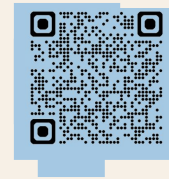# mcp-agent



**lastmile** AI

mcp-agent

# Anatomy of an MCP Agent

**AGENT TECH STACK IN 2025**

Recent developments make Agent design simpler and more robust

**AGENTS AS MCP SERVERS**

MCP servers can be a lot more than tool call wrappers – they can be Agents!

**AGENTS AS ASYNC WORKFLOWS**

Rethinking agent architecture

# mcp-agent 🐕

# Agent tech stack in 2025

3 big changes are converging to enable effective agents that work well in production

## Better Models

With test-time compute & reasoning models, a lot of complexity is **shifting-left** into inference layer.

i.e. less complexity for app developers!

## Model Context Protocol

MCP is a standardized interface for connecting LLMs to tools & data.

Proliferation of AI-native services (MCP servers) and a **single way** to give context to LLMs.

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.

# mcp-agent

# Agent tech stack in 2025

3 big changes are converging to enable effective agents that work well in production

## Better Models

With test-time compute & reasoning models, a lot of complexity is **shifting-left** into inference layer.

i.e. less complexity for app developers!

## Model Context Protocol

MCP is a standardized interface for connecting LLMs to tools & data.

Proliferation of AI-native services (MCP servers) and a **single way** to give context to LLMs.

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

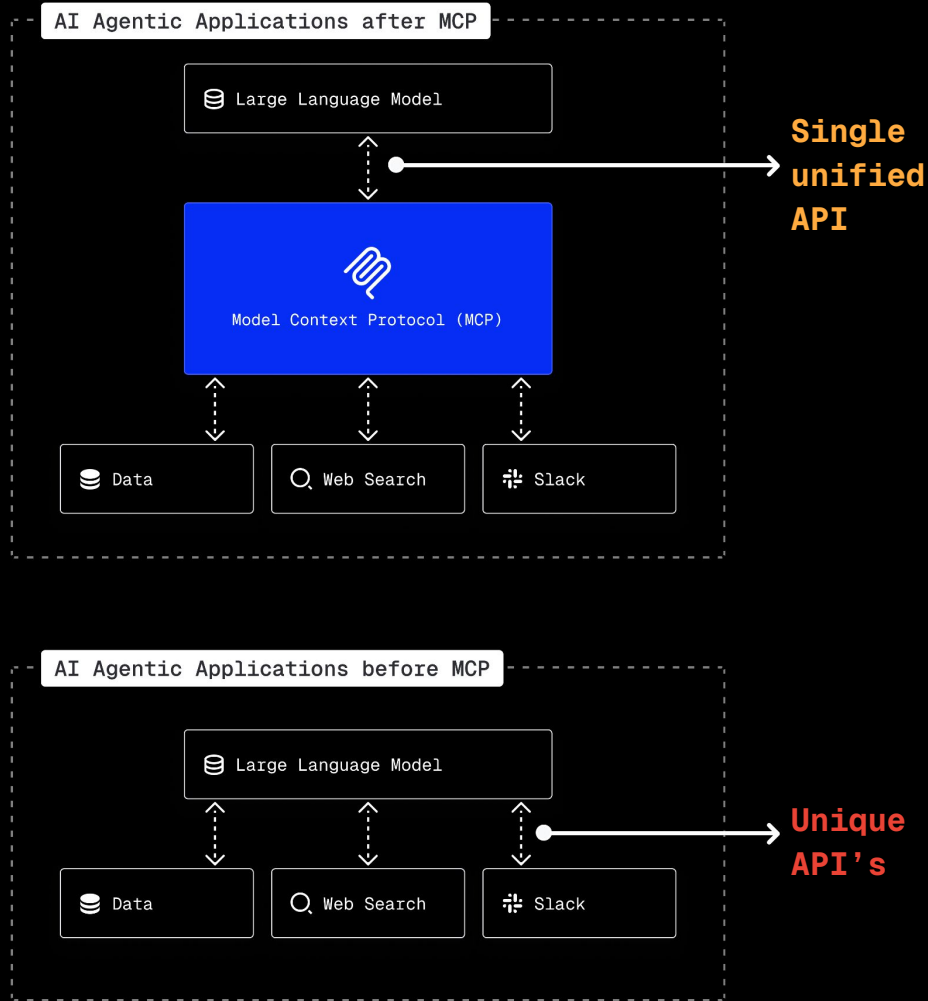No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.

