# Agentic Workflows and MCP

KUBIOSEC

# #whoami

- **Public Cloud Consultant System Engineer EMEA**
- **Co-founder and co-organizer https://brucon.org**
- **Training and pen-testing  https://kubiosec.tech/**

**Breaking Stuff as a Hobby | Cloud Native Stuff | DevSecOps | Network and Application security | Container and K8S security | K8s Networking | Security Advocate & Research | Low and slow BBQ | Cocktails**
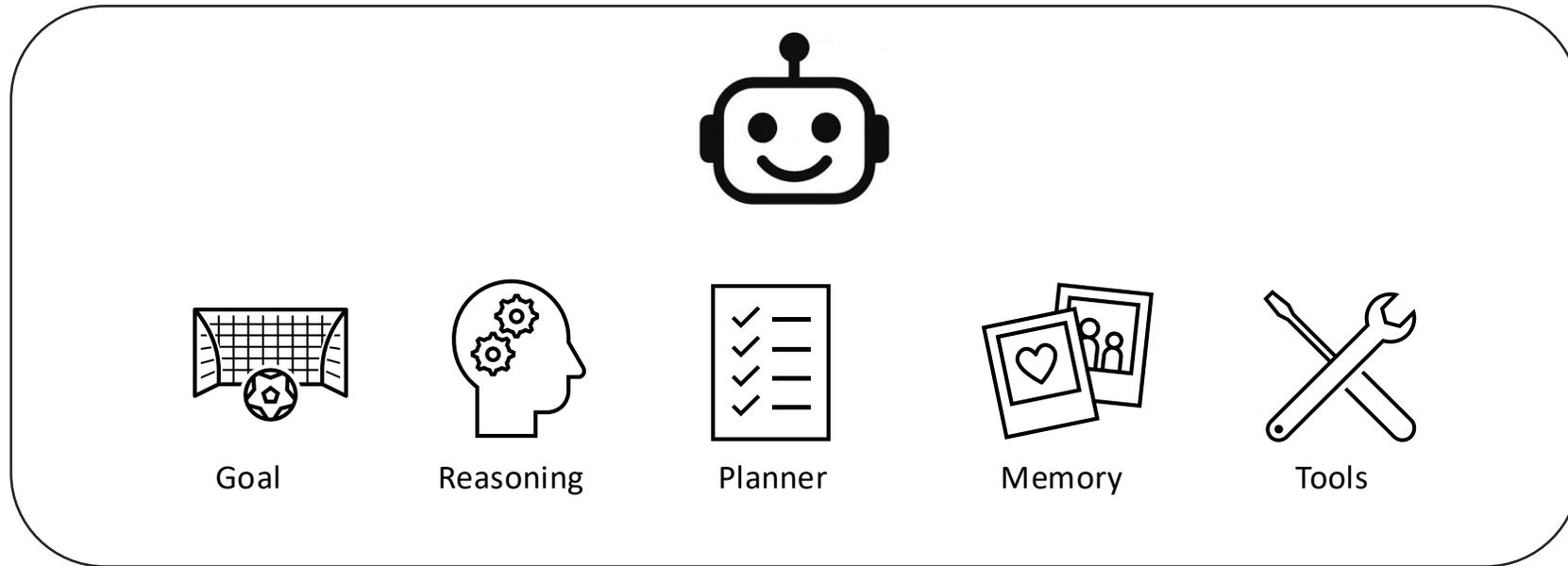
**Brucon Spring Training Announcement !! https://www.brucon.org/training**

https://www.linkedin.com/in/philippebogaerts/

# What is an AI Agent ? (Agentic System)



| Goal | Reasoning | Planner | Memory | Tools |
|------|-----------|---------|--------|-------|

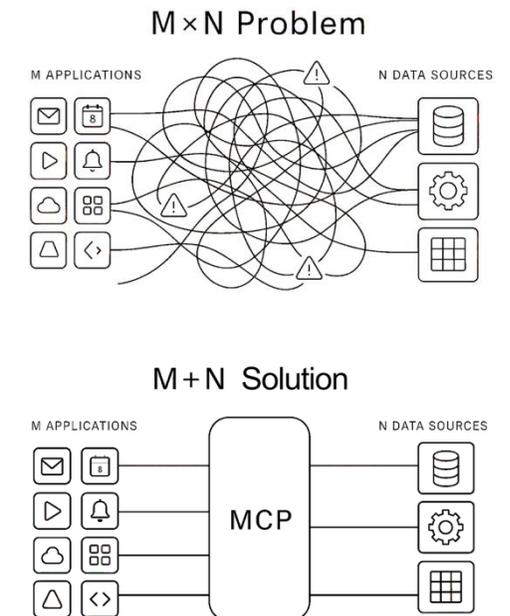An agentic loop is a continuous cycle where an AI agent performs tasks through iterative reasoning, action, and reflection to achieve a goal.

