

BACKSLASH

The Ultimate 2026

# Vibe Coding Security Buyer's Guide



## Table of Contents:

<b>Introduction</b>	<b>3</b>
<b>The 5 maturity Levels of Coding Agents Security</b>	<b>4</b>
<b>Visibility</b>	<b>5</b>
<b>Security Posture</b>	<b>6</b>
<b>Governance</b>	<b>7</b>
<b>Active Tracing</b>	<b>8</b>
<b>Active Protection</b>	<b>9</b>
<b>Summary and Call to Action</b>	<b>10</b>
<b>RFI Template</b>	<b>11</b>

# Introduction: Securing the Vibe

AI adoption within enterprises is quick, but is the most accelerated among software development teams. Tools like Cursor, Windsurf, and Claude Code are reshaping how software is built, shifting development from deliberate, line-by-line coding to intent-driven, agent-assisted and even agent-driven creation.

This new mode of building is commonly referred to as Vibe Coding.

Security leaders face a new mandate: enable AI-native development safely, without slowing velocity or fragmenting developer flow.

This requires a shift from gatekeeping to guardrails. Instead of blocking changes after the fact, security must be embedded directly into the agentic workflows that generate code, execute commands, and deploy systems.

The developer persona is evolving into an orchestrator. Builders now manage agents, MCP servers, prompts, and tools that span IDEs, terminals, CI pipelines, and cloud infrastructure. At the same time, non-traditional builders like product managers are using tools such as Cursor and Claude Code to push minor code changes to production, dramatically expanding the attack surface.

## The new agentic threat landscape introduces new risks:



Autonomous agents executing shell commands, modifying infrastructure, or deploying code without traditional human approval.



MCP servers acting as high-privilege APIs between LLMs and internal systems such as file systems, databases, and services.



Agent internet access and skills that dynamically fetch unverified libraries, browse third-party sites, and introduce supply chain and prompt poisoning risks.



Massive increases in code volume and variability that overwhelm traditional SAST and SCA approaches.



Chaining of components and identities that makes it difficult to assign the right permissions (e.g., agent connecting to an MPC that connects to another agent, and so forth).

- ✖ Securing Vibe Coding is not about retrofitting old tools. It requires a purpose-built, AI-native security approach that covers the entire agentic coding stack.

# The 5 maturity Levels of Coding Agents Security

Organizations that adopt Vibe Coding progress through five levels of security maturity. Each level answers a critical question, moving from awareness to real-time defense.





## LEVEL 1

# Visibility: Which agents are running?

Security teams must start with unified visibility across the entire Vibe Coding ecosystem. Modern developers no longer work with a single tool or a single identity. They routinely use multiple coding agents together, across different environments, often authenticated with different identities at the same time. This reality fundamentally changes the security model.

**This visibility must include a real-time, continuously updated inventory of:**



Local coding agents and AI-native IDEs such as Cursor, Claude Code, and Windsurf.



Remote agents such as CoPilot, Codex, and Devin.



Coding agents in all forms, including CLI tools, AI-native IDEs, and IDE extensions.



The user identity each coding agent is authenticated with, for example whether a personal Gmail account is connected to Cursor or an enterprise SSO is used.



Active MCP servers and the external systems they are connected to.



Which developers and teams are using which agents, and which combinations of agents are used together in practice.



- Without centralized visibility, security teams lack a clear understanding of where AI is generating code, how agents are configured, which identities they operate under, and which internal or external systems they can access. Visibility must be continuous rather than point-in-time, correlated across tools and identities, and presented through a single unified, centralized view that reflects how developers actually work today.



## LEVEL 2

# Security Posture: Do we have critical risks?

Once agents are visible, organizations must move beyond inventory and develop a deep understanding of risk across the AI coding environment. Visibility answers where agents exist. Security posture answers how dangerous they can be and why. This requires continuously evaluating the capabilities, configurations, and privileges of every agent, MCP server, and integration in use.

**Any posture management process must include the identification of agents and MCP servers that have unsecure, overly permissive, or poorly governed configurations, including but not limited to:**



Unrestricted shell or terminal access that allows agents to execute arbitrary commands.



Autonomous execution, deployment, or infrastructure-modifying capabilities without human approval or policy enforcement.



