# RAG Meets Temporal Graphs: Time-Sensitive Modeling and Retrieval for Evolving Knowledge

Jiale Han, Austin Cheung, Yubai Wei, Zheng Yu, Xusheng Wang, Bing Zhu, and Yi Yang

*Abstract*—**Knowledge is inherently time-sensitive and continuously evolves over time. Although current Retrieval-Augmented Generation (RAG) systems enrich LLMs with external knowledge, they largely ignore this temporal nature. This raises two challenges for RAG. First, current RAG methods lack effective time-aware representations. Same facts of different time are difficult to distinguish with vector embeddings or conventional knowledge graphs. Second, most RAG evaluations assume a static corpus, leaving a blind spot regarding update costs and retrieval stability as knowledge evolves. To make RAG time-aware, we propose Temporal GraphRAG (TG-RAG), which models external corpora as a bi-level temporal graph consisting of a temporal knowledge graph with timestamped relations and a hierarchical time graph. Multi-granularity temporal summaries are generated for each time node to capture both key events and broader trends at that time. The design supports incremental updates by extracting new temporal facts from the incoming corpus and merging them into the existing graph. The temporal graph explicitly represents identical facts at different times as distinct edges to avoid ambiguity, and the time hierarchy graph allows only generating reports for new leaf time nodes and their ancestors, ensuring effective and efficient updates. During inference, TG-RAG dynamically retrieves a subgraph within the temporal and semantic scope of the query, enabling precise evidence gathering. Moreover, we introduce ECT-QA, a time-sensitive question-answering dataset featuring both specific and abstract queries, along with a comprehensive evaluation protocol designed to assess incremental update capabilities of RAG systems. Extensive experiments show that TG-RAG significantly outperforms existing baselines, demonstrating the effectiveness of our method in handling temporal knowledge and incremental updates.[1]**

*Index Terms*—**Retrieval-Augmented Generation, Time-Sensitive Question Answering, Temporal Knowledge Graph.**

## I. INTRODUCTION

**R**ETRIEVAL-Augmented Generation (RAG) [1], [2] has emerged as a powerful paradigm that equips large language models (LLMs) [3]–[6] with external knowledge, which can mitigate LLMs' challenges such as hallucination [7], [8], limited domain-specific expertise [9], and outdated knowledge [10]. A typical RAG framework consists of three core components: indexing an external corpus into a vector database, retrieving passages that are semantically similar to the input query, and generating an answer conditioned on both
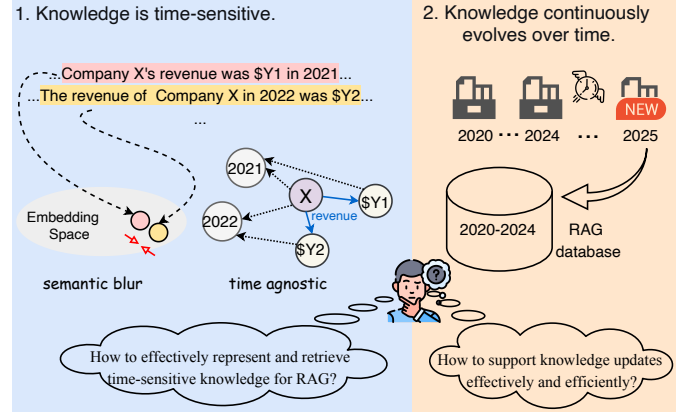


Fig. 1. The temporal challenges for RAG. (1) Time-aware representation: embeddings for time-sensitive facts are often indistinguishable, and knowledge graphs lack explicit temporal attributes. (2) Incremental updates: corpora evolve continuously, but most evaluations assume one-time indexing, creating an evaluation blind spot on update cost and retrieval stability.

the query and the retrieved content. Recent advancements such as GraphRAG [11], [12] further improve the retrieval process by constructing the original corpus as a knowledge graph. This graph-based representation connects isolated facts through shared entities and relations, enabling multi-hop reasoning and deeper understanding of the corpus.

However, existing RAG frameworks overlook a fundamental dimension of knowledge—*time*. In real-world scenarios, knowledge is inherently time-sensitive and continuously evolves over time [13], [14]. For instance, a company's financial reports present varying revenue figures across different fiscal years, and grow over time with the continuous release of new documents. Current RAG systems face significant limitations in represent and retrieve time-sensitive information. For vector-based RAG, factual statements that differ only in the temporal attributes (e.g., "Company X's revenue was $\$Y_1$ in 2021." vs. "Company X's revenue was $\$Y_2$ in 2022.") tend to generate highly similar embeddings, rendering them indistinguishable during retrieval [15], [16]. For graph-based methods, knowledge graphs represent relational facts through static triples (subject, relation, object). Such time-agnostic triples lack native support for temporal attributes. For example, the above facts would be converted as two triples (X, revenue, $\$Y_1$) and (X, revenue, $\$Y_2$), with no indication of the time involved.

Furthermore, real-world corpora evolve continuously, requiring RAG systems to support effective and efficient incremental updates. Although such updates are crucial in real deployments, most studies evaluate RAG systems under a static-

Jiale Han, Austin Cheung, and Yi Yang are with the Hong Kong University of Science and Technology, Hong Kong, China 999077. E-mail: jialehan@ust.hk, mycheungaf@connect.ust.hk, imyiyang@ust.hk.

Yubai Wei is with the University of Turku, Finland FI-20014. E-mail: yubwei@utu.fi.

Zheng Yu, Xusheng Wang, and Bing Zhu are with the HSBC Holdings Plc., Emerging Technology, Innovation, and Ventures, China 200120. E-mail: {matthew.z.yu, atlas.x.wang, bing1.zhu}@hsbc.com.

[1]The dataset and code are available at: https://github.com/hanjiale/Temporal-GraphRAG.

corpus setup, where indexing is built once and subsequent evaluation measures retrieval and answer performance on the fixed database. This leaves a critical evaluation blind spot, including the computational costs of updates and performance changes induced by newly ingested information. As illustrated in Figure 1, these observations highlight two fundamental temporal challenges for RAG: (1) how to effectively represent and retrieve time-sensitive knowledge, and (2) how to support knowledge updates at minimal cost while maintaining robust performance.

To address these challenges, we propose Temporal GraphRAG (TG-RAG), a novel framework that models real-world knowledge with a bi-level temporal graph built from the corpus. The lower layer is a temporal knowledge graph whose nodes are entities mentioned in the corpus and whose edges are relations annotated with timestamps. Relations between the same entity pair at different times are preserved as distinct edges, capturing historical evolution. The upper layer organizes all timestamps into a hierarchical time graph. Cross-layer edges connect each time node to the relation edges that are active at that time. For each time node, we maintain a temporal summary that aggregates the facts attached to that node and the summaries of its finer-grained descendants, yielding a corpus-level, time-scoped view of the knowledge. TG-RAG is update-friendly by design. When new documents arrive, we extract timestamped relations and merge them into the existing graph, then generate summaries only for new time nodes and incrementally propagate updates to their ancestor time nodes, avoiding expensive re-summarization from scratch [11]. For retrieval, we use two time-aware strategies: local retrieval extracts fine-grained facts within a specified time window, and global retrieval leverages temporal summaries to capture significant events or broader trends. In this way, TG-RAG provides effective temporal representation and efficient incremental updates, bridging the gap between conventional RAG and the temporal nature of real-world knowledge.

In addition, due to the scarcity of datasets for complex time-sensitive question answering, we contribute ECT-QA, a high-quality dataset designed to evaluate temporal reasoning in RAG systems. ECT-QA is derived from earnings-call transcripts spanning many companies and time periods, reflecting dynamic factual knowledge within evolving corporate contexts. It contains both specific queries which are fact-oriented and abstract queries that focus on trends and summarization. This allows for comprehensive evaluation of both precise retrieval and deep contextual understanding.

To simulate incremental knowledge updates in real world scenarios, we split the corpus into two distinct parts. The base corpus includes documents from 2020 to 2023, while the new corpus contains documents from 2024. Queries are correspondingly categorized into base questions and new questions based on whether their answers require information from the new corpus. We conduct three evaluation scenarios to thoroughly assess system performance: evaluating base queries on the base corpus, base queries on the fully updated corpus, and new queries on the updated corpus. This comprehensive evaluation framework effectively addresses previously existing blind spots in RAG evaluation, providing insights into both

computational update costs and system stability after updates.

Comprehensive experiments on the ECT-QA dataset show that TG-RAG substantially outperforms current RAG baselines across all evaluation settings, demonstrating the effectiveness of our approach. In particular, when handling evolving knowledge, our method maintains stable performance on historical queries while achieving strong results on new ones with efficient update costs, which highlights its robustness and practicality in dynamic real-world scenarios. Ablation studies and case analyses reveal that these benefits stem from our temporal knowledge graph modeling, which utilizes timestamped edges to encode fine-grained facts and temporal subgraph retrieval to deliver precise evidence gathering. Furthermore, experimental results on widely-used question answering benchmarks and with different LLM backbones confirm the strong generalization capability of our method.

Our contributions are summarized as follows:

- **Make RAG Time-Aware.** We propose Temporal GraphRAG, a novel framework that models external knowledge as a bi-level temporal graph, explicitly representing fine-grained temporal facts while supporting incremental updates for evolving information.
- **Benchmark and Protocol for Incremental Evaluation.** We release ECT-QA, a high-quality benchmark tailored for time-sensitive question answering, along with a comprehensive evaluation protocol specifically designed to assess incremental update capabilities in RAG systems.
- **Effective Retrieval with Efficient Update Cost.** Extensive experiments demonstrate that our approach maintains consistent performance across knowledge updates while achieving low incremental update cost, underscoring its practicality in real-world evolving corpora.

## II. RELATED WORK

### A. Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) has been proposed to connect LLMs with external knowledge sources. Specifically, external corpora are segmented into chunks [17], embedded [18], and stored in a vector database. When a user query arrives, relevant external data is retrieved and combined with the query as input to the LLM, which then generates an informed answer. Beyond this baseline framework, a variety of methods have been developed to further optimize the RAG pipeline. RQ-RAG [19] refines queries for RAG by introducing query rewriting, decomposition, and disambiguation. Self-RAG [20] extends the RAG paradigm by enabling LLMs to adaptively decide when and what to retrieve through self-reflection. LumberChunker [21] leverages LLMs to dynamically chunk documents by identifying semantic shift points. While these approaches enhance the retrieval of explicit facts, they often fall short when addressing global queries that require connecting multiple implicit facts spread across documents.