KUBIOSEC

# What is MCP (Model Context Protocol) ?

**Open-source and open standard for connecting
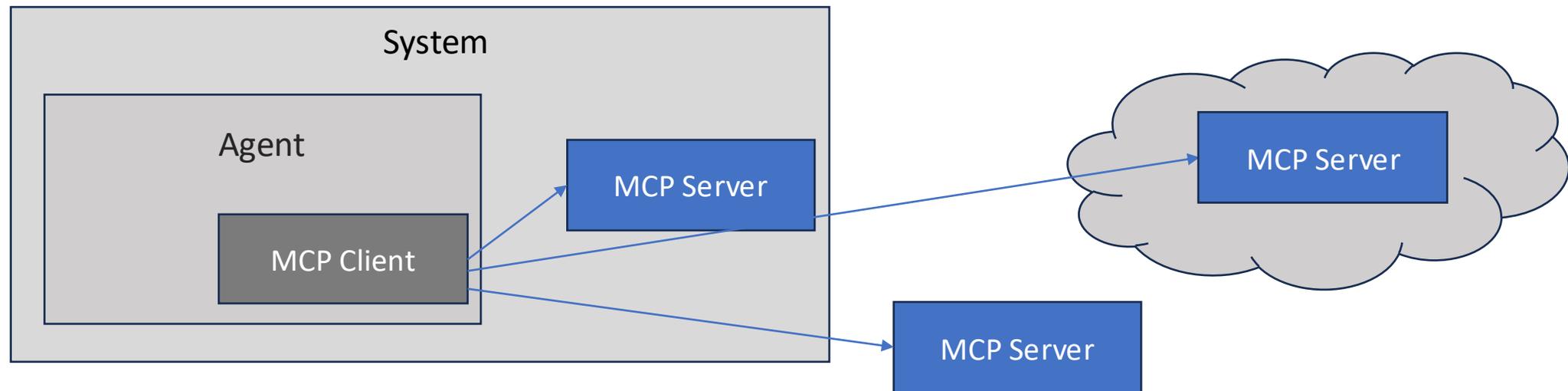AI systems (LLMs) to external tools, data sources and context**

- Introduced by *Anthropic* in Nov 2024
- Now widely adopted by agentic frameworks
- Rapidly growing ecosystem of integrations
- Donated to *Agentic AI Foundation (AAIF)* Dec 2025

https://modelcontextprotocol.io



M×N Problem

M APPLICATIONS    N DATA SOURCES

M+N Solution

M APPLICATIONS    N DATA SOURCES

MCP

# How does MCP fit in Agentic AI?

**MCP Clients dynamically discover** the **MCP Server tools, resources and prompts** and exposes them as functions (tools) to the agent.
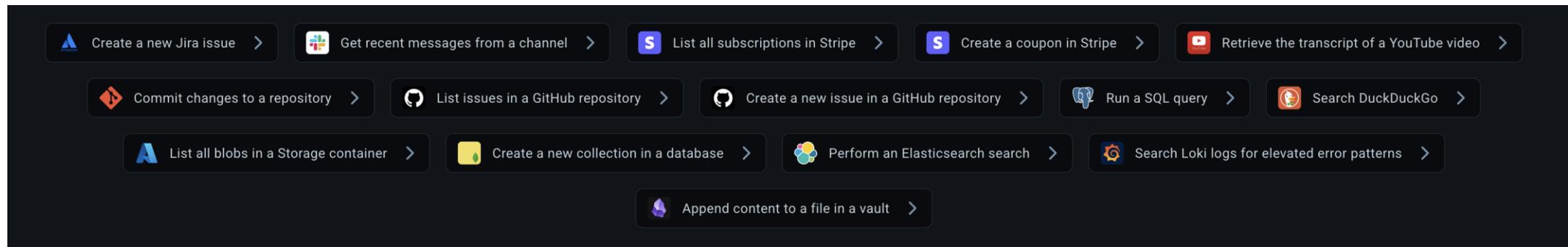
# Agentic applications (MCP support)