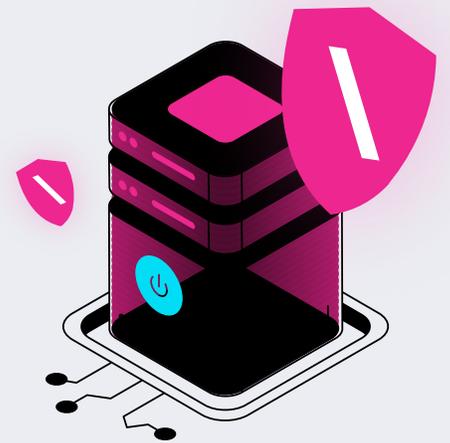
Broad access to file systems, internal networks, cloud resources, source code repositories, or identity and access management services.



Elevated privileges that exceed what is required for the agent's intended purpose.



Unvetted or misconfigured MCP servers that can introduce data leakage, prompt injection, context poisoning, or unintended access to internal organizational services, APIs, and sensitive data.



- Security posture management focuses on understanding effective power, not just presence. It correlates what agents can do, what they are connected to, and which identities they operate under. This allows security teams to compare risk across developers, tools, teams, and environments, identify dangerous patterns early, and prioritize remediation before misconfigurations or excessive privileges turn into security incidents.



### LEVEL 3

## Governance: Which policies should be applied?

Governance introduces enforceable guardrails across agentic workflows. At this level, organizations move from understanding risk to actively controlling it through centralized, policy-driven enforcement. Governance ensures that coding agents operate within clearly defined boundaries, aligned with organizational security, compliance, and engineering standards, without relying on manual review or developer judgment alone.

**At this stage, organizations can define and enforce policies such as:**

-  Agents may access workspace files but are explicitly prohibited from accessing sensitive files such as .env files, secrets, or credentials.
-  Agents may interact with internal or external services through MCPs in read-only mode, but are not allowed to modify, delete, or create data.
-  Agents may have access to local files but are restricted from outbound internet access or external network communication.
-  Agents may propose infrastructure or configuration changes, but cannot apply or deploy them directly without human approval.
-  Agents may install third-party dependencies, but only after those dependencies are scanned and approved by the organization's existing SCA tools and security processes.
-  Developers may install MCPs, but those MCPs must not have known vulnerabilities and must meet defined security baselines.
-  Developers may install MCPs only if they are official, trusted, and authenticated using secure mechanisms such as OAuth rather than static credentials.

-  Effective governance standardizes and hardens coding agents, AI-native IDEs, prompts, and MCP servers across all teams and environments. Policies are enforced automatically and consistently, reducing human error while preserving developer velocity. The result is a secure-by-default agentic environment where guardrails are invisible to developers, but risk is actively constrained at every step.





## LEVEL 4

# Activity Tracing: What has happened?

When vulnerabilities, misconfigurations, or incidents appear, organizations need full end-to-end traceability across agentic activity. Modern AI-driven development requires the ability to reconstruct exactly what happened, why it happened, and which entities were involved, across tools, agents, identities, and systems.



Activity tracing links real-world production outcomes back to:

-  The specific agent that generated an executable artifact such as the deploy binary or script.
-  The agent that executed operational commands, for example running a curl command at a specific point in time.
-  The exact prompt, instructions, and conversational context that led the agent to take that action.
-  The user, identity, and inferred agent intent that initiated or approved the action.
-  The MCP servers, tools, plugins, and external systems involved in the execution path.
-  The policies that allowed, restricted, or blocked actions, such as a policy that prevented an MCP server from being installed.

-  This creates a complete, auditable chain from intent to execution. With this level of traceability, security and engineering teams can respond to incidents quickly, perform accurate root cause analysis, troubleshoot a wide range of failures, and produce reliable compliance and audit reports for AI-generated code, actions, and decisions across the entire development lifecycle.



## LEVEL 5

# Active Protection: Can we stop bad things from happening?

The highest level of maturity is real-time prevention. At this stage, organizations move beyond detection, investigation, and governance into active, inline protection of agentic workflows. Security controls operate in real time, intercepting and stopping malicious, risky, or policy-violating actions before they are executed or cause damage.

**Organizations at this level can automatically identify and block:**

-  Dangerous or destructive shell commands before they are run.
-  Unauthorized access to sensitive data or attempts at data exfiltration, which can then be obfuscated.
-  Context poisoning, prompt injection attacks, and malicious interactions through MCPs or agent tools.
-  Third-party risks introduced through malicious or compromised rules, skills, plugins, or MCP servers.
-  Risky deployments, configuration changes, or releases before they reach production environments.
-  Insider misuse of agents or MCP access to obtain data or capabilities beyond approved authorization.



- Active protection transforms security from monitoring and alerting into live, continuous enforcement. Instead of investigating incidents after the fact, security teams prevent them altogether, while developers retain speed, autonomy, and the ability to work naturally with AI agents. This creates a security posture where innovation is enabled by default, and risk is constrained in real time.

