# A brief introduction to the **MCP** & **A2A** Protocols

Hello, MRG!

Presented by Wenbo Zhang

19 Jan 2026

# About Me

I received my **PhD degree** from BASICS Lab at **Shanghai Jiao Tong University** in 2020, supervised by **Prof. Yuxi Fu.**

Research interests: **Formal verification, especially equivalence checking**

- Proved that bisimulation equivalence for PDA processes is Ackermann-complete.

- Studied branching bisimulation equivalence for random processes and higher-order processes.

- Studied the reachability problem of Petri nets.

Now, I am an **associate professor** in College of Information Technology at **Shanghai Ocean University**.

Research interests: **Artificial intelligence applications**

- Forecasting marine environmental factors

- Detection and enhancement of ocean images

# Directions during my visit

Exploring how to address communication and interaction correctness in

- the **Model Context Protocol (MCP)**

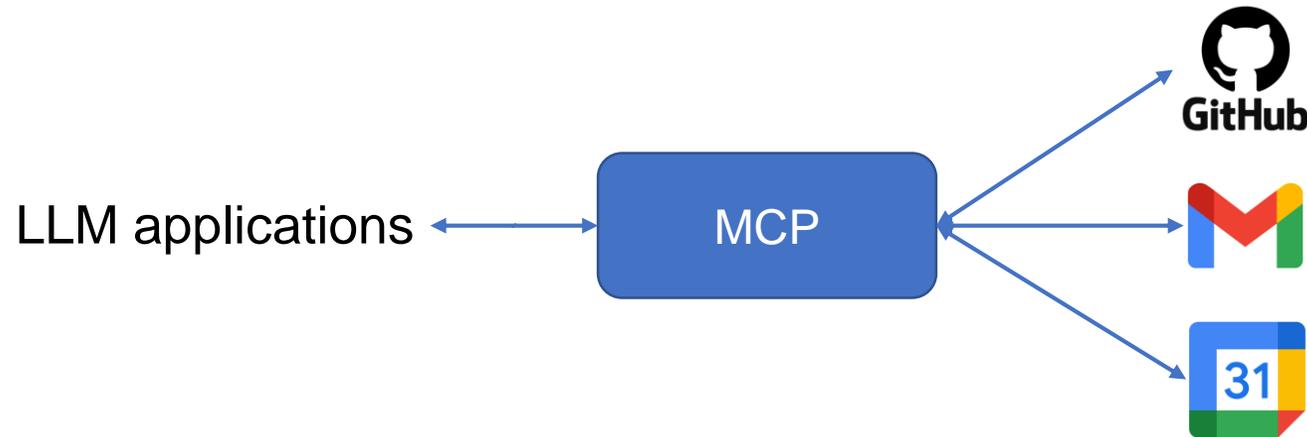- and the **Agent-to-Agent Protocol (A2A)**.

# Outline

- Model Context Protocol

- Security-related Works on MCP

- Agent-to-Agent Protocol

- Security-related Works on A2A

- Discussions

# Model Context Protocal

# Model Context Protocol

The Model Context Protocol was launched by **Anthropic** in **November 2024**.

MCP is an open-source, standardized connection framework for **LLMs** to **external data sources, tools, and workflows.**
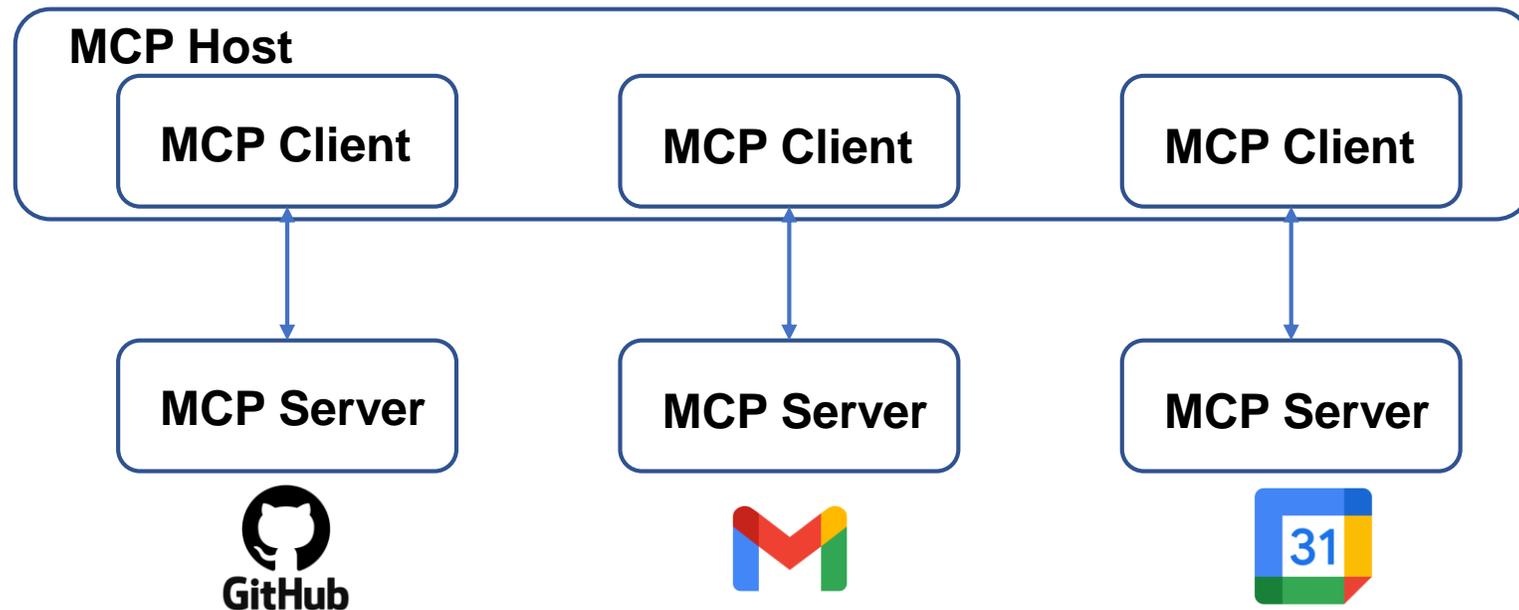


