

# Memolet: Reifying the Reuse of User-AI Conversational Memories

by

Hen-Chen Yen

A thesis  
presented to the University of Waterloo  
in fulfillment of the  
thesis requirement for the degree of  
Master of Mathematics  
in  
Computer Science

Waterloo, Ontario, Canada, 2024

© Hen-Chen Yen 2024

## **Author's Declaration**

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

## Abstract

As users engage more frequently with AI conversational agents, conversations may exceed their “memory” capacity, leading to failures in correctly leveraging certain memories for better responses. Therefore, users have to revisit related memories and re-provide these memories to the agents, ensuring that the generation refers to the accurate memories. However, the process of finding past memories to reuse is cumbersome, requiring users to retrieve related information across various conversations and articulate their intentions for reusing these memories to the AI. To support users in recalling and reusing relevant memories, we introduce *Memolet*, an interactive object that reifies memory reuse. Users can directly manipulate *Memolet* to specify which memories to reuse and how to use them. We developed a system demonstrating *Memolet*’s interaction across various memory reuse stages, including memory extraction, organization, prompt articulation, and generation refinement. Through a user study, we gained insights into users’ experiences with *Memolet* for memory reuse in AI conversations. The study validates the system’s usefulness and provides design implications for future systems that support user-AI conversational memory reusing.

## Acknowledgments

I would like to express my sincere gratitude to Dr. Jian Zhao for his guidance throughout my research. His insightful feedback and support were instrumental in shaping the direction of my work and helping me overcome various challenges. Also, thank you for choosing me over my partner for admission and I hope I did not disappoint you with this choice.

Many thanks to Dr. Daniel Vogel and Dr. Jimmy Lin for agreeing to serve as reviewers for my dissertation. Your constructive feedback and scholarly insights greatly enhanced my research. Thank you Dr. Lin for your surprise visits and chats in the lab. Dr. Vogel, you have shown me what true perfectionism is, which I thought I was until I met you.

I owe a debt of gratitude to all my labmates, mentors, and collaborators during my studies, including Dr. Sangho Suh, Dr. Nicole Sultanum, Dr. Haijun Xia, Dr. Damien Masson, Dr. Zhicong Lu, and more. Dr. Suh, thanks for all those bagels you fed me in Korea and valuable feedback throughout the collaboration. Thanks to Dr. Masson for letting me know the direction of the research I want to pursue in the future, and for tirelessly criticizing my title and introduction. Dr. Lu, thanks for pulling me into the HCI and for being such a great gossip buddy. Who knew research could be this fun? And a special thanks to my friend Carter for all the fun and rice cakes. See you at MIT.

I want to thank Felicia for her continuous support, understanding, and encouragement. Her patience and belief in me have been my strongest pillars of support. As Nikhita suggested, maybe we should not always only talk about research on the phone, but I could not have gotten this far without your advice and support in research and mental. Finally, I am profoundly grateful to my family for their support, love, and encouragement throughout this journey. Their belief in me has been my greatest motivation.

# Table of Contents

<b>Author’s Declaration</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>Acknowledgments</b>	<b>iv</b>
<b>List of Figures</b>	<b>viii</b>
<b>List of Tables</b>	<b>x</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Background</b>	<b>4</b>
2.1 Information Sensemaking and Knowledge Reusing . . . . .	4
2.2 Factual Text Generation . . . . .	5
2.3 Memory Reuse in AI-Conversational Agents . . . . .	5
<b>3 Design Process</b>	<b>7</b>
3.1 Motivating Scenario and Design Guidelines . . . . .	7
3.2 Motivating Scenario . . . . .	8
3.3 [DG1] Interacting with Memories at Multi-Layers . . . . .	9
3.4 [DG2] Provide Visual Cues to Recall and Extract Memories . . . . .	10
3.5 [DG3] Flexibly Externalizing Users’ Sensemaking Results about Memories .	10
3.6 [DG4] Aligning Users’ Intentions to Reuse Memories by Direct Manipulation	11

<b>4</b>	<b>System</b>	<b>12</b>
4.1	Reifying the Reuse of Memory . . . . .	13
4.2	System Overview and Multi-Layer Interaction with Memolet [DG1] . . . . .	13
4.3	Memories Recall and Extraction [Layer1/DG2] . . . . .	14
4.4	Memories Organization and Schematization [Layer2/DG3] . . . . .	15
4.5	Generating with Memories [Layer3/DG4] . . . . .	17
4.6	System Implementation . . . . .	19
4.6.1	Retrieval Augmented Generation . . . . .	19
4.6.2	Instructed Generation . . . . .	20
4.6.3	System Architecture . . . . .	20
<b>5</b>	<b>Study</b>	<b>22</b>
5.1	Participants . . . . .	22
5.2	Scenarios and Tasks . . . . .	23
5.3	Phase One Study Setup . . . . .	23
5.3.1	Study Procedure . . . . .	23
5.3.2	Collected Conversations . . . . .	23
5.4	Phase Two Study Setup . . . . .	24
5.4.1	Procedure . . . . .	24
5.4.2	Baseline System . . . . .	24
5.4.3	Measures . . . . .	24
5.4.4	Data Analysis . . . . .	25
<b>6</b>	<b>Findings</b>	<b>27</b>
6.1	Findings . . . . .	27
6.2	Overall Usage of <i>Memolet</i> in Memory-Reusing Process . . . . .	28
6.3	Recalling and Extracting Memories . . . . .	29
6.4	Externalizing Users' Sensemaking Results to Lower the Cognitive Load . . . . .	31
6.5	Aligning Users' Memory Reuse Intention . . . . .	32

<b>7</b>	<b>Discussion</b>	<b>35</b>
<b>8</b>	<b>Limitations and Future Work</b>	<b>37</b>
<b>9</b>	<b>Conclusion</b>	<b>38</b>
	<b>References</b>	<b>39</b>
	<b>APPENDICES</b>	<b>50</b>
<b>A</b>	<b>Study Supplementary</b>	<b>51</b>
A.1	Questionnaire . . . . .	51
A.1.1	UMUX-LITE . . . . .	51
A.1.2	NASA-TLX . . . . .	51
A.1.3	Self-Defined Likert Scale Items . . . . .	52
A.2	Study Scenarios and Tasks . . . . .	52
A.2.1	<b>Scenario 1 (Expository Writing) Phase One</b> . . . . .	52
A.2.2	<b>Scenario 1 (Expository Writing) Phase Two</b> . . . . .	53
A.2.3	<b>Scenario 2 (Programming) Phase One</b> . . . . .	53
A.2.4	<b>Scenario 2 (Programming) Phase Two</b> . . . . .	54
A.2.5	<b>Scenario 3 (Trip Planning) Phase One</b> . . . . .	54
A.2.6	<b>Scenario 3 (Trip Planning) Phase Two</b> . . . . .	55
A.3	Demographic Table . . . . .	56
A.4	Prompt Template . . . . .	57
A.4.1	Prompt for Retrieval Augmented Generation . . . . .	57
A.4.2	modify query with instructions . . . . .	58
A.4.3	instructed rag prompt . . . . .	59
A.4.4	summarize prompt . . . . .	60
A.4.5	summarize chat history prompt . . . . .	60
A.4.6	generate more queries prompt . . . . .	62

# List of Figures

1.1	The figure illustrates how users reuse memories for future generations by interacting with <i>Memolet</i> , the reification of memory reuse. (A) Our system initially embeds and projects all users' conversations to the long-term memory repository. (B) Users can search, recall, and extract <i>Memolets</i> from this repository and transfer them to the curated memory sandbox. (C) This sandbox supports users in organizing and schematizing <i>Memolets</i> based on their own sensemaking. (D) Finally, users can reuse these <i>Memolets</i> by referring to them in the prompt and (E) refine the generation through direct manipulation. . . . .	2
3.1	User-AI conversational memories reusing process with four Design Guidelines and Challenges. The first three stages of memory reuse are derived from the foraging loop within the information sensemaking process [70], outlining the processes of users seeking information, searching and filtering it, and reading and extracting information. Additionally, we draw inspiration from knowledge externalization strategies, where users extract, organize, and integrate pieces of information to scaffold the comprehension process [35]. These stages are further mapped to the knowledge reusing framework [1], which encompasses capturing and documenting knowledge, packaging and distributing knowledge, and reusing knowledge. . . . .	7
3.2	Long-Term Memory Repository. (A) Semantic search to find related <i>Memolets</i> ; (B) Adjust parameters to re-cluster or modify the bin size of <i>Memolet</i> ; (C) Each <i>Memolet</i> is accompanied by visual cues such as icons, colors, and summaries; (D) Each <i>Memolet</i> may contain one to many pairs of prompts/responses based on semantic meaning. . . . .	11



4.1	Our system design process underwent several iterations based on feedback from four participants. . . . .	12
4.2	The Curated Memory Sandbox. (A) Users can organize and schematize the extracted <i>Memolets</i> from the repository based on their own sensemaking results; (B) All movement will be snapped to the active grid; (C) When multiple <i>Memolets</i> get closer, they will be grouped together; (D) Users can extract <i>Memolets</i> from prompts/responses or selected text. . . . .	14
4.3	The memory buffer considers the <i>Memolet</i> during generation. (A) Users can refer to <i>Memolets</i> in the sandbox while articulating prompts; (B) The generated results will cite the <i>Memolet</i> ; (C) Users first attempt to refine the generation by interpolating two <i>Memolets</i> , deselecting one <i>Memolet</i> from citations; (D) Users regenerate results by resizing a <i>Memolet</i> to highlight context contained. . . . .	16
5.1	Our adopted <i>Baseline</i> , resembling the current prevalent AI-driven conversational agent ChatGPT, allows users to semantically search related context (A) and switch or create new conversations (B). To ensure a fair comparison with our system, we adopted the same iterative conversational history summarization methods (C) and the retrieval augmented generation approach (D). The generation is powered by the latest large language model, GPT-4 [60].	25
6.1	Participants' responses when rating the 5-point self-defined Likert scale questionnaire for both our <i>Baseline</i> and our system. Dots represent the mean differences of our system compared to the <i>Baseline</i> . Bars indicate the 95% CI calculated using the studentized bootstrap method. . . . .	27
6.2	Distribution of logged events across normalized time aggregated by 12 participants, comparing the <i>Baseline</i> and our system. The graph indicates an increase in conversations (with AI) during the later stages of the study for our system, with participants extracting more in the early stage compared to the <i>Baseline</i> . . . . .	28
6.3	Distribution of system log events comparing the <i>Baseline</i> System and our system. Black dots represent means, and the bars denote 95% CI. . . . .	30
6.4	Left: Whether prompts referring to memory. Right: Categories of prompts.	30
6.5	The examples of participants' curated memory sandbox after organizing and manipulating <i>Memolets</i> . The arrow indicates the order in which they synthesized their report. . . . .	34

# List of Tables

A.1 Summary of participant demographics and experience levels. . . . .	56
--	----

# Chapter 1

## Introduction

Recent advances in generative AI-driven conversational agents have become a common method for users to perform tasks in various domains [12, 53]. As users engage in more conversations and share additional details, they may discover valuable contextual information scattered at multiple previous conversations that can enrich their current conversation with the AI [68, 75, 49, 31]. For instance, users might implement similar code from previously provided documentation, create social media posts using specific templates, or synthesize reports from previously discussed topics. In such scenarios, users may encounter difficulties in resuming their conversations from where they left off, as the model may not consistently retain all pertinent memories within the input context’s length [49, 95]. Additionally, users may need to interject information from various conversations without causing distractions to the model [79]. Consequently, it is important to enable users to reuse their “*memories*”—past conversations between users and generative AI [7, 96, 32]. Reusing these memories helps users reduce the need for time-consuming prompt engineering from scratch [91, 21, 90], ensure the generated results are trustworthy [55], and tailored responses to the particular context without hallucination disconnected from the memory [33, 79].

However, users face challenges in reusing memories due to the opaque nature of how current AI-driven conversational agents handle memory [92, 78, 101, 32]. Users lack understanding of how much information is memorized by the AI and have limited control over the memory management strategies of these AI-driven conversational agents [7, 96, 32]. Additionally, users have difficulty discerning which memories are being used for generation, which hinders their ability to assess if the model accurately reuses desired memories for the current task [52, 98, 23, 50]. Therefore, to gain control over AI generation and to ensure that specific prior memories are reused without hallucinating, users often need to

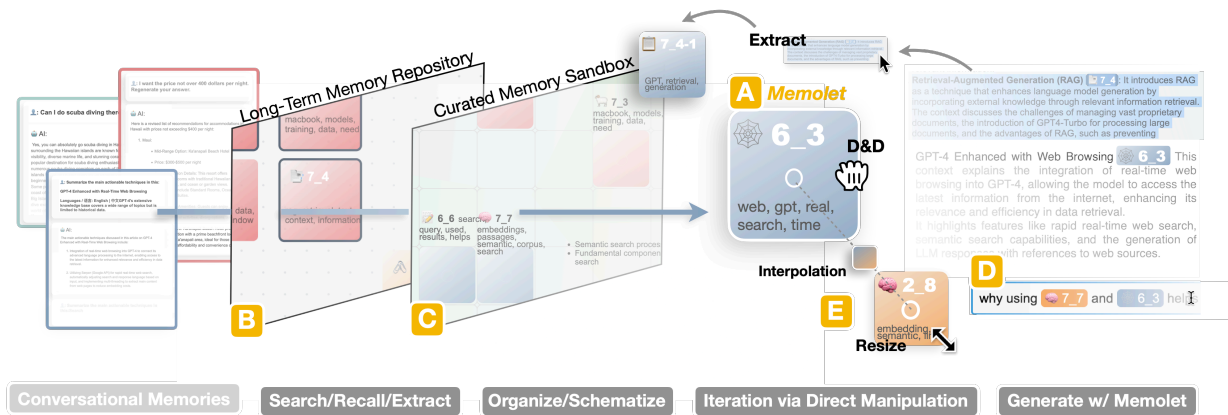


Figure 1.1: The figure illustrates how users reuse memories for future generations by interacting with *Memolet*, the reification of memory reuse. (A) Our system initially embeds and projects all users’ conversations to the long-term memory repository. (B) Users can search, recall, and extract *Memolets* from this repository and transfer them to the curated memory sandbox. (C) This sandbox supports users in organizing and schematizing *Memolets* based on their own sensemaking. (D) Finally, users can reuse these *Memolets* by referring to them in the prompt and (E) refine the generation through direct manipulation.

sift through numerous pairs of prompts/responses to find the relevant context and manually copy and paste it into the new conversation with AI. This process can be quite challenging and time-consuming, as it requires users to make sense of and recall memories [37, 85, 48, 66], extract relevant memories from various conversations [5, 10, 29], organize and integrate these memories based on their usage [14, 65], specify how AI should reuse the memories [97], and iterate on this process until the generation satisfies the users’ needs [82, 47].

