

# From RAG to Agentic system

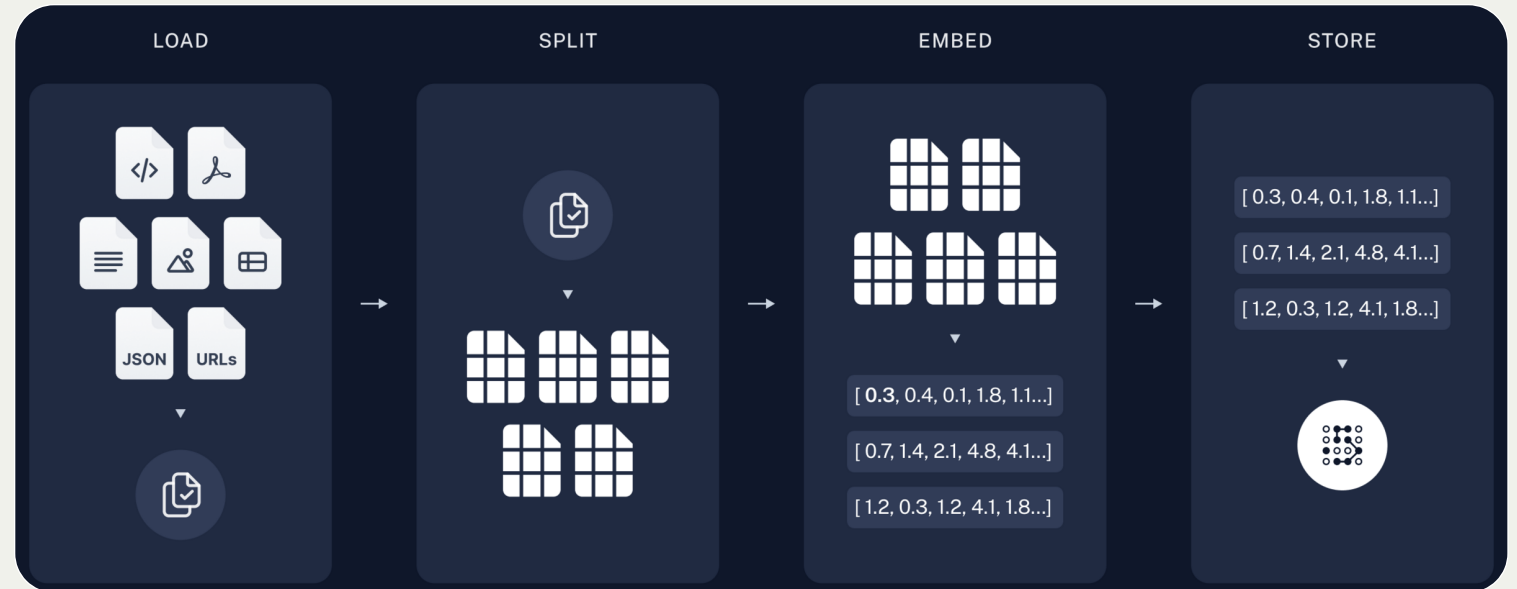


# Intro

# What is a RAG?

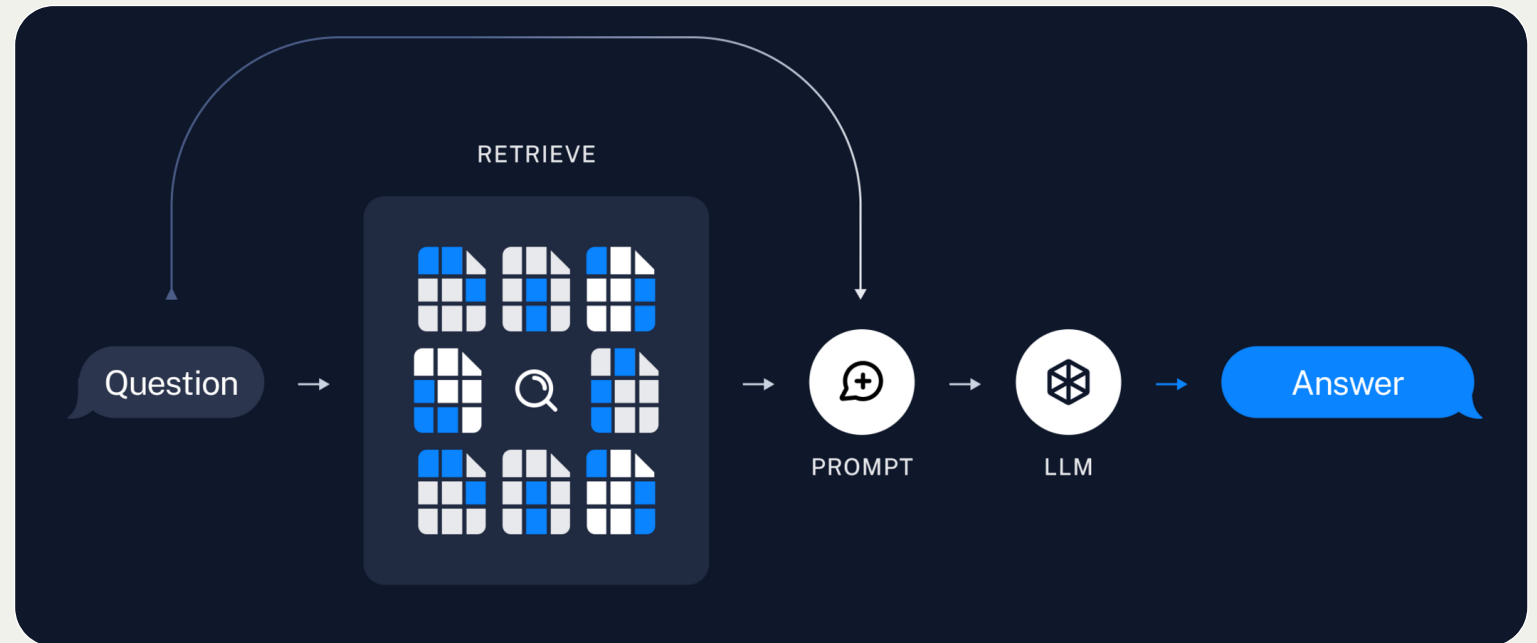
## 1. Indexing

parsing documents  
chunking  
encoding  
indexing in db



