



**CSAI**



**CSAI Foundation**

Cloud Security Alliance AI Safety Initiative

# **MCP Design-Level RCE: Protocol Architecture as Attack Surface**

Critical Design Vulnerabilities in Anthropic's Model Context  
Protocol Threaten the AI Development Supply Chain

Unofficial AI-assisted Research

2026-04-25

**© 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

- Researchers at OX Security disclosed in April 2026 that Anthropic's Model Context Protocol (MCP) contains a systemic, design-level vulnerability in its STDIO transport mechanism that enables remote code execution (RCE) across every supported programming language – Python, TypeScript, Java, and Rust – affecting an estimated 150 million downloads and more than 7,000 publicly exposed servers (across all exposure types) [1][2].
- The flaw is not a traditional implementation bug but an architectural decision: the STDIO interface executes operating system commands passed through configuration parameters without validation, and Anthropic has acknowledged the behavior as intentional rather than committing to a protocol-level fix [3].
- Twelve CVEs have been disclosed across widely deployed AI platforms – including Windsurf, GPT Researcher, LiteLLM, Agent Zero, and LangFlow – ten of which carried high or critical severity ratings at time of publication; the LangFlow and LangBot CVEs were awaiting CVSS assignment at the time of disclosure. These represent a cascading supply chain exposure affecting any downstream tool built on the MCP reference SDK [1][4].
- Beyond the STDIO design flaw, MCP's trust architecture creates a broader attack surface for tool poisoning, prompt injection to RCE, and supply-chain compromise through malicious server registries, compounding risk for enterprises deploying agentic AI systems [5][6].
- Organizations should immediately audit deployed MCP servers, restrict public network access to STDIO endpoints, enforce sandbox isolation for all MCP-enabled services, and treat any externally supplied MCP configuration as untrusted input pending broader protocol reform [2][4].

---

## Background

The Model Context Protocol was introduced by Anthropic in late 2024 as an open standard for connecting large language models (LLMs) to external tools, data sources, and computational resources. MCP defines how AI agents communicate with "MCP servers" – small services that expose capabilities

such as filesystem access, web search, database queries, and API integrations. The protocol rapidly became the dominant integration layer in the agentic AI ecosystem: within months of release, it was integrated into Claude Desktop, VS Code, Cursor, Windsurf, and dozens of open-source agent frameworks, accumulating more than 150 million downloads across package ecosystems [1][7].

The protocol's appeal lies in its simplicity. Developers can register a new tool by standing up an MCP server and pointing an MCP client at it. The STDIO transport – one of two primary communication mechanisms alongside HTTP/SSE – launches MCP servers as local child processes, reading their capabilities from standard output and writing commands to standard input. This design eliminates networking complexity for local deployments and accelerates integration, which has driven its adoption across the AI developer community before corresponding security guidance had been developed for the ecosystem.

That same simplicity, however, encodes a fundamental trust assumption that security researchers have demonstrated is exploitable in multiple deployment contexts. The STDIO transport treats the configuration string specifying which command to launch as authoritative, executing it directly against the host operating system without a sanitization layer. In effect, the configuration field that names the MCP server process doubles as an arbitrary command execution channel – a design choice that has propagated into downstream implementations that use the STDIO transport configuration pathway [3][8].

The security implications became apparent gradually through 2025 as researchers discovered a series of independent vulnerabilities – WhatsApp exfiltration via tool poisoning in April, GitHub MCP prompt injection in May, mcp-remote command injection in July – before the April 2026 disclosure by OX Security revealed that all of these incidents shared a common architectural root [9][10][11].

---

## Security Analysis

### The Core Design Flaw

The fundamental vulnerability in MCP's STDIO transport was described by OX Security researchers as a "configuration-to-command execution" path with no intervening validation [1]. When an MCP client initializes a STDIO server, it constructs a process invocation from parameters supplied in the configuration – parameters that, in most deployment contexts, originate from user input or externally fetched registry data. The MCP SDK executes the configured command regardless of whether the

resulting process is a legitimate MCP server, and critically, execution occurs even when the process returns an error: a malicious command runs before the error is surfaced to the caller, providing no reliable signal to block the attack [3].