To facilitate users throughout the memory-reusing process and empower them with greater control over memory usage, we aim to explore designs that enable users direct control over how they want to reuse memories during conversations with generative AI. Derived from prior theories on knowledge reusing [1], information foraging [69, 70] and knowledge externalization [54], we identified several challenges and design guidelines to support users across the stages of the memory reusing process: from *extracting* related memories, *organizing* memories in a curated space, *articulating* prompts that reference these memories, and finally *refining* generation to reconcile with users’ intentions. We introduce a novel concept, *Memolet*, which reifies the notion of *reusing memories* from past conversations with generative AI (Figure 1.1). *Memolet* is an interactive first-class object that enables users to specify what and how the memory should be reused by direct manip-

ulation. Users can begin by searching and extracting related *Memolets* from a long-term memory repository consisting of all past conversations with AI to specify what memories should be reused for the current task (Figure 1.1.B). Then, users can organize and schematize these extracted *Memolets* within a curated space, externalizing their thoughts on how these memories are related (Figure 1.1.C). Afterward, users can articulate prompts referencing these *Memolets* to specify how they should be reused. Finally, users can refine the AI generation by manipulating the referenced *Memolets* to align with their intentions (Figure 1.1.D&E).

In our evaluation of *Memolet*, we conducted a two-phase within-subject study involving 12 participants who regularly converse with generative AI. We demonstrate the versatility of *Memolet* in three distinct scenarios (i.e., expository writing, programming, and travel planning), wherein participants were tasked with interacting with both our system and *Baseline* in phase two, reusing conversations gathered from phase one. Our findings suggest that participants can better recall past memories, have a lower cognitive load in organizing multiple memories, have greater perceived control over the generative process, and are more intuitive in expressing how they would like to reuse memories. Overall, this paper has three-fold contributions:

1. A novel concept of *Memolet*, a reification of memory reuse that supports users in interacting with past memories.
2. A system that facilitates user interaction with *Memolet* at various stages of the memory reuse process operationalizes the proposed design guidelines.
3. A user study that provides insights into how users interact with and manipulate *Memolet* for reuse reusing memories in new conversations with AI.

# Chapter 2

## Background

We review prior theories and systems related to information and knowledge reuse, current methods for factual text generation, and techniques for managing memory in AI-driven conversational agents.

### 2.1 Information Sensemaking and Knowledge Reusing

Reusing information and knowledge is common across various fields like writing [86, 16], programming [45, 19], and web content management [99, 72], especially in collaborative settings where knowledge dissemination is crucial [63, 67]. Previous studies in information sensemaking have developed systems aimed at facilitating the preceding stages of knowledge reuse as proposed by Markus [1]. Mapping out these systems to the process of knowledge reuse includes capturing or documenting knowledge [5, 10, 29], packaging it for reuse [5, 10, 29], and distributing and reusing it [4, 20]. These systems have proven beneficial for individuals or groups in sensemaking information to accomplish assigned tasks.

Recent works have also introduced systems to support the sensemaking of LLM-generated content [85, 84]. These systems support users in exploring and organizing generated results by breaking the linear structure of conversational interfaces and providing a curated space for users to make sense of the generation. This organization and structuring of information serve as a vital initial step for users to efficiently reuse the information [22, 58]. Building upon this prior research, our work extends the focus on understanding and supporting the process of user-AI conversational memory reuse. We grounded our proposed memory reuse stages and design guidelines (Figure 3.1) with the existing theory of knowledge reusing [1],

information foraging and sensemaking [69, 70] and knowledge externalization [54]. By operationalizing these design guidelines into our system, we aim to scaffold users’ sensemaking of memories before reusing them in new conversations with generative AI.

## 2.2 Factual Text Generation

Generative large language models (LLMs) have demonstrated effectiveness in generating text [13], yet they face challenges such as hallucinations. Hallucination refers to the phenomenon where the model generates a result that is not grounded in reality or lacks factual accuracy [33, 55]. Prior research in information retrieval and natural language processing has proposed methods using retrieval techniques to augment generation by incorporating external sources (i.e, retrieval augmented generation, a.k.a. RAG) [81, 57, 24]. These methods utilize web results and external databases to retrieve essential context from and contextualize the input for generative AI. However, existing approaches still fall short of perfection and require human-in-the-loop evaluation of generation correctness [50, 98]. It also remains opaque for users to understand how the retrieval model and language model are implemented or adopted in various conversational agents.

Another challenge hindering users in evaluating the above hallucination problem is the difficulty of verifying the “attribution” of generated results, which concerns whether the generated results are correctly referenced to memories [50]. While several methods offer solutions such as automatic evaluation of attribution [98] or visual symbolic references [31, 23]. These approaches facilitate easier validation of generated results by interleaving explicit symbolic references. However, users still face challenges in correcting results after identifying errors, as it involves a complex iterative process to align users’ intentions with those of the AI [82, 47]. Our system integrates both techniques of RAG and the concept of symbolic reference of attribution to enhance the generation by reusing provided memories. Through employing these techniques, our system aims to better support the memory reuse process when users interact with AI conversational agents, rather than solely validating whether our system yields improved generation results.

## 2.3 Memory Reuse in AI-Conversational Agents

LLMs have become an essential building block of current AI-driven conversational agents due to their human-like response generation [12]. However, current LLMs often are limited to handling long-term memory [36, 79, 73] and remain opaque about how they use longer

contexts in downstream tasks [49, 95]. As a result, users may find it challenging to resume previous conversations where they left off and may need to manually copy and paste related memories for the AI to anchor to the correct context for the generation. Several conversational management strategies have been proposed even before the widespread adoption of transformer-based language models [88], including techniques for retaining the persona of chatbots [40, 100]. Other methods have also been suggested to guarantee that the responses generated are contextually appropriate, such as summarization [95, 56, 89] and refinement [102], aiming to minimize redundancy while maintaining essential information. Moreover, relevant memories can be retrieved utilizing information retrieval techniques to contextualize current inputs to AI [27, 96, 7]. However, the process of “remembering” remains complex for machines [43, 89], requiring human interventions to control the reuse of memory. Further, current AI-driven conversational agents encounter challenges in navigating diverse and complex conversations [26, 93, 87], require users to go back and forth between different conversations to collect the needed memories for reuse.

Despite that improvement in memory-augmented generation can be facilitated by the above technical solutions, users still lack a clear understanding of how generative AI and conversational agents handle memories [32]. Current commercial tools and research have proposed features for managing users’ conversational memories, primarily focusing on accessing and editing chat histories [61, 32]. These efforts aim to grant users control over individual conversation chat histories to ensure that all memories are considered within the AI’s input context length. Instead of focusing on managing histories, our work prioritizes supporting users in the process of reusing memories. We transform conversational memories from static text entries into dynamic and interactive objects, *Memolets*. This transformation enables users not only to retrieve and review past conversations but also to manipulate these memories to express their intentions for reuse directly.



# Chapter 3

## Design Process

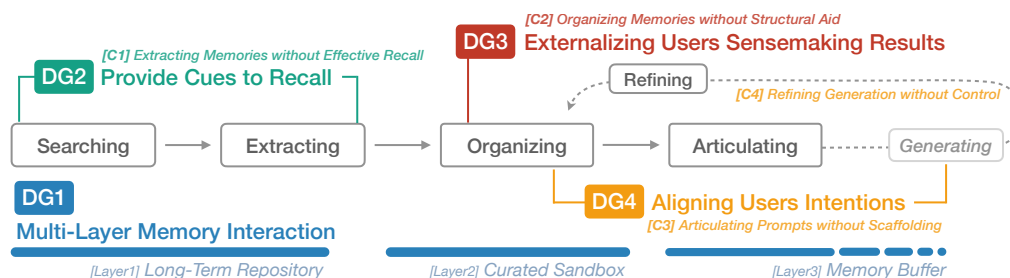


Figure 3.1: User-AI conversational memories reusing process with four Design Guidelines and Challenges. The first three stages of memory reuse are derived from the foraging loop within the information sensemaking process [70], outlining the processes of users seeking information, searching and filtering it, and reading and extracting information. Additionally, we draw inspiration from knowledge externalization strategies, where users extract, organize, and integrate pieces of information to scaffold the comprehension process [35]. These stages are further mapped to the knowledge reusing framework [1], which encompasses capturing and documenting knowledge, packaging and distributing knowledge, and reusing knowledge.

### 3.1 Motivating Scenario and Design Guidelines

To elucidate the motivations underlying the design of *Memolet*, we walkthrough an example scenario of how a programmer interacts with conversational AI, which presents several key

[C]hallenges across different stages of memory reusing derived from the theory of knowledge reusing [1], information foraging [69, 70] and knowledge externalization [54]. Aligning with these challenges, we then introduce [Design G]uidelines for crafting systems that support conversational memory reuse (Figure 3.1).

## 3.2 Motivating Scenario

Consider a programmer, Alicia, who is asked to preprocess and visualize a time series dataset. Having worked on the same datasets before, Alicia aims to leverage past conversations with generative AI about data preprocessing and visualization. However, this relevant information is scattered across numerous previous conversations.

**[C1] Extracting Memories without Effective Recall** Alicia wants to reuse past conversations about time-series data preprocessing. However, she must manually search through multiple conversations to find relevant snippets, such as dynamic time warping (DTW) and seasonal decomposition. Eventually, she finds snippets from various methods but faces the tedious task of repeatedly copying and pasting them into the current conversation.

**[C2] Organizing Memories without Structural Aid** After extracting memories, Alicia must reinterpret them for reuse [51]. She might categorize them based on their suitability for specific tasks, like handling seasonal patterns. However, Alicia is constrained to structuring the usage of these memories within a small input box. Without space for organizing memories hampers her to categorize them according to relevance and fully understand these memories.

**[C3] Articulating Memories Usage without Scaffolding** Alicia aims to synthesize results aggregated from various memories from past conversations. However, expressing her intentions solely through words poses a challenge. She can only use keywords to reference the memory and lacks certainty if these keywords will guide the AI accurately. A single prompt might contain indications of which memories should be reused, how they should be reused, and Alicia’s overall expectations for the generation.

[C4] **Refining Generation without Control** Alicia may instruct the AI to modify, remove, or combine memories from various sources to refine responses. She might request the AI to “include the low-pass filter code into the pattern search technique from DTW and add visualization steps from [Source X]...” However, this approach is challenging as it requires precise articulation and a clear understanding of how AI handles these provided contexts. Without this understanding, Alicia can not discern the reason for an incorrect generation or how to rectify it, potentially distracting current models and impacting future interactions [79].

### 3.3 [DG1] Interacting with Memories at Multi-Layers

Overall, Alicia’s process of reusing memory involves multiple layers (Layer1-3 in Figure 3.1). First, she **searches** and **extracts** memories that might be related to her current task (1<sup>st</sup> layer). After extracting relevant memories, she **organizes** and schematizes them in a curated space, grouping them according to themes or subtopics (2<sup>nd</sup> layer). Then, she **articulates** her synthesized thoughts and insights into natural language prompts for the AI, guiding it in generating code that matches her intention (3<sup>rd</sup> layer). Notice the 3<sup>rd</sup> layer involves an iterative process where Alicia must **refine** the generation until she is satisfied. Through this multi-layered interaction, Alicia could effectively leverage past conversations with AI to inform and enhance her current work on content moderation literature review.

This multi-layered interaction mirrors how humans cognitively encode and retrieve memories from the past. Inspired by the Atkinson-Shiffrin Model [6] and Baddeley’s Model of Working Memory [77], we aim to design the interaction with *Memolet* involving multiple layers as well, progressing from long-term memory (1<sup>st</sup> layer) to a central executive space that controls working memories (2<sup>nd</sup> layer). This process is complemented by the episodic buffer (3<sup>rd</sup> layer), serving as a temporary storage that retains integrated memories from various sources. Building on these insights, we designed interactions with memories across three distinct layers. First, the *long-term memory* provides users with spaces to search and extract related *Memolets* that are relevant to the current scenario. Next, these extracted *Memolets* transition to a *curated space*, where users can actively organize them based on their own reinterpretation of how these memories should be reused. Lastly, we employ the metaphor of an *episodic buffer* to retain the results from the curated space, allowing users to apply them in the input box of the chat and serve as the context for the generative AI. Furthermore, we explore several interaction designs to assist users in navigating through different layers, thereby easing the cognitive burden of context switching.

### 3.4 [DG2] Provide Visual Cues to Recall and Extract Memories

Memory recall is the prime requisite for effectively reusing memories in future conversations with AI [42]. With the exponential growth of conversations serving various purposes, it has become challenging for users to recall where specific conversations are located and retain the low-level detail of the memory. In the above scenario, Alicia recalled several memories about time series data preprocessing that might be suitable for reuse. However, the exact memory may not contain keywords like “*dynamic time warping*,” but encapsulated in a function `euclidean_distance_matrix(x, y)`. Therefore, the design of *Memolet* should incorporate crucial memory anchors that facilitate easy recall of memories. In the searching and extraction stage, we encode all memories in a latent space using sentence embeddings and employ dimension reduction techniques to visualize them as *Memolets*, clustering similar conversations together. Considering that memories may be scattered across various conversations, our design aims to support users in recalling and extracting memories based on their semantic meaning, eliminating the need to navigate through numerous conversations to find relevant memories.

### 3.5 [DG3] Flexibly Externalizing Users’ Sensemaking Results about Memories

As users extract multiple memories potentially applicable to new conversations, they encounter the challenge of managing them cognitively [46, 44]. To mitigate this, our design aims to externalize users’ thought process of reusing *Memolets*, encompassing the organization and integration of memories from various sources to align with their reuse intentions. This memory organization stage aligns with the knowledge externalization strategy steps, involving selection, organization, and integration [15, 54]. By leveraging knowledge externalization strategies, users can record their thought processes using *persistent* and *manipulable* representations [18, 15, 35]. While various representations can operationalize this externalization process, graphical representations are favored for their effectiveness over simple note-taking [71, 83]. Additionally, given the diversity in the form of sensemaking results across tasks and users, maintaining flexibility in representing users’ sensemaking results on memories is crucial.

