# Beyond Fact Retrieval: Episodic Memory for RAG with Generative Semantic Workspaces

**Shreyas Rajesh, Pavan Holur, Chenda Duan, David Chong, Vwani Roychowdhury**

University of California, Los Angeles

{ shreyasrajesh38, pholur, chenda, davidchong13807, vwani}@ucla.edu

## Abstract

Large Language Models (LLMs) face fundamental challenges in long-context reasoning: many documents exceed their finite context windows, while performance on texts that do fit degrades with sequence length, necessitating their augmentation with external memory frameworks. Current solutions, which have evolved from retrieval using semantic embeddings to more sophisticated structured knowledge graphs representations for improved sense-making and associativity, are tailored for fact-based retrieval and fail to build the space-time-anchored narrative representations required for tracking entities through episodic events. To bridge this gap, we propose the **Generative Semantic Workspace** (GSW), a neuro-inspired generative memory framework that builds structured, interpretable representations of evolving situations, enabling LLMs to reason over evolving roles, actions, and spatiotemporal contexts. Our framework comprises an *Operator*, which maps incoming observations to intermediate semantic structures, and a *Reconciler*, which integrates these into a persistent workspace that enforces temporal, spatial, and logical coherence. On the Episodic Memory Benchmark (EpBench) (Huet, Houidi, and Rossi 2025) comprising corpora ranging from 100k to 1M tokens in length, GSW outperforms existing RAG based baselines by up to **20%**. Furthermore, GSW is highly efficient, reducing query-time context tokens by **51%** compared to the next most token-efficient baseline, reducing inference time costs considerably. More broadly, GSW offers a concrete blueprint for endowing LLMs with human-like episodic memory, paving the way for more capable agents that can reason over long horizons.

## Introduction

Large Language Models (LLMs) have transformed natural language understanding, but their ability to reason over long contexts is still limited by finite input windows. Even with token limits in the millions, large document collections can easily exceed these bounds. Performance can also degrade with context length due to phenomena like "context rot" and "lost-in-the-middle" effects (Liu et al. 2023; Hong, Troynikov, and Huber 2025). A common workaround is Retrieval-Augmented Generation (RAG), which supplements the LLM's input with only the most relevant retrieved content at query time. Standard RAG pipelines split documents into smaller chunks, encode them into dense embeddings, and retrieve the top-matching chunks based on semantic similarity to the query—allowing the LLM to focus on a relevant subset of the corpus during inference.

A key limitation of standard RAG methods is that each text chunk is embedded independently, which can lead to incomplete retrieval when a query depends on information spread across multiple chunks. Because similarity scores are computed in isolation, essential context may be missed. To address this, more recent approaches have adopted structured representations — such as knowledge graphs — that explicitly model relationships between chunks across the corpus. At query time, these graphs are traversed or queried to retrieve semantically connected chunks, enabling LLMs to perform more effective multi-hop reasoning and question answering (Gutiérrez et al. 2025a,b; Edge et al. 2025; Guo et al. 2025).

These methods have primarily been evaluated on fact-rich documents such as Wikipedia pages (Yang et al. 2018; Ho et al. 2020; Trivedi et al. 2022). Yet **the vast majority of texts that LLMs encounter are not lists of facts but narratives of evolving real-world situations**. Crime reports, political briefings, corporate filings, legislative records, war dispatches, and multi-day news coverage all describe **actors** (people, organizations, nations) that adopt **roles** (suspect, regulator, bidder, combatant) and transition through **states** (arrested → arraigned → released; startup → unicorn → acquired) while interacting across **space and time**.

We contend that reasoning over such documents would be much more accurate and energy efficient, if one indexed the documents in terms of **an internal world model**— a structured representation that keeps track of *who* is involved, *what* was done, *where* and *when* events occur, *how* roles change, and *what* consequences follow. Indeed, to achieve such a goal, humans possess *episodic memory* (Tulving 1972, 2002) enabling us not only to plan and reason to seamlessly operate in the real world, but also to create new or update existing world models by reasoning across multiple experiences (Schacter, Addis, and Buckner 2007; Hassabis and Maguire 2007).

In this work, we introduce the **Generative Semantic Workspace** (GSW), a unifying computational framework for modeling world knowledge as structured, probabilistic semantics in the era of Large Language Models (LLMs). GSW formalizes how an intelligent agent—human or arti-

ficial—constructs and updates an internal representation of evolving situations from sequential input (e.g., text, video, or dialogue modalities). These representations are interpretable, actor-centric, and predictive: they reflect semantic regularities in the past while projecting likely future outcomes. GSW may be viewed as an instance of *episodic memory* that can be integrated into LLM-based systems as a reasoning and memory module, serving as a symbolic bridge between language and latent world models.

To illustrate how GSW can help LLMs reason accurately, we evaluate it on the Episodic Memory Benchmark (EpBench) (Huet, Houidi, and Rossi 2025), that has recently been introduced as a way to benchmark the episodic memory-like capabilities of LLMs. Following are excerpts from two different documents that relate to an entity, Carter Stewart, in this EpBench dataset:

**Document #1:** The imposing structure loomed before him, its grand facade a testament to both artistry and scientific achievement ...... As he stepped into the **Metropolitan Museum of Art**, the echoing chatter of excited voices ...... The antique clock in the main hall chimed, its resonant tones reminding him of the date: **September 22, 2026** .... found himself particularly engrossed during the third presentation, where **Carter Stewart** explained statistical analysis with a clarity that left the audience spellbound.

**Document #2:** The air crackled with tension as **Carter Stewart** stepped onto the pristine greens of **Bethpage Black Course** on **March 23, 2024** ...... Carter discussed implications of research, his fingers trembling slightly as he adjusted his microphone.

An agent reading the narrative in the first document faces a fundamentally different challenge than traditional fact retrieval. It must understand that "he" refers to a nameless protagonist, who attended a scientific conference where Carter Stewart spoke. The narrator's spatial context (Metropolitan Museum of Art) and temporal context (September 22, 2026), are stated only indirectly and more importantly have to be also assigned to Carter Stewart who is a presenter. GSW is able to create such representations as part of its working memory construction task: "Carter Stewart: **Role:** A presenter at a Scientific Conference; **Date:** September 22, 2026, morning session; **Location:** The Metropolitan Museum of Art, **Topic:** statistical analysis; **Implements Used:** presentation boards and holographic projectors." The second document is more straightforward and GSW creates a memory trace such as: "Carter Stewart; **Role:** a researcher and presenter; **Location:** Bethpage Black Course; **Date:** March 23, 2024, **Did What?:** Presented his research findings at a Scientific Conference." A visualization of the steps of how GSW constructs its working memory is shown in **Fig. 6-8 in Appendix B**.

When presented with a task such as "List all the unique locations and dates where Carter Stewart made presentations at Scientific Conference events." a query resolution module (see Section ) searches through the GSW constructed from all 200 documents and identifies entities mentioned in the query (e.g., Carter Stewart) that match query's intent (e.g., a presenter at scientific conference; another entity named
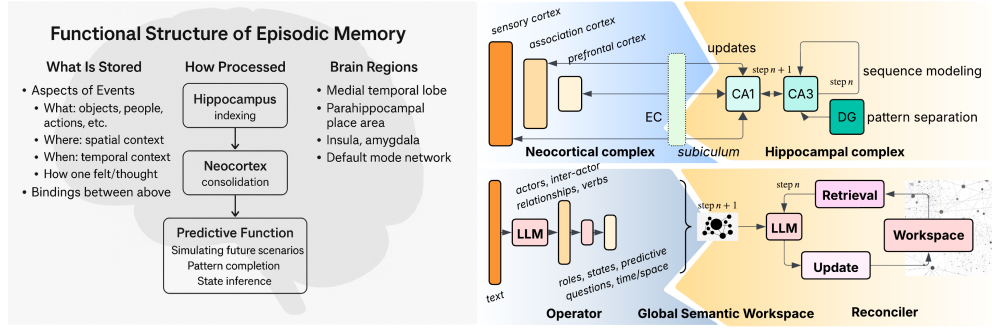
Carter Stewart whose role is that of a baker by profession would be ignored) and then returns just the relevant portion of its memory, as in the preceding paragraph. This results in highly targeted and short texts that an LLM has to reason through to provide an answer. In contrast, current structured RAG methods are designed to facilitate retrieval of either whole chunks or community-level summaries that have different levels of similarity to the entities and other phrases in the query. For example, for this query (see **Appendix** D) GraphRAG's (Edge et al. 2025) summarization missed that Carter Stewart was at the same location as the protagonist in Document #1, and included irrelevant text chunks which led to a list that misses one location and hallucinates two erroneous locations. HippoRAG2 (Gutiérrez et al. 2025b) retrieves the full text of both the relevant documents, along with many other documents, overwhelming the LLM and leading it to hallucinate three erroneous locations. For a more detailed comparison, see Section , and Tables 1, 4, and 3.

In the rest of this paper, we detail the GSW framework (**Section** ) and present a rigorous evaluation on two versions of the EpBench benchmark (**Section** ). The results demonstrate a significant improvement over existing methods. On the EpBench-200 corpus, GSW achieves a state-of-the-art F1-score of 0.85, outperforming strong structured RAG baselines. This advantage is particularly pronounced in the most demanding queries requiring synthesis across as many as 17 different documents, where GSW improves recall by up to **20%** over the next best approach as detailed in Table 2. Furthermore, GSW is efficient, reducing the number of context tokens sent to the LLM by **51%** compared to the most token-efficient baseline, drastically lowers inference costs and reducing the rate of hallucination in question answering (see Table 3). We further show that this powerful combination of accuracy and token efficiency holds at scale; on the EpBench-2000 corpus, a 10x larger dataset, GSW again achieves a state-of-the-art F1-score of 0.773, outperforming the best baseline by more than **15%** on overall recall (Table 4), positioning GSW as a robust and scalable solution for equipping LLMs with effective episodic memory.

Results and discussions are summarized in **Section** and a review of related literature is presented **Section** . Finally, limitations and future work are discussed in Section , and a detailed **Appendix** provides supporting evidence, including manual evaluations performed to validate the power of GSW's episodic memory capabilities.

## The Generative Semantic Workspace (GSW) Framework

In neuroscience, the neocortex is believed to encode hierarchical abstractions of entities, roles, and event templates (George and Hawkins 2009; Botvinick 2008; Felleman and Van Essen 1991). The hippocampus, especially the CA3 module, plays a complementary role by binding these representations into coherent spatiotemporal sequences (Teyler and DiScenna 1986; Rolls 2013; Eichenbaum 2004). During sleep, this neocortical-hippocampal system engages in *experience replay*, a process through which episodic traces are reactivated in reverse or forward order to consolidate memory and refine

Figure 1: **Unifying Brain-Inspired and Generative Semantics for Episodic Memory Modeling** The hippocampal complex (DG, CA3, CA1) and neocortical regions (NC) inspire the *Reconciler* (retrieval, workspace, update) and *Operator* (LLM-driven semantic role extraction), respectively. The neocortical complex, responsible for context-rich consolidation and predictive modeling, aligns with the Operator module's functions. The hippocampal complex, which performs indexing, pattern separation, and sequence modeling, corresponds to the Reconciler. Together, the GSW framework offers a biologically inspired, interpretable model for simulating world knowledge from text inputs.

internal models (Ólafsdóttir, Bush, and Barry 2018; Louie and Wilson 2001; Wilson and McNaughton 1994). This back and forth supports both persistence and prediction of memory (McClelland, McNaughton, and O'Reilly 1995; Rasch and Born 2013), key features of episodic memory.

Motivated by this biological architecture (see Fig 1), an effective memory framework requires a **structured representation** capable of encoding actors along with their evolving roles and states. Crucially, this representation must be capable of spatiotemporal grounding, linking entities and their interactions to specific times and locations, much like the binding function of the hippocampus. Finally, the framework must possess a process for **consolidating and updating these structures** as new information arrives, mirroring the way the neocortical-hippocampal loop constantly refines its world model. This process of building an evolving model is illustrated with a detailed end-to-end example in **Appendix B**.

**From Episodic Memory to Generative Modeling of Situations and Narratives:** The central challenge, therefore, is to create a continuously evolving semantic model, which requires a bidirectional mapping between text and a structured representation. While early symbolic frameworks like Prop-Bank (Kingsbury and Palmer 2002) and FrameNet (Baker, Fillmore, and Lowe 1998) attempted this, they were not designed for this full bidirectional process, relying instead on fixed ontologies that lacked the necessary probabilistic and dynamic interpretation.

*LLMs now make this bidirectional mapping tractable.* They can both infer concise semantic identifiers from text and generate coherent narratives from those identifiers. This enables a new, efficient memory model where compact semantic traces are stored and reactivated in context. The formal model is presented next, and *its approach is validated in **Appendix I**, where a human evaluation shows a strong preference for the GSW semantic maps over those from frameworks like PropBank and FrameNet.*

## A Probabilistic Model for Semantic Memory: The Operator Framework

We now define a minimal schema for encoding these semantic elements—along with predictive cues, spatiotemporal attributes, and utilities—that serves as the foundation of the GSW framework for structured memory in LLMs. The agent must distill and maintain a semantic map from text to build a coherent semantic model.