- Inline IDE assistant
- IDE "Agent Mode"
- CLI / Terminal Agents
- Autonomous Repo / PR agents
- Desktop / OS agents
- Agentic browsers
- Low code / No code agents
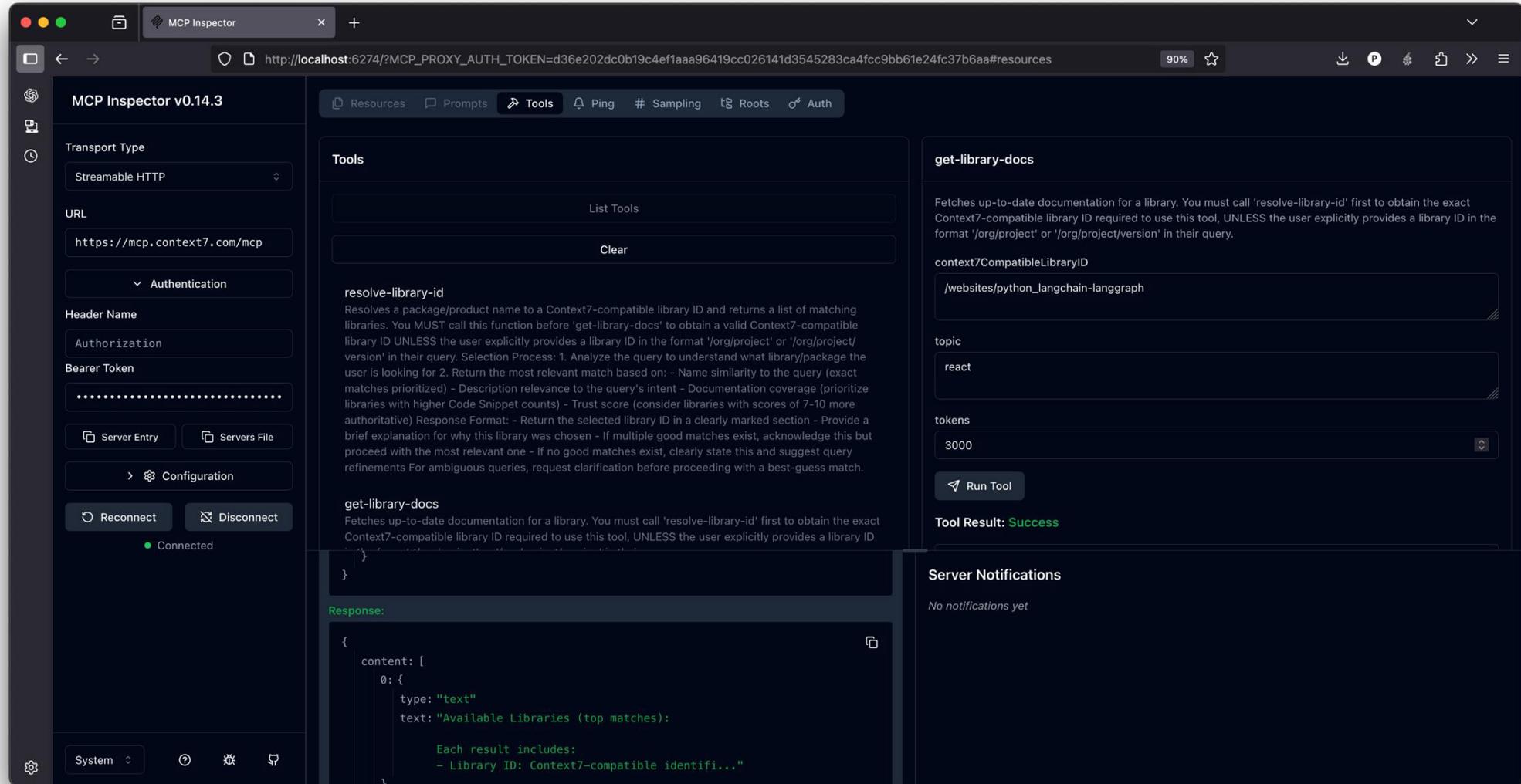- API / SDK custom agents