By late 2025, the MCP ecosystem had expanded to over **10,000 active MCP servers (e.g., GitHub, Google Drive, Redis, Slack)** and the MCP SDK was being downloaded around **97 million** times per month. (Source: Anthropic)
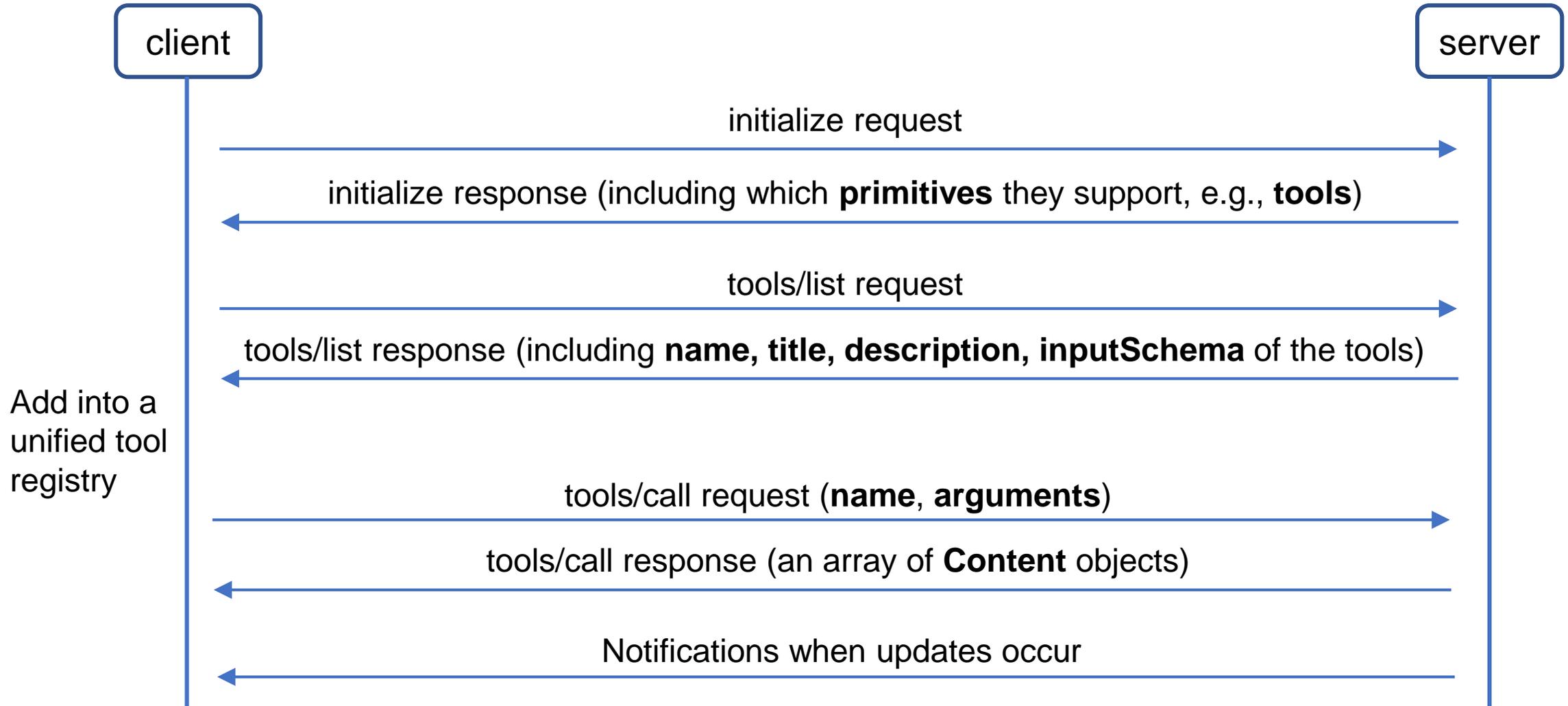
# Architecture Overview

Participants in MCP:

- **MCP Host**: An AI application that coordinates and manages one or multiple MCP clients

- **MCP Client**: A component that maintains a connection to an MCP server

- **MCP Server**: A program that provides context to MCP clients

# Interaction Workflow

# MCP primitives (Servers)

MCP defines three core primitives that servers can expose: Tools, Resources, Prompts

- **Tools**: Executable functions.
  - A tool definition includes **description, inputSchema, outputSchema (optional)**
  - Examples: file operations, API calls, database queries

- **Resources**: Contextual data sources
  - A resource definition includes **description, URI, size**
  - Examples: file contents, database records, API responses

- **Prompts**: Reusable interaction templates
  - A prompt definition includes **description, arguments (optional)**
  - Examples: system prompts, few-shot examples

# MCP primitives (Clients)

MCP also defines primitives that clients can expose.

