

CSAI Foundation | Cloud Security Alliance

MCP Security Crisis: Systemic Design Flaws in AI Agent Infrastructure

Threat Landscape, Incident History, and Defense Guidance for Enterprise Deployments

2026-05-04

 AI-assisted Rapid Research



© 2026 Cloud Security Alliance. Some rights reserved.

You may download, store, display, view, print, redistribute, and link to this document in its original, unmodified form, provided that attribution to the Cloud Security Alliance is maintained and all trademark and copyright notices remain intact.

This document may not be modified or altered. You may quote portions of the document as permitted by the Fair Use provisions of the United States Copyright Act, provided that attribution is given to the Cloud Security Alliance.

This document may be shared on professional and social media platforms in its original form with attribution.

This document was generated with AI assistance and has not undergone official CSA review and approval processes.

Key Takeaways

The Model Context Protocol (MCP), Anthropic's open standard for connecting AI agents to external tools and data sources, has emerged as one of the most rapidly weaponized attack surfaces in agentic AI deployments, given the breadth of its supply chain exposure. A systemic architectural flaw disclosed in April 2026 by OX Security exposes an estimated 200,000 vulnerable instances across a supply chain encompassing more than 150 million package downloads [1][2]. Critically, this is not a coding error in any single product – it is a design default embedded in every official MCP SDK, propagated into downstream projects that trusted the reference implementation. Anthropic has confirmed the behavior is intentional and declined to modify the protocol architecture, leaving remediation responsibility with individual downstream developers [3].

- The MCP STDIO transport executes operating system commands without sanitization or validation, enabling remote code execution on any vulnerable host.
- Beyond the STDIO flaw, the protocol specification provides no native defenses against tool poisoning, rug pull attacks, or cross-server tool shadowing – a set of interconnected attack classes now documented in real-world incidents.
- As of May 2026, at least seven confirmed high- or critical-severity CVEs span major MCP-integrated platforms including MCP Inspector, LiteLLM, Cursor IDE, LibreChat, and Windsurf, with additional vulnerabilities under active tracking [4].
- The MCP authorization specification defines an OAuth 2.1 framework but explicitly marks authorization optional, leaving MCP servers exposed without authentication – a July 2025 internet scan identified at least 1,862 publicly accessible instances responding to unauthenticated requests [20][23].
- Enterprises should treat every MCP server as an untrusted third party, apply zero trust controls at the tool integration layer, and establish ongoing MCP governance programs rather than treating the situation as a one-time patch cycle.

Background

Anthropic introduced the Model Context Protocol in late 2024 as a universal connector standard for AI systems, giving models a structured interface to interact with file systems, databases, APIs, code repositories, and enterprise SaaS applications. Adoption accelerated rapidly. By mid-2025, the MCP ecosystem had accumulated more than 150 million package downloads and thousands of community-built MCP servers spanning developer tooling, productivity applications, financial platforms, and enterprise integrations [6][5]. Claude Desktop, Cursor IDE, Windsurf, and dozens of third-party AI coding assistants shipped with native MCP support, making MCP the most widely adopted standardized integration layer in the early agentic AI ecosystem.

MCP defines two primary transport mechanisms: the Streamable HTTP transport, intended for remote server connections, and the STDIO (standard input/output) transport, the default mechanism for local tool integrations. The STDIO transport requires no network infrastructure and is the mechanism illustrated in Anthropic's official documentation and reference SDKs. Its simplicity contributed to making it the most widely deployed integration path in the ecosystem – which is precisely what makes the architectural flaw discovered by OX Security so far-reaching.

Early security research had already surfaced concerns about MCP's threat model before the OX Security disclosure. Simon Willison documented in April 2025 how tool descriptions visible to the AI model but not displayed in user interfaces could carry hidden adversarial instructions, establishing prompt injection as a structural risk in MCP's design [7]. Invariant Labs demonstrated in a responsible disclosure how a single malicious MCP server could weaponize adjacent trusted servers through cross-server tool shadowing [8]. These warnings were largely framed as implementation challenges. The OX Security research reframed the problem: the risks were not incidental to poor implementation choices, but inherent to what the specification required by default, across every officially supported programming language.

CSA's own prior research notes examined specific components of the emerging MCP risk landscape, including the critical Flowise remote code execution vulnerability (CVE disclosed April 2026) that required emergency patching across hundreds of production AI agent deployments [18], and the systemic STDIO command injection vulnerability disclosed by OX Security [19]. The analysis presented here synthesizes the full threat landscape across the MCP incident history and provides forward-looking strategic guidance for security programs managing agentic AI risk.

