



Model Context Protocol & Elasticsearch

By Moncef Abboud

cefboud.com

What is MCP?

...Like, Is It a Big Deal or Something?

(Spoiler: Yeah, kinda.)

✦ New: Remote MCP Servers

Awesome MCP Servers

A collection of servers for the Model Context Protocol.



Featured

All

Official ✦

Search

Web Scrapping

Communication

Productivity

Development

Database

Cloud Service

File System

Cloud Storage

Version Control

Other

Featured MCPs

[View all](#)

Bright Data

sponsor

Discover, extract, and interact with the web - one interface powering automated access across the public internet.

[View Details](#)

Browserbase

official

Automate browser interactions in the cloud (e.g. web navigation, data extraction, form filling, and more)

[View Details](#)

Chrome DevTools MCP

official

chrome-devtools-mcp lets your coding agent (such as Gemini, Claude, Cursor or Copilot) control and inspect a live Chrome browser

[View Details](#)

Cloudflare

official

Deploy, configure & interrogate your resources on the Cloudflare developer platform (e.g. Workers/KV/R2/D1)

[View Details](#)

Context 7

official

Up-to-date Docs For Any Cursor Prompt

[View Details](#)

DeepWiki by Devin

official

Remote, no-auth MCP server providing AI-powered codebase context and answers

[View Details](#)

E2B

official

Run code in secure sandboxes hosted by E2B

[View Details](#)

Exa

official

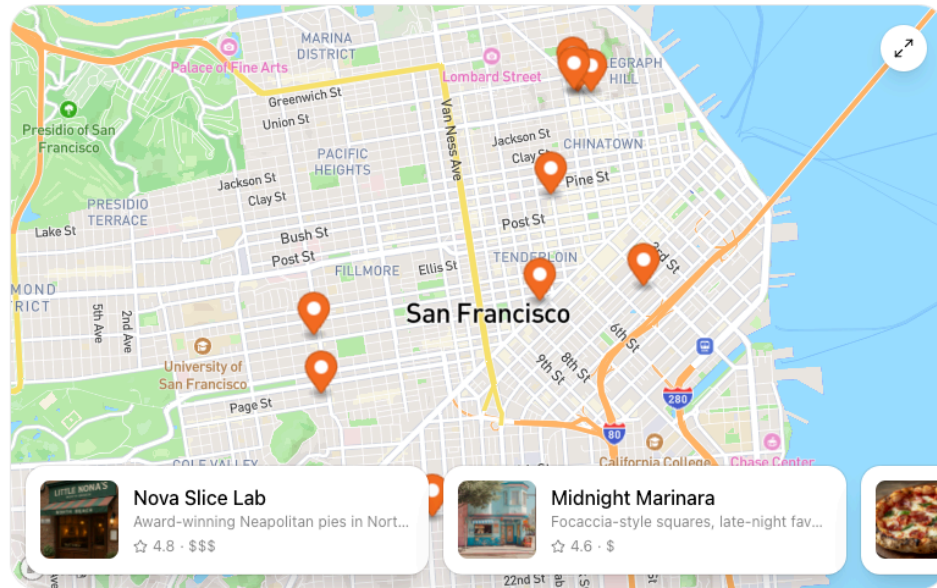
Search Engine made for AIs by Exa

[View Details](#)

Looked for available tools

(=) Served a fresh map

My Pizza App



Here's your **pizza map** 🍕 showing locations themed around **cheese** — enjoy exploring it!



DEVELOPER MODE

+ Ask anything



Let's Travel all the way back to 2022

LLMs are **token-only** systems



The Problem

- LLMs can only output tokens (text, vision)
- They can't **execute actions**
- They can't **query databases**
- They can't **search your logs**

What if they could?

ReAct: Reasoning and Acting

Key Insight: Give LLMs a list of potential actions they can call

```
1  ##### System Prompt #####
2
3  You run in a loop: **Thought → Action → PAUSE → Observation**, then give an **Answer**.
4  **Thought** explains your reasoning.
5  **Action** uses one of the tools below, then returns **PAUSE**.
6  **Observation** is the result of the action.
7
8  **Available actions:**
9
10 * `calculate`: Run a Python-style calculation (e.g., `calculate: 4 * 7 / 3.0`)
11 * `get_weather`: Get current weather for a location (e.g., `get_weather: Tokyo`)
12
13 Example:
14 Q: What's 5 times 3?
15 Thought: I should calculate it.
16 Action: calculate: 5 * 3
17 PAUSE
```