# What is a RAG?

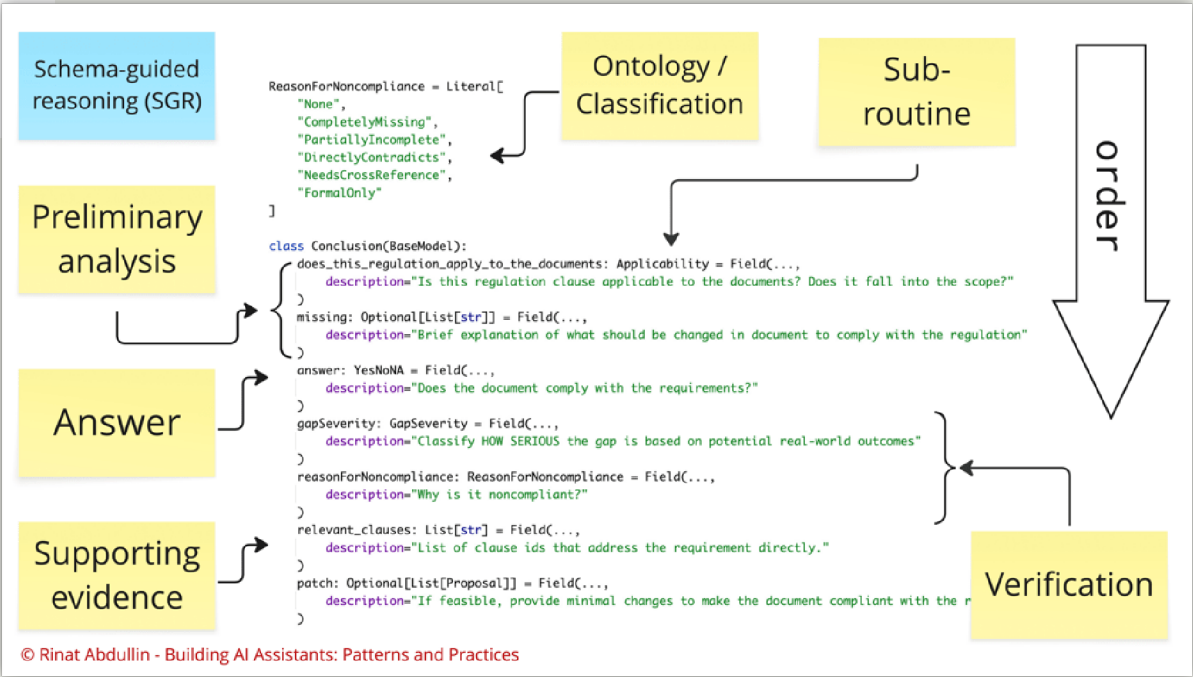
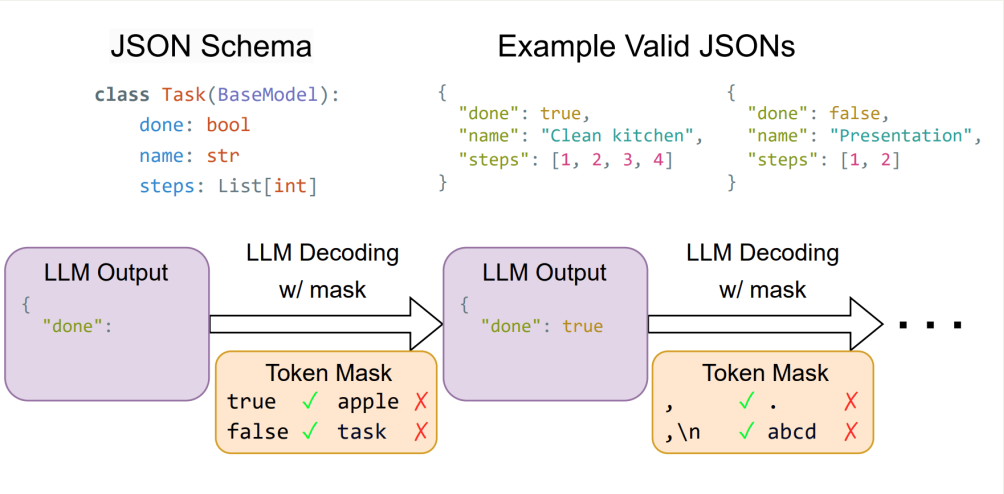
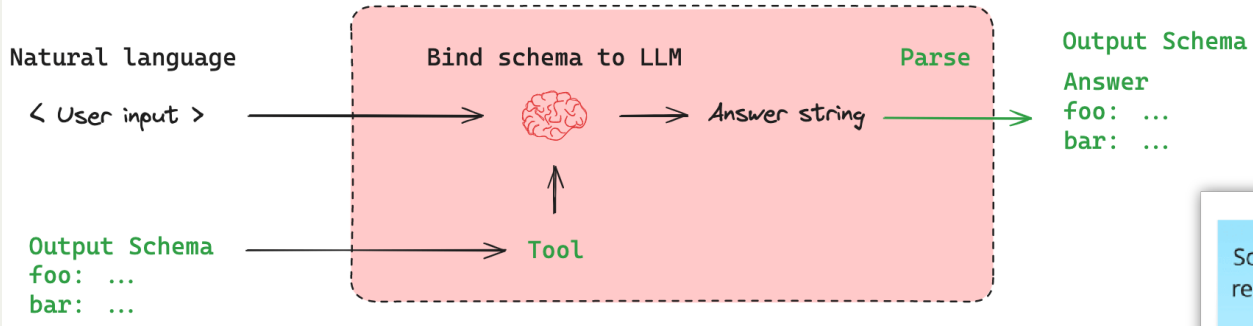
1. Indexing stage
2. **Retrieval & Generation**  
LLM answers question



# What is missing

- **Up-to-date information beyond context:**
  - attach new sources, browse web
- **Memory and state:**
  - continue conversation or update internal parameters
- **Multi-turn:**
  - run retrieval multiple times with different queries on the same topic

