

# Episodic Memory

Previous interactions: [

Human: ...

Assistant: ...

Private Knowledge Base



...

Documentation



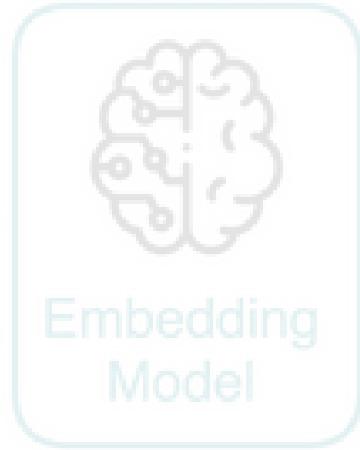
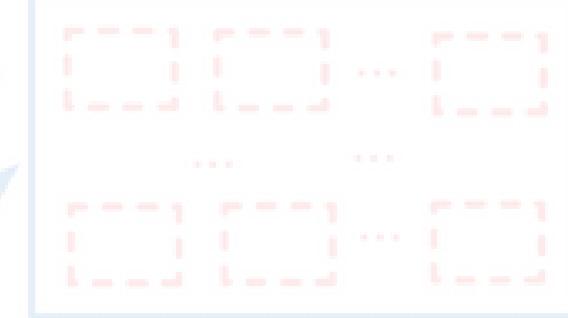
Grounding Context

# Semantic Memory

Approximate  
Nearest  
Neighbour  
search  
(ANN)

Vector  
Database

- Embedding(Latent) Space



Indexing

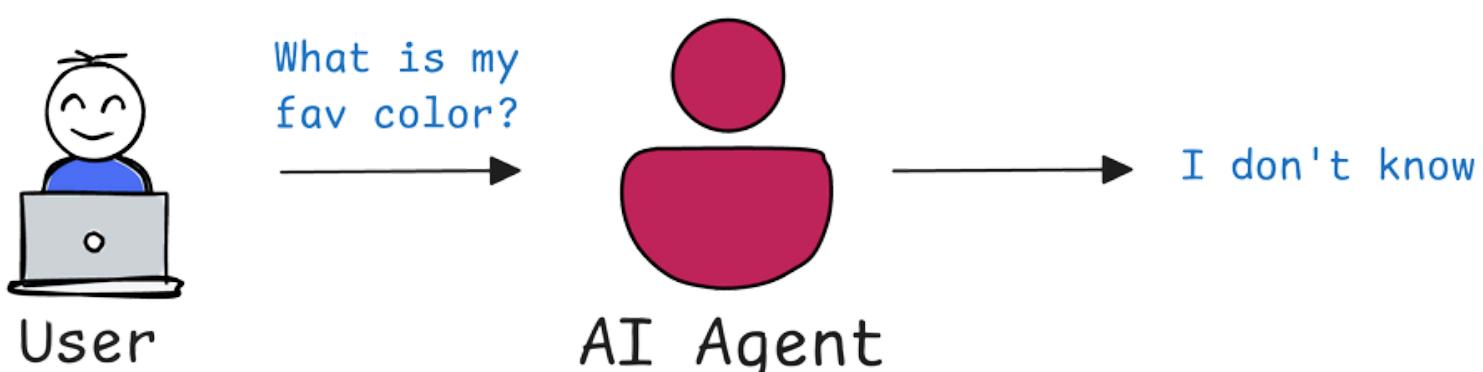
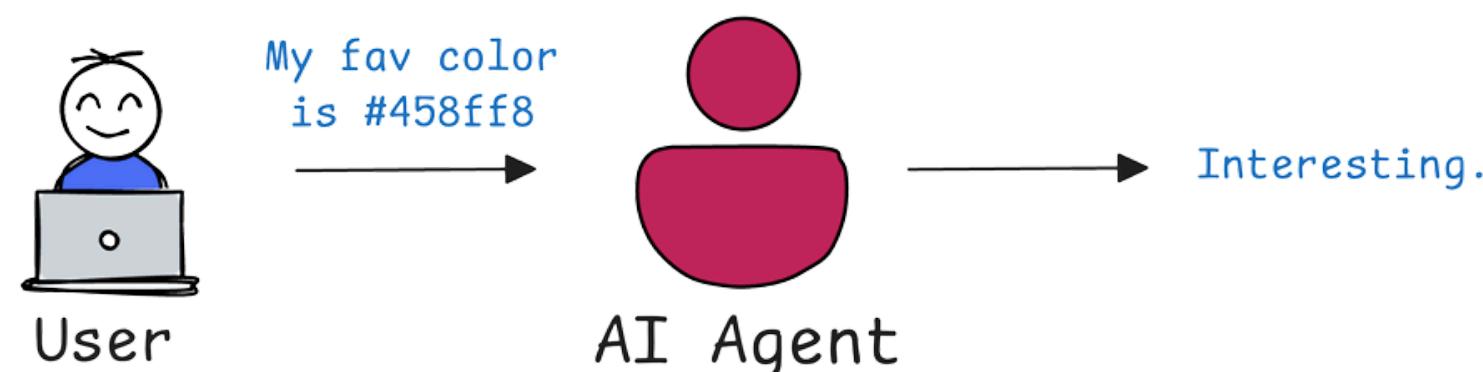
Vector Index

