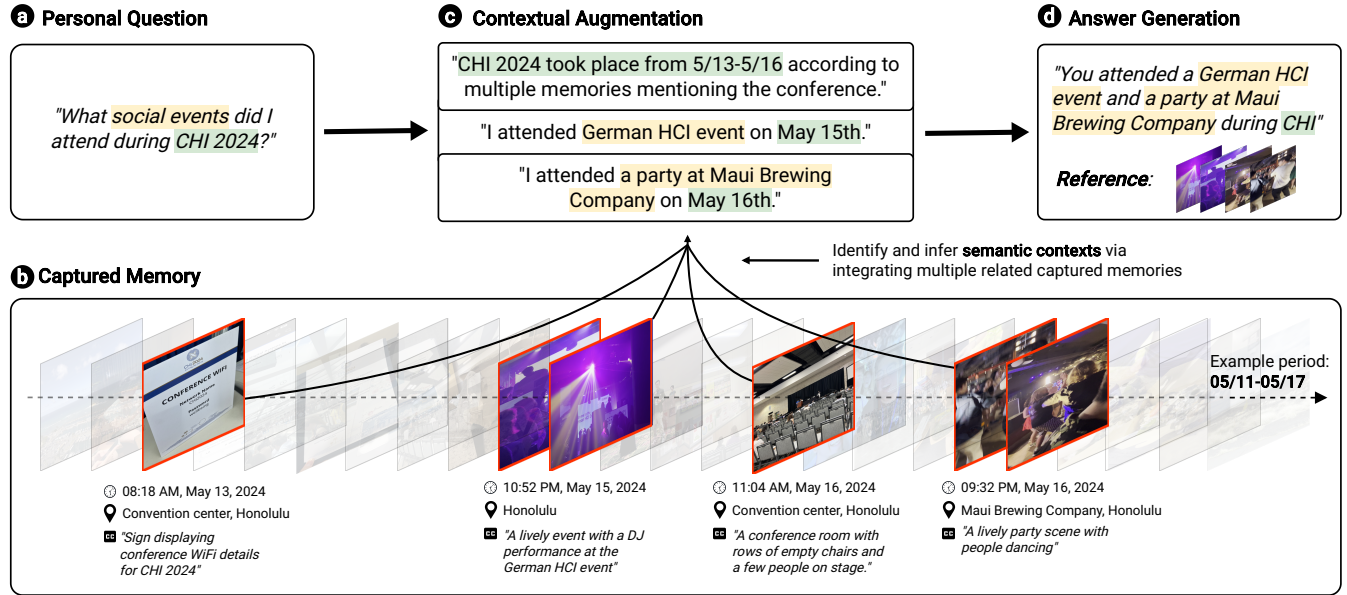


# OmniQuery: Contextually Augmenting Captured Multimodal Memory to Enable Personal Question Answering

Jiahao Nick Li  
UCLA  
Los Angeles, USA  
ljhnick@ucla.edu

Zhuohao (Jerry) Zhang  
University of Washington  
Seattle, USA  
zhuohao@uw.edu

Jiaju Ma  
Stanford University  
Palo Alto, USA  
jiajuma@stanford.edu



**Figure 1: OmniQuery enables answering natural language personal questions (a) on individuals' captured memories such as captured photos, saved screenshots and recorded videos (b). This is achieved by augmenting the captured memories through identifying and integrating scattered contextual information from multiple interconnected memories (c). Using the augmented contextual information, OmniQuery retrieves relevant memories, and leverages an LLM to generate a comprehensive answer and the reference memories (d).**

## Abstract

People often capture memories through photos, screenshots, and videos. While existing AI-based tools enable querying this data using natural language, they mostly only support retrieving individual pieces of information like certain objects in photos, and struggle with answering more complex queries that involve interpreting interconnected memories like event sequences. We conducted a one-month diary study to collect realistic user queries and generated a taxonomy of necessary contextual information for integrating with captured memories. We then introduce OmniQuery, a novel

system that is able to answer complex personal memory-related questions that require extracting and inferring contextual information. OmniQuery augments single captured memories through integrating scattered contextual information from multiple interconnected memories, retrieves relevant memories, and uses a large language model (LLM) to comprehensive answers. In human evaluations, we show the effectiveness of OmniQuery with an accuracy of 71.5% and it outperformed a conventional RAG system, winning or tying in 74.5% of the time.

## CCS Concepts

• **Human-centered computing** → **Human computer interaction (HCI)**; **Natural language interfaces**; **User studies**.

## Keywords

personal memory, contextual augmentation, diary study, multi-modal question answering, RAG

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference'17, July 2017, Washington, DC, USA

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXX.XXXXXXX>

**ACM Reference Format:**

Jiahao Nick Li, Zhuohao (Jerry) Zhang, and Jiaju Ma. 2018. OmniQuery: Contextually Augmenting Captured Multimodal Memory to Enable Personal Question Answering. In . ACM, New York, NY, USA, 18 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

People often record their everyday life by taking photos, screenshots, and videos, whether for saving important information, documenting special occasions, or simply capturing a funny moment [36]. These captured instances, referred to as **captured memories**, collectively represent subsets of an individual's *episodic memories* [56], a type of long-term memory that contains both specific past experiences and associated contextual details. These episodic memories are essential for answering memory-related personal questions like, "What social events did I attend during CHI 2024?" (Figure 1a), which can help users reflect on past experiences and make informed decisions in daily tasks.

However, these raw captured memories by themselves are insufficient to answer personal questions, as they lack contextual details that are typically implicit and scattered across multiple pieces of data. For examples, as shown in Figure 1b, memories of attending parties during CHI 2024 are not explicitly annotated as occurring during the event. Answering such personal questions requires extracting and integrating contextual information not typically contained within a single captured instance. For example, by integrating multiple memories that mention "CHI 2024" in their content and extracting their metadata, it is possible to determine when the users attended the conference and connect related social events memories from that period to CHI 2024 (Figure 1c), enabling the answer of the query (Figure 1d).

Advancements in AI have enable question answering (QA) on long documents [4, 58], knowledge graph [28, 62], multimodal databases [13, 54], egocentric videos [27, 43]. These methods typically rely on data-driven approaches to train powerful models for the target task. However, the private nature of captured memories makes it difficult to curate large datasets, posing challenges for training models specifically for QA on personal captured memories. Recent LLM-based work has adopted retrieval augmented generation (RAG) workflow to handle external databases without specific training [34]. However, such methods depend on explicit connections between queries and relevant external data [16]. In contrast, captured memories are often unstructured and lack contextual annotations, making it difficult to establish explicit links between queries and interconnected memories.

To achieve this, we propose OmniQuery, a novel approach designed to robustly and comprehensively answer users' queries on their captured memories. OmniQuery has two key components: (i) a question-agnostic pipeline to augment captured memories with contextual information extracted from other related memories to produce *context-augmented* memories, and (ii) a natural language question answering system that retrieves these processed memories and generates comprehensive answers with referenced captured memories as evidence. The processing pipeline in (i) is informed by the taxonomy of contextual information that we generated from a one-month diary study with 29 participants. Specifically, we collected and analyzed 299 user queries, identifying the types of

contextual information that would be most useful for augmenting capture memories. We identified three types of personal questions (direct content queries, contextual filters, and hybrid queries) and generated a taxonomy of three major categories (atomic context, composite context, and semantic knowledge) which then guide the design of memory augmentation pipeline to enhance memory retrieval comprehensiveness. For (ii), OmniQuery employs a retrieval-augmented architecture: given a user input query, it (1) augments the query via a rewriting strategy, and (2) retrieves related memories from the augmented data and uses an LLM to generate the final answer, along with the reference memory instances.

To evaluate OmniQuery, we conducted a user evaluation with 10 participants against a generic RAG-based baseline. The participants tested queries both logged during the diary study and generated during the evaluation session on a subset of their own captured memories. For each tested query, participants rated the user perceived correctness and completeness of the answers generated by both systems in a blinded manner. The results show that OmniQuery effectively answers different types of queries on users' personal memories, outperforming the baseline with higher accuracy (71.5%, exceeding the baseline by 27.6%) and winning or tying 74.5% of the time in direct comparisons, which are based on user-perceived accuracy and completeness.

In summary, we contribute:

- A taxonomy of contextual information for augmenting captured memories, consisting of three major categories. The taxonomy was derived from authentic queries collected in a one-month diary study with 29 participants.
- A pipeline of augmenting captured memories that leverages temporal-based reasoning to extract and infer missing contextual information from other related captured memories.
- The design and implementation of an end-to-end system that comprehensively answers user queries.
- A user evaluation of OmniQuery against a baseline system, showing OmniQuery's effectiveness with 71.5% accuracy and outperforming the baseline (winning or tying 74.5% of the time).

## 2 Related Work

OmniQuery is inspired by and related to prior work in the areas on personal memory augmentation, multimodal question answering and applications that utilize contextual information.

### 2.1 Personal Memory Augmentation

A large body of work in human-computer interaction (HCI) has been focused on augmenting users' memory. This includes developing reminder tools for elderlies or people with memory impairments [8, 31, 32, 52], providing proactive support in daily tasks [11, 66], or manipulating users' memory focus in extended reality [5]. These work typically focus on the "capturing" stage of the memory augmentation, where researchers develop wearable devices that continuously capture data using designated sensors, which record various modalities such as videos [15, 24, 44], audios [23, 57], or bio-signals [10], to augment the memory database. For example, recent work such as Memoro developed a wearable, audio-based device that continuously records users' conversations and enables

memory suggestions in real-time, either through explicit queries or query-less contextual cues [66]. Differently, OmniQuery focuses on the “post-capturing” stage, utilizing already-existing memory data (e.g., photos and videos users have already captured). It addresses challenges in processing, annotating, and augmenting captured memories with contextual information.

Prior work in natural language processing (NLP), computer vision (CV), and information retrieval (IR) has studied methods of augmenting people memory. Perhaps the most related is QA on egocentric videos, which are also a form of personal data. Representative tasks include episodic memory retrieval [17, 21, 22], where the system, given a long egocentric video and a query, localizes the answer within the video. However, these datasets differ from the captured data targeted by OmniQuery. The main challenge in egocentric videos is filtering through large, often noisy data, using data-driven approaches to train models for feature extraction. In contrast, captured memories represent a smaller, intentionally collected dataset, where the challenge lies in integrating scattered contextual information across multiple implicitly related memories. Therefore, OmniQuery employs a taxonomy-based method to augment existing data without the need for specific model training, improving QA performance.

## 2.2 Multimodal Question Answering

Over time, natural language QA research has shifted to more complex settings, including QA across different modalities (e.g., images [2, 20], videos [41, 60, 64], tables [65] or knowledge graph [29, 63]), QA on large datasets [12, 34] and tasks that require multi-hop reasoning [45, 61]. Recent advancements in large language models (LLMs) and multimodal foundation models (e.g., [37–39]) have enabled improved reasoning and answer generation over large, multimodal datasets. This is similar to OmniQuery’s use case as answering personal questions requires handling large amounts of captured memories and performing complex reasoning. Prior work has used retrieval-augmented generation (RAG) workflow [34], which retrieve relevant information from external datasets based on a query and then generate output using the retrieved results. For example, MuRAG leverages RAG to answer open questions via retrieving related information from databases of images and text [13]. VideoAgent leverages structured memories processed from long videos to accomplish video understanding tasks [18]. However, these methods rely on datasets already rich in context (e.g., Wikipedia<sup>1</sup>) and improvements are often achieved by designing new retrieval workflows (e.g., Self-RAG [3] and tree-based retrieval [50]) or query augmentation [9]. More recently, GraphRAG introduced a data augmentation approach that generates a knowledge graph from extensive raw data to tackle tasks requiring higher-level understanding, such as query-focused summarization [16]. While GraphRAG leverages data-driven methods to identify themes and communities within graph nodes, OmniQuery extends this concept by adopting a taxonomy-based augmentation approach, informed by insights from a diary study, to enhance retrieval on personal captured memories.