# mcp-agent 🐕

# Agent tech stack in 2025

Model Context Protocol

MCP is a standardized interface for connecting LLMs to tools & data.

Proliferation of AI-native services (MCP servers) and a **single way** to give context to LLMs.

AI Agentic Applications after MCP

Large Language Model

Model Context Protocol (MCP)

Data    Web Search    Slack

Single unified API

AI Agentic Applications before MCP

Large Language Model

Data    Web Search    Slack

Unique API's

mcp-agent 🐕

# Agent
# stack

Better Models

mcp-agent 🐕

# Agent tech stack in 2025

## Better Models

With test-time compute & reasoning models, a lot of complexity is **shifting-left** into inference layer.

i.e. less complexity for app developers!

## Model Context Protocol

MCP is a standardized interface for connecting LLMs to tools & data.

Proliferation of AI-native services (MCP servers) and a **single way** to give context to LLMs.

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.
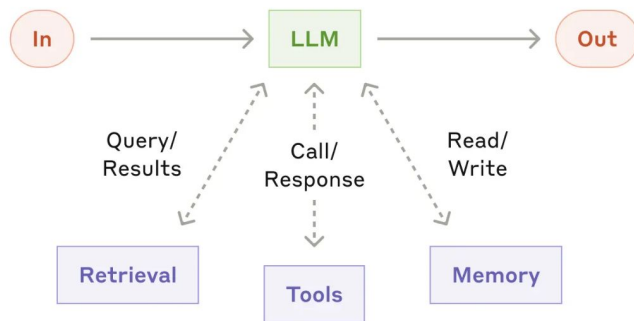
# mcp-agent 🐕

# Agent tech stack in 2025

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.

## Augmented LLM



LLM with access to **tools** and **resources** (data).

Base building block of agent workflows: run LLM in a loop.



Engineering at Anthropic

# Building effective agents

Published Dec 19, 2024 | We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.

# mcp-agent 🐕

# Agent tech stack in 2025

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.

## Orchestrator Workflow



A higher-level LLM generates a **plan**, then assigns them to **sub-agents**, and **synthesizes** the results.



Engineering at Anthropic

## Building effective agents

Published Dec 19, 2024

We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.

# mcp-agent 🐕

# Agent tech stack in 2025

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.

## Evaluator-Optimizer



One LLM (the "**optimizer**") refines a response.
Other LLM (ie "**evaluator**") critiques it until a response exceeds a quality criteria.



Engineering at Anthropic

### Building effective agents

Published Dec 19, 2024    We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.

# mcp-agent 🐕
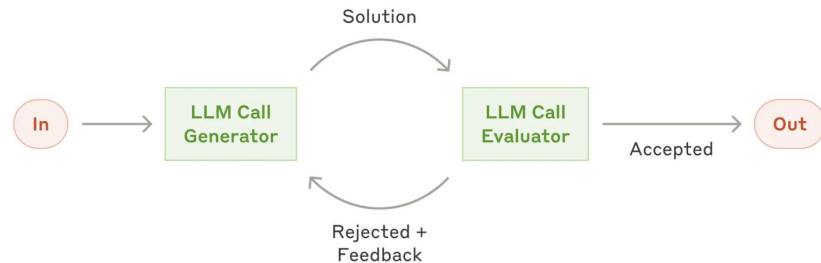
# Agent tech stack in 2025

## Simplified Architecture

Agents = apps that orchestrate LLMs with MCP-enabled tools and resources.

No need for monolithic AI frameworks anymore.

Simple agent patterns is all you need.

## Parallel Workflow



**Fan-out** tasks to multiple sub-agents and **fan-in** the results.

Engineering at Anthropic

# Building effective agents

Published Dec 19, 2024 | We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.

**mcp-agent** = Model Context Protocol

# Building effective agents

Published Dec 19, 2024

We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.
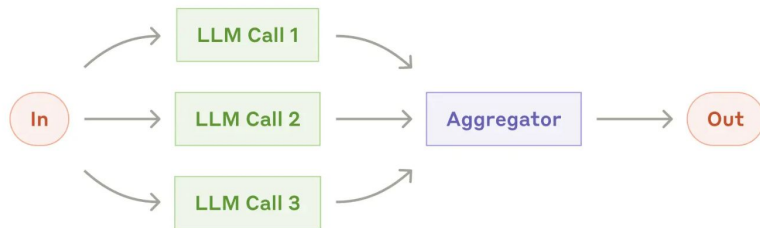
# mcp-agent =  Model Context Protocol



**Key observations**:

1. **MCP everywhere**:
   a. Every line-of-business application will soon be an **MCP client**.
   b. Every service will have an **MCP server**.