To tackle this limitation, Microsoft's GraphRAG [11] constructs an entity knowledge graph from external documents and generates community-level summaries for clusters of closely related entities, thereby enhancing query-focused summarization. MemoRAG [22] leverages a long-context LLM to build a
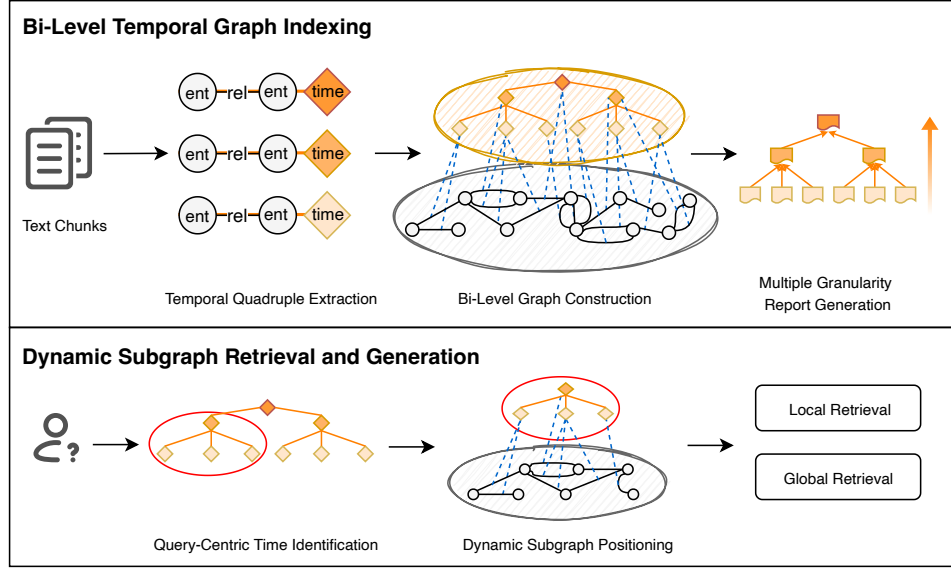
Fig. 2. The overall framework of Temporal GraphRAG.

global memory of the corpus, which produces draft answers and guides retrieval tools to locate useful evidences. LightRAG [23] integrates a graph-based text indexing paradigm with a dual-level retrieval framework, enabling efficient and semantically relevant retrieval to better handle complex queries. HippoRAG [24] converts the corpus into a schemaless knowledge graph, allowing information integration across passages to support reasoning-oriented retrieval. And HippoRAG 2 [25] enhances the original HippoRAG by integrating deeper passage-level context into the Personalized PageRank algorithm. However, few studies explicitly model the temporal attributes of external facts, which limits their effectiveness in domains where knowledge evolves rapidly. To address this, we propose TG-RAG, a novel framework that models the dynamics of knowledge through a bi-level temporal graph, significantly enhancing LLMs' retrieval and reasoning performance in such environments.

### B. Time-Sensitive Question Answering

Time is an important dimension in our physical world. Lots of facts can evolve with respect to time, including those in finance [26], law [27], [28], and healthcare [29]. Therefore, it is important to consider the time dimension and empower the existing question answering (QA) models to reason over time. Son and Oh [30] enhance QA models' temporal understanding through a time-context-dependent span extraction task trained on synthetic temporal data and contrastive time representation learning. Zhu *et al.* [31] reframe time-sensitive QA as programming, using LLMs to translate natural language into executable code that encodes temporal constraints and selects answers via program execution. Yang *et al.* [32] tackle time-sensitive QA by boosting LLMs' temporal sensitivity and reasoning with a temporal information-aware embedding and granular contrastive reinforcement learning. Zhang *et al.* [33] point out that current retrievers struggle with temporal reasoning-intensive questions. The authors further propose a training-free modular retrieval framework that decomposes

questions into content and temporal constraints, retrieves and summarizes evidence accordingly, and ranks candidates with separate semantic and temporal scores.

On the dataset side, a few benchmarks have been introduced to systematically probe models' temporal understanding and reasoning. Chen *et al.* [13] build the first time-sensitive QA dataset TimeQA to investigate whether existing models can understand time-sensitive facts. The dataset is constructed by mining time-evolving facts from Wikidata, aligning them to corresponding Wikipedia pages, having crowdworkers verify and calibrate the temporal boundaries, and finally generating question–answer pairs from the annotated time-sensitive facts. Wei *et al.* [34] propose a dataset called Multi-Factor Temporal Question Answering (MenatQA) containing multiple time-sensitive factors (scope factor, order factor, counterfactual factor), which can be used to evaluate the temporal understanding and reasoning capabilities of QA systems. However, most existing datasets contain only simple single-hop time-sensitive queries, highlighting the need for a more challenging benchmark.

### III. TEMPORAL GRAPH RETRIEVAL-AUGMENTED GENERATION

In this section, we present Temporal Graph Retrieval-Augmented Generation (TG-RAG) and begin by formalizing the task. Given a corpus $\mathcal{D}$ and a time-sensitive query $q$, the goal is to retrieve the set of relevant evidences from $\mathcal{D}$ and generate the answer $a$. TG-RAG operates in two stages, a bi-level temporal graph indexing stage that organizes the corpus into a semantically and temporally structured graph, and a dynamic subgraph retrieval and generation stage where time-relevant facts are retrieved to facilitate response generation. The overall architecture of TG-RAG is illustrated in Figure 2.

## A. Bi-Level Temporal Graph Indexing

In this stage, we reorganize the corpus into a bi-level temporal graph that supports explicit temporal facts representation and deep time-aware understanding, and enables incremental corpus updates.

*1) Temporal Quadruple Extraction:* We first segment the raw corpus $\mathcal{D}$ into multiple chunks $\mathcal{C}$ and prompt an LLM to extract temporal quadruples. For each chunk $c \in \mathcal{C}$, the LLM is asked to first recognize timestamp nodes, then identify the non-temporal entities, and detect the temporal relations linking those entities. The output is a set of quadruples $(v_1, v_2, e, \tau)$, where $v_1$ and $v_2$ are two entities, $e$ denotes the relation, and $\tau$ is the normalized timestamp. The detailed prompt is presented in the Appendix.

*2) Bi-Level Graph Construction:* Based on the extracted temporal quadruples, we construct a temporal knowledge graph $\mathcal{G}_K = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V}$ is the set of non-temporal entities and $\mathcal{E} = \{(v_1, v_2, e, \tau)\}$ is the set of relation edges annotated with timestamps. Multiple facts between the same entity pair at different times are kept as parallel temporal edges, capturing the evolution of their relationship. Furthermore, all unique timestamps are organized into a time hierarchy graph $\mathcal{G}_T$ that partitions the timeline into nested buckets and connects the parent intervals to their immediate sub-intervals (e.g., *year* $\rightarrow$ *quarter* $\rightarrow$ *month* $\rightarrow$ *day*). Cross-layer edges connect each time node in $\mathcal{G}_T$ to the edges in $\mathcal{G}_K$ that are active at that time, yielding a bi-level structure that binds factual relations to their temporal scopes. In addition, we precompute and store vector embeddings for entities and relations based on their textual description: $\mathbf{e}_v$ for each $v \in \mathcal{V}$ and $\mathbf{e}_e$ for each edge in $\mathcal{E}$.

*3) Multi-Granularity Time Report Generation:* For each time node in $\mathcal{G}_T$ from bottom to top, we generate a report that summarizes the activity within its time window. Specifically, we aggregate all entities and relations in $\mathcal{G}_K$ connected to that time node and the time reports of its child nodes, and prompt an LLM to produce a concise report highlighting noteworthy events, interactions, and trends for that time period. These summaries propagate recursively upward, forming a multi-granularity hierarchy of reports, progressively offering broader and richer views of the corpus's temporal dynamics.

*4) Incremental Update of the Temporal Graph:* When new documents arrive, we extract temporal quadruples and merge the resulting timestamped relations into the existing temporal knowledge graph $\mathcal{G}_K$, creating new time nodes in the time hierarchy graph $\mathcal{G}_T$ only when needed. We then generate time reports for the newly created leaf time nodes and incrementally update reports along their ancestor paths in $\mathcal{G}_T$. The reports of unaffected time nodes remain unchanged. Unlike GraphRAG [11] which requires regenerating all summaries upon every update, our approach avoids full recomputation and maintains high efficiency.

## B. Dynamic Subgraph Retrieval and Generation

Given any query $q$, we first identify the query-specific subgraph by semantic and temporal relevance. Then, we design two retrieval modes. Local retrieval ranks entities and time-valid edges in this subgraph and scores their linked chunks to select fine-grained evidence within the query's temporal scope. Global retrieval selects salient time nodes and their summaries to provide global context. The selected contexts are passed to the generator to produce the final answer.

*1) Query-Centric Time Identification:* Given a user query $q$, we prompt an LLM to extract every explicit or implicit temporal expression and determine the full set of timestamps required to answer the question. These timestamps are then aligned to the corresponding nodes in the time hierarchy graph. This process yields a time node set $T^q$.

*2) Dynamic Subgraph Positioning:* We compute the query embedding $\mathbf{e}_q$ and retrieve the top $K$ relation edges from the relation set $\mathcal{E}$, ranked by cosine similarity $\gamma_\varepsilon = \cos(\mathbf{e}_q, \mathbf{e}_\varepsilon)$ with precomputed relation embeddings. The retrieved edges define a query-specific subgraph $\mathcal{G}_K^q = (\mathcal{V}^q, \mathcal{E}^q)$, which serves as the basis for subsequent local and global retrieval. To ensure temporal relevance, we further filter edges whose timestamps lie within the query's temporal scope $T^q$ and gather their neighboring entities, resulting in a temporally focused seed set $\mathcal{V}_t^q$.

*3) Local Retrieval:* We run Personalized PageRank (PPR) on $\mathcal{G}_K^q$ with the temporally filtered seed set $\mathcal{V}_t^q$ as the personalization vector, obtaining a relevance score $s(v)$ for each entity $v \in \mathcal{V}^q$. For each edge $\varepsilon = (v_1, v_2, e, \tau) \in \mathcal{E}^q$, we assign $s(\varepsilon) = \mathbf{1}[\tau \in T^q] (s(v_1) + s(v_2))$. Here, $\mathbf{1}[\cdot]$ denotes the indicator function, which returns 1 when the condition holds and 0 otherwise, so edges with $\tau \notin T^q$ receive zero score. Let $\mathcal{E}(c)$ be the set of edges extracted from chunk $c$ originally, we assign scores for each chunk $c$: $s(c) = w(c) \sum_{\varepsilon \in \mathcal{E}^q} s(\varepsilon)$, where $w(c) = \prod_{\varepsilon \in \mathcal{E}(c)} (1 + \gamma_\varepsilon)$ aggregates the relevance scores of the query and all relation edges to capture the overall semantic importance of all relation edges in chunk $c$. We select the chunks in descending order of $s(c)$ until the token count reaches the predefined context window $L_{\text{ctx}}$ to form the context $\mathcal{C}_q$, and feed it together with the query $q$ to the LLM to generate the final answer $a$. The whole procedure is detailed in Algorithm 1.