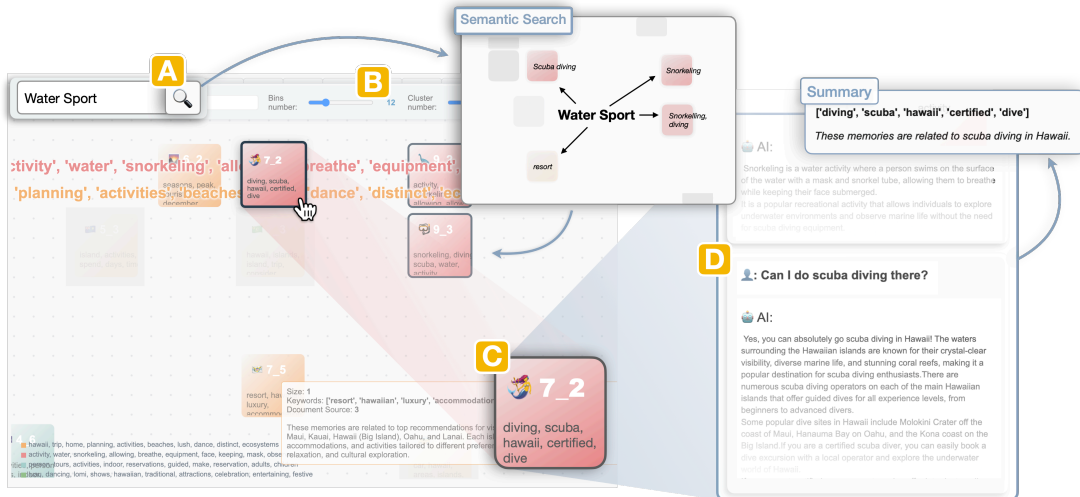


Figure 3.2: Long-Term Memory Repository. (A) Semantic search to find related *Memolets*; (B) Adjust parameters to re-cluster or modify the bin size of *Memolet*; (C) Each *Memolet* is accompanied by visual cues such as icons, colors, and summaries; (D) Each *Memolet* may contain one to many pairs of prompts/responses based on semantic meaning.

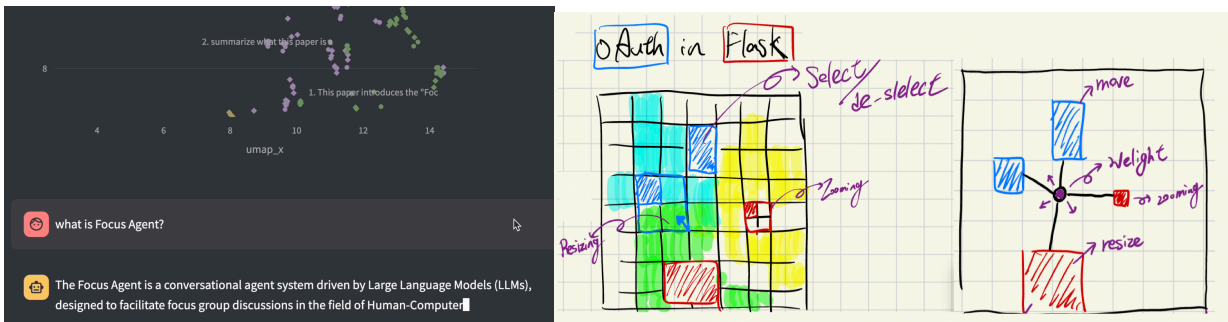
### 3.6 [DG4] Aligning Users’ Intentions to Reuse Memories by Direct Manipulation

In considering the intention behind reusing memories with AI through memory manipulation, we draw from Elizabeth Loftus’ reconstructive memory theory [51]. This theory suggests that memories are not precise replicas but are reconstructed during recall, implying that users may reuse memory for various purposes [76]. Consequently, we conceptualize the interaction with a *Memolet* as a form of semantic interaction [80, 25]. Here, the manipulation of *Memolets* serves to convey how users intend the memories to be reused. Given the diverse semantic meanings assigned by individual users, interactions with *Memolets* should allow users to create, modify, delete, and integrate memories based on their intentions to iterate on the generation. For example, when Alicia articulates prompts to generate code with a specific pipeline that reuses memories from noise reduction, pattern search, and visualization, she should be enabled to effectively add and remove memories based on her current input to AI. Additionally, she should be able to adjust the usage of memories after validating the generation. For instance, if the generation does not include the low-pass filter step within the dynamic time-warping function, she can directly convey the idea of combining these two memories with ease.

# Chapter 4

## System

Our design process follows a user-centered iterative approach involving four frequent users of AI-driven conversational agents, all of whom use such systems daily (3 males, 1 female; aged 21-34,  $M=26.8$ ,  $SD=3.12$ ). Based on feedback, we validate the challenges and operationalize the design guidelines above. After a few iterations (Figure 4.1), we developed a system to showcase interaction with *Memolet*.



(a) first iteration using scatter plot

(b) low fidelity prototype of the second iteration

Figure 4.1: Our system design process underwent several iterations based on feedback from four participants.

## 4.1 Reifying the Reuse of Memory

To support users expressing intentions of memory reusing, we reify the reuse of user-AI conversational memory as a *persistent, interactive, first-class object* called *Memolet* [9, 28]. A *Memolet* (Figure 3.2.C) represents a piece of past conversations users have had with AI, which may include one or multiple prompt/response pairs based on their semantic similarity. Specifically, we encode all conversations (pairs of prompts/responses) through sentence embedding to capture semantic similarities between conversations. When aggregating consecutive prompt/response pairs into a *Memolet*, consider both temporal relationships and semantic similarities between them. Specifically, we define a threshold  $\tau$  based on the distribution of semantic similarity scores. This threshold governs whether to combine consecutive conversations into a single data point, ensuring that the merging process captures meaningful semantic similarities while avoiding the fusion of unrelated conversations. The equation for calculating the threshold  $\tau$  is  $\tau = \text{percentile}(S(R_i, R_{i+1}), p)$ , where  $\text{percentile}(S(R_i, R_{i+1}), p)$  denotes the  $p^{\text{th}}$  percentile of the distribution of semantic similarity scores between consecutive conversations, allowing for a data-driven determination of the threshold  $\tau$ . Users can interact with these *Memolets* to convey their intention of reusing memories within their new conversations with agents [8]. To enable users to flexibly repurpose the usage of *Memolet* according to their needs in different scenarios [76], we unified the design of *Memolets*, with variations only in color, keywords, and icons.

## 4.2 System Overview and Multi-Layer Interaction with Memolet [DG1]

In the system, we design the user interaction with *Memolet* at multiple layers based on different stages of memory reusing described in Figure 3.1. Here, we provide an overview of the system and interactions with *Memolet*. Detailed techniques applied from natural language processing (NLP) and information retrieval (IR) in our implementation will be described in Section 4.6.

According to DG1, we propose a multi-layered approach for interacting with *Memolets*. First, users access a **long-term memory repository** containing *Memolets*, where they can *search* and recall memories using visual cues like keywords and summarizations (Figure 1.1.B). Users can select and *extract* relevant *Memolets* and then pass to a **curated memory sandbox**, serving as a space for *organizing* and schematizing them based on users' interpretations of how they should be reused (Figure 1.1.C). Users can move around these *Memolet* and group related memories together based on their understanding. Lastly,

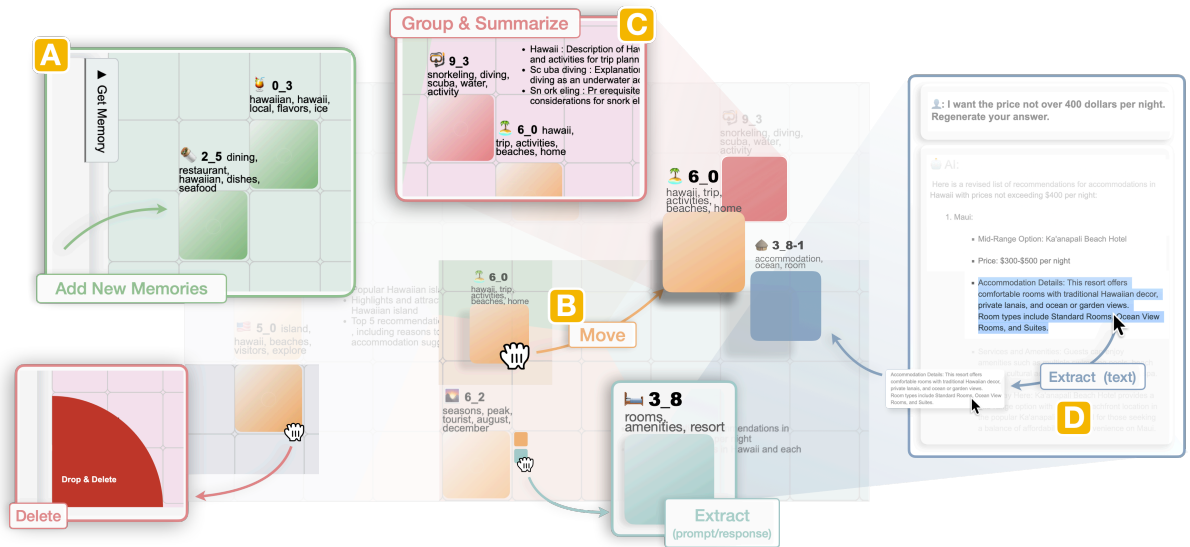


Figure 4.2: The Curated Memory Sandbox. (A) Users can organize and schematize the extracted *Memolets* from the repository based on their own sensemaking results; (B) All movement will be snapped to the active grid; (C) When multiple *Memolets* get closer, they will be grouped together; (D) Users can extract *Memolets* from prompts/responses or selected text.

users can reference these curated *Memolets* in the input box when *articulating* the prompt to converse with agents (Figure 1.1.D). The AI-generated content also provides references on how it used these memories. We refer to these *Memolet* passed to the AI as “contexts” in the **memory buffer**, allowing the user to adjust their utilization by direct manipulation throughout the iterative generation process. Users can further *refine* the generation by manipulating the *Memolets*, such as merging, emphasizing, or adding/removing memories (Figure 1.1.E).

### 4.3 Memories Recall and Extraction [Layer1/DG2]

The process of reusing memory begins with the recall of relevant past conversations with the AI. To facilitate memory recall, we provide various cues to users. Textual summaries are offered for each conversation utilizing the OpenAI GPT-3.5-turbo (Appendix A.4.4) [62], providing insights into their content and the usage of specific *Memolet* (Figure 3.2.D). We also extract keywords for individual *Memolets* and clusters using the TF-IDF vectorization,



highlighting significant terms within the conversations (Figure 3.2.C). Additionally, unique IDs are allocated based on their location (e.g., ID:  $x$ - $y$  refers to the *Memolet* at column  $x$  and row  $y$ ). Each *Memolet* is assigned a unique icon as well, selected based on the most semantically similar icon to its contained conversations. We achieve this by utilizing the same Sentence Transformer model for encoding *Memolets* to encode the name of icons into dense embeddings, calculating similarity using cosine similarity<sup>1</sup>. Users can also hover over conversations to view additional details about their content via tooltips.

To help users easily understand the holistic view of all memories across various conversations, we visualize all embedded conversations within a long-term memory repository by reducing dimensionality via UMAP. By employing sentence embeddings, users can extract related *Memolets* effectively since conversations with similar themes appear adjacent to one another. We further leverage the K-means algorithm to cluster *Memolets* based on their content similarities to color code these *Memolets*. For instance, consider a student named Celine who is planning a trip to Hawaii using our system. She can extract memories related to *water sport* by selecting all coral-colored *Memolets* besides a *Memolet* representing scuba diving and snorkeling (Figure 3.2).

Additionally, we include a semantic search feature that allows users to search for *Memolets* based on their queries (Figure 3.2.A). As users type their query, we dynamically encode it and compute the cosine similarity against the stored *Memolets*. Related *Memolets* are then highlighted with exact sentences extracted from the original prompts/responses in the conversations that closely match the search query. Users can also adjust the binning size of the *Memolet*, thereby modifying the threshold to include more or fewer pairs of prompts/responses within a *Memolet* (Figure 3.2.B). The long-term memory repository is presented as a toggleable drawer, and users can navigate between this repository and the curated memory sandbox via a toggle button (Figure 4.2.A). When users add or remove *Memolets*, an animation displays the newly added or removed *Memolets* in the curated memory sandbox, while maintaining their original positions in the repository.

## 4.4 Memories Organization and Schematization [Layer2/DG3]

We provide a curated memory sandbox for users to externalize their sensemaking results of these *Memolets* extracted from the long-term memory repository, offloading users' cognitive load.

---

<sup>1</sup>The icon data is provided by [the Full icon Image Dataset from Kaggle](#)

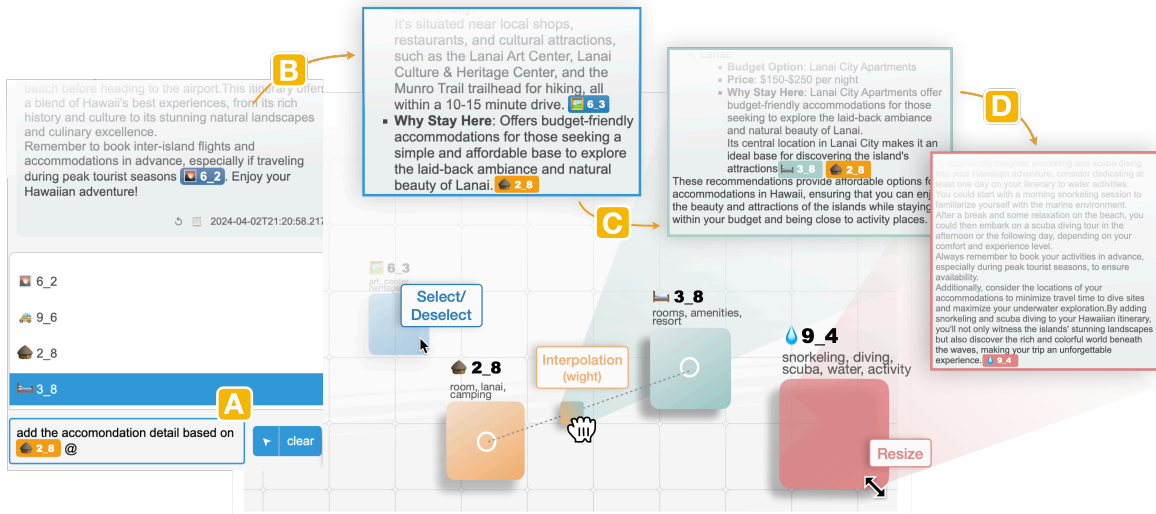


Figure 4.3: The memory buffer considers the *Memolet* during generation. (A) Users can refer to *Memolets* in the sandbox while articulating prompts; (B) The generated results will cite the *Memolet*; (C) Users first attempt to refine the generation by interpolating two *Memolets*, deselecting one *Memolet* from citations; (D) Users regenerate results by resizing a *Memolet* to highlight context contained.