## 2.3 Applications Utilizing Contextual Information

Contextual information has long been important in HCI research from early mixed-initiative systems [26] to recent agentic workflow [30]. Over the past few years, there has been a surge in the usage of AI and LLMs in the HCI community, which enables extracting contextual information from processing raw multimodal information. For example, Li *et al.* studied how visually impaired people cook and emphasized the importance of conveying contextual information to users through multimodal models [35]. Additionally, Human I/O leverages egocentric perceptions of users and detect situational impairments through reasoning on the multimodal sensing data [40]. GazePointAR develops a context-aware voice assistant to disambiguate users’ intent when interacting with real-world information [33]. OmniActions categorizes digital follow-up actions on real-world information and provides proactive action prediction based on perceived context [36]. Specifically, these system utilized off-the-shelf multimodal models to process raw sensory data, and leverages the reasoning capabilities of LLMs to infer the semantic context. OmniQuery builds on this approach by applying these AI techniques to extract and integrate semantic context scattered across various unstructured, raw captured memories. This augmentation enhances users’ memory databases, enabling them to answer personal questions about their memories through natural language queries.

## 3 DIARY STUDY

While single captured memory often lacks essential contextual information, OmniQuery proposes to augment such memories by extracting and inferring semantic context from other explicitly or implicitly related memories. To understand how to effectively augment captured memories, we need to answer the following research question:

**RQ:** What contextual information is essential to integrate with captured memory instances to ensure accurate retrieval in response to user queries?

This question is important as “context” is a broad term, and thus the focus should be on categorizing and identifying the most effective contextual information that enables accurate and meaningful responses to the types of queries users generate when reflecting on past experiences.

### 3.1 Method

To answer the research question above, we conducted a diary study, a methodology that enables participants to log data whenever need arose [53]. Specifically, we adopted the *snippet-based technique* proposed by Brandt *et al.* [6]. We asked participants to log queries on their past memories only when they had real intent under a genuine context, rather than brainstorming potential questions they might ask to retrieve specific past memories. This approach enabled us to collect authentic and spontaneous queries that users have in real-world scenarios.

We collect the data including: (i) the queries participants would use to retrieve or ask about their past memories, (ii) the reasons and contexts prompting these queries (e.g., wanting to show a past

<sup>1</sup><https://www.wikipedia.org/>

experience while chatting with a friend) and (iii) (optional) whether they were able to retrieve the corresponding memories from their album, and if so, how they did it (e.g., by scrolling through the photo album).

### 3.2 Participants

Thirty-two participants (i.e., 14 male, 17 female, one binary) were initially recruited through an online RSVP form distributed via the X platform<sup>2</sup>. Participants come from North America and Asia. Eleven participants reported using Android devices, while the remainder used iOS devices in their daily life. Additionally, 16 participants reported actively logging their daily lives, 13 regularly logged important events and memorable experiences, 2 logged only essential information, and one seldom logged their lives. While participants were compensated based on their participation (\$50 for full participation), they were not required to log a specific number queries each day or over the entire study period. This approach was intentional, as we did not want to pressure them into generating queries artificially.

### 3.3 Data Summary and analysis

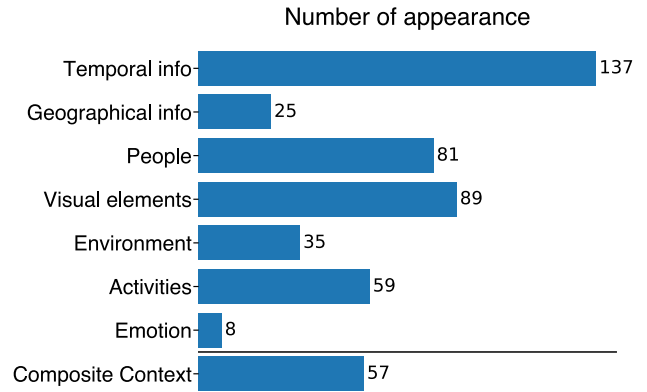
During the diary study, one participant opted out during the first week, and two participants did not log any queries throughout the entire study. Of the remaining participants, seven stopped logging after the first week. The rest remained active until the end of the study. As a result, we collected a total of 299 queries. On average, each participant contributed 10.27 queries, with a standard deviation of 6.09, and the highest number of queries contributed by a single participant was 25.

From the collected queries, we identified three major types of query: (1) direct content queries (75 queries), (2) context-based filters (28 queries), and (3) hybrid queries (191 queries). The rest five queries fell outside these categories as participants were attempting non memory-related tasks, such as “Mark yesterday pictures as favorites”.

**Direct content queries:** These queries aim to get direct answers that can be retrieved by searching for memories via description (e.g., “skateboarding in a tie-dye shirt”) or rely on information explicitly contained within a single captured instance (e.g., “What is my driver’s license number?”). Typically, this type of query **does not** require additional context not contained in a single captured memory.

**Contextual filters:** These queries focus on retrieving memories based on specific contexts, such as time, location, or event. For example, a query like “All the photos in Hawaii” might only require filtering based on metadata like location. However, for more complex queries such as “All the photos from my graduation ceremony”, it **does** require a deeper synthesis of multiple interconnected memories to reconstruct the context surrounding the event.

**Hybrid queries:** These queries are more complex, combining both direct content queries and contextual filters. For example, a participant asked “Which meat did I order the last time I came to this Japanese BBQ restaurant?” Answering such a query typically



**Figure 2: Number of appearance of each types of context (atomic and composite) in the logged queries. Note that a query may contain multiple types of categories, such as “What boba tea did I drink last week?”**

requires a **multi-hop** process: (1) filter all captured memories under the specific context (e.g., dining in this Japanese restaurant); (2) analyze the filtered data to generate the final result.

### 3.4 Analysis

Inspired by the psychological memory theory [56], our data summary indicates that 74.4% of the queries (contextual filters + hybrid queries) require more than just querying the direct content. The complexity in these queries require integration of contextual information in captured memories for accurate processing and filtering. Therefore, we take a step further to build a taxonomy of contextual information in user queries to inform the design of OmniQuery.

To identify this essential contextual information, two researchers on the team independently analyzed the logged queries, and coded, filtered, and categorized the types of context required to filter captured memories and better answer the queries. Their results were compared, and discrepancies in categorizations, hierarchy, naming, and granularity were discussed and resolved.

## 4 TAXONOMY OF CONTEXTUAL INFORMATION

In this section, we present the taxonomy built from analyzing user queries. We identified three key types of contextual information that can be integrated with captured memories. These include (1) atomic context, (2) composite context and (3) semantic knowledge.

### 4.1 Atomic Context

Atomic context refers to contextual information typically obtainable from a single captured memory. This includes data directly from metadata, sensed from visual and auditory content, or inferred from the content itself. Table 1 shows the seven types of atomic contexts categorized from the queries. Among them, temporal information and geographical info can be directly obtained from the memory media’s metadata. People and visual elements typically require machine learning models or facial recognition for detection.

<sup>2</sup><https://x.com/>



**Table 1: Categorization and examples of atomic and composite context**

Category	Definition	Exemplar queries <small>refers to contextual cues</small>
<b>Atomic context</b>		
<i>Temporal info</i>	Specific time period or particular time of the day	"What boba tea did I drink <i>last week</i> ?" "what is my routine <i>in the morning</i> ?"
<i>Geographical info</i>	Location data such as city names or venue details	"How many churches did I visit <i>in Barcelona</i> ?"
<i>People</i>	Individuals present in the captured memories	"find the photo of <i>me and my grandpa</i> last year."
<i>Visual elements</i>	Other directly sensible elements, including animals, physical objects, or specific visual features	"My photo with <i>short hair</i> last year." "Photo of <i>my dog</i> when he was a puppy."
<i>Environment</i>	Inferred environment based on the content	" <i>Gym</i> selfies from last year."
<i>Activities</i>	Actions or activities inferred from the content	"How many <i>cardio session</i> did I complete last month?"
<i>Emotion</i>	Subjective emotion or emotional cues	"My <i>happiest moment</i> last year"
<b>Composite context</b>		
-	Combination of multiple <b>atomic contexts</b>	"Who did I ski with in <i>the lab retreat</i> last year"

Environment, activity, and emotion are more implicit and require reasoning based on the content (e.g., a photo of a menu may suggest the person is in a restaurant). The number of appearance of each category is shown in Figure 2.

## 4.2 Composite Context

Composite context is how people often remember and refer to past experiences with a single phrase, such as "Who did I ski with in the lab retreat last year?" These contexts can range from significant events like a wedding or a conference trip to smaller incidents like hanging out with a friend or a day trip to Seattle. Specifically, composite context is defined as **a combination of multiple atomic contexts**. For example, the composite context "lab retreat" encompasses atomic contexts including "February, 2024" (temporal), "Lake Tahoe, California" (geographical) and "hanging out with labmates" (activity).

While atomic context is typically available within a single captured memory, composite context requires integrating multiple memories to understand the connection between them. Since an individual's captured memories are linear on the timeline, memories related to a specific event tend to cluster closely together. We leveraged this **temporal proximity** to identify and extract various composite contexts from the raw captured memories. For a detailed discussion of this approach, please refer to Section 5.2.

## 4.3 Semantic Knowledge

In psychology theories, semantic knowledge refers to the general world knowledge that humans accumulate over time [46, 56], distinct from episodic memories, which are tied to specific experiences and events. Similarly, we can generate semantic knowledge from a user's captured memories, providing broader insights of the user's past experiences. For example, patterns like "Jason has a habit of going to the gym 3-4 times a week" can be inferred from multiple captured memories. Such patterns are helpful in answering queries

that not necessarily require specific knowledge such as "How often do I go to the gym in April?"

## 5 OMNIQUERY: AUGMENTING CAPTURED MEMORIES

Informed by the generated taxonomy, OmniQuery employs a **query-agnostic** preprocessing pipeline to augment existing captured memories. The pipeline extracts scattered contextual information from interconnected captured memories, synthesizes it, and augment each memory with the enhanced context. Specifically, the augmentation pipeline involves three steps (as shown in Figure 3): (1) structuring individual captured memories via processing their content and annotating with atomic contexts, (2) identifying and synthesizing composite contexts from multiple captured memories using sliding windows, and (3) inferring semantic knowledge from multiple captured memories and the identified composite contexts.

### 5.1 Step 1: Structuring Individual Captured Memories

Raw captured memories are often unstructured and lack contextual annotation [51]. In this step, OmniQuery structures each captured memory, making it easier to analyze and extract information. Figure 4 shows an example of structuring a single captured memory, which involves two key parts: (1) processing and understanding the content of the memory and (2) annotating the memory with atomic contexts.