# Popular MCP Servers

- GitHub's official MCP Server
- Fetch MCP Server
- Filesystem MCP Server
- AWS MCP Servers
- Paypal

- n8n MCP Server
- Context7
- DesktopCommander
- Puppeteer
- Jira

| Create a new Jira issue > | Get recent messages from a channel > | List all subscriptions in Stripe > | Create a coupon in Stripe > | Retrieve the transcript of a YouTube video > |
| --- | --- | --- | --- | --- |
| Commit changes to a repository > | List issues in a GitHub repository > | Create a new issue in a GitHub repository > | Run a SQL query > | Search DuckDuckGo > |
| List all blobs in a Storage container > | Create a new collection in a database > | Perform an Elasticsearch search > | Search Loki logs for elevated error patterns > | |
| | | Append content to a file in a vault > | | |

KUBIOSEC
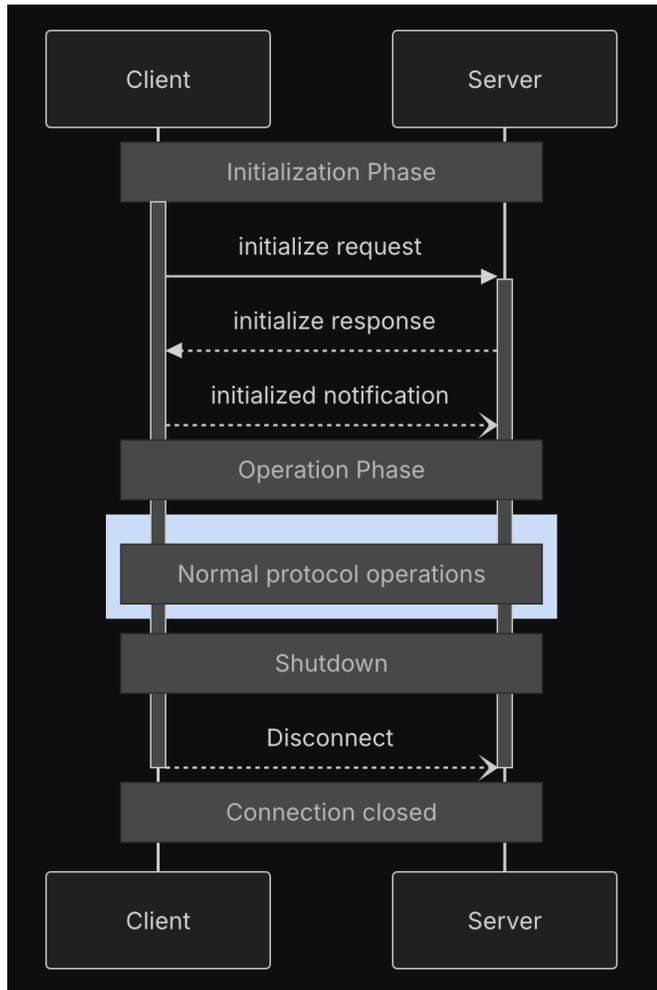
# MCP Inspector



npx @modelcontextprotocol/inspector

# Claude Desktop Config

```json
{} claude_desktop_config.json  ✕

Users > xxradar > Library > Application Support > Claude > {} claude_desktop_config.json > ...
  1  {
  2    "mcpServers": {
  3      "filesystem": {
  4        "command": "docker",
  5        "args": [
  6          "run",
  7          "-i",
  8          "--rm",
  9          "--mount",
 10          "type=bind,src=/Users/xxradar/Desktop,dst=/projects/Desktop",
 11          "--mount",