---

**Algorithm 1** Local Retrieval

---

**Require:** Query $q$; subgraph $\mathcal{G}_K^q = (\mathcal{V}^q, \mathcal{E}^q)$; temporal seed entities $\mathcal{V}_t^q$; time nodes $T^q$; chunks $\mathcal{C}$; context window $L_{\text{ctx}}$; LLM for QA $\mathcal{M}$

**Ensure:** answer $a$

1: **for all** $v \in \mathcal{V}^q$ **do** $\qquad\qquad\qquad$ ▷ entity scores
2: $\qquad s(v) \leftarrow \text{PPR}(\mathcal{G}_K^q, \mathcal{V}_t^q)$ $\qquad\qquad\qquad$ (1)
3: **for all** $\varepsilon = (v_1, v_2, r, \tau) \in \mathcal{E}^q$ **do** $\qquad$ ▷ edge scores
4: $\qquad s(\varepsilon) \leftarrow \mathbf{1}[\tau \in T^q] (s(v_1) + s(v_2))$ $\qquad$ (2)
5: **for all** $c \in \mathcal{C}$ **do** $\qquad\qquad\qquad$ ▷ chunk scores
6: $\qquad \mathcal{E}(c) \leftarrow \{ \varepsilon \in \mathcal{E}^q \mid \varepsilon \text{ extracted from } c \}$
7: $\qquad s(c) \leftarrow \prod_{\varepsilon \in \mathcal{E}(c)} (1 + \gamma_\varepsilon) \cdot \sum_{\varepsilon \in \mathcal{E}^q} s(\varepsilon)$ $\qquad$ (3)
8: $\mathcal{C}^{\downarrow} \leftarrow \text{SORTBYSCOREDESC}(\mathcal{C}, s)$
9: $\mathcal{C}_q \leftarrow \text{GREEDYPACK}(\mathcal{C}^{\downarrow}, L_{\text{ctx}})$
10: $a \leftarrow \mathcal{M}(\mathcal{C}_q, q)$
11: **return** $a$

---

*4) Global Retrieval:* We start by collecting evidences for answering the query. Evidences are formed by: (i) chunks

TABLE I
COMPARISON OF ECT-QA WITH EXISTING QUESTION ANSWERING DATASETS USED FOR RAG.

| Dataset | Source | Creation | time coverage | multi-hop | time-sensitive | abstract |
|---|---|---|---|---|---|---|
| NaturalQuestions [35] | Wikipedia | real Google queries | before 2018 | ✗ | ✗ | ✗ |
| HotPotQA [36] | Wikipedia | crowdsourcing | before 2017 | ✓ | ✗ | ✗ |
| 2WikiMultiHopQA [37] | Wikipedia&Wikidata | logical rules | before 2010 | ✓ | ✗ | ✗ |
| MultiHop-RAG [38] | News Articles | LLM synthesis | Sep-Dec, 2023 | ✓ | Partial | ✗ |
| TimeQA [13] | Wikipedia&Wikidata | crowdsourcing | before 2021 | ✗ | ✓ | ✗ |
| MenatQA [34] | Wikipedia&Wikidata | crowdsourcing | before 2021 | ✗ | ✓ | ✗ |
| UltraDomain [22] | college textbooks | LLM synthesis | Unspecified | ✗ | ✗ | ✓ |
| **ECT-QA** (ours) | Earnings Call Transcripts | LLM synthesis human review | 2020-2024 | ✓ | ✓ | ✓ |

with the highest scores $s(c)$ as detailed in the local retrieval workflow, and (ii) time reports of each node in $T^q$. Then, each piece of evidence is processed independently by an LLM to extract a set of atomic points $P_e$, where each point $p \in P_e$ is a tuple: $p = (\text{description}, \text{score}, \text{confidence})$. In here, description is a statement containing the key information, score reflects the point's perceived importance to the query, and the confidence is the LLM's self-assessed certainty in the point's accuracy. We then iteratively remove lowest confidence points until the input would fit within context window, and order the points by importance score. Finally, an LLM is tasked to synthesize the points into a single, structured answer. Details of the prompts used for each experiment can be found in the Appendix.

## IV. COMPLEX TEMPORAL QUESTION ANSWERING DATASET FOR RAG

To evaluate time-sensitive QA, we construct ECT-QA, a benchmark of specific and abstract temporal questions derived from Earnings Call Transcripts. The dataset is curated through LLM-assisted synthesis and human review to ensure high factual quality. As shown in Table I, ECT-QA fills a key gap in existing RAG benchmarks by addressing time-sensitive multi-hop and abstract temporal reasoning, which is a critical yet underexplored dimension in current datasets.

### A. Temporal Question Definition

Following the definition in Chen *et al.* [13], we consider a question to be time-sensitive if it includes a temporal specifier such as "in 2023" or "before 2024" and altering this specifier would change the correct answer. We design two types of temporal question.

- Specific multi-hop questions, which are fine-grained and fact-based queries that typically require locating and connecting multiple pieces of facts to arrive at an answer. For example, "*Which quarter saw the highest deferred revenue growth for Autodesk Inc from Q3 2022 to Q3 2023?*"
- Abstract questions, which are query-centric summarization that focus on high-level understanding of the dataset rather than isolated facts. An example is "*How did energy companies navigate cost pressures and enhance profitability across 2024?*"

### B. Corpus Collection

Earnings Call Transcripts (ECTs) serve as the corpus for our study, chosen for their rich temporal characteristics. These transcripts record detailed quarterly financial information for each company and are publicly available for listed companies. We crawl ECTs released between 2020 and 2024 from the Motley Fool platform[2] and retain the companies with transcripts for every quarter in this five-year window. For each transcript, we only keep the "prepared remarks" section and discard the rest. The final corpus contains 480 ECTs from 24 companies in 5 different sectors, with a total of 1.58 million tokens.

### C. Specific Question-Answer Synthesis

We generate a diverse and high-quality temporal multi-hop question-answer dataset with golden evidence through the below pipeline.

*1) Temporal Event Extraction and Alignment:* We begin by prompting GPT-4o-mini[3] to extract temporal events from each document, represented as keyword–timestamp pairs with associated evidence sentences. To align semantically similar but syntactically different keywords, we encode them using SentenceTransformer[4] and apply HDBSCAN clustering [39], enabling alignment and cross-document linking for constructing multi-hop questions.

*2) Multi-Hop Question Generation by Temporal and Reasoning Types:* To construct diverse specific queries, we categorize question generation based on both temporal scope and reasoning type. Specifically, we group 3–8 evidence sentences referring to the same event entity into three temporal scopes: single-time (within one time point), multi-time (across multiple periods), and relative-time (before/after a time point). These grouped evidences are then fed into GPT-4o-mini to generate either enumeration or comparison questions, along with a direct and factual answer. In addition, we synthesize unanswerable questions that involve time periods, entities, or evidence not present in the corpus. This process yields a rich set of time-sensitive questions requiring different levels of multi-hop reasoning.

*3) Automatic and Manual Quality Assurance:* To ensure quality, we filter generated questions using four LLM-assessed criteria: reasoning type match, temporal specificity, evidence

[2]https://www.fool.com/earnings-call-transcripts/
[3]https://platform.openai.com/docs/models/gpt-4o-mini
[4]https://sbert.net/

## TABLE II
## Statistics and Examples of ECT-QA.

| Category | Subcategory | Count |
|----------|-------------|-------|
| Corpus | Documents | 480 |
| | Tokens | 1.58 Million |
| Specific Questions | Total | 1,005 |
| | *Temporal Scope* | |
| | Single-time (in) | 483 |
| | Multi-time (between) | 321 |
| | Relative-time (before/after) | 201 |
| | *Reasoning Type* | |
| | Comparison | 282 |
| | Enumeration | 462 |
| | Unanswerable | 261 |
| Examples | What was EPAM Systems, Inc.'s utilization in each quarter after 2024 Q1? (Relative-time & Enumeration & 3 hops) | |
| | In which quarter did EPAM Systems Inc. achieve the highest GAAP gross margin between 2021 Q2 and 2022 Q1? (Multi-time & Comparison & 4 hops) | |
| Abstract Questions | Total | 100 |
| | Single-time (in) | 43 |
| | Multi-time (between) | 57 |
| Examples | How did companies in the information technology sector, such as Baidu and EPAM, navigate macroeconomic challenges and sector-specific headwinds in 2023 Q1? | |
| | Why did Skechers U.S.A., Inc. achieve record revenue achievements between 2020 and 2022? | |

necessity, and evidence sufficiency. Only questions meeting all criteria are retained. Answers are verified for factual correctness and regenerated if needed. A final round of manual review and refinement ensures the overall quality. This pipeline yields a temporal multi-hop QA dataset grounded in golden evidence.

### D. Abstract Question Synthesis

To synthesize realistic and diverse queries, we construct abstract questions by first simulating potential user profiles with information needs, and then generating temporal queries that reflect authentic analytical intents.

*1) Potential User Simulation:* Inspired by prior work [11], we simulate 10 potential user personas who might interact with the corpus. Given a description of the corpus, we ask the LLM to generate descriptive user profiles, each associated with a distinct information need and analytical objective.

*2) User-Guided Abstract Question Generation:* We group multiple ECT documents based on shared metadata, such as time period, company, or sector. On average, each group contains 18.22 documents. To address input length constraints, we first use LLM to generate concise summaries for each individual document. Given a user profile and the corresponding set of document summaries, we encourage LLM to synthesize deep questions that such a user might pose when analyzing the content, with a particular focus on temporal "how" and "why" questions.

*3) Automatic and Manual Quality Assurance:* Similar to the above procedure, we first perform an automatic quality check by prompting the LLM to evaluate each question along

four dimensions: clarity, temporal grounding, analytical depth, and answerability. Questions that pass this validation are then manually reviewed to ensure quality. Finally, we construct 100 abstract questions. Table II summarizes the overall dataset statistics and provides illustrative examples.

## V. Experimental Setup

This section presents the experimental setup, covering the incremental evaluation protocol, baseline methods, evaluation metrics, and implementation details.