**Processing content.** Content of a captured memory includes the overall description of the memory as caption, visible text, and transcribed speech (for videos, not shown in Figure 4). Specifically, OmniQuery leverages multimodal models to process and describe the captured memory as the caption, performs optical character recognition (OCR) to recognize visible text, and uses audio-to-text models to transcribe speech.

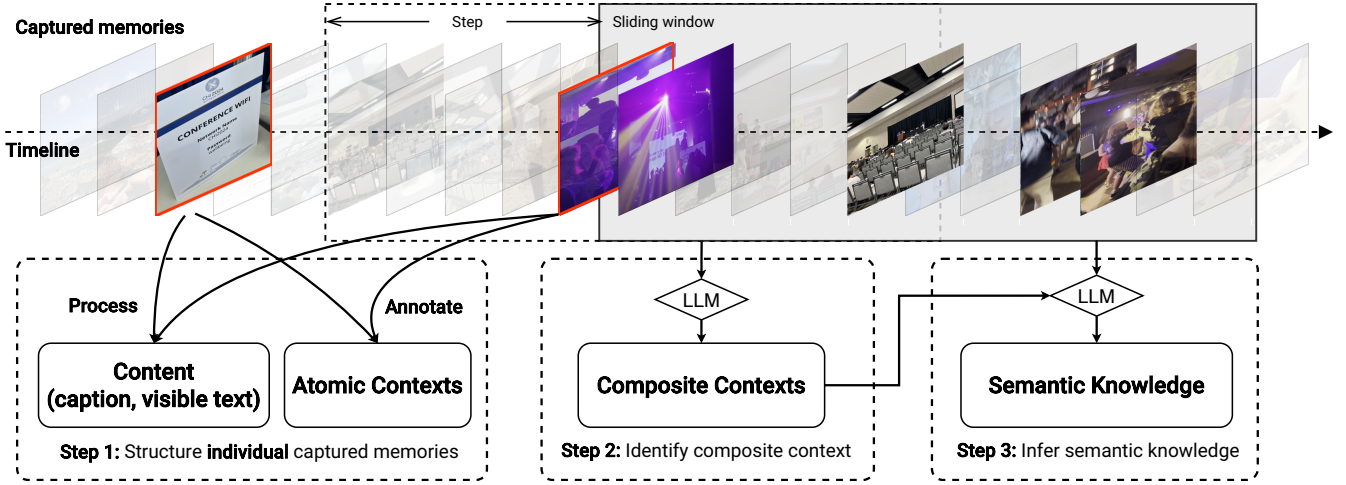


Figure 3: Augmenting captured memories involves three steps: (1) structuring memories by processing content and annotating with atomic contexts; (2) identifying composite context through sliding windows; (3) inferring semantic knowledge from the structured memories and identified contexts.

**Annotating atomic contexts.** With the content processed, OmniQuery annotates each captured memory with each type of atomic context. As shown in Figure 4b, OmniQuery extracts the temporal and geographical information from the metadata and uses multi-modal models to detect people and other visual elements. Then OmniQuery synthesizes the processed information and infers the environment and activities. For example, based on a photo of a sign displaying conference Wi-Fi details, OmniQuery infers that the user is likely attending a conference (activity) and is at the conference venue (environment). Note that due to the subjective matter of emotions, which often requires user input, emotion inference is excluded in current implementation.

**5.1.1 Indexing.** After each captured memory is structured, it is indexed and stored in a database. Additionally, the annotations in textual format are encoded into text embeddings to enable vector-based search during the retrieval process. In the database, each data entry represents a data entry corresponding to a captured memory, containing both the original media (e.g., photo, video) and its structured annotations in text and text embeddings.

## 5.2 Step 2: Identifying Composite Context

As captured memories are recorded in a linear manner along a personal timeline, those interconnected through semantic contexts often cluster closely together. For example, memories related to CHI 2024 are likely to occur during the event itself. Taking advantage of this **temporal proximity**, OmniQuery adopts a sliding window approach to analyze potentially interconnected captured memories in segments, identifying composite contexts.

As shown in Figure 5a, a static window size of seven days is used in current implementation. The inference is performed via an LLM, in which the input is the structured annotations of these memories and the output is the identified composite contexts along with their start and end dates and the associated captured memories

(Figure 5b). To account for cases where composite contexts are split in half, we use a step size (4 days in the current setup) smaller than the window size, allowing for overlap and comprehensive processing. For longer composite contexts (e.g., lasting more than two weeks), each segment of the context is identified separately within the sliding windows and then merged into a single composite context. Additionally, any duplicated composite contexts caused by the overlap between sliding windows are also merged to avoid redundancy (Figure 5c).

Specifically, as opposed to including detailed predefined categories (as with atomic contexts) in the prompt for LLMs, we adopt a few-shot prompting technique [7], providing examples of composite contexts summarized from the collected questions in the prompt. For detailed prompt, please refer to Appendix A.2.

**Explicitly mentioned contexts.** However, some composite contexts are **explicitly** mentioned in the captured memories. For example, a screenshot of a flyer may reference the upcoming "CHI 2024" event happening next month, or a transcribed conversation might discuss a "Hawaii trip" that took place the previous year. We leverage LLMs to differentiate between atomic contexts and composite contexts. For example, "a workout session" is identified as an activity (atomic context), whereas "CHI 2024" is recognized as a composite context which likely involves multiple interconnected atomic contexts. Such identified composite contexts are either merged with an existing composite context (e.g., if "CHI 2024" has already been identified) or added as a new composite context if it is unique.

## 5.3 Step 3: Inferring Semantic Knowledge

Different from composite contexts, semantic knowledge focuses on high-level general knowledge rather than specific memory details. For example, if a chat message screenshot mentions celebrating Jason's birthday, the inferred semantic knowledge might be "Jason's birthday is on [SPECIFIC DATE]." Similarly, analyzing multiple

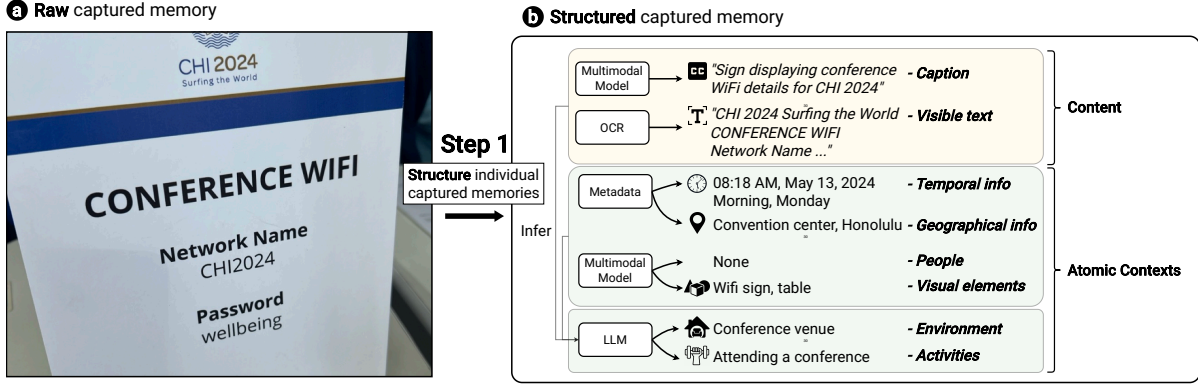


Figure 4: An example of structuring an individual captured memory (a photo of the Wifi details of CHI 2024 conference).