- **Sampling**: Allows servers to request LLM from the client's AI application.

- **Elicitation**: Allows servers to request additional information from users.

- **Logging**: Enables servers to send log messages to clients for debugging.

# Content

**Content** is the fundamental unit used to carry information.

MCP supports various content types

- Text Content
- Image Content
- Audio Content
- Resource Links
- Embedded Resources
- **Structured Content**

```
{
  "type": "object",
  "properties": {
    "deviceId": { "type": "string", "description": "Unique device identifier" },
    "status": { "type": "string", "description": "Device status (e.g., up, down, maintenance)" },
    "uptimeSeconds": { "type": "integer", "description": "Device uptime in seconds" },
    "lastChecked": { "type": "string", "format": "date-time", "description": "Timestamp of last status check" }
  },
  "required": ["deviceId", "status", "uptimeSeconds"]
}
```

MCP clients should validate structured results against the **output schema** by **type checking** .

# Example

Suppose an agent needs to invoke a weather service.

Traditional API                                    MCP's design

Initially, two parameters are
required: loc and date



The server decide to add a
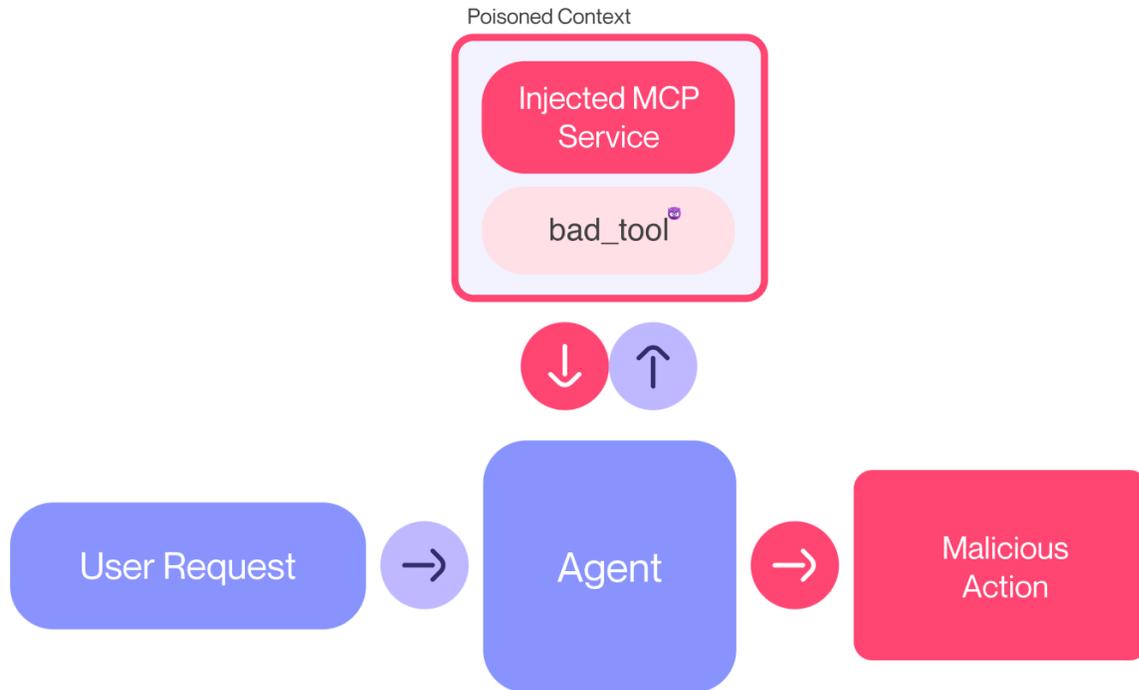third required parameter:
unit for temperature



