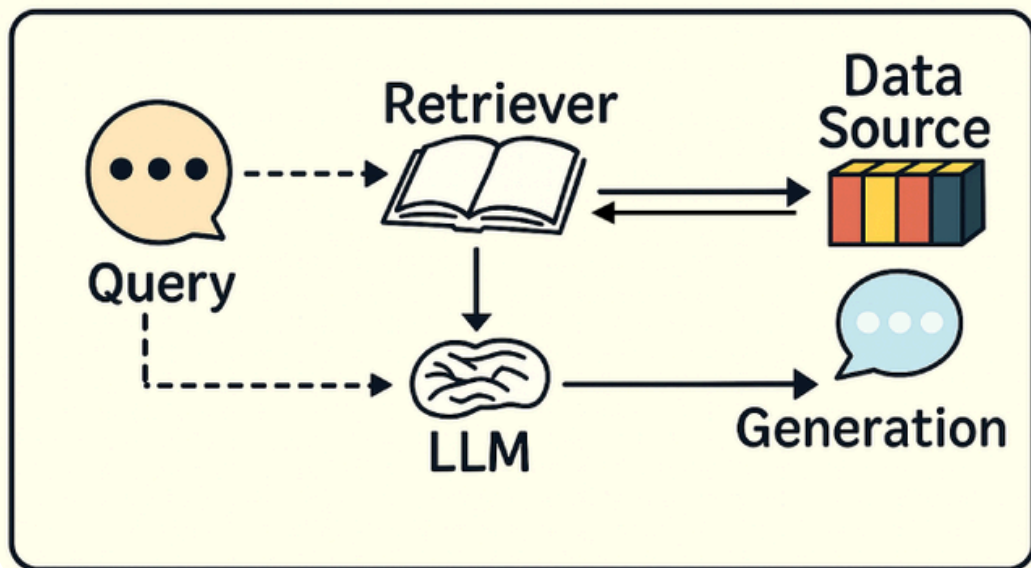
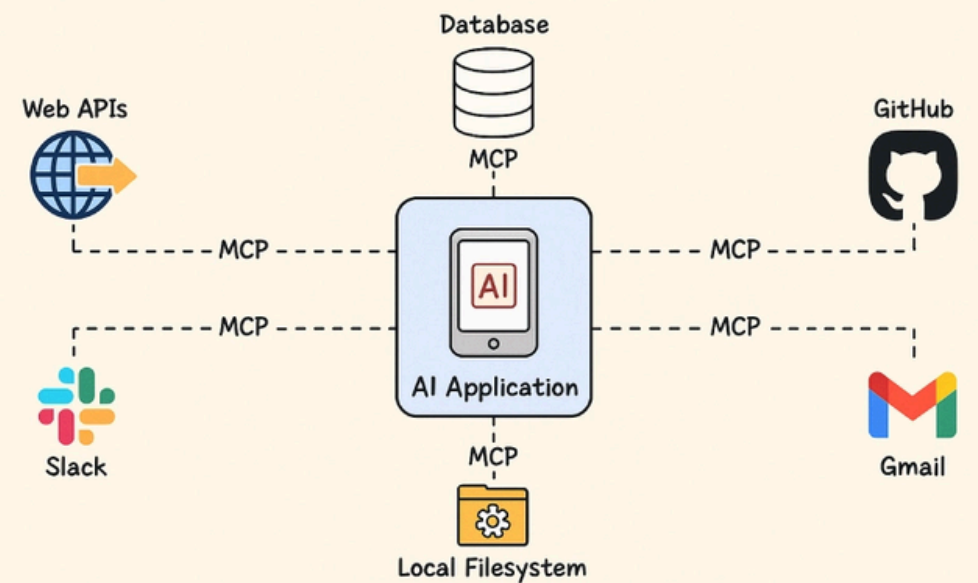