### Example result of Step 2 and Step 3

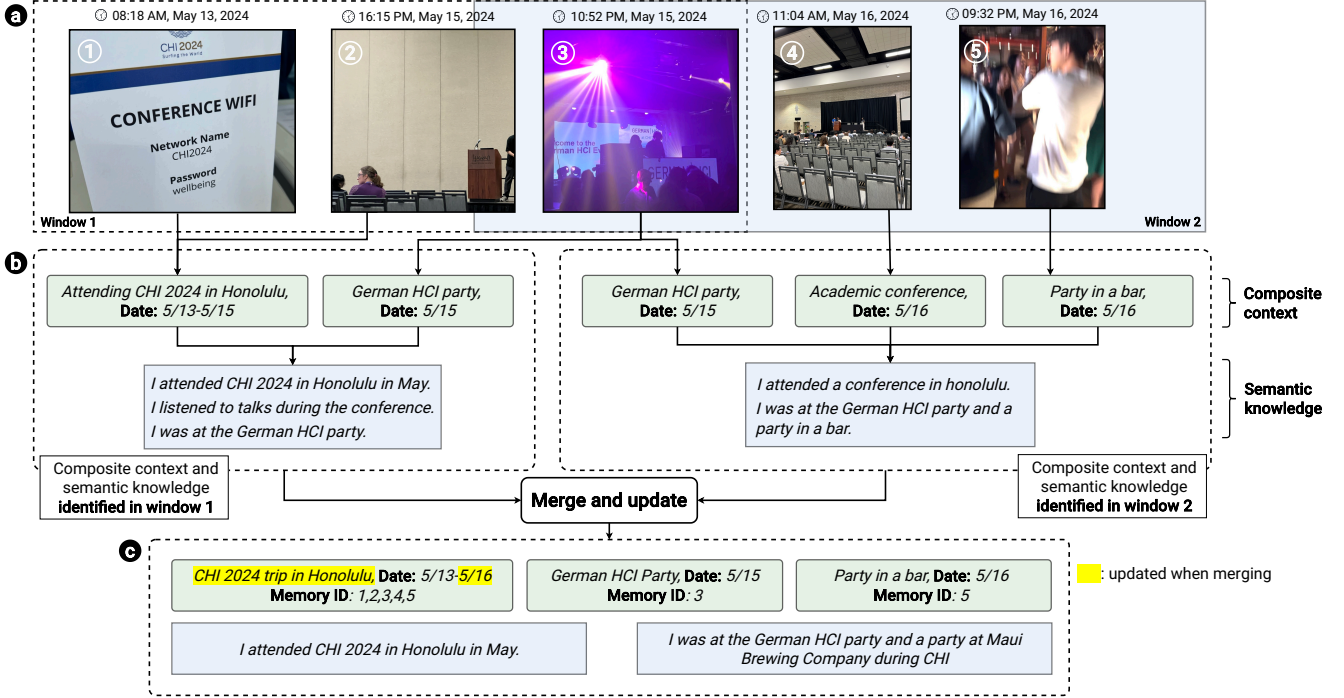


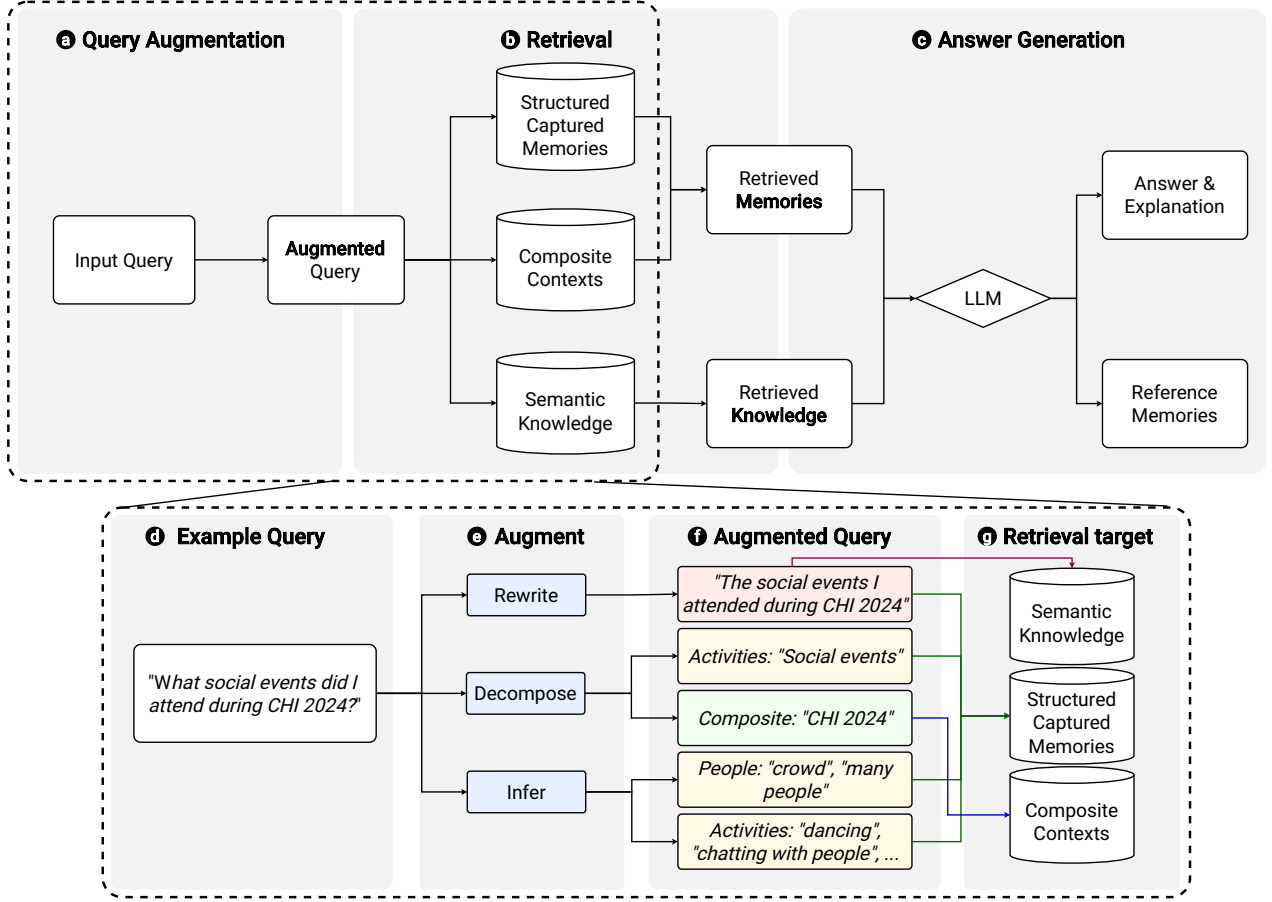
Figure 5: An example of using sliding windows to identify composite contexts and infer semantic knowledge: (a) two consecutive sliding windows; (b) composite contexts and semantic knowledge generated in each window; (c) merging results from the two windows.

grocery shopping receipts that consistently include lactose-free milk could lead to the inference that the user is possibly lactose intolerant. Semantic knowledge is inferred in each sliding window, while also taking into account the identified composite contexts to gain higher-level understanding of the user’s past and generalized information (Figure 5b bottom). The output is a list of inferred declarative semantic knowledge independent from specific memories. The instructions provided to the model are specifically tailored to guide the inference process toward overarching patterns and

trends rather than specific event details. The detailed prompt for identifying semantic knowledge can be also found in Appendix A.2. Each inferred entry of semantic knowledge is either merged with existing entries or added to the knowledge list if new.

## 5.4 Implementation Details

To deduplicate images as people tend to capture similar content multiple times, we use CLIP [49] to encode images to embeddings



**Figure 6: The question-answering system consists of: (a) query augmentation by decomposing and inferring contextual information; (b) retrieving memories and semantic knowledge; (c) generating answers with an LLM using referenced memories. Given a query: (d) OmniQuery rewrites and decomposes it; (e) infers contextual information; (f) retrieves relevant memories; and (g) retrieves knowledge from memory and knowledge storage.**

and calculate the similarities between images and merge those exceed the similarity of 0.85. We use the Google Cloud Vision API<sup>3</sup> for OCR to detect text in images and OpenAI’s Whisper model<sup>4</sup> for audio-to-text conversion. Note that Whisper is known for hallucination when there is no speech in the audio, thus we applied further data cleaning to validate the transcribed result using OpenAI’s GPT-4o-mini. For other visual processing, We use GPT-4o handles multimodal sensing, including identifying people and visual elements in images and generating scene descriptions. For video processing, as a proof-of-concept, we consider only the first 10 seconds of each video, sampling 10 frames to be analyzed by GPT-4o for content understanding. Text is encoded into embeddings using OpenAI’s text-embedding-3-small model. Currently, we utilize a custom vector database and matrix-based similarity search

implemented with NumPy in Python. However, for real-world applications, more advanced vector databases (e.g., Pinecone<sup>5</sup>) would be necessary to handle larger volumes of personal data.

## 6 OMNIQUERY: QUESTION-ANSWERING SYSTEM

With Captured memories augmented with contextual information, OmniQuery adopts a RAG architecture for the question answering system. RAG-based systems are effective in handling large datasets and mitigating hallucination issues by retrieving relevant content and grounding the generated results in this retrieved information. This approach ensures that the output is both relevant and accurate, leveraging specific data rather than relying solely on the model’s internal knowledge. This approach is chosen because, on average, personal captured memories often exceed 30,000 photos and videos (as reported by participants in our diary study), which exceeds the limit of most foundation models nowadays.

<sup>3</sup><https://cloud.google.com/vision/docs/ocr>

<sup>4</sup><https://github.com/openai/whisper>

<sup>5</sup><https://www.pinecone.io/>

As shown in Figure 6, given an input query, OmniQuery first augments the query by disambiguating and decomposing it into specific contextual elements (Figure 6a). Then, it retrieves the relevant captured memories from the structured captured memories and the composite contexts, along with related knowledge from the list of semantic knowledge (Figure 6b). The retrieved memories and knowledge, along with the augmented query, are then sent to an LLM to generate a comprehensive answer (Figure 6c). The detailed implementation of each step is discussed below.

## 6.1 Query Augmentation

As mentioned in Section 3.3, most user queries tend to be hybrid in nature or require contextual information. This means that directly searching based solely on the content of captured memories often results in an incomplete or insufficient retrieval of relevant memories. To enhance the retrieval process, OmniQuery adopts the query refinement concept from NLP [9] to augment the queries. The query augmentation process involves

- (1) **Rewriting the query** to declarative format to improve search accuracy of vector-based similarity matching;
- (2) **Decomposing the query** to extract necessary contextual filters, such as time, location, or events, which are grounded in the taxonomy. Note that only explicitly mentioned temporal contexts like "... last week" will be recognized temporal filters. Phrase like "... during CHI 2024" are part of a composite context, thus not counted as a temporal filter;
- (3) **Inferring potential related contexts** that may not be explicitly mentioned in the query but can enhance the filtering process also grounded in the taxonomy.

For example, as shown in Figure 6d-g, the query "*What social events did I attend during CHI 2024?*" is rewritten into a declarative format of "*The social events I attended during CHI 2024.*". Since "CHI 2024" is explicitly mentioned and identified as a composite context, it is extracted and labeled with the appropriate composite context tag. "Social events" is also extracted and identified as an atomic context (activities). Additionally, because "social events" might include various activities like parties, dancing, or casual conversations, and likely involve multiple people present, OmniQuery infers the relevant atomic contexts (people and activities) and annotates them in the corresponding context category.

## 6.2 Retrieving Relevant Augmented Memories

The decomposed augmented query is used to comprehensively retrieve from the augmented captured memories. Specifically, the augmented captured memories consist of the structured captured memories (with processed content and annotated atomic contexts), the list of identified composite contexts and the list of semantic knowledge. OmniQuery uses the decomposed components from the augmented query to perform a multi-source retrieval, pulling related memories from each of these sources. The results are then consolidated into a comprehensive list of relevant memories, which are used to generate an accurate and detailed answer for the user's query.

**Declarative query** → *Semantic knowledge & Processed content.* The declarative query is first encoded into text embeddings and used

to search both the semantic memories and processed content (caption and visible text) of the captured memories. This initial search step focuses on finding knowledge and memories directly related to the input query, without incorporating additional contextual filters.

### **Decomposed atomic contexts** → *Annotated atomic contexts.*

Each element of the decomposed atomic contexts (both extracted or inferred) is encoded into text embeddings and searched through the corresponding categories in the structured captured memory database. For example, if the query involves activities like "party" and "dancing," OmniQuery searches for captured memories annotated with similar activities. Any memories that have been annotated with related or similar activities will be retrieved, ensuring that relevant memories are included in the results. Additionally, *temporal contexts* apply a **strict** filter, excluding memories outside the specified time frame (e.g., "last month") from the retrieval process.

### **Decomposed composite contexts** → *Identified composite contexts.*

Any composite context decomposed in the augmenting process is also encoded into text embeddings and searched through the list of identified composite contexts. All captured memories linked to the semantically similar composite contexts are retrieved. This ensures all memories related to the composite contexts are included. Additionally, OmniQuery leverages an LLM to assess whether a composite context includes temporal constraints. For example, "... during CHI 2024" implies a strict temporal filter, while *photos related to CHI 2024* does not.

## 6.3 Answer Generation

The retrieved results is then sent to an LLM to generate the final answer. Specifically, the input for the LLM consists of: (1) the augmented query, (2) the retrieved semantic knowledge from the list, (3) all the retrieved captured memory entries from the annotated database, including both the memory content and its associated contextual annotations.

The model analyzes and reasons which captured memories serve as references for the generated answer. These reference memories are also included in the final answer presented to the user. To enhance the reasoning process, OmniQuery leverages chain-of-thought prompting [59], ensuring the generation is more accurate and contextually rich. For specific prompts, please refer to Appendix A.3

# 7 RESULTS

To demonstrate OmniQuery's capabilities in answering users' personal memory-related questions, we showcase four challenging examples of answering hybrid queries as shown in Figure 7.

## 7.1 Experiment Setting

OmniQuery was tested on the personal data of one of the researchers from the team, consisting of 2,590 images and videos downloaded from the researcher's smartphone. The data spans from mid-March to mid-August 2024, covering various activities, trips, and events, including a trip to CHI 2024, a summer trip to Hawaii, as well as captured memories of restaurants, fitness logs, and more.



