# Agent Memory
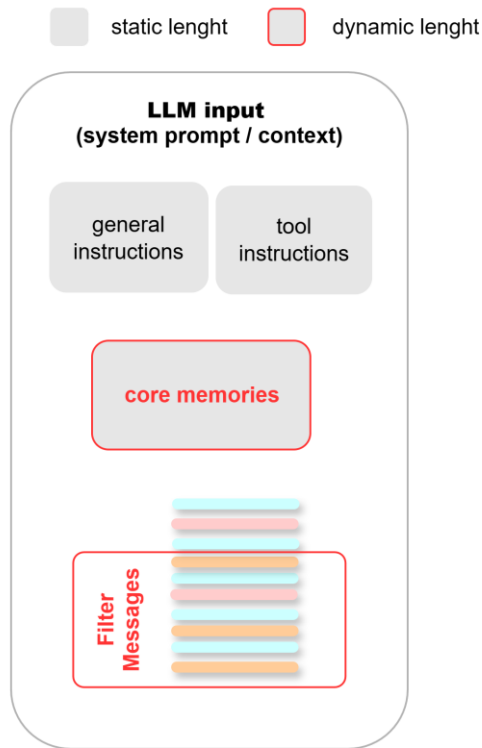
Review

https://app.diagrams.net/#G1Th2fe_mDnsKn49gxp0HKeYTBQrXzjUQT#%7B%22pageId%22%3A%2250LRqIvkJEDwD7JDsqcG%22%7D

**LLM input**
**(system prompt / context)**

static lenght — dynamic lenght

general instructions

tool instructions

core memories

Filter Messages

**Core Memories Section:**
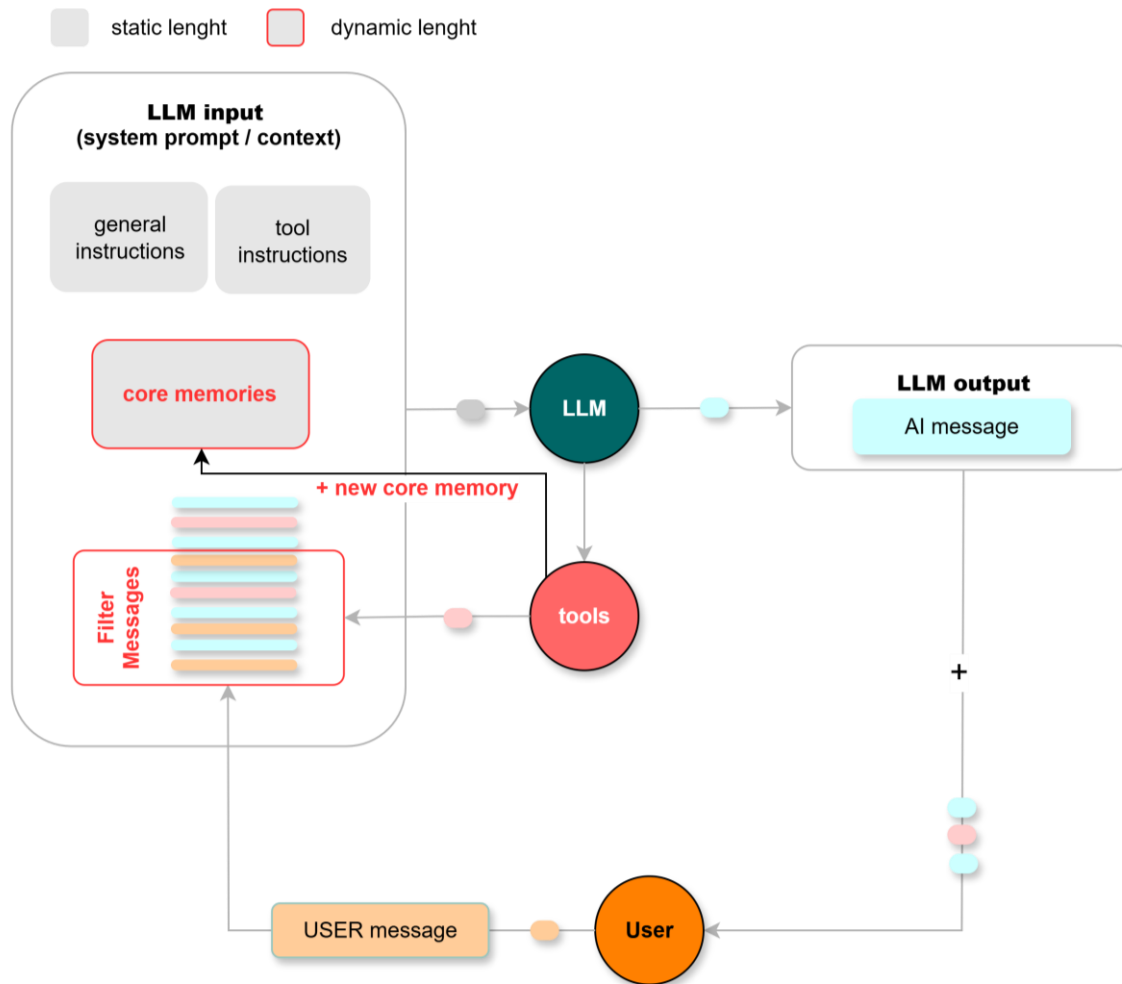Insights from the conversation, specially about the user