# WHY MEMORY MATTERS?

Agentic AI isn't just generating responses - it's solving tasks, over time.

To do that, it needs memory to:

## Interaction without memory



Source

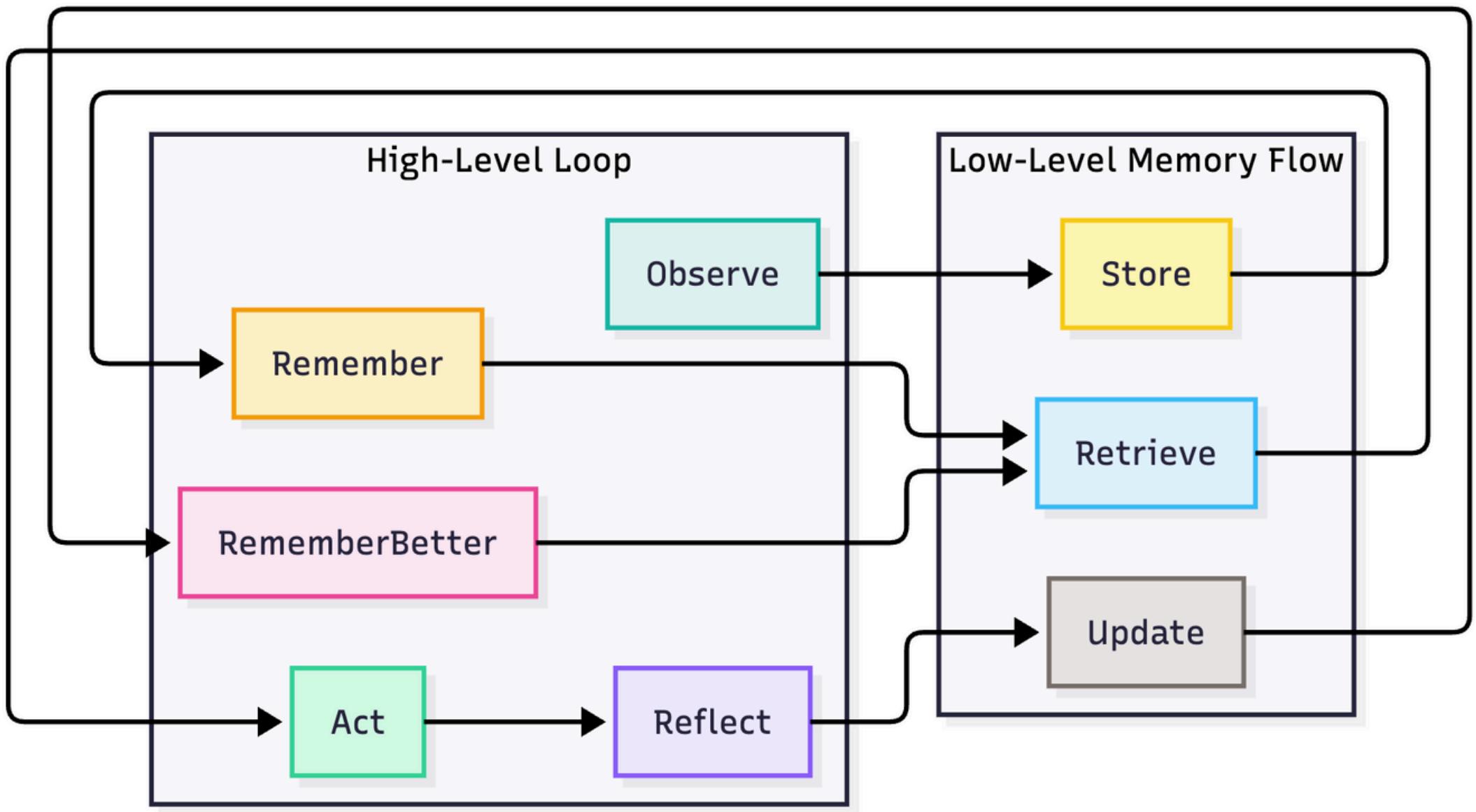
- Track what's already been done.
- Stay aligned with user goals.
- Avoid repeating mistakes.
- Remember relevant knowledge.
- Reflect and improve on the fly

