

---

# DyG-RAG: Dynamic Graph Retrieval-Augmented Generation with Event-Centric Reasoning

---

**Qingyun Sun<sup>1</sup> Jiaqi Yuan<sup>1</sup> Shan He<sup>1</sup> Xiao Guan<sup>1</sup> Haonan Yuan<sup>1</sup>**  
**Xingcheng Fu<sup>2</sup> Jianxin Li<sup>1✉</sup> Philip S. Yu<sup>3</sup>**

<sup>1</sup>Beihang University <sup>2</sup>Guangxi Normal University <sup>3</sup>University of Illinois Chicago

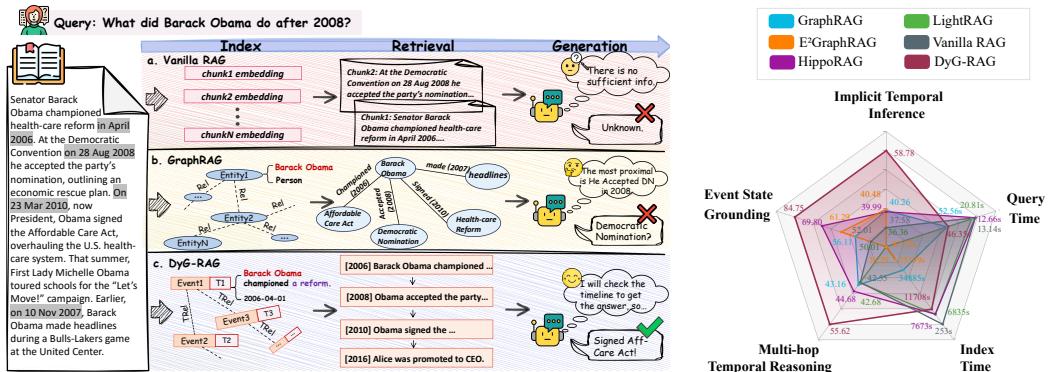
{sunqy,yuanjq,heshan25,xiaoguan,yuanhn,lijx}@buaa.edu.cn,  
fuxc@gxnu.edu.cn, psyu@uic.edu

## Abstract

Graph Retrieval-Augmented Generation has emerged as a powerful paradigm for grounding large language models with external structured knowledge. However, existing Graph RAG methods struggle with temporal reasoning, due to their inability to model the evolving structure and order of real-world events. In this work, we introduce DyG-RAG, a novel event-centric dynamic graph retrieval-augmented generation framework designed to capture and reason over temporal knowledge embedded in unstructured text. To eliminate temporal ambiguity in traditional retrieval units, DyG-RAG proposes Dynamic Event Units (DEUs) that explicitly encode both semantic content and precise temporal anchors, enabling accurate and interpretable time-aware retrieval. To capture temporal and causal dependencies across events, DyG-RAG constructs an event graph by linking DEUs that share entities and occur close in time, supporting efficient and meaningful multi-hop reasoning. To ensure temporally consistent generation, DyG-RAG introduces an event timeline retrieval pipeline that retrieves event sequences via time-aware traversal, and proposes a Time Chain-of-Thought strategy for temporally grounded answer generation. This unified pipeline enables DyG-RAG to retrieve coherent, temporally ordered event sequences and to answer complex, time-sensitive queries that standard RAG systems cannot resolve. Extensive experiments on temporal QA benchmarks demonstrate that DyG-RAG significantly improves the accuracy and recall of three typical types of temporal reasoning questions, paving the way for more faithful and temporal-aware generation. DyG-RAG is available at <https://github.com/RingBDStack/DyG-RAG>.

## 1 Introduction

Recent advances in Large Language Models (LLMs) have demonstrated impressive capabilities in many applications [1]. However, despite their generalization ability, LLMs remain limited by their parametric knowledge, often struggling with factual consistency and adaptability to evolving information. Retrieval-Augmented Generation (RAG) [2, 3] has emerged as a promising paradigm to mitigate these limitations by incorporating external knowledge through retrieval mechanisms. Existing RAG approaches typically rely on static retrieval methods that treat documents as isolated and unstructured items, failing to capture the latent relational structure among retrieved knowledge. To address the limitations of treating retrieved documents as unstructured text, Graph Retrieval-Augmented Generation (Graph RAG) methods [4] have been proposed. These approaches represent knowledge as a structured graph, where nodes correspond to information units (e.g., entities, sentences, or passages) and edges encode explicit relationships such as co-reference, semantic similarity, or knowledge base links. By enabling multi-hop traversal and relational reasoning, Graph RAG models enhance the LLMs’ ability to answer complex queries that require structured inference.



(a) Comparison of RAG pipeline

(b) Comparison of multi-dimensions

Figure 1: Comparison and evaluation of representative RAG methods and DyG-RAG. In Figure 1(b), three axes including Event State Grounding, Implicit Temporal Inference, and Multi-hop Temporal Reasoning are accuracy metrics (%), while Query Time and Index Time represent efficiency metrics. As illustrated, DyG-RAG exhibits superior performance with comparable efficiency, demonstrating a clear advantage over prior RAG methods.

However, both RAG and Graph RAG methods share a fundamental limitation: they struggle with modeling the temporal dynamics of knowledge. In real-world scenarios, information is often not static but unfolds over time. Understanding such temporal evolution is crucial for tasks like timeline generation [5] and historical question answering [6]. Existing retrieval-augmented methods typically treat knowledge as a static snapshot. Vanilla RAG, for instance, relies heavily on semantic similarity between the query and candidate documents, often incorporating timestamps only as an auxiliary dimension in the vector space. Even Graph RAG models, while introducing structural representations, fall short in capturing temporal dynamics. These approaches lack an explicit representation of temporal ordering or causal progression, making it ill-suited for queries that depend on the relative sequence of events. In practice, this leads to several critical failure modes:

- ① Semantic retrievers struggle to distinguish temporal differences.** For example, the questions “*What did Barack Obama do before 2008?*” and “*What did Barack Obama do after 2008?*” yield nearly identical embeddings, as semantic similarity dominates, while temporal intent is ignored.
- ② Temporal constraints are inherently difficult to express and enforce during the retrieval process.** Queries like “*What happened after Obama became president?*” require conditioning on a temporal anchor, but standard retrievers lack mechanisms for relative or conditional retrieval, often returning temporally mismatched results.
- ③ Lack of temporal compositionality.** Questions such as “*How did Obama’s foreign policy evolve after his first term?*” require reasoning over a sequence of events, but standard retrievers return isolated chunks without modeling event progression or temporal dependencies, making it difficult to synthesize coherent answers grounded in temporal chains. A natural extension is to incorporate temporal knowledge graphs (TKGs), where edges are annotated with timestamps to encode when a relation holds. However, these graphs are fundamentally limited in their expressiveness for temporal reasoning: timestamped edges capture only the duration of relations, not the evolution of entity states or event sequences. For example, a TKG may represent that “*Barack Obama–PresidentOf–USA*” holds from 2009 to 2017, but it cannot directly model how his roles, actions, or decisions changed over time, or how they relate to prior and subsequent events.

In this paper, we present **DyG-RAG**, a Dynamic Graph Retrieval-Augmented Generation framework that models the knowledge as a dynamic graph from an event-centric perspective. To organize temporally grounded knowledge, DyG-RAG proposes the Dynamic Event Unit (DEU), which restructures raw text into atomic temporally explicit knowledge units in the event granularity, ensuring that subsequent retrieval is temporally aligned and semantically focused. To capture temporal and causal dependencies across events, we construct an event graph by linking DEUs that share entities and occur close in time, enabling efficient and meaningful multi-hop reasoning. To ensure temporally consistent generation, we introduce an event timeline retrieval pipeline that retrieves event sequences via time-aware vector search and traversal, and guides LLMs using a Time Chain-of-Thought strategy for temporally grounded answers. The contributions are summarized as follows:

Table 1: Representative GraphRAG Methods Comparison

Method	Graph Unit	Edge Type	Retrieval Strategy	Reasoning Mechanism	Dynamic
GraphRAG	Entity + Community	KG Relations + Community Links	Local + Global	Community summary	✗
LightRAG	Chunk Entities	Intra-chunk KG Relations	Dual-level keywords	Shallow path merge	✗
E <sup>2</sup> GraphRAG	Summary Tree + Entity	Semantic + Hierarchical Links	Adaptive Local/Global	Chunk ranking	✗
HippoRAG	Concept Nodes	Concept Associations	PPR-guided multi-hop	PPR subgraph rank	✗
HybridRAG	KG + Chunks	KG relations	Hybrid merge	Evidence voting	✗
<b>DyG-RAG</b>	<b>Dynamic Event Units</b>	<b>Temporal-Semantic Links</b>	<b>Time-aware graph walk</b>	<b>Time-CoT</b>	<b>✓</b>

- We propose a novel framework DyG-RAG, the first dynamic graph retrieval-augmented generation framework that structures and reasons over knowledge from an event-centric perspective.
- By explicitly modeling temporal knowledge as an event graph, DyG-RAG enables time-sensitive retrieval via event timeline reconstruction and enhances reasoning via a Time Chain-of-Thought prompting mechanism, resulting in more faithful and grounded generation.
- Extensive experimental results demonstrate that DyG-RAG significantly improves the performance of three typical types of temporal QA tasks.