### A. Incremental RAG Evaluation

To simulate the dynamic nature of real-world knowledge environments where RAG systems must continuously adapt to evolving information, we design an incremental evaluation protocol by partitioning both the corpus and the query set into two temporal slices. Specifically, for our ECT-QA dataset, we establish a base corpus consisting of all 384 documents from 2020 to 2023, and a new corpus from 2024 containing 96 documents as incremental updates. In parallel, we partition queries according to the temporal scope of the facts required for answering them: base queries, whose answers rely only on facts from 2020 to 2023, and new queries, whose answers require facts spanning 2020 to 2024. In total, the dataset includes 1,105 specific queries with 656 base queries and 349 new queries, and 100 abstract queries containing 72 base abstract queries and 28 new abstract queries. This design allows us to evaluate RAG performance under corpus growth from three perspectives:

- **Base queries on the base corpus**: evaluate system performance under the base corpus setting.
- **Base queries on the updated corpus**: assess consistency and robustness of RAG systems when new corpus is injected.
- **New queries on the updated corpus**: measure adaptability of RAG systems in leveraging newly added knowledge.

Overall, this incremental evaluation captures both the stability of RAG systems on previously seen queries and their ability to handle new queries grounded in evolving knowledge, closely reflecting real-world deployment scenarios.

### B. Baselines

We compare our approach with the following state-of-the-art methods:

- LLM-GT, which concatenates the question with its gold evidence passages and feeds them to the LLM to generate the answer. This approximates an oracle upper bound for the given LLM generator.
- NaiveRAG [2], a standard RAG baseline which chunk the corpus, encode passages with an embedding model, and store embedding vectors into a vector base. At test time, the query is embedded and top-$k$ passages are retrieved by similarity, and the LLM generates the answer from the retrieved context.
- QD-RAG [40], which first decomposes the original question into simpler sub-queries, and then runs retrieval

independently for each sub-query. The retrieved paragraphs are aggregated together with the original question to the LLM to produce the final answer.

- GraphRAG [11], a graph-enhanced RAG pipeline that uses an LLM to extract entities and relations from the corpus, clusters nodes into communities, and generates community reports to capture global context for retrieval and generation.
- LightRAG [23], integrates a graph-based text indexing paradigm with a dual-level retrieval scheme to capture both low- and high-level signals for efficient retrieval.
- HippoRAG2 [25], first builds an open KG by extracting triples from passages with an LLM, and then applies personalized PageRank over the KG to retrieve nodes and maps them back to passages for the final LLM answer.

### C. Evaluation Metrics

To assess the performance of specific queries, we design three LLM-based metrics to automatically judge the factual accuracy of model responses. A powerful LLM (specifically, GPT-4o-mini) serves as the judge, performing a fine-grained, element-wise comparison between the model's prediction and the ground-truth answer, with reference to the provided evidence. Specifically, we utilize this automated framework to measure the proportion of: (1) **Correct** elements, factual claims that are accurately supported by the provided evidence and match the required temporal scope; (2) **Refusal** elements, instances where the model explicitly acknowledges its inability to answer due to lack of evidence; and (3) **Incorrect** elements, responses containing wrong, unsupported, or hallucinated information. These three metrics sum to 1 for each query evaluation. A good QA system is expected to achieve high Correct score and low Incorrect score. Crucially, we expect the model to appropriately refuse answering rather than providing incorrect information [41]. This refusal behavior represents a safety-aware approach that is particularly valuable in high-stakes domains. Besides, we adopt two widely used non-LLM metrics, including **ROUGE-L** and $\mathbf{F_1}$. ROUGE-L[5] evaluates the overlap of the longest common subsequence between generated and reference answers, and $F_1$ measures the harmonic mean of precision and recall at the token level.

Following previous studies [11], [23], we conduct an LLM-based multi-dimensional comparison method for evaluating abstract queries. Specifically, given the answers from two different methods, we employ GPT-4o-mini as a judge to evaluate answer pairs across three key dimensions: **Comprehensiveness**, the depth and detail of information provided in the answer, **Diversity**, the variety of perspectives and insights offered, and **Temporal Coverage**, the accuracy and completeness in handling time-related aspects of the query. For each dimension, the LLM judge selects a winner between two answers and provides detailed explanations. Based on these criteria, an **Overall Winner** is determined to provide a holistic assessment of answer quality. We calculate win rates accordingly, ultimately leading to the final results. All prompts used for evaluation are provided in the Appendix to ensure reproducibility.

[5]https://github.com/google-research/google-research/tree/master/rouge

### D. Implementation Details

In our experiments, the chunk size is set to 1,200 tokens with an overlap of 100 tokens. For local retrieval, we retrieve the top-$K$=20 candidate relations. To ensure a fair comparison, our approach and all baseline models utilize the Gemini-2.5-flash-lite[6] model for the indexing phase, and the GPT-4o-mini model for the query answering phase, unless otherwise specified. For text embedding, we adopt the text-embedding-3-small[7] model. The total length of retrieved content is limited to 12,000 tokens for local retrieval and 24,000 tokens for global retrieval, with 10% of the token budget allocated to chunks and 90% to time reports.

### VI. RESULTS AND DISCUSSION

In this section, we conduct comprehensive experiments to answer the following questions:

- How does our method perform under different evaluation settings? We compare our approach with baselines on both specific queries and abstract queries across three evaluation settings, which demonstrates the effectiveness and efficiency of our method. Detailed in Section VI-A.
- Why does our method perform well? We conduct ablation studies and graph visualizations to better understand the performance gains of our method, demonstrating the importance of temporal graph modeling and time-aware retrieval. The findings are discussed in Section VI-B.
- How well does our method generalize across different LLMs and datasets? We evaluate the generalization of our approach across different LLMs for question answering, and further extend the analysis to generic QA benchmarks beyond the ECT-QA dataset to verify its applicability. See Section VI-C.

### A. How does our method perform?

This sections present and analyze the main results on both specific and abstract question answering tasks, demonstrating our method's superior performance across different question complexities.

*1) Results of Specific Question Answering:* Following the evaluation settings introduced in Section V-A, we first evaluate the performance of different methods on the base queries using the base corpus (2020–2023). Table III summarizes the performance of all methods based on LLM-based factual accuracy metrics (Correct, Refusal, and Incorrect), traditional lexical overlap metrics (ROUGE-L and $F_1$), as well as indexing cost measured by the number of prompt and completion tokens consumed during LLM-based indexing. We observe that LLM-GT achieves very high performance across all metrics, with a Correct ratio of 0.902 and an Incorrect ratio of 0.075. This demonstrates the high factual quality of the ECT-QA dataset, confirming its suitability for reliable evaluation. Compared to the other baselines, our method achieves the highest Correct score of 0.599 and the lowest Incorrect score of 0.210,

[6]https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-flash-lite

[7]https://platform.openai.com/docs/models/text-embedding-3-small

TABLE III
THE PERFORMANCE OF SPECIFIC QUESTION ANSWERING: BASE QUERIES ON THE BASE CORPUS.

| Model | Index Token Cost | | LLM Metrics | | | Non-LLM Metrics | |
| | Prompt ↓ | Completion ↓ | Correct ↑ | Refusal | Incorrect ↓ | ROUGE-L ↑ | $F_1$ ↑ |
|---|---|---|---|---|---|---|---|
| LLM-GT | – | – | 0.902 | 0.023 | 0.075 | 0.626 | 0.647 |
| NaiveRAG | – | – | 0.385 | 0.325 | 0.290 | 0.375 | 0.366 |
| QD-RAG | – | – | 0.380 | 0.329 | 0.291 | 0.369 | 0.359 |
| GraphRAG [11] | 37.1M | 17.7M | 0.405 | 0.280 | 0.315 | 0.371 | 0.375 |
| LightRAG [23] | 17.1M | 9.0M | 0.406 | 0.160 | 0.434 | 0.350 | 0.359 |
| HippoRAG2 [25] | 3.2M | 1.1M | 0.410 | 0.345 | 0.245 | 0.382 | 0.385 |
| **TG-RAG (Ours)** | 6.3M | 7.1M | **0.599** | 0.191 | **0.210** | **0.493** | **0.490** |

TABLE IV
THE PERFORMANCE OF SPECIFIC QUESTION ANSWERING: BASE QUERIES AND NEW QUERIES ON THE UPDATED CORPUS.

| Model | Index Token Cost | | Base Queries | | | | | New Queries | | | | |
| | | | LLM Metrics | | | Non-LLM Metrics | | LLM Metrics | | | Non-LLM Metrics | |
| | Prompt ↓ | Completion ↓ | Correct ↑ | Refusal | Incorrect ↓ | ROUGE-L ↑ | $F_1$ ↑ | Correct ↑ | Refusal | Incorrect ↓ | ROUGE-L ↑ | $F_1$ ↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LLM-GT | – | – | 0.902 | 0.023 | 0.075 | 0.626 | 0.647 | 0.874 | 0.044 | 0.082 | 0.631 | 0.649 |
| NaiveRAG | – | – | 0.366 | 0.347 | 0.287 | 0.352 | 0.343 | 0.390 | 0.406 | 0.205 | 0.384 | 0.356 |
| QD-RAG | – | – | 0.362 | 0.374 | 0.264 | 0.354 | 0.344 | 0.407 | 0.355 | 0.238 | 0.413 | 0.386 |
| GraphRAG | 30.0M | 7.8M | 0.380 | 0.327 | 0.293 | 0.370 | 0.369 | 0.398 | 0.368 | 0.234 | 0.392 | 0.389 |
| LightRAG | 4.6M | 2.6M | 0.386 | 0.203 | 0.412 | 0.345 | 0.355 | 0.382 | 0.229 | 0.389 | 0.354 | 0.348 |
| HippoRAG2 | 0.8M | 0.3M | 0.399 | 0.347 | 0.253 | 0.376 | 0.372 | 0.372 | 0.492 | 0.136 | 0.367 | 0.354 |
| **TG-RAG (Ours)** | 1.6M | 2.0M | **0.587** | 0.193 | **0.220** | **0.483** | **0.475** | **0.617** | 0.216 | **0.167** | **0.524** | **0.487** |

TABLE V
THE PERFORMANCE OF ABSTRACT QUESTION ANSWERING. *Abbreviations: COMP.=COMPREHENSIVENESS, DIV.=DIVERSITY, TEMP. =TEMPORAL COVERAGE.*

| Model | GraphRAG | **Ours** | HippoRAG | **Ours** | LightRAG | **Ours** |
|---|---|---|---|---|---|---|
| *(1) Base queries on the base corpus* | | | | | | |
| Comp. | 0.167 | 0.833 | 0.167 | 0.833 | 0.014 | 0.986 |
| Div. | 0.486 | 0.514 | 0.431 | 0.569 | 0.097 | 0.903 |
| Temp. Cov. | 0.111 | 0.889 | 0.236 | 0.764 | 0.014 | 0.986 |
| Overall | 0.167 | 0.833 | 0.181 | 0.819 | 0.014 | 0.986 |
| *(2) Base queries on the updated corpus* | | | | | | |
| Comp. | 0.222 | 0.778 | 0.167 | 0.833 | 0.028 | 0.972 |
| Div. | 0.472 | 0.528 | 0.375 | 0.625 | 0.042 | 0.958 |
| Temp. Cov. | 0.111 | 0.889 | 0.264 | 0.736 | 0.069 | 0.972 |
| Overall | 0.236 | 0.764 | 0.181 | 0.819 | 0.028 | 0.972 |
| *(3) New queries on the updated corpus* | | | | | | |
| Comp. | 0.107 | 0.893 | 0.214 | 0.786 | 0.000 | 1.000 |
| Div. | 0.357 | 0.643 | 0.536 | 0.464 | 0.071 | 0.929 |
| Temp. Cov. | 0.000 | 1.000 | 0.321 | 0.679 | 0.107 | 0.893 |
| Overall | 0.107 | 0.893 | 0.250 | 0.750 | 0.000 | 1.000 |

indicating superior factual accuracy. In terms of cost efficiency during indexing phase, our method maintains competitive token cost with 6.3 million prompt tokens and 7.1 million completion tokens, which are significantly lower than GraphRAG and LightRAG. It is worth noting that HippoRAG2's low indexing cost comes at the expense of semantic richness, as it constrains the LLM to output only entity and relation names, omitting detailed descriptions and high-level summaries. In contrast, our method achieves a more balanced trade-off, which avoids exhaustive graph construction while achieving high retrieval performance.