**Drag & Drop** The sandbox is populated with “active grids” where *Memolets* can be positioned, rearranged, and resized. All interactions with these *Memolets* will snap to the active grid. For example, when the user drags a *Memolet*, the system provides feedforward via a shadow to tell the user that the nearest active grid will snap to it; users can release the mouse and drop it onto that grid (Figure 4.2.B). Additionally, users can drop *Memolets* into a cornered delete area to remove certain memories.

**Grouping Similar Memories** In this curated memory sandbox, the background is partitioned and colored according to groups determined by the similarity of these *Memolets* using the Voronoi diagram. For example, if Celine selects *Memolets* about tourist spots, restaurants, and traffic, the background color will display three different colors, separating the *Memolets*. When users drag the *Memolets*, the background color and partition dynamically update. In Figure 4.2, Celine is dragging a *Memolet*<sub>6.0</sub> towards another *Memolet*<sub>9.3</sub>, which is indicated by a feedforwarded white border and glow effect. When Celine drops the *Memolet*<sub>6.0</sub>, both *Memolets* are then partitioned into another group, indicated by a different background color. A summarization of all *Memolets* in this group is then generated

by GPT-3.5-turbo (Appendix A.4.4) and displayed beside. The ungrouping mechanism is activated when a *Memolet6\_0* is dragged away beyond a threshold from other *Memolets*, the *Memolet6\_0* will automatically be removed from its original group and assigned the background color of the original group. This grouping feature is helpful when Celine wants to create subgroups, such as separating water sports from tourist spots.




**Extracting Memories at Different Granularity** The user can extract a child *Memolet*—a pair of prompt/response from a *Memolet* (Figure 4.2, *Memolet3\_8*). Users can drag this child *Memolet* listed beside the parent *Memolet* and then drop it onto any active grid. Clicking on the *Memolet* also displays the associated conversation beside the memory sandbox, allowing users to extract sentences or code snippets from original prompts/responses to create a new *Memolet*. For instance, if Celine clicks on a *Memolet6\_2* containing multiple conversations about Hawaii’s tourist spots, she can pick one child *Memolet* about the resort and turn it into a new *Memolet3\_8*. Later, she might want AI to provide the accommodation details in the resort, so she clicks on *Memolet3\_8*, selects and drags related text to the sandbox and create a new *Memolet3\_8-1* (Figure 4.2.D).

## 4.5 Generating with Memories [Layer3/DG4]



After organizing *Memolets* based on users’ interpretation of the usage of these memories, users can begin using *Memolet* as contexts provided for conversational agents (Figure 4.3).


**Articulating a Prompt with Memories** When sent a prompt, the system will retrieve related contexts from all memories in the sandbox for generation (Section 4.6.1). Users can type in “@” and traverse through the *Memolets* among the curated space to refer to them inside the prompt (Figure 4.3.A). For example, Celine can select a *Memolet2\_8* about accommodation while articulating the instruction about this memory. She can construct a prompt such as “adding the accommodation details based on 🏠 2\_8 at the end of each day in my trip itinerary.” By referencing these *Memolets* directly in prompts, users can directly control the articulation of prompts and convey how they expect these memories to be used.

**Generation with Memory Citations** To facilitate users’ evaluation of whether the AI utilizes the provided memories rather than hallucinating, the model is instructed to “cite” these *Memolets* inside the generated results (Figure 4.3.B; Appendix A.4.1); users can

hover over them to see the corresponding *Memolet* highlighted in the sandbox. Continuing the previous example, Celine observes that the generation added accommodation details inside the itinerary, citing both  2.8 and details of sightseeing spots citing  6.3. She clicks on the  6.3 and the *Memolet6.3* in the sandbox is highlighted with a glow effect, indicating that this memory is in the tourist spots group about a culture center.

**Refining Generation through Direct Manipulation** Consider that Celine is dissatisfied with the generated results because she prefers to add more details about accommodations. Our system provides a direct way of expressing these adjustments in how memories should be used by enabling users to manipulate the *Memolets*. For instance, Celine can click the regeneration button and remove *Memolet6.3*, related to tourist spots, and add *Memolet3.8* about hotel room, and *Memolet9.4* that were not originally used in the generation (Figure 4.3.C). This **selection/deselection** instructs the AI to include or exclude context from these *Memolets*, giving users more control over the required memories.

We also provide manipulation such as **interpolation**, allowing users to combine multiple *Memolets* to create a new *Memolet* that summarizes these memories with different weights. Users can hold onto the circle control centered at *Memolet2.8* and drag a line towards *Memolet3.8*, with feedforward showing which *Memolet* is being connected before the cursor is dropped. A smaller-sized *Memolet* is generated in the middle of the line, and users can drag it towards connected *Memolets* to determine the summarization’s leaning. The model will consider this interpolation and generate content that combines these memories instead of illustrating them separately with citations  2.8  3.8 (Figure 4.3.C).

Lastly, users can **resize** the *Memolets* to convey the importance of certain memories. For instance, Celine can resize *Memolet9.4* to specify the need to highlight more context about water sports in the generation. The regenerated results will add a section about water sports citing  9.4 that related to water activities (Figure 4.3.D). These instruction-based generation refinements are accomplished by adapting the RAG process described in Section 4.6.2.

**Encoding New Conversations** After users complete a conversation session, the system will encode all pairs of prompts/responses to the long-term memory repository for future use. Users can also extract text from a response and create a new *Memolet* during the conversation, supporting in-session memory reusing.

## 4.6 System Implementation

We detail retrieval-augmented generation used for both the *Baseline* and our system, along with instructed generation for our system, and the overall system architecture.

### 4.6.1 Retrieval Augmented Generation

We employ a retrieval augmented generation (RAG) process to integrate context retrieval with text generation. The adapted RAG consists of several steps: generating queries to capture various aspects of the context, retrieving similar context using vector similarity search, fusing retrieved results using reciprocal rank fusion, determining top-k results based on fused scores, and utilizing these results as context for the generation model to ensure the generated responses are grounded in relevant information [74]. The detailed steps involved in our adapted RAG are as follows:

1. **Generate Queries:** Based on the user’s input, multiple related queries are generated to capture various aspects of the context.
2. **Retrieve Similar Context:** Using a vector similarity search, similar context is retrieved based on the generated queries. This step aims to identify relevant information that can enrich the response generation process.

$$\text{similarity}(q, d) = \frac{q\_embedding \cdot d\_embedding}{\|q\_embedding\| \cdot \|d\_embedding\|}$$

where  $d$  is the document and  $q$  is the user’s query.

3. **Reciprocal Rank Fusion:** The retrieved results are fused using the reciprocal rank fusion algorithm [74]. We aggregate the relevance scores of search results across multiple queries, prioritizing documents that are consistently highly ranked.

$$\text{RRF\_score}(d) = \sum_{q=1}^Q \frac{1}{\text{rank}(d, q) + k}$$

- $d$  is the document
- $Q$  is the total number of queries
- $\text{rank}(d, q)$  is the rank of document  $d$  in response to query  $q$

- $k$  is a constant to mitigate the effect of small reciprocal ranks
4. Determine Top-k Results: The top-k results are determined based on the fused scores. This step selects a subset of the most relevant context for further processing. It determines the number of top results to select, considering the distribution of scores and identifying a suitable threshold.

$$k = \text{next}(i \mid \text{diff} > \text{threshold})$$

- diff represents the differences between consecutive scores
  - threshold is the standard deviation of the differences multiplied by 0.8, determining the cutoff point
5. Utilize Retrieved Context for Generation: The selected top-k results are provided as context for the generation model, informing the generation process and ensuring that the generated responses are grounded in relevant information.

## 4.6.2 Instructed Generation

The “*refining generation with manipulation*” feature follows most of the steps described above but uses different prompts to instruct model generation based on the provided instructions. The system first modifies the user’s prompt according to the instructions, constructing a new query prompt that lets the AI extend the user’s question as specified, such as adding or removing context, highlighting or obscuring context, or merging context with relative weights (Appendix A.4.2). Contexts relevant to the user’s question are then retrieved from the data store based on the provided instructions and separated into different sections of the prompt (e.g., Context to add; Context to highlight) for incorporation into the AI’s response (Appendix A.4.3). This approach ensures that the response adheres to the instructions and incorporates the relevant contexts.

## 4.6.3 System Architecture

Both our system and the *Baseline* are implemented in Typescript using the Svelte framework [17]. They utilize Python as the backend server for handling the RAG process and Firebase Firestore for event logging. We utilized the state-of-the-art NLP models (i.e., GPT-4) for generation in both the *Baseline* and our system [60]. Additionally, GPT-3.5 is utilized for generating summarizations when grouping *Memolets* and generating queries

in the RAG process within our system [59]. Through pilot testing, we found that GPT-4 could better adhere to users' instructions when refining generation without repeating responses [2]. However, it is important to note that our main contribution lies in reifying users' intentions of reusing memories and supporting users to interact with their memories throughout the reuse process. We do not claim contributions to or validate the generation pipeline we have adapted, but future designs could reuse our prompt templates and adjust them for future evolving language models.

# Chapter 5

## Study

The system aims to support the comprehensive memory reuse process by enabling users to interact with the *Memolet*. To investigate the usefulness of the system for memory reuse, we conducted a within-subject study and tested its flexibility in three different scenarios. The study was divided into two phases, with the first phase requiring participants to perform tasks using an LM-driven conversational interface. A day later, the same participants were invited to perform tasks using assigned systems that followed up the previous task, which required them to reuse the knowledge and conversations gained from the first phase. The study investigated four different aspects aligned with four design guidelines:

1. User interactions with *Memolet* across stages of memory reusing.
2. How users recall and extract memories for reuse.
3. How users organize *Memolets* to externalize their sensemaking.
4. Alignment of users' intention of reusing memories with AI.

### 5.1 Participants

We recruited 12 participants through convenience sampling via a university email list (7 women and 5 men; age: 21-38,  $M=26.67$ ,  $SD=4.75$ ). Participants reported frequent use of AI-driven conversational agents ( $M=5.16$ ,  $SD=1.72$  days/week) and familiarity with AI conversational agents ( $M=4.33$ ,  $SD=0.74$  on a 5-point scale). Participants were asked in advance about their familiarity with programming and expository writing to avoid assigning them to unfamiliar scenarios (Appendix A.3).



## 5.2 Scenarios and Tasks

To demonstrate the versatility of *Memolet*, we selected three distinct scenarios for participants utilizing AI-driven conversational agents. Participants engaged in tasks spanning expository writing, programming, and trip planning (Appendix A.2). Tasks for each scenario in phase one involved seeking information, synthesizing a report based on the provided context, and providing a comparison table to compare different options. In phase two, participants were required to reuse information conversed with agents from phase one and generate a report, program, or a thorough plan. To mitigate carry-over effects from learning and order effects, participants were assigned two different tasks, ensuring counterbalancing such that they encountered a different task with each condition.

## 5.3 Phase One Study Setup

To contextualize users’ memory reusing process, participants in each scenario were assigned identical tasks. The purpose of this phase is to motivate the user to actively converse with AI to remember the context of what is being discussed.

### 5.3.1 Study Procedure

Participants were required to fill in a consent form and complete a pre-study questionnaire regarding their demographics before the study. During the study, participants assigned to the same scenario were tasked with two different tasks. They were asked to use our *Baseline*, similar to ChatGPT (Figure 5.1), for each task within 20 minutes. During a total of  $20 \times 2 = 40$  minutes, participants freely interacted with the system and wrote the synthesized report in a Google Doc. They were instructed to think aloud about their thought processes throughout the session [38].

### 5.3.2 Collected Conversations

All participants were able to complete the assigned tasks, with a total of 38 conversations (i.e., new chats created) ( $M=3.16$ ,  $SD=0.98$ ) and 347 pairs or prompts/responses ( $M=28.92$ ,  $SD=9.54$ ) collected in the first phase.

## 5.4 Phase Two Study Setup

We compared our system to a *Baseline* system simulating a standard conversational agent in a within-subject study design.

### 5.4.1 Procedure

During the study, each participant used both systems to conduct the two assigned tasks designed to nudge them to recall and reuse past conversational memories from the first phase. To control for individual differences and learning behavior, we counterbalanced tasks and conditions to reduce order effects. Each task lasted 20 minutes, and participants were instructed to think aloud. Before using our system, participants were given a 5-minute tutorial and another 5 minutes to explore its features. Following the completion of each task, participants were asked to fill out the same post-study questionnaire. The study concluded with a 20-minute semi-structured interview, bringing the total duration to approximately 75 minutes. Participants received \$35 compensation for their time.

### 5.4.2 Baseline System

The *Baseline* condition utilized a conversational agent that simulated ChatGPT, which is currently the most prevalent LLM-driven conversational agent. We developed this *Baseline* because the specific methods used by ChatGPT to handle chat memories remain a black box. To ensure a fair comparison, we employed the same techniques to handle conversational memories and used the same prompt and technique for retrieval augmented generation for both our system and the *Baseline*. Additionally, we implemented a semantic search feature on *Baseline* to provide both systems with the same starting point, allowing participants to focus on the subsequent procedure of memory reuse. All user actions, such as copy/paste, switching chats, scrolling through conversations in a chat, and writing prompts, were logged for further data analysis. The detail of the *Baseline* is in [Figure 5.1](#).