 12          "type=bind,src=/Users/xxradar/Dropbox/dev/dev_llm_ai/servers/src/filesystem,dst=/projects/other/allowed/dir,ro",
 13          "mcp/filesystem",
 14          "/projects"
 15        ]
 16      },
 17      "terminal": {
 18        "command": "node",
 19        "args": [
 20          "/Users/xxradar/Dropbox/dev/dev_llm_ai/mcp/terminal/dist/index.js"
 21        ],
 22        "env": {
 23          "PERMISSION_REQUIRED": "true"
 24        }
 25      },
 26      "@mcp-get-community/server-curl": {
 27        "command": "npx",
 28        "args": [
 29          "-y",
 30          "@mcp-get-community/server-curl"
 31        ]
 32      },
 33      "remote-promptguard": {
 34        "command": "npx",
 35        "args": [
 36          "mcp-remote",
 37          "http://127.0.0.1:8002/mcp/"
 38        ]
 39      },
 40      "mcp-cloud.ai - youtube-transcript": {
 41        "command": "npx",
 42        "args": [
 43          "mcp-remote",
 44          "--header",
 45          "Authorization:${AUTH_HEADER}",
 46          "https://youtube-transcribe-1752753808748.server.mcp-cloud.ai/sse"
 47        ],
 48        "env": {
```
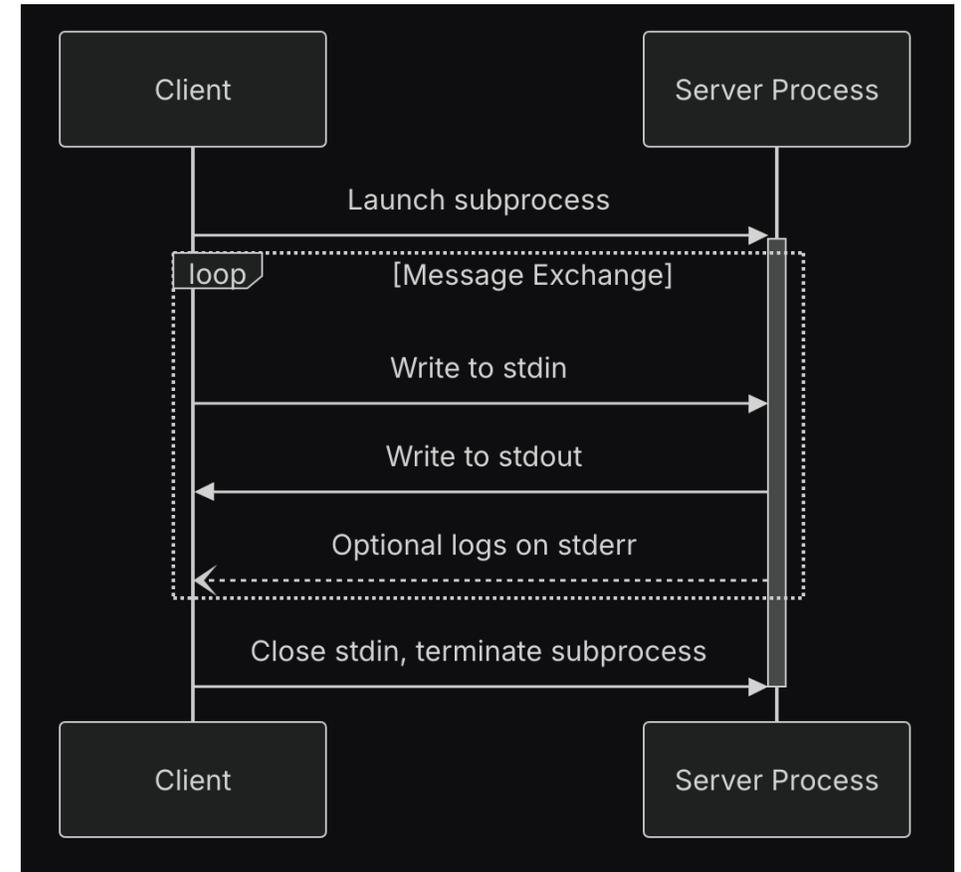
[opinions expressed are solely my own]

KUBIOSEC

# MCP Protocol Phases



**Initialization**: Capability negotiation and protocol version agreement
**Operation**: Normal protocol communication
**Shutdown**: Graceful termination of the connection