Next, we inject the new corpus from 2024 into the existing database to simulate the real-world scenario where RAG systems have to incrementally adapt to newly emerging knowledge. We record the update cost and evaluate all methods on both the base queries and the new queries. The results are summarized in Table IV. TG-RAG achieves the best performance on both base and new queries, and maintains almost consistent performance on base queries between the original corpus and the updated corpus, with only a negligible drop compared to other methods. This indicates that our approach effectively models time-sensitive knowledge in scenarios with continuously evolving information. Moreover, in contrast to GraphRAG which incurs a high indexing cost of 30.0 million prompt and 7.8 million completion tokens due to complete graph summary regeneration, our method requires only 1.6 million prompt and 2.0 million completion tokens by updating the summaries for new time points and their corresponding ancestor nodes, resulting in a significantly lower update cost. Overall, these findings confirm that our method adapts well to evolving knowledge, preserving retrieval accuracy while efficiently integrating new information, which is a crucial capability for real-world QA systems that evolve over time.

*2) Results of Abstract Question Answering:* We evaluate the performance of different methods on abstract question answering, as shown in Table V. For the dimension of Temporal Coverage, our approach achieves remarkable win rates of 0.889 against GraphRAG, 0.764 against HippoRAG2, and 0.986 against LightRAG on base queries over the base corpus. This strong performance in temporal aspects directly reflects the effectiveness of our temporal-aware graph construction and retrieval framework. For Comprehensiveness and Diversity dimension, our method also shows substantial advantages with consistently high win rates, highlighting its ability to generate detailed and diverse answers.

TABLE VI
A REPRESENTATIVE CASE COMPARING GRAPHRAG TO OUR METHOD ON ABSTRACT QUESTION ANSWERING.

| Query | Why did Iron Mountain Incorporated maintain consistent EBITDA margin growth from 2020 Q1 to 2022 Q4? |
|---|---|
| **GraphRAG** | This performance is supported by a diversified business model encompassing physical storage, records and information management (RIM), data centers, digital solutions, and asset lifecycle management (ALM). ### Key Drivers of Performance **Data Center Business Growth:** IRM's Global Data Center Business has been a significant contributor to its financial success. This segment has experienced robust growth, with strong leasing activity and increased demand for bandwidth, partly driven by remote working trends. The adjusted EBITDA margin for the data center business has been reported at a healthy 45.8%, with year-on-year growth reaching 21% and 25% in different periods. Financial strategies such as capital recycling and sale-leaseback transactions have been employed to fuel the expansion of this business. **Productivity and Mix Improvements:** Cost productivity and a strong mix of business have been identified as key drivers for improving the Adjusted EBITDA margin across operations. These factors directly influence profitability and are central to the company's strategy for enhancing financial health... **Challenges and Offsetting Factors:** While the company has shown improvements, certain factors have partially offset gains in the Adjusted EBITDA margin. These include fixed cost deleverage, which occurs when lower service revenue leads to less efficient utilization of fixed costs, and higher bonus compensation accrual... |
| **TG-RAG** | Various factors played a vital role, including capital recycling programs, digital transformations, and innovative projects that enhanced operational efficiency. **In 2020**, Iron Mountain faced the initial pandemic challenges but simultaneously began laying the groundwork for future growth. The company emphasized cost management strategies and operational efficiency improvements. ### Key Initiatives **Cost Management**: Strategic projects that focused on reducing operational costs started to take effect. **Capital Recycling**: A strong approach to capital recycling set the stage for improved financial flexibility. **By 2021**, Iron Mountain had escalated its initiatives focusing on its data center business and digital transformation. This included the launch of innovative products like Smart Sort and Secure ITAD... **As 2022 progressed**, Iron Mountain reported notable financial performance improvements, achieving $1.28 billion in revenue during Q4 2022... the adjusted EBITDA was positively influenced by innovations and strategic partnerships, such as those established with a global content delivery network. The continued enhancement of core services, including data center offerings, supported margin expansion with notable year-on-year growth in service revenues. ## **Conclusion**: From 2020 Q1 to 2022 Q4, Iron Mountain Incorporated's consistent EBITDA margin growth can be attributed to a multifaceted strategy that combined operational excellence, strategic project execution, and a robust pipeline of innovative products and customer contracts. By focusing on cost management, capital recycling, and expanding digital infrastructures, ..., resulting in a strong EBITDA growth narrative. |
| **LLM Decision** | **Comprehensiveness**: Winner: TG-RAG; **Diversity**: Winner: TG-RAG; **Temporal Coverage**: Winner: TG-RAG; **Overall Winner**: Winner: TG-RAG |

TABLE VII
ABLATION STUDY ON SPECIFIC QUESTION ANSWERING (BASE QUERIES ON THE BASE CORPUS).

| Temporal Retrieval | PPR Ranking | Temporal Indexing | LLM Metrics | | | Non-LLM Metrics | |
|---|---|---|---|---|---|---|---|
| | | | Correct ↑ | Refusal | Incorrect ↓ | ROUGE-L ↑ | $F_1$ ↑ |
| ✓ | ✓ | ✓ | **0.599** | 0.191 | **0.210** | **0.493** | **0.490** |
| ✓ | ✗ | ✓ | 0.580 | 0.223 | 0.197 | 0.483 | 0.472 |
| ✗ | ✓ | ✓ | 0.382 | 0.423 | 0.195 | 0.376 | 0.356 |
| ✗ | ✗ | ✓ | 0.482 | 0.294 | 0.223 | 0.434 | 0.416 |
| ✗ | ✗ | ✗ | 0.381 | 0.458 | 0.161 | 0.359 | 0.345 |

To illustrate the differences qualitatively, Table VI presents a representative case comparing GraphRAG and our method. When answering "Why did Iron Mountain Incorporated maintain consistent EBITDA margin growth from 2020 Q1 to 2022 Q4?", our TG-RAG method demonstrates clear superiority in temporal reasoning by explicitly organizing the explanation along the temporal progression from 2020 to 2022. This temporal structuring enables comprehensive coverage of the entire growth narrative while maintaining logical coherence. The LLM judge ranks our answer higher across all dimensions, proving the effectiveness of our method.

### B. Why does our method perform well?

To understand the strong performance of our method, we conduct a detailed analysis from both quantitative and qualitative perspectives. First, we perform ablation studies to quantitatively evaluate the contribution of each core component in our framework. Second, we present a graph visualization with statistical comparisons to illustrate how our bi-level temporal structure effectively represents temporal facts while maintaining structural compactness.

*1) Ablation Study:* To evaluate the contribution of each core component, we conduct an ablation study on specific question answering using the base queries over the base corpus. Our model consists of three key modules: Temporal retrieval performs query-centric time filtering to select facts within the temporal scope $T^q$, PPR Ranking applies Personalized PageRank on the subgraph to propagate relevance scores from the seed entities and obtain more informative entity-level rankings; and Temporal Indexing, which builds time-stamped relation quadruples for precise evidence representation and retrieval. Four ablation variants are compared against the full model.

- *w/o PPR Ranking*, which removes graph propagation and directly assigns edge scores using relation and query similarity. Specifically, the original edge scoring function in Eq. (2) is replaced by $s(\varepsilon) = \mathbf{1}[\tau \in T^q] \gamma_\varepsilon$, for each $\varepsilon = (v_1, v_2, r, \tau) \in \mathcal{E}^q$.
- *w/o Temporal Retrieval*, which disables query-centric time

Question 1: What was Western Digital Corporation's revenue in each quarter from 2023 Q1 to Q3?

Question 2: What were Western Digital Corporation's operating cash flow, gross debt outstanding, and earnings per share in 2020 Q3?
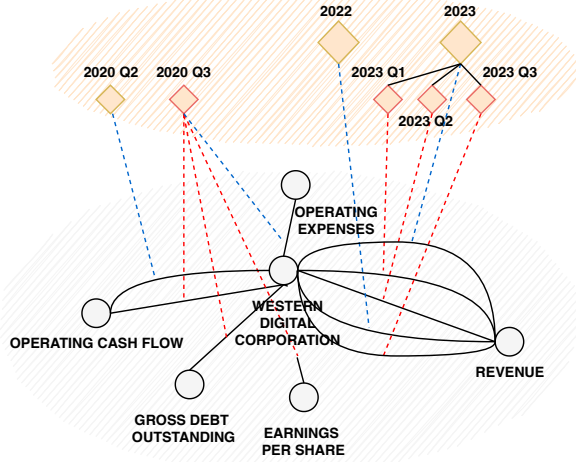


Fig. 3. Visualization of a subgraph sampled from our constructed bi-level temporal graph. Dashed lines represent the connections between upper-layer temporal nodes and lower-layer relational edges, where red lines indicate the temporal clues used to answer given query examples, and blue lines denote others present in the graph.

identification, running standard PageRank over all nodes to calculate entity scores and defining edge scores without temporal constraints. The original Eq. (1) and Eq. (2) is modified to $s(v) = \mathrm{PR}\left(\mathcal{G}_K^q\right)$, $s(\varepsilon) = s(v_1) + s(v_2)$.

- *w/o Temporal Retrieval + PPR*, which removes both time filtering and graph propagation, relying solely on relation and query similarity. The original edge scoring function in Eq. (2) is replaced by $s(\varepsilon) = \gamma_\varepsilon$.
- *w/o All*, which further disables the temporal graph indexing module, degrading the graph construction from temporal to static. Retrieval in this setting is performed over a non-temporal graph without timestamped relations.