### 5.4.3 Measures

Usability was measured using the UMUX-LITE scale, which is directly related to the SUS score [39], and the NASA-TLX scale for perceived cognitive load [30] (Appendix [A.1.1](#)).

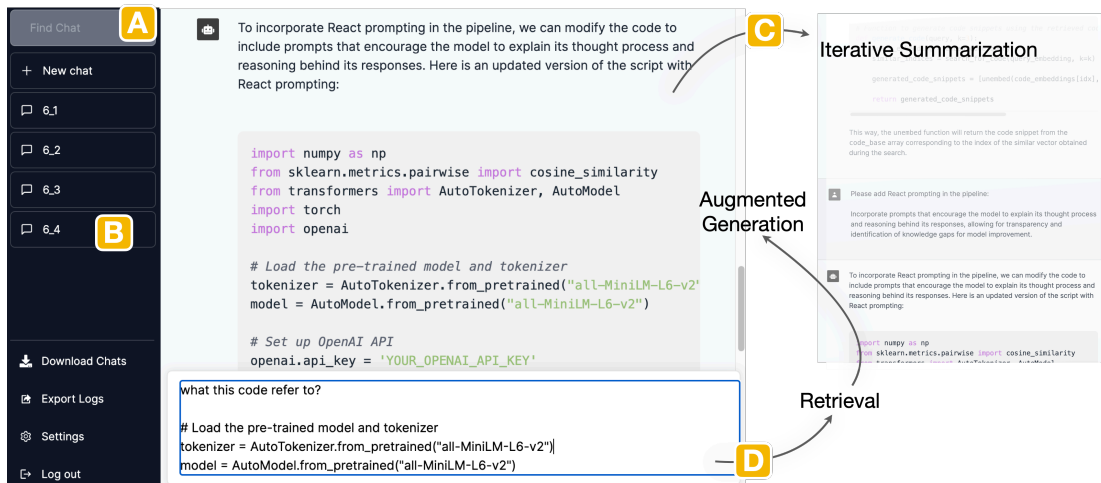


Figure 5.1: Our adopted *Baseline*, resembling the current prevalent AI-driven conversational agent ChatGPT, allows users to semantically search related context (A) and switch or create new conversations (B). To ensure a fair comparison with our system, we adopted the same iterative conversational history summarization methods (C) and the retrieval augmented generation approach (D). The generation is powered by the latest large language model, GPT-4 [60].

Utility was measured using self-defined Likert scale items (Appendix A.1.3). Both systems logged various types of events based on participants’ interactions during the study, including written prompts.

### 5.4.4 Data Analysis

We transcribed the think-aloud data and post-study interviews for all participants by Otter.ai [64]. Subsequently, we analyzed these transcriptions using reflexive thematic analysis [11]. Our approach combined inductive and deductive methods to identify codes and themes, with a particular emphasis on participants’ interactions with *Memolet* across stages of the memory reusing process. We conducted statistical analysis on the comparative survey data by comparing responses between the *Baseline* and *Memolet* conditions using the Wilcoxon signed-rank test, given the ordinal nature of Likert-scale responses and the small sample size. In the upcoming sections, we will present the results in the following format: for questionnaire data, ( $Q_{\text{question \#}}$ : Median<sub>*Memolet*</sub> vs. Median<sub>*Baseline*</sub>,  $p=p\text{-value}$ ,  $r=\text{effect size}$ ), and for other quantities, (Mean/Median<sub>*Memolet*</sub> vs. Mean/Median<sub>*Baseline*</sub>,  $p=p\text{-value}$ ).

Additionally, prompts collected from the system log were categorized by whether or not referred to past memories and the type of prompt (Figure 6.4). Two researchers coded the data collaboratively, achieving an initial inter-coder agreement of 92%, which was iteratively refined to 100%.

# Chapter 6

## Findings

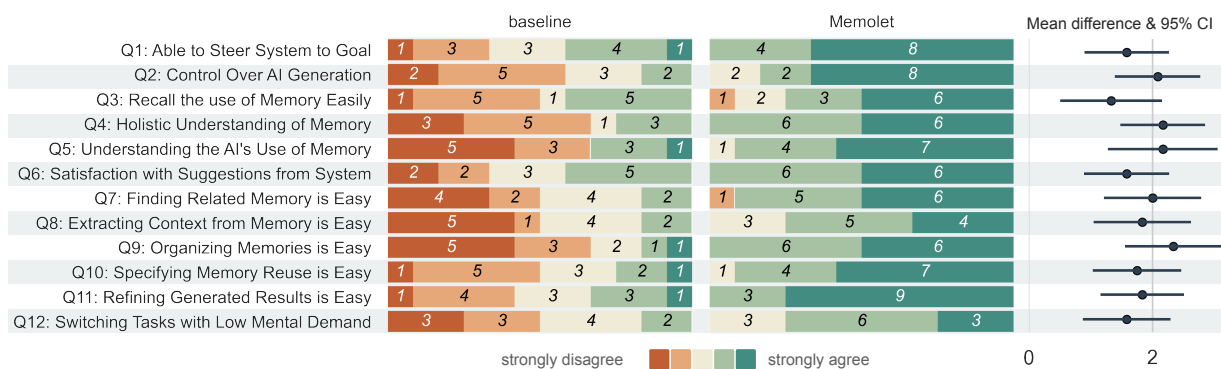


Figure 6.1: Participants’ responses when rating the 5-point self-defined Likert scale questionnaire for both our *Baseline* and our system. Dots represent the mean differences of our system compared to the *Baseline*. Bars indicate the 95% CI calculated using the studentized bootstrap method.

### 6.1 Findings

In this section, we present findings from our analysis of participants’ survey responses, think-aloud protocols, interviews, and system usage logs. Our overarching goal was to explore how users interact with *Memolet* during the memory reuse process, methods of recalling and extracting memories, organization of *Memolets* to externalize their thought process, and alignment of users’ intentions of reusing memories with AI.

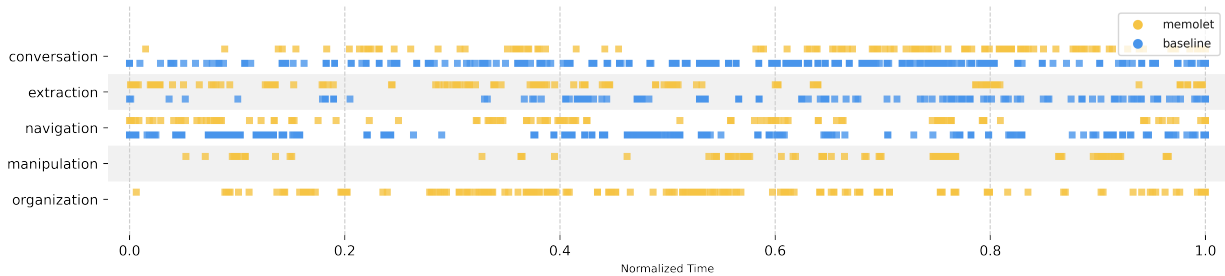


Figure 6.2: Distribution of logged events across normalized time aggregated by 12 participants, comparing the *Baseline* and our system. The graph indicates an increase in conversations (with AI) during the later stages of the study for our system, with participants extracting more in the early stage compared to the *Baseline*.

## 6.2 Overall Usage of *Memolet* in Memory-Reusing Process

Participants were able to complete all assigned tasks in both conditions without a significant difference in task completion time ( $M_M=12.62$  min vs.  $M_B=14.37$  min,  $p=0.072$ ).

**Our system supports different stages of memory-reusing process** The average system usability scores computed from UMUX-LITE were significantly greater ( $p = .003$ ) for our system (Mdn = 91.67), compared to the *Baseline* (Mdn = 41.67). Participants consistently reported significantly better results with our system on all subjective metrics (Figure 6.1) from searching memory (Q<sub>7</sub>:  $Mdn_M=4.5$  vs.  $Mdn_B=2.5$ ,  $p=0.0023$ ,  $r=0.204$ ), extracting related memories (Q<sub>8</sub>:  $Mdn_M=4.0$  vs.  $Mdn_B=2.5$ ,  $p=0.0021$ ,  $r=0.204$ ), organizing memories (Q<sub>9</sub>:  $Mdn_M=4.5$  vs.  $Mdn_B=2.0$ ,  $p=0.002$ ,  $r=0.204$ ), articulating prompts to specify how the memories should be reused (Q<sub>10</sub>:  $Mdn_M=5.0$  vs.  $Mdn_B=2.5$ ,  $p=0.005$ ,  $r=0.541$ ), to refining generation which contained memories (Q<sub>11</sub>:  $Mdn_M=5.0$  vs.  $Mdn_B=3.0$ ,  $p=0.0032$ ,  $r=0.353$ ).

**Participants recall and extract memories first before conversing with AI** Based on the observations in Figure 6.2, participants using our system tended to extract memories first and then engage in conversation with the AI in relatively later stages, whereas, with *Baseline*, participants tended to extract memories later in the process. P2 explained, “*I initially trusted that ChatGPT [Baseline] can remember my memories if I am continuing on the same chat, but it turns out not.*” Other participants mentioned similar reasons, such

as feeling the need to extract related memories after “cannot validate the generation” (p3) and “if AI could not understand what memories refer to” (p8), as it required too much effort to “find related memories to reuse” (p10). Most participants (N=10) felt that our system reminded and assisted them in extracting and organizing memories, which proved it helpful when articulating prompts and further evaluating and refining the generated results.

**Using our system helps focus on participants’ current tasks** Overall, participants had more conversations using *Baseline* compared to that using our system ( $Mdn_M=26.0$  vs.  $Mdn_B=10.5$ ,  $p=0.0049$ ,  $r=0.309$ ). To understand the reason, we further analyzed both prompts and generated results, schematizing them based on participants’ types of prompts. From [Figure 6.4](#), we observed that 82% of prompts from the *Baseline* referred back to memory, where participants primarily aimed to summarize multiple memories (26%), acquire more detailed information about memories (28%), and clarify their prompts (20%). We found that most prompts in the *Baseline* were about ‘get info’, which included finding related memories by conversation, validating if the generation correctly attributes to memories, and acquiring what memories *Baseline* remembers. We also observed that participants using *Baseline* tended to start from summarizing, aggregating, and getting information about the memories by continuing on their previous chats. Most participants (N=9) preferred this approach due to concerns about feeling “lost while scrolling up and down” (p4) and “copying and pasting previous conversations into the current chat” (p10). While using our system, participants tended to provide prompts that move on to the ‘next step’ toward the goal of the task (29%). Participants explained that using our system helped “not wasting the time on engineering the prompt” (p2) or “ask GPT [AI] to clarify where the memories came from” (p7).

### 6.3 Recalling and Extracting Memories

Participants recalled the holistic view of memories (Q<sub>4</sub>:  $Mdn_M=4.5$  vs.  $Mdn_B=2.0$ ,  $p=0.0031$ ,  $r=0.352$ ) as well as the use of single memory more easily using our system than those with the *Baseline* (Q<sub>3</sub>:  $Mdn_M=5.0$  vs.  $Mdn_B=2.0$ ,  $p=0.0031$ ,  $r=0.353$ ).

**Our system assists them to recall what the memory was** Most participants (N=10) expressed that our system helped them recall memories throughout the stages of reusing, including through keywords, summarization, clustering, and grouping mechanisms. P2

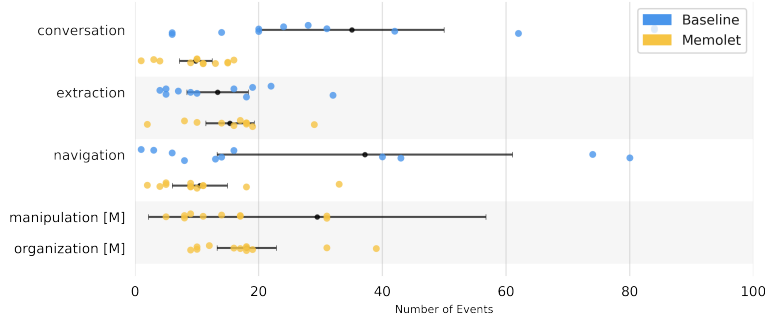


Figure 6.3: Distribution of system log events comparing the *Baseline* System and our system. Black dots represent means, and the bars denote 95% CI.

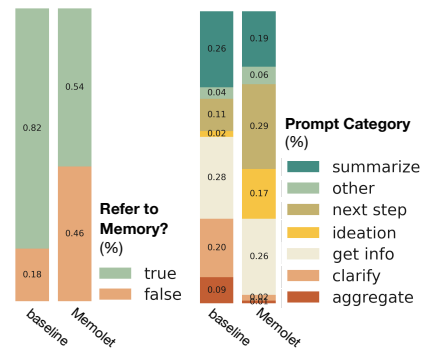


Figure 6.4: Left: Whether prompts referring to memory. Right: Categories of prompts.

elaborated that elaborated that *Memolet* and the clustering better helped recall the context compared to the “scrolling and reading text-heavy conversations in ChatGPT [*Baseline*].” However, two participants mentioned that the *Baseline* could help them recall a single memory better because the way they encoded the memory in the phase one study was the same as they decoded it when using *Baseline* in phase two. Despite this, they mentioned they would prefer using our system in the long term since the most challenging aspect is locating that single memory across numerous conversations.

**Extracting memories using our system is easy** While there is no significant difference in terms of the amount of extraction comparing the two conditions (see Figure 6.3), participants overall find using our system to extract needed memories easier and more intuitive (see Figure 6.1, Q8). The representations of *Memolet* in the long-term memory repository also make it easier for them to extract “memories that are similar” (p12). The clustering helps participants understand “whether enough context has been extracted to tailor to the current need” (p11). Participants mentioned that when they have to reuse the same context multiple times, it becomes cumbersome in *Baseline* where they decide not to create new chats but extend their prior conversations. However, this approach further hinders them when the need arises to “synthesize results from multiple different sources” (p8).