To make this concrete, let's consider a single text input $C_n$ at some time step $n$ : *Yesterday, in a swift response to a reported robbery, law enforcement officers apprehended Jonathan Miller, a 32-year-old resident of Greenview Avenue, in the downtown area.*

Explicit information in $C_n$ typically specifies a configuration of participating actors $a_1, \ldots, a_K$ and the relations or interactions among them. The agent must distill and maintain a semantic map from these clues to build a coherent semantic model. Let's represent this interaction pattern at time step $n$ as (here each entry denotes an interaction from actor $a_i$ to $a_j$ as inferred from $C_n$):

$$
\mathcal{C}_n \approx \begin{pmatrix} (a_1 \to a_1)^n & \cdots & (a_1 \to a_K)^n \\ \vdots & \ddots & \vdots \\ (a_K \to a_1)^n & \cdots & (a_K \to a_K)^n \end{pmatrix};
$$

### Actors, Roles and States

The word 'Miller', in isolation, corresponds to a broad, unconditioned distribution over possible behaviors of a human. If 'Miller' is likely to commit a crime, the agent would probably refer to Miller with a label 'Criminal'. We call these labels *roles*.

**Role:** An identifier that specifies a distribution over potential actions that an actor $a_i \in \mathcal{A}$ may take toward other actors $a_j \in \mathcal{A}$:

$$
\pi_r : \mathcal{A} \times \mathcal{A} \to [0, 1] \tag{1}
$$

where $\pi_r(a_i \to a_j)$ denotes the probability of $a_i$ acting on $a_j$ in role $r$. For example, assigning the role of 'criminal'

to Miller increases the *likelihood* that he will engage in actions such as *committing a crime* against another actor or increasing the chances that Miller will *attempt to flee* from 'law enforcement'.

The agent would also *know* that in addition to Miller being a *criminal*, Miller has been *caught*. Or perhaps he *escaped*. We call these labels *states*.

**State:** An identifier that induces a contextual attribute that modulates the probability distribution over actions available to an actor within a given role. Given an actor $a_i$ with role $r$, a state $s \in \mathcal{S}_r$ constrains the role-induced action distribution $\pi_r$:

$$\pi_{r,s}(a_i \to a_j) = \pi_r(a_i \to a_j \mid s), \qquad (2)$$

where $\pi_{r,s}$ denotes the subset of actions available to actor $a_i$ in state $s$. For instance, a *criminal* in the state *captured* may be limited to passive or compliant interactions, precluding actions such as fleeing or committing further crimes. Thus, states act as dynamic modifiers of an actor's interaction profile within a given situation.

**Verbs and Valences**

Verbs encode structured semantic attributes helping the agent to structure an event by drawing on prior experience, as verbs tend to generalize across contexts more reliably than nouns. They provide causal certificates for roles/states of actors. For example, understanding why Miller transitions from being *free* to *captured* relies on identifying the underlying interaction – such as being arrested – that bridges those states. A verb's valences are efficient means of capturing information needed for reasoning about future outcomes. Verbs can be modeled similar to roles and states:

$$v(a_i \to a_j) : \mathcal{A} \times \mathcal{A} \to \mathcal{L}_v, \qquad (3)$$

where the valences $\ell_k \in \mathcal{L}_v$ signal the change in roles and states of the actors interacting via the verb. When Miller is running from the police, the *next* state for Miller might be *escaped* or *caught*: a distribution of potential *future* roles and states.

**Time and Space Continuity**

Spatiotemporal continuity constraints are crucial to capture world models, not only for individual actors but especially as interactions/verbs couple their coordinates. For instance, if Officers are actively apprehending Johnathan Miller in the Downtown area, then it enforces a shared location and time among the actors. Moreover, if the next day Miller is found in a city a thousand miles away, it would constrain his unobserved action to that of having flown and lead the agent to narrow down events that could have led to such a spatial shift. In effect, the flow of time and space regularizes the semantic map, biasing verb selection toward contextually coherent transitions. If the position information derived from $\mathcal{C}_n$ at time step $n$ is $\mathcal{X}_n$ and the temporal information is $\mathcal{T}_n$, then:

**Temporal continuity:** $\mathcal{T}_{n+1} - \mathcal{T}_n$ must be consistent with the expected temporal scope of $v$,
**Spatial proximity:** $\|\mathcal{X}_n(a_i) - \mathcal{X}_n(a_j)\|$ must fall within a valid range for the verb (e.g., *tackle* requires physical closeness)

**Forward-Falling Questions to Capture Potential Outcomes and Actions**

The collection of roles/states, verbs, and spatiotemporal coordinates constrain the space of future progression and can be efficiently encoded as a set of questions $\mathcal{Q}_n$. For example, given that Miller has been arrested, "When would Miller be indicted," "where and when would the trial happen?" "Will he be free on bail?" A prosecutor agent, for example, would need to start strategizing about such potential outcomes.

A complete workspace instance can be written as a sampled distribution from an underlying "Workspace" generative process:

$$\mathcal{M}_n \sim p(\mathcal{A}, \mathcal{R}, \mathcal{S}, \mathcal{V}, \mathcal{T}, \mathcal{X}, \mathcal{Q} \mid \mathcal{C}_{0:n}) \qquad (4)$$

where $\mathcal{M}_n \mapsto q(\mathcal{M}_{n+1} \mid \mathcal{M}_n)$ models the likelihood of generating the next workspace instance.

**Enabling Recursive Updates: A State Space Approach (The Reconciler Framework)**

Given a single text input $\mathcal{C}_0$, GSW models the workspace instance $\mathcal{M}_0$ as $P(\mathcal{M}_0|\mathcal{C}_0)$. We seek to compute: $P(\mathcal{M}_n|\mathcal{C}_{0:n})$. For $\mathcal{M}_1$, we introduce $\mathcal{W}_1$, an intermediate representation to decompose $P(\mathcal{M}_1|\mathcal{C}_0, \mathcal{C}_1)$ into parts:

$$\begin{aligned} P(\mathcal{M}_1 \mid \mathcal{C}_0, \mathcal{C}_1) \\ = \sum_{\mathcal{M}_0, \mathcal{W}_1} P(\mathcal{M}_1 \mid \mathcal{M}_0, \mathcal{W}_1) \\ \times P(\mathcal{M}_0 \mid \mathcal{C}_0) P(\mathcal{W}_1 \mid \mathcal{C}_1) \qquad (5) \end{aligned}$$

Here, we assume conditional independence between the workspace state $\mathcal{M}_0$ and the intermediate representation $\mathcal{W}_1$ given the context sequence, such that:

$$\begin{aligned} P(\mathcal{M}_0, \mathcal{W}_1 \mid \mathcal{C}_0, \mathcal{C}_1) \\ = P(\mathcal{M}_0 \mid \mathcal{C}_0) P(\mathcal{W}_1 \mid \mathcal{C}_1) \qquad (6) \end{aligned}$$

where we define $\mathcal{W}_1$ to depend solely on the current context $\mathcal{C}_1$, and $\mathcal{M}_0$ solely on the initial context $\mathcal{C}_0$. For an arbitrary step $n$:

$$\begin{aligned} P(\mathcal{M}_n|\mathcal{C}_{0:n}) = \sum_{\mathcal{M}_{n-1}, \mathcal{W}_n} P(\mathcal{M}_n|\mathcal{M}_{n-1}, \mathcal{W}_n) \\ \times P(\mathcal{M}_{n-1}|\mathcal{C}_{0:(n-1)}) P(\mathcal{W}_n|\mathcal{C}_n) \qquad (7) \end{aligned}$$

Estimating a workspace instance $\mathcal{M}_n$ involves learning parameterized models for three components: the transition model, the prior workspace, and the context-derived augmentation. The prior workspace $\mathcal{M}_{n-1}$ is recursively computed from previous steps. The augmentation step produces an intermediate representation of the current context $\mathcal{C}_n$. We refer to the model estimating this distribution as the **Operator**. The transition model uses a Markovian assumption to produce the updated workspace instance by reconciling existing workspace semantic maps with new semantic information. We refer to this module as the **Reconciler**. Together, the Operator and Reconciler implement a sequential inference mechanism where the Operator maps each new context $\mathcal{C}_n$ to an intermediate state $\mathcal{W}_n$, and the Reconciler performs a structured update $\mathcal{M}_{n-1} \to \mathcal{M}_n$.

## Question Answering with GSW

Figure 2 illustrates this process: memory construction via Operator and Reconciler modules, followed by retrieval, reranking and QA. As described in the caption, once a working memory instance is constructed, answering a query involves the following steps: the system first matches entities from the query to the GSW, then generates contextual summaries for those matched entities from the workspace, re-ranks the summaries for relevance, and finally passes the top-ranked summaries to an LLM to synthesize the answer.

## EpBench: An Episodic Memory Benchmark

Our experiments utilize the Episodic Memory Benchmark (EpBench) (Huet, Houidi, and Rossi 2025), a benchmark specifically designed to evaluate the capabilities of LLMs for episodic memory recall and reasoning over long narratives. Unlike many standard Question Answering (QA) benchmarks (Kočiskỳ et al. 2018; Zhang et al. 2024; Yang et al. 2018) – focusing on localized factual retrieval – EpBench targets core episodic capabilities: remembering specific events situated in unique spatiotemporal contexts and distinguishing between recurring events involving the same actors.

| Statistic | Value |
|---|---|
| Number of Chapters | 200 |
| Total Tokens | 102,870 |
| Total Queries (QA Pairs) | 686 |
| Queries by Event Category | |
| (0 / 1 / 2 / 3-5 / 6+ Cues) | 180 / 180 / 108 / 128 / 90 |
| Max. Chapters Referenced per Query | 17 |
| Min. Chapters Referenced per Query | 0 |

Table 1: **EpBench-200 Dataset Statistics.** The statistics for around 1M token Epbench-2000 are presented in **Appendix F**

EpBench documents are structured as synthetic books generated chapter-by-chapter from event templates (detailing date, location, entity, content) sampled from a larger universe, ensuring recurring elements that necessitate disambiguation and temporal tracking. Chapters are generated via LLM prompts and verified for coherence. Moreover, the same time/location/actors (collectively referred to as cues) appear across multiple chapters. For our evaluation, we use both the standard 200-chapter version and the extended 2000 chapter version of the dataset and report its Statistics in Table 1 **Appendix F**.

## Evaluation Metrics

To evaluate model performance on the EpBench dataset's queries (detailed in Section ), we adopt the LLM-as-a-Judge evaluation paradigm (Zheng et al. 2023). For consistency, we strictly follow the LLM-based answer processing and extraction procedure outlined by the EpBench benchmark authors. This approach accounts for the possibility that model responses might be longer or more elaborate than the typically concise ground truth answers. These LLM extracted answers are then used to compute Precision, Recall and F1 scores which we report in Table 2

## Baseline Methods

We compare GSW against several baseline approaches: **Vanilla LLM**, standard **Embedding-based RAG** (Karpukhin et al. 2020; Ram et al. 2023) for which we utilized the **Voyage-03**[1] embedding model selected for its strong performance on retrieval benchmarks (Thakur et al. 2021; Muennighoff et al. 2022) , and the structured RAG methods **GraphRAG** (Edge et al. 2025), **HippoRAG2**(Gutiérrez et al. 2025b), and **LightRAG** (Guo et al. 2025). We detail the hyper-parameter settings for all baselines in **Appendix E** .

## Implementation Details

The GSW **Operator** (Section ) and **Reconciler** (Section ) were implemented by prompting GPT-4o (Hurst et al. 2024) according to task-specific instructions, using temperature set to 0 for deterministic behavior. To ensure fair comparison, we standardized both the maximum context utilization (limited to 17 chapters per query, matching the maximum relevant chapters per query) and the answer generation model (GPT-4o) across all evaluated methods. The complete prompts are provided in **Appendix A**, and API interactions were managed using the Bespoke Curator library (Marten et al. 2025), indexing costs are reported in **Appendix H**. To generate an answer for a given query, we first identify named entities within the query text. These entities are then matched to corresponding nodes within the current GSW memory ($\mathcal{M}_n$) using simple string matching. Summaries for the matched entities – aggregated from the GSW structure – are then retrieved and re-ranked based on semantic similarity to the query. The final re-ranked summaries are provided to the LLM to answer the query as illustrated in Figure 2. An end-to-end QA example is provided in **Appendix C**.

## Results and Discussion

**QA Performance:** Table 2 presents a comparative analysis of GSW against the baseline methods detailed in Section  across Precision (P), Recall (R), and F1-Score (F1) metrics, categorized by the number of matching cues per query. Across the aggregated metrics, GSW achieves the highest overall F1-Score (0.850), Precision (0.865), and Recall (0.894), improving overall metrics by more than **10%** over the next-best method. GSW also demonstrates consistent performance across the various Cue categories, achieving the highest score in **16 out of 18** individual metric computations, and ranking second in the remaining two, highlighting its robust performance across varying levels of episodic recall complexity. Particularly noteworthy is GSW's performance in the '6+ Cues' category. *This is the most demanding scenario*, where correct responses can require reasoning across information spanning up to 17 distinct chapters (see Table 1). Even in this complex setting, GSW demonstrates robust efficacy and achieves the highest performance over all metrics:

---

[1]https://blog.voyageai.com/2024/09/18/voyage-3/

[2]Cost calculated using GPT-4o pricing of $2.50 per million tokens.
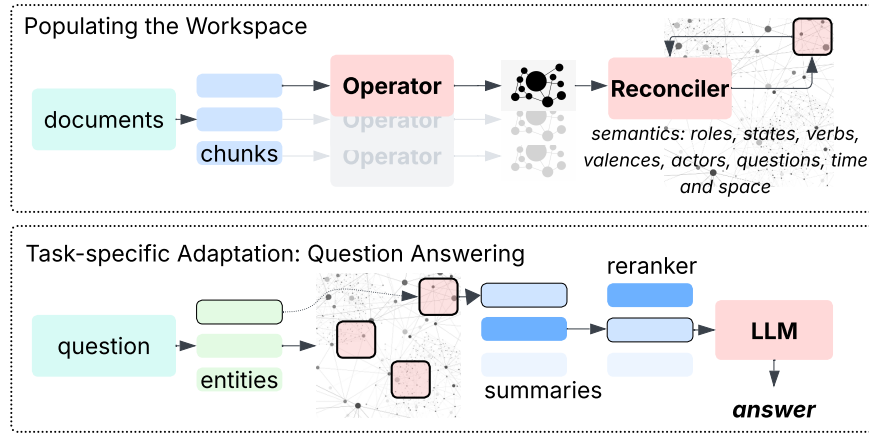
Figure 2: **Episodic Memory Creation and QA:** Figure illustrates the end-to-end process of constructing a workspace and question answering from the workspace. *(top)* Large-scale text is segmented into semantically coherent chunks. Each chunk is processed by the *Operator* model to generate a local workspace instance, represented as a semantic graph. These instances are incrementally integrated by the *Reconciler* resulting in a unified Global Memory. *(bottom)* During question answering, the system retrieves relevant portions of this memory by matching named entities in the query to identifiers in the semantic network. For each match, it reconstructs episodic summaries—contextual recreations of past situations—which are re-ranked and passed to an LLM to generate the final answer.