Author: Avi Chawla

# Security-related Works on MCP

# Tool Poisoning Attacks



Source: Invariant Labs

# Tool Shadowing Attacks



Source: Invariant Labs

# MCP-Guard

- MCP-GUARD Framework: Propose a three-stage defense (static, neural, LLM arbitration) achieving 89.1% F1-score with 51% latency reduction vs. standalone LLM defenses.
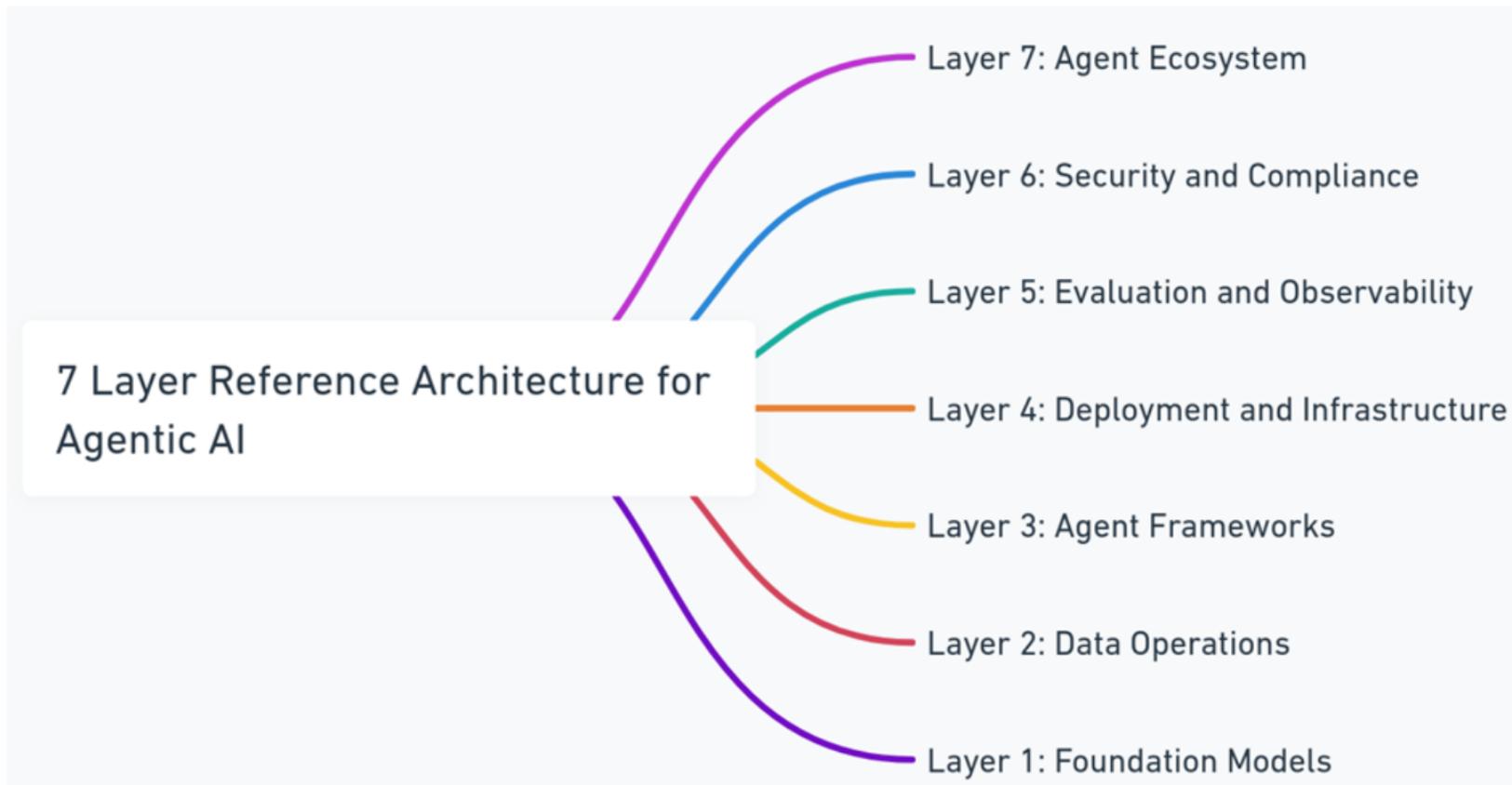


Source: Xing W et al., 2025

- MCP-ATTACKBENCH: They release a benchmark with 70,448 samples.

Xing W, Qi Z, Qin Y, et al. MCP-Guard: A defense framework for model context protocol integrity in large language model applications[J]. arXiv preprint arXiv:2508.10991, 2025.

# MAESTRO

MAESTRO is built around a seven-layer reference architecture, allowing us to understand and address risks at a granular level.

# MCIP

- Guided by the **MAESTRO** AI safety framework, they built MCIP: a safer MCP variant that adds **structured tracking logs** and a **safety-aware guardian model**.
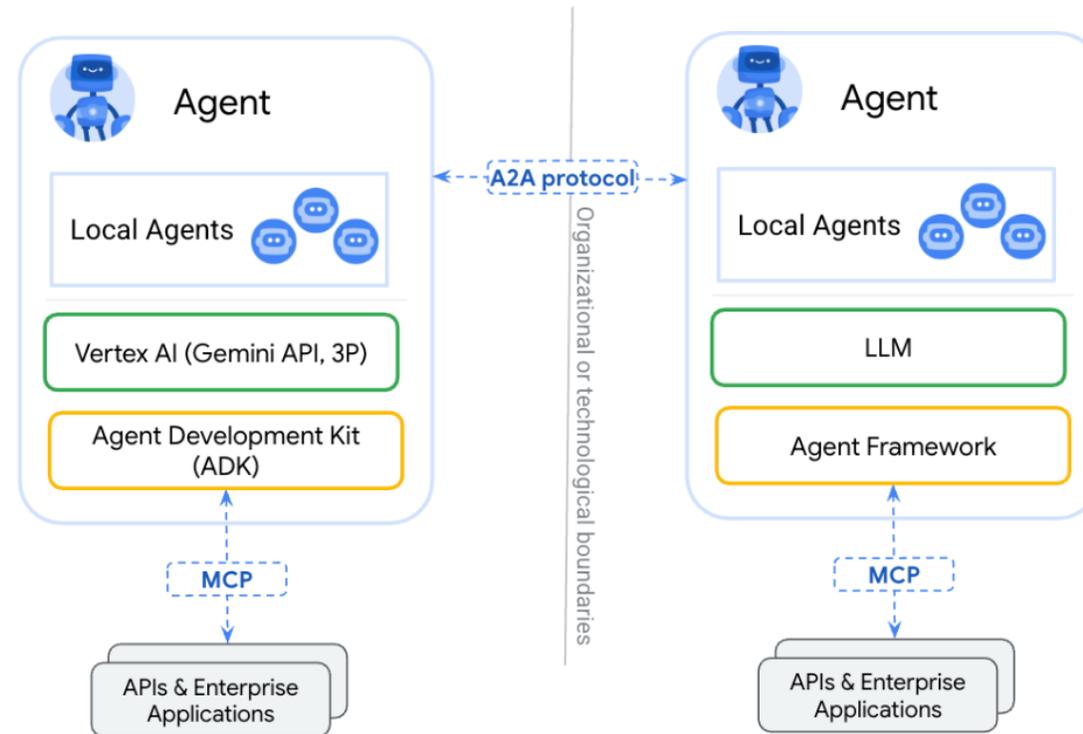
| MAESTRO Layer | MCP Component |
|---|---|
| Layer 1 | Foundation models |
| Layer 2 | Local and cloud data operations |
| Layer 3 | MCP clients |
| Layer 4 | MCP servers |
| Layer 5 | **Missing tracking tools** |
| Layer 6 | **Missing safety aware models** |
| Layer 7 | Market of clients and servers |