The experimental results are shown in Table VII. First, removing only PPR Ranking results in a relatively modest performance drop compared to the full model, indicating that graph propagation helps refine local relevance. Second, when Temporal Retrieval is disabled, Correct score drops substantially from 0.599 to 0.382 and Refusal rate increases significantly to 0.423, highlighting the importance of temporal filtering for handling time-sensitive queries. Without proper time scoping, the retrieval system is overwhelmed by temporally irrelevant evidence, which subsequently misleads the answer generation process. Interestingly, the variant removing both Temporal Retrieval and PPR Ranking performs better than removing only Temporal Retrieval. This suggests that without time filtering, regular pagerank would propagate noisiness and dilute relevance, leading to poorer chunk scores and worse answers. The consistent superiority of the full model demonstrates that the three components work synergistically: temporal retrieval ensures relevant fact positioning, PPR propagates these relevant signals effectively, and temporal indexing provides the neces-

TABLE VIII
COMPARISON OF GRAPH SIZES CONSTRUCTED BY DIFFERENT METHODS, MEASURED BY THE NUMBER OF ENTITIES AND RELATIONS EXTRACTED FROM THE BASE CORPUS AND THE UPDATED CORPUS.

| Model | Base Corpus | | Updated Corpus | |
|---|---|---|---|---|
| | #entity | #relation | #entity | #relation |
| GraphRAG | 45,540 | 59,679 | 54,854 | 71,989 |
| LightRAG | 35,003 | 42,892 | 42,514 | 53,383 |
| HippoRAG2 | 19,318 | 24,205 | 22,128 | 28,153 |
| Ours | 16,817 | 28,157 | 20,110 | 34,671 |

sary structural foundation for precise evidence representation.

*2) Graph Visualization and Statistics:* To intuitively illustrate how our method repesents and leverages temporal information, we visualize a small subgraph extracted from the constructed bi-level temporal graph, as shown in Figure 3. The detailed edge descriptions are provided in the Appendix. The upper layer (orange region) represents time nodes, while the lower layer (gray region) depicts entities and their associated relations. A key advantage of this design is that the graph can represent the same factual relation across different time periods via multiple time-stamped edges. For example, the relation between Western Digital Corporation and Revenue can appear under 2023 Q1, Q2, and Q3 as distinct temporal instances. This fine-grained temporal modeling effectively captures the factual evolution of the same knowledge across time, allowing the model to distinguish between time-specific evidence and thus prevent temporal confusion during retrieval and answer generation. For instance, given the query *"What was Western Digital Corporation's revenue in each quarter from 2023 Q1 to Q3?"*, our model first identifies the temporal nodes 2023 Q1, 2023 Q2, and 2023 Q3, and retrieves the respective company–revenue relations, yielding accurate results. Similarly, for the query *"What were Western Digital Corporation's operating cash flow, gross debt outstanding, and earnings per share in 2020 Q3?"*, our model first identifies the temporal node 2020 Q3 and successfully locates the correct facts. These visualized examples clearly demonstrate that the proposed temporal graph structure enables effective time-aware retrieval and precise retrieval for time sensitive question answering.

As shown in Table VIII, our method encodes the same corpus with substantially fewer entities and relations, resulting in lower storage costs and faster retrieval. This compactness stems from our temporal design, which shares entity nodes across time while introducing timestamped connections to capture evolution, thus effectively eliminating structural redundancy.

## C. Applicability to Alternative LLMs and Generic QA Datasets

To further examine the generalization of our approach, we apply it to different LLMs and extend the evaluation to generic QA datasets. Table IX reports the performance of various QA LLMs on base queries over the base corpus. Our method maintains strong performance across diverse model architectures, notably achieving a 0.625 Correct score with DeepSeek-R1. It is worth noting that even when using the open-source Llama-3.3-70B-Instruct, our approach still outperforms other baselines built upon GPT-4o-mini. The effectiveness

TABLE IX
EVALUATION OF DIFFERENT LLMs FOR QUERY ANSWERING: BASE
QUERIES ON THE BASE CORPUS.

| QA LLM | Correct | Refusal | Incorrect | ROUGE-L | $F_1$ |
|---|---|---|---|---|---|
| GPT-4o-mini | 0.599 | 0.191 | 0.210 | 0.493 | 0.490 |
| Llama-3.3-70B-Instruct | 0.466 | 0.373 | 0.160 | 0.419 | 0.403 |
| Qwen3-235B-A22B-Instruct-2507 | 0.619 | 0.200 | 0.180 | 0.500 | 0.482 |
| DeepSeek-R1 | 0.625 | 0.174 | 0.201 | 0.493 | 0.468 |

TABLE X
QA PERFORMANCE ($F_1$ SCORES) ACROSS DIFFERENT RAG DATASETS
USING LLAMA-3.3-70B-INSTRUCT AS THE INDEX AND QA LLMs.
NV-EMBED-V2 ARE ADOPTED AS THE EMBEDDING MODEL, FOLLOWING
THE SETTING IN [25]. RESULTS MARKED WITH † ARE REPORTED BY [25].

| Model | NaturalQuestions | 2WikiMultihopQA | HotpotQA |
|---|---|---|---|
| NaiveRAG† | 0.619 | 0.615 | 0.753 |
| GraphRAG† | 0.469 | 0.586 | 0.686 |
| LightRAG† | 0.166 | 0.116 | 0.024 |
| HippoRAG2† | 0.633 | 0.710 | 0.755 |
| Ours | **0.648** | **0.719** | **0.757** |

across both proprietary and open-source models indicates that our approach does not critically depend on specific LLM implementations, making it flexible and accessible for different practical deployments.

In addition, we evaluate the general applicability of our approach on three widely-used RAG datasets, NaturalQuestions [35], 2WikiMultihopQA [37], and HotpotQA [36]. To ensure a consistent and fair comparison, we follow the configuration of HippoRAG2 [25] by using Llama-3.3-70B-Instruct for both indexing and question answering, and NV-Embed-v2 [42] as the embedding model. For graph indexing, we prompt the model to extract not only temporal quadruples but also non-temporal factual triples. During retrieval, we remove temporal filtering when the query lacks a clearly identifiable time scope, enabling our method to operate robustly across generic QA tasks. As shown in Table X, our method achieves $F_1$ scores of 0.648 on NaturalQuestions, 0.719 on 2WikiMultihopQA, and 0.757 on HotpotQA, confirming its broad applicability and robustness beyond time-sensitive QA.

## VII. CONCLUSION

In this paper, we identify and address a critical yet often overlooked challenge in retrieval-augmented generation: managing the temporal dynamics of real-world knowledge. We propose Temporal GraphRAG (TG-RAG), a novel framework that represents external corpora as bi-level temporal graphs, where timestamped relations capture factual evolution and hierarchical time summaries encode multi-granular trends. The framework incorporates time-aware retrieval strategies that dynamically select relevant temporal contexts to enhance answer accuracy. Furthermore, we contribute a new benchmark ECT-QA and an incremental evaluation protocol that simulates realistic corpus growth, allowing assessment of retrieval accuracy and update efficiency. Extensive experiments demonstrate the effectiveness and efficiency of our method. We believe this work provides a concrete step toward making RAG systems temporally aware and practically deployable in dynamic environments.

## REFERENCES

[1] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, "Retrieval-augmented generation for knowledge-intensive NLP tasks," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2020.

[2] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, Q. Guo, M. Wang, and H. Wang, "Retrieval-augmented generation for large language models: A survey," *arXiv preprint arXiv:2312.10997*, 2023.

[3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2020.

[4] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[5] D. Guo, D. Yang, H. Zhang, J. Song, P. Wang, Q. Zhu, R. Xu, R. Zhang, S. Ma, X. Bi, and et al., "Deepseek-r1 incentivizes reasoning in llms through reinforcement learning," *Nature*, vol. 645, no. 8081, pp. 633–638, September 2025.

[6] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, and et al., "Qwen3 technical report," *arXiv preprint arXiv:2505.09388*, 2025.

[7] K. Shuster, S. Poff, M. Chen, D. Kiela, and J. Weston, "Retrieval augmentation reduces hallucination in conversation," in *Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2021*, 2021, pp. 3784–3803.

[8] J. Song, X. Wang, J. Zhu, Y. Wu, X. Cheng, R. Zhong, and C. Niu, "RAG-HAT: A hallucination-aware tuning pipeline for LLM in retrieval-augmented generation," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 1548–1558.

[9] K. Bhushan, Y. Nandwani, D. Khandelwal, S. Gupta, G. Pandey, D. Raghu, and S. Joshi, "Systematic knowledge injection into large language models via diverse augmentation for domain-specific RAG," in *Proceedings of the Findings of the Association for Computational Linguistics: NAACL*, 2025, pp. 5922–5943.

[10] J. Ouyang, T. Pan, M. Cheng, R. Yan, Y. Luo, J. Lin, and Q. Liu, "Hoh: A dynamic benchmark for evaluating the impact of outdated information on retrieval-augmented generation," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2025, pp. 6036–6063.

[11] D. Edge, H. Trinh, N. Cheng, J. Bradley, A. Chao, A. Mody, S. Truitt, and J. Larson, "From local to global: A graph RAG approach to query-focused summarization," *arXiv preprint arXiv:2404.16130*, 2024.

[12] S. Shang, X. Cheng, Y. Xiong, F. Guo, S. Gao, X. Chen, F. Wang, Y. Wang, D. Zhao, and R. Yan, "Personalized Review Summarization by Using Graph-based Retrieval Augmemted Generation," *IEEE Transactions on Knowledge and Data Engineering*, no. 01, pp. 1–16, 2025.

[13] W. Chen, X. Wang, and W. Y. Wang, "A dataset for answering time-sensitive questions," in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, 2021.

[14] F. Zhang, Z. Zhang, F. Zhuang, Y. Zhao, D. Wang, and H. Zheng, "Temporal knowledge graph reasoning with dynamic memory enhancement," *IEEE Transactions on Knowledge and Data Engineering*, vol. 36, no. 11, pp. 7115–7128, 2024.

[15] N. Kanhabua and A. Anand, "Temporal information retrieval," in *Proceedings of the International ACM SIGIR conference on Research and Development in Information Retrieval*, 2016, pp. 1235–1238.

[16] A. Abdallah, B. Piryani, J. Wallat, A. Anand, and A. Jatowt, "Extending dense passage retrieval with temporal information," *arXiv preprint arXiv:2502.21024*, 2025.

[17] A. Jimeno-Yepes, Y. You, J. Milczek, S. Laverde, and R. Li, "Financial report chunking for effective retrieval augmented generation," *arXiv preprint arXiv:2402.05131*, 2024.