# MCP vs RAG

## What is RAG?



## What is MCP?



### User Query

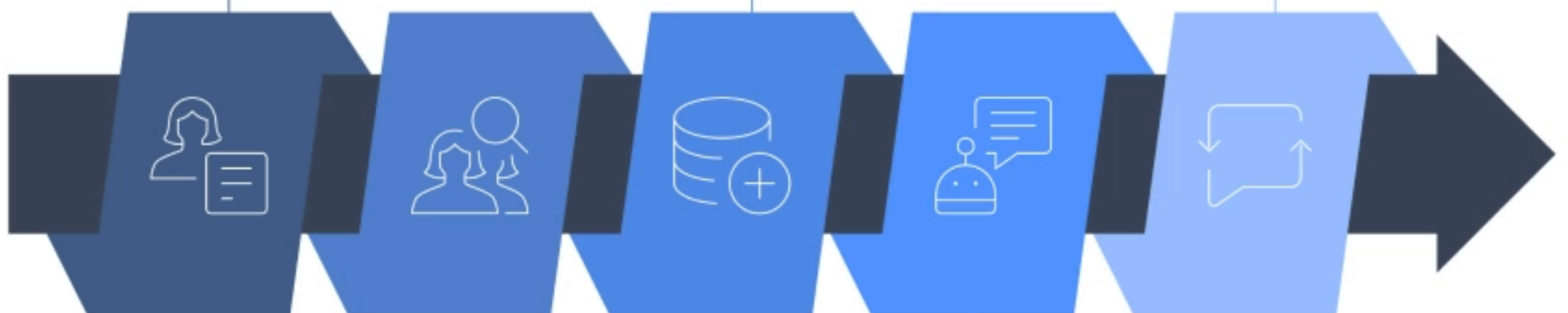
The user inputs a query to the system.

### Augmentation

The retrieved information is added to the original query.

### Response to User

The generated response is shared with the user.