F1:0.834 P:0.891, R:0.822. In particular when compared to HippoRAG2, next most performant in this category, GSW outperforms it by approximately **20%** in recall. *Recall, in particular, measures a framework's ability to map queries to the correct chapter and context*, and it is revealing that for all competing frameworks recall decreases as the number of matching cues increases, whereas the GSW maintains consistently strong performance, highlighting the strength of its structured representation in storing episodic information. Finally, the Vanilla LLM is consistently the poorest performing baseline (e.g overall F1 Score of 0.642) reaffirming the inherent difficulty of the episodic QA task and the necessity of specialized memory frameworks like the GSW.

**Scalability on EpBench-2000**: To assess the scalability of our method, we evaluate GSW on the EpBench-2000 dataset, which increases the corpus size by 10 fold. The results, presented in Table 4, show that GSW maintains its performance lead by achieving an overall F1-score of 0.773, which is **15% higher** than the strongest baseline (embedding RAG), and **22% higher** than other structured RAG methods. Thus, GSW's advantages in recall and reasoning persist even at a significantly larger scale. Due to space constraints, the full breakdown table by cue category is provided in **Appendix F**.

**Token Efficiency:** Beyond query performance, GSW demonstrates substantial improvements in token efficiency, as detailed in Table 3, which presents the average number of context tokens supplied to the LLM per query, and the corresponding cost for all compared methods. GSW achieves a remarkable **51%** reduction in token usage when compared to the next most token-efficient baseline (GraphRAG). This advantage is even more pronounced when compared to stronger performing baselines such as Embedding RAG and HippoRAG2, against which GSW offers a token reduction of

nearly **59%**. GSW's efficient approach to query resolution contributes to the reduction in token count: Rather than passing entire chapters or raw document chunks, GSW utilizes its semantic structure to generate entity-specific summaries (Prompt in **Appendix A**), thereby providing only targeted query-specific information to the LLM as illustrated in **Appendix C**. This focused contextual information also reduces hallucinations as supported by the GSW's leading performance in the '0 Cues' category, where no matching cues are present in the source document.

Several additional **ablation studies** are presented in **Appendix F**, including the removal of identifier types (e.g., temporal and spatial tags), evaluations on a shortened version of the EpBench dataset, and comparisons across different retrieval strategies. These experiments highlight the contribution of each component in the GSW architecture and underscore the importance of principled memory querying. For qualitative insights into GSW's behavior and outputs, see **Appendix D**.

## Related Work

The relevant literature has been discussed in the Introduction, and a detailed literature review is included in **Appendix G**. To summarize, Retrieval-Augmented Generation (RAG) (Lewis et al. 2021; Gao et al. 2024; Karpukhin et al. 2020) retrieves relevant chunks from indexed documents using dense (Devlin et al. 2019; Reimers and Gurevych 2019; Lee et al. 2025), sparse (Robertson and Zaragoza 2009), or hybrid (Cormack, Clarke, and Buettcher 2009) embeddings. While effective for fact-based QA, standard RAG struggles to connect dispersed information due to its reliance on chunk-based retrieval (Chen et al. 2023; Merola and Singh 2025). Structured approaches like GraphRAG(Edge et al. 2025), LightRAG(Guo et al. 2025) and HippoRAG(Gutiérrez et al.

| Metric | Method | 0 Cues (N=180) | 1 Cue (N=180) | 2 Cues (N=108) | 3-5 Cues (N=128) | 6+ Cues (N=90) | Overall (N=686) |
|---|---|---|---|---|---|---|---|
| **P** | Vanilla LLM | $0.840 \pm 0.019$ | $0.734 \pm 0.021$ | $0.735 \pm 0.026$ | $0.703 \pm 0.021$ | $0.806 \pm 0.028$ | $0.766 \pm 0.010$ |
| | Embedding RAG | $0.906 \pm 0.021$ | $\underline{0.745} \pm 0.026$ | $0.803 \pm 0.028$ | $0.823 \pm 0.025$ | $0.886 \pm 0.029$ | $\underline{0.832} \pm 0.012$ |
| | GraphRAG (Edge et al. 2025) | $\underline{0.950} \pm 0.016$ | $0.657 \pm 0.029$ | $0.677 \pm 0.034$ | $0.753 \pm 0.028$ | $0.816 \pm 0.035$ | $0.781 \pm 0.013$ |
| | HippoRAG2 (Gutiérrez et al. 2025b) | $0.829 \pm 0.027$ | $0.704 \pm 0.029$ | $\mathbf{0.817} \pm 0.026$ | $\underline{0.839} \pm 0.026$ | $\mathbf{0.940} \pm 0.020$ | $0.812 \pm 0.013$ |
| | LightRAG (Guo et al. 2025) | $0.946 \pm 0.017$ | $0.668 \pm 0.029$ | $0.615 \pm 0.036$ | $0.695 \pm 0.031$ | $0.822 \pm 0.037$ | $0.763 \pm 0.014$ |
| | GSW (Ours) | $\mathbf{0.978} \pm 0.011$ | $\mathbf{0.755} \pm 0.026$ | $\underline{0.810} \pm 0.027$ | $\mathbf{0.878} \pm 0.019$ | $\underline{0.890} \pm 0.024$ | $\mathbf{0.865} \pm 0.010$ |
| **R** | Vanilla LLM | $0.840 \pm 0.019$ | $0.781 \pm 0.021$ | $0.526 \pm 0.021$ | $0.419 \pm 0.017$ | $0.229 \pm 0.014$ | $0.616 \pm 0.011$ |
| | Embedding RAG | $0.906 \pm 0.021$ | $\underline{0.863} \pm 0.025$ | $0.773 \pm 0.033$ | $0.746 \pm 0.027$ | $0.624 \pm 0.036$ | $\underline{0.807} \pm 0.012$ |
| | GraphRAG (Edge et al. 2025) | $\underline{0.950} \pm 0.016$ | $0.764 \pm 0.031$ | $0.686 \pm 0.035$ | $0.645 \pm 0.026$ | $0.537 \pm 0.030$ | $0.748 \pm 0.014$ |
| | HippoRAG2 (Gutiérrez et al. 2025b) | $0.829 \pm 0.027$ | $0.823 \pm 0.026$ | $\underline{0.800} \pm 0.029$ | $\underline{0.749} \pm 0.026$ | $\underline{0.675} \pm 0.030$ | $0.787 \pm 0.013$ |
| | LightRAG (Guo et al. 2025) | $0.946 \pm 0.017$ | $0.716 \pm 0.033$ | $0.628 \pm 0.035$ | $0.559 \pm 0.029$ | $0.458 \pm 0.029$ | $0.699 \pm 0.015$ |
| | GSW (Ours) | $\mathbf{0.978} \pm 0.011$ | $\mathbf{0.863} \pm 0.025$ | $\mathbf{0.869} \pm 0.023$ | $\mathbf{0.893} \pm 0.015$ | $\mathbf{0.822} \pm 0.022$ | $\mathbf{0.894} \pm 0.009$ |
| **F1** | Vanilla LLM | $0.840 \pm 0.019$ | $0.709 \pm 0.022$ | $0.585 \pm 0.021$ | $0.476 \pm 0.017$ | $0.325 \pm 0.017$ | $0.629 \pm 0.010$ |
| | Embedding RAG | $0.906 \pm 0.021$ | $\underline{0.726} \pm 0.026$ | $0.723 \pm 0.030$ | $0.745 \pm 0.026$ | $0.680 \pm 0.035$ | $\underline{0.771} \pm 0.013$ |
| | GraphRAG (Edge et al. 2025) | $\underline{0.950} \pm 0.016$ | $0.625 \pm 0.029$ | $0.625 \pm 0.034$ | $0.657 \pm 0.026$ | $0.607 \pm 0.032$ | $0.714 \pm 0.013$ |
| | HippoRAG2 (Gutiérrez et al. 2025b) | $0.829 \pm 0.028$ | $0.676 \pm 0.028$ | $\underline{0.762} \pm 0.028$ | $\underline{0.754} \pm 0.025$ | $\underline{0.746} \pm 0.027$ | $0.753 \pm 0.013$ |
| | LightRAG (Guo et al. 2025) | $0.946 \pm 0.017$ | $0.594 \pm 0.030$ | $0.587 \pm 0.032$ | $0.579 \pm 0.028$ | $0.561 \pm 0.030$ | $0.678 \pm 0.014$ |
| | GSW (Ours) | $\mathbf{0.978} \pm 0.011$ | $\mathbf{0.744} \pm 0.026$ | $\mathbf{0.807} \pm 0.024$ | $\mathbf{0.868} \pm 0.016$ | $\mathbf{0.834} \pm 0.022$ | $\mathbf{0.850} \pm 0.010$ |

Table 2: **GSW performance on Epbench-200 (200-Chapters Book)** Performance is grouped by metric (Precision, Recall, F1-Score) across different numbers of matching cues per query. (N=X) indicates questions per category. Error bars are estimated via bootstrap resampling. Best score in each column for each metric group is **bold**; second best is underlined.

| Method | Avg. Tokens | Avg. Cost[2] |
|---|---|---|
| Vanilla LLM | $\sim101{,}120$ | $\sim\$0.2528$ |
| Embedding RAG | $\sim8{,}771$ | $\sim\$0.0219$ |
| GraphRAG (Edge et al. 2025) | $\sim\underline{7{,}340}$ | $\sim\underline{\$0.0184}$ |
| HippoRAG2 (Gutiérrez et al. 2025b) | $\sim8{,}771$ | $\sim\$0.0219$ |
| LightRAG (Guo et al. 2025) | $\sim40{,}476$ | $\sim\$0.1012$ |
| GSW (Ours) | $\sim\mathbf{3{,}587}$ | $\sim\mathbf{\$0.0090}$ |

Table 3: **GSW's Efficiency**: Average context tokens passed to the LLM per query on EpBench-200, and the estimated cost to answer that query. GSW achieves the best performance (detailed in Table 2) with the significantly lowest token count and cost, as highlighted below. Best score in each column is **bold**; second best is underlined.

| Method | Precision | Recall | F1 |
|---|---|---|---|
| Embedding RAG | $\underline{0.827} \pm 0.014$ | $\underline{0.688} \pm 0.015$ | $\underline{0.675} \pm 0.015$ |
| GraphRAG | $0.761 \pm 0.017$ | $0.548 \pm 0.017$ | $0.544 \pm 0.017$ |
| HippoRAG2 | $0.759 \pm 0.016$ | $0.648 \pm 0.016$ | $0.635 \pm 0.015$ |
| LightRAG | $0.649 \pm 0.018$ | $0.497 \pm 0.017$ | $0.494 \pm 0.016$ |
| GSW (Ours) | $\mathbf{0.830} \pm 0.010$ | $\mathbf{0.796} \pm 0.009$ | $\mathbf{0.773} \pm 0.009$ |

Table 4: **Overall performance on Epbench-2000 (2000-Chapters Book).** The same convention as in Table 2 is followed. For a more descriptive full table please refer to **Appendix F**.

interpretable alternative to long-context or retrieval-based systems.

Nevertheless, we identify key limitations and avenues for future work. Firstly, GSW's evaluation, while utilizing Ep-Bench for its strengths in spatiotemporal assessment, is constrained by the limited scope of current episodic memory benchmarks in thoroughly probing the complex evolution of actor roles and states within extended narratives; we are developing a more comprehensive benchmark to specifically address this gap. Secondly, the present GSW implementation relies on a strong closed-source LLM (GPT-4o). Empirical validation of promising open-source alternatives (Yang et al. 2024; Grattafiori et al. 2024; Team et al. 2025) within our Operator-Reconciler architecture is therefore essential. Expanding GSW to process and integrate information from diverse data modalities beyond text also presents an important direction for future development, aiming to broaden its applicability to more complex, real-world experiential inputs.

2025a,b) mitigate this by modeling relationships and supporting multi-hop reasoning.

## Concluding Remarks and Limitations

In this work, we introduced the Generative Semantic Workspace (GSW) as a framework for equipping LLMs with human-like episodic memory. Its two core components—the Operator, which interprets local semantics within short context windows, and the Reconciler, which integrates and updates these representations over time—jointly construct a persistent, structured memory. This memory maps raw text into evolving configurations of roles, states, and interactions within a coherent global workspace. On the Episodic Memory Benchmark, GSW outperforms existing approaches in both accuracy and token efficiency, offering a scalable and

# References

Aguilar, J.; Beller, C.; McNamee, P.; Van Durme, B.; Strassel, S.; Song, Z.; and Ellis, J. 2014. A Comparison of the Events and Relations Across ACE, ERE, TAC-KBP, and FrameNet Annotation Standards. In Mitamura, T.; Hovy, E.; and Palmer, M., eds., *Proceedings of the Second Workshop on EVENTS: Definition, Detection, Coreference, and Representation*, 45–53. Baltimore, Maryland, USA: Association for Computational Linguistics.

