# vijil diamond

## Why Vijil

Vijil helps you ship AI agents to enterprises 4x faster while lowering operational risks. Vijil provides a layer of trust between agent development frameworks and agent runtime platforms to make AI agents more reliable, secure, and safe for enterprises from development to operations, and back.

## Why You

**You balance innovation and risk:** You're the senior leader at a fast-moving company who wants to use AI agents in production quickly. You want to see evidence of reliability, security, and safety throughout development and operations.

**You use special agents:** Your AI team is building or buying custom domain-specific agents using (open or closed) LLMs with access to confidential data and restricted tools.

## What Diamond Does

- Quantifies risks of your custom agent with bespoke tests
- Shortens time-to-trust™ (time-to-value at lower risk)
- Lowers cost of governance, risk mitigation, compliance

## Why Diamond is Better

- **Comprehensive**: Includes continuously evolving tests based on our 7-level taxonomy of trustworthiness
- **Customizable**: Generates bespoke tests tailored to your agent specification, user personas, and org policies
- **Policy-driven**: Includes tests based on international, federal, and state regulations, and industry standards
- **Auditable**: Uniquely, verifies agent code (SAST) and behavior (DAST) inside a trusted exec environment (TEE)
- **Multi-level**: Tests multi-agent systems, multi-turn interactions, and multi-component (LLM, knowledge base, tools using MCP or A2A) agents
- **Fast**: Produces Vijil Trust Score and Report in minutes
- **Actionable**: Recommends test-driven remediation

## Root of Trust

- Runs as a SOC2 Type II service on Google Cloud
- Runs as an agent, in a TEE using confidential computing to ensure integrity of testing and verification

> Vijil made it easy to test our AI assistant thoroughly, throughout its development, so we could deploy it into production 4x sooner.
> – Benedikt Klinglmayer, **Autonoma Cloud**



Trust Score



Trust Report

Gartner COOL VENDOR 2025

CBINSIGHTS AI 100 2025