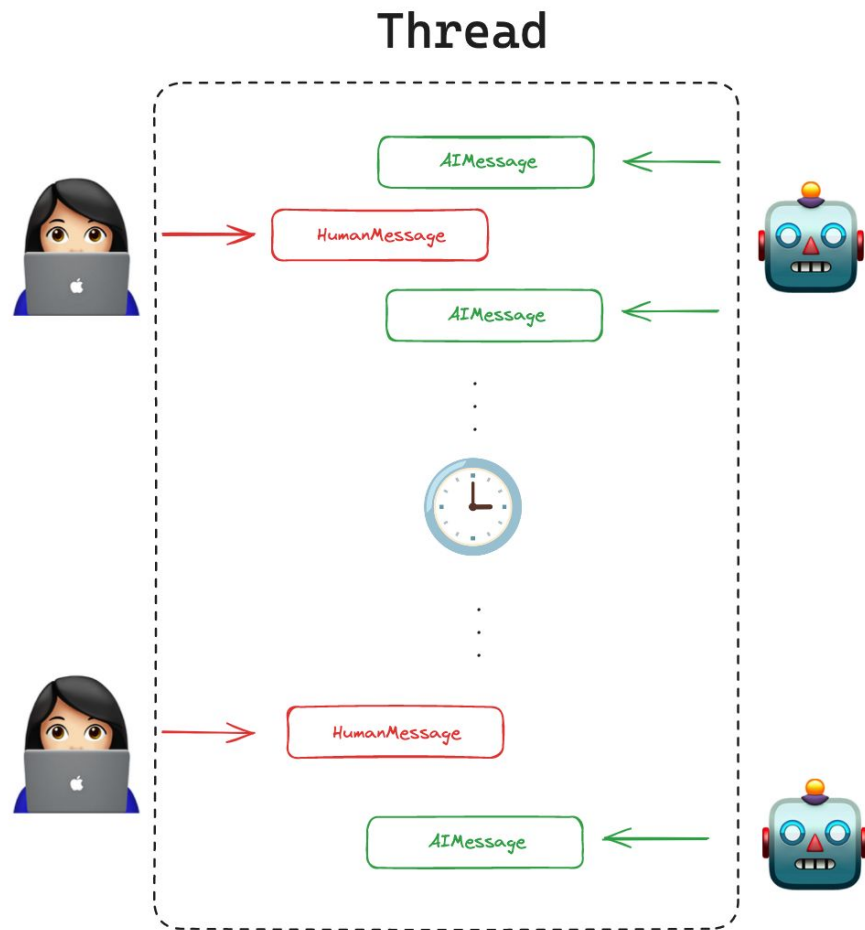


LangChain Academy

Introduction to  LangGraph

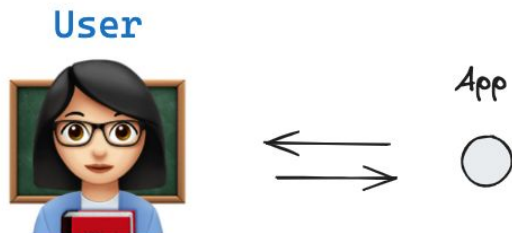
Memory is a cognitive function that allows people to store, retrieve, and use information to understand their present and future.

