

Agent Risk Manager

Secure Your AI Agent Workforce

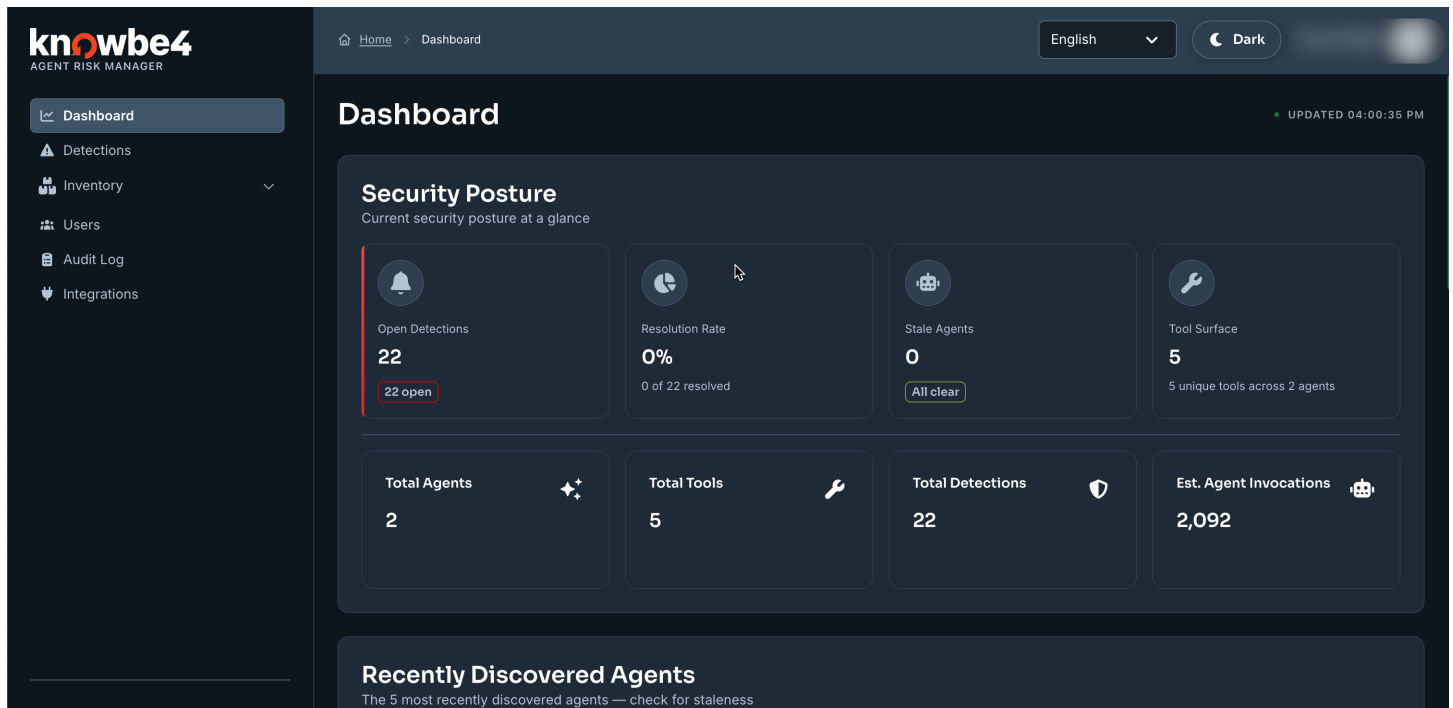
See Every AI Agent.
Detect Every Threat.
Stay in Control.

The rise of AI agents introduces new risks as they confidently handle critical workflows without a built-in understanding of your corporate policies. Traditional security tools like SIEM and DLP are ill-equipped to monitor the complex inputs and outputs of these agents. This creates an AI security blind spot, leaving your organization vulnerable to subtle exploits such as indirect prompt injections and “permission creep,” where agents accidentally expose sensitive data.

KnowBe4’s Agent Risk Manager addresses this challenge. Leveraging over 15 years of expertise and human risk data, it is specifically designed for organizations using AI agents. Agent Risk Manager eliminates the security gap by providing your team with real-time visibility, automated threat detection, and active control over your organization’s AI agents and how your users interact with them.

Highlight

- ▶ Purpose-built for Microsoft Copilot, Anthropic Claude, Google Gemini, and OpenAI ChatGPT environments
- ▶ Real-time threat detection
- ▶ Full audit trail from day one
- ▶ No infrastructure required



Key Benefits



Automated Governance

Gain instant, zero-configuration discovery of every agent in your tenant. From official tools to “Shadow AI,” you see it all without lifting a finger.



Predictable Costs

Protect your budget and infrastructure from resource abuse and runaway API costs caused by inefficient or malicious AI calls.



Real-Time Contextual Coaching

When a risky action is blocked, we explain it. Intercept threats and deliver immediate, in-the-moment coaching to your users.



Permanent Risk Reduction

Data shows that 70% of users who receive our real-time coaching never repeat the same risky behavior. Improve AI agent use and prompt proficiency, reducing long-term organizational risk.