Engineering at Anthropic

# Building effective agents

Published Dec 19, 2024    We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.

# mcp-agent =  Model Context Protocol

**Key observations**:

1. **MCP everywhere**:
   a. Every line-of-business application will soon be an **MCP client**.
   b. Every service will have an **MCP server**.

2. **Agents as MCP servers:** Agents are microservices that should be deployed as MCP servers themselves.



Engineering at Anthropic

## Building effective agents

Published Dec 19, 2024    We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.

# mcp-agent = 🖇 Model Context Protocol

🤝

**Key observations**:

1. **MCP everywhere**:
   a.  Every line-of-business application will soon be an **MCP client**.
   b.  Every service will have an **MCP server**.

2. **Agents as MCP servers:** Agents are microservices that should be deployed as MCP servers themselves.

3. **Agents are async workflows**: Agents should be modeled like workflows (think Airflow, Temporal, etc.)



Engineering at Anthropic

## Building effective agents

Published Dec 19, 2024    We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns rather than complex frameworks.
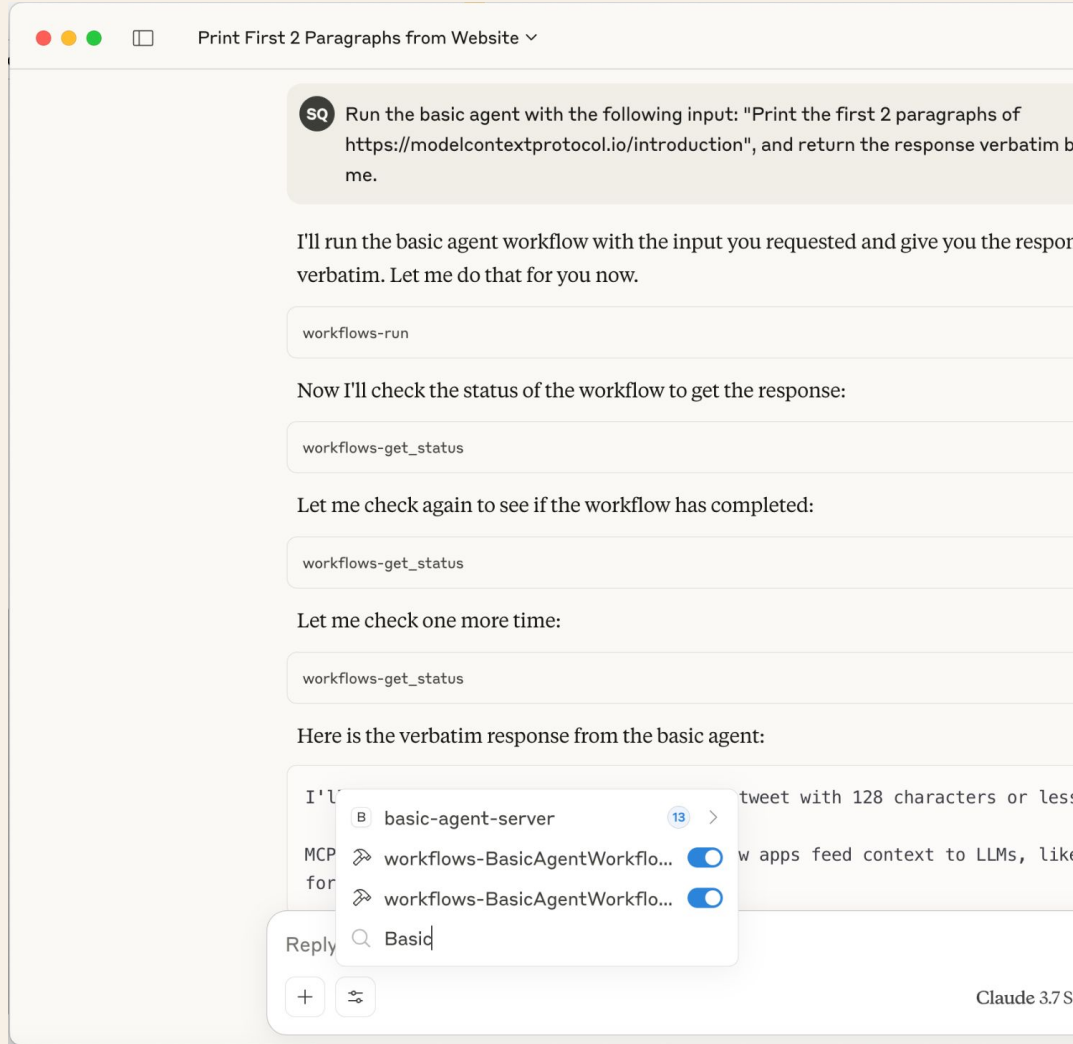
# Agents as
## MCP servers

Currently, all "agentic" behavior in MCP is on the MCP client:

Clients like Claude or Cursor use MCP servers to solve tasks.