1 Within session (thread) memory



1 Across session (thread) memory

First chat session

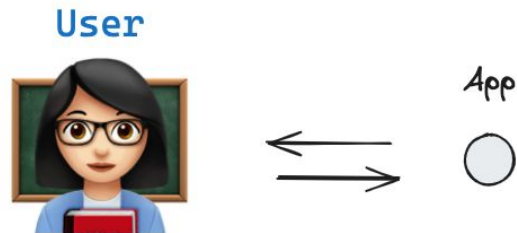


Save user specific
details, facts, etc

Memories



Second chat session



Use saved user specific
information to improve
application experience

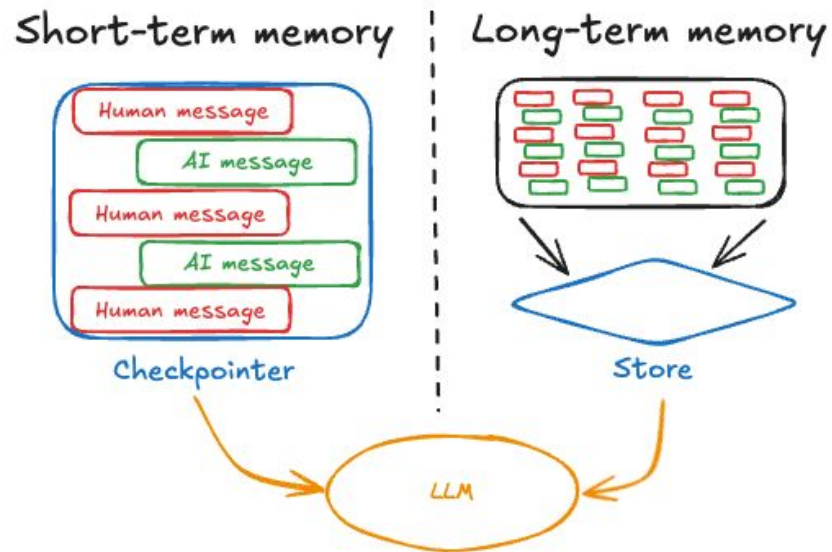
Memories



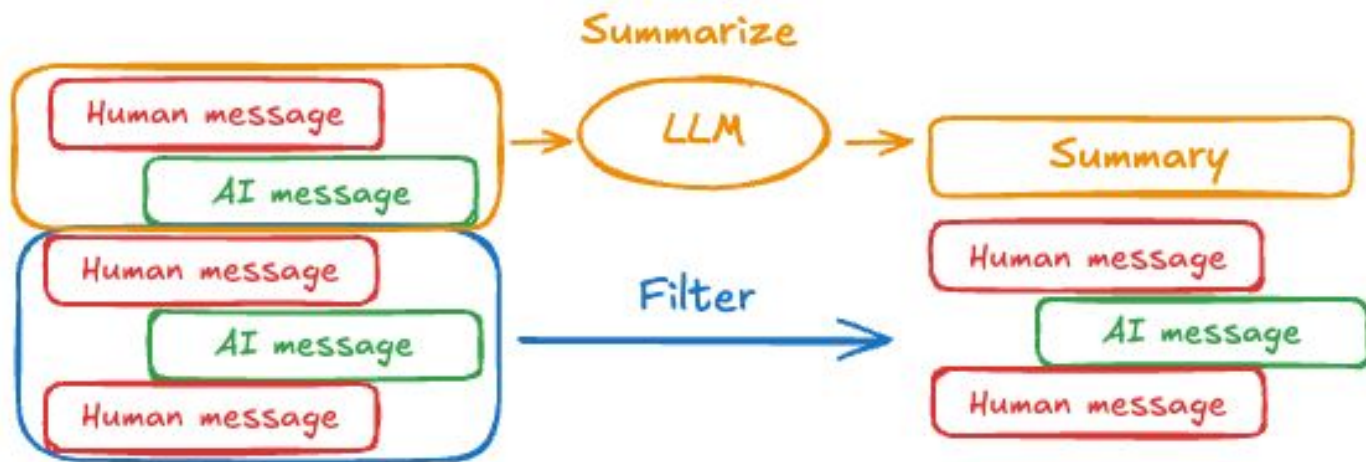
1 Comparison

	Short-term	Long-term
Scope	Within session (thread)	Across session (thread)
Example use-case	Persist conversational history, allow interruptions in a chat (e.g., if user is idle or to allow human-in-the-loop)	Remember information about a specific user across all chat sessions
LangGraph usage	Checkpointers	Store

