# XGraphRAG: Interactive Visual Analysis for Graph-based Retrieval-Augmented Generation

Ke Wang, Bo Pan, Yingchaojie Feng, Yuwei Wu, Jieyi Chen, Minfeng Zhu, and Wei Chen
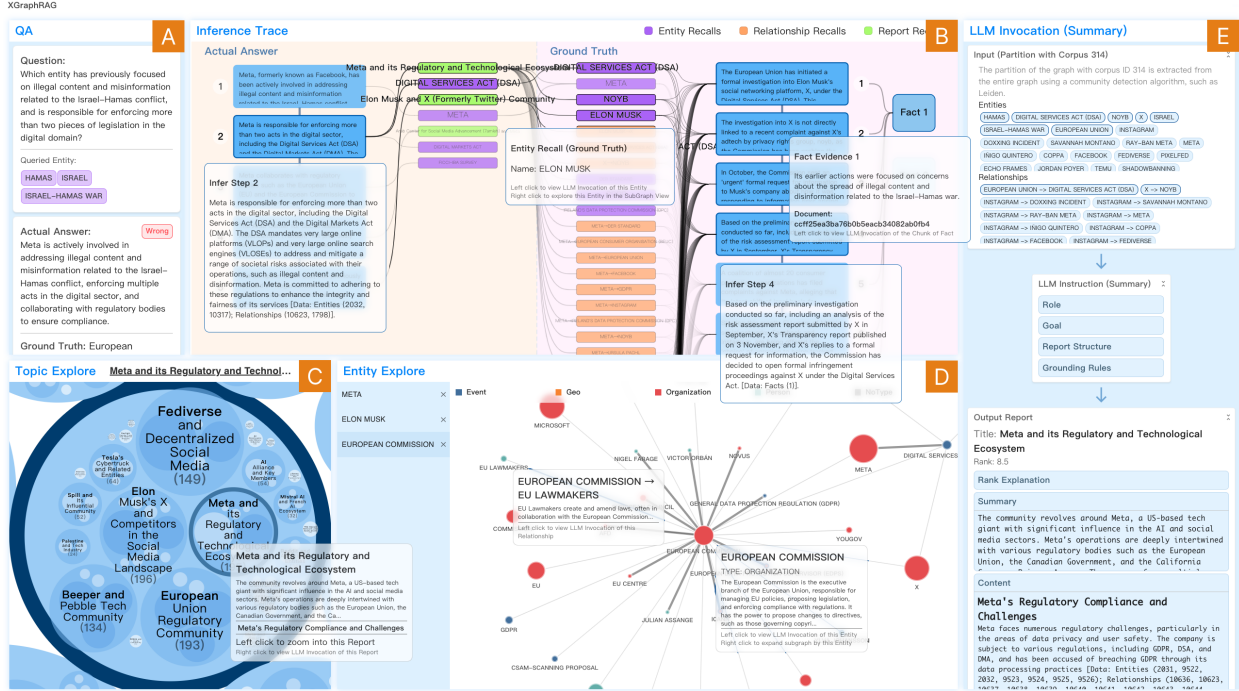
Fig. 1: *XGraphRAG* contains four views: (A&B) The QA & Inference Trace View allows users to locate suspicious recalls. (C) The Topic Explore View supports global relevance analysis on the graph for suspicious retrievals. (D) The Entity Explore View supports local relevance analysis on the graph for suspicious retrievals. (E) The LLM Invocation View presents details for each LLM invocation.