Security Analysis

The STUDIO Design Flaw

OX Security's April 2026 research, titled "The Mother of All AI Supply Chains," identified a command execution vulnerability embedded in Anthropic's official MCP SDKs for Python, TypeScript, Java, and Rust [1]. The STUDIO transport processes incoming configuration by passing parameters directly to the host operating system's shell without input sanitization or validation. Any process command passed to the interface executes on the host system regardless of whether it initializes a valid MCP server. An attacker who can influence an MCP configuration file – through a compromised software repository, a malicious package registry entry, or a social engineering scenario – can achieve arbitrary code execution on the target machine.

OX researchers catalogued four distinct exploitation families, each demonstrated against live production systems [2][14]. Unauthenticated command injection via direct STUDIO configuration allows an attacker to deliver a malicious MCP configuration to a system that has no session authentication, exploiting scenarios where configuration files are fetched from external sources. Authenticated command injection via MCP STUDIO affects deployments where valid users configure MCP servers from attacker-controlled packages. Hardening-bypass variants circumvent common mitigations such as shell restrictions or process sandboxing. A fourth family chains MCP command execution with secondary privilege escalation techniques to achieve persistence or lateral movement across the target environment.

The affected ecosystem is broad. OX Security identified vulnerabilities in more than ten downstream projects including LiteLLM, LangChain, LangFlow, Flowise, LettaAI, LangBot, LibreChat, DocsGPT, Bisheng, and Windsurf [4]. The research produced over thirty responsible disclosures and catalogued at least seven confirmed high- or critical-severity CVEs, with additional vulnerabilities under active tracking. Several projects shipped patches in the weeks following the disclosure. As of May 2026, Windsurf remains in reported state [4].

| CVE | Affected Project | Severity | Status (May 2026) |
|----------------|---------------------------------|----------|-------------------|
| CVE-2025-49596 | MCP Inspector | High | Patched |
| CVE-2025-54136 | Cursor IDE (MCPoison) | High | Patched (v1.3) |
| CVE-2025-54994 | @akoskm/create-mcp-server-stdio | High | Patched |
| CVE-2026-22252 | LibreChat | High | Patched |

| CVE | Affected Project | Severity | Status (May 2026) |
|----------------|------------------|----------|-------------------|
| CVE-2026-22688 | WeKnora | High | Patched |
| CVE-2026-30623 | LiteLLM | Critical | Patched |
| CVE-2026-30615 | Windsurf | High | Unpatched |

Sources: OX Security, NVD, Tenable, LiteLLM security advisory [1][4][9][10][15]

Anthropic's response has been to confirm the behavior is intentional and to update its SECURITY.md file – a change made nine days after OX Security's initial contact [3] – with a note that STDIO adapters "should be used with care." No architectural changes were made to the SDK. The OX research team reports having asked Anthropic multiple times to modify the protocol's defaults before going public [5]. The disclosure and Anthropic's response received immediate coverage across the security community [13]. Anthropic's position assigns sanitization responsibility to downstream developers, which creates a structural liability: the reference implementation ships without guard rails, and developers building on it inherit the exposure with minimal signal – a SECURITY.md update does not substitute for documentation, warnings in the SDK, or mandatory sanitization in the reference implementation.

The Broader Attack Taxonomy

The STDIO command execution flaw – carrying the highest CVSS severity of the documented CVE classes – is one vector in a broader taxonomy of MCP attack patterns that have matured over the past year.

Tool poisoning exploits the structural asymmetry between what an AI model processes and what a user sees. MCP tool descriptions – the text that tells an agent what a tool does and how to invoke it – are consumed by the language model but are rarely displayed in the user interface at runtime. Adversaries can embed hidden instructions within tool descriptions that direct the agent to perform unauthorized actions, exfiltrate data, or suppress notifications about its own behavior. Invariant Labs demonstrated this attack class in April 2025 by showing how a malicious trivia-game MCP server embedded instructions in its tool description targeting a legitimate WhatsApp MCP server connected to the same agent session [8]. The agent followed the embedded instructions to extract WhatsApp message history and route it outbound through the trusted server, appearing as ordinary traffic. End-to-end encryption at the transport layer provided no protection because the exfiltration occurred above the encryption boundary, through the agent's legitimate access.

