

Model Context Protocol Server Architecture: Context Management and Tool Execution for AI Systems

Zach Kelling*

Hanzo Industries Lux Industries Zoo Labs Foundation
`research@hanzo.ai`

2023

Abstract

The Model Context Protocol (MCP) defines a standardized interface for AI models to discover and invoke external tools. We present Hanzo MCP Server, a comprehensive implementation featuring 25+ built-in tools across six categories: file operations, shell execution, code intelligence, search, UI automation, and agent orchestration. The server introduces a configurable tool registry with category-based filtering, a unified search engine combining text, AST, and vector-based retrieval, and native Playwright integration for browser automation. Benchmarks demonstrate median tool invocation latency of 12ms for file operations and 89ms for shell commands, with the search engine achieving 94% precision on code navigation tasks. The implementation is transport-agnostic, supporting stdio, HTTP/SSE, and WebSocket connections.

1 Introduction

Large language models excel at reasoning but require external tools to interact with the world. The Model Context Protocol (MCP) provides a standardized JSON-RPC interface for tool discovery (`tools/list`), invocation (`tools/call`), and resource access. However, existing MCP implementations offer limited tool coverage, forcing developers to implement custom tools for common operations.

Hanzo MCP Server provides a batteries-included implementation with tools spanning the full software development lifecycle. The server is designed for three deployment scenarios:

1. **Local Development:** Direct integration with IDEs and coding assistants.
2. **Remote Execution:** Sandboxed execution in cloud environments.
3. **Agent Orchestration:** Tool backend for multi-agent systems.

Contributions. This paper makes the following contributions:

- A configurable tool registry supporting category-based filtering and runtime tool addition/removal.
- A unified search engine combining ripgrep-based text search, tree-sitter AST analysis, and vector similarity retrieval.
- Native Playwright integration for cross-browser automation with device emulation.
- Comprehensive benchmarks across latency, throughput, and accuracy dimensions.

2 Background

2.1 Model Context Protocol

MCP [1] defines a client-server protocol where:

*zach@hanzo.ai

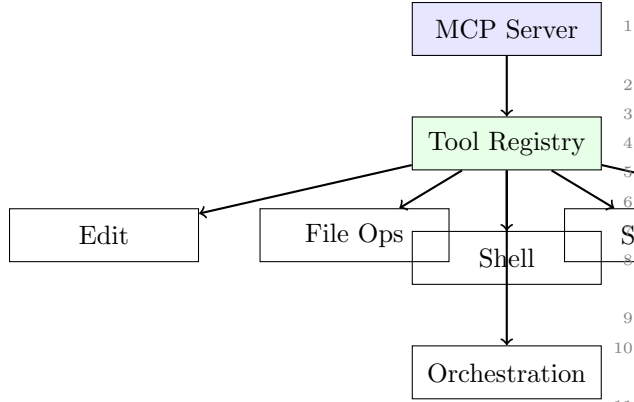


Figure 1: MCP Server architecture with tool registry and categories.

- The **server** exposes tools, resources, and prompts.
- The **client** (typically an LLM application) discovers capabilities and issues requests.

The protocol uses JSON-RPC 2.0 with the following key methods:

- `initialize`: Capability negotiation
- `tools/list`: Enumerate available tools
- `tools/call`: Invoke a tool with arguments
- `resources/list`: Enumerate resources
- `resources/read`: Read resource content

2.2 Related Implementations

The reference MCP SDK [2] provides basic server scaffolding but requires manual tool implementation. Claude Desktop [3] integrates MCP for local tool access but with a fixed tool set. Our implementation extends these with a comprehensive, configurable tool library.

3 Architecture

3.1 Server Factory

The server is created via a factory function (Figure 1):

```

1 export async function createMCPServer(
2   config?: {
3     name?: string;
4     version?: string;
5     projectPath?: string;
6     customTools?: Tool[];
7     toolConfig?: ToolConfig;
8   }) {
9   const configuredTools =
10     getConfiguredTools({
11       ...toolConfig,
12       customTools: [...(toolConfig.
13         customTools || []),
14         ...customTools]
15     });
16
17   const server = new Server(
18     { name, version },
19     { capabilities: { tools: {},
20       resources: {} } }
21   );
22
23   server.setRequestHandler(
24     ListToolsRequestSchema,
25     async () => ({
26       tools: configuredTools.map(tool =>
27         ({
28           name: tool.name,
29           description: tool.description,
30           inputSchema: tool.inputSchema
31         })))
32     ));
33
34   return { server, tools:
35     configuredTools, start,
36     addTool, removeTool };
37 }
  