1 Comparison



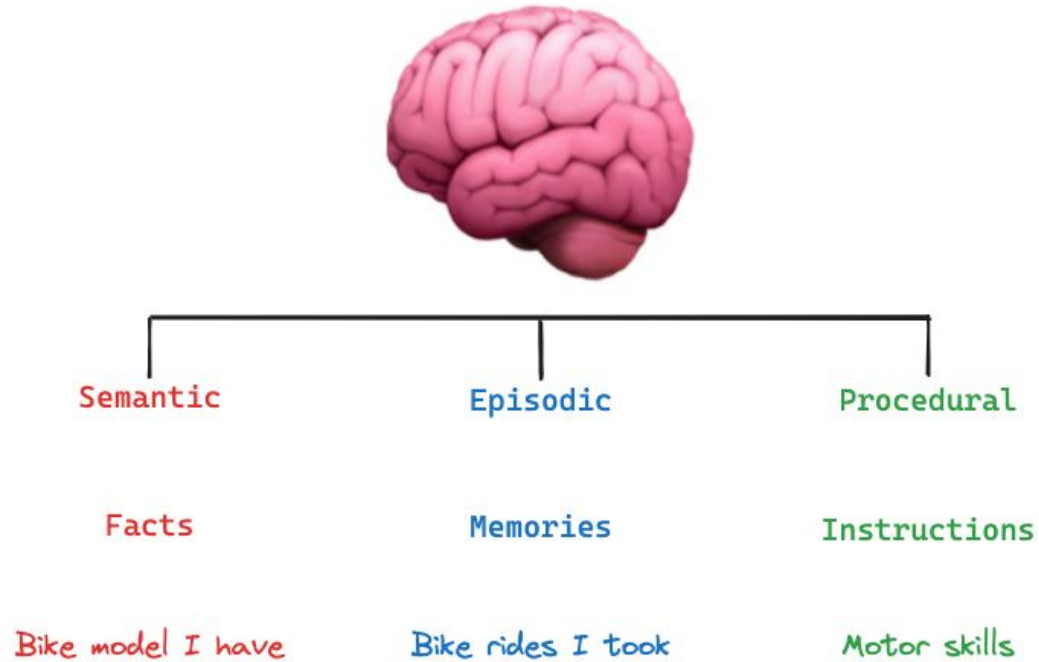
2 Short-term memory



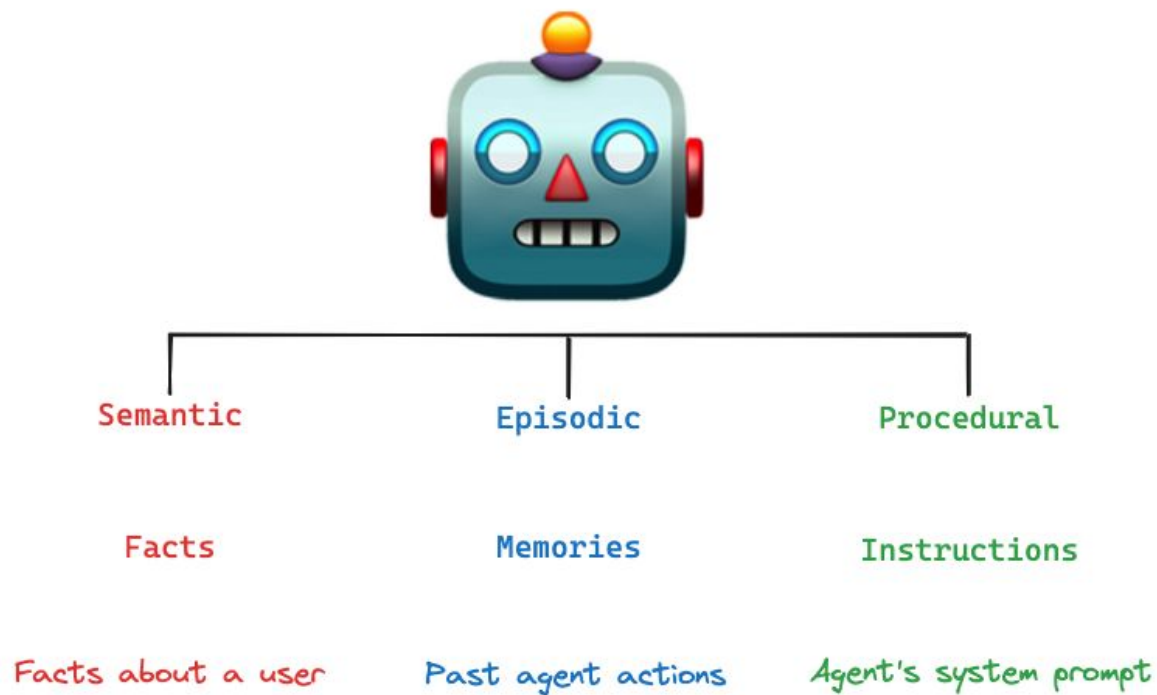
What is the type of memory?