# Structured output & Schema Guided Reasoning

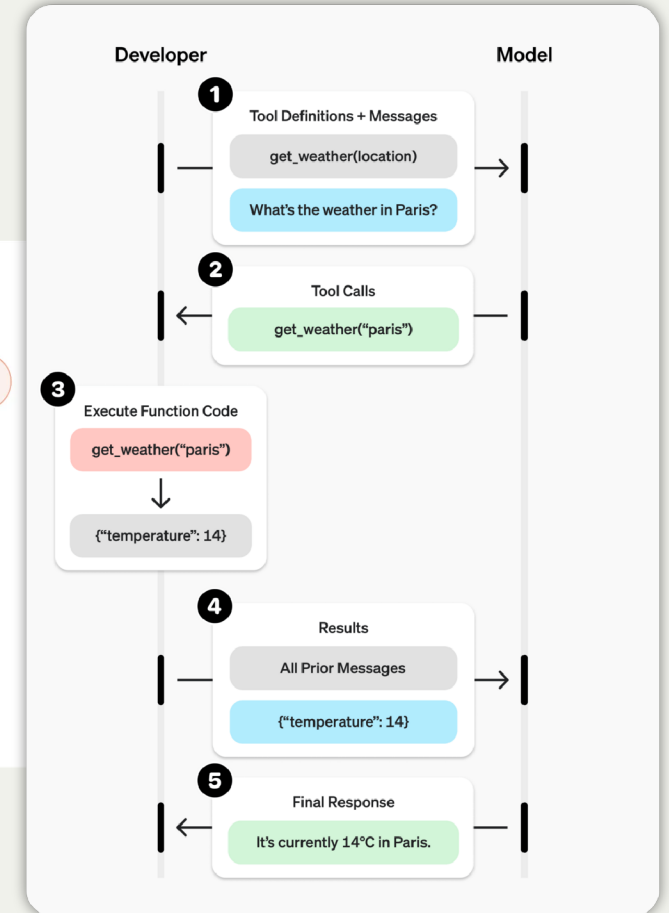
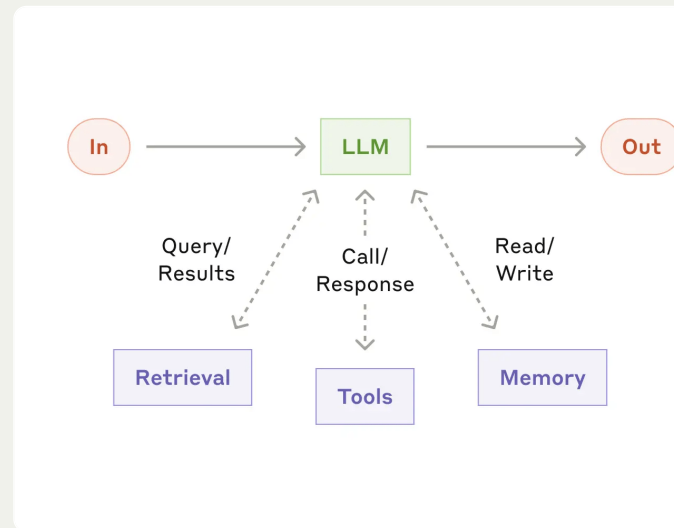