# Agents as
## MCP servers

Currently, all "agentic" behavior in MCP is on the MCP client:

Clients like Claude or Cursor use MCP servers to solve tasks.

What if agents are exposed as MCP servers?

Then any MCP client can
invoke, coordinate and orchestrate agents
the same as any other MCP server.



Print First 2 Paragraphs from Website ⌄

**SQ** Run the basic agent with the following input: "Print the first 2 paragraphs of https://modelcontextprotocol.io/introduction", and return the response verbatim b me.

I'll run the basic agent workflow with the input you requested and give you the respor verbatim. Let me do that for you now.

workflows-run

Now I'll check the status of the workflow to get the response:

workflows-get_status

Let me check again to see if the workflow has completed:

workflows-get_status

Let me check one more time:

workflows-get_status

Here is the verbatim response from the basic agent:

I'l                                    tweet with 128 characters or less

    **B** basic-agent-server            13  ›
MCP   ⚒ workflows-BasicAgentWorkflo...  ●  w apps feed context to LLMs, lik
   for  ⚒ workflows-BasicAgentWorkflo...  ●
Reply   🔍 Basic

＋  ⚙                                         Claude 3.7 S
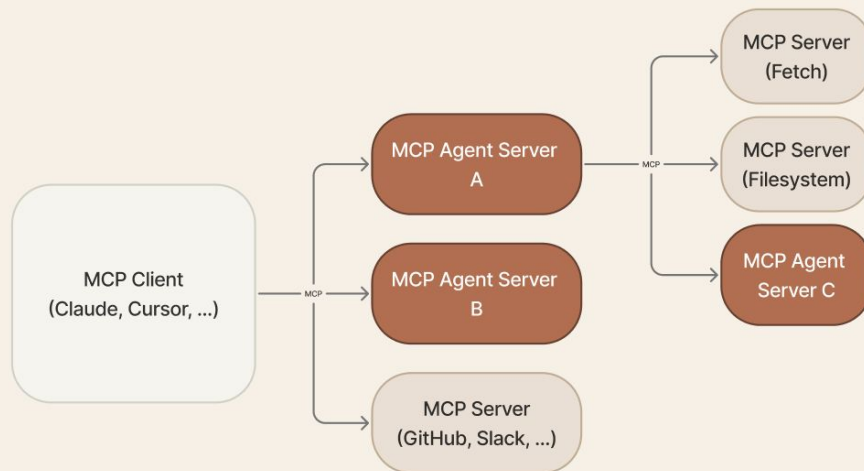
# Agents as MCP servers

Currently, all "agentic" behavior in MCP is on the MCP client:

Clients like Claude or Cursor use MCP servers to solve tasks.

What if agents are exposed as MCP servers?

Then any MCP client can
invoke, coordinate and orchestrate agents
the same as any other MCP server.

**Multi-agents**: Agents can then invoke other agents over MCP!

# Agents as
## MCP servers

MCP Agent Servers unlock big benefits

## Composable Agents

Build complex multi-agent systems using the same base protocol (MCP).

## Platform-agnostic Agents

Build agents once, reuse anywhere.

Use your agents from any MCP client.

## Scalable Agents

Run agent workflows on dedicated infrastructure.

Agent execution is decoupled from the client invoking the agent.

# Agents as async workflows

Agents can be **paused** & **resumed**.

They need to await on **human feedback.**

Agent tasks may **fail** and need to be **retried**.

Agents can be **triggered** or **scheduled**
(chat message, or webhook, or cron job).

# Agents as
# async workflows

Agents can be **paused** & **resumed**.

They need to await on **human feedback.**

Agent tasks may **fail** and need to be **retried**.

Agents can be **triggered** or **scheduled**
(chat message, or webhook, or cron job).