When do you want to update memories?

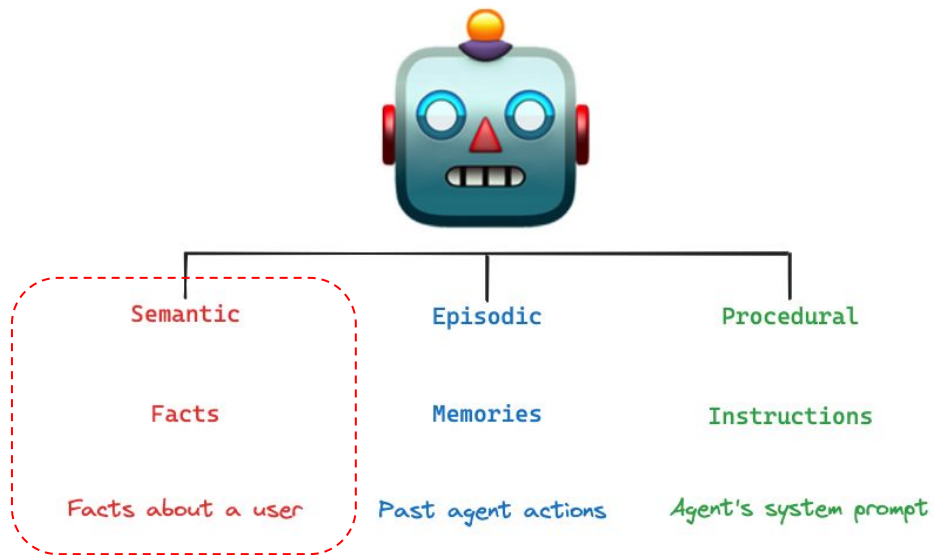
3 Long-term memory: Humans



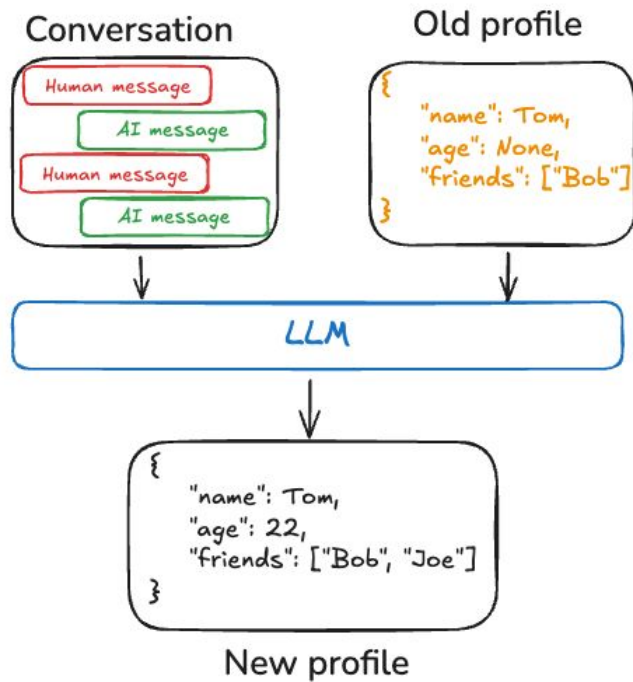
3 Long-term memory: Agents



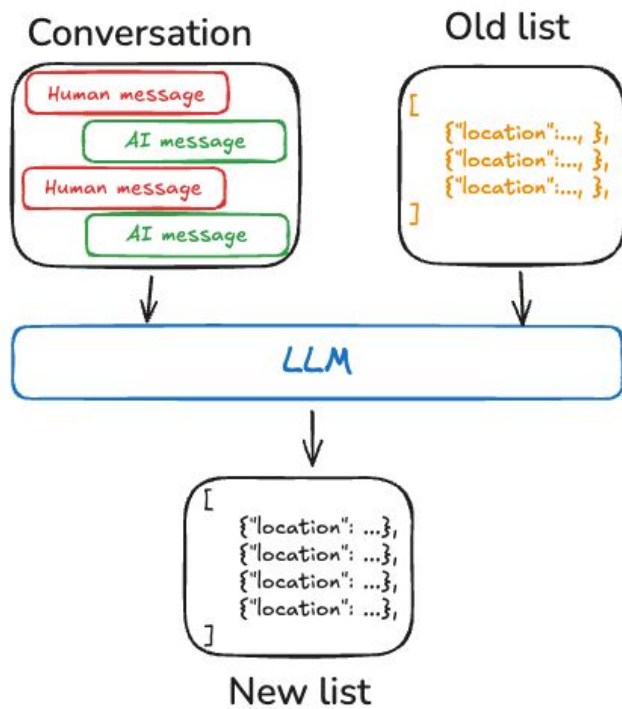
How to structure facts?



3 Long-term memory: Semantic (Profile)



3 Long-term memory: Semantic (Collection)

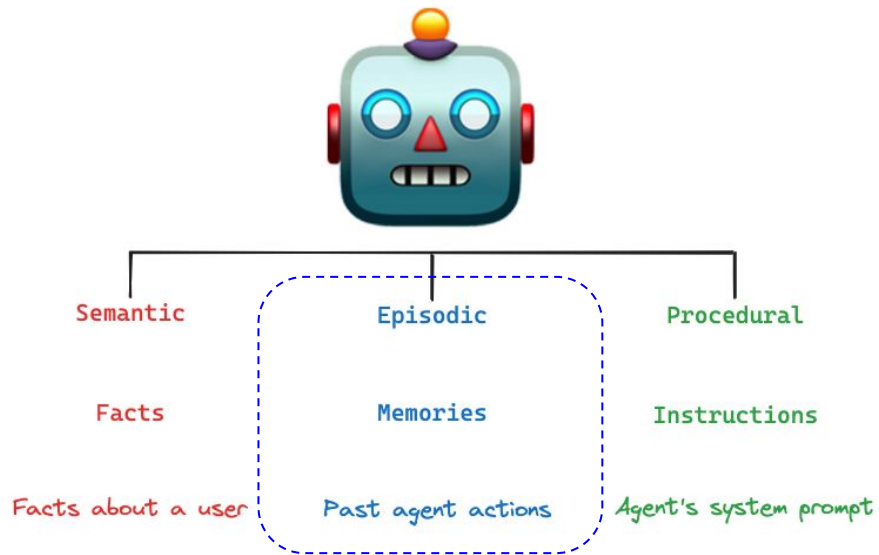


3 Long-term memory: Semantic (Comparison)