- <https://til.simonwillison.net/llms/python-react-pattern>

ReAct (2)

The Process:

1. LLM outputs action name + input
2. Runtime executes the action
3. Result returned as "observation"
4. LLM continues with new context

The Problem with Non-Standardized Tools

Each AI application has its **own format** for tool integration

```
1  # App A format
2  {"action": "search_logs", "params": {"query": "error", "time": "1h"}}
3
4  # App B format
5  <tool name="search_logs"><param name="query">error</param></tool>
6
7  # App C format
8  USE_TOOL(search_logs, query=error, timeframe=1h)
```

Result: Every tool needs custom integration for every AI app

The Schema Fragmentation Problem

Each LLM Provider Has Different Tool Formats

The Challenge:

Every tool developer must maintain **N different implementations** of the same tool

OpenAI Function Calling Schema

Tool Definition

```
1  {
2    "type": "function",
3    "function": {
4      "name": "search_elasticsearch",
5      "description": "Search Elasticsearch",
6      "parameters": {
7        "type": "object",
8        "properties": {
9          "index": {
10             "type": "string",
11             "description": "Index name"
12           },
13          "query": {
14            "type": "string"
15          }
16        },
17        "required": ["index", "query"]
18      }
19    }
20  }
```

API Call

```
1  from openai import OpenAI
2
3  client = OpenAI()
4
5  response = client.chat.completions.create(
6      model="gpt-4",
7      messages=[
8          {"role": "user",
9           "content": "Search for errors"}
10     ],
11     tools=[tool_definition],
12     tool_choice="auto"
13 )
```

Nested structure with "function" wrapper and "parameters" object

Anthropic (Claude) Tool Schema

Tool Definition

```
1  {
2    "name": "search_elasticsearch",
3    "description": "Search Elasticsearch",
4    "input_schema": {
5      "type": "object",
6      "properties": {
7        "index": {
8          "type": "string",
9          "description": "Index name"
10       },
11       "query": {
12         "type": "string"
13       }
14     },
15     "required": ["index", "query"]
16   }
17 }
```

API Call

```
1  import anthropic
2
3  client = anthropic.Anthropic()
4
5  response = client.messages.create(
6    model="claude-3-5-sonnet-20241022",
7    max_tokens=1024,
8    messages=[
9      {"role": "user",
10       "content": "Search for errors"}
11    ],
12    tools=[tool_definition]
13  )
```

Flatter structure with "input_schema" instead of "parameters"

Google (Gemini) Function Schema

Tool Definition

```
1  {
2    "name": "search_elasticsearch",
3    "description": "Search Elasticsearch",
4    "parameters": {
5      "type": "object",
6      "properties": {
7        "index": {
8          "type": "string",
9          "description": "Index name"
10       },
11       "query": {
12         "type": "string"
13       }
14     },
15     "required": ["index", "query"]
16   }
17 }
```

API Call

```
1  import google.generativeai as genai
2
3  model = genai.GenerativeModel(
4      'gemini-1.5-pro',
5      tools=[tool_definition]
6  )
7
8  response = model.generate_content(
9      "Search for errors"
10 )
```

Similar to OpenAI but without the "function" wrapper, uses "parameters" directly

Schema Differences Summary

Provider	Wrapper	Args Key	Other Differences
OpenAI	function object	parameters	Nested structure
Claude	None	input_schema	Flat structure
Gemini	None	parameters	Similar to OpenAI

The Problem

```
1  # Without a standard, you need:
2
3  def tool_for_openai():
4      return {"type": "function", ...}
5
6  def tool_for_claude():
7      return {"name": ..., "input_schema": ...}
8
9  def tool_for_gemini():
10     return {"name": ..., "parameters": ...}
```

3x the code, 3x the bugs, 3x the maintenance

Enter MCP

Model Context Protocol

An **agreed-upon standard** for AI applications and tools to communicate

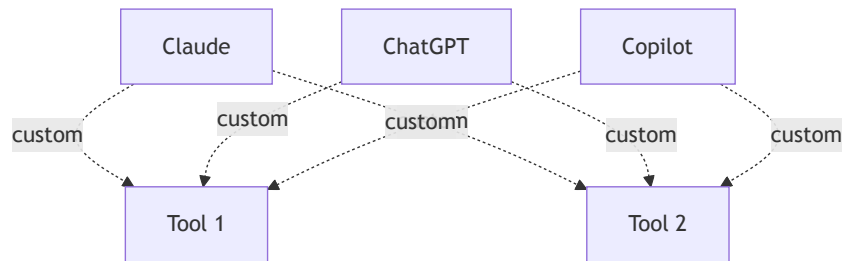
Think of it like:

- **HTML + HTTP** for the web
- **JDBC/ODBC** for databases

One standard. Universal compatibility.