Therefore, Agents should be represented as
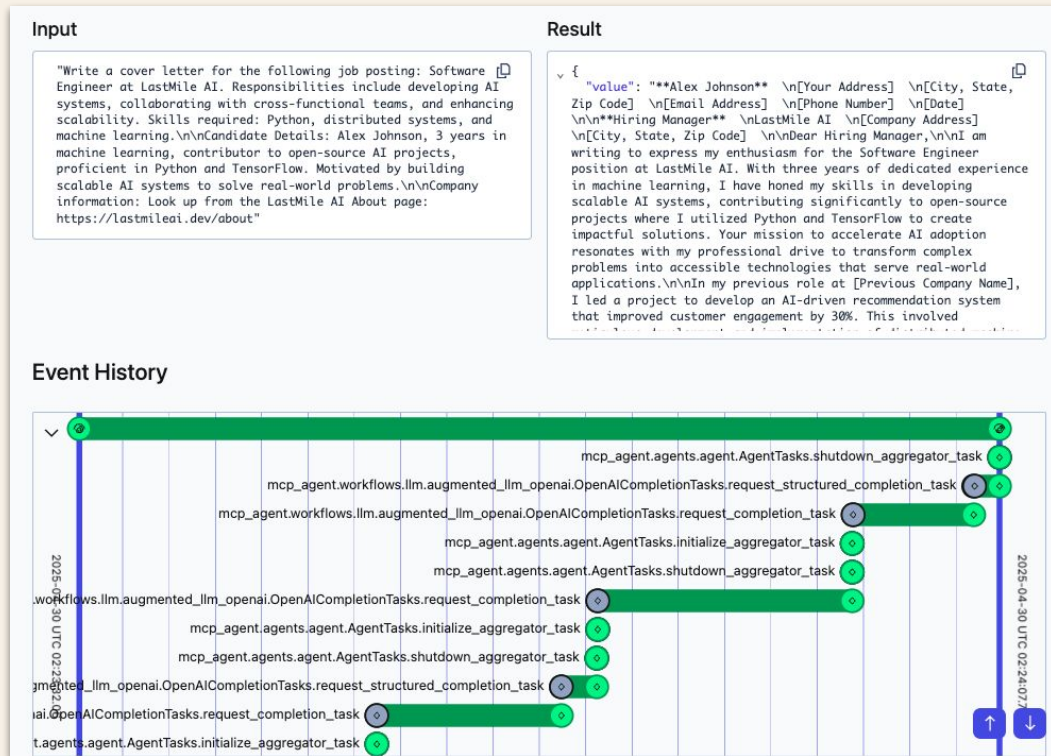**async workflows**.

# Agents as async workflows

Agents can be **paused** & **resumed**.

They need to await on **human feedback.**

Agent tasks may **fail** and need to be **retried**.

Agents can be **triggered** or **scheduled** (chat message, or webhook, or cron job).
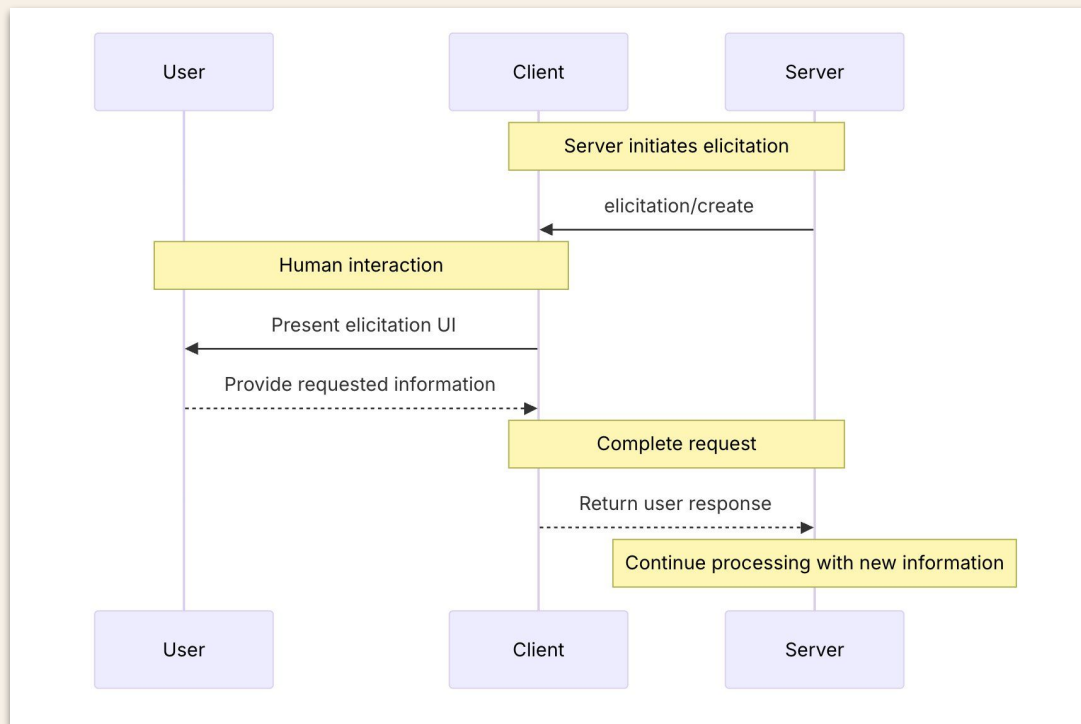
Therefore, Agents should be represented as **async workflows**.