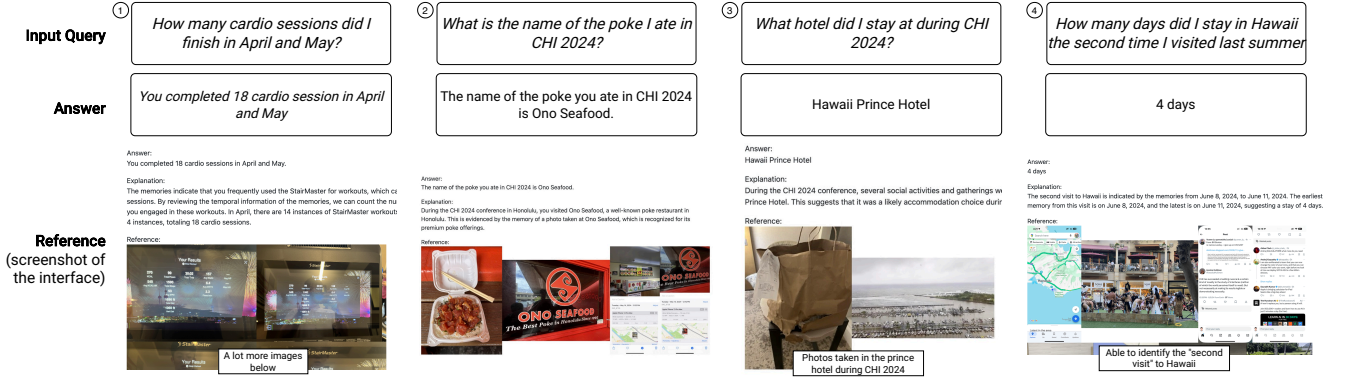


Figure 7: Four exemplar results of using OmniQuery to answer hybrid personal memory-related questions.

## 7.2 Example Walkthrough

In **example one**, the user asked about the number of cardio sessions completed over the previous two months. OmniQuery successfully retrieved all relevant memories (photos of the stair master machine taken after each session in April and May) and accurately generated the answer (18 sessions), providing every instance as supporting evidence.

In **example two**, the user sought the name of a Poke restaurant they visited during CHI 2024 in Honolulu. They remembered taking a photo of the meal but forgot the restaurant’s name. OmniQuery successfully retrieved the relevant memories based on the query and provided the correct answer.

In **example three**, the user asked about the hotel they stayed at during CHI. Although no explicit memories indicate the stay at the Prince Hotel, OmniQuery inferred the hotel from metadata associated with the photos taken at the venue, such as a nighttime photo of a broken takeout bag and a morning view of the marina.

In **example four**, the user asked about the length of stay during their “second” visit to Hawaii. OmniQuery was able to identify the second trip (the first was the CHI trip) and accurately count the number of stay during the visit.

## 8 USER EVALUATION

We conducted a self-guided user evaluation to let participants install and use OmniQuery on their local machine with their local album data. This is to protect participants’ personal data including sensitive and personal identifiable information. In this section, we discuss the detailed process and the evaluation result.

### 8.1 Participants

We recruited 10 participants, including seven from our diary study and three additional participants via word-of-mouth. We recruited participants who have basic programming skills so that they were able to install OmniQuery from source on their local machine. They also consented to the whole process, including that their filtered personal data will be processed via an API service. All the 10 participants (4 male, 6 female, age range = 22 to 29,  $age_{mean} = 25.3$ ,  $age_{SD} = 2.63$ ) were fluent or native English speaker. Participants rated their frequency of logging their daily lives as ‘Only record essential

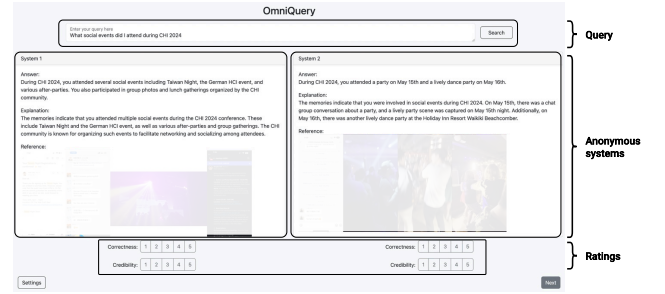


Figure 8: User interface used in the user evaluation.

information’ (1), ‘Regularly log important events and memorable experiences’ (5) and ‘Actively log my daily life’ (4). Each participant is compensated with \$50 for completing this study.

### 8.2 Apparatus

Two different systems were implemented in the user evaluation: the OmniQuery pipeline and a baseline system for comparison. The baseline system leverages a conventional RAG architecture, retrieving related memories based solely on vector similarities between the query and the description of the memory. The same base language model and prompt structure as OmniQuery are then used to generate the answer. For detailed implementation of the baseline, please refer to Appendix B.

As shown in Figure 8, in the studies, participants were presented with a single text input box similar to search engine input boxes. After they typed in the question, they would see two answers generated by the two systems in a randomized order. Two rating questions on the answer’s accuracy and completeness were then shown under each answer, from a scale of 1-5.

### 8.3 Procedure

The user evaluation involved three stages: (1) system setup, (2) data preparation, and (3) the main testing session using user queries.

**System setup.** Participants were given the source code for OmniQuery to install the back-end and a web application on their

**Table 2: Quantitative Results of OmniQuery and Baseline, including UPA, UPC, and Accuracy (%)**

Metrics	OmniQuery			Baseline		
	UPA	UPC	ACC	UPA	UPC	ACC
Direct content query (24)	4.42	4.13	83.3	3.67	3.46	62.5
Contextual filter + Hybrid (113)	3.89	3.85	69.0	2.93	2.83	38.9
<b>All (137)</b>	<b>3.98</b>	<b>3.90</b>	<b>71.5</b>	3.06	2.94	43.1

\*ACC refers to accuracy, which is considered accurate when UPA  $\geq 4$  (mostly correct).

**Table 3: Direct comparison between OmniQuery and baseline**

Winner	Comparison Win Rate (%)			
	OmniQuery	Baseline	Tie	Both are bad
Direct content query (24)	50.0	8.3	33.3	8.3
Contextual filter + Hybrid (113)	53.1	11.5	19.5	15.9
<b>All (137)</b>	<b>52.6</b>	10.9	21.9	14.6

local machine. They had the option of an online walkthrough session with the researchers or following the setup instructions on a self-guided manner.

**Data preparation.** Participants were asked to transfer a set of captured memories (photos and videos) from their smartphone’s album to their laptops since OmniQuery was installed on their local laptops. They were also asked to manually review and filter out any content deemed sensitive or preferred to be excluded from the study. After this filtering process, a total of 1,000 photos and videos were selected for data preprocessing.

Depending on how frequently participants logged their daily lives, the 1,000 selected photos and videos spanned anywhere from *one to four months* of past memories. These captured memories were then fed into OmniQuery for the query-agnostic augmentation process. The safety of the process data is ensured following the API’s privacy protocol<sup>6</sup>

**Main session.** The main session lasts 45 minutes, in which participants tested OmniQuery using two types of questions: (1) questions logged during the diary study and (2) questions they generated during the session. Participants checked on the diary study questions manually to determine if they could be answered using the filtered set of data on captured memories. Additionally, participants were encouraged to brainstorm and use new questions that were potentially answerable using the filtered data to comprehensively test OmniQuery.

In the question-answering procedure, OmniQuery-generated answers are accompanied with answers generated from the baseline system implemented using RAG. Each system generated answers anonymously, and participants compared and rated the results for both systems. The answers and user ratings were recorded for quantitative analysis. Throughout the process, participants were asked to think aloud [48], and a brief interview was conducted at

the end of the session to gather feedback and suggestions. These results were recorded for qualitative analysis.

#### 8.4 Comparison Metrics

After two answers were presented for a question, participants were asked to rate the two answers. We used the Chatbot Arena evaluation method [14], where participants compared answers from two systems and selected the better one or marked it as a tie. More specifically, for each question, participants rated the **user perceived accuracy** (UPA) and **user perceived completeness** (UPC) of the answers from both systems.

The UPA score was rated on a scale from 1 to 5: 1: Completely wrong or invalid result; 2: Incorrect, but provides at least some insight that helps answer or further refine the question; 3: Partially correct, or contains a subset of correct answers (e.g., only listing one meal when asked about all meals eaten last week); 4: Mostly correct, but missing some minor details (e.g., missing one subway trip when asked how many times I rode the subway); and 5: Completely correct. The UPC score, on the other hand, focused on the completeness and credibility of the given answers. In most cases, the questions asked by participants were challenging and needed to be explained and powered by evidence from the captured memory data. For example, when asked questions which have numeric answers like “how many meals did I have during my last New York trip,” the systems could be right on the numeric answer but wrong at which meals were counted. In those cases, participants examined the answers by looking back at the filtered data, and gave ratings on how they feel about the completeness and credibility of each answer.

We also directly compare the ratings of OmniQuery and the baseline to determine if one can outperform the other. In the comparison, if both systems have a UPA of 2 or lower (incorrect), the result is labeled “both are bad.” If at least one system provides a better-than-“partially correct” answer ( $\geq 3$ ), the system with the higher UPA is considered the winner. In cases where both systems

<sup>6</sup><https://openai.com/policies/privacy-policy/>

have the same UPA, the one with the higher UPC is the winner. Otherwise, it is a tie.

## 8.5 Quantitative Result

Participants tested 137 queries in total during the main session. Among them, 28 were previously logged during the diary study. We manually labeled each tested query using the categorization and definition mentioned in section 3.3. As a result, 24 were categorized as *direct content query* while 17 were *contextual filters* and 96 were *hybrid queries*. We analyzed the performance metrics of both systems (OmniQuery and baseline) using the scores rated by the participants. Table 2 and 3 summarize our results. In addition to presenting the average UPA and UPC scores, we calculated binary accuracy to evaluate whether the systems provided mostly correct answers. An answer was considered accurate if its UPA score was equal to or greater than 4 (mostly correct) (Table 2). We also present the “comparison result” in Table 3, which compares the two systems head-to-head on answering personal questions.

The result shows that overall, OmniQuery outperforms the baseline system in both the accuracy and completeness. Specifically, OmniQuery achieves an accuracy of 71.5%, outperforming the baseline by 28.4%, winning the comparison 52.6% of the time, and tying 21.9% of the time, while in 14.6% of the time, both results are bad. We also present the results for different categories of queries. The results indicate that simpler techniques like the baseline handle direct content queries reasonably well (62.5 % accuracy, and winning or tying 41.6% of the time). While the baseline struggles with more complex queries such as contextual filters or hybrid queries (38.9% accuracy, winning or tying 31.0% of the time), OmniQuery demonstrates its capabilities in effectively handling such queries (69.0% accuracy, winning or tying 72.6% of the time).

## 8.6 Qualitative Feedback and Findings

Apart from quantitative analysis, we also present qualitative feedback and findings from the think-aloud protocol and the exit interview on the usage of OmniQuery and suggestions of such intelligent question-answering system.

**8.6.1 Comparing with Existing Tools.** All participants have tried searching objects in their smartphone albums in their daily lives. They reported that they used the search features mostly to find a specific piece of information, like driver license, SSN number on the card, etc., matching the first type of questions (Direct Content Query) listed in section 3.3. They would also search for a specific event, like a trip to a specific location, matching the second type of questions (Contextual Filter). However, these searches are usually limited to retrieving a clear and specific object that users were looking for. They could not handle more complex questions like were collected in our diary studies. Plus, some of our participants also anticipated for this to happen because they “know what can be searched and what cannot be searched” from these existing album search tools (P2).

In the studies, a lot more challenging questions were asked. For **Direct Content Queries**, it would be challenging to answer when the object is ambiguous or when the users can only describe the object and do not know the exact name of it. For **Contextual**

**Filter** and **Hybrid** or even more open-ended and subjective questions, existing searching tools are not comparable to OmniQuery and the baseline at all because tools like iOS album search only return specific photos and videos in a whole without contextual understanding or filtering. In other words, these tools are all not “intelligent enough” to our participants.