This behavior is not an accident of any single implementation. OX Security confirmed the vulnerability across official SDKs in Python, TypeScript, Java, and Rust, meaning any organization that built an MCP-enabled product using STDIO transport through Anthropic's reference libraries inherited the exposure unless additional input sanitization was applied at the application layer [2]. Anthropic, upon notification, confirmed the behavior is intentional and stated that input sanitization is the developer's responsibility, declining to modify the protocol's architecture [3]. The company updated its security guidance to recommend using STDIO MCP adapters "with caution," effectively converting a systemic protocol defect into individual developer liability.

OX Security's investigation identified four distinct exploitation families that map to real-world deployment patterns. In the first, attackers supply malformed JSON configuration through publicly accessible web UI interfaces – configuration screens that accept MCP server specifications – allowing arbitrary command injection without authentication in some products. In the second, applications that implement allowlists restricting commands to known executables like `python`, `npm`, or `npx` can be bypassed by injecting commands through argument flags: `npx -c <malicious_command>` passes the allowlist while executing arbitrary code [4]. The third family exploits AI coding assistants that process attacker-controlled HTML content, causing the assistant to autonomously modify local MCP configuration files and register malicious STDIO servers – prompt injection escalating to RCE through the agent's own tool use. The fourth involves modifying network requests in transit to change a server's transport type from HTTP to STDIO and inject a malicious command, exploiting backend systems that process STDIO configurations not visible in application UIs [4].

## CVE Inventory

The OX Security investigation resulted in twelve publicly disclosed CVEs across widely deployed platforms. CVE-2026-30615 affecting the Windsurf IDE is notable as the only instance among the batch requiring zero user interaction for exploitation – a fully automated attack path [4]. CVE-2026-30623 in LiteLLM and CVE-2026-30624 in Agent Zero both carry critical severity ratings. GPT Researcher (CVE-2025-65720), Jaz (CVE-2026-30616), Langchain-Chatchat (CVE-2026-30617), DocsGPT (CVE-2026-26015), and Bisheng (CVE-2026-33224) round out the critical-severity disclosures. Flowise (CVE-2026-40933) and Upsonic (CVE-2026-30625) were rated high severity, and LangFlow and LangBot received unassigned CVEs at the time of publication [4].

These represent only the formally disclosed cases from a single research effort. Researchers confirmed live RCE on six production platforms and estimated up to 200,000 vulnerable instances across the broader ecosystem, with more than 200 open-source projects carrying the inherited flaw [1][2]. Trend Micro independently identified 492 MCP servers exposed to the internet with no authentication layer, providing unobstructed access to the underlying command execution pathway [16].

## Earlier Incidents: A Converging Pattern

The April 2026 disclosure crystallizes a pattern visible throughout 2025. In July 2025, JFrog disclosed CVE-2025-6514 (CVSS 9.6) in `mcp-remote`, an OAuth proxy with 437,000 downloads used to connect Claude Desktop and similar clients to remote MCP servers [11][17]. The vulnerability exploited the OAuth authorization flow: a malicious server responded with a weaponized `authorization_endpoint` value that `mcp-remote` passed to the `open()` npm package, which on Windows encoded the value as a PowerShell command and executed it. Researchers demonstrated command execution using the PowerShell subexpression operator to bypass URL encoding restrictions. The fix landed in `mcp-remote` version 0.1.16. While OX Security did not explicitly link CVE-2025-6514 to the April 2026 disclosure, the core pattern – unsanitized external data flowing into a command execution context – anticipates the architectural flaw by nine months.

In June 2025, CVE-2025-49596 was issued for MCP Inspector's `inspector-proxy` component, allowing unauthenticated RCE on developer workstations – exposing entire filesystems, API keys, and environment secrets [9]. A September 2025 supply-chain incident involved a malicious package impersonating the Postmark MCP server, which silently BCC'd all outbound email communications to an attacker-controlled address [9]. That same month, a path-traversal vulnerability in the Smithery MCP hosting platform allowed attackers to compromise build credentials controlling more than 3,000 downstream applications [9].

## Tool Poisoning and Prompt Injection

The design-level RCE vulnerabilities exist alongside a related but distinct class of attacks that exploit MCP's trust model at the semantic layer. Tool poisoning embeds malicious instructions in an MCP server's tool metadata – its name, description, or parameter definitions – which the LLM reads and interprets as behavioral guidance. Because LLMs rely on these descriptions to understand how to invoke a tool, a sufficiently crafted description can redirect the model's behavior regardless of what the server actually does. Invariant Labs demonstrated in April 2025 that a malicious MCP server registered alongside a legitimate WhatsApp MCP server could silently exfiltrate a user's entire message history by embedding exfiltration instructions in its tool description, bypassing data-loss prevention controls that examined only the visible conversation [9].