## 2 Related Work

### 2.1 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) aims to enhance LLMs by combining relevant information retrieved from external sources, enabling factual grounding and reducing hallucinations [2, 3]. Subsequent works have extended RAG in various directions, including retriever design [7], retrieval strategies [8], and generation techniques [9]. How to enable RAG to retrieve task-relevant and domain-specific knowledge and to perform better reasoning has become a key focus of current research. For example, RAG-end2end [10] extends the vanilla RAG by updating all components with domain-specific knowledge during training. Self-RAG [11] retrieves and reflects using reflection tokens, enabling the generation to tailor to diverse task requirements. Despite these advances, existing RAG methods employ static text chunks and focus on semantic relevance, often overlooking temporal structure and event-level reasoning, which our DyG-RAG framework aims to address.

### 2.2 Graph Retrieval-Augmented Generation

Graph Retrieval-Augmented Generation (Graph RAG) shows superior performance in reasoning by representing knowledge in a more structural way [4, 12]. GraphRAG [13] builds a graph index by deriving an entity knowledge graph from the source documents and uses community summaries to generate the partial responses. LightRAG [14] proposes a dual-level retrieval mechanism from both low-level and high-level knowledge discovery to improve response times. E<sup>2</sup>GraphRAG [15] constructs a summary tree and an entity graph for knowledge modeling and then constructs bidirectional indexes and an adaptive retrieval strategy to capture the many-to-many relationships efficiently. HippoRAG [16] is inspired by human long-term memory and uses personalized PageRank to identify relevant subgraphs for multi-hop reasoning. HybridRAG [17] combines Graph RAG and vanilla Vector RAG techniques to enhance question-answer systems to generate more accurate and contextually relevant answers. Some works focus on exploiting the domain-specific knowledge to enhance the LLMs. However, current Graph RAG approaches typically rely on static or schema-constrained knowledge structures and overlook the evolving nature of real-world information. This makes it difficult to capture temporally sensitive knowledge or reason over event dynamics, which our proposed DyG-RAG aims to address by modeling knowledge as dynamic event graphs. The summary of representative RAG methods and Graph RAG methods is shown in Table 1.

### 2.3 Enhancing LLMs with Temporal Knowledge

Prior works have shown that LLMs often possess outdated or temporally inconsistent knowledge, which limits their ability to reason about time-sensitive or evolving facts [18]. To address this, recent efforts have explored augmenting LLMs with temporal knowledge. A common line of work uses explicit time representations in the generation process of LLMs. Chain-of-Timeline [19] introduces

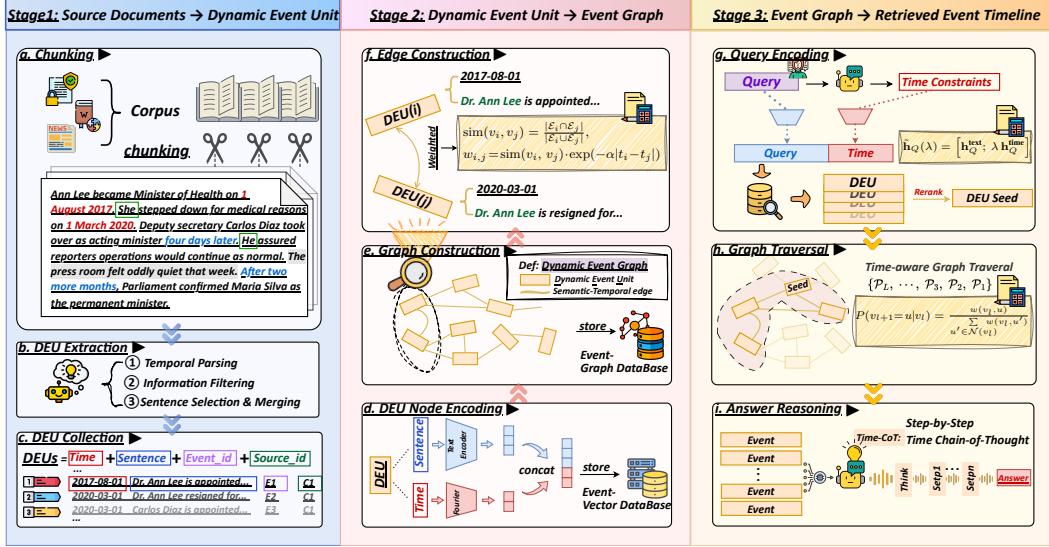


Figure 2: Overall framework of DyG-RAG. The framework consists of three stages: (1) Source documents are parsed into structured Dynamic Event Units (DEUs) via an LLM; (2) A dynamic event graph is constructed from the DEUs and stored in both vector and graph databases; (3) Given a query, the system performs bi-encoding, DEU seed retrieval, and time-aware graph traversal to generate temporally coherent timeline events and finally reason the answer.

topic-relevant event timelines in SQL-style formats to align model attention with time identifiers. Recently, another popular line of work incorporates the pre-defined Temporal Knowledge Graphs (TKGs) to enhance LLMs’ temporal reasoning ability. These methods typically enhance LLMs via two mechanisms: (1) integrating TKG facts into prompts [20, 21] or (2) aligning language model representations with temporal knowledge graph embeddings [22]. TimeR<sup>4</sup> [20] exploits the knowledge in the TKGs in a ‘‘Retrieve-Rewrite-Retrieve-Rerank’’ way. GenTKGQA [22] retrieves subgraphs from TKGs and uses the subgraph embeddings to enhance answers. ARI [23] extracts relevant information and provides it to LLMs via abstract reasoning induction.

### 3 DyG-RAG

In this section, we systematically introduce the DyG-RAG framework, a novel solution that enhances the Graph RAG with dynamic event extraction and reasoning. The overall framework of DyG-RAG is shown in Figure 2.

#### 3.1 Preliminary

Following the setting of general Graph RAG [24], the steps are as follows. First, given a text data source  $D = \{d_1, d_2, \dots, d_n\}$ , Graph RAG transforms the text data into a graph-structured data source  $G$  by a constructor. Second, the users’ query  $Q$  is processed by a query processor and passed to the retriever to obtain the relevant content  $C$ . Finally, the LLMs generate the final answer using the retrieved content  $C'$  as triggers.

#### 3.2 DyG-RAG Workflow

As an overview, the workflow of DyG-RAG converts unstructured text into a dynamic graph by extracting event units, building an interconnected index, and retrieving relevant event sequences. Each stage transitions naturally into the next: from detecting events to indexing relationships to resolving queries, ensuring that the evolving dynamics of knowledge are consistently captured and preserved throughout the entire process.

**Stage 1: Source Documents → Dynamic Event Unit (Section 3.3).** DyG-RAG decomposes raw text in the source documents  $D = \{d_1, d_2, \dots, d_n\}$  into precisely anchored Dynamic Event Units {DEU} that capture essential semantic and temporal information, forming the foundational storage schema for subsequent processing.

**Stage 2: Dynamic Event Unit → Event Graph (Section 3.4).** DyG-RAG organizes extracted dynamic event units into a knowledge-enriched dynamic event graph  $G$  where edges encode both entity co-occurrence and temporal proximity, creating a navigable structure that models narrative flows and causal dependencies.

**Stage 3: Event Graph → Retrieved Event Timeline (Section 3.5).** DyG-RAG executes a four-phase pipeline: coarse retrieval with time-enhanced embeddings, cross-encoder-based semantic filtering, multi-seed graph traversal, and chronological timeline construction, to reconstruct coherent event sequences and resolve complex event-centric queries.

### 3.3 Dynamic Event Unit Extraction

Existing RAG systems typically index paragraph-level chunks, which often blur temporal boundaries and embed multiple time points, leading to imprecise temporal matching and reduced interpretability. To enable temporally grounded retrieval, we replace traditional retrieval units (*e.g.*, paragraphs or knowledge graph triples) with **Dynamic Event Units (DEUs)** defined as follows:

**Definition 1 (Dynamic Event Unit).** A *Dynamic Event Unit (DEU)* is a self-contained factual statement that describes a discrete event or stable state occurring at a specific time point or over a clearly defined interval.

$$\text{DEU} = \{s_i, t_i, \text{ID}_{\text{event}}, \text{ID}_{\text{source}}\}, \quad (1)$$

where  $s_i$  is the sentence,  $t_i$  is the normalized timestamp,  $\text{ID}_{\text{event}}$  and  $\text{ID}_{\text{source}}$  are the event ID and source ID.