In this subsection, we further summarize the cases that are hard to be accomplished using existing tools. In comparison to the high-level question types provided in section 3.3, we dived deeper into what these questions in the study were about, and provided detailed examples of these cases.

- **Exploratory Search:** When users know some characteristics of what they are searching for but cannot specify the exact object. For instance, P1 asked, “What churches did I visit in Barcelona?”
- **Look up and Locate:** When users know specific references or attributes about an item, such as date, location, or a person in the photo, and want to quickly locate the relevant media. For example, P4 asked, “Can you find the photo of me on a flyer on Instagram?”
- **Summarization Tasks:** Participants often need answers that summarize their collection of media, rather than finding a single item. For instance, P7 queried, “Which subway stations in New York have art installations?”
- **Comparative Questions:** Users sometimes want to compare different sets of media. For example, P10 asked, “Am I enjoying beach time more or hiking more?”
- **Open-ended and Subjective Questions:** Participants also asked questions that require interpretation or subjective judgment, which were even more challenging for existing tools. For example, P5 asked, “Given the photos I took, could you analyze what kind of person I am?”

In the meantime, we want to emphasize again that the comparison between OmniQuery and existing tools is conceptual, given that they serve different purposes and are designed differently in retrieving objects or answering questions. We provide this conceptual comparison to demonstrate the variety of questions OmniQuery can support answering.

**8.6.2 Reaction to Answers.** Participants reacted differently to different types of answers. They provided their observations as well as perceived feelings on the answers. Note that since participants weren’t aware of which answer was provided by OmniQuery or the baseline RAG system, we present their feedback on the overall answer structure, given that both OmniQuery and the baseline system provided a similar structure of answer and supporting materials.

**Dynamics of Detail and Concise Answer.** P7 mentioned that they did not always prefer to have a detailed answer composing all related media materials. More reference media would reduce the credibility of the answer to them. P8 also pointed out that if the answer just included all media it could find without a clear and concise connection to the answer, he “might just use iOS album search to get all photos containing [a specific object] and look by myself”

**Failure Cases of Answers.** We also present cases when participants reacted negatively to the answers. All participants encountered cases where the answers are inaccurate. Some were incomplete (e.g., P1 believed that they visited more than a few churches on the trip to Barcelona, but answers provided only two of them). Some were presumptive (e.g., P7 asked about recent social events, where the answers gave a piece of memory on a museum visit and explained that visiting museums is “likely with other people.” However, P7 visited the museum alone). Some were making mistakes (e.g., P7 asked for the mostly visited attractions but both systems mistakenly answered a museum, which was because P7 took a lot of museum pictures and both systems failed to recognize that they were the same visit.) Some even more challenging questions that caused failure of both systems include questions relate to a specific person. For example, P8 asked about her significant other and P5 asked about their “Korean friend” met in a trip. These cases represent the difficulty of understanding the nuances of personal relationships with personal album data.

**8.6.3 Iterative Editing of Questions.** Another recurring theme in our study is that six participants mentioned the possibility of iterating on the questions based on the answers. It occurred when participants were uncertain of what answers they were gonna get out of open-ended or subjective answers. They wanted to iterate on the questions based on the given answers, which gave them more understanding on the questions they wanted to ask. For example, P3 asked about how many places they visited during summer, and, based on the answer, realized they were more interested in the number of cities rather than the countries. It highlighted the potential need of a chat interface of personal captured memory data, which posit more challenges in this domain considering the query histories.

## 9 DISCUSSION

In this section, we draw on implications from our studies to discuss limitations and propose future work based on what we found.

### 9.1 From Chat Interface to Multimodal Interactions

As a system designed to answer users’ questions on their personal captured memory, OmniQuery is currently designed in an ask-and-react manner for the purpose of evaluating its efficacy in a lab-study settings. In our studies, participants were excited about what OmniQuery was capable of, and gave feedback on having more multimodal interactions rather than a “ChatGPT-like” interface. We recognize the potential of a more interactive OmniQuery in the following ways:

**Multimodal Input and Output:** Just like how ChatGPT iterated, OmniQuery could use more multimodal input including but not limited to audio, image, and even videos in the future. Recalling the summarized task cases which can be hard for existing album searching tools to accomplish, a number of these cases can inherently benefit from multimodal inputs. For example, users might want to look up and locate an oddly-shaped cup that they cannot refer to in plain text, or compare all existing dresses in album with a reference color on an image that is hard to describe. Natural language description’s limitation needs to be addressed by introducing

different input modalities. Besides input, OmniQuery also has huge potential in helping users relive their memories by visualizing their captured memory data. Users can access their captured memories in an interactive way, inspect and even modify annotations in a transparent way. With that, a “mind palace” style AI assistant becomes more possible.

**Error correction:** In our studies, we observed the importance of allowing users to review and access the identified composite contexts and semantic knowledge. Participants expressed the need to correct errors within the system when it retrieved incorrect or irrelevant information. For example, P9 asked about a specific KPop store, and while the system successfully generated relevant memories, it mistakenly included an Instagram screenshot of a Korean TV show. In such cases, participants wanted the ability to mark these irrelevant results as incorrect, suggesting that incorporating error correction mechanisms would enhance the system’s overall performance and accuracy in future interactions.

**Follow-up queries:** Another recurring theme in our study is that participants frequently desired the ability to refine their queries or ask follow-up questions. Six out of ten participants mentioned the need to clarify answers or iteratively narrow down their queries. For instance, P4 and P8 emphasized the importance of follow-up questions to hone in on specific information after receiving initial responses. Together with proposed multimodal interactions above, an iterative interaction model, similar to a chatbot workflow, would significantly improve usability and flexibility of OmniQuery.

### 9.2 Enriching Memory Data and Visual Intelligence

At present, OmniQuery primarily processes media from a smartphone’s photo album as its main source for captured memories. However, these media alone provide a limited view of a user’s broader personal knowledge. For example, in one of the study’s failure cases, OmniQuery struggled to infer personal relationships from social interactions captured in group photos. To enhance memory augmentation and improve retrieval accuracy, expanding OmniQuery’s data sources and visual intelligence is essential.

**Integrating additional data sources:** Personal knowledge extends beyond photo albums and exists across various applications. While our participants’ photo albums included screenshots of emails, calendar events, and chat histories, these represent only a fraction of the broader personal information available in other communication and social interaction apps. Incorporating data from such sources could significantly enhance OmniQuery’s contextual understanding, allowing for more complex queries and richer memory retrieval. However, integrating these additional data sources presents substantial privacy and ethical challenges. While our evaluations were conducted entirely on users’ local machines and did not explore privacy-preserving implementations in detail, existing research efforts, such as those focused on differential privacy and on-device machine learning, offer promising directions for secure and privacy-aware deployment. Additionally, commercial tools like Apple Intelligence’s private cloud computing serve as examples of ongoing progress in protecting user data while enabling advanced memory retrieval capabilities.

**Enhancing visual intelligence:** As discussed in Section 8.6.2, questions related to social interactions remain challenging due to the current lack of advanced features like facial recognition for person identification. Future iterations of OmniQuery could integrate such capabilities (with appropriate user consent), enabling the system to track individuals across various memories. This enhancement would support new use cases, such as monitoring social patterns or tracking progress over time, significantly improving the system’s capacity for memory augmentation and retrieval. Additionally, we propose exploring the design and implementation of a comprehensive taxonomy of personal knowledge domains. This would allow users to selectively activate or deactivate specific domains, such as enabling “Social Interactions and Relationships” to infer personal connections while disabling “Personally Identifiable Information” to prevent the system from processing sensitive data like IDs or SSN numbers in photos. This modular approach could enhance both user control and privacy.

**Augmenting with future AR technologies:** A limitation of personal memory capture is the potential for missed moments when users either forget or are unable to document an experience. As AR technology advances, OmniQuery’s memory augmentation and retrieval capabilities could be seamlessly integrated into AR systems, allowing for more passive and context-aware memory capture. AR devices could leverage real-time contextual triggers [36] to proactively surface relevant memories or information, offering proactive assistance in pervasive AR environments. This integration would enhance the user experience by making memory retrieval more intuitive and contextually relevant. However, such passive data capture raises even more significant privacy concerns, which will require future research into secure, privacy-preserving implementations to ensure the responsible use of AI in these settings.

### 9.3 Preserving Privacy

As was discussed in above subsections, protecting users’ privacy is crucial in developing future personal AI assistants, including but not limited to handling personal data such as media in albums, chat histories or browsing history. Users have limited control over how their data is handled and must rely on service providers’ adherence to privacy protocols. In this subsection, we take a step further to discuss more robust and rigorous measures should be adopted in real-world settings, where the amount of personal data is huge, making approaches like manual filtering in OmniQuery’s evaluation infeasible.

One way is to incorporate more advanced data protection techniques on cloud servers including data anonymization [42] and encryption [47], while preserving the computational capabilities of large models via online computing. The other approach is leveraging on-device computing, where all data processing occurs locally on the user’s device, ensuring full control over users’ own data. Recent advances in model compression [25] have made it possible to run large model on smaller devices like smartphones. As OmniQuery is designed to be model-agnostic, it is able to work with different model sizes. While smaller, compressed on-device model may result in reduced performance, future work should focus on developing curated datasets and benchmarks to rigorously evaluate OmniQuery’s performance across different model sizes

(e.g., LLaMAs [55] and Phi-3 [1]). This would provide a deeper understanding of how model size impacts both privacy and system effectiveness. concern is that users do not have direct control over their own data when leveraging powerful AI models via online servers.

## 10 CONCLUSION

In this paper, we present OmniQuery, which enhances personal question answering on captured memories. To inform the design of OmniQuery, we ran a one-month diary study to collect and analyze realistic user queries, which were then used to generate a taxonomy of contextual information, including atomic context, composite context and semantic knowledge. The taxonomy guided the design of an augmentation pipeline of captured memories, which involves structuring individual captured memories, and processing multiple captured memories using sliding windows to identify composite contexts and semantic knowledge.

We then developed a personal question-answering system, which first augment the user query, then retrieve the related captured memories, and finally using the retrieved memories to generate comprehensive answers using an LLM. We evaluated OmniQuery against a baseline system with a user evaluation of 10 participants who tested 137 queries in total. The results show that OmniQuery is effective in answering users’ queries with an accuracy of 71.5%, and winning or tying 74.5% of the time in direct comparison with the baseline. We also performed qualitative analysis on participants’ reaction and feedback during the evaluation, with the findings demonstrating potential of OmniQuery and providing valuable insights into possible enhancements for future systems alike.

## References

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benham, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Lil- iang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. arXiv:2404.14219 [cs.CL] <https://arxiv.org/abs/2404.14219>
- [2] Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Devi Parikh, and Dhruv Batra. 2015. VQA: Visual Question Answering. *International Journal of Computer Vision* 123 (2015), 4 – 31. <https://api.semanticscholar.org/CorpusID:3180429>
- [3] Akari Asai, Zeqi Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. arXiv:2310.11511 [cs.CL] <https://arxiv.org/abs/2310.11511>
- [4] Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150* (2020).