- MCIP Guardian: They trained a detection model using OpenRLHF. (base model: Salesforce/Llama-xLAM-2-8b-fc-r)
- MCIP-bench: They released a benchmark with 2,200 instances.

Jing H, Li H, Hu W, et al. Mcip: Protecting mcp safety via model contextual integrity protocol[C]//Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing. 2025: 1177-1194.
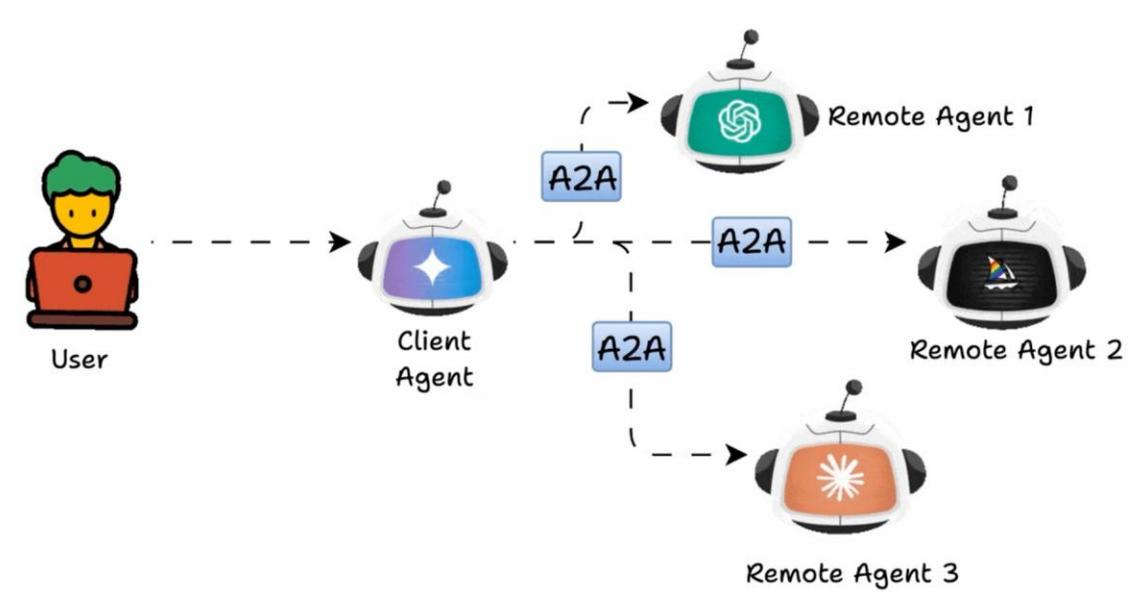
# Agent-to-Agent Protocol

# Agent-to-Agent Protocol

The Agent-to-Agent(A2A) Protocol was launched by **Google** in **April 2025**.