	Type	Pro	Con
Profile	Single document	Easily retrieved single representation	Challenging to maintain as it grows larger
List	List of documents	Allows for smaller, narrowly scoped memories w/ easy addition of new information	Retrieval can be challenge as the list grows larger

3 Long-term memory: Episodic

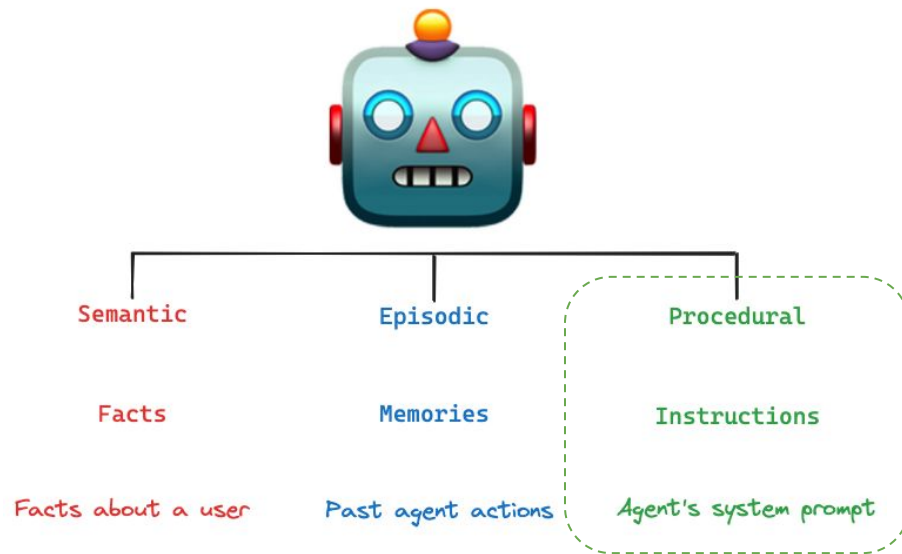
How to present past actions to an agent?