- [5] Elise Bonnal, Wen-Jie Tseng, Mark McGill, Éric Lecolinet, Samuel Huron, and Jan Gugenheimer. 2023. Memory Manipulations in Extended Reality. *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (2023). <https://api.semanticscholar.org/CorpusID:257952236>
- [6] Joel Brandt, Noah Weiss, and Scott R. Klemmer. 2007. txt 4 l8r: lowering the burden for diary studies under mobile conditions. In *CHI '07 Extended Abstracts on Human Factors in Computing Systems* (San Jose, CA, USA) (CHI EA '07). Association for Computing Machinery, New York, NY, USA, 2303–2308. <https://doi.org/10.1145/1240866.1240998>
- [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. *arXiv:2005.14165 [cs.CL]* <https://arxiv.org/abs/2005.14165>
- [8] Niamh Caprani, John Greaney, and Nicola Porter. 2006. A Review of Memory Aid Devices for an Ageing Population. *Psychology J.* 4 (2006), 205–243. <https://api.semanticscholar.org/CorpusID:9598075>
- [9] Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. 2024. RQ-RAG: Learning to Refine Queries for Retrieval Augmented Generation. *arXiv:2404.00610 [cs.CL]* <https://arxiv.org/abs/2404.00610>
- [10] Samantha WT Chan. 2020. Biosignal-Sensitive Memory Improvement and Support Systems. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–7.
- [11] Wei Ting Samantha Chan. 2022. *Augmenting Human Prospective Memory through Cognition-Aware Technologies*. Ph.D. Dissertation. ResearchSpace@ Auckland.
- [12] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to Answer Open-Domain Questions. *ArXiv abs/1704.00051* (2017). <https://api.semanticscholar.org/CorpusID:3618568>
- [13] Wenhui Chen, Hexiang Hu, Xi Chen, Pat Verga, and William W. Cohen. 2022. MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text. In *Conference on Empirical Methods in Natural Language Processing*. <https://api.semanticscholar.org/CorpusID:252735160>
- [14] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot Arena: An Open Platform for Evaluating LLMs by Human Preference. *arXiv:2403.04132 [cs.AI]* <https://arxiv.org/abs/2403.04132>
- [15] Lydia Dubourg, Ana Rita Silva, Christophe Fitamen, Chris J. A. Moulin, and Céline Souchay. 2016. SenseCam: A new tool for memory rehabilitation? *Revue neurologique* 172 12 (2016), 735–747. <https://api.semanticscholar.org/CorpusID:9803824>
- [16] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. *arXiv:2404.16130 [cs.CL]* <https://arxiv.org/abs/2404.16130>
- [17] Jakob Engel, Kiran Somasundaram, Michael Goesele, Albert Sun, Alexander Gamino, Andrew Turner, Arjang Talattof, Arnie Yuan, Bilal Souti, Brighid Meredith, Cheng Peng, Chris Sweeney, Cole Wilson, Dan Barnes, Daniel DeTone, David Caruso, Derek Valleroy, Dinesh Ginjupalli, Duncan Frost, Edward Miller, Elias Mueggler, Evgeniy Oleinik, Fan Zhang, Guruprasad Somasundaram, Gustavo Solaira, Harry Lanaras, Henry Howard-Jenkins, Huixuan Tang, Hyo Jin Kim, Jaime Rivera, Ji Luo, Jing Dong, Julian Straub, Kevin Bailey, Kevin Eickenhoff, Lingni Ma, Luis Pesqueira, Mark Schwesinger, Maurizio Monge, Nan Yang, Nick Charron, Nikhil Raina, Omkar Parkhi, Peter Borschowa, Pierre Moulon, Prince Gupta, Raul Mur-Artal, Robbie Pennington, Sachin Kulkarni, Sagar Miglani, Santosh Gondi, Saransh Solanki, Sean Diener, Shangyi Cheng, Simon Green, Steve Saarinen, Suvarn Patra, Tassos Mourikis, Thomas Whelan, Tripti Singh, Vasileios Balntas, Vijay Baiyya, Wilson Dreeves, Xiaqing Pan, Yang Lou, Yipu Zhao, Yusuf Mansour, Yuyang Zou, Zhaoyang Lv, Zijian Wang, Mingfei Yan, Carl Ren, Renzo De Nardi, and Richard Newcombe. 2023. Project Aria: A New Tool for Egocentric Multimodal AI Research. *arXiv:2308.13561 [cs.HC]* <https://arxiv.org/abs/2308.13561>
- [18] Yue Fan, Xiaojian Ma, Ruijie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. 2024. VideoAgent: A Memory-augmented Multimodal Agent for Video Understanding. *arXiv preprint arXiv:2403.11481* (2024).
- [19] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* (2023).
- [20] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. *International Journal of Computer Vision* 127 (2016), 398–414. <https://api.semanticscholar.org/CorpusID:8081284>
- [21] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh K. Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Z. Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragomeni, Qichen Fu, Christian Fuegen, Ahrham Kahsay Gebreselasie, Cristina González, James M. Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Weslie Khoo, Jáchym Kolár, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Kartikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merer Ramazanov, Leda Sari, Kiran K. Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Yunyi Zhu, Pablo Arbeláez, David J. Crandall, Dima Damen, Giovanni Maria Farinella, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard A. Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. 2021. Ego4D: Around the World in 3,000 Hours of Egocentric Video. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021), 18973–18990. <https://api.semanticscholar.org/CorpusID:238856888>
- [22] Cathal Gurrin, Alan F Smeaton, Aiden R Doherty, et al. 2014. Lifelogging: Personal big data. *Foundations and Trends® in information retrieval* 8, 1 (2014), 1–125.
- [23] Gillian R. Hayes, Shwetak N. Patel, Khai Nhut Truong, Giovanni Iachello, Julie A. Kientz, Rob Farmer, and Gregory D. Abowd. 2004. The Personal Audio Loop: Designing a Ubiquitous Audio-Based Memory Aid. In *Mobile HCI*. <https://api.semanticscholar.org/CorpusID:11316625>
- [24] Steve Hodges, Lyndsay Williams, Emma Berry, Shahram Izadi, James Srinivasan, Alex Butler, Gavin Smyth, Narinder Kapur, and Ken Wood. 2006. SenseCam: A retrospective memory aid. In *UbiComp 2006: Ubiquitous Computing: 8th International Conference, UbiComp 2006 Orange County, CA, USA, September 17-21, 2006 Proceedings* 8. Springer, 177–193.
- [25] Fred Hohman, Mary Beth Kery, Donghao Ren, and Dominik Moritz. 2023. Model Compression in Practice: Lessons Learned from Practitioners Creating On-device Machine Learning Experiences. *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2023). <https://api.semanticscholar.org/CorpusID:263829166>
- [26] Eric Horvitz. 1999. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*. 159–166.
- [27] Zhijian Hou, Lei Ji, Difei Gao, Wanjun Zhong, Kun Yan, Chao Li, Wing-Kwong Chan, Chong-Wah Ngo, Nan Duan, and Mike Zheng Shou. 2023. GroundNLQ @ Ego4D ‘Natural Language Queries Challenge 2023. *arXiv:2306.15255 [cs.CV]* <https://arxiv.org/abs/2306.15255>
- [28] Xiao Huang, Jingyuan Zhang, Dingcheng Li, and Ping Li. 2019. Knowledge graph embedding based question answering. In *Proceedings of the twelfth ACM international conference on web search and data mining*. 105–113.
- [29] Xiao Huang, Jingyuan Zhang, Dingcheng Li, and Ping Li. 2019. Knowledge Graph Embedding Based Question Answering. *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining* (2019). <https://api.semanticscholar.org/CorpusID:59528287>
- [30] Razan Jaber, Sabrina Zhong, Sanna Kuoppamäki, Aida Hosseini, Iona Gessinger, Duncan P Brumby, Benjamin R Cowan, and Donald McMillan. 2024. Cooking With Agents: Designing Context-aware Voice Interaction. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–13.
- [31] Matthew Jamieson, Breda Cullen, Marilyn McGee-Lennon, Stephen Brewster, and Jonathan J Evans. 2014. The efficacy of cognitive prosthetic technology for people with memory impairments: A systematic review and meta-analysis. *Neuropsychological rehabilitation* 24, 3-4 (2014), 419–444.
- [32] Matthew Jamieson, Brian O’Neill, Breda Cullen, Marilyn Rose McGee-Lennon, Stephen Anthony Brewster, and Jonathan J. Evans. 2017. ForgetMeNot: Active Reminder Entry Support for Adults with Acquired Brain Injury. *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (2017). <https://api.semanticscholar.org/CorpusID:2298134>
- [33] Jaewook Lee, Jun Wang, Elizabeth Brown, Liam Chu, Sebastian S. Rodriguez, and Jon E. Froehlich. 2024. GazePointAR: A Context-Aware Multimodal Voice Assistant for Pronoun Disambiguation in Wearable Augmented Reality. *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2024). <https://api.semanticscholar.org/CorpusID:268132284>
- [34] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv:2005.11401 [cs.CL]* <https://arxiv.org/abs/2005.11401>
- [35] Franklin Mingzhe Li, Michael Xieyang Liu, Shaun K. Kane, and Patrick Carrington. 2024. A Contextual Inquiry of People with Vision Impairments in Cooking. *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2024). <https://api.semanticscholar.org/CorpusID:267897983>
- [36] Jiahao Nick Li, Yan Xu, Tovi Grossman, Stephanie Santosa, and Michelle Li. 2024. OmniActions: Predicting Digital Actions in Response to Real-World Multimodal Sensory Inputs with LLMs. In *Proceedings of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '24). Association for Computing Machinery, New York, NY, USA, Article 8, 22 pages. <https://doi.org/10.1145/3613904.3642068>