This formulation aligns directly with temporal questions (*e.g.*, “What happened?”, “When did something change?”) and serves as the minimal and coherent unit for time-aware retrieval.

The extraction pipeline of DEU consists of four key components: document chunking, temporal parsing, information filtering, and sentence selection and merging.

**(1) Document Chunking.** We first split each source document  $d_i \in D$  into overlapping segments of fixed length to preserve contextual coherence and bound computational cost. To mitigate topic drift and maintain semantic consistency, we prepend the document title to each chunk. Documents shorter than the predefined chunk length are retained as single segments. Subsequently, we extract candidate events from each chunk using an LLM for downstream processing.

**(2) Temporal Parsing.** We then identify temporal expressions within each candidate event and normalize them to create a consistent time anchor. Absolute timestamps (*e.g.*, “March 2008”, “2021-06-15”) are identified and prioritized by granularity. The finest-grained date is assigned as  $t_i$  and stored in a time stack for anchoring future events. Relative or vague expressions (*e.g.*, “earlier that year”, “recently”) are resolved by referencing the most recent absolute date within the same paragraph or context window. Time intervals (*e.g.*, “from 2010 to 2015”) retain their full span in text, but use the earliest point as  $t_i$  for indexing. If no reliable temporal anchor can be extracted,  $t_i$  is assigned a static value, indicating timeless background facts.

**(3) Information Filtering.** To ensure retrieval-relevant content, we compute an information score for each candidate event sentence based on the presence of key attributes: ① Contains named entities or coreferable referents. ② describes a state change or eventive predicate (*e.g.*, “became”, “resigned”, “launched”). ③ includes results or quantitative indicators. ④ is anchored in time with month-level precision or higher. We assign one point for each criterion that the candidate sentence satisfies. Only candidates with  $\text{score}(s) \geq 1$  are retained as DEUs. This step removes generic, underspecified, or off-topic sentences.

**(4) Sentence Selection and Merging.** Finally, we normalize, disambiguate, and aggregate valid DEUs to ensure appropriate granularity and semantic coherence. For coreference resolution, ambiguous pronouns are replaced with explicit entity mentions. Regarding event granularity, we assign one DEU per sentence, unless multiple tightly-related actions share the same time anchor, in which case their predicates are merged into a coordinated DEU. For the temporal resolution constraint, events with different day-level timestamps are not merged to preserve temporal precision.

This process yields a clean, structured set of DEUs, each representing a unique, time-localized event ready for graph-based indexing and retrieval. By matching the natural unit of human temporal questions (e.g., “what happened?”, “when did something change?”), DEUs serve as the fundamental building block for downstream dynamic graph construction and temporal reasoning.

### 3.4 Event Graph Construction and Indexing

To support structured retrieval and temporal reasoning, we organize the extracted DEUs into a dynamic Event Graph.

**Definition 2** (Event Graph). *An Event Graph  $G = (\mathcal{V}, \mathcal{E})$  is weighted graph, where nodes  $v_i \in \mathcal{V}$  represent DEUs and edges  $e_{ij} \in \mathcal{E}$  encode the degree of temporal and semantic relevance between events DEU $_i$  and DEU $_j$  with edge weight  $w_{i,j} \in [0, 1]$ .*

**(1) DEU Node Encoding.** Each DEU is encoded into a dense vector representation that fuses semantic and temporal information. Specifically, given a DEU  $v_i$  with sentence text  $s_i$  and timestamp  $t_i$ , we compute its embedding as:

$$\mathbf{z}_i = \text{Concat}(\mathbf{h}_i^{\text{text}}, \mathbf{h}_i^{\text{time}}), \mathbf{h}_i^{\text{text}} = \text{Encoder}_{\text{text}}(s_i), \mathbf{h}_i^{\text{time}} = \phi(t_i) \quad (2)$$

Where  $\text{Concat}(\cdot, \cdot)$  is the concatenation operation,  $\mathbf{h}_i^{\text{text}}$  is the sentence embedding generated by a pretrained encoder, and  $\mathbf{h}_i^{\text{time}}$  is the time embedding.  $\mathbf{h}_i^{\text{time}}$  is obtained via a Fourier time encoder  $\phi(\cdot)$  that maps the timestamp  $t_i$  into a smooth periodic representation capturing relative time distances.

**(2) Edge Construction and Weighting.** The edges between DEU nodes in the event graph are constructed based on two core criteria: *Entity Co-occurrence* and *Temporal Proximity*. For Entity Co-occurrence, two DEUs must mention at least one common named entity or co-referent entity. For Temporal Proximity, the absolute time difference between their timestamps must be within a threshold window  $\Delta t$ . Formally, an undirected edge  $e_{i,j}$  is added between nodes  $v_i$  and  $v_j$  only if:

$$\text{EntityOverlap}(e_{i,j}) > 0 \text{ and } |t_i - t_j| \leq \Delta t. \quad (3)$$

To capture semantic closeness between two events, we leverage their involved entities for computing the overlap of their respective entity sets:

$$\text{sim}(v_i, v_j) = \frac{|\mathcal{E}_i \cap \mathcal{E}_j|}{|\mathcal{E}_i \cup \mathcal{E}_j|}, \quad (4)$$

where  $\mathcal{E}_i$  and  $\mathcal{E}_j$  denote the sets of named entities extracted from events  $v_i$  and  $v_j$ .

This entity-based similarity ensures that multi-hop traversals remain consistently anchored around common entities, thereby reducing semantic drift in longer reasoning paths. Moreover, successive overlaps among entities across events naturally reveal meaningful event chains (e.g., *Person → Organization → Policy*), effectively enhancing the ability to uncover complex, indirect relationships.

To further incorporate temporal information into the edge weighting, each edge is assigned a weight  $w_{i,j}$  that combines semantic similarity and temporal closeness:

$$w_{i,j} = \text{sim}(v_i, v_j) \cdot \exp(-\alpha|t_i - t_j|), \quad (5)$$

where  $\alpha$  is a decay hyperparameter controlling the sensitivity of edge weights to temporal proximity.

To ensure sparsity and focus traversal on meaningful connections, we restrict each node  $v_i$  to connect with at most  $K$  most related nodes  $\mathcal{N}_i$  that satisfy both conditions, where  $K$  is a tunable parameter.

$$\mathcal{N}_i = \text{Top-}K_{j \neq i}(w_{i,j}). \quad (6)$$

When introducing a new event node, entities extracted from it determine connections to existing nodes, naturally embedding the event into meaningful semantic and causal contexts. This incremental insertion enriches the event graph dynamically, enabling coherent growth and the continuous extension of event chains, effectively capturing the evolving dynamics intrinsic to the DyG.

**(3) Event Graph Indexing.** To balance semantic relevance and temporal proximity, we extend NANOVECTORDB<sup>1</sup> with a `TimestampEnhancedVectorStorage` layer that concatenates a sinusoidally encoded timestamp  $\mathbf{h}_i^{\text{time}} \in \mathbb{R}^{d_\tau}$  to each DEU semantic embedding  $\mathbf{h}_i^{\text{text}} \in \mathbb{R}^d$ , producing a

<sup>1</sup><https://github.com/gusye1234/nano-vectordb>

time-enhanced embedding  $\mathbf{z}_i \in \mathbb{R}^{d+d_\tau}$  that is stored in the vector index. The event graph  $G = (\mathcal{V}, \mathcal{E})$  is managed using NETWORKX<sup>2</sup>, where edge weights  $w_{i,j}$  are preserved for traversal and a bidirectional ID map links vector search results to graph nodes.

By combining semantic and temporal signals, the event graph reflects narrative flow with meaningful path semantics: each hop represents a grounded, interpretable transition between events. This design also supports parallel, multi-path traversal strategies in retrieval, enabling DyG-RAG to reconstruct temporally consistent and semantically coherent event chains for complex queries (as introduced in Section 3.5).

### 3.5 Event Timeline Retrieval and Prompting with Time CoT

**(1) Query Parse with Temporal Intent.** To enable time-aware retrieval over the event graph, we design a query processing module that transforms natural language questions into joint semantic-temporal embeddings. This allows queries to align not only with topical content but also with temporal intent.