Rug pull attacks exploit the absence of any mechanism in the MCP specification for tracking tool definition changes or requiring re-approval when they occur. A malicious server presents benign, useful tools to earn initial user approval, then silently modifies tool definitions or behavior in subsequent sessions. The agent transacts with the modified server session after session with no mechanism to detect that the tool it approved on day one is no longer the tool it is invoking today. The MCPoison vulnerability in Cursor IDE (CVE-2025-54136) demonstrated the real-world consequence of this pattern: an attacker could commit a benign MCP configuration to a shared code repository, wait for a developer to approve it, and then replace it with a malicious payload in a later commit. Every subsequent Cursor session would silently execute the attacker's commands without any re-approval prompt [11]. Check Point Research, which disclosed this issue to Cursor in July 2025, noted that any credential accessible to the Cursor process – API keys, cloud credentials, repository tokens – fell within reach of a successful exploit. Cursor's version 1.3 remediation now requires explicit user re-approval for any configuration change, including whitespace modifications.

Cross-server tool shadowing extends the blast radius of a single compromised MCP server across the entire set of tools an agent has access to in a given session. Because a single MCP session can connect an agent to dozens of servers simultaneously – file system, GitHub, Slack, Stripe, database connectors – a malicious server can inject tool descriptions that redefine the agent's understanding of adjacent trusted tools. Instructions embedded in one server's schema can direct the agent to route requests through a different server in ways the user never authorized, effectively turning legitimate tool integrations into exfiltration conduits without any change to those integrations themselves.

The authentication gap compounds every other risk category. A July 2025 internet-wide scan identified 1,862 MCP servers exposed to the public internet responding to unauthenticated tool-listing requests [20]. The MCP authorization specification defines an OAuth 2.1 framework but explicitly marks authorization optional for implementations [23]. In practice, many deployed MCP servers accept connections and expose full tool listings – including potentially sensitive capability descriptions – without any credential validation, meaning an attacker with network access to such a server can enumerate capabilities, probe for exploitable behavior, and attempt to trigger tool execution without valid credentials.

Documented Incidents

The MCP attack surface has moved from theoretical to operational, with a documented history of real-world exploits and active malicious tooling [12]. The first confirmed in-the-wild malicious MCP server appeared in September 2025, when a package named `postmark-mcp` – an impersonation of a legitimate Postmark email integration – shipped fifteen clean versions to a public npm registry to establish a trust baseline before version 1.0.16 silently added a single line that blind-carbon-copied every outgoing email

to an attacker-controlled address. Koi Security's disclosure, as reported by Pipelab, estimated approximately 300 organizations had integrated the package into active workflows at time of removal [4].

The Supabase-Cursor incident, documented in July 2025, illustrated how privilege amplification compounds MCP risk. A Cursor AI agent operating with privileged service-role database access was processing customer support tickets that included user-supplied text. Attackers embedded SQL directives within support ticket content, directing the agent to read sensitive integration tokens and exfiltrate them through a public support thread [21][22]. The agent's elevated database permissions, combined with no boundary between user-supplied input and executable instruction, produced a data breach from what was intended as a routine support automation workflow.

These incidents collectively establish three patterns that security teams should internalize. MCP as a protocol surface is actively targeted, not merely theoretical. The documented incidents to date suggest supply chain vectors – package impersonation, configuration file poisoning, repository compromise – as an early and observable initial access path in MCP attacks. And the sensitivity of what MCP agents routinely access – credentials, communications, code repositories, financial transaction data – makes this protocol a high-value target that warrants treatment commensurate with other privileged infrastructure components.

Recommendations

Immediate Actions

Security teams should immediately conduct an inventory of all MCP servers in use across the organization, including those introduced through developer tooling such as Cursor, Windsurf, or Claude Desktop. For each deployed MCP server, verify the patch status against the CVE table above and apply available updates. Any MCP server whose configuration is sourced from an external repository or package registry should be treated as potentially compromised until the provenance and integrity of that configuration can be independently verified.

For organizations running MCP STUDIO integrations, the immediate mitigation is process isolation. MCP server processes should run in dedicated containers or sandboxes with no access to host-level credentials, API keys, or sensitive file paths. Host-level shell command execution from within the MCP server process context should be blocked or logged at the container boundary. This does not resolve the underlying design flaw, but it limits the blast radius of a successful exploitation by confining what any attacker can reach from an MCP server process.