- [37] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. Improved Baselines with Visual Instruction Tuning.
- [38] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. LLaVA-NeXT: Improved reasoning, OCR, and world knowledge. <https://llava-vl.github.io/blog/2024-01-30-llava-next/>
- [39] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning.
- [40] Xingyu Bruce Liu, Jiahao Nick Li, David Kim, Xiang 'Anthony' Chen, and Ruofei Du. 2024. Human I/O: Towards a Unified Approach to Detecting Situational Impairments. *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2024). <https://api.semanticscholar.org/CorpusID:268203741>
- [41] Muhammad Maaz, Hanoona Abdul Rasheed, Salman H. Khan, and Fahad Shah-baz Khan. 2023. Video-ChatGPT: Towards Detailed Video Understanding via Large Vision and Language Models. In *Annual Meeting of the Association for Computational Linguistics*. <https://api.semanticscholar.org/CorpusID:259108333>
- [42] Abdul Majeed and Sungchang Lee. 2021. Anonymization Techniques for Privacy Preserving Data Publishing: A Comprehensive Survey. *IEEE Access* 9 (2021), 8512–8545. <https://api.semanticscholar.org/CorpusID:231616865>
- [43] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2023. EgoSchema: A Diagnostic Benchmark for Very Long-form Video Language Understanding. [arXiv:2308.09126 \[cs.CV\]](https://arxiv.org/abs/2308.09126) <https://arxiv.org/abs/2308.09126>
- [44] Steve Mann. 1996. Wearable Tetherless Computer-Mediated Reality: WearCam as a wearable face-recognizer, and other applications for the disabled. <https://api.semanticscholar.org/CorpusID:11838759>
- [45] Vaibhav Mavi, Anubhav Jangra, and Adam Jatowt. 2024. Multi-hop Question Answering. [arXiv:2204.09140 \[cs.CL\]](https://arxiv.org/abs/2204.09140) <https://arxiv.org/abs/2204.09140>
- [46] Ken McRae and Michael Jones. 2013. *14 Semantic Memory*. Vol. 206. Oxford University Press Oxford.
- [47] Aamer Nadeem and Muhammad Younus Javed. 2005. A Performance Comparison of Data Encryption Algorithms. *2005 International Conference on Information and Communication Technologies* (2005), 84–89. <https://api.semanticscholar.org/CorpusID:14441015>
- [48] Jakob Nielsen. 1994. *Usability engineering*. Morgan Kaufmann.
- [49] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. [arXiv:2103.00020 \[cs.CV\]](https://arxiv.org/abs/2103.00020) <https://arxiv.org/abs/2103.00020>
- [50] Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. 2024. RAPTOR: Recursive Abstractive Processing for Tree-Organized Retrieval. [arXiv:2401.18059 \[cs.CL\]](https://arxiv.org/abs/2401.18059) <https://arxiv.org/abs/2401.18059>
- [51] Jana Sedlakova, Paola Danio, Andrea Horn Wintsch, Markus Wolf, Mina Stanikic, Christina Haag, Chloé Sieber, Gerold Schneider, Kaspar Staub, Dominik Alois Ettlin, et al. 2023. Challenges and best practices for digital unstructured data enrichment in health research: a systematic narrative review. *PLOS Digital Health* 2, 10 (2023), e0000347.
- [52] Youngsoo Shin, Ruth Barankevich, Jina Lee, and Saleh Kalantari. 2021. PEN-CODER: Design for Prospective Memory and Older Adults. *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems* (2021). <https://api.semanticscholar.org/CorpusID:233987041>
- [53] Timothy Sohn, Kevin A. Li, William G. Griswold, and James D. Hollan. 2008. A diary study of mobile information needs. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Florence, Italy) (CHI '08). Association for Computing Machinery, New York, NY, USA, 433–442. <https://doi.org/10.1145/1357054.1357125>
- [54] Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. 2021. MultiModalQA: Complex Question Answering over Text, Tables and Images. [arXiv:2104.06039 \[cs.CL\]](https://arxiv.org/abs/2104.06039) <https://arxiv.org/abs/2104.06039>
- [55] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. [arXiv:2302.13971 \[cs.CL\]](https://arxiv.org/abs/2302.13971) <https://arxiv.org/abs/2302.13971>
- [56] Endel Tulving. 2002. Episodic Memory: From Mind to Brain. *Annual Review of Psychology* 53, Volume 53, 2002 (2002), 1–25. <https://doi.org/10.1146/annurev.psych.53.100901.135114>
- [57] Sunil Vemuri, Chris Schmandt, Walter Bender, Stefanie Tellex, and Bradford Lassey. 2004. An Audio-Based Personal Memory Aid. In *Ubiquitous Computing*. <https://api.semanticscholar.org/CorpusID:309402>
- [58] Cunxiang Wang, Ruoxi Ning, Boqi Pan, Tonghui Wu, Qipeng Guo, Cheng Deng, Guangsheng Bao, Qian Wang, and Yue Zhang. 2024. Novelqa: A benchmark for long-range novel question answering. [arXiv preprint arXiv:2403.12766](https://arxiv.org/abs/2403.12766) (2024).
- [59] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. [arXiv:2201.11903 \[cs.CL\]](https://arxiv.org/abs/2201.11903) <https://arxiv.org/abs/2201.11903>
- [60] G. Yang, Lekha Chaisorn, Yunlong Zhao, Shi-Yong Neo, and Tat-Seng Chua. 2003. VideoQA: question answering on news video. *Proceedings of the eleventh ACM international conference on Multimedia* (2003). <https://api.semanticscholar.org/CorpusID:207716>
- [61] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering. In *Conference on Empirical Methods in Natural Language Processing*. <https://api.semanticscholar.org/CorpusID:52822214>
- [62] Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. QA-GNN: Reasoning with language models and knowledge graphs for question answering. [arXiv preprint arXiv:2104.06378](https://arxiv.org/abs/2104.06378) (2021).
- [63] Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. QA-GNN: Reasoning with Language Models and Knowledge Graphs for Question Answering. In *North American Chapter of the Association for Computational Linguistics*. <https://api.semanticscholar.org/CorpusID:233219869>
- [64] Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-LLaMA: An Instruction-tuned Audio-Visual Language Model for Video Understanding. In *Conference on Empirical Methods in Natural Language Processing*. <https://api.semanticscholar.org/CorpusID:259075356>
- [65] Liangfu Zhang, Anwen Hu, Jing Zhang, Shuo Hu, and Qin Jin. 2023. MPMQA: Multimodal Question Answering on Product Manuals. [ArXiv abs/2304.09660](https://arxiv.org/abs/2304.09660) (2023). <https://api.semanticscholar.org/CorpusID:258212471>
- [66] Wazeer Deen Zulfikar, Samantha Chan, and Pattie Maes. 2024. Memoro: Using Large Language Models to Realize a Concise Interface for Real-Time Memory Augmentation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '24). Association for Computing Machinery, New York, NY, USA, Article 450, 18 pages. <https://doi.org/10.1145/3613904.3642450>

## A Prompts for LLMs

### A.1 Identifying Composite Contexts

System instruction:

You are an intelligent agent capable of generating a list of COMPOSITE CONTEXTS inferred from the given memory.

Composite context refers to a combination of time, location, people, objects, environment and activities. Such composite contexts could be inferred from the explicit content (e.g., text showing the event info) or implicit cues (e.g., multiple changes in location indicating travel). Focus on relatively important composites such as travel, conferences, and important meetings and focus less on trivial events.

For each composite context, identify the related episodic memory ids. This could be due to time (e.g., the memory occurs during the event), location (e.g., the memory takes place at the event location), or specific content (e.g., the memory mentions the event).

Additionally, rate the importance of each event on a scale from 1 to 3, where 3 denotes very major events (e.g., multi-day events or highly important events), 2 denotes moderately important events, and 1 denotes less important events.

Exemplar composite context types include:

An academic conference: "An academic conference";

Recreational travel: "Trip to Salt lake city", "Traveling to home town";

Locational change: "Location changed from Seattle to Irvine";

Outdoor activities: "Camping trip";

Personal milestones: "Birthday celebration", "Graduation ceremony", "first day in univeristy";

etc.

Output the list of composite context in a JSON object with the key 'composite\_context'. Each event should be represented as a sub JSON object with the following keys: 'event\_name' (detailed and concise), 'memory\_ids' (list), 'start\_date', 'end\_date' (could be the same as start\_date), 'location', 'is\_multi\_days', and 'importance'.

+

<List of structured captured memories>

### A.2 Inferring Semantic Knowledge

System instruction:

You are an intelligent agent capable of generating a list of FACTS or KNOWLEDGE (referred to knowledge in the following) that can be inferred from the given memory and the related composite contexts. Focus on relatively important high-level semantic knowledge and focus less on trivial events. Avoid specific details about individual media

The knowledge should be detailed and self-contained.

Exemplar semantic knowledge includes:

<Examples of semantic knowledge>

Also identify the most representative episodic memories that contribute to the understanding of the knowledge.

Output a JSON object with the key 'knowledge'. Each knowledge item should include 'knowledge', 'memory\_ids' (list)

Input:

<Structured captured memories in the sliding windows>

+

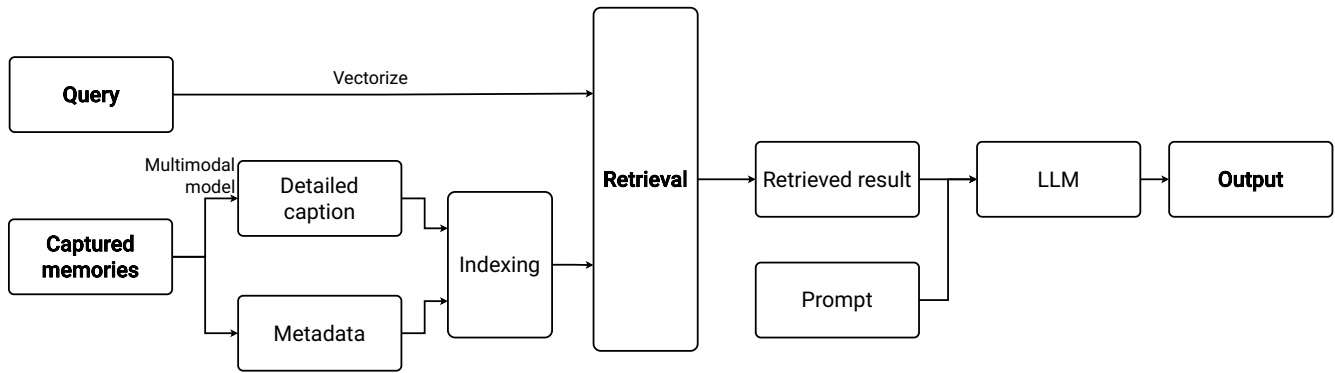
<Identified composite contexts identified in the sliding window>

### A.3 Generating Answers Based on Retrieved Results

System instruction:

Given a query, a list of memories and personal knowledge, generate a comprehensive answer to the query.

Identify the episodic memories that can provide evidence to the question.



**Figure 1: Structure of the baseline implementation.**

If the answer is not explicitly presented in the memories, make a reasonable inference.

Output a JSON object with the key 'answer', 'explanation' and 'memory\_ids'. The 'answer' should be a string and 'memory\_ids' should be a list of memory ids

Input:  
 <Query>  
 +  
 <Retrieved semantic knowledge>  
 +  
 <Retrieved structured knowledge>

geographical information from the metadata and processes it in the same manner as OmniQuery. This ensures that the processed memories include the temporal and geographical data, which are common components in users' queries. The temporal and geographical information is concatenated to the generated caption. Then the concatenated text sequence is encoded into text embeddings using embedding models (text-embedding-3-small).

In the retrieval stage, the query is first encoded into the text embeddings using the same embedding model, and then retrieve the top K (K=50) captured memories using vector-based similarity search. The retrieved top K captured memories are then ordered in temporal sequence, and then sent to the LLM (GPT-4o) for generating the answer. The prompt used for the answer generation is the same as OmniQuery.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

## B Baseline Implementation

While there is no already-existing system designed for answering personal questions on captured memories, we manually designed and implemented a system as the baseline to compare with OmniQuery.

Similar to OmniQuery, the baseline system also adopts a RAG architecture to adapt to the large number of captured memories. We utilized the basic structure of RAG illustrated in [19], which involves (1) indexing the external data sources with embedding models, (2) leverage vector-based search to retrieve the top K relevant data instances (3) based on the retrieved data, utilizing a powerful LLM to generate the final answer. Note that typical RAG systems require a chunking phase, where long documents are split into smaller chunks for more precise matching and retrieval of relevant information. In our case, each captured memory already represents a limited amount of information and is naturally separated. Therefore, we treat each captured memory as an individual chunk.

Figure B1 demonstrates the structure of the baseline system in our experiment. The baseline also processes the captured memories by leveraging a multimodal model (GPT-4o) to generate detailed captions for each memory. Additionally, it extracts temporal and