Given a user query  $Q$ , we perform two steps: temporal information extraction and query embedding. For the temporal information extraction, we leverage an LLM to extract temporal constraints  $t_Q$  from the query text. Then we encode the query by separately extracting its semantic and temporal representations:

$$\mathbf{h}_Q^{\text{text}} = \text{Encoder}_{\text{text}}(Q), \quad \mathbf{h}_Q^{\text{time}} = \phi(t_Q), \quad (7)$$

where  $\mathbf{h}_Q^{\text{text}}$  is the semantic embedding produced by the same encoder used for DEUs and  $\mathbf{h}_Q^{\text{time}}$  is the temporal encoding generated using the same Fourier-based mapping  $\phi(\cdot)$  as in the event graph. Before retrieval, we reweight the temporal component of the query embedding by a tunable factor  $\lambda$ , producing:

$$\tilde{\mathbf{h}}_Q(\lambda) = [\mathbf{h}_Q^{\text{text}}; \lambda \mathbf{h}_Q^{\text{time}}], \quad \text{sim}(Q, i) = \cos(\tilde{\mathbf{h}}_Q(\lambda), \mathbf{h}_i), \quad (8)$$

Where  $\lambda \in [0, 1]$  controls the trade-off between semantic and temporal components, and  $i \in V$  indexes each candidate DEU node with embedding  $\mathbf{h}_i$ .

By dynamically tuning the temporal scaling factor, DyG-RAG retriever strikes an optimal balance between semantic similarity and temporal proximity, substantially boosting retrieval precision for time-sensitive queries. It also guards against selecting events that merely coincide in time but bear no semantic relation to the query.

**(2) Graph Traversal and Timeline Construction.** Then the resulting vector  $\tilde{\mathbf{h}}_Q$  is used to retrieve temporally aligned and semantically relevant event nodes from the vector database, forming the entry point for subsequent graph traversal and answer generation. To refine this initial set, a cross-encoder reranker scores each retrieved event against the query to filter out irrelevant or weakly aligned nodes, ensuring higher-quality seeds. DyG-RAG then selects a set of top-ranked nodes  $\mathcal{V}_{\text{seed}}$  to serve as the initial event nodes for graph traversal. Given the seed node set  $\mathcal{V}_{\text{seed}}$  retrieved from the vector database, DyG-RAG performs graph traversal to collect supporting evidence in the form of a coherent event sequence. A weighted random walk mechanism is used to explore the event graph while respecting semantic and temporal locality. For each seed node  $v_i \in \mathcal{V}_{\text{seed}}$ , we initiate multiple random walks with fixed length  $L$ . At each step, the next node  $v_{l+1}$  is sampled from the neighbors of  $v_l$  according to an edge-weighted transition probability:

$$P(v_{l+1} = u | v_l) = \frac{w(v_l, u)}{\sum_{u' \in \mathcal{N}(v_l)} w(v_l, u')}, \quad (9)$$

where  $w(v_l, u) \in [0, 1]$  denotes the edge weight encoding semantic and temporal relevance between nodes  $v_l$  and  $u$ , as defined in Equation (5).

This stochastic traversal allows the model to discover multi-hop paths that encode causal chains, narrative progression, or temporally adjacent developments beyond immediate neighbors. The random walk process produces a diverse set of paths  $\{\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_L\}$ , each representing a candidate sequence of related events.

We then convert the path set into a structured timeline. Specifically, we separate static events from timestamped ones to ensure that temporally ambiguous information and unchanged facts are isolated

---

<sup>2</sup><https://networkx.org/>

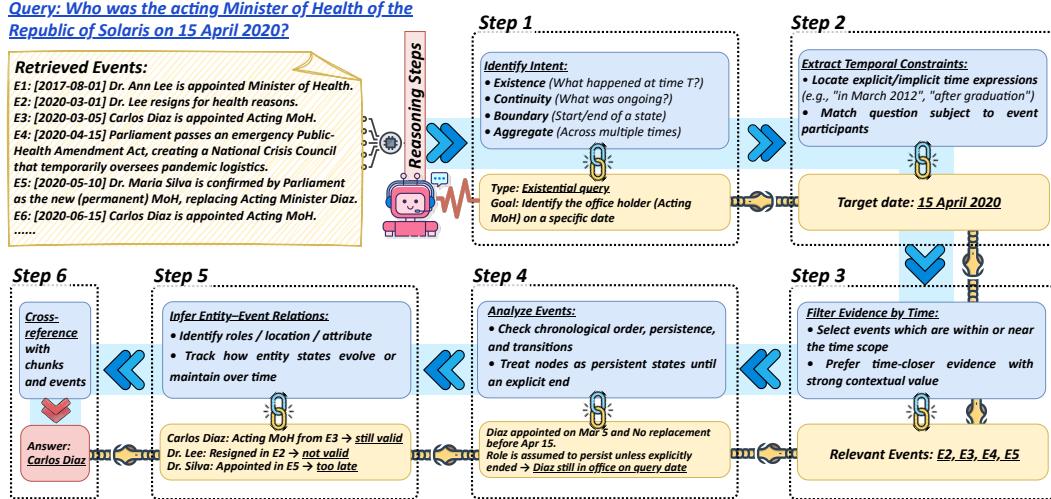


Figure 3: Illustration of Time-CoT.

from concrete temporal anchors. The timestamped events are then chronologically sorted to preserve the natural order of occurrence. Both static and temporal events are merged into a unified, coherent sequence, with each entry explicitly formatted as “Event # [index] [timestamp]: sentence”, making temporal cues transparent to the model. This event timeline is passed to the LLM as structured context, enabling it to perform fine-grained temporal reasoning (e.g., identifying intervals, detecting overlaps, and inferring state persistence across time points) within a clear and interpretable reasoning chain.

**(3) Time Chain-of-Thought for Prompt.** In standard RAG workflows, the retrieved chunks are concatenated into a single prompt without any specific order. As a result, LLMs must implicitly infer and reorder the timeline, which can lead to errors such as temporal inversion or redundant references. To address these challenges, we introduce Time Chain-of-Thought (**Time-CoT**), a specialized chain-of-thought prompt that explicitly encodes temporal relations (e.g., intervals, containment, and causal succession), guiding LLMs through temporal verification and reasoning before answer generation.

Chain-of-Thought prompting is known to boost multi-step reasoning by forcing the model to “think” step-by-step [25], and our Time-CoT extends this benefit to temporal logic, explicitly guiding the LLM through interval comparison and persistence inference. Time-CoT enriches the LLM input with two key components: *Structured Event Timeline* and *Temporal Reasoning Template*. As illustrated in Figure 3, after obtaining the Structured Event Timeline, we apply a Temporal Reasoning Template to structure the prompt and explicitly guide the LLM through symbolic temporal inference. This template encodes procedural reasoning steps: The template first directs the model to identify evidence within the question’s time scope, using both explicit timestamps and relative expressions. It then encourages checking event order, tracking state continuity (e.g., “*office incumbency*”), and distinguishing between event types such as instantaneous actions or ongoing processes. To further enhance precision, the template decomposes question semantics into predefined classes (e.g., “*boundary*”, “*continuity*”, “*aggregate*”), and couples each with tailored reasoning heuristics. Finally, the prompt encourages cross-referencing between events and text chunks, enforcing justifications that cite specific timestamps and chains of events. This structured format converts flat retrieval results into a temporally aware reasoning process grounded in discrete, interpretable steps.

Time-CoT not only improves the model’s ability to handle complex temporal reasoning, mitigating temporal hallucination, but also yields a more interpretable reasoning process, as the chronological event chain makes the temporal logic transparent and grounded in the retrieved evidence.

### 3.6 Discussion: DyG-RAG vs. TKG enhanced LLM Generation.

Temporal Knowledge Graphs (TKGs) are often employed to inject temporal information into LLM generation by encoding facts as timestamped tuples “ $\langle \text{subject}, \text{relation}, \text{object}, \text{time} \rangle$ ” as introduced in Section 2. These symbolic structures are often integrated into generation pipelines via entity linking, temporal fact retrieval, or TKG embeddings. While effective in certain tasks such as complex QA and explainable temporal reasoning, TKG-enhanced generation faces several key limitations when compared with our dynamic graph retrieval-augmented generation framework, DyG-RAG.

**(1) Knowledge Granularity and Expressiveness.** TKGs are fundamentally relation-centric and represent knowledge as atomic facts spanning time intervals. While effective for modeling persistent relations (e.g., “ $\langle \text{Obama}, \text{presidentOf}, \text{US}, 2009\text{--}2017 \rangle$ ”), they lack expressiveness for capturing transient, state-changing events, causal chains, or multi-clause descriptions. In contrast, DyG-RAG introduces Dynamic Event Units (DEUs), defined as self-contained, time-anchored factual statements, offering much finer granularity and aligning naturally with how humans ask and reason about temporal questions.