Developers building on MCP SDKs should audit every location where user-supplied or externally sourced data is passed to MCP STUDIO configuration parameters and apply strict allowlist validation before any process execution occurs. The sanitization responsibility that Anthropic has assigned to downstream developers is real, and no value sourced from a user, a package registry, or a network-retrieved configuration file should be treated as safe input to the STUDIO transport without explicit validation.

Short-Term Mitigations

Over the next thirty to sixty days, security teams should extend runtime behavioral monitoring to cover MCP server processes specifically. Legitimate MCP servers operate with well-defined process profiles; unexpected child processes, outbound network connections to previously unseen destinations, and file system access outside expected paths are all indicators of exploitation or tool definition manipulation. Endpoint detection and response tooling should include MCP server process behavior in its detection logic.

Organizations should establish MCP server provenance standards analogous to existing software supply chain controls. Package indexes and tool registries used to source MCP servers should require published security advisories, maintained changelogs, and provenance metadata. Configuration management systems should pin MCP server versions explicitly and alert on definition changes – implementing the rug-pull resistance that the protocol specification itself omits. The Enhanced Tool Definition Interface (ETDI), documented in academic research, offers a formal approach to cryptographically binding tool definitions to approved versions using OAuth token-based attestation [16]. Organizations with mature security programs should evaluate ETDI as a near-term architectural mitigation for environments where MCP tool integrity is a priority control.

For environments where multiple MCP servers share an agent session, session isolation should become the default architectural pattern. Each MCP server should operate in a separate process context with distinct credentials scoped only to the resources that specific tool integration requires. Cross-server interactions should be mediated through an explicit policy layer rather than relying on the agent's internal reasoning to enforce access boundaries – reasoning that, as the Invariant Labs research demonstrated, is susceptible to manipulation through tool description injection [8].

Strategic Considerations

The MCP security situation reflects a recurring pattern in infrastructure standardization: a protocol gains rapid adoption by optimizing for developer experience and capability, with security properties treated as derivative concerns for implementers to address. Anthropic's position that STUDIO sanitization is downstream developer responsibility is not unreasonable in isolation, but it creates systemic exposure

when the reference implementation ships without guard rails and the majority of developers have no security baseline to reason from. The result is an ecosystem where every project that trusted the reference implementation inherited a vulnerability class that was invisible to them at time of adoption.

Enterprises should not treat the current situation as a bounded CVE remediation exercise. The MCP ecosystem continues expanding while its security posture lags significantly behind its adoption curve. New MCP servers are currently published without any mandatory security review process, authentication requirement, or coordinated disclosure mechanism visible in the ecosystem. The attack surface will grow as MCP integration deepens into sensitive enterprise workflows: production financial systems, regulated data environments, privileged infrastructure management, and customer-facing AI services. CISOs should establish MCP governance as an ongoing program requirement – with defined policies for which categories of MCP tools are permitted, under what authentication and sandboxing requirements, and with what monitoring obligations attached.

The banking sector has been specifically flagged by researchers as particularly exposed, given the depth of financial data and transaction authority being granted to MCP-enabled agentic systems [17]. Any organization granting AI agents access to production financial systems, customer data, or regulated information should apply heightened scrutiny to every MCP integration in those environments and assess whether the current state of MCP protocol security is consistent with applicable regulatory obligations.

CSA Resource Alignment

The MCP security crisis maps directly to multiple CSA frameworks that organizations can use to structure their response and communicate risk to stakeholders.

CSA's MAESTRO threat modeling framework for agentic AI systems identifies tool interface abuse, supply chain compromise, and privilege escalation as primary threat vectors for agent architectures. The STDIO command injection flaw, tool poisoning, and rug pull patterns documented in this note correspond precisely to MAESTRO's tool manipulation and supply chain threat categories. Organizations applying MAESTRO threat models to their agentic AI deployments should explicitly incorporate MCP server trust boundaries as a modeling primitive and assign threat scenarios to each attack class documented here.

The AI Controls Matrix (AICM), CSA's comprehensive framework for AI system security controls, addresses the foundational requirements that MCP deployments currently fail to satisfy. AICM controls covering input validation, least-privilege access, software supply chain integrity, and runtime behavioral monitoring each map to specific defensive actions in the recommendation section above. Organizations

using the AICM as a compliance baseline should add MCP server configuration management and tool definition integrity verification as explicit control objectives, distinct from traditional software dependency management.