Without memory, the agent restarts every time like a GPS that forgets the destination after every turn.

# HOW DOES MEMORY WORK?

Memory in agents works like a dynamic knowledge system.

Here's what happens under the hood:



- **Store**: The agent logs relevant data during task execution (like goals, decisions, errors).
- **Retrieve**: When needed, it pulls past info to guide current actions.
- **Update**: After each task or feedback cycle, it revises memory just like learning.

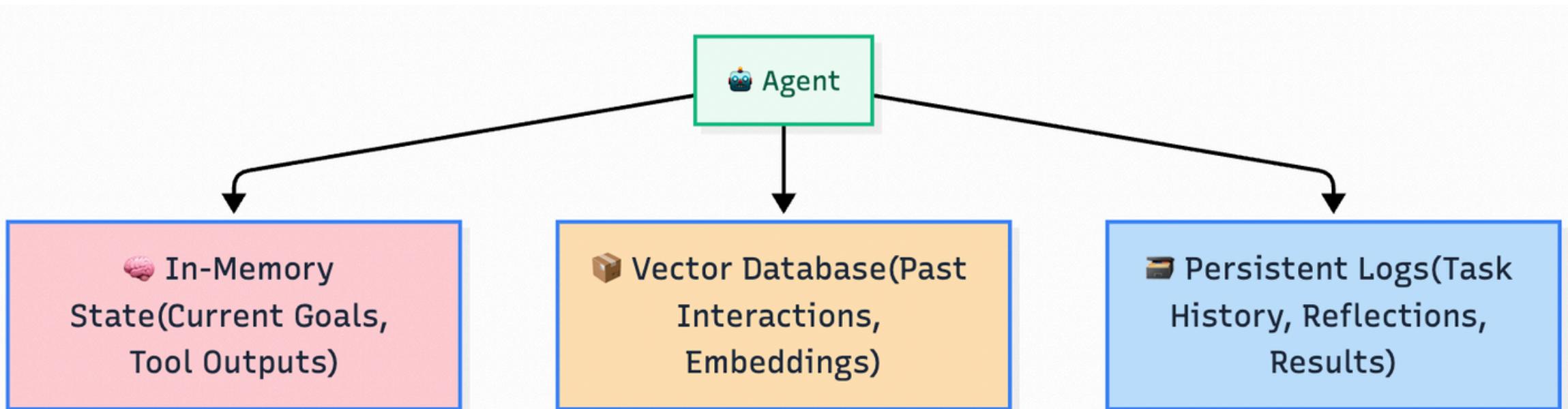
Think of it as a loop:

- Observe → Remember → Act → Reflect → Remember Better

# WHERE IS MEMORY STORED?

In agentic systems, memory isn't magic it's data, carefully organized and stored.