Baker, C. F. 2017. Framenet: Frame semantic annotation in practice. *Handbook of Linguistic Annotation*, 771–811.

Baker, C. F.; Fillmore, C. J.; and Lowe, J. B. 1998. The berkeley framenet project. In *COLING 1998 Volume 1: The 17th International Conference on Computational Linguistics*.

Botvinick, M. 2008. Hierarchical models of behavior and prefrontal function. *Trends in cognitive sciences*, 12: 201–8.

Chanin, D. 2023. Open-source Frame Semantic Parsing. *arXiv preprint arXiv:2303.12788*.

Chen, H.; Pasunuru, R.; Weston, J.; and Celikyilmaz, A. 2023. Walking Down the Memory Maze: Beyond Context Limit through Interactive Reading. ArXiv:2310.05029 [cs].

Cormack, G. V.; Clarke, C. L.; and Buettcher, S. 2009. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, 758–759.

Das, P.; Chaudhury, S.; Nelson, E.; Melnyk, I.; Swaminathan, S.; Dai, S.; Lozano, A.; Kollias, G.; Chenthamarakshan, V.; Jiří; Navrátil; Dan, S.; and Chen, P.-Y. 2024. Larimar: Large Language Models with Episodic Memory Control. ArXiv:2403.11901 [cs].

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv:1810.04805.

Doddington, G. R.; Mitchell, A.; Przybocki, M. A.; Ramshaw, L. A.; Strassel, S. M.; and Weischedel, R. M. 2004. The automatic content extraction (ace) program-tasks, data, and evaluation. In *Lrec*, volume 2, 837–840. Lisbon.

Edge, D.; Trinh, H.; Cheng, N.; Bradley, J.; Chao, A.; Mody, A.; Truitt, S.; Metropolitansky, D.; Ness, R. O.; and Larson, J. 2025. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. ArXiv:2404.16130 [cs].

Eichenbaum, H. 2000. A cortical–hippocampal system for declarative memory. *Nature Reviews Neuroscience*, 1(1): 41–50. Publisher: Nature Publishing Group.

Eichenbaum, H. 2004. Hippocampus: cognitive processes and neural representations that underlie declarative memory. *Neuron*, 44(1): 109–120.

Felleman, D. J.; and Van Essen, D. C. 1991. Distributed hierarchical processing in the primate cerebral cortex. *Cerebral cortex (New York, NY: 1991)*, 1(1): 1–47.

Fountas, Z.; Benfeghoul, M.; Oomerjee, A.; Christopoulou, F.; Lampouras, G.; Ammar, H. B.; and Wang, J. ????. Human-like Episodic Memory for Infinite Context LLMs.

Gao, Y.; Xiong, Y.; Gao, X.; Jia, K.; Pan, J.; Bi, Y.; Dai, Y.; Sun, J.; Wang, M.; and Wang, H. 2024. Retrieval-Augmented Generation for Large Language Models: A Survey. ArXiv:2312.10997 [cs].

George, D.; and Hawkins, J. 2009. Towards a mathematical theory of cortical micro-circuits. *PLoS computational biology*, 5(10): e1000532.

Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Guo, Z.; Xia, L.; Yu, Y.; Ao, T.; and Huang, C. 2025. LightRAG: Simple and Fast Retrieval-Augmented Generation. ArXiv:2410.05779 [cs].

Gutiérrez, B. J.; Shu, Y.; Gu, Y.; Yasunaga, M.; and Su, Y. 2025a. HippoRAG: Neurobiologically Inspired Long-Term Memory for Large Language Models. ArXiv:2405.14831 [cs].

Gutiérrez, B. J.; Shu, Y.; Qi, W.; Zhou, S.; and Su, Y. 2025b. From RAG to Memory: Non-Parametric Continual Learning for Large Language Models. ArXiv:2502.14802 [cs].

Hassabis, D.; and Maguire, E. A. 2007. Deconstructing episodic memory with construction. *Trends in cognitive sciences*, 11(7): 299–306.

Ho, X.; Duong Nguyen, A.-K.; Sugawara, S.; and Aizawa, A. 2020. Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, 6609–6625. Barcelona, Spain (Online): International Committee on Computational Linguistics.

Hong, K.; Troynikov, A.; and Huber, J. 2025. Context Rot: How Increasing Input Tokens Impacts LLM Performance. Technical report, Chroma.

Hsieh, C.-P.; Sun, S.; Kriman, S.; Acharya, S.; Rekesh, D.; Jia, F.; Zhang, Y.; and Ginsburg, B. 2024. RULER: What's the Real Context Size of Your Long-Context Language Models? ArXiv:2404.06654 [cs].

Huet, A.; Houidi, Z. B.; and Rossi, D. 2025. Episodic Memories Generation and Evaluation Benchmark for Large Language Models. ArXiv:2501.13121 [cs].

Hurst, A.; Lerer, A.; Goucher, A. P.; Perelman, A.; Ramesh, A.; Clark, A.; Ostrow, A.; Welihinda, A.; Hayes, A.; Radford, A.; et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Johansson, R.; and Nugues, P. 2008. Dependency-based semantic role labeling of PropBank. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, 69–78.

Karpukhin, V.; Oguz, B.; Min, S.; Lewis, P.; Wu, L.; Edunov, S.; Chen, D.; and Yih, W.-t. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 6769–6781. Online: Association for Computational Linguistics.

Kingsbury, P.; and Palmer, M. 2002. From TreeBank to Prop-Bank. In González Rodríguez, M.; and Suarez Araujo, C. P., eds., *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC'02)*. Las Palmas, Canary Islands - Spain: European Language Resources Association (ELRA).

Kočiský, T.; Schwarz, J.; Blunsom, P.; Dyer, C.; Hermann, K. M.; Melis, G.; and Grefenstette, E. 2018. The narrativeqa reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6: 317–328.

Lee, C.; Roy, R.; Xu, M.; Raiman, J.; Shoeybi, M.; Catanzaro, B.; and Ping, W. 2025. NV-Embed: Improved Techniques for Training LLMs as Generalist Embedding Models. ArXiv:2405.17428 [cs].

Leetaru, K.; and Schrodt, P. A. 2013. GDELT: Global data on events, location, and tone. *ISA Annual Convention*.

Leng, Q.; Portes, J.; Havens, S.; Zaharia, M.; and Carbin, M. 2024. Long context rag performance of large language models. *arXiv preprint arXiv:2411.03538*.

Levin, B. 1993. *English verb classes and alternations: A preliminary investigation*. University of Chicago press.

Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; Yih, W.-t.; Rocktäschel, T.; Riedel, S.; and Kiela, D. 2021. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. ArXiv:2005.11401 [cs].

Li, J.; Sun, A.; Han, J.; and Li, C. 2020. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1): 50–70.

Li, M.; Li, S.; Wang, Z.; Huang, L.; Cho, K.; Ji, H.; Han, J.; and Voss, C. 2021. The Future is not One-dimensional: Complex Event Schema Induction by Graph Modeling for Event Prediction. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 5203–5215. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.

Liu, J.; Chen, Y.; Liu, K.; Bi, W.; and Liu, X. 2020. Event extraction as machine reading comprehension. In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, 1641–1651.

Liu, N. F.; Lin, K.; Hewitt, J.; Paranjape, A.; Bevilacqua, M.; Petroni, F.; and Liang, P. 2023. Lost in the Middle: How Language Models Use Long Contexts. ArXiv:2307.03172 [cs].

Louie, K.; and Wilson, M. A. 2001. Temporally structured replay of awake hippocampal ensemble activity during rapid eye movement sleep. *Neuron*, 29(1): 145–156.

Lowe, J. B. 1997. A frame-semantic approach to semantic annotation. In *Tagging Text with Lexical Semantics: Why, What, and How?*

Lu, Y.; Lin, H.; Xu, J.; Han, X.; Tang, J.; Li, A.; Sun, L.; Liao, M.; and Chen, S. 2021. Text2Event: Controllable Sequence-to-Structure Generation for End-to-end Event Extraction. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2795–2806. Online: Association for Computational Linguistics.

Marten, R.; Vu, T.; Ji, C. C.-J.; Sharma, K.; Pimpalgaonkar, S.; Dimakis, A.; and Sathiamoorthy, M. 2025. Curator: A Tool for Synthetic Data Creation. https://github.com/bespokelabsai/curator.

McClelland, J. L.; McNaughton, B. L.; and O'Reilly, R. C. 1995. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological review*, 102(3): 419.

Mellor, A. 2017. The temporal event graph. *Journal of Complex Networks*, 6(4): 639–659.

Merola, C.; and Singh, J. 2025. Reconstructing Context: Evaluating Advanced Chunking Strategies for Retrieval-Augmented Generation. *arXiv preprint arXiv:2504.19754*.

Muennighoff, N.; Tazi, N.; Magne, L.; and Reimers, N. 2022. MTEB: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.

Ng, V.; and Cardie, C. 2002. Improving machine learning approaches to coreference resolution. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 104–111.

Nivre, J. 2010. Dependency parsing. *Language and Linguistics Compass*, 4(3): 138–152.

Ólafsdóttir, H. F.; Bush, D.; and Barry, C. 2018. The role of hippocampal replay in memory and planning. *Current Biology*, 28(1): R37–R50.

Palmer, M. 2009. Semlink: Linking propbank, verbnet and framenet. In *Proceedings of the generative lexicon conference*, 9–15. GenLex-09, Pisa, Italy.

Palmer, M.; Gildea, D.; and Xue, N. 2011. *Semantic role labeling*. Morgan & Claypool Publishers.

Ram, O.; Levine, Y.; Dalmedigos, I.; Muhlgay, D.; Shashua, A.; Leyton-Brown, K.; and Shoham, Y. 2023. In-Context Retrieval-Augmented Language Models. ArXiv:2302.00083 [cs].

Rasch, B.; and Born, J. 2013. About sleep's role in memory. *Physiological reviews*.

Reimers, N.; and Gurevych, I. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. ArXiv:1908.10084 [cs].

Robertson, S.; and Zaragoza, H. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends® in Information Retrieval*, 3(4): 333–389.

Rolls, E. T. 2013. A quantitative theory of the functions of the hippocampal CA3 network in memory. *Frontiers in cellular neuroscience*, 7: 98.

Sarthi, P.; Abdullah, S.; Tuli, A.; Khanna, S.; Goldie, A.; and Manning, C. D. 2024. RAPTOR: Recursive Abstractive Processing for Tree-Organized Retrieval. ArXiv:2401.18059 [cs].

Schacter, D. L.; Addis, D. R.; and Buckner, R. L. 2007. Remembering the past to imagine the future: the prospective brain. *Nature reviews neuroscience*, 8(9): 657–661.

Schuler, K. K. 2005. *VerbNet: A broad-coverage, comprehensive verb lexicon*. University of Pennsylvania.

Shi, L.; and Mihalcea, R. 2005. Putting Pieces Together: Combining FrameNet, VerbNet and WordNet for Robust Semantic Parsing. In Gelbukh, A., ed., *Computational Linguistics and Intelligent Text Processing*, 100–111. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-540-30586-6.

Shi, P.; and Lin, J. 2019. Simple BERT Models for Relation Extraction and Semantic Role Labeling. *ArXiv*, abs/1904.05255.

Spaulding, E.; Conger, K.; Gershman, A.; Uceda-Sosa, R.; Brown, S. W.; Pustejovsky, J.; Anick, P.; and Palmer, M. 2023. The DARPA Wikidata Overlay: Wikidata as an ontology for natural language processing. In Bunt, H., ed., *Proceedings of the 19th Joint ACL-ISO Workshop on Interoperable Semantics (ISA-19)*, 1–10. Nancy, France: Association for Computational Linguistics.

Speer, R.; Chin, J.; and Havasi, C. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.

Team, G.; Georgiev, P.; Lei, V. I.; Burnell, R.; Bai, L.; Gulati, A.; Tanzer, G.; Vincent, D.; Pan, Z.; Wang, S.; et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.

Team, G.; Kamath, A.; Ferret, J.; Pathak, S.; Vieillard, N.; Merhej, R.; Perrin, S.; Matejovicova, T.; Ramé, A.; Rivière, M.; et al. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.

Teyler, T. J.; and DiScenna, P. 1986. The hippocampal memory indexing theory. *Behavioral neuroscience*, 100(2): 147.

Thakur, N.; Reimers, N.; Rücklé, A.; Srivastava, A.; and Gurevych, I. 2021. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models. *arXiv preprint arXiv:2104.08663*.

Trivedi, H.; Balasubramanian, N.; Khot, T.; and Sabharwal, A. 2022. MuSiQue: Multihop Questions via Single-hop Question Composition. arXiv:2108.00573.

Tulving, E. 1972. Episodic and semantic memory. In *Organization of memory*, xiii, 423–xiii, 423. Oxford, England: Academic Press.

Tulving, E. 2002. Episodic memory: From mind to brain. *Annual review of psychology*, 53(1): 1–25.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

von Ahn, L.; Kedia, M.; and Blum, M. 2006. Verbosity: a game for collecting common-sense facts. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '06, 75–78. New York, NY, USA: Association for Computing Machinery. ISBN 1595933727.

Vrandečić, D.; and Krötzsch, M. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57(10): 78–85.

Wadden, D.; Wennberg, U.; Luan, Y.; and Hajishirzi, H. 2019. Entity, relation, and event extraction with contextualized span representations. *arXiv preprint arXiv:1909.03546*.