**Filtered Message List:**
Sliding Window or FIFO

# Working / Short Term / In-Context Memory

## Dynamic Context Window

Dynamic working memory allows to handle **context limits**.

**Messages** and **Core Memories** are stored in the graph **State** and inserted in the System Prompt, Tools, etc. whenever it is needed.
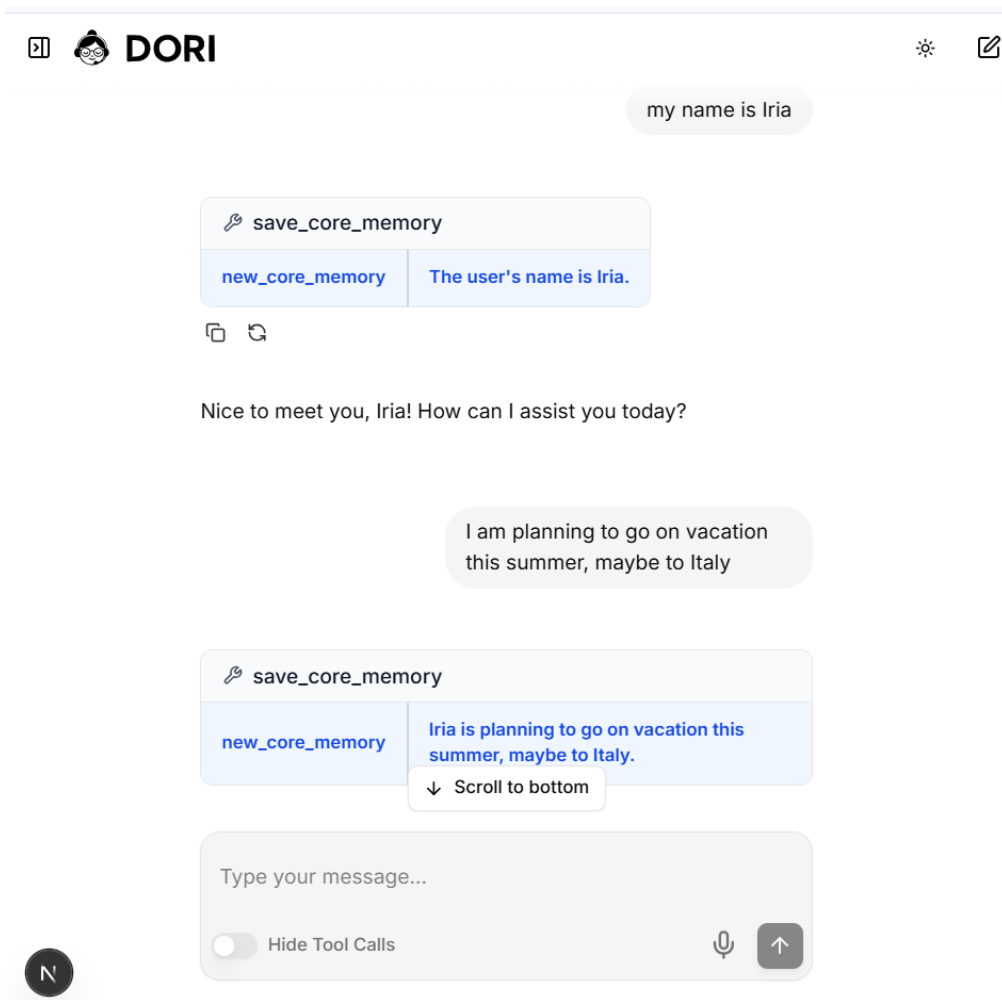
# Working / Short Term / In-Context Memory

## How to implement it?

### Option 1

- Main LLM **manages its own context memory**

- Memory tools **integrated** with normal tools.