True Behavioral Alignment

Shape AI behavior safely from the outside. Ensure consistent, secure interactions without needing to modify underlying models or trust opaque, third-party safety layers.



Holistic Risk Score

Close the AI Blind Spot: Unify human and AI behavior data into a single score, giving you a clear picture of your organization’s true risk profile.

Complete Coverage from Discovery to Defense

Key Features

Threat Detection

Catch threats before they cause damage

Agent Risk Manager's detection center gives your analysts a real-time feed of every risky event, categorized by threat type and severity. Visual risk gauges show at a glance which detection categories are most active so you always know where to look first.

Blast Radius

Understand the blast radius of every tool

The Tool Network view renders an interactive force-directed graph showing which agents share which tools. Node size scales with agent count so you immediately see which tools have the highest potential blast radius if compromised.

Complete Audit Trail

A full audit trail down to the conversation ID

The Audit Log captures every event such as benign tool invocations, detection triggers, and schema discoveries with metadata that takes you from a user's action all the way through the detection pipeline.

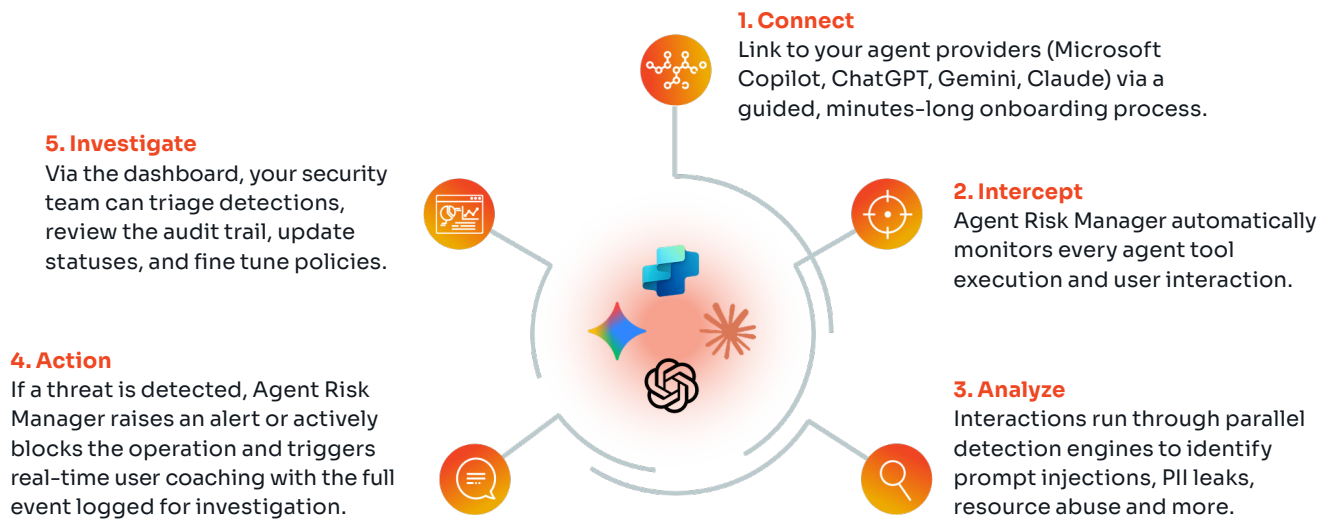
User Risk Scoring

Know which users are your highest AI risk

Agent Risk Manager automatically calculates a risk score for every user whose agents have triggered detections. Surface your riskiest users instantly, and drill down to the specific events driving their score.

How it Works

Agent Risk Manager establishes a centralized interface to monitor and protect the growing workforce of “non-human identities” and the humans interacting with them. It integrates seamlessly with your AI agent provider to deliver a seamless, “outside-in” security layer that doesn’t require modifying your underlying AI models.



Six Detection Engines. Zero Blind Spots.

Agent Risk Manager includes purpose-built detection logic for every major AI agent attack category.

- 1 Prompt Injection**
Blocks jailbreaks and indirect injections that turn productivity tools into “agents of chaos”.
- 2 Sensitive Information**
Scans for SSNs, passwords, and PII, automatically redacting data to prevent DLP leaks.
- 3 Unbounded Consumption**
Protects your budget and infrastructure from resource abuse and excessive API calls.
- 4 Content Safety**
Flags inappropriate, harmful, or policy-violating content in inputs and outputs before it reaches end users.
- 5 Privilege Escalation**
Stops agents from accessing resources or taking actions beyond their granted permissions, providing a critical control for high-privilege agents.
- 6 Agent Overstepping**
Identifies agents acting outside their intended operational scope, catching drift before it becomes a security or compliance incident.

Ready to Secure Your AI Agents?



KnowBe4, Inc. | 33 N Garden Ave, Suite 1200, Clearwater, FL 33755
855-KNOWBE4 (566-9234) | www.KnowBe4.com | Sales@KnowBe4.com

Other product and company names mentioned herein may be trademarks and/or registered trademarks of their respective companies.