Wang, H.; Shi, H.; Tan, S.; Qin, W.; Wang, W.; Zhang, T.; Nambi, A.; Ganu, T.; and Wang, H. 2024. Multimodal needle in a haystack: Benchmarking long-context capability of multimodal large language models. *arXiv preprint arXiv:2406.11230*.

Wilson, M. A.; and McNaughton, B. L. 1994. Reactivation of hippocampal ensemble memories during sleep. *Science*, 265(5172): 676–679.

Xiang, W.; and Wang, B. 2019. A survey of event extraction from text. *IEEE Access*, 7: 173111–173137.

Yang, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

Yang, Z.; Qi, P.; Zhang, S.; Bengio, Y.; Cohen, W. W.; Salakhutdinov, R.; and Manning, C. D. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

Zhan, Q.; Li, S.; Conger, K.; Palmer, M.; Ji, H.; and Han, J. 2023. GLEN: General-Purpose Event Detection for Thousands of Types. *arXiv preprint arXiv:2303.09093*.

Zhang, X.; Chen, Y.; Hu, S.; Xu, Z.; Chen, J.; Hao, M.; Han, X.; Thai, Z.; Wang, S.; Liu, Z.; and Sun, M. 2024. ∞Bench: Extending Long Context Evaluation Beyond 100K Tokens. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 15262–15277. Bangkok, Thailand: Association for Computational Linguistics.

Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing, E. P.; Zhang, H.; Gonzalez, J. E.; and Stoica, I. 2023. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. ArXiv:2306.05685 [cs].

# Appendix

Our technical appendix is structured as follows:

1. Appendix A: Prompts to LLM.
2. Appendix B: Example GSW Instance.
3. Appendix C: GSW QA Example.
4. Appendix D: Qualitative Analysis of GSW performance.
5. Appendix E: Further Implementation Details.
6. Appendix F: Ablation studies.
7. Appendix G: Related work on Memory Augmentation for LLMs.
8. Appendix H: Computational Costs and Resources for Building the GSW.
9. Appendix I: Related Computational Models of Workspaces.

# A    Prompts to the LLM

In the following section, we describe the prompts used by each component of our GSW framework.

## Operator

We present the prompt to generate operator representation in Fig 3 and 4. The full prompt is considerably longer and includes detailed instructions for each task. For brevity, we have included the introduction and first task in full, with summaries of the remaining tasks. The complete prompt is available in our code repository.

## Reconciler

We present the prompt to reconcile unanswered queries with incoming context in Fig 5

## Question Answering

We present the prompt to generate entity summaries which are passed to the answering agent in Fig 6 and the prompt used by the answering agent is presented in Fig 7

# B    Example GSW instance

In this section we present an example instance of the GSW, highlighting the functionality of both the Operator and the Reconciler. Fig 8 illustrates the operator representations for two separate chunks. Fig 9 presents the result of reconciling the two representations presented in Fig 8. Finally Fig 10 presents a portion of the final reconciled workspace with reconciled entities, space/time coupling and answered forward falling questions.

# C    GSW QA Example

Figure 11 illustrates the end-to-end question answering (QA) pipeline of the GSW framework, showcasing how a sample query from the EpBench dataset is processed through each stage.

# D    Qualitative Analysis of GSW performance

This section presents a qualitative analysis of selected queries to further illustrate GSW's superior performance and token efficiency compared to baseline methods, as detailed in Table 6. The chosen queries, whose full text and ground truth answers are provided in Table 5, are representative of varying complexity, with answers requiring the synthesis of information linked to two to seven distinct contextual cues. This detailed examination reveals specific failure modes in baseline approaches that GSW is naturally suited to overcome.

For instance, GraphRAG, which generates summaries of varying detail from source documents, frequently struggles with information loss and often provides an excessive volume of irrelevant context to the LLM, increasing the likelihood of hallucinations. This limitation is particularly noticeable in its handling of queries Q3 and Q4 (see Table 6). These

---

[3]Background context generated according to contextual chunking by Anthropic, see https://www.anthropic.com/news/contextual-retrieval.

---

queries demand precise spatial and temporal understanding of events, aspects that GraphRAG's summarization process does not natively or consistently capture, leading to missing information or inaccuracies in its responses.

HippoRAG2, on the other hand, processes every query through its knowledge graph –constructed by connecting semantically similar phrases across triples derived from all the chapters– to identify the relevant chapters, and then provides full texts of these chapters as context to the LLM for a final answer. The strength of this approach is that they do not need to perform fine-grained analysis of the text -for example for dates and locations;- as long as their retrieval process identifies the right chapter, the onus is on the LLM to retrieve the relevant spatio-temporal information. This is an effective approach if the documents themselves are short and the number of documents needed to answer a query are few. In the EpBench data set the document size is around 500 tokens and the number of documents needed to answer some of the questions is 17; since QA cannot know the number of documents needed for any given question, 17 documents (chapters) were sent for each query across all evaluated methods. As observed for queries Q2 and Q4 in Table 6, this strategy of providing full documents can overwhelm the LLM, leading to hallucinations or the failure to pinpoint the correct answer even when the right document with the necessary information is present in the retrieved context. Furthermore, there were instances (e.g., Q3, Q5) where HippoRAG2 failed to retrieve all the pertinent documents required to comprehensively answer the query.

In contrast, GSW's structured representation and targeted summary generation (as detailed in Table 6 showing 'None' for errors and lower token counts) effectively mitigate these issues. The ability of our GSW framework to collate and then structure spatio-temporal information scattered across the length of a document (via reconciliation) is aptly captured in the entity-level summary for Carter Stewart that is retrieved in response to Q2 (first three sentences are shown below):

> On September 22, 2026, during the morning sessions of a scientific conference at the Metropolitan Museum of Art, Carter Stewart took on the role of a presenter, delivering a final presentation that included statistical analysis using presentation boards and holographic projectors.

The necessary information – Carter Stewart, location, and time –in the original document came from three different paragraphs; in fact, Carter Stewart is referred to as "He" until after location and time information is given:

> The imposing structure loomed before him, its grand facade a testament to both artistry and scientific achievement ...... As he stepped into the **Metropolitan Museum of Art**, the echoing chatter of excited voices ...... The antique clock in the main hall chimed, its resonant tones reminding him of the date: **September 22, 2026** .... found himself particularly engrossed during the third presentation, where **Carter Stewart** explained statistical analysis with a clarity that left the audience spellbound."

Figure 3: LLM prompt for Operator extraction.[3]

Figure 4: LLM prompt for Space Time coupling.

| Query ID | Query Text | Ground Truth Answer |
|---|---|---|
| Q1 | Consider all events that Jackson Ramos has been involved in. List all the locations where these events took place, without mentioning the events themselves. | High Line, Snug Harbor Cultural Center, Central Park, One World Trade Center, Ellis Island |
| Q2 | Reflect on the experiences of Carter Stewart related to Scientific Conference. List all the unique locations where these events took place, without mentioning the events themselves. | Bethpage Black Course, Metropolitan Museum of Art |
| Q3 | Consider all events that Ezra Edwards has been involved in. List all the locations where these events took place, without mentioning the events themselves. | Water Mill Museum, Port Jefferson, Yankee Stadium, New York Botanical Garden, Brooklyn Bridge, Bethpage Black Course, One World Trade Center |
| Q4 | Recall the events related to Tech Hackathon that occurred on March 23, 2025. List all the locations where these events took place, without describing the events themselves. | Yankee Stadium, Water Mill Museum, Woolworth Building, Queensboro Bridge |
| Q5 | Recall the events related to Tech Hackathon that occurred on November 13, 2026. List all the locations where these events took place, without describing the events themselves. | Trinity Church, Woolworth Building, Statue of Liberty, Fire Island National Seashore |

Table 5: Selected Queries and Ground Truth Answers for Qualitative Analysis

Figure 5: LLM prompt for QA reconciliation.

| Query ID | Method | Token Count | Error Description | Analysis/Reason |
|---|---|---|---|---|
| | GSW | 2011 | None | NA |
| Q1 | HippoRAG2 | 9289 | None | NA |
| | GraphRAG | 8189 | Missing 1 location | Info not available in retrieved context. |
| | GSW | 1568 | None | NA |
| Q2 | HippoRAG2 | 8225 | Hallucinated 3 extra locations. | Too much irrelevant information resulted in LLM hallucination. |
| | GraphRAG | 8220 | Missed 1 location and Hallucinated 2 | All required info present in context but LLM hallucinated. |
| | GSW | 1726 | None | NA |
| Q3 | HippoRAG2 | 8475 | Missed 1 location | Info not available in retrieved context. |
| | GraphRAG | 7058 | Missed 1 location | Info not available in retrieved context. |
| | GSW | 5530 | None | NA |
| Q4 | HippoRAG2 | 8614 | Missed 2 locations | All required info present in context but LLM hallucinated. |
| | GraphRAG | 7936 | Missed 3 locations and Hallucinated 1 | All required info present in context but LLM hallucinated. |
| | GSW | 6452 | None | NA |
| Q5 | HippoRAG2 | 8355 | Missed 1 location | Info not available in retrieved context |
| | GraphRAG | 7936 | Missed 2 locations | Info not available in retrieved context. |

Table 6: Qualitative Performance Comparison on Selected Queries (referencing Query IDs from Table 5)

Figure 6: LLM prompt for entity summary generation.

# E  Implementation Details

In this section, we provide further implementation details for the GSW as well as baselines implemented.

## Operator

The operator representations are obtained by prompting GPT-4o with the prompt presented in Fig. 3 with a temperature of 0 to reduce stochasticity. Prior to obtaining the operator representations, we perform co-reference resolution at a Chapter level resolution. Chapters are then chunked into smaller text chunks each containing three sentences without overlap between consecutive chunks. Space-Time coupling is performed after the operator representations are obtained by prompting GPT-4o with the prompt presented in Fig. 4 with temperature set to 0 and max generation tokens set to 1000.

## Reconciler

Reconciliation is performed on consecutive chunks of operator representations; for our study, we reconcile all chunks of a particular chapter to produce one reconciled GSW representation per chapter. Roles and states for reconciled entities are time-stamped and stored, and this historical information is subsequently utilized during the generation of entity-level summaries.

When a reconciled entity provides new space/time information, its associated space/time nodes are updated accordingly. All previously recorded space/time information is also time-stamped and preserved to enrich these entity-level summaries. Furthermore, it is important to note that if an update to a space/time node is triggered by one entity, this new spatio-temporal information is propagated to all other entities

coupled with that same node; this dynamic is illustrated in Figures 9 and 10.

Finally, the reconciliation process also addresses *forward-falling questions* —queries identified by previous Operator instances that can now be answered using the integrated information from the reconciled GSW as detailed in Section 2.1 of the main paper. These questions are resolved by prompting GPT-4o with the instructions detailed in Figure 5. For this QA resolution task, the temperature is set to 0 and maximum generation tokens are set to 500.

## QA

Prior to the final question answering (QA) stage, entity-specific summaries are generated using the GSW structure. For each entity, a prompt is constructed incorporating its roles, states, associated spatio-temporal information, and the questions it addresses through verb phrases (as captured in its GSW representation). This summarization prompt, detailed in Figure 6, is processed by GPT-4o with a temperature of 0 and a maximum of 500 generation tokens.

The question answering (QA) process unfolds as follows: First, Named Entity Recognition (NER) is performed on the input question to identify relevant entities for querying the GSW. Based on these extracted entities, basic string matching is used to find corresponding entities within the consolidated GSW representations. Next, the **entity-specific summaries** (generated as described previously) for these matched entities are retrieved and then re-ranked. This re-ranking is based on the cosine similarity between the embeddings of the entity summaries and the embedding of the input query. To ensure consistency, the Voyage-03 model is employed as the embedding model for both the summaries and the query. Fi-

**User Prompt:**

*You are a  question answering agent that only uses provided information to answer questions.*
*Your task is to answer questions based exclusively on the knowledge graph information provided. Do not use any external knowledge.*
*The information provided is extracted from a Generative Semantic Workspace (GSW) representation, which captures:*
*- Entities: People, places, objects, and concepts*
*- Verb Phrases: Actions or events involving the entities*
*- Spatial Relationships: Locations of entities*
*- Temporal Relationships: Time periods of entities*
*We use the GSW to extract entity summaries, and you will be provided with these summaries along with graph structure for the GSW for each relevant chapter in order to answer the question.*
*Always ground your answer in the provided information, and only provide answers for which there is clear evidence in the information provided. If the information needed is not available,*
*state that you cannot answer based on the available information.*
*Please answer the following question using ONLY the information provided in the knowledge base extract below.*

*First determine which chapters are most likely to contain relevant information based on the question, then based on the entity summaries and the graph structure for those chapters, determine the most likely answer.*
*Answers will always be a SINGLE entity representing a person, event, location or time period. It will not be a description or a concept.*
*QUESTION: questions*
*KNOWLEDGE BASE INFORMATION: gsw summaries*
*First provide a reasoning for which chapters are most likely to contain relevant information based on the question.*
*Then provide a reasoning for which entity is most likely to be the correct answer based on the entity summaries and the graph structure for those chapters.*

**Inputs:**

*Question: "Question to be answered"*
*GSW Summaries: "Summaries produced by the GSW relevant to answer questions"*

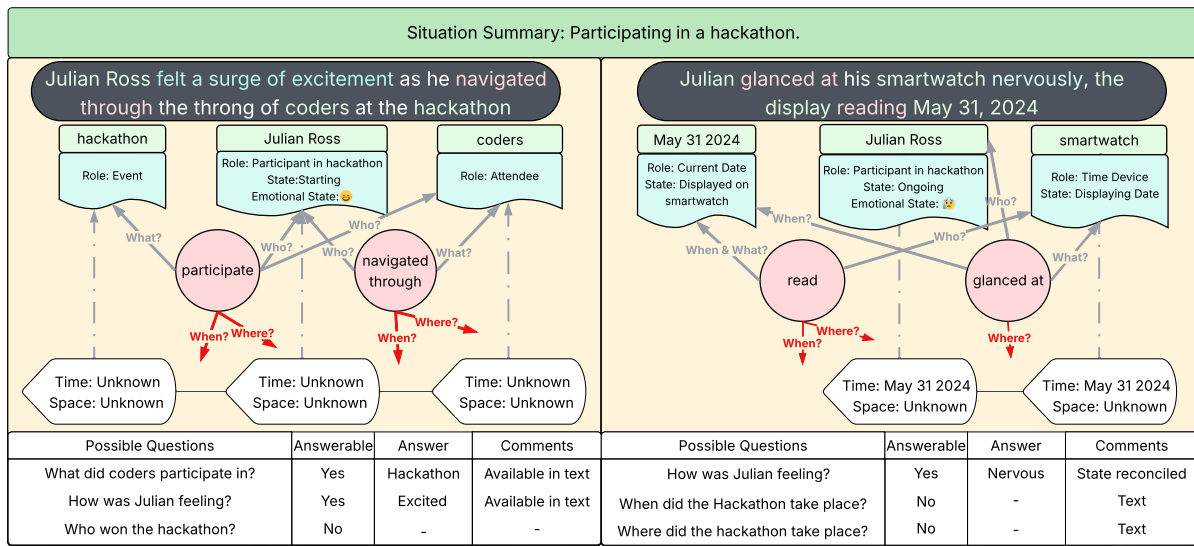Figure 7: LLM prompt for final Question Answering.