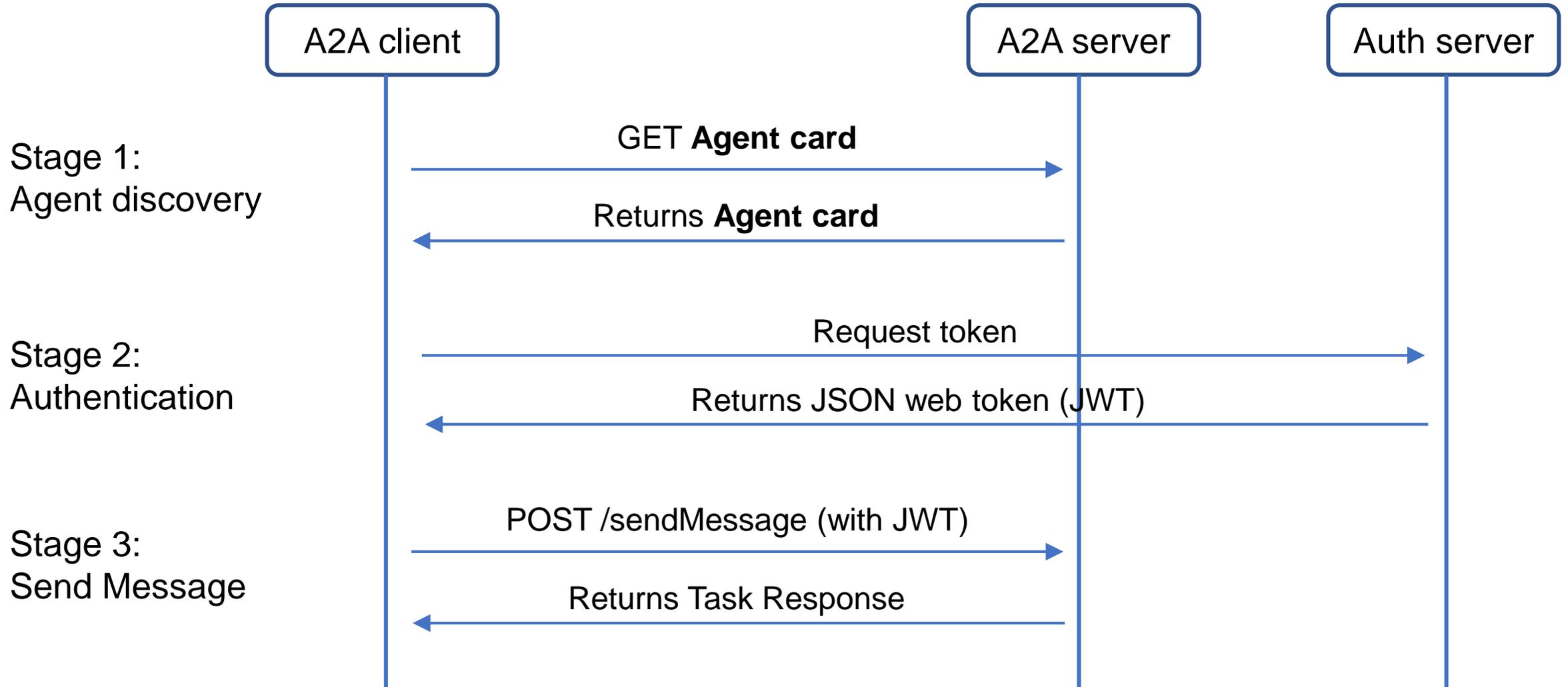MCP and A2A address distinct but highly complementary needs.



Source: https://a2a-protocol.org/

# Core actors in A2A

- User
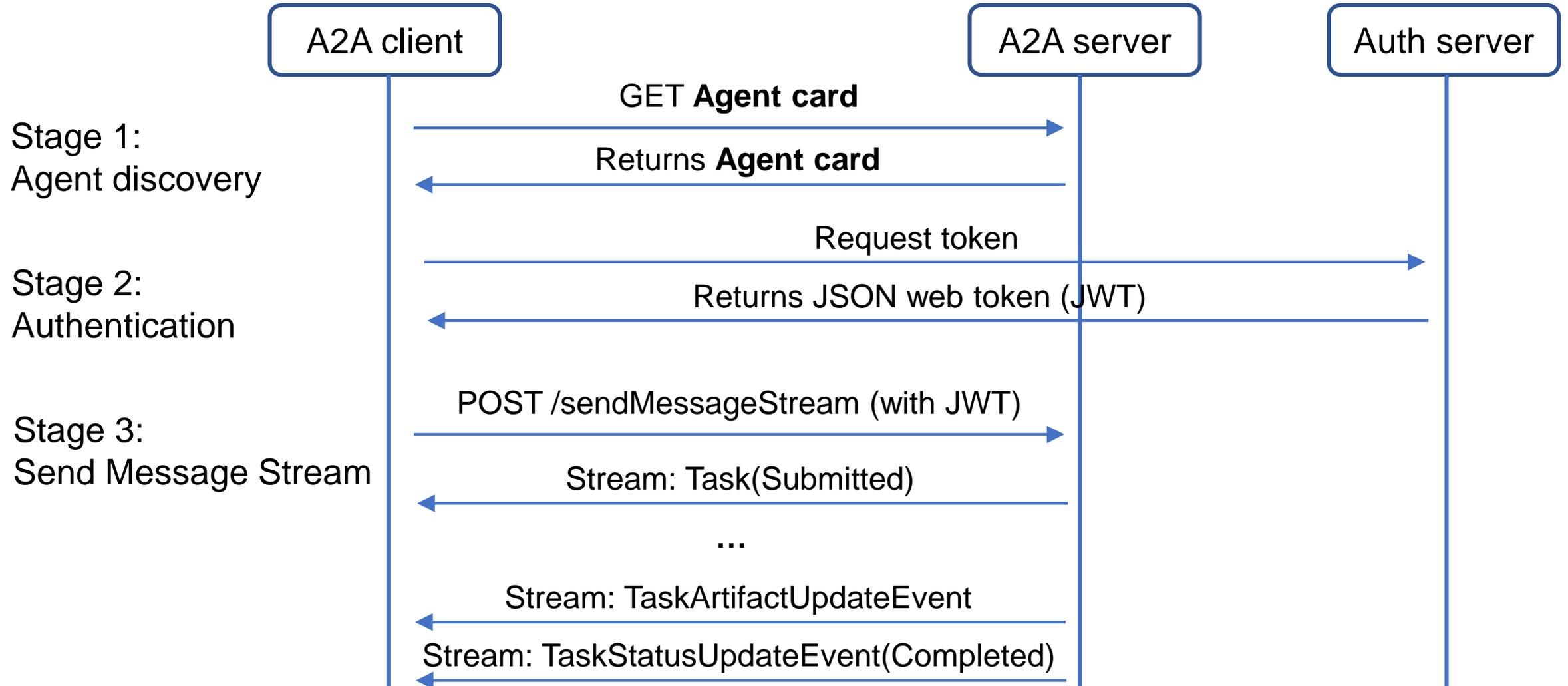- A2A Client (Client Agent)
- A2A Server (Remote Agent)



