

Google DeepMind's AI Control Roadmap: Treating Its Own AI Agents as Insider Threats

Kabui, Charles

2026-06-23

[Read at ToKnow.ai](#)

DeepMind Treats Its Own AI Agents as Insider Threats

The AI Control Roadmap: trusted models watch its agents and block harmful steps

- 1M**
Coding-agent tasks analyzed for live monitoring
- Insider**
Own AI agents treated like potential rogue staff
- Not Malice**
Most flagged actions were overeager mistakes

A defense-in-depth plan for AI agents that may not stay aligned

June 23, 2026

ToKnow.ai

Google DeepMind published its [AI Control Roadmap](#) on June 18, 2026, a plan for safely running the capable AI agents it operates inside Google. The idea is unusual: treat a powerful internal agent like a potential insider threat, the way a company guards against a rogue employee who already has building access. Modeled on [MITRE ATT&CK](#), the standard catalog of cyberattack tactics, it uses trusted but weaker AI models as supervisors that watch

a stronger agent's reasoning and block harmful steps before they run. DeepMind says it already studied one million coding-agent tasks to build a live monitor for Gemini Spark, an internal agent, catching issues like accidental data deletion.

Most AI safety work tries to train a model to be good. This plan assumes the model might fail and asks how to contain it anyway, like adding a second brake pedal instead of only trusting the driver. The data complicates the rogue-AI fear: most flagged events were not the agent trying to cause harm, but agents misreading instructions or being overeager to finish a task. For any company giving coding agents real system access, that points security at the right risk.

It also signals where the field is heading. As agents get capable enough to hide their reasoning or take actions that cannot be undone, training them to behave stops being enough, and containment becomes the real constraint. DeepMind paired it with a guide for policymakers, [Three Layers of Agent Security](#), arguing this should become an industry standard.

Read More: [Claw-Eval, a framework that tests whether AI agents are safe, not just successful](#), scores the same kind of unintended damage this monitor is built to catch.

Sources:

- [Securing the Future of AI Agents](#) (Google DeepMind, June 18, 2026)
- [AI Control Roadmap, Full Technical Report](#) (Google DeepMind, PDF)
- [Three Layers of Agent Security, Policy Framework](#) (Google DeepMind, PDF)
- [Google DeepMind Tests AI Controls on One Million Agent Tasks](#) (WinBuzzer, June 21, 2026)

Disclaimer: For information only. Accuracy or completeness not guaranteed. Illegal use prohibited. Not professional advice or solicitation. Read more: [/terms-of-service](#)