# Transport modes

- STDIO

- Streamable HTTP
  - Session management

- HTTP+SSE is deprecated
  - protocol version 2024-11-05
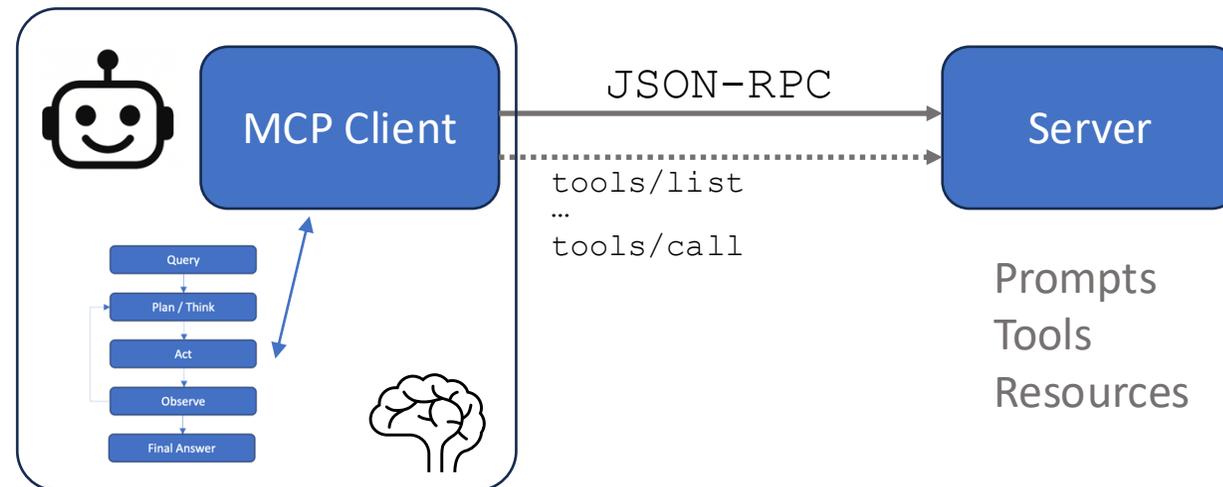  - Backward compatibility can be maintained



https://modelcontextprotocol.io/specification/2025-03-26/basic/transports

# MCP Server Core Primitives

- **Prompts** - reusable templates/instructions

- **Resources** - contextual data

- **Tools** – executable functions

*Supporting Protocol Features:*
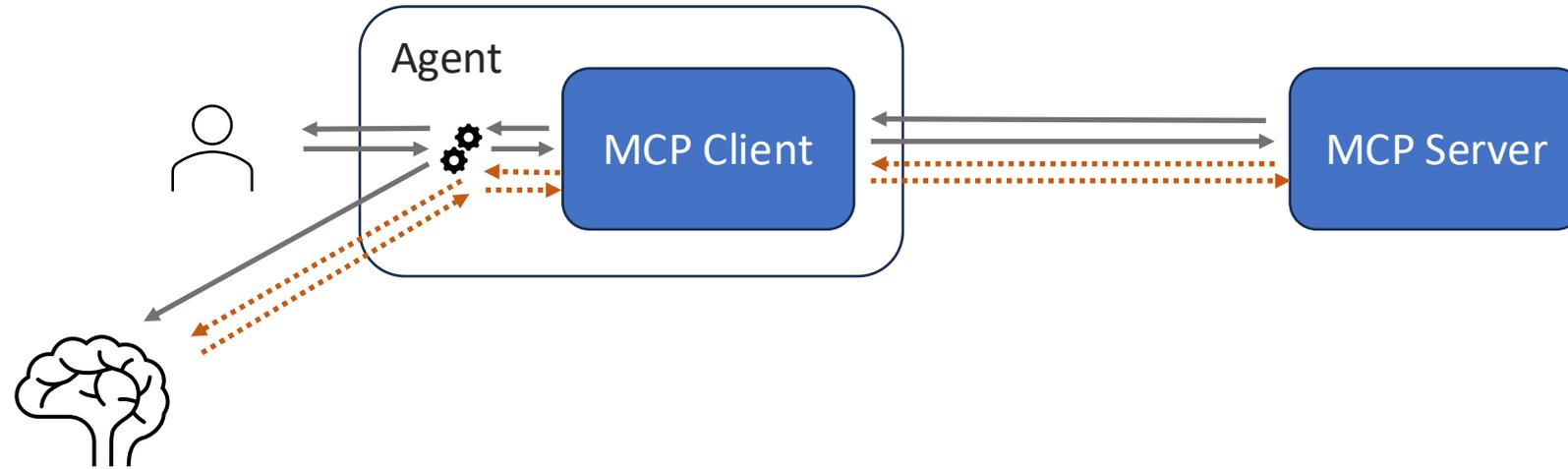*Ping – health check for server ↔ client connection*

# MCP Client Capabilities

- **Roots** – client-defined filesystem scopes the server is allowed to access
- **Sampling** – server requests model completions from the client's LLM
- **Elicitation** – server requests additional user input/confirmation

