# Week 6 → MCP (Model Context Protocol)

→ MCP is a standardized framework for AI agents to interact with external systems and data, enabling them to operate with enhanced memory reasoning, and multi-task execution.

## What MCP is not →
① A framework for building agent
② A fundamental change to how agents works
③ A way to code agent

## What MCP is actually →
① A protocol - a standard
② A simple way to integrate tools, resources, prompts
⑤ "A USB-C port for AI application"

## Reason not be excited →
① It's just a standard, it not a tool themselves.
② Langchains already has big Tool ecosystem
③ You can already make any function into any tool

## Reason to be excited →
① Makes it easy to integrate tools.
② It's taking off! Exploding ecosystem
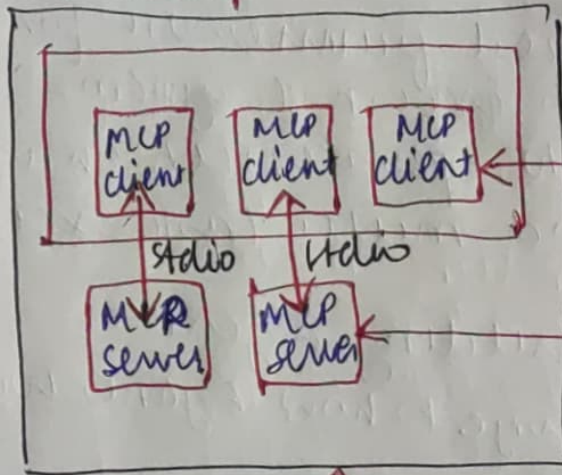③ HTML was a standard too.

## MCP core concept →
The three components →
① (Host) is an LLM app like claude or our Agent architecture.
② (MCP client) lives inside Host and connects 1:1 to mcp server
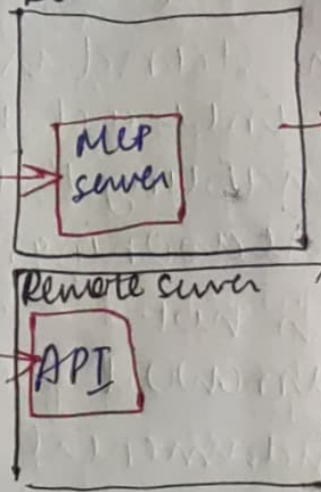③ (mcp server) provides tools, context and prompts

example →
fetch → MCP server that searches the web via a headless browser

**Your computer** — Remote server diagram

- Your computer: MCP client, MCP client, MCP client — SSE → MCP server → alternative architecture
- Stdio / Stdio connections to MCP server, MCP server
- Remote server: API

→ In most of example or cases we will be seeing that MCP server are running Outside the host and your connecting it , it still running on your machine

→ MCP servers can takes advantages of functionality over internet.

Two different Transport machenism→

① stdio → (standard Input output) → The AI runs and mcp server as a local program, communication via standard input/output.

② SSE → (Server-Sent-Events (SSE)) → The AI connects to a remote server via HTTP and exchanges message asynchronously.

**commands** →

① fetch - params = {"command": "(uvx)", 'args':
   { "mcp-server-fetch" }
   
   *This also uses playwright*
   
   This commands is going to create an MCP client running within OpenAI Agents SDK.

→ for running mcp server by using JS node npm
   playwright -param = {"command": "npx", 'args':
   { "@playwright mcp @' latest"}}

## Distributed agent Runtime?

→ A distributed agent runtime handles life cycle and communication lifecycle and communication accross process boundaries, consisting of:

(a) **host service** → A host service connected to worker runtimes, handling message delivery and sessions for direct message

(b) **worker runtime** → A. worker wruntime advertises agent to the host service and handles execution their code.

→ Grpc is an host to run and send message on local host according to the port no. specified.

## Which framework to select?

→ It's not the most important question and it doesn't really matter!

→ Pick the framework that suits your styles and the skill of your team.

→ ~~my go to OpenAI~~

→ Other frameworks → Google ADK, Huggingface smeragent and Pydantic AI - but they will seem familiar.

## What matters →

① Start with the problem, not the solution.

② Have a metric to evaluate success

③ favour workflow over autonomy finally.
(ao try to do the code by hardwriting)

④ work bottom up, not top down
(do not start with a big problem, instead start with a small thing).

⑤ Start simple, then add [small piece → big problem]

⑥ Start with large frontier model, then reduce

⑦ Think context rather than memory.

⑦ [Try giving LLM the right context to be able to answer its question].

⑧ most problems are solved with better prompts.

⑨ look at the traces

⑩ Be a scientist, no shortcut to R&D.