Common storage places include:



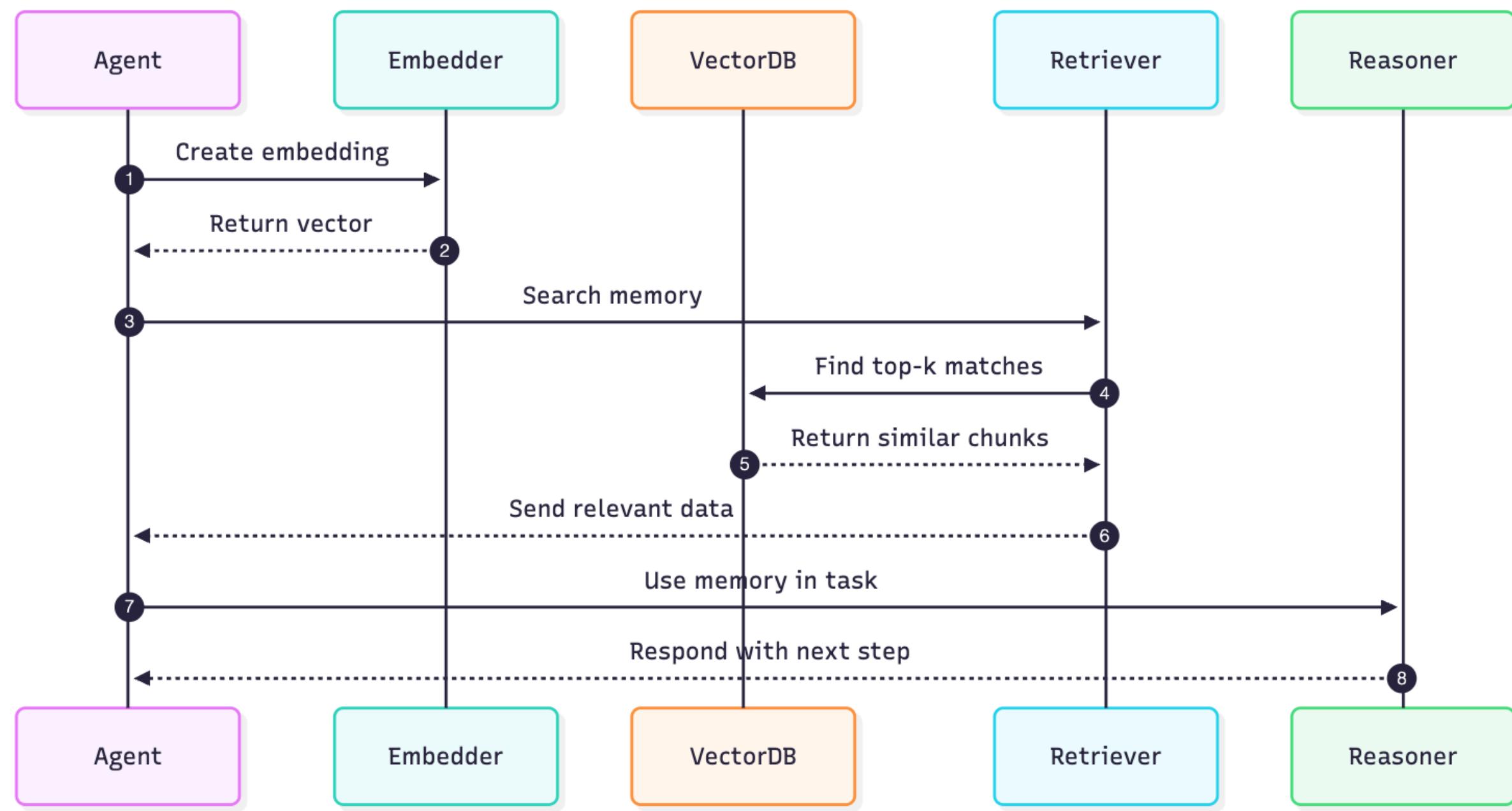
- **In-memory state:** Temporary info held during a task (like current goals or tool outputs)
- **Vector databases:** Store past interactions as embeddings, so agents can retrieve similar experiences later
- **Persistent logs or files:** Used to save long-term events, reflections, and task outcomes across sessions

Different memory types live in different layers but together, they give agents the power to remember, reason, and improve.