Prior reasoning trajectories

```
system: You are requested to solve math questions in an alternate  
universe. The operations have been altered to yield different  
than expected. Do not guess the answer or rely on your innate  
of math. Use the provided tools to answer the question. When  
associativity and commutativity apply, distributivity does  
Answer the question using the fewest possible tools. Only  
numeric response without any clarifications. Here are some  
conversations of the user interacting with the AI until the  
answer is reached:  
  
user: evaluate the negation of -100  
assistant:  
    tool_calls: [{"name": "negate", "args": {"a": -100}}]  
tool (negate): -100  
assistant: So the answer is 100.  
user: 100 is incorrect. Please refer to the output of your tool call.  
assistant:  
    content: You're right, my previous answer was incorrect. Let  
using the tool output  
    tool_calls: [{"name": "negate", "args": {"a": -100}}]  
tool (negate): -100  
assistant: The answer is -100.0
```


How to update agent instructions?



LARGE LANGUAGE MODELS ARE HUMAN-LEVEL PROMPT ENGINEERS

**Yongchao Zhou^{1,2,*}, Andrei Ioan Muresanu^{2,3,*}, Ziwon Han^{1,2,*}, Keiran Paster^{1,2},
Silviu Pitis^{1,2}, Harris Chan^{1,2}, Jimmy Ba^{1,2}**

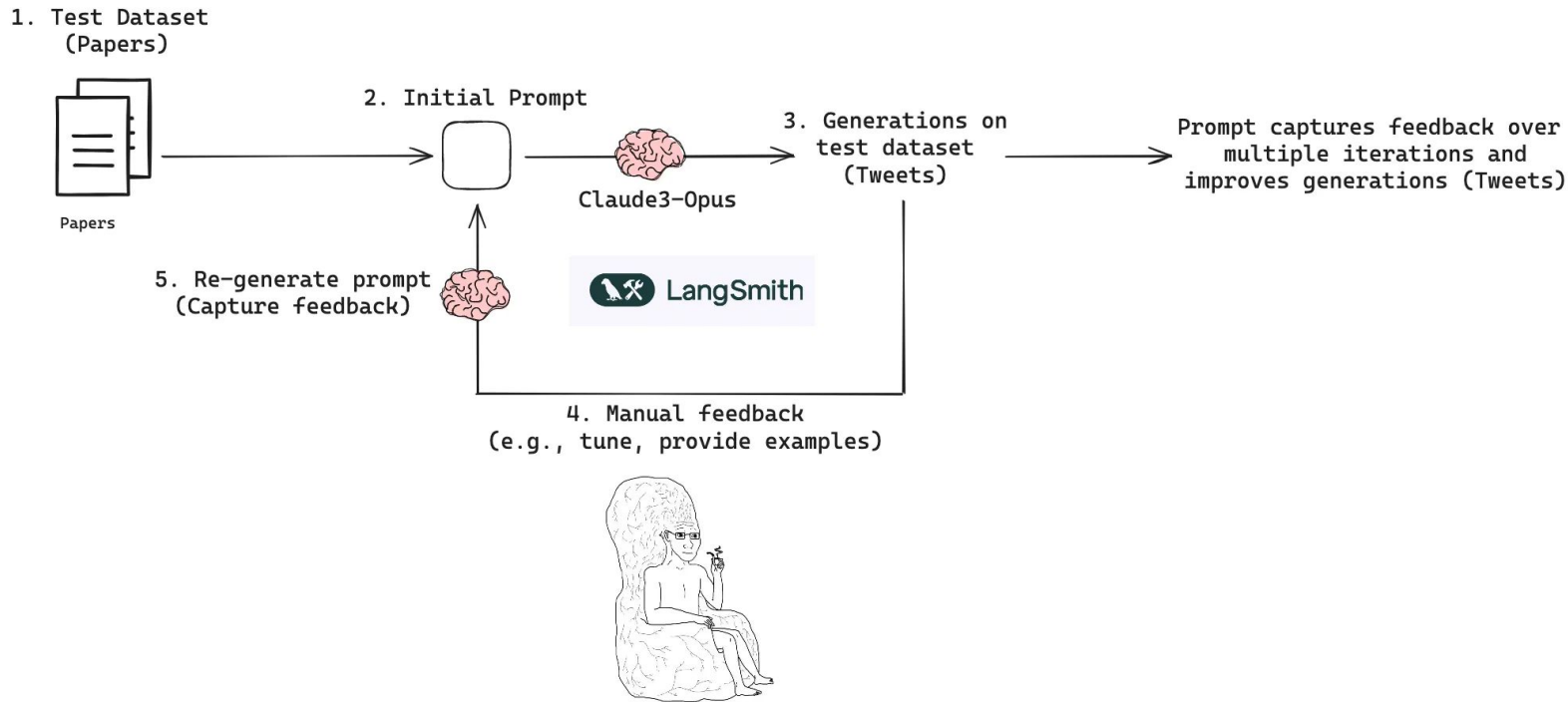
¹University of Toronto ²Vector Institute ³University of Waterloo *Equal contribution
{yczhou, hanzhiwen, keirp, spitis, hchan, jba}@cs.toronto.edu
{andrei.muresanu}@uwaterloo.ca