## Summary

Backslash takes a completely fresh and comprehensive approach to vibe coding security. Instead of retrofitting legacy shift-left tools with AI assistants, or offering generic AI security disconnected from real developer behavior, Backslash is built for how modern software is actually written today. Developers work with multiple agents, tools, identities, and workflows at the same time, across local and remote environments. Backslash addresses this reality head-on by securing the full vibe coding ecosystem, not just isolated tools or point-in-time scans.

The Backslash platform unifies visibility, posture management, governance, traceability, and real-time prevention across both the ecosystem layer and the code layer. By combining deep code context with live insight into agents, IDEs, MCP servers, identities, and actions, Backslash delivers continuous protection across the entire vibe coding stack, from intent to execution to production.



## Next Steps

- ✓ If your organization is adopting AI coding agents, now is the moment to secure them intentionally. Backslash provides security and engineering teams with the visibility, control, and active protection needed to enable vibe coding at scale, without sacrificing developer speed or autonomy. Talk to our team to move from reactive security to real-time protection across your AI-driven development workflows.

[GET A LIVE DEMO](#)

# RFI Template



[Download RFI Excel Template](#)

## Inventory & Visibility

Can the solution identify local MCP clients?	
Can the solution identify remote MCP servers?	
Can the solution identify MCPs invoked via IDE plugins?	
Can the solution identify MCPs invoked via local daemon processes?	
Can the solution identify MCPs invoked via remote endpoints?	
Does the solution provide visibility into MCP parameters?	
Does the solution provide visibility into LLM usage inside IDEs?	
Can the solution classify MCPs by purpose?	
Can the solution classify MCPs by source such as local or remote?	
Can the solution classify MCPs by trust level such as signed?	
Can the solution capture workspace rules?	
Can the solution capture sub-agents?	
Can the solution capture agent extensions?	
Can the solution capture agent skills?	
Does the solution provide visibility into model selection per tool?	
Does the solution log the model name or ID?	
Does the solution log the model provider?	
Does the solution log the model version?	

## Coding agent Compatibility

Compatibility with Cursor?	
Compatibility with Claude code?	
Compatibility with Github Copilot?	
Compatibility with Google Antigravity?	
Compatibility with Gemini Code Assist?	
Compatibility with Windsurf?	
Compatibility with VSCode?	
Compatibility with JetBrains?	
Compatibility with OpenAI Codex?	
Compatibility with CodeGPT?	
Compatibility with Devin?	
Compatibility with Amazon Kiro?	

## Risk Posture

Can the solution detect agents connected with private user account?	
Can the solution detect malicious MCP behavior?	
Can the solution detect insecure dependencies?	
Can the solution detect suspicious network usage?	
Can the solution detect privilege escalation patterns?	
Can the solution detect unsigned MCP packages?	
Can the solution detect dependency confusion risks?	
Can the solution detect pulling models or tools from untrusted locations?	

## Usage Tracing & Observability

Does the solution provide visibility into MCP tool calls?	
Can the solution capture all tool calls including parameters and file paths?	
Can the solution report on the frequency of tool invocation?	
Can the solution report on tool chains and sequences?	
Can the solution detect file write or delete actions?	
Can the solution detect network access?	
Can the solution detect shell invocation?	
Can the solution detect repository modification?	
Can the solution export audit logs to SIEM or SOAR?	
Can the solution capture developer prompts?	
Can the solution capture agent messages?	
Does the solution provide visibility into agent mode actions?	
Does the solution provide visibility into autorun steps?	
Does the solution provide visibility into autonomous tool execution?	

## Guardrails & Active Protection

Can the solution distribute rule patterns?	
Can the solution distribute ignore file patterns?	
Can the solution prevent Data leakage?	
Can the solution prevent prompt injections?	
Does the solution support policy-based restrictions on tools?	
Does the solution support policy-based restrictions on prompts?	
Does the solution support policy-based restrictions on code generation types?	
Can the solution detect insecure code suggestions?	
Can the solution block unauthorized MCPs?	
Can the solution block remote or unknown model downloads?	
Can the solution block malicious tool invocation?	
Can the solution block known-risk LLM actions?	
Does the solution support model whitelisting?	
Does the solution support MCP whitelisting?	
Does the solution support tool whitelisting?	
Does the solution support endpoint whitelisting?	
Can the solution inject and enforce enterprise prompts in IDEs?	
Can the solution inject and enforce enterprise prompts in agents?	