# HOW RETRIEVAL WORKS?

When faced with a new situation, agents don't scroll through all memory. Instead, they search semantically.

Here's how it works:



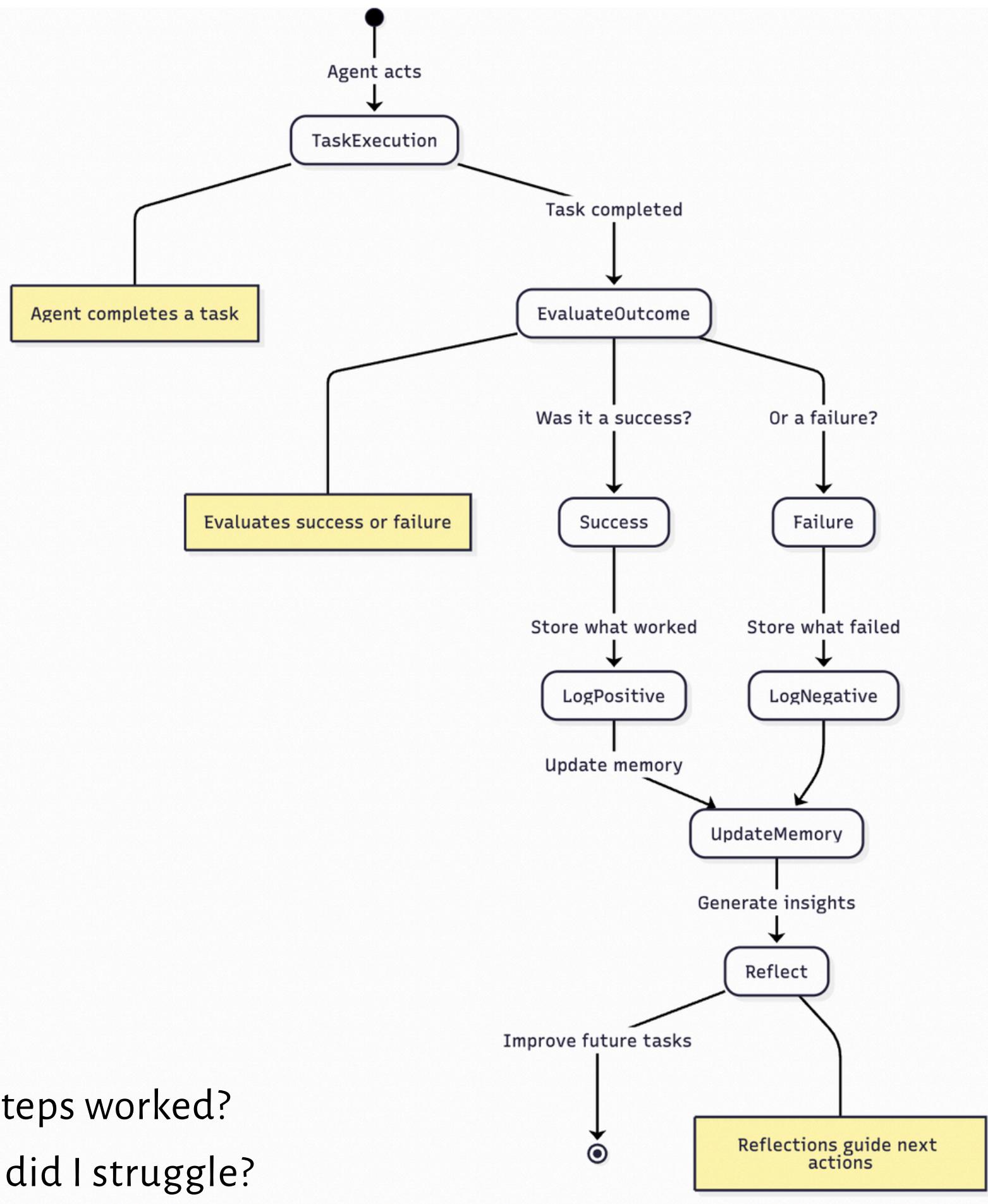
- The agent turns the current task or question into an embedding
- It compares this to past memory chunks stored in a vector database
- Then it pulls back the most relevant ones like similar errors, prior results, or past user goals

It's fast, fuzzy, and context-aware more “find what's useful” than “remember everything.”