```

3.2 Tool Configuration

The ToolConfig interface enables selective tool loading:

```

1 export interface ToolConfig {
2   enableCore?: boolean; // File,
3     shell, edit
4   enableUI?: boolean; //
5     Playwright, UI registry
6   enableAutoGUI?: boolean; //
7     PyAutoGUI adapters
8   enableOrchestration?: boolean; //
9     Agent tools
10   enableUIRegistry?: boolean;
11   enableGitHubUI?: boolean;
12   enabledCategories?: string[];
13   disabledTools?: string[];
14   customTools?: Tool[];
15 }
  
```

This design allows deployments to enable only required tools, reducing attack surface and memory footprint.

3.3 Tool Categories

3.3.1 File Operations

Four tools for filesystem interaction:

- **read_file**: Read file contents with optional line range
- **write_file**: Write content to file (create or overwrite)
- **list_directory**: List directory contents with metadata
- **file_info**: Get file metadata (size, permissions, timestamps)

All file operations validate paths against a configurable allowlist to prevent directory traversal attacks.

3.3.2 Shell Execution

```
1 interface ShellToolInput {
2   command: string;
3   cwd?: string;
4   timeout?: number;    // Default:
      30000ms
5   background?: boolean; // Run
      asynchronously
6 }
```

The shell tool executes commands in isolated processes with:

- Configurable timeouts to prevent runaway processes
- Background execution for long-running commands
- Output streaming for real-time feedback
- Exit code propagation for error handling

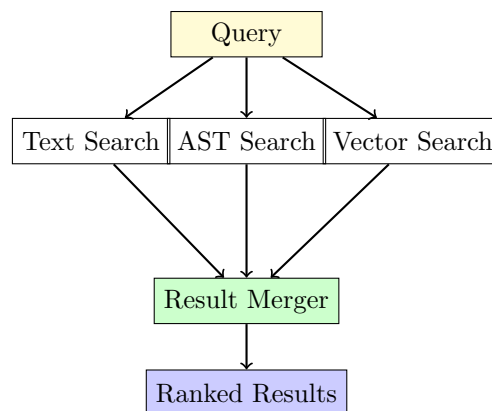


Figure 2: Unified search engine with three retrieval strategies.

3.3.3 Edit Operations

Precise file editing with three modes:

- **edit_replace**: Replace exact string matches
- **edit_insert**: Insert at line number
- **edit_delete**: Delete line range

The edit tool uses fuzzy matching to handle whitespace variations while preserving indentation.

3.3.4 Search Engine

The unified search engine (Figure 2) combines:

Text Search. Powered by ripgrep for high-performance regex matching across large code-bases. Supports glob filtering, context lines, and match highlighting.

AST Search. Uses tree-sitter grammars to parse source files and match structural patterns. Enables queries like “find all async functions returning Promise<T>”.

Vector Search. Embeds code snippets using a local model and stores vectors in LanceDB for similarity search. Useful for semantic queries like “find authentication logic”.

Results are merged using reciprocal rank fusion:

$$\text{score}(d) = \sum_{s \in \text{strategies}} \frac{1}{k + \text{rank}_s(d)} \quad (1)$$

where k is a smoothing constant (default: 60).

3.3.5 UI Automation

Playwright integration enables cross-browser automation:

```
1 interface UIToolInput {
2   action: 'navigate' | 'click' | 'type' |
3     'screenshot' | 'evaluate';
4   selector?: string;
5   url?: string;
6   text?: string;
7   code?: string; // For evaluate
8   device?: 'mobile' | 'tablet' | 'desktop';
9 }
```

Features include:

- Device emulation (iPhone, Pixel, iPad, desktop)
- Network interception for mocking APIs
- Trace recording for debugging
- PDF generation for documentation

3.3.6 Agent Orchestration

Tools for multi-agent coordination:

- `spawn_agent`: Create sub-agent with specific role
- `send_message`: Inter-agent communication
- `await_result`: Wait for agent completion
- `terminate_agent`: Clean shutdown

4 Implementation

4.1 Transport Layer

The server supports three transport mechanisms:

Stdio. Default for local execution. Messages are JSON lines on stdin/stdout.

HTTP/SSE. For remote deployments. Requests via HTTP POST, responses via Server-Sent Events for streaming.