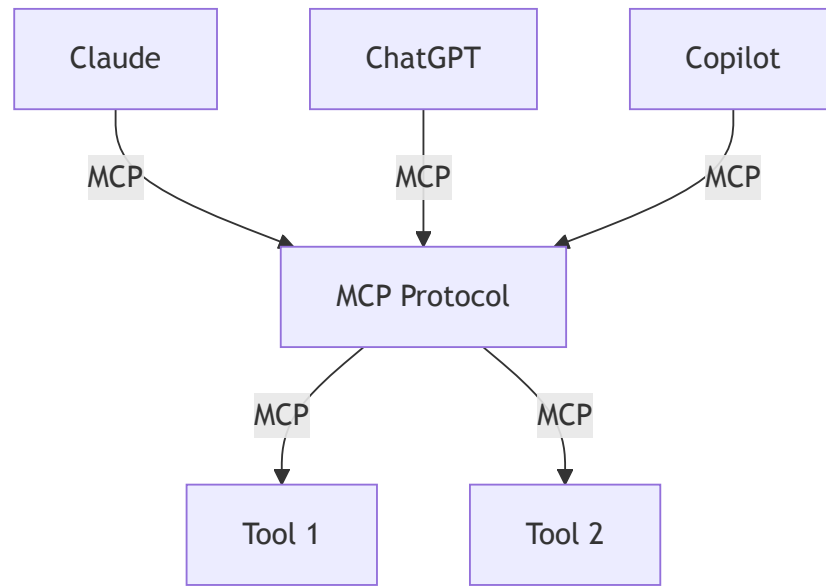
MCP = Interoperability

Without MCP



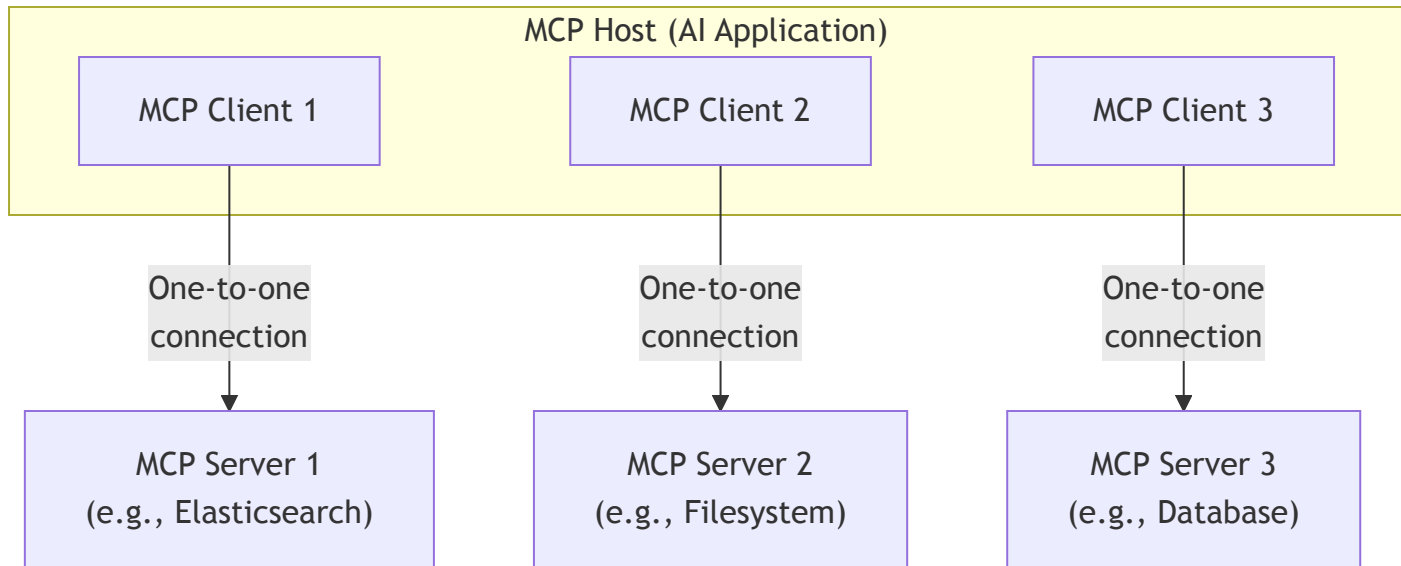
N×M integrations needed

With MCP



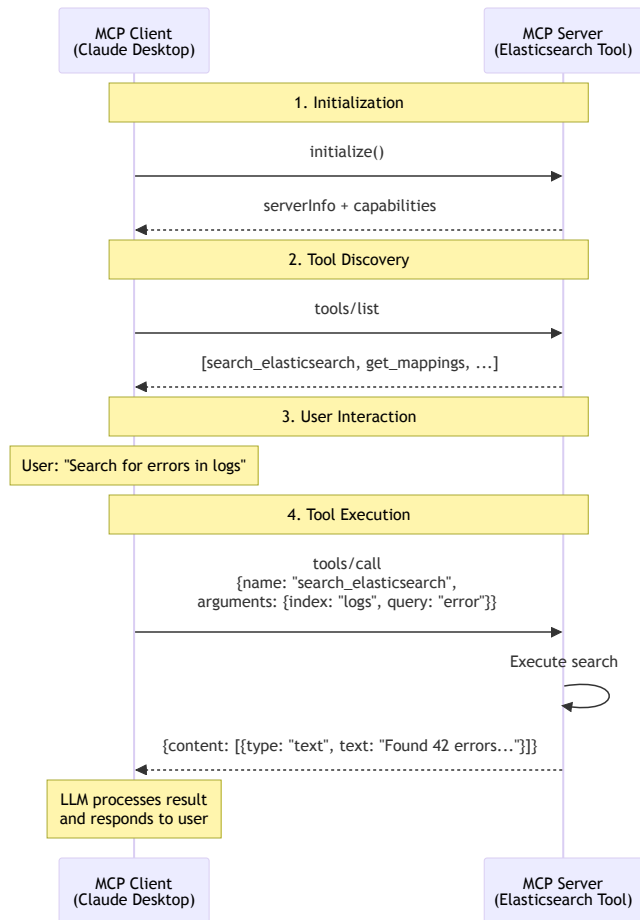
Plug and play compatibility

The Three Participants



- ◆ MCP Host: Claude Desktop, VS Code, your future AI coworker
- MCP Client: 1 per MCP server, Bridges host ↔ server.
- MCP Server: **Elasticsearch**, Figma, Github, any service with an MCP API

MCP Protocol Flow



MCP's Two Layers



Data Layer (Inner)

What messages are exchanged

Messages like:

- `tools/list` - Get available tools
- `tools/call` - Execute a tool
- `resources/list` - Get resources
- Session initialization



Transport Layer (Outer)

How messages are transmitted and exchanged.

JSON-RPC

JSON Remote Procedure Call

```
1  // Client Request
2  {
3    "jsonrpc": "2.0",
4    "method": "subtract",
5    "params": { "minuend": 42, "subtrahend": 23 },
6    "id": 1
7  }
8
9  // Server Response
10 {
11   "jsonrpc": "2.0",