Figure 8: **Operator example:** Operator instances of two different chunks, as the GSW framework processes a story.
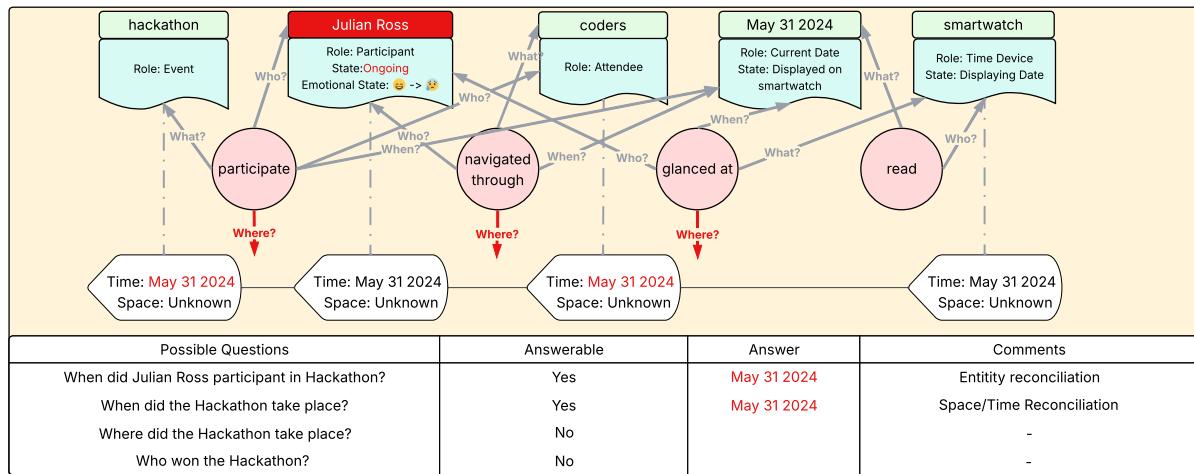


Figure 9: **Reconciler example:** Reconciled result of the two chunks presented in Fig 8
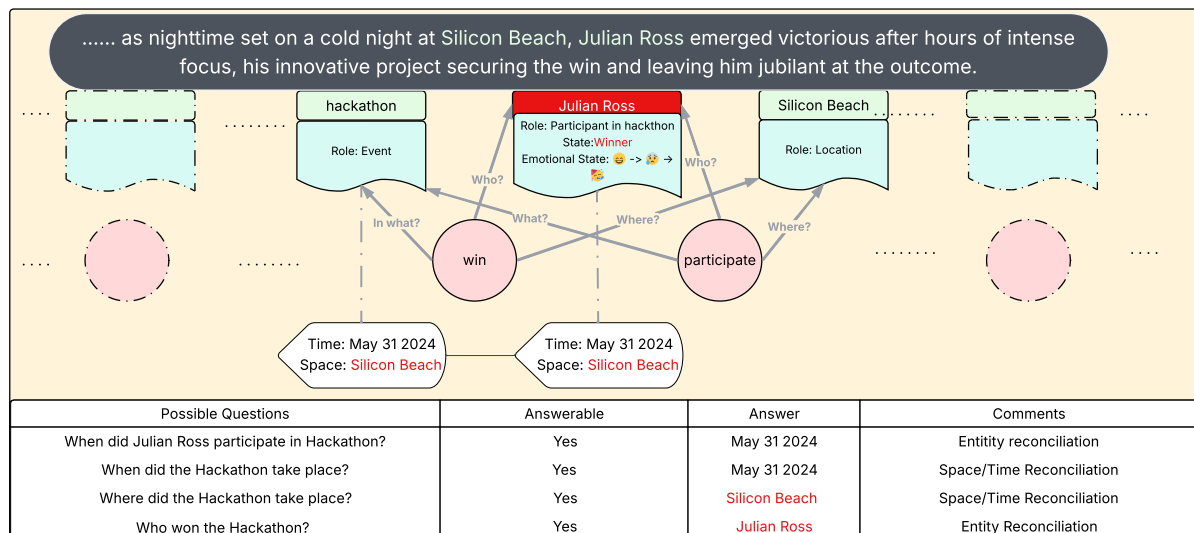


Figure 10: **Final GSW:** A portion of the final reconciled GSW

**Input Query:**

*Reflect on the experiences of Carter Stewart related to Scientific Conference. List all the unique locations where these events took place, without mentioning the events themselves.*

**Named Entities :**

*Carter Stewart, Scientific Conference*

**Retrieved Summaries:**

*Chapter 29:*

Entity: Carter Stewart - Summary: On January 3, 2026, at Yankee Stadium, Carter Stewart, a performer and mime artist, was preparing for a significant performance....

*Chapter 49:*

Entity: Carter Stewart - Summary: Entity: Carter Stewart - Summary: On September 22, 2026, during the morning sessions of a scientific conference at the Metropolitan Museum of Art, Carter Stewart took on the role of a presenter,....

Entity: The scientific conference - Summary: The scientific conference, held on September 22, 2026, was a pivotal moment that took place at the Metropolitan Museum of Art. This event was attended by various individuals

*Chapter 134:*

Entity: Carter Stewart - Summary: On December 25, 2025, Carter Stewart organized a literary-themed festival at Yankee Stadium, stepping onto the field with a sense of pride as both an organizer and participant...

Entity: Carter Stewart's pocket watch - Summary: On December 25, 2025, Carter Stewart's pocket watch, a timekeeping device adorned with intricate clockwork gears and miniature constellations....

*Chapter 166:*

Entity: Carter Stewart - Summary: On March 23, 2024, Carter Stewart, a researcher and presenter, stepped onto the Bethpage Black Course to present his research findings at a Scientific Conference...

Entity: Scientific Conference - Summary: The Scientific Conference, held at the Bethpage Black Course on March 23, 2024, was an event that buzzed with anticipation as it unfolded in a unique setting—a golf course....

**Reranked Summaries:**

*Chapter 166:*

Entity: Carter Stewart - Summary: On March 23, 2024, Carter Stewart, a researcher and presenter, stepped onto the **Bethpage Black Course** to present his research findings at a Scientific Conference...

Entity: Scientific Conference - Summary: The Scientific Conference, held at the **Bethpage Black Course** on March 23, 2024, was an event that buzzed with anticipation as it unfolded in a unique setting—a golf course....

*Chapter 49:*

Entity: Carter Stewart - Summary: Entity: Carter Stewart - Summary: On September 22, 2026, during the morning sessions of a scientific conference at the **Metropolitan Museum of Art**, Carter Stewart took on the role of a presenter,....

Entity: The scientific conference - Summary: The scientific conference, held on September 22, 2026, was a pivotal moment that took place at the **Metropolitan Museum of Art**. This event was attended by various individuals

.
.
. } *Other reranked*
.     *summaries*

**Final Answer:**

**Bethpage Black Course, Metropolitan Museum of Art**

Figure 11: **Illustrative example of the GSW QA framework:** First, NER is performed on the input query to identify key entities. In this version of QA implementation these extracted entities are matched to the relevant GSW instances of chapters via string matching, and the entity-specific summaries (see Appendix E) from the GSWs are retrieved. Subsequently, these retrieved entity summaries are re-ranked based on their semantic similarity to the input query—a score calculated using cosine similarity between their embeddings and the query's embedding. The figure displays a selection of initially retrieved summaries followed by the top re-ranked summaries. Finally, these re-ranked summaries are passed to an answering LLM, which then produces the final answer. As our considerably smaller average token count shows, our entity summaries are already concise, and only entity-relevant chapters are retrieved. Future implementations could leverage several avenues for further reduction in token counts without compromising performance. For example, in a query involving multiple entities, GSWs that have all the entities could be retrieved and sent to the LLM for a final answer; currently our re-ranking step ranks them at the top but we send summaries from other chapters as well, which is not necessary.

nally, these re-ranked summaries are passed to the answering agent (GPT-4o ). The context provided to the agent is limited to summaries derived from a maximum of 17 diverse chapters, a constraint applied to maintain consistency across all evaluated methods and to ensure all dataset questions can be addressed. A detailed example of the QA process is presented in Appendix C.

## Baselines

For HippoRAG2 (Gutiérrez et al. 2025b), GraphRAG (Edge et al. 2025), and LightRAG (Guo et al. 2025), we adhere to each method's default hyperparameters and prompt formats as provided in their respective official implementations. To ensure consistency across baselines, we modify the answering model in HippoRAG2 to use GPT-4o, aligning it with other evaluated methods. Additionally, we set top-k to 17 for HippoRAG2 to retrieve the top 17 relevant documents to align with the QA settings. The detailed configurations for each baseline are listed in Tables 7–9.

| Setting | Value |
|---|---|
| Mode | Local |
| LLM Model | gpt-4o |
| Embedding Model | text-embedding-3-small |
| Response Type | Multiple paragraphs |
| Max Context Tokens | 12000 |
| Text Unit Proportion | 0.5 |
| Community Report Proportion | 0.1 |
| Top-K Entities | 10 |
| Top-K Relationships | 10 |
| Include Entity Rank | True |
| Include Relationship Weight | True |
| Include Community Rank | False |

Table 7: GraphRAG Baseline Parameter

| Setting | Value |
|---|---|
| LLM Indexing Model | gpt-4o-mini |
| LLM Answering Model | gpt-4o |
| Embedding Model | NV-Embed-v2 |
| QA Top-K | 17 |
| Linking Top-K | 5 |
| Retrieval Top-K | 200 |

Table 8: HippoRAG2 Baseline Parameter

| Setting | Value |
|---|---|
| LLM Model | gpt-4o |
| Embedding Model | text-embedding-3-small |
| Retrieval Mode | Hybrid |
| Chunk Token Size | 1200 |
| Chunk Overlap Size | 100 |

Table 9: LightRAG Baseline Parameter

## Bootstrapping for Evaluation

In our main evaluation for EpBench-200 and EpBench-2000, we represent error bars computed via bootstrap resampling on 1,000 iterations. For each evaluation, we resampled the test set predictions with replacement and computed performance metrics on each bootstrap sample. The LLM judge operated with temperature=0 for deterministic outputs. These standard deviations indicate the variability in scores when different combinations of test examples are weighted through resampling

## F    Ablation Studies

We present the results of ablation studies we performed on our GSW framework.

### Evaluating the GSW on the Short Book Dataset

Table 10 presents results comparing GSW against Vanilla LLM on the shorter 20-chapter variant of EpBench. Both GSW and Vanilla LLM demonstrate strong performance on this smaller dataset. The Vanilla LLM performs particularly well on this version because the entire context length is approximately 10,000 tokens, which easily fits within the model's context window. Notably, even with this shorter context, we observe that Vanilla LLM begins to struggle relative to GSW as the number of matching cues increases, particularly in the 3-5 cue category where GSW shows superior recall (0.910 vs 0.781) and F1-score (0.857 vs 0.777).

This finding further supports our main results presented in Table 2 of the main paper, as it demonstrates how Vanilla LLM's performance deteriorates with increased context length. While performing competitively on short narratives, Vanilla LLM struggles with the 200-chapter version where context exceeds 100,000 tokens. In contrast, GSW maintains robust performance across both short and long narratives , highlighting the value of our approach.

### Detailed Results on EpBench-2000

The detailed statistics for EpBench-2000 are presented in Table 11. Although the maximum number of chapters referenced per query in the EpBench-2000 dataset reaches 138, we choose to limit the maximum context utilization to 17 chapters per query, maintaining the same configuration applied to EpBench-200 in the main paper. This choice is based on the fact that the 138-chapter scenario represents an extreme outlier, while 17 chapters suffice to address the majority of queries effectively. Furthermore, processing 138 chapters per query would introduce significant computational overhead and inefficiencies, as it requires feeding an excessive volume of text to the model, which could negatively impact both performance and resource utilization. Since we use the same number of chapters per query as in EpBench-200, we therefore expect a very similar token usage.