**(2) Graph Construction and Context Adaptability.** TKGs are typically constructed using a pre-defined schema, a fixed set of entity types and relation predicates curated from knowledge base ontologies. While this schema-driven approach ensures consistency, it also imposes significant limitations. The expressivity of TKGs is constrained by prior assumptions in the ontologies, making it difficult to capture emergent and domain-specific temporal information outside the schema. In contrast, DyG-RAG builds the event graph directly from free-form text, using data-driven DEUs without reliance on a pre-fixed schema. This allows the event graph to naturally adapt to the distribution of knowledge in the corpus and flexibly support context-dependent retrieval, enabling finer alignment with the LLM’s generation needs.

**(3) Grounding and Text Fidelity.** A critical limitation of many TKGs is that they are pre-constructed independently of the raw corpus, typically storing abstracted entity-relation triples without links to the original sentences. This detachment hinders provenance tracking, fine-grained context retrieval, and answer justification. In contrast, DyG-RAG preserves raw text grounding and context, which can be surfaced during generation to support transparency and trustworthiness.

In summary, while traditional TKGs offer structured storage of time-stamped facts, they struggle with expressiveness, adaptability, and grounding with text fidelity. In contrast, DyG-RAG builds an event graph from raw text using fine-grained dynamic event units, captures both semantic and temporal relations, and supports efficient multi-hop retrieval. This makes it more flexible, interpretable, and better suited for generating temporally grounded, context-aware answers.

## 4 Experiments

This section empirically evaluates the proposed DyG-RAG on the temporal question-answering (QA) task. The experiments focus on the following research questions:

- **Q1.** How does DyG-RAG perform in the temporal QA task? (Section 4.2)
- **Q2.** How critical is the event graph for temporally grounded generation tasks? (Section 4.3)
- **Q3.** Can DyG-RAG effectively retrieve and reason over temporal event knowledge? (Section 4.4)
- **Q4.** How efficient is DyG-RAG in terms of index and query? (Section 4.5)

### 4.1 Experimental Settings

#### 4.1.1 Datasets

To evaluate DyG-RAG’s ability to retrieve and reason over temporal information, we focus on three types of temporal questions: **Implicit Temporal Inference**, **Event State Grounding**, and **Multi-hop Temporal Reasoning**. Specifically, we conduct experiments on three open-source temporal QA benchmarks for the above three question types: TimeQA [26], TempReason [27], and ComplexTR [28]. These datasets, derived from Wikipedia-based sources [29], closely simulate temporal RAG settings by providing structured temporal annotations and factual content essential for constructing realistic question-answer pairs.

Table 2: Statistics of Temporal QA Datasets.

Dataset	Documents	Tokens	Avg. Tokens	Questions	Question Type
TimeQA [26]	3,159	1,819,999	576	2,613	Implicit Temporal Inference
TempReason [27]	2,721	2,667,820	980	5,397	Event State Grounding
ComplexTR [28]	112	240,165	2,144	312	Multi-hop Temporal Reasoning

- **TimeQA** [26]. Focused on detailed reasoning with implicit temporal references and cross-sentence cues, we adopt the hard mode, where questions demand deeper inference rather than keyword matching.
- **TempReason** [27]. Centered on event state grounding, we use the L2 mode to evaluate the model’s retrieval ability in determining event states at specified timestamps.
- **ComplexTR** [28]. This benchmark targets multi-hop temporal reasoning, which requires combining multiple temporal cues across events. We use the golden test set to assess the model’s chaining capability.

As no specialized open-source benchmark exists for temporal RAG, we adapt these datasets into a unified document corpus and structured question–answer pairs, reflecting real-world RAG pipelines for temporal retrieval and generation. We have made our processed datasets<sup>3</sup> publicly available to facilitate further research. The statistics of the dataset are shown in Table 2.

#### 4.1.2 Baselines

To the best of our knowledge, there are no existing Dynamic Graph RAG works. For comparison, we select several RAG methods and Graph RAG methods. **RAG methods** include *Vanilla RAG*. **Graph RAG methods** include *GraphRAG* [13], *LightRAG* [14], *E<sup>2</sup>GraphRAG* [15], *HippoRAG* [16].

- *Vanilla RAG*<sup>4</sup> [30]: A standard baseline in simple RAG systems that store the chunked texts in a vector database with text embeddings and use representation vectors to directly retrieve text chunks based on the similarity for the query.
- *GraphRAG*<sup>5</sup> [13]: It transforms the retrieved nodes into communities and traverses the communities to capture the global information. We use two variants of GraphRAG, including local and global retrieval, denoted as *GraphRAG-L* and *GraphRAG-G*.
- *LightRAG*<sup>6</sup> [14]: It uses a dual-level retrieval mechanism to capture both low-level and high-level information of the index graph. We use three variants of LightRAG, including local, global, and hybrid retrieval, denoted as *LightRAG-L*, *LightRAG-G*, and *LightRAG-H*.
- *E<sup>2</sup>GraphRAG*<sup>7</sup> [15]: It constructs a summary tree and an entity graph based on document chunks and then constructs bidirectional indexes to capture their many-to-many relationships.
- *HippoRAG*<sup>8</sup> [16]: It first constructs a knowledge graph for offline indexing and then uses personalized PageRank to identify relevant subgraphs for multi-hop reasoning.

#### 4.1.3 Settings

To ensure consistency across baselines and fair comparison, we adopt the following unified settings for model architecture, document chunking, and retrieval parameters. For the backbone models,

<sup>3</sup><https://github.com/RingBDStack/DyG-RAG/datasets>

<sup>4</sup><https://huggingface.co/facebook/rag-token-nq>

<sup>5</sup><https://github.com/Graph-RAG/GraphRAG/>

<sup>6</sup><https://github.com/HKUDS/LightRAG>

<sup>7</sup><https://github.com/YiboZhao624/E-2GraphRAG>

<sup>8</sup><https://github.com/OSU-NLP-Group/HippoRAG>

Table 3: Results on temporal QA task, where Acc. means accuracy, **Crimson** denotes the best results and **Orange** denotes the runner-ups.

Datasets	TimeQA		TempReaon		ComplexTR	
	Acc. (%)	Recall (%)	Acc. (%)	Recall (%)	Acc. (%)	Recall (%)
Vanilla RAG	37.58	44.14	52.01	64.38	42.55	51.46
GraphRAG-L	40.26	46.28	56.11	67.55	43.16	54.29
GraphRAG-G	10.10	13.83	8.81	11.34	20.97	30.81
LightRAG-L	34.94	39.89	54.23	66.82	41.03	49.46
LightRAG-G	6.66	8.12	13.54	15.27	12.46	13.56
LightRAG-H	36.36	43.06	50.01	62.95	42.68	53.59
HippoRAG	39.99	45.39	<b>69.80</b>	<b>80.54</b>	<b>44.68</b>	<b>55.28</b>
E <sup>2</sup> GraphRAG	<b>40.48</b>	<b>50.19</b>	61.29	73.58	38.29	54.99
<b>DyG-RAG (ours)</b>	<b>58.78</b>	<b>67.02</b>	<b>84.75</b>	<b>91.47</b>	<b>55.62</b>	<b>69.88</b>

we employ Qwen2.5-14B [31] as the large language model for all graph-construction and generation tasks. And all retrieval modules use the BGE-M3 [32] encoder, a multilingual general-purpose embedding model, to encode both queries and text chunks into dense vector representations.

For our DyG-RAG, entity extraction from document chunks is performed using the `dslim/bert-base-NER` model [33], which enables high-quality span-level tagging of person, organization, and location entities. And we use the lightweight but effective TinyBERT-L-2-v2 [34] as the initial retrieval candidate events reranker. All documents are segmented using a sliding window approach with a fixed chunk size of 1,200 tokens and an overlap of 64 tokens. During retrieval, we select the Top-20 candidate chunks or nodes per query, and the input to the language model is truncated to a maximum of 16,384 tokens per instance.

We conduct all experiments using 4 NVIDIA V100 GPUs. Each GPU is equipped with 32 GB of HBM2 memory and delivers up to 15.7 TFLOPS of FP32 performance and 125 TFLOPS via Tensor Core FP16 operations. We deploy LLM inference using the vLLM framework to ensure efficient batching and low-latency throughput. We fix the maximum request concurrency to 32 in the graph construction stage if the method implementation supports concurrent LLM calls.

## 4.2 Results on Temporal QA

To assess DyG-RAG’s effectiveness on the temporal-aware QA task, we follow the token-level evaluation protocol [35], reporting Accuracy and Recall based on the overlap between predicted and gold-standard answer tokens. The accuracy and recall on three benchmark datasets are shown in Table 3. As we can see, DyG-RAG consistently achieves superior performance across all datasets, showing significant improvements in both accuracy and recall. Specifically, compared with the strongest baselines, DyG-RAG achieves absolute accuracy gains of approximately 18.30%, 14.95%, and 10.94% on TimeQA, TempReason, and ComplexTR datasets, respectively. Correspondingly, recall improvements are about 16.83%, 10.93%, and 14.60%.