12   "result": 19,
13   "id": 1
14 }
```

- Method name + parameters
- Requests have IDs (mapped to responses)
- Multiple requests can be in flight
- **Notifications** have no ID (no response expected)

JSON-RPC Notifications

Fire-and-forget messages with no ID

```
1  {
2    "jsonrpc": "2.0",
3    "method": "notifications/tools/list_changed"
4  }
```

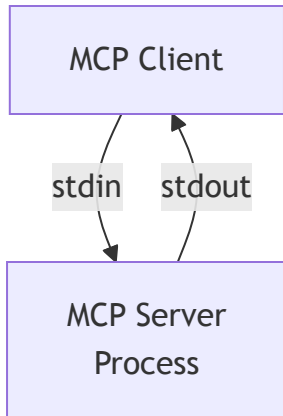
- Initialization complete
- Tools list updated
- Status changes
- Log messages

Transport Layer



Standard I/O

For local servers

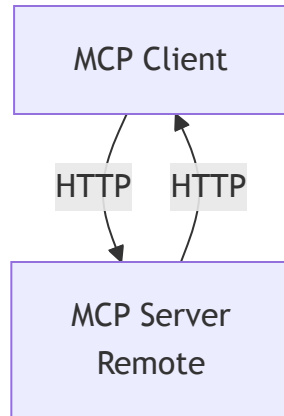


- Same machine
- Fast & secure
- No network needed
- Process communication



Streamable HTTP

For remote servers



- Over the network
- HTTP POST for client-to-server messages
- SSE for streaming and server notifications

MCP Primitives

The building blocks of functionality

Server Primitives

What MCP Servers offer to Clients

Tools

Actions/Functions Server operations with names, descriptions, parameters, and results.