Prompt injection through MCP differs from conventional prompt injection in its persistence: once a tool description is poisoned, any agent session that loads that tool may inherit the malicious instruction – a one-time compromise with potentially broad reach absent active monitoring and tool-registry governance [6]. The GitHub MCP incident in May 2025 illustrated how over-privileged personal access tokens, combined with a prompt injection payload embedded in a repository file, could direct an agent to exfiltrate private repository contents, salary data, and project details to a remote server [9].

Research published in early 2025 examining publicly available MCP server implementations found that 43% contained command injection flaws and 30% permitted unrestricted URL fetching – attack surface characteristics that enable both direct exploitation and secondary data exfiltration [6]. These figures suggest the vulnerabilities cataloged by name represent a fraction of the actual exposed population.

## Supply Chain Amplification

What distinguishes the MCP situation from a conventional software vulnerability is the supply chain amplification mechanism. The vulnerability propagates through a layer of abstraction: developers who build agentic applications do not themselves write MCP SDK internals, they consume them. A security defect in Anthropic's reference SDK becomes a defect in LangChain, LiteLLM, LangFlow, and the hundreds of open-source frameworks that extend those libraries, which in turn become defects in the enterprise products and custom deployments built on top of those frameworks. OX Security's estimate of 200,000 vulnerable instances reflects not seven or twelve software products but the geometric expansion of a foundational flaw across an entire development ecosystem [1][2].

The supply chain risk is compounded by the emerging ecosystem of MCP server registries. Developers routinely install MCP servers from public repositories and third-party hosting platforms, creating a distribution channel that, unlike established package registries, currently lacks standardized security review, signing requirements, or abuse detection tooling. The February 2026 fake Oura MCP package incident demonstrated that threat actors have already weaponized this channel, cloning legitimate projects and substituting StealC malware to harvest developer credentials, browser passwords, and cryptocurrency wallets from developer workstations [9].

---

# Recommendations

## Immediate Actions

Organizations running MCP-enabled services should treat the STDIO transport endpoint as an unauthenticated remote code execution surface until proven otherwise. Network isolation is the first priority – as it addresses the broadest attack surface while patches are assessed – requiring that organizations block all public IP access to services that accept MCP STDIO configuration input and ensure MCP-enabled developer tools such as IDEs, agent frameworks, and coding assistants are not exposing STDIO endpoints to untrusted network segments. Applying all vendor patches against the disclosed CVEs (LiteLLM, Windsurf, GPT Researcher, and others) reduces risk for the specific platforms covered, though organizations should recognize that patched implementations have addressed their individual instance of the flaw while the underlying protocol design remains unchanged.

## Short-Term Mitigations

Any input that specifies an MCP server command, transport type, or configuration parameter must be treated as untrusted. Validated allowlists of permitted executables and paths should be enforced at the application layer, with explicit prohibition on argument-flag bypasses (for example, by validating the full tokenized command rather than the executable name alone). MCP-enabled services should run inside sandboxed environments – containers or VMs with minimal filesystem access, no access to host credentials, and restricted network egress – so that a successful RCE exploit yields a constrained foothold rather than host compromise. Tool invocations should be logged comprehensively, with anomaly detection configured to alert on unexpected subprocess creation, outbound connections, or filesystem modifications originating from MCP server processes. Enterprises should audit all installed MCP servers against known-good registry entries and remove any packages not sourced from verified publishers.

## Strategic Considerations

The MCP disclosure is reminiscent of early web-era open-redirect and SSRF vulnerabilities: a design decision that seemed convenient in a low-threat environment became a systemic liability as adoption scaled and adversarial interest grew. The appropriate long-term response is protocol-level reform – specifically, separating the MCP server configuration channel from the command execution pathway, introducing cryptographically verified server identity at the registry layer, and standardizing a sandboxed

execution model as a protocol requirement rather than a vendor option. CSA encourages organizations engaged with Anthropic's MCP working group or contributing to open-source MCP implementations to advocate for these architectural changes.

At the enterprise governance level, MCP server provenance should be treated as a software supply chain control analogous to open-source dependency management. Organizations adopting agentic AI systems should require provenance attestation for any MCP server entering their environment, apply the same SCA and vulnerability scanning tooling used for software packages, and establish a documented approval process before new MCP servers are registered in production deployments. MCP servers, like other open-source dependencies, are subject to abandonment – the Smithery incident demonstrated that even hosted MCP infrastructure can become a liability when maintenance lapses [9]. Organizations should plan for regular re-verification of installed servers' maintenance status and security posture.

---

## CSA Resource Alignment

The vulnerabilities described in this note map directly to threats modeled within CSA's MAESTRO framework for agentic AI security [12]. MAESTRO's Layer 3 (Agent Frameworks) addresses vulnerabilities in the tool invocation and orchestration layer – precisely the layer where MCP's STUDIO design flaw resides – while Layer 7 (Agent Ecosystem) addresses analogous threat categories including marketplace manipulation, agent impersonation, and supply chain compromise of agent tools [12]. Organizations using MAESTRO for threat modeling should incorporate the four exploitation families identified by OX Security into their Layer 3 and Layer 7 threat scenarios.

CSA's Agentic AI Red Teaming Guide, co-developed with OWASP AI Exchange, provides specific test procedures for supply chain and dependency attacks relevant to the MCP context [13]. The guide's coverage of multi-agent exploitation and knowledge base poisoning also applies to tool poisoning scenarios where a compromised MCP server's metadata influences an agent's behavior across multiple sessions.

CSA's AI Controls Matrix (AICM) – a control framework for AI systems and a superset of the Cloud Controls Matrix – provides the compliance scaffolding for implementing the mitigations described above [14]. Relevant AICM control domains include Supply Chain Security (requiring provenance verification for AI components), Application and Interface Security (input validation at tool integration boundaries), and Audit Assurance and Compliance (logging and monitoring of tool invocations). Organizations pursuing CSA STAR certification for AI-enabled services should document their MCP server governance processes under these domains.

The AI Organizational Responsibilities guidance on core security responsibilities addresses the developer liability question at the center of the Anthropic response: when a protocol designer declines to implement safety controls and shifts responsibility to implementers, the implementers must have clear documented standards governing how that responsibility is discharged [15]. Organizations should not rely on implied or ad-hoc practices for MCP security; they should codify configuration validation requirements, sandbox standards, and supply chain approval processes as formal security controls.

# 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] The Hacker News. "[Anthropic MCP Design Vulnerability Enables RCE, Threatening AI Supply Chain.](#)" The Hacker News, April 2026.
- [4] SecurityWeek. "['By Design' Flaw in MCP Could Enable Widespread AI Supply Chain Attacks.](#)" SecurityWeek, April 2026.
- [5] Infosecurity Magazine. "[Systemic Flaw in MCP Protocol Could Expose 150 Million Downloads.](#)" Infosecurity Magazine, April 2026.
- [6] Elastic Security Labs. "[MCP Tools: Attack Vectors and Defense Recommendations for Autonomous Agents.](#)" Elastic Security Labs, 2025.
- [7] eSentire. "[Model Context Protocol Security: Critical Vulnerabilities Every CISO Must Address in 2025.](#)" eSentire, 2025.
- [8] Practical DevSecOps. "[MCP Security Vulnerabilities: How to Prevent Prompt Injection and Tool Poisoning Attacks in 2026.](#)" Practical DevSecOps, 2026.
- [9] AuthZed. "[A Timeline of Model Context Protocol \(MCP\) Security Breaches.](#)" AuthZed Blog, 2026.
- [10] Simon Willison. "[Model Context Protocol has prompt injection security problems.](#)" simonwillison.net, April 2025.
- [11] JFrog Security Research. "[Critical RCE Vulnerability in mcp-remote: CVE-2025-6514 Threatens LLM Clients.](#)" JFrog, July 2025.
- [12] Cloud Security Alliance. "[Agentic AI Threat Modeling Framework: MAESTRO.](#)" CSA Blog, February 2025.
- [13] Cloud Security Alliance. "[Agentic AI Red Teaming Guide.](#)" CSA AI Organizational Responsibilities Working Group, 2025.

[14] Cloud Security Alliance. ["AI Controls Matrix \(AICM\)."](#) CSA, 2025.

[15] Cloud Security Alliance. ["AI Organizational Responsibilities: Core Security Responsibilities."](#) CSA, 2025.

[16] Trend Micro. ["MCP Security: Network-Exposed Servers Are Backdoors to Your Private Data."](#) Trend Micro, 2025.

[17] The Hacker News. ["Critical mcp-remote Vulnerability Exposes LLM Clients to Remote Code Execution."](#) The Hacker News, July 2025.