ABSTRACT

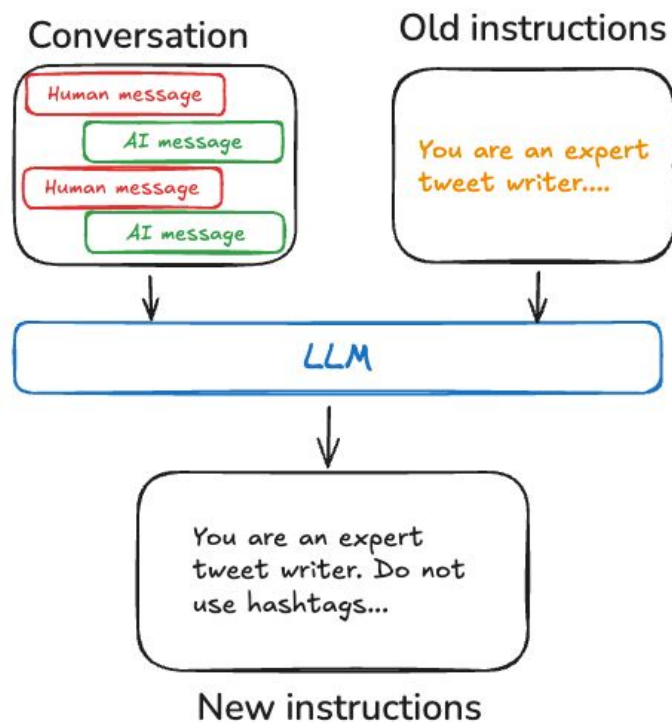
By conditioning on natural language instructions, large language models (LLMs) have displayed impressive capabilities as general-purpose computers. However, task performance depends significantly on the quality of the prompt used to steer the model, and most effective prompts have been handcrafted by humans. Inspired by classical program synthesis and the human approach to prompt engineering, we propose *Automatic Prompt Engineer*¹ (APE) for automatic instruction generation and selection. In our method, we treat the instruction as the “program,” optimized by searching over a pool of instruction candidates proposed by an LLM in order to maximize a chosen score function. To evaluate the quality of the selected instruction, we evaluate the zero-shot performance of another LLM following the selected instruction. Extensive experiments show that our automatically generated instructions outperform the prior LLM baseline by a large margin and achieve better or comparable performance to the instructions generated by human annotators on 24/24 Instruction Induction tasks and 17/21 curated BIG-Bench tasks. We conduct extensive qualitative and quantitative analyses to explore the performance of APE. We show that APE-engineered prompts are able to improve few-shot learning performance (by simply prepending them to standard in-context learning prompts), find better zero-shot chain-of-thought prompts, as well as steer models toward truthfulness and/or informativeness.²