The TimeQA dataset evaluates the model’s capability to handle complex temporal spans and hierarchical inclusion relations, which need reasoning in implicit temporal cues across sentences. DyG-RAG demonstrates superior performance primarily due to its dynamic event units (DEUs), which systematically organize temporal relationships through semantic-temporal links. This structured temporal arrangement significantly enhances the model’s capability to infer implicit temporal connections. And the auxiliary Time-CoT reasoning guides step-wise disambiguation of complex temporal connections. In contrast, baselines such as GraphRAG show unsatisfactory performance, as their static structures lack the granularity and adaptability required for effectively capturing temporal information.

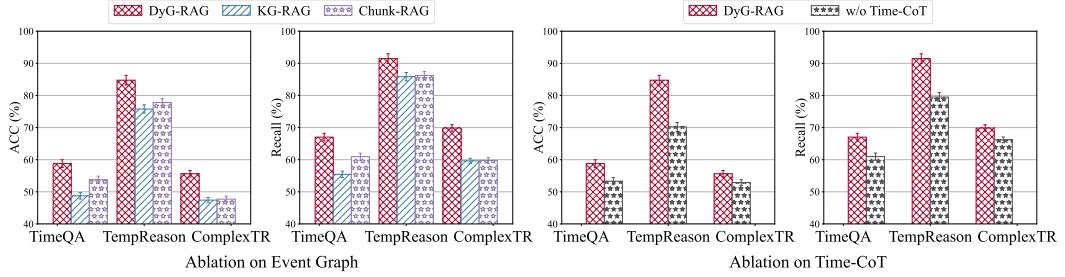


Figure 4: Ablation Study

The TempReason dataset assesses the accuracy of event state identification at specified timestamps. DyG-RAG excels in this task primarily due to its fine-grained DEU retrieval mechanism and precise time embeddings, enabling accurate retrieval of events aligned specifically to queried timestamps. By contrast, HippoRAG organizes information using concept nodes and employs PPR-guided multi-hop retrieval, which can diffuse relevance over loosely connected events and thus retrieve extraneous or time-mismatched states. Other methods that segment events at either overly fine-grained (entity-level) or overly coarse (chunk-level) granularity, which hinder the encoding of temporal discriminative information, omit essential context, or retrieve extraneous information, leading to retrieval errors and degraded event state grounding.

The ComplexTR dataset evaluates the capability for multi-hop temporal reasoning. DyG-RAG achieves superior performance by seamlessly integrating DEUs into coherent semantic-temporal structures, effectively enabling robust multi-event reasoning across complex temporal contexts. Furthermore, the structured reasoning approach of Time-CoT supports coherent narrative chains over multiple temporal spans. Baseline methods such as LightRAG, which rely on keyword-based retrieval and shallow path aggregation, often miss important long-range event associations and lack dynamic re-weighting of emerging temporal links. Consequently, these approaches struggle to reliably construct robust multi-hop event sequences, resulting in inferior performance.

### 4.3 Investigation on Event Graph Construction

To better understand the role of event graph construction in DyG-RAG, we conduct an ablation study focusing on the impact of different knowledge structuring paradigms. Specifically, we compare DyG-RAG with two alternative designs: **KG-based RAG**, which uses a static knowledge graph and temporal prompting, and **Chunk-based RAG**, which performs standard chunk-level retrieval without graph structure. We adopt different knowledge construction strategies and corresponding retrieval units, explicitly forcing the LLM to construct a timeline from the retrieved content before performing step-by-step reasoning to answer the question.

As Figure 4 shows, DyG-RAG consistently outperforms two knowledge construction methods across the three datasets, especially on ComplexTR, which requires multi-hop temporal reasoning. Unlike KG-RAG, which uses a static KG lacking explicit chronological relations, and Chunk-RAG, which retrieves isolated textual chunks without inherent temporal connectivity, our DyG-RAG inherently encodes events within a dynamic and temporally structured graph. This intrinsic event-level structuring enables the LLM to organize retrieved content into coherent and temporally consistent timelines, thereby substantially enhancing its capability to perform temporal reasoning.

### 4.4 Investigation on Event Timeline and Time CoT

To assess the contribution of temporal reasoning in DyG-RAG, we conduct an ablation study in DyG-RAG where the components of Time-CoT are removed. Specifically, we exclude the explicit event timeline construction from the retrieval outputs and omit the structured temporal reasoning instructions in the prompt.

As Figure 4 shows, incorporating Event Timelines and Time-CoT prompting leads to consistently better performance across all datasets. This suggests that prompting the LLM to reason step by step over retrieved timeline events significantly enhances its ability to handle complex temporal

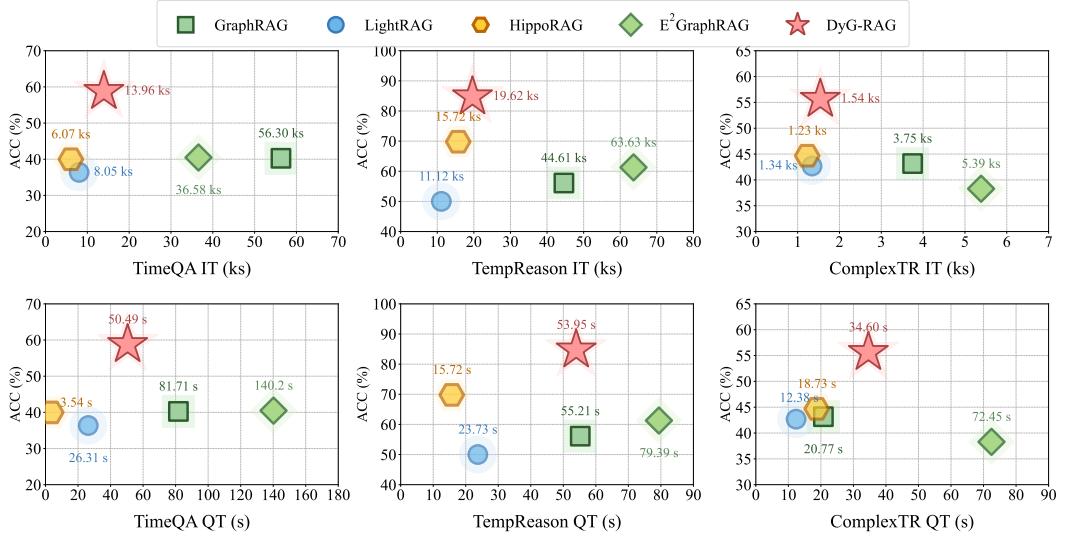


Figure 5: Efficiency Comparison

questions. Instead of treating retrieved information as isolated facts, Time-CoT encourages LLM to construct a coherent chronological narrative and reason explicitly within this structure. Such guided reasoning aligns with the core principles of chain-of-thought prompting, helping LLM more effectively capture temporal dependencies and improve answer accuracy in time-sensitive tasks.

#### 4.5 Efficiency of DyG-RAG

To assess the efficiency of DyG-RAG, we report the **Indexing Time (IT)** and **Querying Time (QT)** across three temporal QA datasets. We include a representative set of GraphRAG-based baselines, all of which adopt structured graph construction for retrieval-augmented generation. For consistency, indexing time is defined as the total time required for graph construction, while querying time refers to the average time per query, encompassing both retrieval from the constructed index and answer generation by LLM.

As Figure 5 shows, efficiency across evaluated methods varies significantly, reflecting differences in both algorithmic complexity and engineering maturity. HippoRAG and LightRAG achieve the fastest querying time and indexing time, primarily due to their mature, optimized code implementations and efficient pipeline design. In contrast, GraphRAG exhibits notably slower indexing performance, mainly attributed to the expensive community report generation phase, which involves intensive graph clustering and serialization steps. E<sup>2</sup>GraphRAG introduces valuable acceleration strategies, particularly in entity extraction and retrieval; however, its end-to-end efficiency is hindered by the lack of concurrent LLM request implementation and lengthy answer-generation prompts. Our DyG-RAG demonstrates moderate yet acceptable latency, ranking consistently third-best overall. Considering its substantial accuracy improvements, the current efficiency trade-off is reasonable. Nonetheless, further speed optimization remains an important open issue for future work.

#### 4.6 Case Study

To gain a deeper understanding of the underlying mechanisms behind retrieval-based temporal reasoning, we perform case studies of three representative methods: DyG-RAG, GraphRAG, and Naive RAG on the three Temporal QA datasets. These analyses illustrate how each method retrieves, structures, and reasons over temporal information, providing interpretable insights into their relative behaviors and failure modes.