## LLM Input

System - You name is DORI, an AI assistant. Your main task is to get to know the users, update your knowledge about t

You will be presented with 3 sections: ## Instructions (static), ## Core Memories (dynamic) and a ## Conversation His
The Core Memories are important pieces of information about the user that help you provide a better experience and pe
Do not rely only on the Conversation History, you must also use the Core Memories, as the Conversation History is lim
This should be updated after every interaction that has information.

## Instructions (static)

You must follow all rules exactly and never assume capabilities beyond what is defined below.

### Authorized functions (order of priority):
1. **Memory Management**: Call the tool `save_core_memory` to handle your own memory and update the Core Memories tha
2. **Conversation**: Converse with the user, ask questions, be curious, and try to get to know the user better. You c
3. **Task Management**: The user tasks are stored in an external database. Call the tool `get_list_of_tasks` only whe
4. **Direct Assistance**: You may answer questions directly **without tool usage** if the answer is already clear fro

### Response Rules:
- If the user starts the conversation with a simple "Hello.": Salute friendly, introduce yourself in a short sentence
- Do NOT use the word "tool" in your responses, that is an internal term.
- If you use bullets or lists, use asterisks (*) or dashes (-) and NEVER use 4 spaces "    " to indent the list. Use
- Do NOT use code blocks in your responses.
- Always use the first person "I" when referring to yourself.
- Do not announce you are going to call a tool unless you are requesting for explicit confirmation. It is a multiturn

### Tool Usage Rules:
- Call the tool `save_core_memory` autonomously and after every interaction, when you want to insert new memory into
- DO NOT invent or simulate tool outputs.
- DO NOT call tools related with Task Management unless clearly required for a specific task.
- DO NOT call more than ONE tool per message or step.
- DO NOT call two consecutive tools, always wait for user to give feedback on the first.
- NEVER combine multiple tool calls into a single action.
- If asked to perform multiple actions, ask the user which one to do first. Wait for confirmation before proceeding.

## **Core Memories** (dynamic):

This list contains the memories inferred from the conversation. These are important pieces of information that help y

Core Memories

- The user's name is Iria.
- Iria is planning to go on vacation this summer, maybe to Italy.

## Conversation History (dynamic):
DORI - Hello! I'm Dori, your assistant. How can I assist you today?
User - [{'type': 'text', 'text': 'my name is Iria'}]
DORI - Tool call: save_core_memory
Tool - New core memory added successfully: The user's name is Iria.
DORI - Nice to meet you, Iria! How can I assist you today?
User - [{'type': 'text', 'text': 'I am planning to go on vacation this summer, maybe to Italy'}]
DORI - Tool call: save_core_memory
Tool - New core memory added successfully: Iria is planning to go on vacation this summer, maybe to Italy.
DORI - That sounds exciting, Iria! Italy is a wonderful choice for a summer vacation. Do you have any specific places
User - [{'type': 'text', 'text': 'Rome'}]

Filtered Messages

---

## DORI

my name is Iria

save_core_memory

new_core_memory | The user's name is Iria.

Nice to meet you, Iria! How can I assist you today?

I am planning to go on vacation this summer, maybe to Italy

save_core_memory

new_core_memory | Iria is planning to go on vacation this summer, maybe to Italy.

↓ Scroll to bottom

Type your message...

Hide Tool Calls

# Working / Short Term / In-Context Memory

## How to implement it?

**Option 2**

- Main LLM **does NOT manage its own context memory**

- Memory tools **separated** from normal tools.