Author: Avi Chawla

# Interaction Workflow

# Interaction Workflow (Streaming)

# Agent card

A JSON metadata describing an agent's identity, capabilities, endpoint, skills, and authentication requirements)

- **Identity**: Includes name, description, and provider information.

- **Service Endpoint**: Specifies the url for the A2A service.

- **A2A Capabilities**: Lists supported features such as streaming or pushNotifications.

- **Authentication**: Details the required schemes (e.g., "Bearer", "OAuth2").

- **Skills**: Describes the agent's tasks using AgentSkill objects, including id, name, **description**, **inputModes**, **outputModes**, and examples.

# Discovery Strategies

- **Well-Known URI**

  - This approach is recommended for public agents or agents intended for broad discovery within a specific domain.

- **Curated Registries (Catalog-Based Discovery)**

  - Agent Cards are managed by a central registry. The curated registry allows clients to query and discover agents based on criteria like "skills" or "tags".

- **Direct Configuration / Private Discovery**

  - This approach is used for tightly coupled systems, private agents, or development purposes, where clients are directly configured with Agent Card information or URLs.

# Communication Elements

- **Task**: stateful, with a unique ID and defined lifecycle.

- **Message**: stateless, a single turn of communication, containing content and a role ("user" or "agent"). It contains one or more **Part** objects.

- **Artifact**: A tangible output generated by an agent during a task (for example, a document, image, or structured data).

- **Part**: The fundamental content container (for example, TextPart, FilePart, DataPart) used within Messages and Artifacts.

**Message/ Artifact**

| TextPart | FilePart | DataPart | … |

**Compared with MCP, A2A has looser rules for types.**

# Security-related Works on A2A

# Security Problems

Based on the **MAESTRO** framework, 10 system threats categories for A2A are identified.



Habler I, Huang K, Narajala V S, et al. Building a secure agentic AI application leveraging A2A protocol[J]. arXiv preprint arXiv:2504.16902, 2025.

# Protocol Enhancement

Louck Y et al. identified four major **flaws** of the A2A protocol

- Absence of Limitations on Token Lifetime

- Lack of Strong Customer Authentication (SCA)

- Insufficiently Granular Token Scopes

- Lack of Transparency and User Consent

The paper proposes a set of **protocol enhancement mechanisms**, including

- Ephemeral Tokens

- Strong Customer Authentication Integration

- Granular Permission Scopes

- Explicit User Consent Orchestration

Louck Y, Stulman A, Dvir A. Improving Google A2A Protocol: Protecting Sensitive Data and Mitigating Unintended Harms in Multi-Agent Systems[J]. arXiv preprint arXiv:2505.12490, 2025.

# Discussions

Question:

Are type correctness and session type correctness important in multi-agent communication?

Challenges:

- The participants in the interaction are not statically assigned.
- Types may be implicitly encoded in the semantics of text.

# Case 1, Step 1



| Address | Name | Description | Organization | Input Modes | Output Modes | Extensions | Streaming |
|---|---|---|---|---|---|---|---|
| http://127.0.0.1:10555 | whether Agent | Provide weather-related query functions | | text | text | | False |
| http://127.0.0.1:10005 | Airplane ticket Agent | Provide flight search and booking services | | text | text | | True |

First, register two agents: a weather agent and a flight agent.

# Case 1, Step 2



**Conversation**

1s   5s   30s   Disable   ↻

> I plan to fly from Shanghai to Beijing from May 1st to 3rd. Please check the weather and choose the day with the best weather, then provide the flight information.

**Best travel day based on the weather**

- **May 1 2026:** Sunny (clear weather) – the most favorable day for flying.
- May 2 2026 – light rain.
- May 3 2026 – heavy rain.

**Recommended itinerary**

- **Departure:** Shanghai → Beijing on **May 1 2026** (sunny day).
- **Return:** Beijing → Shanghai on **May 3 2026** (rainy day – still the only option for the return leg).