Figure 6 illustrates the outputs of three representative methods on the TimeQA dataset, which involves tasks requiring implicit temporal reasoning. The specific query, “*Which team did Bruno Pereirinha play for in May 2013?*”, exemplifies implicit temporal inference as correctly answering the question requires reasoning over two temporally distinct events rather than directly locating ex-

<b>Query:</b> Which team did Bruno Pereirinha play for in May 2013?
<b>Answer:</b> Lazio
<b>a. DyG-RAG:</b>
From the events:
- Event #8: Pereirinha made his professional debut for Sporting CP in 2006–07.
- Event #14: On 22 June 2010, Pereirinha was loaned to fellow Primeira Liga side Vitória de Guimarães.
- Event #17: In early January 2011, Pereirinha joined Kavala F.C. of the Super League Greece also on loan.
- Event #20: On 28 July 2011, Pereirinha signed with S.S. Lazio in Italy for three and a half years.
- Event #28: Pereirinha left S.S. Lazio on 15 July 2015, after buying out his contract.
From the text chunks:
- Chunk 8 mentions that Pereirinha signed with S.S. Lazio in Italy for three and a half years in 2012.
- Chunk 8 also notes that Pereirinha left the Romans (S.S. Lazio) on 15 July 2015, after buying out his contract.
Therefore, the timeline shows that Pereirinha signed with S.S. Lazio in 2012 and stayed with them until July 2015. Thus, in May 2013, Pereirinha was playing for S.S. Lazio.
<b>b. GraphRAG:</b>
Based on the available information, there is no explicit mention of the question. The timeline of his career shows:
- 2006-2007: Professional debut with Sporting CP.
- 2010: Loaned to Vitória de Guimarães.
- 2011: Signed with S.S. Lazio.
- 2012-2013: Returned to Sporting CP.
- 2015: Joined Atlético Paranaense in Brazil.
While it is plausible that Pereirinha was with Sporting CP in May 2013, the exact details for that month are not provided. Therefore, we cannot conclusively determine his team affiliation for May 2013 based on the available data.
<b>c. Naive RAG:</b>
Unknown. The provided context does not specify which team Bruno Pereirinha played for in May 2013.

Figure 6: Case Study on TimeQA

plicit mentions of the year 2013. DyG-RAG correctly identified “*S.S. Lazio*” as the team due to its capability to build a coherent event chain, explicitly capturing both contract initiation and termination events and thus accurately inferring the implicit timeframe. In contrast, GraphRAG failed by hallucinating an incorrect event (“*Returned to Sporting CP*”), caused by incomplete temporal attributes in its KG structure. Meanwhile, Naive RAG lacked the ability to retrieve or infer the necessary implicit temporal connections, leaving the query unanswered. This case illustrates DyG-RAG’s effectiveness in reliably modeling implicit temporal reasoning through dynamic event graph construction.

Figure 7 illustrates the outputs of three representative methods on the TempReason dataset. This example specifically demonstrates event state grounding, as the query “*Which position did Anne-Marie Descôtes hold in Mar 2012?*” requires determining a person’s exact professional state at a given timestamp, achievable by pinpointing a single relevant event interval. In this case, DyG-RAG successfully answers “*director*” by leveraging its explicit time embeddings to accurately retrieve the event interval (2008–2013) closely aligned with the queried timestamp (*March 2012*). DyG-RAG benefits further from structured event triples paired with textual chunks, allowing LLM to validate and reinforce its reasoning effectively. In contrast, GraphRAG incorrectly selects the later-held position (“*Ambassador of France to Germany*”), primarily due to overly dense connectivity within its predefined KG around *Anne-Marie Descôtes*. Its KG-based graph is constructed with nodes representing entities and positions, and edges encoding their relations, but without explicit temporal scopes or directional indicators of state progression. As a result, it cannot capture whether a state is being maintained or has changed over time, leading to confusion between temporally distant roles and misalignment with the queried time. Meanwhile, Naive RAG, despite correctly retrieving a text chunk containing relevant temporal information (“*director from 2008 to 2013*”), fails because isolated and unstructured passages hinder LLM’s ability to reliably infer persistence of states over time. Consequently, Naive RAG’s lack of explicit temporal structuring causes LLM to discard the correct evidence, demonstrating the necessity of structured temporal grounding as provided by DyG-RAG.

Figure 8 illustrates the outputs of three representative methods on the ComplexTR dataset. This query exemplifies multi-hop temporal reasoning, as it requires bridging two semantically distinct event types: an academic milestone (“*pursuing a second doctorate at Université libre de Bruxelles in 1970*”) and a subsequent employment event occurring 3 years and 11 months later. Among the compared methods, only DyG-RAG correctly identified (“*Autonomous University of Madrid*”) as

<b>Query:</b> Which position did Anne-Marie Descôtes hold in Mar, 2012?
<b>Answer:</b> director
<b>a. DyG-RAG:</b>
Based on the provided events and text chunks, let's analyze:
From the events:
- Event #40 [2008-01]: Anne-Marie Descôtes holds the position of director from Jan, 2008 to Jan, 2013.
- Event #80 [2017-06]: Anne-Marie Descôtes holds the position of ambassador of France to Germany from Jun, 2017 to Aug, 2022.
From the text chunks:
- Chunk 6: Anne-Marie Descôtes holds the position of director from Jan, 2008 to Jan, 2013.
- Chunk 6: Anne-Marie Descôtes holds the position of ambassador of France to Germany from Jun, 2017 to Aug, 2022.
Given the temporal scope of the question, we can see that Anne-Marie Descôtes held the position of director from January 2008 to January 2013. This means that in March 2012, she was still in the position of director.
<b>b. GraphRAG:</b>
In March 2012, Anne-Marie Descôtes did not hold the position of Ambassador of France to Germany, as she assumed this role on June 6, 2017.
Explanation Why:
The records show that Anne-Marie Descôtes became the Ambassador of France to Germany on June 6, 2017. There is no information provided that indicates she held any specific position in March 2012.
Entities:
- ID 0: Anne-Marie Descôtes (PERSON)
- Description: Anne-Marie Descôtes is the Ambassador of France to Germany, representing France in the country.
- ID 2: ambassadorship of France to Germany (ROLE)
- Description: Ambassadorship of France to Germany is the diplomatic position held by Anne-Marie Descôtes since June 6, 2017.
Relationships:
- ID 5: Anne-Marie Descôtes -> France (RELATIONSHIP)
- Description: Anne-Marie Descôtes represented France in various capacities, including ambassadorial roles and leadership within the French Foreign Affairs Ministry.
<b>c. Naive RAG:</b>
Unknown. There is no specific mention of Anne-Marie Descôtes holding any particular position in March 2012 within the provided context. The closest dates mentioned are her roles as the ambassador of France to Germany from June 2017 and director from January 2008 to January 2013, but neither of these positions align with the date asked about.

Figure 7: Case Study on TempReason

the answer, benefiting from its temporal-semantic graph structure that naturally supports reasoning across temporally distant events. By explicitly encoding both temporal and semantic relations, DyG-RAG enables the construction of a coherent event timeline, allowing the LLM to perform step-by-step inference over long-range dependencies. In contrast, both GraphRAG and Naive RAG failed for a common reason: they retrieved the nearest temporal event (“1969–1971, University of Texas at Austin”) and prematurely terminated the reasoning process, exhibiting early stopping behavior. GraphRAG, despite its multi-hop capability, lacks temporal awareness in its edge construction, limiting traversal to semantically close nodes while ignoring chronological progression. As a result, both baselines missed critical events spanning a broader temporal window, hindering the LLM’s ability to resolve the query through temporal chaining.

## 5 Conclusion

In this work, we propose DyG-RAG, a dynamic graph retrieval-augmented generation framework that addresses the limitations of existing RAG methods in temporal reasoning. By transforming unstructured text into a dynamic event graph, DyG-RAG captures the evolving structure of knowledge and supports precise, time-aware retrieval. It identifies semantically and temporally grounded event nodes through event unit extraction, encodes entity and temporal relations to model narrative and causal flows via graph indexing, and reconstructs coherent event sequences through structured retrieval. By enabling fine-grained temporal understanding and event-level reasoning, DyG-RAG lays a foundation for more faithful, context-aware language models in dynamic, real-world domains.