# Tools & Function calling

Structured output enables  
more stable function calling

**Input:** tool definitions + schemas

**Output:** selected tool + arguments



# Tools: Model Context Protocol

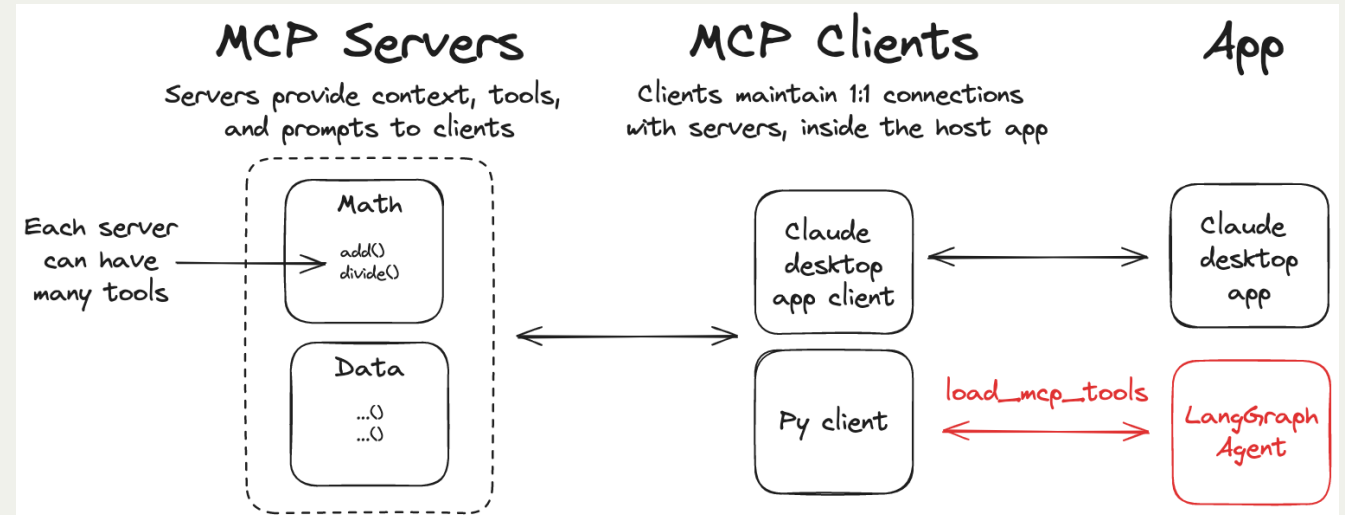
Proposed way for services to expose **tools** to LLMs

```
@mcp.tool()
async def get_forecast(latitude: float, longitude: float) -> str:
    """Get weather forecast for a location.

    Args:
        latitude: Latitude of the location
        longitude: Longitude of the location
    """
    # First get the forecast grid endpoint
    points_url = f"{NWS_API_BASE}/points/{latitude},{longitude}"
    points_data = await make_nws_request(points_url)

    if not points_data:
        return "Unable to fetch forecast data for this location."

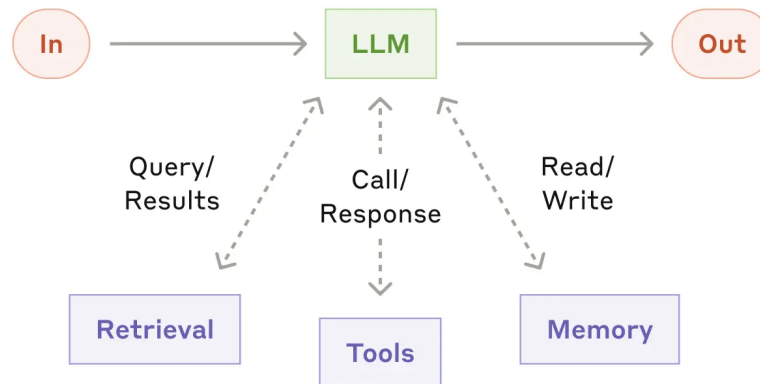
    # Get the forecast URL from the points response
    forecast_url = points_data["properties"]["forecast"]
    forecast_data = await make_nws_request(forecast_url)
```





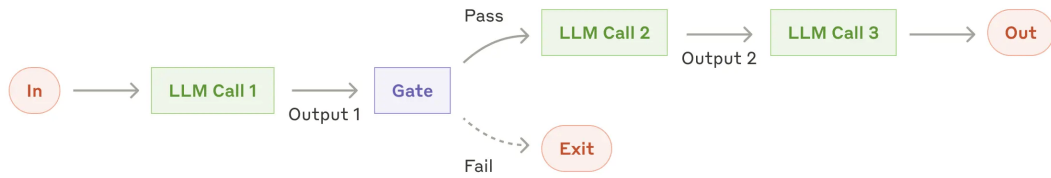
# What is an Agent?