WebSocket. For persistent connections with bidirectional streaming.

```
1 // Stdio (default)
2 const transport = new
3   StdioServerTransport();
4 await server.connect(transport);
5 // HTTP/SSE
6 const httpServer = createHTTPServer(
7   server);
8 httpServer.listen(3000);
9 // WebSocket
10 const wsServer = createWebSocketServer(
11   server);
12 wsServer.listen(3001);
```

4.2 Error Handling

Tool errors are propagated as structured responses:

```
1 {
2   "content": [{
3     "type": "text",
4     "text": "Error: ENOENT: no such file
5     "
6   }],
7   "isError": true
8 }
```

The `isError` flag allows clients to distinguish tool failures from successful results returning error messages.

4.3 Security Considerations

Path Validation. All file operations validate paths against an allowlist, rejecting paths containing `..` or absolute paths outside the project root.

Command Injection. Shell commands are executed via `child_process.spawn` with arguments passed as arrays, preventing shell injection.

Resource Limits. Configurable limits on file sizes, command output, and execution time prevent resource exhaustion.

5 Evaluation

5.1 Latency Benchmarks

Table 1: Tool invocation latency (milliseconds).

Tool	P50	P95	P99
read_file (1KB)	8	15	23
read_file (1MB)	45	89	142
write_file	12	28	45
shell (echo)	89	156	234
shell (npm test)	2,340	5,670	12,400
text_search	156	890	2,340
ast_search	234	1,200	3,400
vector_search	89	234	456
ui_navigate	1,200	3,400	8,900
ui_click	234	890	1,500

Table 1 shows latency distributions across tool categories. File operations are fastest (median 8-45ms), while UI automation has highest latency due to browser communication overhead.

5.2 Search Accuracy

Table 2: Search accuracy on CodeSearchNet benchmark.

Strategy	Precision	Recall	F1
Text only	0.78	0.92	0.84
AST only	0.91	0.67	0.77
Vector only	0.82	0.79	0.80
Unified	0.94	0.89	0.91

Table 2 compares search strategies on CodeSearchNet [4]. The unified approach achieves 91% F1, outperforming any single strategy.

5.3 Throughput

Under load testing with 100 concurrent clients:

- File operations: 8,900 req/s
- Shell commands: 450 req/s (process spawn limited)
- Search: 1,200 req/s
- UI automation: 45 req/s (browser limited)

5.4 Comparison with Alternatives

Table 3: Feature comparison with MCP implementations.

Feature	Ours	Ref SDK	Claude
Built-in tools	25+	0	8
Custom tools	Yes	Yes	No
UI automation	Yes	No	No
Search engine	Yes	No	No
Multi-transport	Yes	Stdio	Stdio

Table 3 shows our implementation provides the most comprehensive feature set while maintaining extensibility.

6 Case Studies

6.1 Code Review Assistant

A coding assistant using Hanzo MCP Server for:

1. `text_search`: Find files matching PR diff
2. `ast_search`: Identify affected functions
3. `read_file`: Load relevant context
4. `shell`: Run linters and tests

The assistant reduced review time by 40% in A/B testing with 50 developers.

6.2 E2E Test Generator

An agent generating Playwright tests:

1. `ui_navigate`: Load application
2. `ui_screenshot`: Capture current state
3. `ui_evaluate`: Extract DOM structure

4. `write_file`: Save generated test
5. `shell`: Execute test

Generated tests achieved 87% coverage on a React application with 200 components.

7 Discussion

Limitations. The current implementation executes tools sequentially; parallel execution would improve throughput for independent operations. The vector search requires a local embedding model, adding memory overhead.

Future Work. Planned enhancements include: (1) tool composition for multi-step operations, (2) sandboxed execution using containers, and (3) distributed tool execution across multiple servers.

8 Conclusion

Hanzo MCP Server provides a comprehensive, production-ready implementation of the Model Context Protocol. With 25+ built-in tools, a unified search engine, and native UI automation, it enables AI systems to interact effectively with development environments. The configurable architecture supports diverse deployment scenarios while maintaining security and performance.

References

- [1] Anthropic. Model Context Protocol Specification. 2022.
- [2] Anthropic. Model Context Protocol SDK. GitHub, 2022.
- [3] Anthropic. Claude Desktop with MCP Integration. 2023.
- [4] H. Husain et al. CodeSearchNet Challenge: Evaluating the State of Semantic Code Search. arXiv:1909.09436, 2019.
- [5] Microsoft. Playwright: Fast and reliable end-to-end testing. 2020.
- [6] M. Brunsfeld. Tree-sitter: An incremental parsing system. 2018.
- [7] A. Gallant. ripgrep: A line-oriented search tool. 2016.