# REFLECTION & LEARNING

Once an agent finishes a task, it doesn't just move on. It reflects. This self-assessment is stored back into memory not just as data, but as lessons.

It asks:

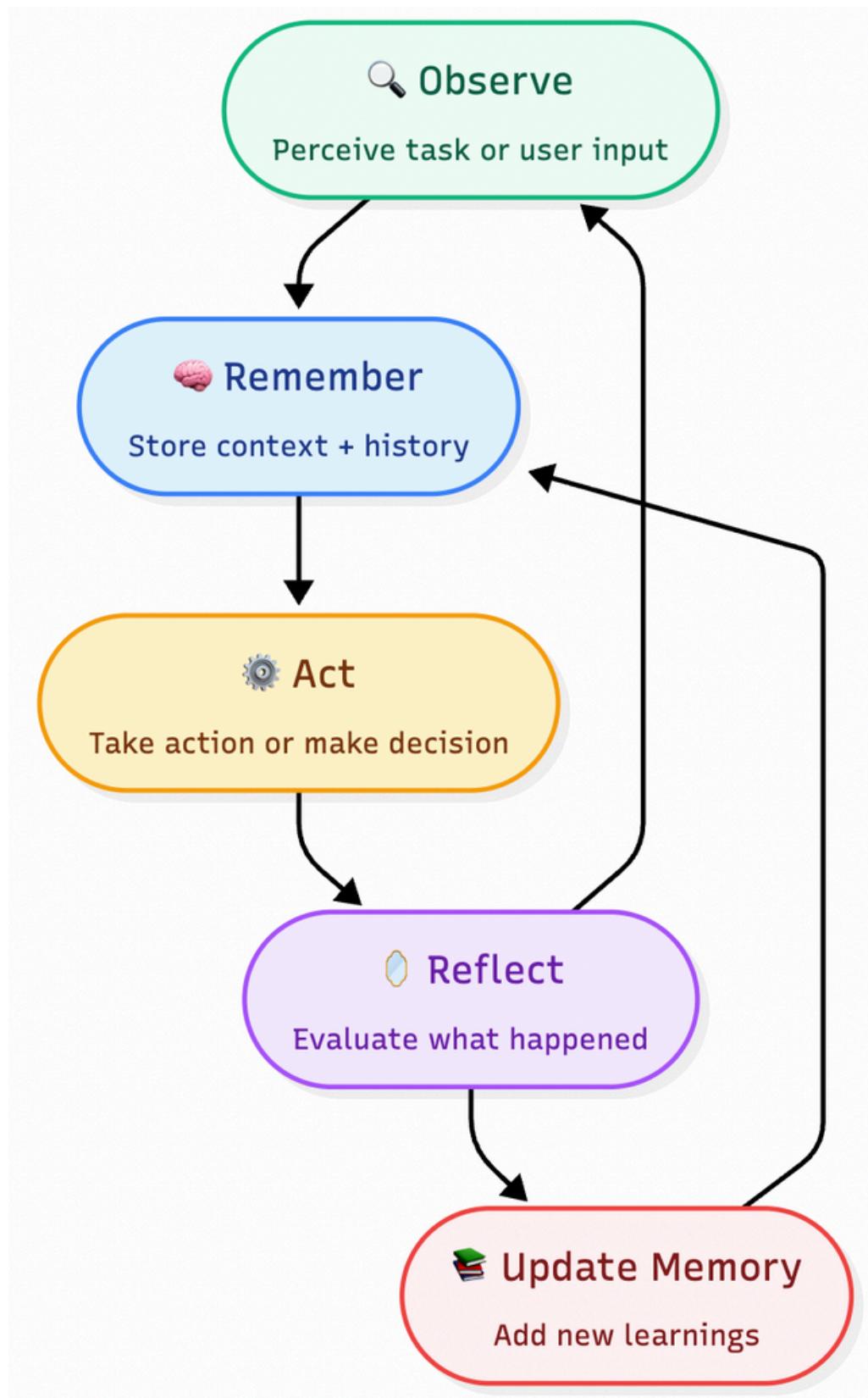


- What steps worked?
- Where did I struggle?
- Should I try a different approach next time?

This reflection loop is what separates a one-shot LLM from a truly agentic system.

# THE MEMORY LOOP

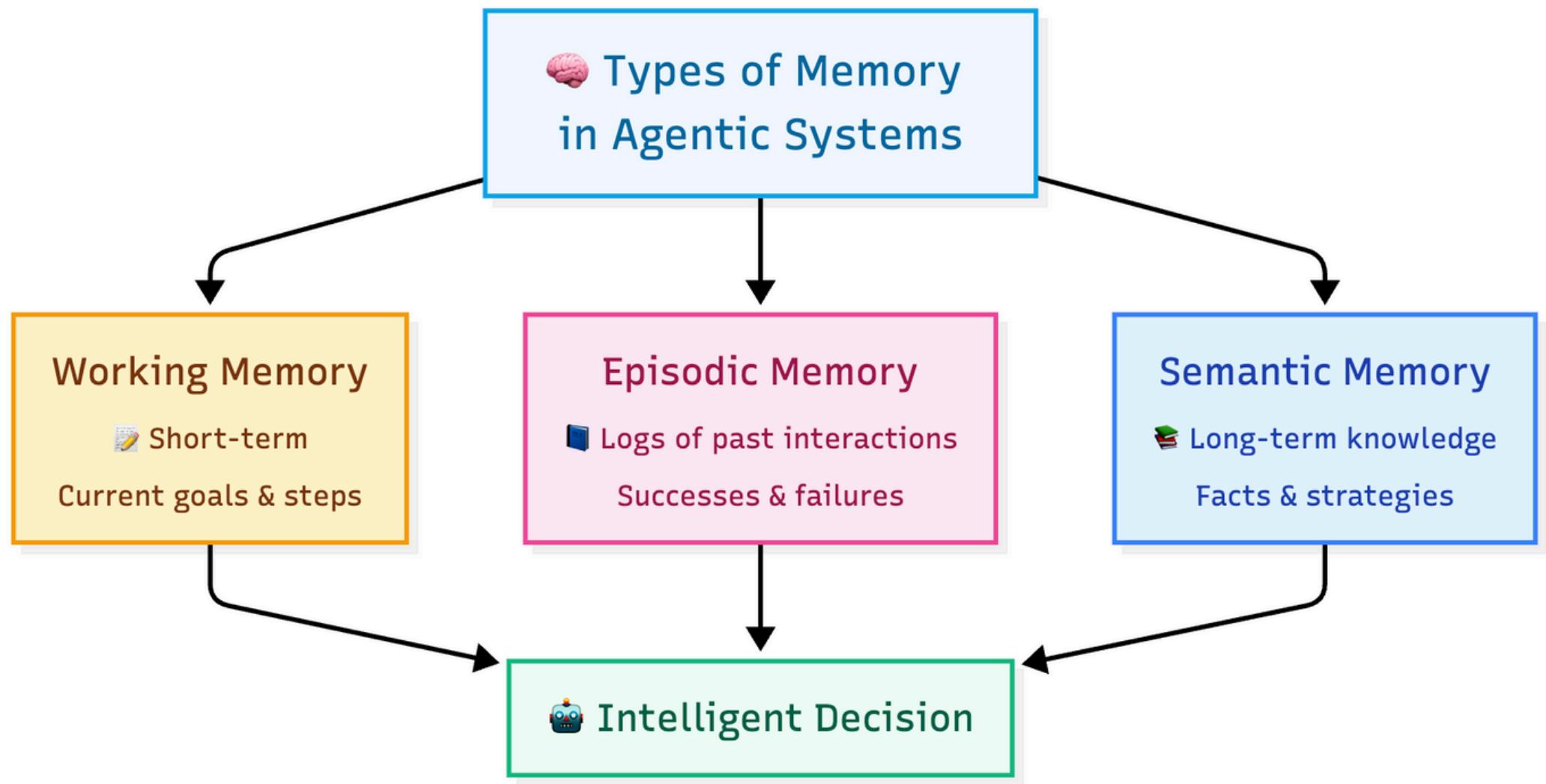
Memory isn't a one-time thing it's a continuous cycle. Here's how memory actually works inside an agent:



