's intended scope, purpose, task parameters, and tiered permission levels.<sup>157</sup> This practice is imperative because it reduces the risk of unintended consequences, improves agentic reliability, and mitigates the potential of agentic drift, misalignment, and overreach.<sup>158</sup> Organizations should also update internal policies that outline acceptable agentic design and deployment strategies and train employees on how agents can be responsibly used for their individual roles and responsibilities.<sup>159</sup>

#### Deploy AI Agents Incrementally with Built-In Evaluation and Rollback Protocols

Responsible and secure deployment of AI agents requires iterative testing and continuous monitoring.<sup>160</sup> Organizations and end users should always aim to introduce agents gradually, starting with sandboxed environments and escalating through staged pilot programs, while applying regular red-teaming, running automated stress tests, and preparing rollback protocols at each phase of agentic deployment.<sup>161</sup> This allows developers and end users alike to continuously observe an AI agent's real-world behavior, identify surface novel risks in real-time, calibrate agent performance over time, and course correct effectively if an agent starts completing harmful tasks or expanding beyond its defined purpose and scope of work.<sup>162</sup>

#### Conclusion

The rise of AI agents signals a marked shift in how emerging technologies interact with, interpret, and influence our digital world. Increasingly described as the "third wave" of AI innovation, AI agents represent a departure from passive models that rely on continuous human oversight and intervention.<sup>163</sup> With the ability to act autonomously, reason, and learn through experience, AI agents are poised to redefine the contours of human–machine collaboration.

These agentic advancements also introduce complex governance questions and cybersecurity challenges. As AI agents take on more decision-making

158. Ibid.





<sup>155.</sup> Canon, "Why cyber hygiene remains critical in the era of Al-driven threats," CSO, Feb. 26, 2025. https://www.csoonline.com/article/3820782/why-cyber-hygieneremains-critical-in-the-era-of-ai-driven-threats.html; "It May be Time to Review Your Cyber Hygiene," Security Magazine, May 8, 2024. https://www.securitymagazine. com/articles/100646-it-may-be-time-to-review-your-cyber-hygiene.

**<sup>156.</sup>** Korolov. https://www.cio.com/article/3489045/ai-agents-will-transform-business-processes-and-magnify-risks.html.

<sup>157. &</sup>quot;Guidelines for defining goals and instructions for AI agent," Webex Help Center, Feb. 7, 2025. https://help.webex.com/en-us/article/nelkmxk/Guidelines-for-defininggoals-and-instructions-for-AI-agent.

<sup>159.</sup> Lynch. https://hai.stanford.edu/news/predictions-for-ai-in-2025-collaborative-agents-ai-skepticism-and-new-risks.

<sup>160.</sup> Xuhui Zhou et al., "HAICOSYSTEM: An Ecosystem for Sandboxing Safety Risks in Human-Al Interactions," arXiv, Sept. 24, 2024. https://arxiv.org/abs/2409.16427. 161. Ibid.

<sup>163.</sup> Caldwell. https://www.forbes.com/councils/forbestechcouncil/2024/10/22/the-evolution-of-ai-agents-in-the-third-wave-of-ai.



roles and responsibilities, their actions may increasingly reflect the assumptions, priorities, and constraints embedded into their underlying models. Ensuring these systems are secure, reliable, and aligned with welldefined objectives requires more than technical measures alone; it demands coordinated efforts from stakeholders across industry, government, and the public. Our regulatory frameworks and risk management strategies must also evolve in parallel with AI agents to enable effective oversight, support responsible design and deployment, and reduce emerging risks.

While we cannot perfectly anticipate all the long-term impacts AI agents will have on how we work, learn, and live, what remains certain is that our governance and cybersecurity responsibilities are only beginning. With balanced, flexible safeguards in place, agentic systems can be deployed in ways that maximize their benefits while proactively mitigating the cybersecurity risks they may introduce or amplify. Beyond the recommendations outlined in this study, there are still many opportunities for future research and development, ranging from scalable audit mechanisms and real-time monitoring tools to infrastructure-agnostic safeguards and standards for agent-to-agent interactions.<sup>164</sup> Moreover, as AI agents are developed and deployed across borders, the need for voluntary and shared global governance norms for acceptable use and permissible development strategies will only grow.<sup>165</sup> Ultimately, the broader imperative is not simply to keep pace with emerging technologies but to guide and shape their trajectory, ensuring that they augment human talent and skills, reinforce America's technological leadership and economic competitiveness, and remain grounded in our founding values.<sup>166</sup>

#### R Street Policy Study No. 325 May 2025



Ultimately, the broader imperative is not simply to keep pace with emerging technologies but to guide and shape their trajectory, ensuring that they augment human talent and skills, reinforce America's technological leadership and economic competitiveness, and remain grounded in our founding values.

#### About the Authors

**Haiman Wong** is a resident fellow in the Cybersecurity and Emerging Threats team at the R Street Institute. Her research examines the intersection of cybersecurity and emerging technologies, including Al, edge computing, and connected vehicles.

**Tiffany Saade** is a master's candidate in the Stanford Ford Dorsey Program in International Policy, with a concentration in Cyber Policy and Security. Her research explores the intersection of AI and the evolving cybersecurity threat landscape, focusing on how AI models can be exploited in offensive cyber operations and strategies to secure them against misuse.

<sup>164. &</sup>quot;The Rise of Al Agents: Transforming Citizen Development and User Experiences," Windows Forum, Feb. 27, 2025. https://windowsforum.com/threads/the-rise-of-aiagents-transforming-citizen-development-and-user-experiences.353983/?amp=1.

<sup>165.</sup> Betsy Morris, "Beyond Intelligence: The Impact of AI Agents," Stanford University McCoy Family Center for Ethics in Society, Aug. 16, 2024. https://ethicsinsociety. stanford.edu/news/beyond-intelligence-impact-advanced-ai-agents.

<sup>166.</sup> Richard Reisman and Richard Whitt, "New Perspectives on AI Agentiality and Democracy: 'Whom Does It Serve?," Tech Policy Press, Dec. 6, 2024. https://www. techpolicy.press/new-perspectives-on-ai-agentiality-and-democracy-whom-does-it-serve; Eleonore Fournier-Tombs, "An Ethical Grey Zone: AI Agents in Political Deliberations," Carnegie Council for Ethics in International Affairs, Nov. 13, 2024. https://www.carnegiecouncil.org/media/article/ethical-grey-zone-ai-agents-politicaldeliberation.