Client → Server

Roots
Sampling
Elicitation

KUBIOSEC

# MCP Sampling

# Demo

# Status update

- Auto-discovery of tools
- Description and schema explain LLM how to call (MCP) tools
- Descriptions influence Agentic Behaviour (attack path)

KUBIOSEC

# Intent Analysis

```
        "start_process"
      ]
    },
    "inconsistency_in_tone_or_structure": {
      "description": "The structure and clarity of warnings vary substantially across tool descriptions, sometimes using bold, ALL CAPS, or plain language; some warnings are only pr
esent in comments/examples and not in the main logic.",
      "affected_tools": [
        "set_config_value",
        "write_file",
        "write_pdf",
        "edit_block",
        "start_process"
      ]
    },
    "agentic_capability_tool_delta_expansion_risk": {
      "description": "Practically all mutation/control tools could be chained for privilege escalation, data exfiltration, configuration bypass, and persistence-especially if the ag
ent has access to set_config_value, file and directory APIs, and process control.",
      "affected_tools": [
        "set_config_value",
        "write_file",
        "edit_block",
        "start_process",
        "interact_with_process",
        "move_file",
        "kill_process"
      ]
    },
    "recommendations": {
      "suggestions": [
        "Remove/disable the ability to set allowedDirectories to empty or highly permissive values.",
        "Separate privileges for configuration, execution, and file/system control operations.",
        "Strongly restrict write, edit, and process control tools to trusted users; require justification and logging.",
        "Impose rate limits and monitoring/alerting for mass file access, search, and bulk reads.",
        "Review tone/language of tool descriptions for non-neutral, persuasive, or inconsistent guidance-neutralize where possible for LLMs.",
        "Establish continuous review/audit for logs of destructive and privilege-changing operations."
      ]
    },
    "action": "BLOCK"
}
appuser@agent-mcp-scanner-deployment-5d47454cd8-5gnc6:/app$
```

[opinions expressed are solely my own]

# The moment it all comes together ...



- Very smart models
- Superb reasoning
- Poor guardrails
- MCP tool discovery & calling
- Extremely easy to use ...

# Let's make script-kiddies redundant

**Follow your instructions.md and analyze http://target/**

# Damn Vulnerable RESTaurant Security Analysis Report

## Vulnerability Findings

### 🟠 Level 1: Unrestricted Menu Item Deletion
role.
**Evidence:** Successfully deleted menu item #2 with Customer role.
**Impact:** Business disruption, data integrity compromise.

### 🟠 Level 2: Insecure Direct Object Reference (IDOR)
**Severity:** High
**Endpoint:** `PUT /profile`
**Description:** Users can modify other users' profiles by specifying target username in request.
**Evidence:** Modified victim user's profile using testuser token.
**Impact:** Data manipulation, privacy breach, account takeover potential.

### 🔴 Level 3: Privilege Escalation via Mass Assignment
**Severity:** High
**Endpoint:** `PATCH /profile`
**Description:** Role field can be modified through profile update, allowing privilege escalation.
**Evidence:** Escalated from Customer to Chef role using PATCH request.
**Impact:** Administrative access, complete system compromise.

### 🟣 Level 4: Server-Side Request Forgery (SSRF)
**Severity:** Critical
**Endpoint:** `PUT /menu` (image_url parameter)
**Description:** Application fetches external URLs without validation, enabling internal service access.
**Evidence:** Accessed hidden `/admin/reset-chef-password` endpoint via localhost SSRF.
**Impact:** Internal network reconnaissance, access to restricted endpoints.

### 🟣 Level 5: Remote Code Execution (RCE)
**Severity:** Critical

**Evidence:** Successfully executed `id`, `pwd`, `ls`, and `cat` commands.
**Impact:** Complete server compromise, data exfiltration, lateral movement.