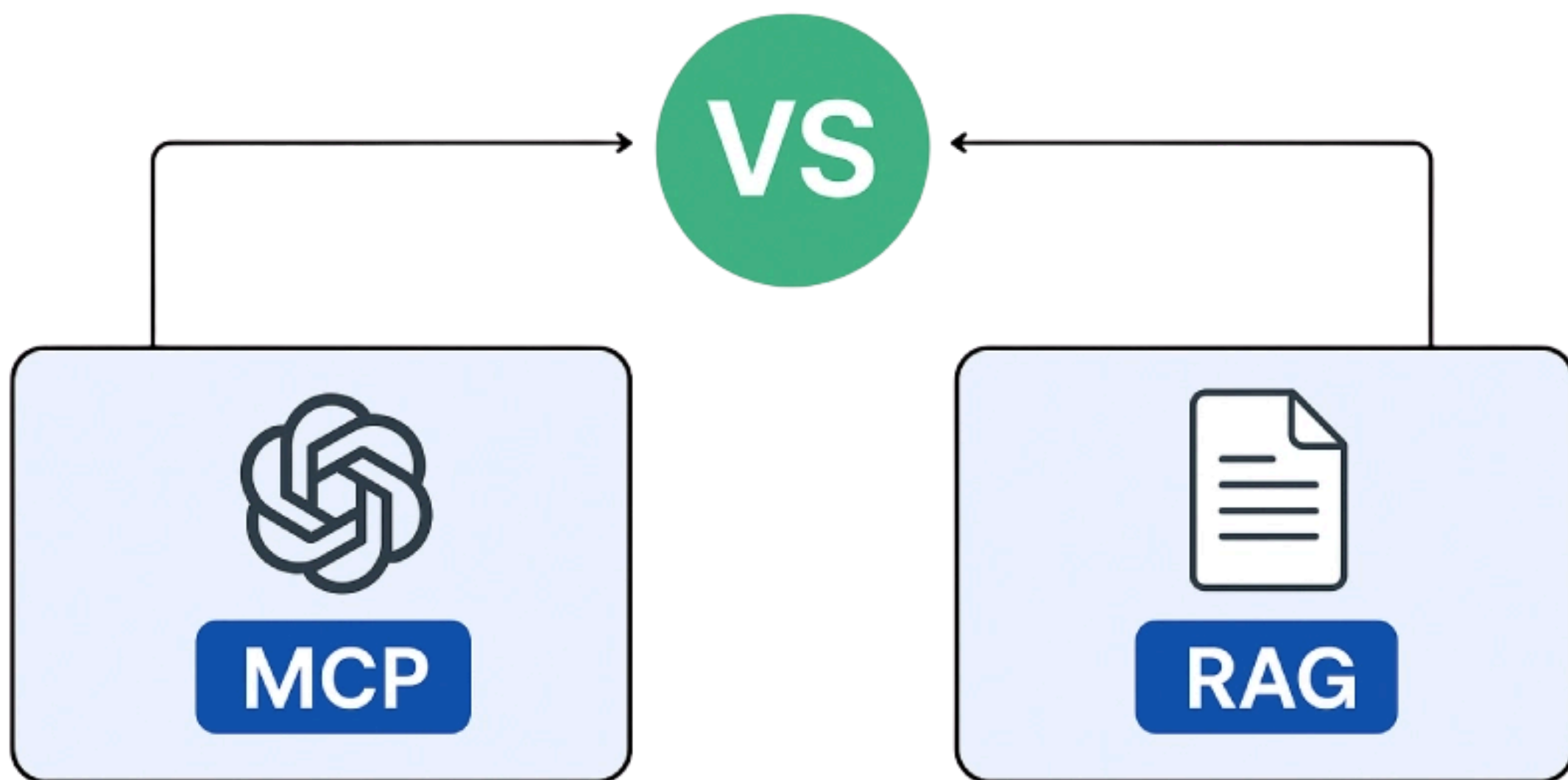
### Retrieval Process

The system searches for relevant information in the knowledge base.

### LLM Generation

The LLM generates a response using the combined input.

# MCP vs RAG: Competitors?



No, MCP and RAG are not competitors in the way they work or the tasks they perform.

As we have discussed in the previous sections, MCP and RAG perform different tasks and empower LLMs in different ways.

RAG powers LLMs with additional data while MCP grants LLMs the ability to act.

The key differences between MCP and RAG are summarised in the table below:

Feature	RAG (Retrieval-Augmented Generation)	MCP (Model Context Protocol)
<b>Purpose</b>	Enhances knowledge of LLMs by retrieving relevant external data	Extends the capabilities of LLMs to use tools and perform actions
<b>Function</b>	Pulls info from documents, databases, or search APIs	Connects to tools, APIs, software, and real-time systems
<b>Use Case Type</b>	Improves response accuracy and context relevance	Enables real-world actions, tool use, and automation
<b>How It Works</b>	Retrieves relevant documents → augments the prompt → generates output	Uses structured tool schemas → selects tool → executes action
<b>Data Access</b>	Typically works with textual or vector data	Works with functional endpoints (e.g., APIs, plugins, webhooks)
<b>Execution</b>	Passive: Only retrieves and informs	Active: Can take actions like submitting forms or updating systems
<b>Example Task</b>	“What is our refund policy?” → fetches from policy doc	“Cancel my subscription” → triggers refund API
<b>Model Input Impact</b>	Expands the prompt with more content for better grounding	Doesn’t always expand the prompt, focuses on decision and execution
<b>Complexity</b>	Requires vector DB, chunking, and embedding logic	Requires tool definitions, security layers, and execution control
<b>Best Used For</b>	Knowledge-based Q&A, grounding, and content generation	Workflow orchestration, automation, and tool-augmented agents

# Can MCP and RAG work together?

Yes, MCP and RAG can work together to help us design highly sophisticated AI workflows. RAG allows LLMs to pull relevant information while MCP executes tasks based on retrieved knowledge. Using these two together, we can create the following workflows:

## 1. RAG as a tool within the MCP framework

In this case, an LLM operating with MCP can have RAG as one of its tools, which it can use to fetch the required information.