[18] Y. Wei, J. Han, and Y. Yang, "Adapting general-purpose embedding models to private datasets using keyword-based retrieval," in *Proceedings of the Findings of the Association for Computational Linguistics, ACL*, 2025, pp. 6856–6870.

[19] C. Chan, C. Xu, R. Yuan, H. Luo, W. Xue, Y. Guo, and J. Fu, "RQ-RAG: learning to refine queries for retrieval augmented generation," *arXiv preprint arXiv:2404.00610*, 2024.

[20] A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi, "Self-rag: Learning to retrieve, generate, and critique through self-reflection," in *Proceedings of the International Conference on Learning Representations*, 2024.

[21] A. V. Duarte, J. Marques, M. Graça, M. Freire, L. Li, and A. L. Oliveira, "Lumberchunker: Long-form narrative document segmentation," in *Proceedings of the Findings of the Association for Computational Linguistics: EMNLP*, 2024, pp. 6473–6486.

[22] H. Qian, P. Zhang, Z. Liu, K. Mao, and Z. Dou, "Memorag: Moving towards next-gen RAG via memory-inspired knowledge discovery," *arXiv preprint arXiv:2409.05591*, 2024.

[23] Z. Guo, L. Xia, Y. Yu, T. Ao, and C. Huang, "Lightrag: Simple and fast retrieval-augmented generation," *arXiv preprint arXiv:2410.05779*, 2024.

[24] B. J. Gutierrez, Y. Shu, Y. Gu, M. Yasunaga, and Y. Su, "Hipporag: Neurobiologically inspired long-term memory for large language models," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2024.

[25] B. J. Gutiérrez, Y. Shu, W. Qi, S. Zhou, and Y. Su, "From RAG to memory: Non-parametric continual learning for large language models," in *Proceedings of the International Conference on Machine Learning*, 2025.

[26] Y. Pei, J. Zheng, and J. Cartlidge, "Dynamic graph representation with contrastive learning for financial market prediction: Integrating temporal evolution and static relations," in *Proceedings of the International Conference on Agents and Artificial Intelligence*, 2025, pp. 298–309.

[27] L. A. Khan, "Temporality of law," *McGeorge Law Review*, vol. 40, p. 55, 2009.

[28] B. M. Stewart, "Chronolawgy: A study of law and temporal perception," *University of Miami Law Review*, vol. 67, p. 303, 2012.

[29] S. Turner and D. P. Fernandez, ""everything was much more dynamic": Temporality of health system responses to covid-19 in colombia," *Plos one*, vol. 19, no. 9, p. e0311023, 2024.

[30] J. Son and A. Oh, "Time-aware representation learning for time-sensitive question answering," in *Proceedings of the Findings of the Association for Computational Linguistics: EMNLP*, 2023, pp. 70–77.

[31] X. Zhu, C. Yang, B. Chen, S. Li, J. Lou, and Y. Yang, "Question answering as programming for solving time-sensitive questions," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 12 775–12 790.

[32] W. Yang, Y. Li, M. Fang, and L. Chen, "Enhancing temporal sensitivity and reasoning for time-sensitive question answering," in *Proceedings of the Findings of the Association for Computational Linguistics: EMNLP*, 2024, pp. 14 495–14 508.

[33] Z. Siyue, X. Yuxiang, Z. Yiming, W. Xiaobao, L. A. Tuan, and Z. Chen, "MRAG: A modular retrieval framework for time-sensitive question answering," *arXiv preprint arXiv:2412.15540*, 2024.

[34] Y. Wei, Y. Su, H. Ma, X. Yu, F. Lei, Y. Zhang, J. Zhao, and K. Liu, "Menatqa: A new dataset for testing the temporal comprehension and reasoning abilities of large language models," in *Proceedings of the Findings of the Association for Computational Linguistics: EMNLP*, 2023, pp. 1434–1447.

[35] T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. P. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, K. Toutanova, L. Jones, M. Kelcey, M. Chang, A. M. Dai, J. Uszkoreit, Q. Le, and S. Petrov, "Natural questions: a benchmark for question answering research," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 452–466, 2019.

[36] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning, "Hotpotqa: A dataset for diverse, explainable multi-hop question answering," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 2369–2380.

[37] X. Ho, A. D. Nguyen, S. Sugawara, and A. Aizawa, "Constructing A multi-hop QA dataset for comprehensive evaluation of reasoning steps," in *Proceedings of the International Conference on Computational Linguistics*, 2020, pp. 6609–6625.

[38] Y. Tang and Y. Yang, "Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries," in *Proceedings of the Conference on Language Modeling*, 2024.

[39] L. McInnes and J. Healy, "Accelerated hierarchical density based clustering," in *Proceedings of the IEEE International Conference on Data Mining Workshops*, 2017, pp. 33–42.

[40] P. J. L. Ammann, J. Golde, and A. Akbik, "Question decomposition for retrieval-augmented generation," *arXiv preprint arXiv:2507.00355*, 2025.

[41] A. T. Kalai, O. Nachum, S. S. Vempala, and E. Zhang, "Why language models hallucinate," *arXiv preprint arXiv:2509.04664*, 2025.

[42] C. Lee, R. Roy, M. Xu, J. Raiman, M. Shoeybi, B. Catanzaro, and W. Ping, "Nv-embed: Improved techniques for training llms as generalist embedding models," in *Proceedings of the International Conference on Learning Representations*, 2025.

## TEMPORAL EXTRACTION PROMPT

```
1  -Goal-
2  Given a user query that is potentially ask a time-
       related question, identify the timestamp
       entities in the query, follows the structure:
3  - entity_name: standard format of the timestamp
       entity identified in context, following {
       timestamp_format}
4  - entity_type: {timestamp_types}
5  - temporal_logic: one of <at, before, after, between
       >
6  If temporal_logic is <between>, format a pair of
       timestamp entities as ("entity"{tuple_delimiter
       }<between>{tuple_delimiter}<entity_name>{
       tuple_delimiter}<entity_name>{tuple_delimiter}<
       entity_type>)
7  If temporal_logic is <at, before, after>, format a
       pair of timestamp entities as ("entity"{
       tuple_delimiter}<temporal_logic>{tuple_delimiter
       }<entity_name>{tuple_delimiter}<entity_type>)
8
9  When finished, output {completion_delimiter}
10
11 #####################
12 -Examples-
13 #####################
14 Example 1:
15 User Query:
16 Who was the CEO of DXC Technology on January 1,
       2022?
17 ################
18 Output:
19 ("entity"{tuple_delimiter}"at"{tuple_delimiter}"
       2022-01-01"{tuple_delimiter}"date"){
       completion_delimiter}
20 ############################
21 Example 2:
22 User Query:
23 What strategic decisions were made between Q2 and Q4
        2022?
24 #############
25 Output:
26 ("entity"{tuple_delimiter}"between"{tuple_delimiter}
       "2022-Q2"{tuple_delimiter}"2022-Q4"{
       tuple_delimiter}"quarter"){completion_delimiter}
27
28 ###########################
29 Example 3:
30 User Query:
31 What has changed in Aon's leadership after the NFP
       acquisition in 2023?
32 #############
33 Output:
34 ("entity"{tuple_delimiter}"after"{tuple_delimiter
       }"2022"{tuple_delimiter}"year"){
       completion_delimiter}
35 ###########################
36
37 Output:
38 -Real Data-
39 ####################
40 User Query: {input_text}
41 ####################
42 Output:
```

## LOCAL QUERY PROMPT

```
1  ---Role---
2  You are a helpful assistant responding to questions
       about temporal data in the relevant chunks
       provided.
3
4  ---Goal---
```

```
5  Answer the user's question based on the available
       information in the given chunks. Your response
       should:
6
7  1. **Answer the question directly** - Provide the
       specific information requested
8  2. **Be temporally accurate** - Ensure temporal
       information matches the question's scope
9  3. **Be evidence-based** - Only provide answers when
        you have clear, unambiguous evidence
10
11 **Critical Guidelines:**
12 - **ONLY answer if you have clear, unambiguous
       evidence** from the provided chunks
13 - **Refuse to answer** if the evidence is unclear,
       conflicting, incomplete, or uncertain
14 - **Refuse to answer** if you cannot find specific
       information requested in the question
15 - **Refuse to answer** if the chunks contain related
        information but not the exact answer needed
16 - For temporal queries, be flexible with temporal
       expressions (e.g., "2023 Q4" vs "fourth quarter
       of 2023")
17 - Use the given chunks as your primary source of
       information
18 - The chunks are ordered by relevance, so focus on
       the most relevant chunks
19
20 **When to refuse (respond with "No explicit evidence
        for the question"):**
21 - No relevant information found in chunks
22 - Information is present but unclear or ambiguous
23 - Information is incomplete for the specific
       question asked
24 - You are uncertain about the answer
25 - Cannot find the specific information requested
26
27 ---Chunks--- (Ordered by relevance)
28
29 {"".join(processed_chunks)}
30
31 ---Query---
32 {query}
33
34 ---Target response length and format---
35
36 {response_format}
```

## GLOBAL QUERY PROMPT

```
1  ---Role---
2  You are a helpful assistant responding to questions
       about a dataset by synthesizing perspectives
       from multiple analysts.
3
4
5  ---Goal---
6  Generate a response of the target length and format
       that responds to the user's question, summarize
       all the reports from multiple analysts who
       focused on different parts of the dataset.
7
8  Note that the analysts' reports provided below are
       ranked in the **descending order of importance
       **.
9
10 If you don't know the answer or if the provided
       reports do not contain sufficient information to
        provide an answer, just say so. Do not make
       anything up.
11
12 The final response should remove all irrelevant
       information from the analysts' reports and merge
        the cleaned information into a comprehensive
       answer that provides explanations of all the key
```

```
        points and implications appropriate for the
            response length and format.
13
14 Add sections and commentary to the response as
            appropriate for the length and format. Style the
            response in markdown.
15
16 The response shall preserve the original meaning and
            use of modal verbs such as "shall", "may" or "
            will".
17
18 Do not include information where the supporting
            evidence for it is not provided.
19
20 Create a comprehensive analysis that demonstrates
            diversity through:
21
22 **STRUCTURAL DIVERSITY:**
23 - Use varied section headers (mix of descriptive,
            analytical, and thematic headers)
24 - Alternate between different organizational
            patterns within the same response
25 - Combine narrative flow with analytical rigor
26 - Mix chronological and thematic organization
27
28 **PRESENTATION DIVERSITY:**
29 - Vary paragraph styles (some analytical, some
            narrative, some data-focused)
30 - Use different evidence presentation methods (
            direct quotes, summaries, bullet points, tables)
31 - Mix formal analysis with accessible explanations
32 - Alternate between macro and micro perspectives
33
34 **CONTENT DIVERSITY:**
35 - Present multiple analytical frameworks and
            perspectives
36 - Include both quantitative and qualitative insights
37 - Balance strategic, operational, and financial
            viewpoints
38 - Show both short-term and long-term implications
39
40 **STYLISTIC DIVERSITY:**
41 - Vary sentence structures and lengths
42 - Mix technical precision with engaging narrative
43 - Use different transition styles between sections
44 - Combine data-driven analysis with strategic
            insights
45
46 ---Target response length and format---
47 {response_type}
48
49
50 ---Analyst Reports---
51 {report_data}
52
53 ---Target response length and format---
54 {response_type}
55
56 Add sections and commentary to the response as
            appropriate for the length and format. Style the
            response in markdown.
```