Table 12 reports the complete set of metrics for GSW and all baselines on the EpBench-2000 dataset, broken down by cue complexity. These results expand upon the summary in the main text, demonstrating that GSW retains its lead across all levels of episodic complexity, and outperforming the strongest baseline by more than **15%** in F1-score and

| Metric | Method | 0 Cues (N=180) | 1 Cue (N=180) | 2 Cues (N=72) | 3-5 Cues (N=24) | Overall (N=456) |
|---|---|---|---|---|---|---|
| **P** | Vanilla LLM | 0.889 | **0.781** | **0.900** | 0.799 | **0.843** |
|  | GSW (Ours) | **0.939** | 0.751 | 0.804 | **0.854** | 0.841 |
| **R** | Vanilla LLM | 0.889 | **0.919** | 0.813 | 0.781 | 0.883 |
|  | GSW (Ours) | **0.939** | 0.856 | **0.819** | **0.910** | **0.886** |
| **F1** | Vanilla LLM | 0.889 | **0.812** | **0.821** | 0.777 | **0.842** |
|  | GSW (Ours) | **0.939** | 0.745 | 0.784 | **0.857** | 0.834 |

Table 10: **Full EpBench (20-Chapter Book) Performance by Event Categories:** Precision (P), Recall (R), and F1-Score for Vanilla LLM vs. GSW across different event category complexities. (N=X) indicates questions per category.

**14%** in recall. The EpBench-2000 experiment further highlights GSW's ability to scale effectively while maintaining strong performance in long-context, high-recall settings.

| Statistic | Value |
|---|---|
| Number of Chapters | 1967 |
| Total Tokens | 1,012,097 |
| Total Queries (QA Pairs) | 623 |
| Queries by Event Category |  |
| (0 / 1 / 2 / 3-5 / 6+ Cues) | 90 / 165 / 114 / 124 / 130 |
| Max. Chapters Referenced per Query | 138 |
| Min. Chapters Referenced per Query | 0 |

Table 11: **EpBench-2000 Dataset Statistics.**

## Ablating components of the GSW for Question Answering

Table 13 presents the results of ablating both components of the GSW as well as approaches to retrieval, highlighting the importance of each component and our string matching + reranking retrieval mechanism. We note that while naive string matching achieves almost similar performance to our retrieval method, it consumes almost double the number of tokens.

## G    Related work on Memory Augmentation for LLMs

Enabling LLMs to effectively process long narratives requires capabilities akin to human episodic memory – constructing and maintaining a dynamic, coherent understanding of events unfolding over space and time (Tulving 1972; Eichenbaum 2000). Key to this is the ability to accurately track entities, including their evolving states and roles, and to ground events and answer queries based on specific spatial and temporal contexts established within the narrative (Huet, Houidi, and Rossi 2025). While LLMs possess remarkable core abilities, achieving this level of sophisticated, stateful reasoning over extended sequences remains a significant challenge. The following sections analyze inherent limitations in common approaches used to provide context to LLMs, evaluating why they often fall short of systematically delivering these specific episodic memory capabilities.

**Leveraging Long context LLMs**

One approach to providing LLMs with relevant context is to leverage their increasingly large context windows, potentially feeding the entire long narrative along with a query into the prompt. The rapid expansion of context lengths, now reaching millions of tokens, has certainly broadened the scope of tasks LLMs can handle by allowing more raw information to be processed simultaneously (Team et al. 2024).

However, relying solely on this native processing mechanism faces significant hurdles when evaluated against the demands of episodic memory. Firstly, while context windows are growing, they are not infinite, and extremely long narratives may still exceed even the largest available limits. Secondly, even when a narrative technically fits, processing vast amounts of text for every query is computationally expensive, impacting latency and cost. More fundamentally, processing quality often degrades with extreme context lengths (Leng et al. 2024; Hsieh et al. 2024; Wang et al. 2024). Research has shown that LLMs can struggle to consistently access and utilize information spread across very long contexts, with performance notably dipping for information located in the middle (*lost in the middle* phenomenon) (Liu et al. 2023). Feeding potentially large amounts of irrelevant text alongside the crucial details for a specific episodic query can distract the model and hinder its ability to pinpoint and reason over the necessary information.

Finally, perhaps the most critical limitation for systematic episodic tracking is the inherently unstructured nature of the input context. Even with all the necessary details about entity states, roles, locations, and times present within the text, the LLM lacks explicit mechanisms to structure this information dynamically. It must rely solely on its attention mechanism and in-context learning to piece together relationships, track state changes, and maintain temporal coherence across potentially thousands of tokens. This makes the reliable, systematic tracking required for robust episodic memory challenging and often brittle when relying only on the native context window (Huet, Houidi, and Rossi 2025).

**Memory Augmentation for LLMs**

To overcome the challenges of static parametric knowledge and the inefficiencies of processing entire documents in context, Retrieval-Augmented Generation (RAG) has become a standard technique (Lewis et al. 2021; Gao et al. 2024). The typical RAG pipeline involves pre-processing a

| Metric | Method | 0 Cues (N=90) | 1 Cue (N=165) | 2 Cues (N=114) | 3-5 Cues (N=124) | 6+ Cues (N=130) | Overall (N=623) |
|---|---|---|---|---|---|---|---|
| **P** | Embedding RAG | $0.789 \pm 0.043$ | $\underline{0.751} \pm 0.028$ | $\mathbf{0.845} \pm 0.026$ | $\underline{0.840} \pm 0.031$ | $\mathbf{0.911} \pm 0.025$ | $\underline{0.827} \pm 0.014$ |
| | GraphRAG (Edge et al. 2025) | $\mathbf{0.943} \pm 0.025$ | $0.747 \pm 0.038$ | $0.639 \pm 0.040$ | $0.692 \pm 0.038$ | $0.795 \pm 0.043$ | $0.761 \pm 0.017$ |
| | HippoRAG2 (Gutiérrez et al. 2025b) | $0.620 \pm 0.051$ | $0.638 \pm 0.032$ | $0.803 \pm 0.032$ | $0.824 \pm 0.028$ | $\underline{0.893} \pm 0.021$ | $0.759 \pm 0.016$ |
| | LightRAG (Guo et al. 2025) | $0.790 \pm 0.042$ | $0.534 \pm 0.039$ | $0.560 \pm 0.040$ | $0.593 \pm 0.035$ | $0.787 \pm 0.039$ | $0.649 \pm 0.018$ |
| | GSW (Ours) | $\underline{0.867} \pm 0.0025$ | $\mathbf{0.761} \pm 0.0020$ | $\underline{0.841} \pm 0.0019$ | $\mathbf{0.841} \pm 0.0019$ | $0.870 \pm 0.0019$ | $\mathbf{0.830} \pm 0.0010$ |
| **R** | Embedding RAG | $0.789 \pm 0.043$ | $\underline{0.764} \pm 0.032$ | $\underline{0.795} \pm 0.033$ | $0.637 \pm 0.031$ | $0.480 \pm 0.028$ | $\underline{0.688} \pm 0.015$ |
| | GraphRAG (Edge et al. 2025) | $\mathbf{0.943} \pm 0.025$ | $0.492 \pm 0.037$ | $0.587 \pm 0.039$ | $0.538 \pm 0.036$ | $0.321 \pm 0.025$ | $0.548 \pm 0.017$ |
| | HippoRAG2 (Gutiérrez et al. 2025b) | $0.620 \pm 0.050$ | $0.703 \pm 0.034$ | $0.769 \pm 0.031$ | $\underline{0.647} \pm 0.029$ | $\underline{0.491} \pm 0.026$ | $0.648 \pm 0.016$ |
| | LightRAG (Guo et al. 2025) | $0.790 \pm 0.042$ | $0.525 \pm 0.038$ | $0.549 \pm 0.038$ | $0.440 \pm 0.033$ | $0.270 \pm 0.017$ | $0.497 \pm 0.017$ |
| | GSW (Ours) | $\underline{0.867} \pm 0.025$ | $\mathbf{0.844} \pm 0.019$ | $\mathbf{0.864} \pm 0.016$ | $\mathbf{0.792} \pm 0.017$ | $\mathbf{0.633} \pm 0.017$ | $\mathbf{0.796} \pm 0.009$ |
| **F1** | Embedding RAG | $0.789 \pm 0.043$ | $\underline{0.644} \pm 0.031$ | $\underline{0.758} \pm 0.032$ | $0.679 \pm 0.031$ | $0.561 \pm 0.029$ | $\underline{0.675} \pm 0.015$ |
| | GraphRAG (Edge et al. 2025) | $\mathbf{0.943} \pm 0.025$ | $0.436 \pm 0.035$ | $0.547 \pm 0.038$ | $0.541 \pm 0.036$ | $0.405 \pm 0.027$ | $0.544 \pm 0.017$ |
| | HippoRAG2 (Gutiérrez et al. 2025b) | $0.620 \pm 0.050$ | $0.583 \pm 0.031$ | $0.732 \pm 0.031$ | $\underline{0.681} \pm 0.027$ | $\underline{0.578} \pm 0.026$ | $0.635 \pm 0.015$ |
| | LightRAG (Guo et al. 2025) | $0.790 \pm 0.042$ | $0.436 \pm 0.034$ | $0.514 \pm 0.037$ | $0.463 \pm 0.033$ | $0.375 \pm 0.021$ | $0.494 \pm 0.016$ |
| | GSW (Ours) | $\underline{0.867} \pm 0.025$ | $\mathbf{0.741} \pm 0.020$ | $\mathbf{0.818} \pm 0.017$ | $\mathbf{0.789} \pm 0.017$ | $\mathbf{0.698} \pm 0.016$ | $\mathbf{0.773} \pm 0.009$ |

Table 12: **GSW performance on Epbench (2000-Chatpers Book):** Performance is grouped by metric (Precision, Recall, F1-Score) across different numbers of matching cues per query. (N=X) indicates questions per category. Error bars are estimated via bootstrap resampling. Best score in each column for each metric group is **bold**; second best is underlined.

knowledge corpus (e.g., the entire narrative document) into smaller chunks. These chunks are then indexed, commonly using dense vector embeddings obtained from encoder style LLMs(Devlin et al. 2019; Reimers and Gurevych 2019; Lee et al. 2025), though sparse methods like BM25(Robertson and Zaragoza 2009) or hybrid approaches are also employed (Cormack, Clarke, and Buettcher 2009). At inference time, the user query is used to retrieve the top-k most relevant chunks from the index based on a similarity metric (e.g., cosine similarity for dense vectors). These retrieved chunks are then presented as augmented context to an LLM, which generates the final response based on both its parametric knowledge and the retrieved information.(Ram et al. 2023)

This approach has proven effective for many knowledge-intensive tasks, particularly fact-based question answering where retrieving specific evidence snippets is sufficient (Karpukhin et al. 2020). However, when evaluated against the requirements of robust episodic memory recall over long narratives, the limitations of standard RAG become apparent (Huet, Houidi, and Rossi 2025). Firstly, the process of retrieving discrete, potentially disconnected chunks based on local query relevance often **fragments the narrative flow**. This makes it exceedingly difficult for the LLM to reliably follow evolving storylines or track the **changing states and roles of entities** over time, as the necessary context may be spread across multiple chunks that are not retrieved together(Chen et al. 2023).

Moreover, this fragmentation problem is compounded by the framework being highly sensitive to the initial chunking strategy(Merola and Singh 2025). Arbitrary chunk boundaries can split crucial information related to an event or an entity's state, leading to information loss during retrieval. For instance, if a character's state changes within a passage, but the chunking algorithm divides this passage at an inoppor-

tune point, the complete context of this state change may not be captured in any single retrieved chunk. Optimal chunking is non-trivial and can significantly impact the ability to reconstruct the necessary context for complex episodic reasoning. Consequently, while standard RAG offers efficiency gains over naive long-context processing, its inherent lack of structure and narrative coherence makes it ill-suited for systematically addressing the dynamic, stateful, and context-dependent nature of episodic memory tasks.

Additionally, standard RAG mechanisms based on semantic similarity often struggle with incorporating specific spatio-temporal constraints that are essential for episodic memory. Embeddings typically capture semantic content but may not adequately encode the nuances of time and location, making it difficult to retrieve context relevant to a specific point in time or place mentioned in a query or implied by the narrative history.

## Structured Representations as Memory

Recognizing the limitations of standard RAG, particularly its tendency to fragment narratives and struggle with temporal coherence, recent work has explored incorporating more explicit structure into the retrieval and augmentation process. Instead of treating the source narrative as a flat sequence of independent chunks, these methods attempt to build richer representations that capture relationships or hierarchies within the text, aiming to provide more contextually relevant information to the LLM.