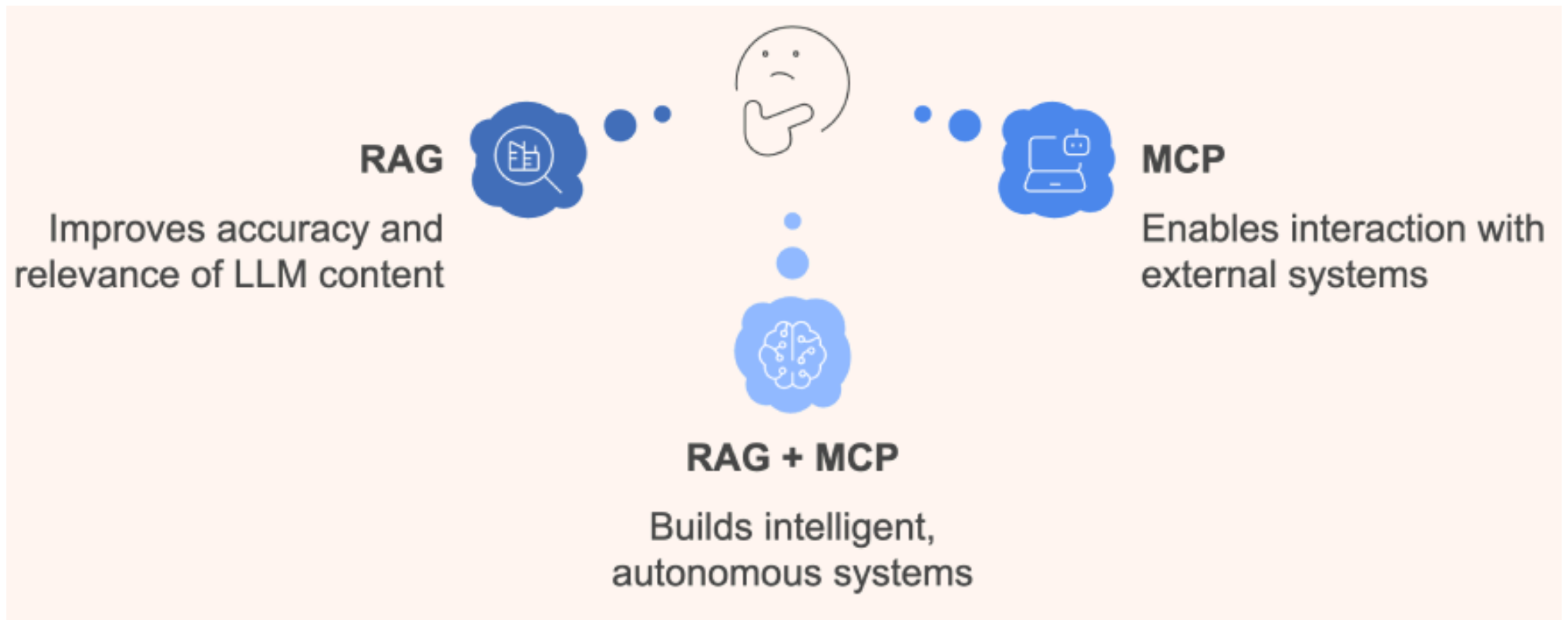
Example: An MCP-powered AI system for a Marketing Campaign. It uses RAG to retrieve information regarding previous campaigns and competitor information. Then, using MCP-powered tools, it creates social media posts and schedules them across different platforms.

## 2. MCP for guiding RAG-Powered Agents

In systems involving multi-agents, each agent can have its own RAG pipeline and MCP can act as a coordinator for the system.



# Which one should you pick?



The choice between RAG, MCP, or RAG + MCP depends on the task. Each of the frameworks has its unique strengths. Here is how you can decide which approach to take:

- **RAG:** If your main goal is to improve the accuracy, relevance, and factual grounding of LLM-generated content, then “RAG” should be your choice.
- **MCP:** If your main goal is to allow your LLM to interact with external systems, perform actions, or leverage tools to complete its tasks, then “MCP” is your go-to path.
- **RAG + MCP:** If your goal is to build an intelligent, autonomous system that can better understand and act decisively, then the combination of RAG and MCP is your go-to option.