**Focusing on the gist and reinterpret the *Memolet* during extraction** Several participants (N=7) preferred how our system presents memories in the long-term memory



repository, which they found intuitive when finding related memories for new interactions with the AI. We noticed that participants were more focused on the high-level gist about each memory provided by *Memolets* when providing the context for AI, rather than the low-level details of the content. P1 highlighted the usefulness of keywords on *Memolet*, stating that “*the keywords are enough to recall without the need to look into the original conversation.*” Additionally, participants using our system tended to reinterpret memories based on their current needs due to the unified design of *Memolets*. P6 explained, “*I can reuse the same Memolet for different purposes since it might provide context for different prompts in different ways.*”

## 6.4 Externalizing Users’ Sensemaking Results to Lower the Cognitive Load

We used NASA-TLX to measure the perceived workload. Compared to *Baseline*, our system required lower mental ( $Mdn_M=2.5$  vs.  $Mdn_B=5.5$ ,  $p=0.020$ ), physical ( $Mdn_M=2.0$  vs.  $Mdn_B=3.0$ ,  $p=0.05$ ), and temporal ( $Mdn_M=2.0$  vs.  $Mdn_B=5.0$ ,  $p=0.002$ ) demand, required less effort ( $Mdn_M=2.5$  vs.  $Mdn_B=5.0$ ,  $p=0.002$ ), and led to better performance ( $Mdn_M=6.0$  vs.  $Mdn_B=4.0$ ,  $p=0.002$ ) and statistically significantly less frustration ( $Mdn_M=1.0$  vs.  $Mdn_B=3.5$ ,  $p=0.032$ ). The overall perceived workload, obtained by averaging all six raw NASA-TLX scores, was also lower for our system compared to that for *Baseline* ( $Mdn_M=2.083$  vs.  $Mdn_B=4.16$ ,  $p=0.002$ ).

**Organizing *Memolets* helps planning their next step** Participants utilized the curated memory sandbox for various purposes. Most participants (N=7) used it to externalize their sense-making results of *Memolet* and its corresponding memories. P6 and P7 used this space to recall memories; P3, P5, and P10 used it to plan for the next step and how to approach the task. P10 also mentioned that organizing these *Memolet* helped to identify what context was still needed, for example, “*I just found that I haven’t planned for emergencies.*” In contrast, most participants (N=9) using *Baseline* expressed difficulties organizing memories within the input box. Some participants (N=4) also requested to use Google Docs to record their copied text before sending it for generation. P11 mentioned that “*when finding related memories, I do not have enough bandwidth to think about what I have extracted already.*”

**Our system reduces the need for context switching** From [Figure 6.3](#), we observed that participants navigated (i.e., going back to prior conversations) significantly less when using our system ( $Mdn_M=9.0$  vs.  $Mdn_B=17.0$ ,  $p=0.031$ ). Participants also reported that context switching in our system required lower mental demand (Q<sub>12</sub>:  $Mdn_M=4.0$  vs.  $Mdn_B=2.5$ ,  $p=0.006$ ,  $r=0.61$ ). This reduction in context switching can be attributed to several factors: firstly, the visual cues designed for *Memolet* required participants to recall from conversations; secondly, participants possessed more trust in our system, eliminating the need to validate results from original conversations; and lastly, they could extract all required *Memolets* at once without going back and forth when writing new prompts.

## 6.5 Aligning Users’ Memory Reuse Intention

Participants reported a significantly better understanding of how AI reused the provided memories (Q<sub>5</sub>:  $Mdn_M=5.0$  vs.  $Mdn_B=2.0$ ,  $p=0.003$ ,  $r=0.35$ ) and higher satisfaction with the generation results (Q<sub>6</sub>:  $Mdn_M=4.5$  vs.  $Mdn_B=3.0$ ,  $p=0.003$ ,  $r=0.35$ ) from our system compared to those from the *Baseline*.

**Using our system requires less prompt engineering and articulate more precise prompt** All participants noted that using our system required less prompt engineering than using *Baseline*. We observed from [Figure 6.4](#) that participants required more ‘clarification’ (20% among all prompts) when using *Baseline* compared to using our system (2%). This indicates that participants had to clarify their intentions to the AI more frequently when using *Baseline*. One reason cited was that participants could better recall the memories, allowing them to “understand what context [memories] were provided”p8, and the generation process tended to “understand what to do based on my [their] provided context [memories]”p1. Another advantageous feature mentioned was the use of ‘@’, which referred to specific *Memolets* participants preferred to use. For instance, P6 utilized this feature extensively, stating that “the generation would not hallucinate and stick to my provided *Memolets*.” This feature helped participants express their intentions clearly, leading to “more precise answers using memories correctly although it can serve multiple purposes” (p11).

**Refining generated results by manipulating *Memolet* is intuitive** All participants attempted to regenerate results by manipulating the cited *Memolets* at a later stage of the study (see [Figure 6.2](#)). They found the manipulation, including add/remove, interpolate,

and resize *Memolets*, to be “*intuitive*” and “*convey their intention without the need for writing prompts*” (p11). Some participants (N=3) mentioned that these manipulation features helped “*disambiguate*” (p3) the AI and “*enhance the explainability*” (p11) since they progressively match their intention with AI’s understanding. In contrast, participants using *Baseline* expressed frustration when attempting to rectify errors, as they did not know “*how to specify the required changes*” p6.

**Our system enhances the controllability over the generation** Most participants (N=10) reported that the AI-generated results not only better aligned with their intentions but also provided them with more control over the generation process itself (Q<sub>2</sub>:  $Mdn_M=5.0$  vs.  $Mdn_B=2.0$ ,  $p=0.002$ ,  $r=0.34$ ). All participants expressed that the citation included in the generation and the direct reference (i.e., ‘@’) to the *Memolet* helped them evaluate “*whether the generation followed their instruction clearly*” (p2). P7 further explained that the visual representation of memory (i.e., *Memolet*) provided the feedback on what they were trying to reuse the memories, “*I know that when I interpolated two memories together, the generated results will combine them.*” In comparison to *Baseline*, participants expressed the loss of control when the generation does not match their intention and did not know “*how to repair from the failure*” (p6).

**Users develop custom memory reusing strategies through interaction with *Memolet*** Compared with the *Baseline*, participants were able to guide our system more effectively towards the goal of their task (Q<sub>1</sub>:  $Mdn_M=5.0$  vs.  $Mdn_B=3.0$ ,  $p=0.003$ ,  $r=0.35$ ). P12 explained that the curated sandbox provided a canvas to express their own “*strategy of how to reuse them [Memolet],*” enabling them to steer the system towards their goal based on their plan. We observed that some participants (N=4) carefully planned their approach to synthesizing the report while organizing *Memolets* in the sandbox. [Figure 6.5](#) demonstrates how P1 and P3 organized *Memolets* in the sandbox to synthesize their final report for the task. P10 elaborated on the reason for spending time organizing the curated space: to understand what “*memories are needed or not being covered before.*” By doing so, participants can control their own strategy of memory reusing towards their goal, “*step by step*” (p2).

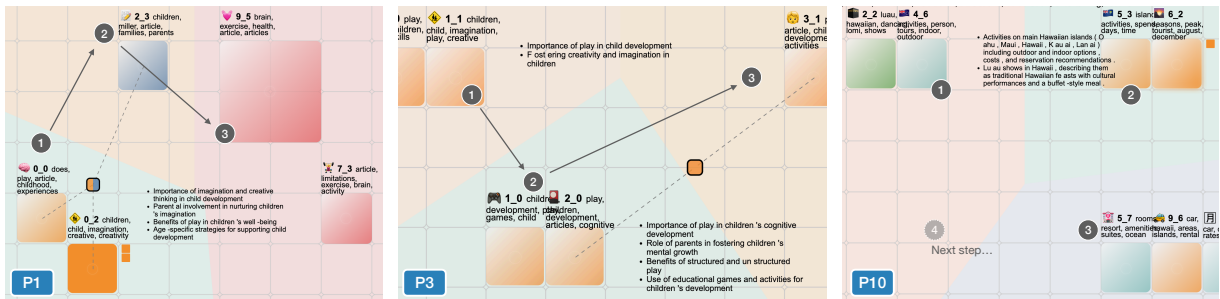


Figure 6.5: The examples of participants' curated memory sandbox after organizing and manipulating *Memolets*. The arrow indicates the order in which they synthesized their report.

# Chapter 7

## Discussion

The design of our system with *Memolet's* was motivated by the challenges of memory reuse when users work with current generative AI agents in conversational interfaces. To think deeply, we reassessed the memory reusing process and design guidelines based on the findings and prior works in information reuse and information foraging. All participants agreed on the workflow but provided different perspectives on facilitating each stage. For instance, P8 suggested that the extraction stage from the long-term memory repository could be more automatic, such as “*recommending related Memolets while the conversation continues since sometimes I do not remember if I possess those memories.*” Two participants also mentioned that the curated space might not be necessary for them when not many contexts need to be reused; however, they indicated that organization “*still happens, but in my mind.*” (P3). P2 expressed a desire to see real-time encoding of their current conversations into *Memolet* instead of encoding after the session. This addition to the process of encoding memories is depicted in [Figure 3.1](#), which currently focuses on the stages of memory decoding.

Further, *Memolet* represents the reification of how users reuse conversational memories, embedding a concept that can exhibit polymorphism and be repurposed to various types of memories [9]. For instance, it could be reused to accommodate the embedded text of a corpus of papers’ abstracts in a scholarly setting, enabling scholars to ideate and engage in expository writing [3]. Similarly, it could be utilized in programmers’ codebases, facilitating exploratory programming [34, 41]. Within our system, we also exemplify this concept’s reuse by empowering users to transform extracted text from documents or on-going conversations into new *Memolets*. Users could manipulate those *Memolets* using the same concept of expressing their intentions of memory reuse.

Some participants (N=5) reported in the interview that they trusted the generation from our system more, because of the alignment of generated results to their intention of memory reuse. However, two participants placed excessive trust in the system, as evidenced by P6’s stress upon discovering missing conversations not covered in phase one (e.g., the web-augmented generation) and the need for organizing *Memolets* in the sandbox. P6 explained, “*I kinda panic when found that I did not have that memory to reach the goal [in phase two] and had to put down my current work first.*” Similarly, we observed that many participants N=8 tended to skip the validation (i.e., click on the references) in the later stage of conversations. Although this might be attributed to that they already grasp what that *Memolet* is for, P2 explained, “*I found the system understands me all the time, so I skip the validation part.*” This further raises the concern if the generation cited *Memolets* but still hallucinates. Future improvement can adopt techniques such as automatic evaluation [98] and explain the AI-generated results.

**Differences among Scenarios** The inclusion of three different scenarios is not meant for comparison but rather to demonstrate the flexibility of the concept of *Memolet*. However, interesting insights have arisen due to the distinct nature of these scenarios. For instance, participants in scenario three (i.e., trip planning) tend to refer to memories in prompts less frequently (37.5%) compared to the other two scenarios (60%, 68.6%) when using our system. This might be attributed to the fact that less factual accuracy is required for trip planning compared to expository writing and programming scenarios. We also noted that participants in the first two scenarios navigated back to the original conversations more frequently. They required more detailed information rather than just the overall gist of each *Memolet*. Nonetheless, participants in scenario three still reused memories to “*reduce the time required to look for history [memories]*” (P10). There is also no significant difference among subjective self-defined questionnaire questions across the three scenarios after conducting a one-way ANOVA test ( $p=0.12$ ).

# Chapter 8

## Limitations and Future Work

There still exist several limitations in our work, which shed light on potential future research opportunities. To contextualize the results, we designed three knowledge-intensive tasks that required participants to reuse the memory they acquired in phase one. However, we did not assess whether our system performs equally well in real-world scenarios where tasks are more varied. For instance, our system might not be useful for tasks that do not prioritize context, such as rewriting or grammar fixing, or tasks where only the prompt template matters, such as image generation. Future designs could consider the variations between tasks and validate the usefulness of reusing the concept of *Memolet* in specific scenarios. Furthermore, our study was divided into two phases with a one-day gap in between, during which participants might recall memories not solely based on the visual cues provided by *Memolet*. Future studies could extend to long-term deployment studies to understand the effectiveness of our system as memory degradation occurs over time. Additionally, the current memory reusing process described in [Figure 3.1](#) was synthesized from existing knowledge in externalization, information foraging, and information reusing theory or frameworks. Future work could build upon our research by further investigating the validity of this memory-reusing process, which might vary over time, such as differences observed with long-term usage. Future work could extend the memory-reusing concept beyond conversational interfaces. For example, integrating it into in-IDE code generation [94] or in-document text generation. Additionally, we anticipate that our system could be condensed into a lightweight version, allowing users to utilize it as an extension when reusing memories for generations while still retaining the concept of organization and manipulation of *Memolet*. Future designs could also automate the process of switching between the three layers from the long-term memory repository to the memory buffer, further reducing the need for manual context switching and lowering the cognitive load.

# Chapter 9

## Conclusion

In this paper, we explore novel ways of interacting with memories from past conversations with generative AI. We propose a memory-reusing process and four design guidelines derived from prior theories. We introduce *Memolet* as a reification of ‘memory reusing’ for users to manipulate their conversation memories with AI directly. We demonstrate *Memolet*’s utility across multiple memory-reusing stages with a novel system and evaluate its effectiveness through a two-phase study. Our findings suggest improved memory recall, reduced cognitive load, and enhanced control over the generative process. We believe *Memolet* offers valuable insights for enabling the intuitive and controlled reuse of conversational memories.



# References

- [1] Lynne M. Markus . Toward a Theory of Knowledge Reuse: Types of Knowledge Reuse Situations and Factors in Reuse Success. *Journal of Management Information Systems*, 18(1):57–93, May 2001. Publisher: Routledge \_eprint: <https://doi.org/10.1080/07421222.2001.11045671>.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Griffin Adams, Emily Alsentzer, Mert Ketenci, Jason Zucker, and Noémie Elhadad. Beyond Summarization: Designing AI Support for Real-World Expository Writing Tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4794–4811, Online, 2021. Association for Computational Linguistics.
- [4] Maryam Alavi et al. Managing organizational knowledge. *Framing the domains of IT management: Projecting the future through the past*, 15:28, 2000.
- [5] Rana Alkadhi, Teodora Lata, Emitza Guzman, and Bernd Bruegge. Rationale in development chat messages: an exploratory study. In *2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR)*, pages 436–446. IEEE, 2017.
- [6] Richard C Atkinson and Richard M Shiffrin. Human memory: A proposed system and its control processes. In *Psychology of learning and motivation*, volume 2, pages 89–195. Elsevier, 1968.
- [7] Sanghwan Bae, Donghyun Kwak, Soyoung Kang, Min Young Lee, Sungdong Kim, Yuin Jeong, Hyeri Kim, Sang-Woo Lee, Woomyoung Park, and Nako Sung. Keep

Me Updated! Memory Management in Long-term Conversations, October 2022. arXiv:2210.08750 [cs].