While these structured approaches offer advantages over standard RAG by preserving more relational or hierarchical context and enabling more sophisticated information integration (like multi-hop reasoning or global summarization), they still face challenges when viewed through the lens of episodic memory (Huet, Houidi, and Rossi 2025).

| Metric | GSW Configuration / Ablation | 0 Events (N=150) | 1 Event (N=150) | 2 Events (N=90) | 3-5 Events (N=98) | 6+ Events (N=60) | Overall (N=548) |
|---|---|---|---|---|---|---|---|
| **P** | w/o Space/Time Linking | 0.978 | 0.799 | 0.814 | 0.851 | 0.854 | 0.868 |
| | QA Input: Verb Phrases | 0.939 | 0.839 | 0.807 | 0.896 | 0.874 | 0.878 |
| | Retrieval: Str.Match, No reranking | 0.967 | 0.773 | 0.860 | 0.872 | 0.932 | 0.879 |
| | Retrieval: Emb. Match, No reranking | 0.922 | 0.792 | 0.797 | 0.874 | 0.876 | 0.855 |
| | Retrieval: NER emb, no reranking | 0.944 | 0.747 | 0.823 | 0.872 | 0.854 | 0.854 |
| | **GSW (Full)** | **0.978** | **0.755** | **0.810** | **0.878** | **0.891** | **0.865** |
| **R** | w/o Space/Time Linking | 0.978 | 0.800 | 0.810 | 0.738 | 0.723 | 0.827 |
| | QA Input: Verb Phrases | 0.939 | 0.766 | 0.644 | 0.674 | 0.551 | 0.747 |
| | Retrieval: Str.Match, No reranking | 0.967 | 0.834 | 0.850 | 0.819 | 0.822 | 0.867 |
| | Retrieval: Emb. Match, No reranking | 0.922 | 0.820 | 0.833 | 0.825 | 0.781 | 0.845 |
| | Retrieval: NER emb, no reranking | 0.944 | 0.710 | 0.750 | 0.721 | 0.624 | 0.768 |
| | **GSW (Full)** | **0.978** | **0.863** | **0.868** | **0.892** | **0.822** | **0.894** |
| **F1** | w/o Space/Time Linking | 0.978 | 0.731 | 0.764 | 0.762 | 0.761 | 0.811 |
| | QA Input: Verb Phrases | 0.939 | 0.733 | 0.633 | 0.719 | 0.621 | 0.754 |
| | Retrieval: Str.Match, No reranking | 0.967 | 0.748 | 0.826 | 0.823 | 0.859 | 0.846 |
| | Retrieval: Emb. Match, No reranking | 0.922 | 0.726 | 0.788 | 0.827 | 0.810 | 0.817 |
| | Retrieval: NER emb, No reranking | 0.944 | 0.629 | 0.717 | 0.748 | 0.693 | 0.756 |
| | **GSW (Full)** | **0.978** | **0.745** | **0.806** | **0.867** | **0.834** | **0.850** |

Table 13: **Ablation Study of GSW Components on EpBench (200-Chapter Book):** Performance across different event categories (Precision, Recall, F1-Score). (N=X) indicates questions per category. Full GSW model results (at the bottom) are for reference from Table 2 in the main paper.

Graph-based methods like GraphRAG(Edge et al. 2025), LightRAG(Guo et al. 2025), HippoRAG(Gutiérrez et al. 2025b,a) and RAPTOR(Sarthi et al. 2024) suffer from two broad limitations. First, they lack mechanisms to track entity state/role changes across time—they represent entities as static nodes without modeling how attributes or relationships evolve throughout a narrative. Second, they provide no specific framework to ground the evolving narrative in space and time, making it difficult to represent sequential developments or causal relationships. These methods typically represent semantic relationships or summarize community structures within a potentially static corpus, but they are not explicitly designed to model the temporal flow of events within a single narrative or to meticulously track the dynamic changes in entity states and roles as the narrative unfolds sequentially. Their structure captures connections, but not necessarily the chronological progression and state transitions required for recalling specific episodes.

Other research efforts have targeted episodic memory more directly. For instance, Larimar (Das et al. 2024) proposes modifications to the LLM's attention mechanism, while EM-LLM (Fountas et al.) introduces specific memory components integrated with openweight models. While promising, these approaches often require fundamental changes to the LLM architecture or are designed specifically for openweight models, limiting their applicability. In contrast, our GSW framework is proposed as a plug-and-play episodic memory module compatible with any underlying LLM (including closed-source models like GPT-4o via API) and, critically, requires no specialized training or fine-tuning of model parameters, relying instead on the LLM's capabilities for its operator and reconciliation functions.

## H    Computational Costs and Resources for Building the GSW

The primary computational costs for the Generative Semantic Workspace (GSW) framework are associated with its initial, one-time indexing process. To index the 200 chapters of the Epbench dataset, the total expense is approximately $15 when utilizing GPT-4o. This cost covers all stages of GSW construction, including the generation of operator representations, reconciliation, and the creation of entity-specific summaries. By leveraging parallel calls to the OpenAI API, managed via the Bespoke Curator library (Marten et al. 2025), this entire indexing task for 200 chapters can be completed in roughly 1 hour. Alternatively, the OpenAI Batch API can be used to reduce costs, with indexing taking hours.

Our primary experiments leverage API-based models (e.g., GPT-4o) and therefore do not necessitate dedicated local computing infrastructure. However, for tasks such as running the baseline method evaluations reported in this study, and for broader experimentation involving various dense retriever models or locally-hosted chat models, we utilized a single server node equipped with four A6000 GPUs.

## I    Related Computational Models of Semantics

Semantic representation frameworks have a rich history in NLP, yet as we explore below, their design choices create inherent limitations for tracking evolving actor states and relationships—a critical requirement for episodic memory. Among the most influential frameworks are PropBank (Kingsbury and Palmer 2002) and FrameNet (Baker, Fillmore, and Lowe 1998), which attempt to define correspondences between (a) the syntactic "realizations" of semantics *explicit*

within language structure, and (b) finite, discrete sets of semantic "roles" (Levin 1993). These approaches rely heavily on manually-annotated lexicon ontologies developed by expert linguists. While valuable for understanding individual sentences, they were not designed for the dynamic, interconnected tracking that episodic memory demands. Below, we detail these frameworks and their limitations for serving as memory systems:

**PropBank:** PropBank utilized a *bottom-up* approach: (1) Dependency Parse Trees (Nivre 2010) were applied to a large text corpus to distill shared syntactic patterns ("Framesets") specific to each verb (a process known as "lexical sampling"). (2) For each Frameset, the corresponding sentences were manually annotated with an enumerated set of *arguments* ARG:0,..., ARG:N. These arguments were later associated to verb-specific definitions using VerbNet (Schuler 2005). The semantic roles are identified as corresponding *spans* within the sentence (commonly a NP, NNP subtree in the dependency parsing). For example, the sentence (A):

**Officers** captured **Sarah** at the
**Sepulveda on-ramp** of the **405**.

would be annotated with the arguments:

*Agent*: **officers**, *Predicate*: **captured**, *Patient*: **Sarah**.

Perhaps the greatest benefit of PropBank was that its syntactic "grounding" made it possible for rule-based and early ML models (Johansson and Nugues 2008; Shi and Lin 2019) to *learn the task of distilling the semantics* given a sentence, albeit within the confines of a *limited* ontology of $> 3000$ verbs and $> 4000$ Framesets

*Event Databases:* PropBank evolved in several directions, including efforts to unify it with related semantic lexicon such as VerbNet and FrameNet (Palmer 2009; Shi and Mihalcea 2005), or augment it via the DWD overlay (Spaulding et al. 2023) to WikiData (Vrandečić and Krötzsch 2014). The latter of these efforts now manifests as "Event" databases (Liu et al. 2020; Xiang and Wang 2019) such as the ACE (Doddington et al. 2004) and ERE (Aguilar et al. 2014) datasets, and led to the DARPA initiative of Event identification/extraction challenges. Events are best motivated by their related identification tasks: Given a sentence, identify the event(s) – from a set of *hundreds* of events in a pre-annotated schema (Zhan et al. 2023; Wadden et al. 2019; Lu et al. 2021) – that the sentence is referring to. For example, (A) would be annotated with the *Capture* event.

**FrameNet:** In contrast to PropBank and related Event ontologies, FrameNet[4] utilizes a *top-down* approach that is not tethered to the syntax structure. Rather, expert linguists aggregated roles (redefined as "Frame Elements" (FE)) from a large corpus of sentences, which are *known to co-exist* under a conceptual gestalt, or "Frame". Each frame additionally comprises a set of "Lexical Units" (LU) - valences (mostly verbs and nouns) whose occurrence in a sentence increases the likelihood of a frame. For example,

the *Frame*: **Taking Captive**

would contain the following frame elements and lexical units:

---

FE: *Agent*, *Captive*, *Cause*
LU: *capture.v*, *secure.v*

FrameNet (thousands of frames and tens of thousands of FE) is a substantially larger and more comprehensive ontology (Baker 2017) compared to Propbank. When originally constructed, automated systems could not effectively identify the frames implied by a sentence; today, however, Transformer models (Vaswani et al. 2017; Devlin et al. 2019; Chanin 2023) have demonstrated success at accurately modeling the sentence-to-frame mapping.

Despite the enormous success and wide adoption of PropBank, FrameNet, and their descendant works, the explicit, finite, and discrete lexicons they employ raise the question: *When is an explicit lexicon ontology complete?*. While FrameNet provides Frame-Frame precedence and subset relationships, these are coarse-grained and do not adequately answer the question: *How can we track the evolution of semantics across a stream of sentences?* - a key requirement for any semantic model to serve as a memory.

More recent work (Speer, Chin, and Havasi 2017) has attempted to assemble a comparatively larger (and less stringent) open-schema semantic ontology of concepts using game-play based crowd-sourcing techniques (von Ahn, Kedia, and Blum 2006). However, such efforts to scale manual annotation ultimately do not address how a complete ontology can be constructed. Event Graph Models (EGM) (Mellor 2017) generate event networks to describe the dynamics of events in a text corpus, often using a combination of submodules such as Coreference Resolution (Ng and Cardie 2002), Named Entity Recognition (NER) (Li et al. 2020) and Semantic Role Labeling (Palmer, Gildea, and Xue 2011). Extensions (Li et al. 2021) generate the *most likely* event *template* sequences. These methods rely on predefined event schema to enumerate the set of possible events. However, while EGMs both track the evolution of semantics across sentences and offer an unsupervised approach to extending existing ontologies, these often marginalize across individual contexts in the training corpus, and generate the most likely event *schema* that follows a current event schema network. As a result, these works have yet to design methods to track the semantics across a specific document. The GSW, particularly through the operator, is designed to overcome these challenges by generating actor-centric, evolving semantic maps that are not constrained by predefined, static lexicons and can capture the nuances of unfolding situations. To demonstrate the GSW Operator's effectiveness in producing these comprehensive semantic maps from complex, real-world text, we conducted a qualitative human evaluation.

## Comparing Existing Semantic Models to the Operator

To empirically validate the Operator's capability in generating these comprehensive semantic maps, particularly its proficiency in interpreting narrative-rich texts where actor, roles and states undergo significant evolution, we leveraged news reports, as they are a popular resource for sampling semantics-rich stories that belong to universally recognized situation patterns. We query GDELT (Leetaru and Schrodt

2013), a Jigsaw-powered news-indexing platform, with Situation identifiers to retrieve a small set of situation-conditioned `en_US` articles. Table 14 presents statistics about the data. These situations were manually selected as an initial seed set – similar to FrameNet's early versions containing few frames (Lowe 1997) – to assess the validity of the GSW framework. Situation-specific seeds are assembled using a bootstrapped method that invokes FrameNet: (a) Frames are linked using subframe and precedence relationships to create weakly connected components; (b) Headers/labels of the frames in each component form the seed search phrases. We evaluate our framework on situations like "Crime and Justice", "Fire Fighting", "Healthcare", and "Technology Development".

| Situation Label | Documents | Sentences | Tokens |
|---|---|---|---|
| *crime and justice* | 80 | 1209 | 100,635 |
| *fire fighting* | 79 | 1116 | 87,901 |
| *technology development* | 81 | 1334 | 122,493 |
| *healthcare* | 81 | 1259 | 117,962 |
| *economy* | 78 | 1264 | 110,605 |

Table 14: **Data Statistics:** Situation-specific news reports are sampled from GDELT. Each document (or article) is split into short contexts $C_1, \ldots, C_N$ (of 3 sentences) before being passed into the Operator to generate the sematic representation.

Table 15 presents these results across five diverse situations, showing strong human preference for the Operator generated representations compared to existing semantic frameworks.

## Annotator Guidelines

Annotators who exceeded $50K$ in total gross pay were recruited from UpWork, a talent resource. These candidates were first interviewed in a 10-minute session to verify that they were proficient in English and those that had prior experience in annotating large-scale AI/ML data – listed as a verified skill on the platform – were selected to move on to the next round. Following this, they were given a set of 10 task prototype examples and 10 unanswered labeling tasks. Those that got 9 out of the 10 annotations right moved on to the first round of labeling. Each task was labeled twice – by the *annotator* and a *verifier* – to ensure quality of the results. Annotators were paid $5/40$ samples which was estimated to take them about 30 minutes at most, or at the rate of $10/$hour, which was confirmed to exceed the federal minimum wage where the annotators were situated. Annotator guidelines are presented in Fig. 12

| Situation | Ours vs. Baseline | | |
| --- | --- | --- | --- |
| | vs. Zhan et al. (GLEN) | vs. Shi & Lin (BERT-SRL) | vs. Chanin (FST) |
| Crime & Justice | 0.90 (0.10) | 0.96 (0.04) | 0.70 (0.30) |
| Economy | 0.98 (0.02) | 0.96 (0.04) | 0.86 (0.14) |
| Firefighting | 0.98 (0.02) | 0.98 (0.02) | 0.79 (0.21) |
| Healthcare | 1.00 (0.00) | 0.96 (0.04) | 0.94 (0.06) |
| Tech. Development | 0.96 (0.04) | 0.96 (0.04) | 0.86 (0.14) |

Table 15: **Operator Evaluations**: Comparison with Existing Frameworks: Given a short context, English-speaking annotators are shown the unlabeled outputs of the Operator and a baseline framework (GLEN, BERT-SRL, FST) and asked to select the one which best summarizes the semantics in the text. The Operator is preferred over baselines across situations.



Figure 12: **Annotator instructions for UpWork Task:** Annotators are asked to compare the outputs of the Operator to the Semantic map output by a baseline framework (either GLEN, BertSRL, FST) given a shared input text context. During annotation, one random baseline map and the Operator output are presented in random order and the annotator is asked to pick the representation of the Semantics that best reflects the information in the context.