Resources

Context Data to enrich LLM understanding Files, databases, or documents that enrich LLM context.

Prompts

Packaged Workflows and best practices Packaged workflows for guided task execution with built-in expertise.

Client Primitives

What MCP Clients offer to Servers

Sampling

LLM calls are made by the client on demand; server requests and user approves each interaction for control and transparency.

Elicitation

Server engages the user to fill in gaps—confirm actions, provide missing data, or set preferences before proceeding.

Security

Critical for remote servers and your ES cluster!

MCP emphasizes security with multiple authentication schemes:

```
1  Authorization: Bearer <token>
2  Authorization: ApiKey <elasticsearch-api-key>
3  Authorization: Basic <base64-credentials>
```

For Elasticsearch MCP Server:

- `ES_API_KEY` - Recommended approach
- `ES_USERNAME` + `ES_PASSWORD` - Basic auth

OAuth2.1

Latest MCP authorization relies builds on top of OAuth2

- Authorization Server Metadata (RFC8414)
- Dynamic Client Registration Protocol (RFC7591)
- Protected Resource Metadata (RFC9728)
- Resource Indicator (RFC8707)

Elasticsearch MCP in Action

Live demo time!

Introducing Elastic's MCP Server

Official implementation from Elastic

```
1  # Docker image
2  docker.elastic.co/mcp/elasticsearch
3
4  # Supports MCP protocols
5  - stdio (local)
6  - SSE (deprecated)
7  - streamable-HTTP (recommended)
8
9  # Elasticsearch versions
10 - 8.x ✅ Officially supported
11 - 9.x ✅ Officially supported
12 - 7.x ⚠️ May work, no guarantees
```

Status: ⚠️ EXPERIMENTAL (but fully functional!)

Setup: Docker Command

Running the Elasticsearch MCP Server

```
1  # For stdio mode (local)
2  docker run -i --rm \
3    -e ES_URL="https://my-cluster.es.io:9200" \
4    -e ES_API_KEY="your-api-key-here" \
5    docker.elastic.co/mcp/elasticsearch \
6    stdio
7
8  # For HTTP mode (remote)
9  docker run --rm \
10    -e ES_URL="https://my-cluster.es.io:9200" \
11    -e ES_API_KEY="your-api-key-here" \
12    -p 8080:8080 \
13    docker.elastic.co/mcp/elasticsearch \
14    http
```

Environment Variables:

- `ES_URL` - Your cluster URL
- `ES_API_KEY` - API key (recommended) or...
- `ES_USERNAME` + `ES_PASSWORD` - Basic auth

Claude Desktop Configuration

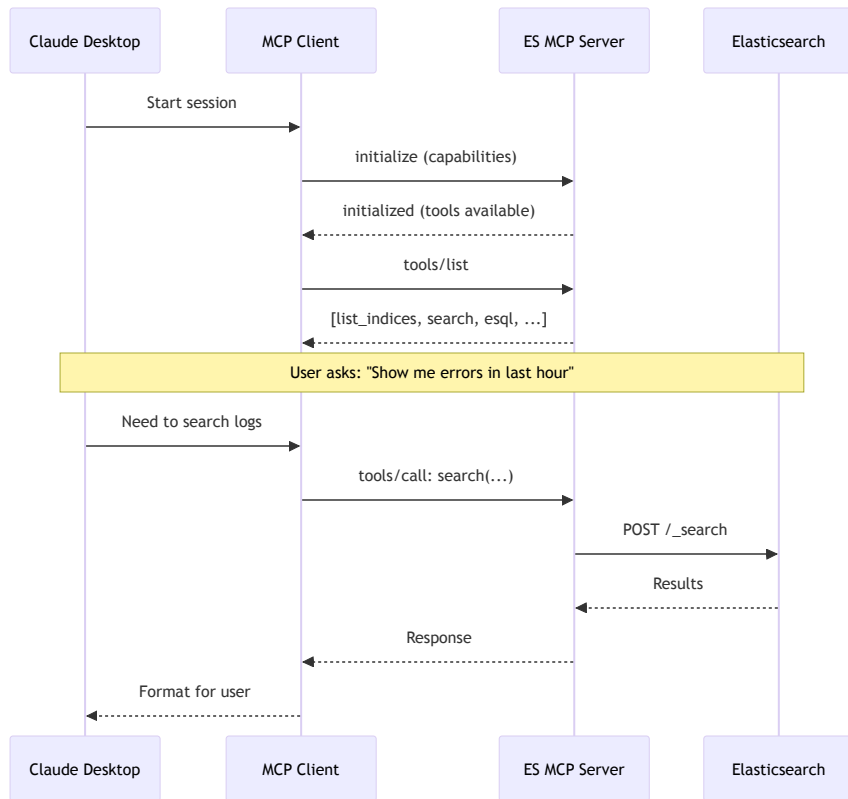
Connecting Claude to Elasticsearch

```
1  {
2    "mcpServers": {
3      "elasticsearch": {
4        "command": "docker",
5        "args": [
6          "run", "-i", "--rm",
7          "-e", "ES_URL",
8          "-e", "ES_API_KEY",
9          "docker.elastic.co/mcp/elasticsearch",
10         "stdio"
11       ],
12       "env": {
13         "ES_URL": "https://my-cluster.es.io:9200",
14         "ES_API_KEY": "VnVhQ2ZHY0JDZGJrU ... "
15       }
16     }
17   }
18 }
```