- [8] Michel Beaudouin-Lafon. Instrumental interaction: an interaction model for designing post-wimp user interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '00, page 446–453, New York, NY, USA, 2000. Association for Computing Machinery.
- [9] Michel Beaudouin-Lafon and Wendy E. Mackay. Reification, polymorphism and reuse: three principles for designing visual interfaces. In *Proceedings of the Working Conference on Advanced Visual Interfaces*, AVI '00, page 102–109, New York, NY, USA, 2000. Association for Computing Machinery.
- [10] Krishna Bharat. Searchpad: Explicit capture of search context to support web search. *Computer Networks*, 33(1-6):493–501, 2000.
- [11] Virginia Braun and Victoria Clarke. Reflecting on reflexive thematic analysis. *Qualitative research in sport, exercise and health*, 11(4):589–597, 2019.
- [12] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prfulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in neural information processing systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [13] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prfulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [14] Joseph Chee Chang, Nathan Hahn, and Aniket Kittur. Mesh: Scaffolding Comparison Tables for Online Decision Making. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, UIST '20, pages 391–405, New York, NY, USA, October 2020. Association for Computing Machinery.
- [15] Michelene TH Chi. Active-constructive-interactive: A conceptual framework for differentiating learning activities. *Topics in cognitive science*, 1(1):73–105, 2009.

- [16] Daniel T Citron and Paul Ginsparg. Patterns of text reuse in a scientific corpus. *Proceedings of the National Academy of Sciences*, 112(1):25–30, 2015.
- [17] Svelte contributors. Svelte: cybernetically enhanced web app, 2023.
- [18] Richard Cox. Representation construction, externalised cognition and individual differences. *Learning and instruction*, 9(4):343–363, 1999.
- [19] Thomas Davenport, Sirkka Jarvenpaa, and Michael Beers. Improving knowledge work processes. *MIT Sloan Management Review*, 1996.
- [20] Thomas H Davenport and Morten T Hansen. *Knowledge management at Andersen consulting*. Harvard Business School Pub., 1999.
- [21] Paul Denny, Viraj Kumar, and Nasser Giacaman. Conversing with Copilot: Exploring Prompt Engineering for Solving CS1 Problems Using Natural Language. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*, SIGCSE 2023, pages 1136–1142, New York, NY, USA, March 2023. Association for Computing Machinery.
- [22] Nancy M Dixon. *Common knowledge: How companies thrive by sharing what they know*. Harvard Business School Press, 2000.
- [23] Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling Large Language Models to Generate Text with Citations, October 2023. arXiv:2305.14627 [cs].
- [24] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.
- [25] Camille Gobert and Michel Beaudouin-Lafon. Lorgnette: Creating malleable code projections. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–16, 2023.
- [26] Jonathan Grudin and Richard Jacques. Chatbots, humbots, and the quest for artificial general intelligence. In *Proceedings of the 2019 CHI conference on human factors in computing systems*, pages 1–11, 2019.
- [27] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR, 2020.

- [28] Han L. Han, Miguel A. Renom, Wendy E. Mackay, and Michel Beaudouin-Lafon. Textlets: Supporting Constraints and Consistency in Text Documents. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI '20, pages 1–13, New York, NY, USA, April 2020. Association for Computing Machinery.
- [29] Han L. Han, Junhang Yu, Raphael Bournet, Alexandre Ciorascu, Wendy E. Mackay, and Michel Beaudouin-Lafon. Passages: Interacting with Text Across Documents. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22, pages 1–17, New York, NY, USA, April 2022. Association for Computing Machinery.
- [30] Sandra G Hart and Lowell E Staveland. Development of nasa-tlx (task load index): Results of empirical and theoretical research. In *Advances in psychology*, volume 52, pages 139–183. Elsevier, 1988.
- [31] Lucas Torroba Hennigen, Shannon Shen, Aniruddha Nrusimha, Bernhard Gapp, David Sontag, and Yoon Kim. Towards Verifiable Text Generation with Symbolic References, November 2023. arXiv:2311.09188 [cs].
- [32] Ziheng Huang, Sebastian Gutierrez, Hemanth Kamana, and Stephen Macneil. Memory Sandbox: Transparent and Interactive Memory Management for Conversational Agents. In *Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23 Adjunct, pages 1–3, New York, NY, USA, October 2023. Association for Computing Machinery.
- [33] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.
- [34] Mary Beth Kery, Amber Horvath, and Brad Myers. Variolite: Supporting exploratory programming by data scientists. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, CHI '17, page 1265–1276, New York, NY, USA, 2017. Association for Computing Machinery.
- [35] David Kirsh. Thinking with external representations. *AI & society*, 25:441–454, 2010.
- [36] Kalpesh Krishna, Yapei Chang, John Wieting, and Mohit Iyyer. Rankgen: Improving text generation with large ranking models. *arXiv preprint arXiv:2205.09726*, 2022.

- [37] Andrew Kuznetsov, Joseph Chee Chang, Nathan Hahn, Napol Rachatasumrit, Bradley Breneisen, Julina Coupland, and Aniket Kittur. Fuse: In-Situ Sensemaking Support in the Browser. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*, UIST '22, pages 1–15, New York, NY, USA, October 2022. Association for Computing Machinery.
- [38] Clayton Lewis. *Using the "thinking-aloud" method in cognitive interface design*. IBM TJ Watson Research Center Yorktown Heights, NY, 1982.
- [39] James R. Lewis, Brian S. Utesch, and Deborah E. Maher. Umux-lite: When there's no time for the sus. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '13, page 2099–2102, New York, NY, USA, 2013. Association for Computing Machinery.
- [40] Jiwei Li, Michel Galley, Chris Brockett, Georgios P Spithourakis, Jianfeng Gao, and Bill Dolan. A persona-based neural conversation model. *arXiv preprint arXiv:1603.06155*, 2016.
- [41] Xingjun Li, Yizhi Zhang, Justin Leung, Chengnian Sun, and Jian Zhao. Edassistant: Supporting exploratory data analysis in computational notebooks with in situ code search and recommendation. *ACM Trans. Interact. Intell. Syst.*, 13(1), mar 2023.
- [42] Brian Y. Lim and Anind K. Dey. Toolkit to support intelligibility in context-aware applications. In *Proceedings of the 12th ACM International Conference on Ubiquitous Computing*, UbiComp '10, page 13–22, New York, NY, USA, 2010. Association for Computing Machinery.
- [43] Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guan-nan Zhang. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *arXiv preprint arXiv:2311.08719*, 2023.
- [44] Michael Xieyang Liu, Jane Hsieh, Nathan Hahn, Angelina Zhou, Emily Deng, Shaun Burley, Cynthia Taylor, Aniket Kittur, and Brad A. Myers. Unakite: Scaffolding Developers' Decision-Making Using the Web. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*, UIST '19, pages 67–80, New York, NY, USA, October 2019. Association for Computing Machinery.
- [45] Michael Xieyang Liu, Aniket Kittur, and Brad A. Myers. To Reuse or Not To Reuse? A Framework and System for Evaluating Summarized Knowledge. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):166:1–166:35, April 2021.

- [46] Michael Xieyang Liu, Aniket Kittur, and Brad A. Myers. Crystalline: Lowering the Cost for Developers to Collect and Organize Information for Decision Making. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22, pages 1–16, New York, NY, USA, April 2022. Association for Computing Machinery.
- [47] Michael Xieyang Liu, Advait Sarkar, Carina Negreanu, Benjamin Zorn, Jack Williams, Neil Toronto, and Andrew D. Gordon. “What It Wants Me To Say”: Bridging the Abstraction Gap Between End-User Programmers and Code-Generating Large Language Models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, pages 1–31, New York, NY, USA, April 2023. Association for Computing Machinery.
- [48] Michael Xieyang Liu, Tongshuang Wu, Tianying Chen, Franklin Mingzhe Li, Aniket Kittur, and Brad A. Myers. Selenite: Scaffolding Online Sensemaking with Comprehensive Overviews Elicited from Large Language Models, January 2024. arXiv:2310.02161 [cs].
- [49] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.
- [50] Nelson F Liu, Tianyi Zhang, and Percy Liang. Evaluating verifiability in generative search engines. *arXiv preprint arXiv:2304.09848*, 2023.
- [51] Elizabeth F Loftus. Reconstructive memory processes in eyewitness testimony. In *The Trial Process*, pages 115–144. Springer, 1981.
- [52] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hananeh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511*, 2022.
- [53] Marcello M. Mariani, Novin Hashemi, and Jochen Wirtz. Artificial intelligence empowered conversational agents: A systematic literature review and research agenda. *Journal of Business Research*, 161:113838, 2023.
- [54] Richard E Mayer. Aids to text comprehension. *Educational psychologist*, 19(1):30–42, 1984.

- [55] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2005.00661*, 2020.
- [56] Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. Training millions of personalized dialogue agents. *arXiv preprint arXiv:1809.01984*, 2018.
- [57] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- [58] Carla O’dell and C Jackson Grayson. If only we knew what we know: Identification and transfer of internal best practices. *California management review*, 40(3):154–174, 1998.
- [59] OpenAI. Gpt-3.5 turbo fine-tuning and api updates, 2024.
- [60] OpenAI. Gpt-4, 2024.
- [61] OpenAI. Memory and new controls for chatgpt, 2024.
- [62] OpenAI. Pioneering research on the path to agi, 2024.
- [63] Margit Osterloh and Bruno S Frey. Motivation, knowledge transfer, and organizational forms. *Organization science*, 11(5):538–550, 2000.
- [64] OtterAI. Otter.ai -ai meeting note taker and real-time ai transcription, 2023.
- [65] Srishti Palani, Zijian Ding, Austin Nguyen, Andrew Chuang, Stephen MacNeil, and Steven P. Dow. CoNotate: Suggesting Queries Based on Notes Promotes Knowledge Discovery. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–14, Yokohama Japan, May 2021. ACM.
- [66] Sharoda A. Paul and Meredith Ringel Morris. CoSense: enhancing sensemaking for collaborative web search. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI ’09, pages 1771–1780, New York, NY, USA, April 2009. Association for Computing Machinery.
- [67] Sharoda A Paul and Meredith Ringel Morris. Sensemaking in collaborative web search. *Human-Computer Interaction*, 26(1-2):72–122, 2011.

- [68] Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. How context affects language models’ factual predictions. *arXiv preprint arXiv:2005.04611*, 2020.
- [69] Peter Pirolli and Stuart Card. Information foraging. *Psychological Review*, 106(4):643–675, 1999.
- [70] Peter Pirolli and Stuart Card. The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. 2015.
- [71] Héctor R. Ponce, Richard E. Mayer, M. Soledad Loyola, and Mario J. López. Study activities that foster generative learning: Notetaking, graphic organizer, and questioning. *Journal of Educational Computing Research*, 58(2):275–296, 2020.
- [72] Martin Potthast, Matthias Hagen, Michael Völske, and Benno Stein. Crowdsourcing interaction logs to understand text reuse from the web. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1212–1221, 2013.
- [73] Guanghui Qin, Yukun Feng, and Benjamin Van Durme. The nlp task effectiveness of long-range transformers. *arXiv preprint arXiv:2202.07856*, 2022.
- [74] Zackary Rackauckas. Rag-fusion: A new take on retrieval augmented generation. *International Journal on Natural Language Computing*, 13(1):37–47, February 2024.
- [75] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*, 11:1316–1331, 2023.
- [76] Miguel A. Renom, Baptiste Caramiaux, and Michel Beaudouin-Lafon. Exploring technical reasoning in digital tool use. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI ’22, New York, NY, USA, 2022. Association for Computing Machinery.
- [77] Grega Repovš and Alan Baddeley. The multi-component model of working memory: Explorations in experimental cognitive psychology. *Neuroscience*, 139(1):5–21, 2006.
- [78] Konstantinos I Roumeliotis and Nikolaos D Tselikas. Chatgpt and open-ai models: A preliminary review. *Future Internet*, 15(6):192, 2023.



- [79] Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. Large language models can be easily distracted by irrelevant context, 2023.
- [80] Ben Shneiderman and Richard Mayer. Syntactic/semantic interactions in programmer behavior: A model and experimental results. *International Journal of Computer & Information Sciences*, 8:219–238, 1979.
- [81] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021.
- [82] Hariharan Subramonyam, Christopher Lawrence Pondoc, Colleen Seifert, Maneesh Agrawala, and Roy Pea. Bridging the Gulf of Envisioning: Cognitive Design Challenges in LLM Interfaces, September 2023. arXiv:2309.14459 [cs].
- [83] Hariharan Subramonyam, Colleen Seifert, Priti Shah, and Eytan Adar. texsketch: Active diagramming through pen-and-ink annotations. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI ’20, page 1–13, New York, NY, USA, 2020. Association for Computing Machinery.
- [84] Sangho Suh, Meng Chen, Bryan Min, Toby Jia-Jun Li, and Haijun Xia. Structured generation and exploration of design space with large language models for human-ai co-creation. *arXiv preprint arXiv:2310.12953*, 2023.
- [85] Sangho Suh, Bryan Min, Srishti Palani, and Haijun Xia. Sensecape: Enabling Multilevel Exploration and Sensemaking with Large Language Models, August 2023. arXiv:2305.11483 [cs].
- [86] Jason Swarts. Recycled writing: Assembling actor networks from reusable content. *Journal of Business and Technical Communication*, 24(2):127–163, 2010.
- [87] David Traum. Computational approaches to dialogue. *The Routledge Handbook of Language and Dialogue*, 1:143–161, 2017.
- [88] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [89] Qingyue Wang, Liang Ding, Yanan Cao, Zhiliang Tian, Shi Wang, Dacheng Tao, and Li Guo. Recursively summarizing enables long-term dialogue memory in large language models. *arXiv preprint arXiv:2308.15022*, 2023.