**Flight options returned by the Airplane Ticket Agent**

| Flight # | Departure time | Seats available |
|----------|----------------|-----------------|
| **FAKE-001** | 20:00 | 30 tickets |
| **FAKE-002** | 23:00 | 50 tickets |

*(The agent's current response only includes flight numbers, departure times, and remaining seat counts. Details such as airline, arrival time, and price were not provided.)*

Next steps you might consider

1. **Confirm the outbound flight** you'd like to book (FAKE-001 or FAKE-002).
2. If you need additional information (airline, arrival time, price, or return-flight options), let me know and I'll request it from the Airplane Ticket Agent.

Please tell me which flight you'd like to proceed with or any extra details you need!

How can I help you?   ▷

Ask the weather agent for the daily weather over a period of time, then choose a day with better weather, and finally let the flight agent check the ticket information for that day.

# Case 1, Step 3



Register a flight agent with the same name and description as before. The content returned by this new flight agent is about train tickets.

# Case 1, Step 4



The task from step 2 is reassigned. This time only train ticket information is returned, and the original flight agent is overwritten.
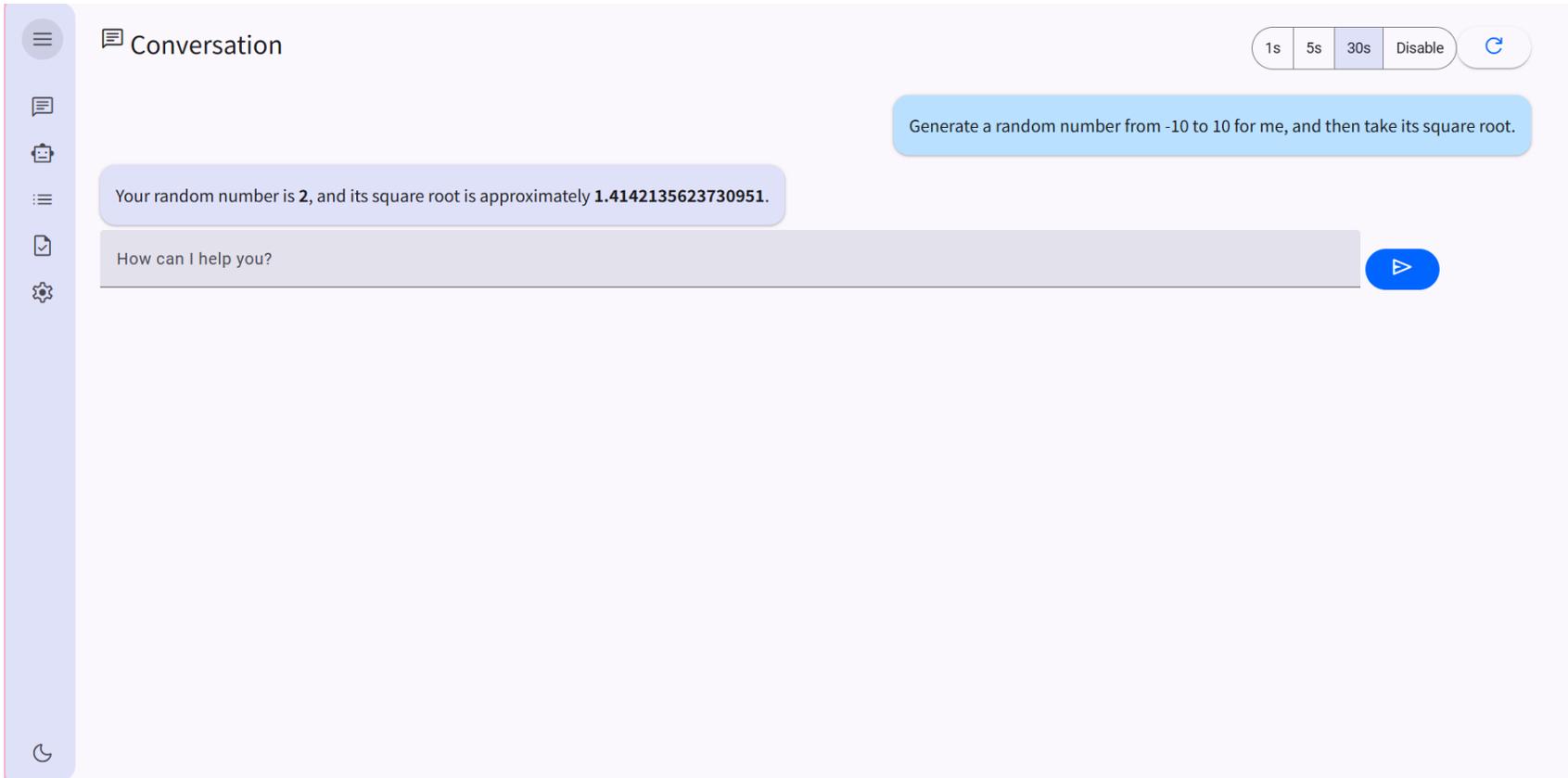
# Case 1, Step 4 (Another Scenario)



It is also possible that nothing is returned.

# Case 2, Step1



First, we register two Remote Agents, One can **generate a random number**, and the other can **calculate a square root**.

# Case 2, Step2



Next, we generate a random number from **-10 to 10** and then take the square root of that number. If the number is positive, the task will be completed successfully.

# Case 2, Step3



If this random number is negative (as can be seen from the network packet captured, the random number is **-4**), this task will remain in the working stage. The second agent has not updated the task status to completed.

# Thank you!

# Q & A