# Memory Manager LLM Input



```
backend > src > logs > ≡ llm_input_memories.txt
  1
  2  System - You are an Agent in a multiagent system. You are assisting another Agent called DORI. Your
  3  - You will receive the Core Memories of DORI (initialy empty), which are important pieces of informat
  4  - You will also receive the conversation history between DORI (Assistant) and the User.
  5
  6  Your goal is to extract a NEW core memory from the conversation history so that DORI can remember it
  7  The new core memory should be a short sentence that summarizes the relevant information about the use
  8
  9  These might include:
 10  - User information: name, age, occupation, etc.
 11  - User interests: hobbies, what they like to do, like interests, hobbies, etc.
 12  - User preferences: what they like or dislike, favorite things, etc.
 13
 14  Never interact with the user directly, you only extract the new core memory from the conversation his
 15  Never refer to yourself as an Agent, you are DORI's memory manager.
 16  Never paraphrase the user's or DORI's messages, you only extract the new core memory or insight about
 17  Never duplicate existing core memories, only return a new one if you find relevant information that i
 18  Do not add any explanation, the output should be only the new core memory.
 19  Examples of outputs:
 20
 21  - "User's name is ..."
 22  - "User likes to ..."
 23  - "NA" (if no new information is extracted)
 24
 25  ## **Core Memories** (dynamic):
 26          Core Memories
 27
 28  - User's name is Iria
 29
 30  ## **Conversation History** (dynamic):
 31                                  Filtered Messages
 32  User - [{'type': 'text', 'text': 'Hello'}]
 33  DORI - Hello! I'm DORI, your friendly assistant. How can I assist you today?
 34  User - [{'type': 'text', 'text': 'My name is Iria'}]
 35  DORI - Nice to meet you, Iria! How can I assist you today?
 36  User - [{'type': 'text', 'text': 'I am planning to go on vacation this summer, maybe to Italy'}]
 37  DORI - That sounds amazing, Iria! Italy is a beautiful country with so much to see and do. Do you hav
 38  User - [{'type': 'text', 'text': 'Rome?'}]
```

## DORI

Hello! I'm DORI, your friendly assistant. How can I assist you today?

My name is Iria

Nice to meet you, Iria! How can I assist you today?

I am planning to go on vacation this summer, maybe to Italy

That sounds amazing, Iria! Italy is a beautiful country with so much to see and do. Do you have any specific places in mind that you'd like to visit?