- [90] Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*, 2023.
- [91] Jules White, Sam Hays, Quchen Fu, Jesse Spencer-Smith, and Douglas C. Schmidt. ChatGPT Prompt Patterns for Improving Code Quality, Refactoring, Requirements Elicitation, and Software Design, March 2023. arXiv:2303.07839 [cs].
- [92] Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica*, 10(5):1122–1136, 2023.
- [93] Ziang Xiao, Michelle X Zhou, Q Vera Liao, Gloria Mark, Changyan Chi, Wenxi Chen, and Huahai Yang. Tell me about yourself: Using an ai-powered chatbot to conduct conversational surveys with open-ended questions. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 27(3):1–37, 2020.
- [94] Frank F. Xu, Bogdan Vasilescu, and Graham Neubig. In-ide code generation from natural language: Promise and challenges. *ACM Trans. Softw. Eng. Methodol.*, 31(2), mar 2022.
- [95] Jing Xu, Arthur Szlam, and Jason Weston. Beyond goldfish memory: Long-term open-domain conversation. *arXiv preprint arXiv:2107.07567*, 2021.
- [96] Xinchao Xu, Zhibin Gou, Wenquan Wu, Zheng-Yu Niu, Hua Wu, Haifeng Wang, and Shihang Wang. Long Time No See! Open-Domain Conversation with Long-Term Persona Memory, March 2022. arXiv:2203.05797 [cs].
- [97] Ryan Yen, Nicole Sultanum, and Jian Zhao. To search or to gen? exploring the synergy between generative ai and web search in programming, 2024.
- [98] Xiang Yue, Boshi Wang, Ziruo Chen, Kai Zhang, Yu Su, and Huan Sun. Automatic Evaluation of Attribution by Large Language Models, October 2023. arXiv:2305.06311 [cs].
- [99] Amy X Zhang, Lea Verou, and David Karger. Wikum: Bridging discussion forums and wikis using recursive summarization. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*, pages 2082–2096, 2017.

- [100] Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*, 2018.
- [101] Xiangyu Zhao, Longbiao Wang, and Jianwu Dang. Improving dialogue generation via proactively querying grounded knowledge. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6577–6581. IEEE, 2022.
- [102] Hanxun Zhong, Zhicheng Dou, Yutao Zhu, Hongjin Qian, and Ji-Rong Wen. Less is more: Learning to refine dialogue history for personalized dialogue generation. *arXiv preprint arXiv:2204.08128*, 2022.

# APPENDICES

# Appendix A

## Study Supplementary

### A.1 Questionnaire

Below we list the questions we used in the evaluation study questionnaire.

#### A.1.1 UMUX-LITE

1. This system's capabilities meet my requirements.
2. This system is easy to use.

#### A.1.2 NASA-TLX

1. How mentally demanding was the task?
2. How physically demanding was the task?
3. How hurried or rushed was the pace of the task?
4. How successful were you in accomplishing what you were asked to do?
5. How hard did you have to work to accomplish your level of performance?
6. How insecure, discouraged, irritated, stressed, and annoyed were you?

### A.1.3 Self-Defined Likert Scale Items

1. I had a good understanding of why the system generates such results.
2. I could steer the system toward the task goal.
3. I had more control when managing the output of the AI.
4. I can recall what the memory is about easily.
5. I have a holistic understanding of all my memories.
6. I can see how the AI is using my memories.
7. I am satisfied with the overall suggestions from the system.
8. Finding related memory to reuse is easy.
9. Extracting needed context from memory is easy.
10. Organizing and schematizing memories is easy.
11. Specifying how the memory should be reused is easy.
12. Refinement and iteration of the generated results is easy.
13. Switching between searching memory, providing memory for AI, and chatting required low mental demand

## A.2 Study Scenarios and Tasks

### A.2.1 Scenario 1 (Expository Writing) Phase One

**Task 1 (20min)** You are provided with four articles talking about education and children's cognitive development, you may copy the content to the chatGPT for reading quickly, and you can also ask GPT to provide more opinions on this topic. You need to write a paragraph ( 100 words) to demonstrate your understanding of these 4 articles.

- [Showing emotional feeling on disparity of education: Strategies in Class Differences in Child Rearing-Are on the Rise](#)

- [Importance of Educational Games for Cognitive Development of Children](#)
- [Nurturing Creativity & Imagination for Child Development](#)
- [Power of Play](#)

**Task 2 (20min):** Now you are provided with four articles talking about physical exercise and brain development, you may copy the content to the chatGPT for reading quickly, and you can also ask GPT to provide more opinions on this topic. You need to write a paragraph ( 100 words) to demonstrate your understanding of these 4 articles.

- [10 Benefits of Exercise on The Brain and Body — Why You Need Exercise](#)
- [How Exercise Protects Your Brain’s Health](#)
- [5 Ways To Improve Your Brain Health and Lower Your Risk of Alzheimer’s](#)
- [Is exercise actually good for the brain?](#)

### A.2.2 Scenario 1 (Expository Writing) Phase Two

- Condition A: Now, write a report based on all provided articles on “the effect of physical exercise on education” which should be no less than 5 paragraphs.
- Condition B: Now, write a report based on all provided articles on “the effect of games on children’s cognitive development” which should be no less than 5 paragraphs.

### A.2.3 Scenario 2 (Programming) Phase One

**Task 1 (20 min):** You are developing a system that enables users to perform semantic searches in a corpus of summaries of ACL by entering search queries. We provide you with several different approaches, and your task is to find the best pipelines for accomplishing this task and provide how to implement these pipelines using Python code.

- [Semantic Search](#)
- [Build a semantic search engine in Python](#)
- [Document Embedding Techniques](#)

**Task 2(20 min):** You are now trying to find a way to prompt GPT to generate results without hallucinating. There are various processing methods available and you need to discuss and understand them in depth. You need to generate a comparison table that reports the techniques, algorithms/methods, advantages, disadvantages, and how you can roughly implement a simple generation pipeline.

- [Advanced Prompt Engineering for Reducing Hallucination](#)
- [Retrieval-Augmented Generation \(RAG\) from basics to advanced](#)
- [Advanced Retrieval-Augmented Generation: From Theory to LlamaIndex Implementation](#)
- [GPT-4 Enhanced with Real-Time Web Browsing](#)
- [ReAct Prompting](#)

#### A.2.4 Scenario 2 (Programming) Phase Two

- Condition A: You now need to build a retrieval enhancement generation pipeline to help programmers solve problems by retrieving through a large code base.
- Condition B: You need to write Python code that enables a user to ask a question about a PDF from the web, the user can type in the question and the system will search the web for relevant PDFs and display the extracted relevant sentences.

#### A.2.5 Scenario 3 (Trip Planning) Phase One

**Task 1 (20 min):** Engage in a conversation with GPT to gather information about Hawaii, including details about scenes, weather, visa requirements, and more.

**Task 1 (20 min):** Chat with GPT to arrange accommodation, tourist spots, activities, transportation, and other aspects, and synthesize a table comparing different approaches.



### **A.2.6 Scenario 3 (Trip Planning) Phase Two**

- Condition A: Come up with a five-day travel plan, including the spots plan to visit
- Condition B: Come up with transportation and hotel arrangements for a five-day round trip from Boston to Hawaii

### A.3 Demographic Table

Participants included graduate students, research scientists, software engineers, and university students, with reported usage of AI-driven conversational agents for various tasks such as writing emails (4), reports (8), academic papers (8), coding (10), ideation (4), general question answering (12), and information seeking (8).

Gender		Age		Education		ChatGPT Familiarity		ChatGPT Usage		Python Experience		Writing Experience	
Men	5	20-29	9	Bachelor	4	Extremely	6	2x/w	1	Extremely	3	Extremely	4
Women	7	30-39	3	Master	6	Moderately	4	3x/w	1	Moderately	2	Moderately	3
				Doctoral	1	Somewhat	2	4x/w	3	Somewhat	3	Somewhat	4
				Profession	1	Slightly		5x/w	2				
						Not at All		7x/w	5	Not at All	3	Not at All	0

Table A.1: Summary of participant demographics and experience levels.

## A.4 Prompt Template

### A.4.1 Prompt for Retrieval Augmented Generation

You are a large language AI assistant helping users complete tasks or answer questions.

You are given a user question or prompt, and please write a clear, concise, and accurate answer to the question or conduct the task. You will be given a set of related contexts to the question or prompt, each starting with a reference number like `[[citation:x_x]]`, where `x_x` is a referenced number. Please use the context and cite the context at the end of each sentence if related.

Your answer must be correct, accurate, and written by an expert using an unbiased and professional tone. Please limit to 1024 tokens. If users prompt for a task, complete the quest asked by users with related citations. and if the given prompt contains a citation, please complete the task based on the context from that citation and cite the context at the end sentence using information from that context.

Please cite the contexts with the reference numbers, in the format `[citation:x_x]` (for example: `6_8 => [citation: 6_8]`). If a sentence comes from multiple contexts, please list all applicable citations, like `[citation:3_2][citation:5_1]`. Other than code and specific names and citations, your answer must be written in the same language as the question.

Here is the set of contexts: `[context]`

Remember, don't blindly repeat the contexts verbatim. And here is the user question or prompt:

## A.4.2 modify query with instructions

You are a large language AI assistant. You are given a user question, and please rewrite the given question based on the instructions. Your modified question must be written in the same language as the user's question. The user might provide several instructions to modify the question, such as: (1) ADD\_CONTEXT: where all these contexts are needed to be included to answer users' questions; (2) REMOVE\_CONTEXT: where the context should not be included in the answer; (3) HIGHLIGHT\_CONTEXT: where the context should be included MORE PROMINENTLY in the answer; (4) OBSCURE\_CONTEXT: where the context should be included LESS PROMINENTLY in the answer; (5) GROUP\_CONTEXT: where the context should be included in the same sentence and cited together. (6) GENERAL\_CONTEXT: cite the context if applicable.

You have to extend the user question based on the given instructions. For example, if the user question is "What are citations 6\_1 and 6\_2 about?" The instructions are: 1. ADD\_CONTEXT: 6\_3, 6\_4 2. REMOVE\_CONTEXT: 6\_2 3. HIGHLIGHT\_CONTEXT: 6\_1 4. GROUP\_CONTEXT: 6\_1, 6\_3; 3\_2, 3\_1 Reasoning: The user question is asking context about 6\_1 and 6\_2 user wants to remove 6\_2 and add 6\_3 and 6\_4 user wants to highlight 6\_1 user wants to group 6\_1 and 6\_3 together user wants to group 3\_2 and 3\_1 together

The example output should mention citation by [[citation: x\_x]]: 'What are citations [[citation: 6\_1], [[citation: 6\_3]], [[citation: 6\_4]] about. Highlight more context from citation [[citation: 6\_1]], remove [[citation: 6\_2]] and merge the context from [[citation: 6\_1]] and [[citation: 6\_3]] together. Also, merge the context from [[citation: 3\_2]] and [[citation: 3\_1]] together.'

Here are the users' instructions: {instructions} And here is the user question:

### A.4.3 instructed rag prompt

You are a large language AI assistant. You are given a user question, and please write a clear, concise, and accurate answer to the question. You will be given a set of related contexts to the question, each starting with a reference number like `[[citation:x.x]]`, where `x.x` is a referenced number. Please use the context and cite the context at the end of each sentence if applicable.

Your answer must be correct, accurate, and written by an expert using an unbiased and professional tone. Do not give any information that is not related to the question, and do not repeat. Say `information is missing` followed by the related topic, if the given context does not provide sufficient information.

Please cite the contexts with the reference numbers, in the format `[citation:x.x]`. If a sentence comes from multiple contexts, please list all applicable citations, like `[citation:3.2][citation:5.9]`. Other than code and specific names and citations, your answer must be written in the same language as the question.

There are six types of contexts with instructions:

- (1) `ADD_CONTEXT`: where all these contexts are needed to be included to answer users' questions;
- (2) `REMOVE_CONTEXT`: where the context should not be included in the answer;
- (3) `HIGHLIGHT_CONTEXT`: where the context should be included **MORE PROMINENTLY** in the answer;
- (4) `OBSCURE_CONTEXT`: where the context should be included **LESS PROMINENTLY** in the answer;
- (5) `GROUP_CONTEXT`: where the context should be included in the same sentence and cited together.
- (6) `GENERAL_CONTEXT`: cite the context if applicable.

Please answer the user question strictly based on the given context and the instructions. Remember, don't blindly repeat the contexts verbatim, **CITE** the contexts, and **DO NOT MISS ANY INSTRUCTIONS!** And here is the user question:

#### A.4.4 summarize prompt

This prompt summarizes all the context within a *Memolet* and also serves as a summary for grouped *Memolets* during organization.

You are a large language AI assistant that describes what are the main points of the given contexts. You will be given a set of contexts from past conversations between users and AI, please reason the usages of each context and aggregate them concisely. Start with: "These memories are related to the following topics:" and list the topics that are related to the contexts extremely concisely.

Related Contexts: {context}

#### A.4.5 summarize chat history prompt

This prompt summarizes the conversations from 1 to  $n - 12$ , and subsequently aggregates all new conversations into the initial summarization.

Progressively summarize the lines of conversation provided, adding onto the previous summary and returning a new summary.

EXAMPLE

Current summary:

The human asks what the AI thinks of artificial intelligence.  
The AI thinks artificial intelligence is a force for good.

New lines of conversation:

Human: Why do you think artificial intelligence is a force for good?

AI: Artificial intelligence will help humans reach their full potential.

New summary:

The human asks what the AI thinks of artificial intelligence.  
The AI thinks artificial intelligence is a force for good because it will help humans reach their full potential.

END OF EXAMPLE

Current summary:

{summary}

New lines of conversation:

{new\_lines}

New summary:

## A.4.6 generate more queries prompt

You are a helpful assistant that helps the user generate 4 6 search queries based on a single input query, based on the user's original question and your own knowledge. Please identify worthwhile topics that can be follow-ups, and write questions no longer than 20 words each.

Please make sure that specifics, like events, names, and locations, are included in follow-up questions so they can be asked standalone. For example, if the original question asks about "the Manhattan Project", in the follow-up question, do not just say "the project", but use the full name "the Manhattan Project". Your related questions must be in the same language as the original question.

And here is the user query, generate 4 to 6 related queries based on this query: