

# A Critical Security and Architectural Review of the Model Context Protocol (MCP) Ecosystem

Yusuf Talha ARABACI  
M.Sc. Student in Software Engineering  
Karabuk University  
Karabuk, Turkey

**Abstract**—The Model Context Protocol (MCP), introduced in late 2024, represents a pivotal shift in how Large Language Models (LLMs) interact with the external world. By standardizing the connection between AI agents and data sources, MCP aims to resolve the combinatorial “ $N \times M$ ” integration bottleneck that has long hampered the scalability of agentic systems. This paper provides a critical synthesis of the burgeoning MCP ecosystem, moving beyond a simple feature overview to analyze the architectural trade-offs, security implications, and optimization challenges that accompany this new standard. We scrutinize the protocol’s core transport layers, evaluate the industry’s pivot from direct tool calling to a “Code Execution Paradigm” to combat context bloat, and dissect the “Lethal Trifecta” of risks—Data, Internet, and Action—that emerges when agents are granted autonomy. Finally, we explore the “AgentX” pattern for multi-agent orchestration and propose necessary governance frameworks to ensure the ecosystem matures securely.

**Keywords**—Model Context Protocol, MCP, context bloat, code execution paradigm, agentic security, lethal trifecta.

## I. INTRODUCTION: THE INTEGRATION KNOT

The transition of Large Language Models (LLMs) from passive chat interfaces to autonomous “agents” capable of executing complex workflows has hit a significant architectural wall: the lack of a universal connection standard. Historically, connecting a specific LLM ( $N$ ) to a specific data source ( $M$ ) required a bespoke integration. As the number of models and tools exploded, this created an unmanageable “ $N \times M$ ” integration knot, stifling innovation and creating massive maintenance overhead key [7].

Enter the Model Context Protocol (MCP). Functioning effectively as a “USB-C port” for the AI operational layer, MCP standardizes this interface, theoretically reducing the integration complexity from quadratic ( $N \times M$ ) to linear ( $N + M$ ). Since its release, backed by heavyweights like Anthropic, OpenAI, and Google, MCP has rapidly de-facto standardized the field, shifting the paradigm from stateless API calls to stateful, session-based bidirectional communication. However, as this paper argues, widespread adoption brings its own set of distinct challenges, particularly regarding security and context efficiency.

## II. ARCHITECTURAL PARADIGMS

MCP’s architecture is designed to decouple the agent’s reasoning capabilities (the Client/Host) from the specific implementation details of the tools it uses (the Server). This modularity is achieved through three core primitives.

### A. Core Primitives

Rather than treating every interaction as unique, MCP abstracts them into three categories:

- 1) **Tools:** These are the agent’s “hands.” They are executable functions—such as querying a SQL database or sending a Slack message—that the model can invoke to manipulate the external world.
- 2) **Resources:** These represent the agent’s “eyes.” They are passive data streams—like log files, database schemas, or API responses—that the application exposes for the model to read and analyze.
- 3) **Prompts:** These serve as the agent’s “instruction manuals,” providing standardized, user-defined templates (e.g., “Analyze this security log”) to guide common workflows.

### B. Transport Layers: Local vs. Remote

Implementation choices in MCP are dictated by the specific security and latency needs of the deployment:

- **Stdio (Standard Input/Output):** The gold standard for local, high-security environments. Here, the server runs as a subprocess of the host application. Since data never leaves the local machine’s memory space, it offers both ultra-low latency and a minimized attack surface.
- **HTTP with SSE (Server-Sent Events):** The necessary choice for distributed or cloud-native agents. While the client pushes data via standard HTTP POST, the server maintains an open channel via SSE to push asynchronous updates (like “File Changed”). This approach is firewall-friendly but introduces standard web security vectors.

### C. The Automation Revolution

Perhaps the most significant accelerant for MCP adoption has been the emergence of tools like **AutoMCP**. By parsing existing OpenAPI (Swagger) specifications, these compilers can generate fully functional MCP servers with near-perfect accuracy (99.9% success rate), compressing weeks of integration work into mere minutes [5].

## III. PERFORMANCE: THE BATTLE AGAINST BLOAT

Solving connectivity is only half the battle; the other half is efficiency. MCP introduces a phenomenon we term “Context

Bloat,” where the sheer volume of available tools overwhelms the model.

#### A. The Context Bloat Paradox

Giving an agent access to thousands of tools sounds powerful, but in practice, loading thousands of JSON schemas into the context window is disastrous. It not only spikes token costs by up to  $236\times$  but also degrades performance, as models suffer from the “Lost in the Middle” phenomenon when drowning in irrelevant definitions [6].

#### B. The Code Execution Paradigm Shift

To counter this, the industry is pivoting from a “Direct Tool Calling” model to a “Code Execution Paradigm”:

- **Direct Tool Calling (Legacy):** The model painstakingly calls one tool, waits for the result, processes it, and then calls the next. It is chatty, slow, and token-expensive.
- **Code Execution (Modern):** The model acts like a developer. It writes a single Python script that imports the necessary tools, performs the data fetching, filtering, and analysis in one go, and executes it within a sandbox. This approach can reduce token usage by over 98% and significantly lower latency.

## IV. SECURITY: THE LETHAL TRIFECTA

Empowering safe, text-generating models with the ability to act on the world creates what security researchers call the “Lethal Trifecta”: the combination of **Data Access**, **Internet Connectivity**, and **Action Capability**.

#### A. Emerging Threat Vectors

The attack surface for MCP is distinct from traditional web applications:

- **Indirect Prompt Injection (IPI):** This is the most insidious threat. An attacker hides malicious instructions (e.g., in white text on a harvested webpage) that the agent reads. The agent, treating this data as context, unwittingly prioritizes these instructions over the user’s original intent.
- **Tool Poisoning:** Malicious actors may publish MCP servers with subtle traps in their tool descriptions, tricking the model into executing harmful commands under the guise of legitimate operations.
- **Sampling Manipulation:** A compromised server could subtly alter the conversation history during a sampling request, effectively gaslighting the model into a compromised state.

#### B. Defense-in-Depth Strategies

Securing MCP requires a layered approach, moving beyond simple border firewalls:

- 1) **Human-in-the-Loop (HITL):** For high-stakes actions—like deleting files or transferring funds—automation must yield to human authorization.

- 2) **Rigorous Sandboxing:** Every MCP server should be treated as potentially hostile. Running them in isolated containers or microVMs with read-only file system roots is non-negotiable.
- 3) **Least Privilege:** Agents should never run as “root.” Scoped tokens must restrict access strictly to the resources required for the current task.

## V. THE FUTURE: MULTI-AGENT ORCHESTRATION

As we look forward, MCP is evolving from a point-to-point connector into the connective tissue of Multi-Agent Systems (MAS).

#### A. The AgentX Pattern

Are we building one super-agent, or a team of specialists? The trend favors the latter. The **AgentX** pattern orchestrates a team: a Planner to break down tasks, a Researcher to gather data, a Coder to process it, and an Auditor to verify the results. MCP provides the common language that allows these distinct agents to share context and seamlessly hand off resources.

#### B. Hybrid Architectures

The future is not a binary choice between REST and MCP. We envision a hybrid architecture where REST/gRPC remains the backbone of high-performance microservices, while MCP becomes the standard protocol for the unpredictable, dynamic “Agent-System” interface layer.

## VI. CONCLUSION

The Model Context Protocol has successfully solved the integration problem, turning AI agents from isolated chatbots into capable digital colleagues. However, this capabilities upgrade comes with a significant responsibility. The ecosystem’s long-term viability depends less on further protocol tweets and more on rigorous governance. Adopting standards like ISO 42001, embracing secure-by-default architectures, and respecting the “Lethal Trifecta” are the necessary steps to ensure that our agents remain helpful assistants rather than becoming security liabilities.

## KAYNAKÇA

- [1] X. Hou, Y. Zhao, S. Wang, and H. Wang, “Model Context Protocol (MCP): Landscape, Security Threats, and Future Research Directions,” *arXiv preprint arXiv:2503.23278*, 2025.
- [2] N. Krishnan, “Advancing Multi-Agent Systems Through Model Context Protocol: Architecture, Implementation, and Applications,” *arXiv preprint arXiv:2504.21030*, 2025.
- [3] A. Ehtesham, A. Singh, G. K. Gupta, and S. Kumar, “A Survey of Agent Interoperability Protocols,” *arXiv preprint arXiv:2505.02279*, 2025.
- [4] M. M. Hasan et al., “Model Context Protocol (MCP) at First Glance: Studying the Security and Maintainability of MCP Servers,” *arXiv preprint arXiv:2506.13538*, 2025.
- [5] M. Mastouri, E. Ksontini, and W. Kessentini, “Making REST APIs Agent-Ready: From OpenAPI to MCP Servers,” *arXiv preprint arXiv:2507.16044*, 2025.
- [6] W. Song et al., “Help or Hurdle? Rethinking Model Context Protocol-Augmented Large Language Models,” *arXiv preprint arXiv:2508.12566*, 2025.
- [7] A. Singh et al., “A Survey of the Model Context Protocol (MCP): Standardizing Context,” *Preprints 202504.0245*, 2025.
- [8] W. Xing et al., “MCP-Guard: A Defense Framework for Model Context Protocol Integrity,” *arXiv preprint arXiv:2508.10991*, 2025.