## LOCAL QUERY PREDICTION EVALUATION PROMPT

```
1 You are a fact-level evaluation assistant.
2
3 Given:
4 - A user's question
5 - The ground-truth answer (unanswerable)
6 - The model's prediction (which may contain multiple
            factual claims)
7
8 Task:
9 1. Identify all **distinct factual elements** in the
            prediction answer (e.g., numbers, facts, etc.).
```

```
10 2. For each element, compare with the model's
            prediction and classify as one of:
11    - Correctly Refusal: The model explicitly refused
            to answer this element (e.g., said "I don't
            know" or "No explicit evidence").
12    - Incorrect: The model provided a wrong or
            hallucinated value.
13
14 Output:
15 - Count how many elements are correctly refusal or
            incorrect.
16 - Use exact matching unless it's obvious the meaning
            is equivalent.
17 - Only consider content explicitly present in the
            prediction. Do not infer missing values.
18        """
19
20    USR_TEMPLATE = """
21 ### Input
22 **Question:**
23 {question}
24
25 **Ground-Truth Answer:**
26 {answer}
27
28 **Model Prediction:**
29 {prediction}
30
31 ### Assistant
32 Return the result in **strict JSON format** as:
33 {{
34    "correctly refusal": <int>,
35    "incorrect": <int>
36 }}
```

## GLOBAL QUERY PREDICTION EVALUATION PROMPT

```
1 ---Role---
2 You are an expert tasked with evaluating two answers
            to the same question based on three criteria:
            **Comprehensiveness**, **Diversity**, and **
            Temporal Coverage**.
3
4 ---Goal---
5 You will evaluate two answers to the same question
            based on three criteria: **Comprehensiveness**,
            **Diversity**, and **Temporal Coverage**.
6
7 - **Comprehensiveness**: How much detail does the
            answer provide to cover all aspects and details
            of the question?
8
9 - **Diversity**: Evaluate the richness and variety
            of the answer across multiple dimensions:
10    - **Content Diversity**: Does the answer explore
            different aspects, angles, or facets of the
            question rather than repeating similar points?
11    - **Analytical Perspectives**: Does it present
            multiple viewpoints, analytical frameworks, or
            methodological approaches?
12    - **Evidence Variety**: Does it draw from diverse
            sources, data types, or evidence categories?
13    - **Structural Variety**: Does the presentation
            style, organization, or formatting vary
            appropriately to the content?
14    - **Depth vs. Breadth Balance**: Does it strike a
            good balance between detailed analysis and
            broad coverage?
15    - **Innovation**: Does it offer unique insights
            or creative approaches to addressing the
            question?
16    Note: Length alone does not indicate diversity.
            Focus on the richness of perspectives and
            approaches.
17
```

```
18  - **Temporal Coverage**: How well does the answer
       capture the **time dimension** of the question?
19    - Does it clearly reference the relevant periods
       (years, quarters, events) mentioned or implied
       in the query?
20    - Are chronological relationships accurate and
       complete?
21    - Is the timeline explanation easy to follow and
       logically organized?
22
23  For each criterion, choose the better answer (either
       Answer 1 or Answer 2) and explain why. Then,
       select an overall winner based on these three
       categories.
24
25  Here is the question:
26  {question}
27
28  Here are the two answers:
29
30  **Answer 1:**
31  {prediction1}
32
33  **Answer 2:**
34  {prediction2}
35
36  Evaluate both answers using the three criteria
       listed above and provide detailed explanations
       for each criterion.
37
38  Output your evaluation in the following JSON format:
39
40  {{
41      "Comprehensiveness": {{
42          "Winner": "[Answer 1 or Answer 2]",
43          "Explanation": "[Provide explanation here]"
44      }},
45      "Diversity": {{
46          "Winner": "[Answer 1 or Answer 2]",
47          "Explanation": "[Provide explanation here]"
48      }},
49      "Temporal Coverage": {{
50          "Winner": "[Answer 1 or Answer 2]",
51          "Explanation": "[Provide explanation here]"
52      }},
53      "Overall Winner": {{
54          "Winner": "[Answer 1 or Answer 2]",
55          "Explanation": "[Summarize why this answer
       is the overall winner based on the three
       criteria]"
56      }}
57  }}
```

## CASE STUDY OF TEMPORAL GRAPH VISUALIZATION

### Example of Timestamped Relations

Table XI lists the detailed edge descriptions for the visualization example shown in Figure 3. Each record corresponds to a temporal quadruple $(v_1, v_2, r, \tau)$ representing a factual relation grounded at a specific time. This example is constructed from the subgraph of *Western Digital Corporation* and illustrates how our bi-level temporal graph organizes multi-temporal financial facts.

### Case Study Description

Table XII details two representative question answering cases from the ECT-QA benchmark. Each case demonstrates how our temporal graph structure supports accurate retrieval and reasoning over time-dependent financial facts.

*Single-time Specific Query:* Consider the query: *"What were Western Digital Corporation's operating cash flow, gross debt outstanding, and earnings per share in 2020 Q3?"*

Our retrieval process is as follows:

1) Query-Centric Time Identification: The query is analyzed to identify the relevant time point *2020-Q3*.
2) Dynamic Subgraph Positioning: TG-RAG retrieves top-30 relevant relations in the temporal knowledge graph and identifies temporally focused seed set anchored to the *2020-Q3* time node.
3) Local Retrieval: Top 5 nodes with highest PPR score are "WESTERN DIGITAL CORPORATION", "2020-Q3","GROSS MARGIN", "REVENUE", "FREE CASH FLOW", after chunk scoring the chunk with highest score comes from the earnings call transcript from company Western Digital Corporation in 2020-Q3.

The retrieved evidences make the QA LLM to generate the accurate answer: *"Operating cash flow: $142 million, gross debt outstanding: $9.8 billion, earnings per share: $0.85."*

*Multi-Time Query:* Now consider a query spanning multiple time points: *"What was Western Digital Corporation's revenue in each quarter from 2023 Q1 to Q3?"*

Our method handles this as follows:

1) Query-Centric Time Identification: The query is recognized as requiring data from three distinct time points: *2023-Q1*, *2023-Q2*, and *2023-Q3*.
2) Dynamic Subgraph Positioning: TG-RAG retrieves top-20 relevant relations in the temporal knowledge graph and identifies temporally focused seed set anchored to the time nodes.
3) Local Retrieval: Top 5 nodes with highest PPR score are "REVENUE", "WESTERN DIGITAL CORPORATION","2023-Q1","2023-Q3", after chunk scoring the top 3 chunks with highest scores come from the earnings call transcripts from company Western Digital Corporation in 2023-Q1, 2023-Q2 and 2023-Q3 respectively.

The retrieved evidences enable the LLM to generate the accurate answer: *"$142 million, $9.8 billion, and $0.85."*

This case study confirms that our graph structure is not merely a storage format but an effective retrieval backbone. By explicitly modeling time as a first-class citizen, it prevents temporal confusion and provides the necessary granularity to answer complex, time-sensitive questions precisely.

TABLE XI
TIMESTAMPED RELATIONS EXTRACTED FOR THE WESTERN DIGITAL CORPORATION CASE STUDY.

| Entity 1 | Entity 2 | Time | Relation Description |
|---|---|---|---|
| Western Digital Corporation | Operating Cash Flow | 2020 Q2 | In 2020 Q2, Western Digital Corporation's operating cash flow was $257 million. |
| Western Digital Corporation | Operating Cash Flow | 2020 Q3 | In Q3 2020, Western Digital Corporation's operating cash flow was $142 million. |
| Western Digital Corporation | Gross Debt Outstanding | 2020 Q3 | At the end of Q3 2020, Western Digital Corporation's gross debt outstanding was $9.8 billion. |
| Western Digital Corporation | Earnings Per Share | 2020 Q3 | In Q3 2020, Western Digital Corporation's earnings per share was $0.85. |
| Western Digital Corporation | Operating Expenses | 2020 Q3 | In Q3 2020, Western Digital Corporation's operating expenses were $738 million, slightly lower than expected. |
| Western Digital Corporation | Revenue | 2022 | For the full fiscal year 2022, Western Digital Corporation's revenue was $18.8 billion, an 11% increase from fiscal 2021. |
| Western Digital Corporation | Revenue | 2023 | For fiscal year 2024, Western Digital's total revenue was $13 billion, up 6% from fiscal year 2023. |
| Western Digital Corporation | Revenue | 2023 Q1 | In Q1 2023, Western Digital Corporation reported total revenue of $3.7 billion, down 17% sequentially and 26% year over year. |
| Western Digital Corporation | Revenue | 2023 Q2 | In fiscal Q2 2023, Western Digital Corporation delivered revenue of $3.1 billion, which was at the high end of its guidance range. |
| Western Digital Corporation | Revenue | 2023 Q3 | In Q3 2023, Western Digital Corporation reported total revenue of $2.8 billion, which was down 10% sequentially and 36% year over year. |

TABLE XII
CASE STUDY EXAMPLES AND EVALUATION METRICS ON ECT-QA.

| Question | Golden Answer | Prediction | Correct | Refusal | Incorrect | ROUGE-L | $F_1$ |
|---|---|---|---|---|---|---|---|
| What were Western Digital Corporation's operating cash flow, gross debt outstanding, and earnings per share in 2020 Q3? | $142 million, $9.8 billion, and $0.85. | Operating cash flow: $142 million, gross debt outstanding: $9.8 billion, earnings per share: $0.85. | 1.0 | 0.0 | 0.0 | 0.583 | 0.500 |
| What was Western Digital Corporation's revenue in each quarter from 2023 Q1 to Q3? | $3.7 billion, $3.1 billion, and $2.8 billion. | 2023 Q1: $3.7 billion, 2023 Q2: $3.1 billion, 2023 Q3: $2.8 billion. | 1.0 | 0.0 | 0.0 | 0.720 | 0.632 |