- **Store**: Capture goals, inputs, and outcomes
- **Retrieve**: Pull relevant info when needed
- **Use**: Apply that info to act or decide
- **Reflect**: Analyze what happened
- **Update**: Feed new insights back into memory

# TYPES OF MEMORY IN AGENTS

To act intelligently, agents use different types of memory, each serving a unique purpose:



- **Working memory**: Short-term. Stores current instructions, goals, or steps in progress. Think of it like a mental sticky notes.
- **Episodic memory**: Logs past experiences successes, failures, user interactions. Think of it like a task diary.
- **Semantic memory**: Stores long-term knowledge: facts, patterns, strategies, like the agent's internal knowledge base.

These memory types work together just like our own brain does to make decisions smarter and more contextual.

# TOOLS THAT POWER MEMORY IN AGENTS

Modern agent frameworks come with built-in memory support from short-term buffers to long-term retrievers.

Here are some of the most popular:



- **LangChain**: Offers working memory, retrievers, and vector-based storage
- **LangGraph**: Adds stateful memory with dynamic control flows
- **LlamaIndex**: Ideal for complex retrieval, summarization, and long-term memory chunks
- **CrewAI**: Enables shared memory across multiple agents

Each one helps agents store, retrieve, and reflect the core loop of intelligent behavior.

# WHAT IS MEMO?

Memo is a scalable, intelligent memory architecture built for agents not just stacked onto them. Here's what makes it special:

Feature	Traditional LLM Agents	With Memo
Memory Management	Manual or static	Automatic, relevance-informed
Cost Efficiency	High token usage	Saves ~90% tokens & latency
Remembering Smart	Echo prior text	Consolidates, filters, corrects itself
Long-Term Coherence	Session-limited	Persistent, evolving knowledge over time

- **Smart extraction + Update pipeline:** LLMs automatically extract meaningful facts, detect contradictions, and decide to add, update, delete, or skip keeping memory coherent and cost-efficient
- **Fast & Efficient retrieval:** Operating with 91% lower latency and 90% fewer tokens than full-context systems meaning smarter responses, faster.
- **Dynamic memory management:** It filters what to remember, when to forget, and how to consolidate much like human memory.
- **Optional graph memory (Memog):** Supports complex, relational reasoning via knowledge graphs ideal for deeper, multi-hop understanding

# Stay Ahead with Our Tech Newsletter! 🚀

👉 Join 1k+ leaders and professionals to stay ahead in GenAI!

🔗 <https://bhavishyapandit9.substack.com/>

## Join our newsletter for:

- Step-by-step guides to mastering complex topics
- Industry trends & innovations delivered straight to your inbox
- Actionable tips to enhance your skills and stay competitive
- Insights on cutting-edge AI & software development

## WTF In Tech

[Home](#)   [Notes](#)   [Archive](#)   [About](#)

People with no idea about AI saying it will take over the world:



My Neural Network:

## Object Detection with Large Vision Language Models (LVLMs)

Object detection, now smarter with LVLMs

MAR 27 · BHAVISHYA PANDIT

## AI Interview Playbook : Comprehensive guide to land an AI job in 2025

Brownie point: It includes 10 Key AI Interview Questions (With Answers).

MAR 22 · BHAVISHYA PANDIT



WTF In Tech

My personal Substack

💡 Whether you're a developer, researcher, or tech enthusiast, this newsletter is your shortcut to staying informed and ahead of the curve.

Bhavishya Pandit



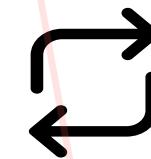
**Follow to stay updated on  
Generative AI**



**LIKE**



**COMMENT**



**REPOST**