Location: ~/Library/Application Support/Claude/claude_desktop_config.json

Restart Claude Desktop and you're ready! 🚀

The Complete Exchange Flow



Real Example: Initialization

Step 1: Client initializes session

```
1  {
2    "jsonrpc": "2.0",
3    "method": "initialize",
4    "id": 1,
5    "params": {
6      "protocolVersion": "2025-06-18",
7      "capabilities": {
8        "roots": { "listChanged": true },
9        "sampling": {}
10     },
11     "clientInfo": {
12       "name": "Claude Desktop",
13       "version": "0.7.4"
14     }
15   }
16 }
```

Server responds:

```
1  {
2    "jsonrpc": "2.0",
3    "id": 1,
4    "result": {
5      "protocolVersion": "2025-06-18",
6      "capabilities": {
7        "tools": { "listChanged": true },
8        "resources": {}
9      },
10     "serverInfo": {
11       "name": "elasticsearch-mcp-server",
12       "version": "0.4.0"
13     }
14   }
15 }
```

Real Example: Getting Tools

Step 2: Client requests available tools

```
1  {  
2    "jsonrpc": "2.0",  
3    "method": "tools/list",  
4    "id": 2  
5  }
```

Real Example: Getting Tools (2)

Server responds with tool definitions:

```
1  {
2    "jsonrpc": "2.0",
3    "id": 2,
4    "result": {
5      "tools": [
6        {
7          "name": "list_indices",
8          "description": "List all available Elasticsearch indices",
9          "inputSchema": { "type": "object", "properties": {} }
10       },
11       {
12         "name": "search",
13         "description": "Perform an Elasticsearch search",
14         "inputSchema": {
15           "type": "object",
16           "properties": {
17             "index": { "type": "string" },
18             "body": { "type": "object" }
19           },
20           "required": ["index", "body"]
21         }
22       }
23     // ... more tools
```

Real Example: User Query

User asks: "Show me error logs from the last hour"

Step 3: LLM decides to list indices first

```
1  {
2    "jsonrpc": "2.0",
3    "method": "tools/call",
4    "id": 3,
5    "params": {
6      "name": "list_indices"
7    }
8  }
```

Response:

```
1  {
2    "jsonrpc": "2.0",
3    "id": 3,
4    "result": {
5      "content": [
6        { "index": "app-logs", "status": "open", "docs.count": 12 },
7        { "index": "products", "status": "open", "docs.count": 13 }
8      ]
9    }
10 }
```

Key Takeaways

1. **MCP is a protocol** - like HTTP/HTML for AI tools
2. **Two layers**: Data (JSON-RPC) + Transport (stdio/HTTP)
3. **Three participants**: Host (Claude) → Client → Server (Elasticsearch)
4. **Primitives**: Tools, Resources, Prompts / Sampling, Elicitation,
5. **Elasticsearch MCP is neat** -Talk to your ES using natural language. Unlock agentic workflows.
6. **Powerful capabilities** - search, ES|QL, mappings.

Get in touch: cefboud.com

—

Slides: github.com/CefBoud/slides