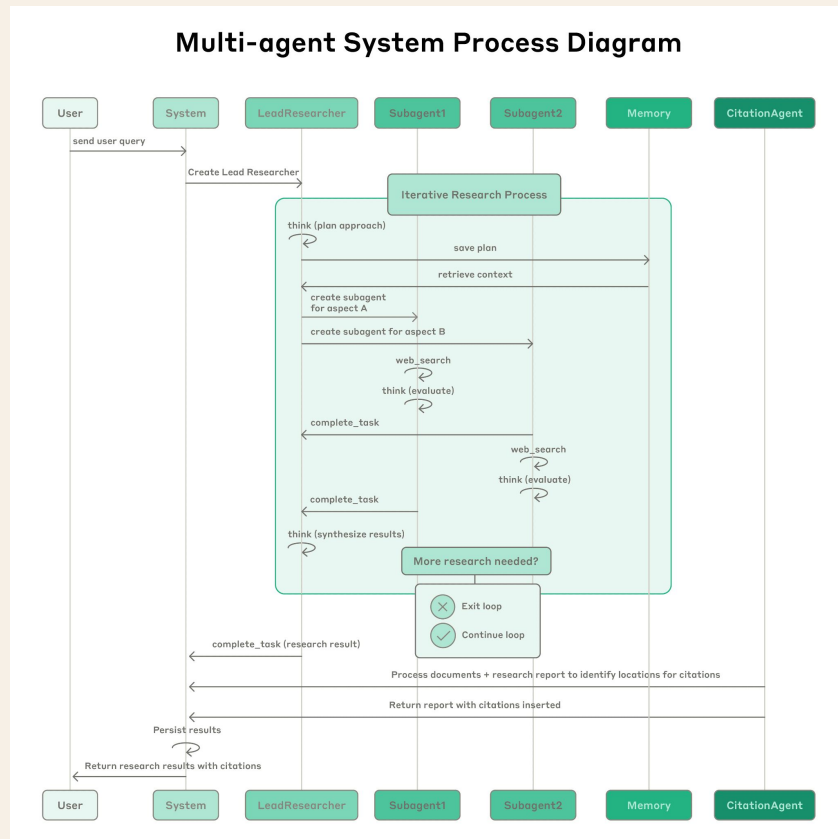mcp-agent 🐕 uses Temporal to execute agent tasks.

# Elicitation

**\*New\*:** MCP servers can request human feedback from clients (and therefore users). This is called "elicitation".

# Emerging agent architecture

1. **Orchestrator:** Planner/task manager agent (main point of contact with human).

2. **Multiple sub-agents** with specialized tasks

3. **MCP access:** Agents connect to tools & data via MCP servers

4. **Durable workflows**: Long-running async tasks (jobs).
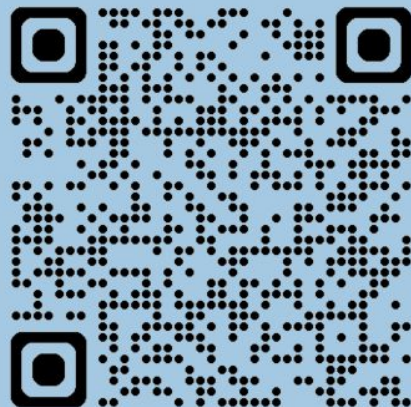
5. **Human-in-the-loop** (not full autonomy)



Multi-agent System Process Diagram

# Problems yet to be solved

1. Agents as MCP servers
   a. Long-running tools
   b. Agent authentication / authorization
2. Context management
3. Tool management
4. Human in the loop for long-running workflows

We are working on solving these!

# THANK YOU!



github.com/lastmile-ai/mcp-agent

@qadri_sarmad

sarmad@lastmileai.dev

lastmile AI