## References

- [1] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-*

<p><b>Query:</b> Which employer did Manuel García Velarde work for 3 years and 11 months after he/she worked for Université libre de Bruxelles?</p> <p><b>Answer:</b> Autonomous University of Madrid</p>
<p><b>a. DyG-RAG:</b></p> <p>Step-by-Step Analysis:</p> <ol style="list-style-type: none"> <li>Identify the timeline for Université libre de Bruxelles involvement: <ul style="list-style-type: none"> <li>- Event #22 [1970]: Manuel García Velarde got his second PhD degree at the Université Libre de Bruxelles in 1970. We interpret this as looking for the relevant period after <b>1970</b>.</li> </ul> </li> <li>Calculate the period after Université libre de Bruxelles involvement: <ul style="list-style-type: none"> <li>- The period of 3 years and 11 months after 1970 would be approximately <b>1974</b>.</li> </ul> </li> <li>Identify Manuel García Velarde's employment timeline around 1974: <ul style="list-style-type: none"> <li>- Event #24 [1971]: Manuel García Velarde started to teach and research at the Autonomous University of Madrid in 1971.</li> <li>- Event #29 [1979]: Manuel García Velarde started to work at the National University of Distance Education in 1979.</li> </ul> </li> </ol> <p>Given these events, there is a gap between 1971 and 1979, which includes 1974.</p> <ol style="list-style-type: none"> <li>Determine the employer around 1974: <ul style="list-style-type: none"> <li>- From the timeline, the period between 1971 and 1979 is not covered by any specific employer changes except for starting at the Autonomous University of Madrid in 1971 and moving to the National University of Distance Education in 1979.</li> </ul> </li> </ol> <p>Since there is no specific event mentioning any change in employment between 1971 and 1979, it is reasonable to conclude that Manuel García Velarde was employed at the <b>Autonomous University of Madrid</b> throughout this period.</p>
<p><b>b. GraphRAG:</b></p> <p>After working at Université libre de Bruxelles from January 1965 to January 1969, Manuel García Velarde worked at the <b>University of Texas at Austin</b> from January 1969 to January 1971, which is a period of 3 years and 11 months. Supporting Evidence:</p> <ul style="list-style-type: none"> <li>- Entity:</li> <li>- Université libre de Bruxelles:</li> <li>- Type: ORGANIZATION</li> <li>- Description: "Université libre de Bruxelles is a prominent university in Belgium that has contributed significantly to scientific advancements."</li> <li>- Relation: <ul style="list-style-type: none"> <li>- Source: Manuel García Velarde</li> <li>- Target: Université Libre de Bruxelles</li> <li>- Description: "Manuel García Velarde worked at Université libre de Bruxelles from January 1965 to January 1969."</li> </ul> </li> <li>- Relation: <ul style="list-style-type: none"> <li>- Source: Manuel García Velarde</li> <li>- Target: University of Texas at Austin</li> <li>- Description: "Manuel García Velarde worked at University of Texas at Austin from January 1969 to January 1971."</li> </ul> </li> </ul>
<p><b>c. Naive RAG:</b></p> <p>Based on the provided information, Manuel García Velarde worked at the University of Texas at Austin for 3 years and 11 months after his time at the Université libre de Bruxelles.</p> <p>Supporting evidence:</p> <p>"In 1969 to 1971 he worked at the University of Texas at Austin, where Prigogine led a research institute."</p>

Figure 8: Case Study on CompexTR

*Confidence Computing*, page 100211, 2024.

- [2] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6491–6501, 2024.
- [3] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2(1), 2023.
- [4] Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. Graph retrieval-augmented generation: A survey. *arXiv preprint arXiv:2408.08921*, 2024.
- [5] Daivik Sojitra, Raghav Jain, Sriparna Saha, Adam Jatowt, and Manish Gupta. Timeline summarization in the era of llms. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2657–2661, 2024.
- [6] Samaa Maged, Asmaa ElMaghraby, Ali Marzban, Mohamed Essawey, Amira Ahmed, Esraa Negm, and Wael H Gomaa. Historyquest: Arabic question answering in egyptian history with llm fine-tuning and transformer models. In *2024 Intelligent Methods, Systems, and Applications (IMSA)*, pages 135–140. IEEE, 2024.
- [7] Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. G-retriever: Retrieval-augmented generation for textual graph un-

- derstanding and question answering. *Advances in Neural Information Processing Systems*, 37:132876–132907, 2024.
- [8] Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonello, and Fabrizio Silvestri. The power of noise: Redefining retrieval for rag systems. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 719–729, 2024.
  - [9] Paulo Finardi, Leonardo Avila, Rodrigo Castaldoni, Pedro Gengo, Celio Larcher, Marcos Piau, Pablo Costa, and Vinicius Caridá. The chronicles of rag: The retriever, the chunk and the generator. *arXiv preprint arXiv:2401.07883*, 2024.
  - [10] Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. Improving the domain adaptation of retrieval augmented generation (rag) models for open domain question answering. *Transactions of the Association for Computational Linguistics*, 11:1–17, 2023.
  - [11] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2023.
  - [12] Haoyu Han, Harry Shomer, Yu Wang, Yongjia Lei, Kai Guo, Zhigang Hua, Bo Long, Hui Liu, and Jiliang Tang. Rag vs. graphrag: A systematic evaluation and key insights. *arXiv preprint arXiv:2502.11371*, 2025.
  - [13] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitansky, Robert Osazuwa Ness, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
  - [14] Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. Lightrag: Simple and fast retrieval-augmented generation. 2024.
  - [15] Yibo Zhao, Jiapeng Zhu, Ye Guo, Kangkang He, and Xiang Li. E<sup>2</sup>graphrag: Streamlining graph-based rag for high efficiency and effectiveness. *arXiv preprint arXiv:2505.24226*, 2025.
  - [16] Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. Hipporag: Neurobiologically inspired long-term memory for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
  - [17] Bhaskarjit Sarmah, Dhagash Mehta, Benika Hall, Rohan Rao, Sunil Patel, and Stefano Pasquali. Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction. In *Proceedings of the 5th ACM International Conference on AI in Finance*, pages 608–616, 2024.
  - [18] Bhawna Piryani, Abdelrahman Abdullah, Jamshid Mozafari, Avishek Anand, and Adam Jatowt. It’s high time: A survey of temporal information retrieval and question answering. *arXiv preprint arXiv:2505.20243*, 2025.
  - [19] Jiaying Wu and Bryan Hooi. Chain-of-timeline: Enhancing llm zero-shot temporal reasoning with sql-style timeline formalization. In *Workshop on Reasoning and Planning for Large Language Models*.
  - [20] Xinying Qian, Ying Zhang, Yu Zhao, Baohang Zhou, Xuhui Sui, Li Zhang, and Kehui Song. Timer4: Time-aware retrieval-augmented large language models for temporal knowledge graph question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6942–6952, 2024.
  - [21] Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. Back to the future: Towards explainable temporal reasoning with large language models. In *Proceedings of the ACM Web Conference 2024*, pages 1963–1974, 2024.

- [22] Yifu Gao, Linbo Qiao, Zhigang Kan, Zhihua Wen, Yongquan He, and Dongsheng Li. Two-stage generative question answering on temporal knowledge graph using large language models. *arXiv preprint arXiv:2402.16568*, 2024.
- [23] Ziyang Chen, Dongfang Li, Xiang Zhao, Baotian Hu, and Min Zhang. Temporal knowledge question answering via abstract reasoning induction. *arXiv preprint arXiv:2311.09149*, 2023.
- [24] Haoyu Han, Yu Wang, Harry Shomer, Kai Guo, Jiayuan Ding, Yongjia Lei, Mahantesh Halappanavar, Ryan A Rossi, Subhabrata Mukherjee, Xianfeng Tang, et al. Retrieval-augmented generation with graphs (graphrag). *arXiv preprint arXiv:2501.00309*, 2024.
- [25] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc., 2022.
- [26] Wenhui Chen, Xinyi Wang, and William Yang Wang. A dataset for answering time-sensitive questions. *arXiv preprint arXiv:2108.05266*, 2021.
- [27] Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards benchmarking and improving the temporal reasoning capability of large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Long Papers)*, pages 14820–14835, 2023.
- [28] Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards robust temporal reasoning of large language models via a multi-hop qa dataset and pseudo-instruction tuning. *arXiv preprint arXiv:2311.09821*, 2023.
- [29] Wikipedia contributors. Wikipedia, the Free Encyclopedia. <https://www.wikipedia.org/>, 2023. [Online; accessed in 2023].
- [30] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Namnan Goyal, Heinrich Kütterer, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.
- [31] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2: The next generation of qwen language models, 2024.
- [32] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge-m3: A multilingual, multi-task, multi-vector embedding model, 2024.
- [33] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL*, 2019.
- [34] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174. Association for Computational Linguistics, November 2020.
- [35] Yingli Zhou, Yaodong Su, Youran Sun, Shu Wang, Taotao Wang, Runyuan He, Yongwei Zhang, Sicong Liang, Xilin Liu, Yuchi Ma, and Yixiang Fang. In-depth analysis of graph-based rag in a unified framework, 2025.