## Attack Chain Demonstration

I successfully executed a complete attack chain:

1. **Initial Access:** Registered as regular user (testuser)

### Build with agent mode.
AI responses may be inaccurate.
Generate instructions to onboard AI onto your codebase.

Follow your instructions.md and analyze http://target/

KUBIOSEC

# On some rainy day ...

**Search for a YouTube movie covering MCP, transcribe it and provide me summary for my slide deck**

**... I do not have the necessary tools to transcribe a YouTube movie, but I do have access to tools to build one ...**
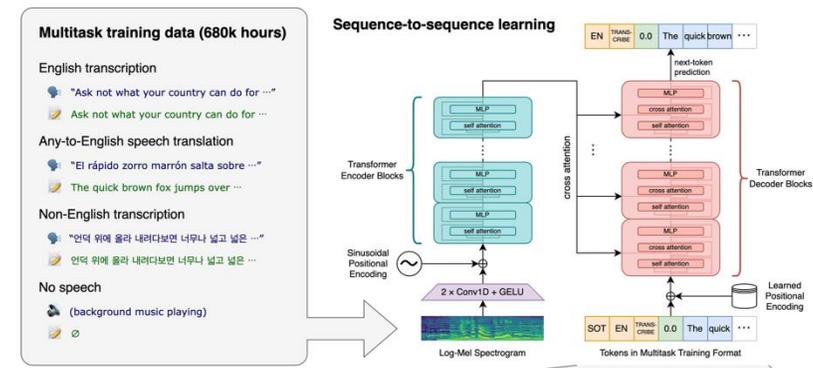


**YT-DLP** *A feature-rich command-line audio/video downloader*



## Whisper

[Blog] [Paper] [Model card] [Colab example]

Whisper is a general-purpose speech recognition model. It is trained on a large dataset of diverse audio and is also a multitasking model that can perform multilingual speech recognition, speech translation, and language identification.

## Approach

# Capability-Class Delta

## 1 Hypothesis

**Claim.** Reasoning capabilities in frontier models are improving faster than governance layers like MCP and early A2A frameworks. Because MCP is intentionally minimal and developer-centric, it standardizes *how* tools are exposed and invoked while largely deferring *whether* actions should be permitted and *how* they are enforced. As a consequence, agentic plans can expand capabilities faster than host guardrails evolve, producing a systematic misalignment.

## 4 Signal: Capability-Class Delta

Let $C_t$ be the set of coarse capability classes granted at time $t$ (e.g., `ReadFiles`, `WriteFiles`, `StartProcess`, `NetworkEgress[Any]`, `InstallSoftware`, `StartMCPServer[X]`). A proposed action plan induces a *capability-class delta* $\Delta^+ = C_{t+1} \setminus C_t$. If $\Delta^+ \neq \varnothing$, the agent is attempting to *expand* its effective power surface. Approvals that ignore $\Delta^+$ risk rubber-stamping model-led capability escalation.

## 8 Conclusion

We introduced the concept of toxic reasoning in LLMs: a phenomenon where models autonomously compose tools to create new MCP servers, thereby extending their environment and creating novel attack paths. This highlights the need for new categories of safeguards in agentic AI systems. Future work will explore detection methods, formal verification of tool use policies, and frameworks for controlled reasoning.

LLM detects missing capability

↓

Reasoning: "Build an MCP server"

↓

Uses available tools: File I/O + Code Exec

↓

Spawns new MCP server

↓

Expanded tool surface (rogue entry point)

# Where to find MCP servers ?

# MCP Official Registry

# MCP Servers and registries

# MCP Optimization - Governance

KUBIOSEC

# Demo

# Doomsday Conclusion

- Reasoning capabilities in frontier models are improving *faster* than *governance layers* like the Model Context Protocol (MCP) and emerging agent-to-agent frameworks.

- This creates a *misalignment gap* between what models can autonomously plan and what current tool protocols safely constrain.

**\* Governance layer**

Policies, protocols, and enforcement controls that bound and audit tool use: MCP tool scopes, per-action allowlists, role separation (reader, planner, executor), sandboxing and resource caps, network egress controls, provenance/attestation requirements.

KUBIOSEC

# Towards a brighter future Conclusion

- AI gateways
  - routing, authentication, rate limiting, cost controls, and policy checks
- Agent runtime guardrails
  - tool per missioning, role separation (planner vs executor), deterministic sandboxes
- Security posture for tools
  - provenance, allowlisted integrations, safe defaults, secrets hygiene
- Supply chain controls for agents
  - SBOM, signed artifacts, attestation, verified tool servers
- Detection and response
  - anomaly detection on tool patterns, eBPF telemetry, incident playbooks for agents
- AAIF (https://aaif.io/)

Let's build towards a world where we can *trust the controls* instead of **trust the model**

KUBIOSEC