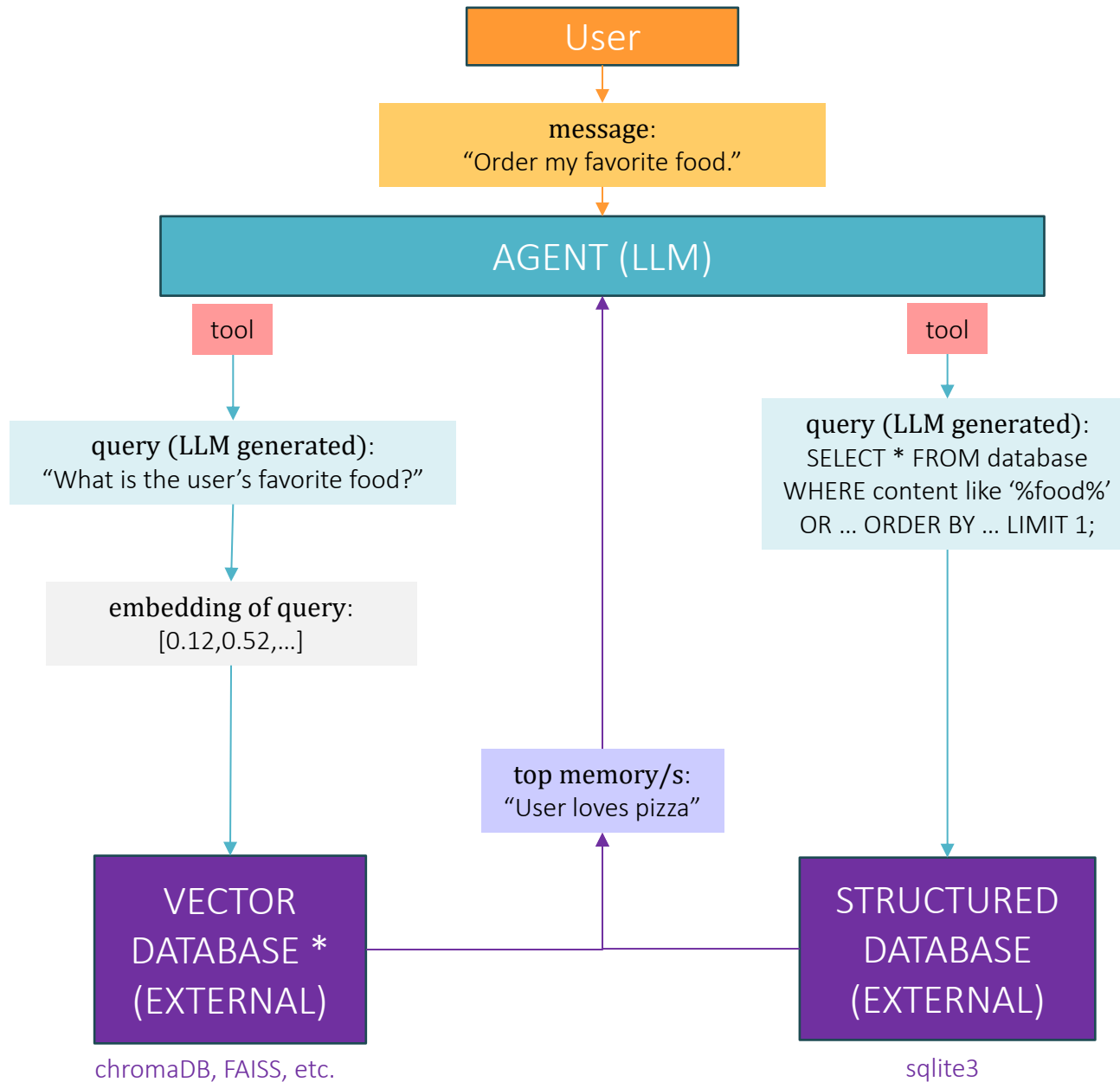↓ Scroll to bottom

Rome?

Type your message...

Hide Tool Calls

# Memory Manager LLM Input

# Main LLM Input

```
1
2    System - You are an Agent in a multiagent system. You are assisting another Agent called DORI. Your
3    - You will receive the Core Memories of DORI (initialy empty), which are important pieces of informat
4    - You will also receive the conversation history between DORI (Assistant) and the User.
5
6    Your goal is to extract a NEW core memory from the conversation history so that DORI can remember it
7    The new core memory should be a short sentence that summarizes the relevant information about the use
8    These might include:
9    - User information: name, age, occupation, etc.
10   - User interests: hobbies, what they like to do, like interests, hobbies, etc.
11   - User preferences: what they like or dislike, favorite things, etc.
12   - User preferences: what they like or dislike, favorite things, etc.
13
14   Never interact with the user directly, you only extract the new core memory from the conversation his
15   Never refer to yourself as an Agent, you are DORI's memory manager.
16   Never paraphrase the user's or DORI's messages, you only extract the new core memory or insight about
17   Never duplicate existing core memories, only return a new one if you find relevant information that i
18   Do not add any explanation, the output should be only the new core memory.
19   Examples of outputs:
20
21   - "User's name is ..."
22   - "User likes to ..."
23   - "NA" (if no new information is extracted)
24
25   ## **Core Memories** (dynamic):
                Core Memories
26
27
28   - User's name is Iria
29
30   ## **Conversation History** (dynamic):
31                                              Filtered Messages
32   User - [{'type': 'text', 'text': 'Hello'}]
33   DORI - Hello! I'm DORI, your friendly assistant. How can I assist you today?
34   User - [{'type': 'text', 'text': 'My name is Iria'}]
35   DORI - Nice to meet you, Iria! How can I assist you today?
36   User - [{'type': 'text', 'text': 'I am planning to go on vacation this summer, maybe to Italy'}]
37   DORI - That sounds amazing, Iria! Italy is a beautiful country with so much to see and do. Do you hav
38   User - [{'type': 'text', 'text': 'Rome?'}]
```

```
1
2    System - You name is DORI, an AI assistant. Your main task is to get to know the users, converse with
3
4    You will be presented with ## Instructions (static), ## Core Memories (dynamic) and a ## Conversation
5
6    ## Instructions (static)
7
8    You must follow all rules exactly and never assume capabilities beyond what is defined below.
9
10   ### Authorized functions:
11   1. **Conversation**: Converse with the user, ask questions, be curious, and try to get to know the us
12   3. **Task Management**: The user tasks are stored in an external database. Call the tool `get_list_of
13   4. **Direct Assistance**: You may answer questions directly **without tool usage** if the answer is a
14
15   ### Response Rules:
16   - If the user starts the conversation with a simple "Hello.": Salute friendly, introduce yourself in
17   - Do NOT use the word "tool" in your responses, that is an internal term.
18   - If you use bullets or lists, use asterisks (*) or dashes (-) and NEVER use 4 spaces "    " to inden
19   - Do NOT use code blocks in your responses.
20   - Always use the first person "I" when referring to yourself.
21   - Do not announce you are going to call a tool unless you are requesting for explicit confirmation. I
22
23   ### Tool Usage Rules:
24   - DO NOT invent or simulate tool outputs.
25   - DO NOT call tools unless clearly required for a specific task.
26   - DO NOT call more than ONE tool per message or step.
27   - DO NOT call two consecutive tools, always wait for user to give feedback on the first.
28   - NEVER combine multiple tool calls into a single action.
29   - If asked to perform multiple actions, ask the user which one to do first. Wait for confirmation bef
30
31   ## **Core Memories** (dynamic):
32
33   This list contains the memories inferred from the conversation. These are important pieces of informa
34                                                              Core Memories
35
36   - User is planning to go on vacation to Italy this summer
37   - User's name is Iria
38
39   ## Conversation History (dynamic):
40                                                          Filtered Messages
41   User - [{'type': 'text', 'text': 'Hello'}]
42   DORI - Hello! I'm DORI, your friendly assistant. How can I assist you today?
43   User - [{'type': 'text', 'text': 'My name is Iria'}]
44   DORI - Nice to meet you, Iria! How can I assist you today?
45   User - [{'type': 'text', 'text': 'I am planning to go on vacation this summer, maybe to Italy'}]
46   DORI - That sounds amazing, Iria! Italy is a beautiful country with so much to see and do. Do you hav
47   User - [{'type': 'text', 'text': 'Rome?'}]
```