3 Long-term memory: Procedural

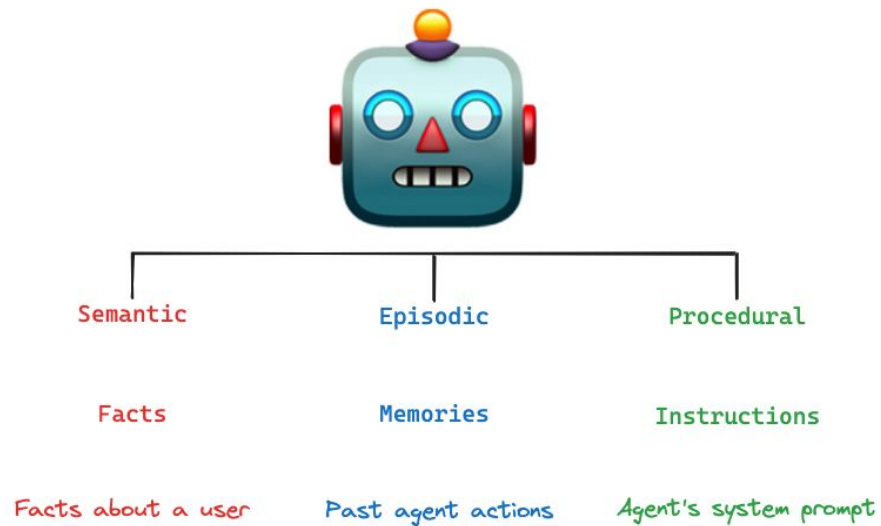
Teaching Claude3 to Tweet like a pro using iterative prompt engineering



3 Long-term memory: Procedural



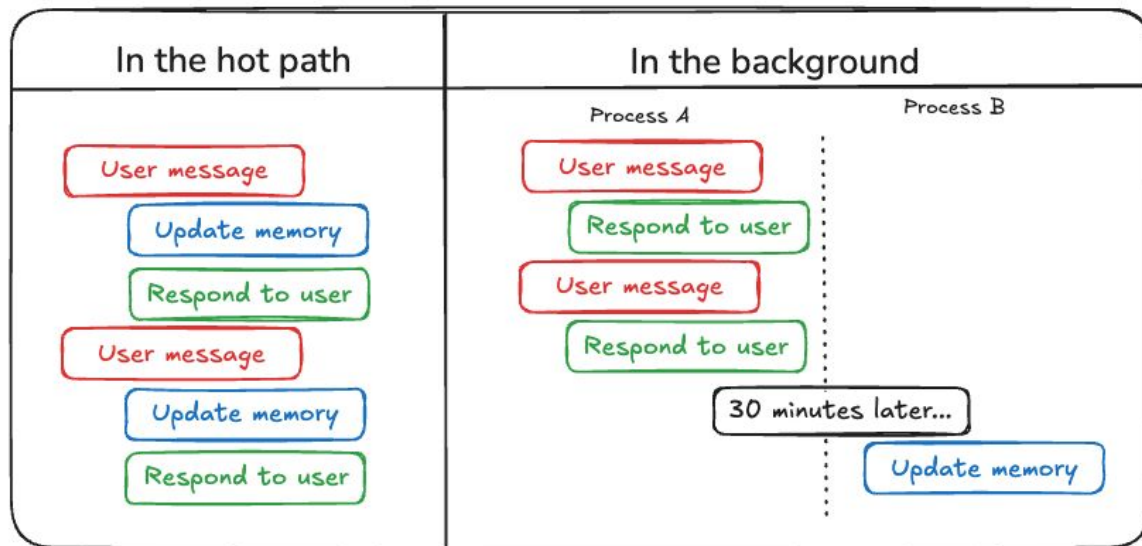
What are the memory types?



What is the type of memory?

When do you want to update memories?

3 Long-term memory updating



3 Long-term memory updating (Comparison)

	Type	Pro	Con
Hot-path	During runtime (ChatGPT)	Real-time updates with transparency for user	Can affect UX / latency and degrade performance
Background	As a separate process	Lower risk of UX / performance degradation	Frequency of memory writing needs to be tuned

4 End goal