## Building block with state and tool usage

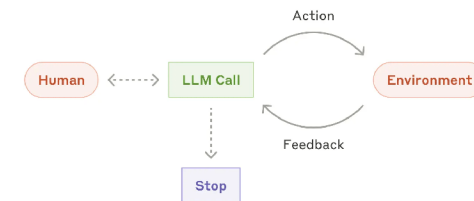
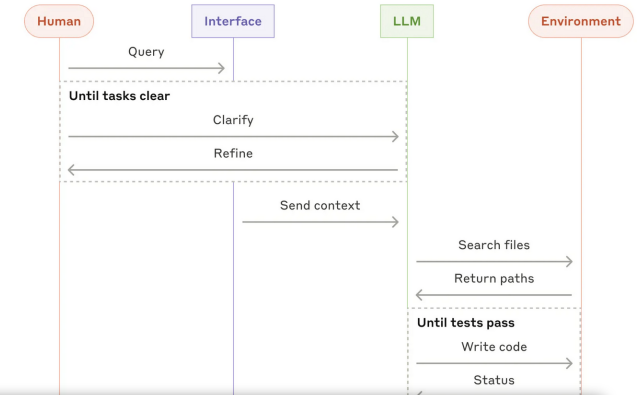


# What is an Agent?

## "Workflows"



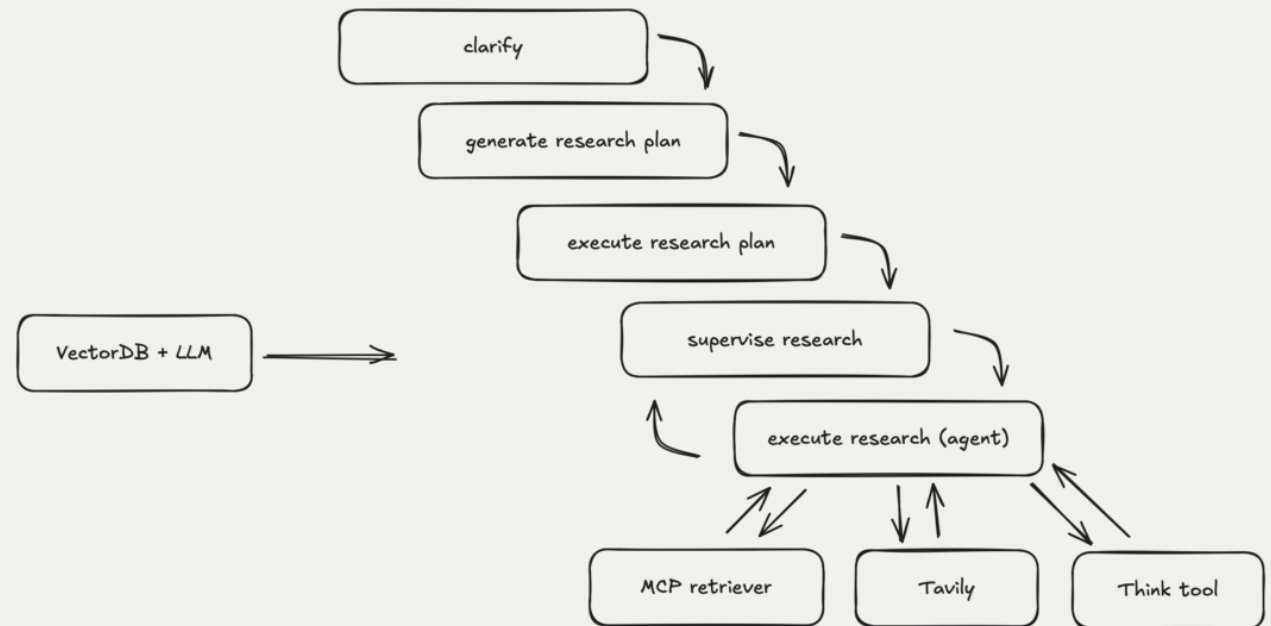
## "Agents"



# We are gonna build

We will walk through:

- structured outputs
- search agent (LLM + Tools)
- retrieval as MCP server
- states and communication between agents



# Thank you!

I am

Vladimir Ageev, DS @ EPAM Systems