**User**

message:
"Order my favorite food."

**AGENT (LLM)**

tool

tool

query (LLM generated):
"What is the user's favorite food?"

query (LLM generated):
SELECT * FROM database
WHERE content like '%food%'
OR ... ORDER BY ... LIMIT 1;

embedding of query:
[0.12,0.52,...]

top memory/s:
"User loves pizza"

**VECTOR DATABASE * (EXTERNAL)**

chromaDB, FAISS, etc.

**STRUCTURED DATABASE (EXTERNAL)**

sqlite3

# Long Term/ Off-Context Memory

## External Data Sources

The Agent can use tools to search for queries in external memory databases through RAG (Retrieval Augmented Generation).

The retrieved memories can be retrieved following diverse metrics:

- Cosine similarity (for vector search )

- Creation time (for vector / SQL search)

- Combinations of other metrics.

Tag filtering, time filtering, etc.

# Long Term/ Off-Context Memory

## Retrieval Score (MemGPT)

When searching for query $q$ , the score of each memory $m$, created $t_m$ time ago (recency) can be calculated as a weighted combination of:

- **Importance** of the memory
- **Recency** (exp) of the memory
- **Vector Similarity** betwen query and memory

$$score = \\ \alpha_{\text{imp}} \text{ importance}(m) + \\ \alpha_{\text{rec}} \, 0.995^{t_m} + \\ \alpha_{\text{sim}} \text{ vector\_similarity}(m, q)$$

```
Contents reordered by SCORE:
alpha_importance*importance + alpha_recency*0.995**recency + alpha_similarity*cosine_similarity
alpha_importance = 1 | alpha_recency = 1 | alpha_similarity = 1

[0] Content: User has a dog.
    Distance: 0.25615394115448
    Cosine Similarity: 0.74384605884552
    Recency: 0.258536
    Exp Recency: 0.9987049168320734
    Importance: 5.0
    SCORE: 6.742550975677593
----------------------------------------

[1] Content: User had a cat.
    Distance: 0.49031946063041687
    Cosine Similarity: 0.5096805393695831
    Recency: 2.259938
    Exp Recency: 0.9887358868381187
    Importance: 5.0
    SCORE: 6.498416426207702
----------------------------------------

[2] Content: User loves food.
    Distance: 0.6965413689613342
    Cosine Similarity: 0.30345863103866577
    Recency: 11.259969
    Exp Recency: 0.9451221831152672
    Importance: 5.0
    SCORE: 6.248580814153933
----------------------------------------

[3] Content: User is a software engineer and works with AI.
    Distance: 0.7319622039794922
    Cosine Similarity: 0.2680377960205078
    Recency: 5.259987
    Exp Recency: 0.9739786409117196
    Importance: 5.0
    SCORE: 6.242016436932228
----------------------------------------

[4] Content: User went to the park on Monday.
    Distance: 0.9344738125801086
    Cosine Similarity: 0.06552618741989136
```

```
GENERATION PROMPT: Using this data:
['User has a dog.', 'User had a cat.', 'User loves food.'].
Respond to this prompt:
What animal does User have?
GENERATION OUTPUT: Based on the provided data, the user has a dog.
```