CSA's Zero Trust guidance – applied specifically to AI system architectures in CSA's agentic AI security publications – frames the core design principle that should govern MCP deployments: never trust, always verify applies equally to AI tool integrations. Unauthenticated MCP management interfaces violate this principle regardless of what the protocol specification permits. The authentication gap documented in this note – where the specification marks authorization optional and a substantial proportion of observable deployments have none – represents one of the clearest violations of zero trust principles in the agentic AI stack. CSA's position is that enterprise MCP deployments must implement explicit authentication and authorization controls as a baseline requirement, not an optional enhancement.

CSA's STAR (Security Trust Assurance and Risk) program provides a vehicle for organizations to assess and report on MCP security controls as part of broader cloud and AI security assurance programs. Vendors and cloud providers offering MCP-enabled services should be assessed against STAR criteria that explicitly address tool supply chain integrity, authentication requirements for all MCP server connections, and sandboxing requirements for agent tool execution environments. Customers procuring MCP-integrated services should include these criteria in vendor security assessments and contractual security requirements.

References

- [1] OX Security. "[The Mother of All AI Supply Chains: Critical, Systemic Vulnerability at the Core of Anthropic's MCP.](#)" OX Security Blog, April 2026.
- [2] OX Security. "[MCP Supply Chain Advisory: RCE Vulnerabilities Across the AI Ecosystem.](#)" OX Security Blog, April 2026.
- [3] VentureBeat. "[200,000 MCP servers expose a command execution flaw that Anthropic calls a feature.](#)" VentureBeat, May 2026.
- [4] Pipelab. "[The State of MCP Security 2026: Incidents, Attack Patterns, and Defense Coverage.](#)" Pipelab, April 2026.
- [5] The Register. "[MCP 'design flaw' puts 200k servers at risk: Researcher.](#)" The Register, April 16, 2026.
- [6] Infosecurity Magazine. "[Systemic Flaw in MCP Protocol Could Expose 150 Million Downloads.](#)" Infosecurity Magazine, 2026.
- [7] Simon Willison. "[Model Context Protocol has prompt injection security problems.](#)" simonwillison.net, April 9, 2025.
- [8] Invariant Labs. "[MCP Security Notification: Tool Poisoning Attacks.](#)" Invariant Labs Blog, April 2025.
- [9] NIST National Vulnerability Database. "[CVE-2025-54136 Detail.](#)" NVD, 2025.
- [10] LiteLLM. "[Security Update: CVE-2026-30623 – Command Injection via Anthropic's MCP SDK.](#)" LiteLLM Documentation, April 2026.
- [11] Check Point Research. "[Cursor IDE's MCP Vulnerability.](#)" Check Point Research, 2025.
- [12] AuthZed. "[A Timeline of Model Context Protocol \(MCP\) Security Breaches.](#)" AuthZed Blog, 2026.
- [13] The Hacker News. "[Anthropic MCP Design Vulnerability Enables RCE, Threatening AI Supply Chain.](#)" The Hacker News, April 2026.
- [14] OX Security. "[The Mother of All AI Supply Chains: Technical Deep Dive.](#)" OX Security Blog, April 2026.
- [15] Tenable. "[FAQ: CVE-2025-54135, CVE-2025-54136 Vulnerabilities in Cursor \(CurXecute and MCPoison\).](#)" Tenable Blog, 2025.

- [16] Tan et al. "[ETDI: Mitigating Tool Squatting and Rug Pull Attacks in Model Context Protocol \(MCP\)](#)." arXiv:2506.01333, 2025.
- [17] American Banker. "[Unpatched AI flaw poses risk to banking sector](#)." American Banker, 2026.
- [18] Cloud Security Alliance Labs. "[Flowise CVSS 10.0 RCE: AI Agent Builders Under Attack](#)." CSA Labs, April 9, 2026.
- [19] Cloud Security Alliance Labs. "[MCP by Design: RCE Across the AI Agent Ecosystem](#)." CSA Labs, April 20, 2026.
- [20] Dark Reading. "[Nearly 2,000 MCP Servers Possess No Security Whatsoever](#)." Dark Reading, July 2025.
- [21] General Analysis. "[Supabase MCP Can Leak Your Entire SQL Database](#)." General Analysis, July 2025.
- [22] Simon Willison. "[Supabase MCP: The Lethal Trifecta](#)." simonwillison.net, July 6, 2025.
- [23] Model Context Protocol. "[MCP Authorization Specification](#)." modelcontextprotocol.io, 2025.