**Abstract**—Graph-based Retrieval-Augmented Generation (RAG) has shown great capability in enhancing Large Language Model (LLM)'s answer with an external knowledge base. Compared to traditional RAG, it introduces a graph as an intermediate representation to capture better structured relational knowledge in the corpus, elevating the precision and comprehensiveness of generation results. However, developers usually face challenges in analyzing the effectiveness of GraphRAG on their dataset due to GraphRAG's complex information processing pipeline and the overwhelming amount of LLM invocations involved during graph construction and query, which limits GraphRAG interpretability and accessibility. This research proposes a visual analysis framework that helps RAG developers identify critical recalls of GraphRAG and trace these recalls through the GraphRAG pipeline. Based on this framework, we develop XGraphRAG, a prototype system incorporating a set of interactive visualizations to facilitate users' analysis process, boosting failure cases collection and improvement opportunities identification. Our evaluation demonstrates the effectiveness and usability of our approach. Our work is open-sourced and available at `https://github.com/Gk0Wk/XGraphRAG`.

**Index Terms**—Retrieved-augmented generation, large language model, visual analysis, interactive visualization

---◆---

## 1 INTRODUCTION

Large Language Models (LLMs) have rapidly gained prominence in recent years, demonstrating significant value and potential across various

- Ke Wang, Bo Pan, Yingchaojie Feng, Yuwei Wu, Jieyi Chen, and Wei Chen are with the State Key Lab of CAD&CG, Zhejiang University. E-mail: {sttotphd | bopan | fycj | 22451008 | chenjieyi_juraws | chenvis}@zju.edu.cn.
- Minfeng Zhu is with Zhejiang University. E-mail: minfeng_zhu@zju.edu.cn.
- Minfeng Zhu is the corresponding author.

sectors such as industry [35, 36], healthcare [1, 26, 34], science [2, 38], and Finance [40, 42]. Despite their unprecedented language comprehension and text generation capabilities, LLMs can face limitations when dealing with domain-specific knowledge, real-time updated information, and proprietary data, which are not in their pre-training corpus [11]. A promising paradigm for addressing this limitation is Retrieval-Augmented Generation (RAG), which integrates a retrieval component within the generation process to leverage the information contained in a customized large text corpus. This process not only enriches the contextual depth of the responses but also boosts their factual accuracy and specificity [7, 23]. Although RAG has achieved impressive results, it faces limitations in real-world scenarios that involve complex structures of relationships among different entities in the customized text corpus [23, 24]. To address this, recent works [6, 12, 39] propose introducing graphs as an intermediate representation during the retrieval process (i.e. Graph-based RAG) to capture better and leverage

the structured relational knowledge contained in the customized text corpus, bringing about new opportunities to elevate the precision and comprehensiveness of generation results in RAG applications.

However, despite the promising prospects of the Graph-based RAG (GraphRAG) paradigm, current RAG application developers often face challenges in analyzing the effectiveness of GraphRAG on their datasets. This reduces their confidence and motivation to deploy GraphRAG solutions in actual production. This difficulty arises from the two additional layers of complexity introduced by GraphRAG compared to traditional RAG: First, the information processing pipeline of a GraphRAG system is much longer and more complex than that of traditional RAG. Starting from raw text, the information goes through a series of intricate processes, including entity and relation extraction, graph construction, information aggregation and summarization, and parallel multi-stage retrievals. Users lack effective tools to track and analyze these processes. Second, the information processing in a GraphRAG system involves frequent and extensive calls to LLMs. It is difficult to understand the context and role of each LLM invocation within the graph during the analysis process.

As such, it is highly desirable to provide analysis support for GraphRAG developers to examine the GraphRAG process on their customized text corpus, enabling them to make informed design decisions and adjustments. To clarify the analysis targets, we survey recent GraphRAG-related papers and open-source projects to abstract a common pipeline for GraphRAG. For the sake of better understanding the underlying challenges in the analysis process, we conducted a formative study with 5 experienced RAG experts to have a grasp of the current analysis process, common analysis strategies, and the challenges they encounter when analyzing the GraphRAG process.

Based on the findings of the formative study, we propose a two-stage visual analysis framework for the process diagnosis of GraphRAG systems. In the first stage, it leverages the power of the LLM to provide automatic evaluation of model answers and identify the suspicious incorrect recalls that mislead the inference process. In the second stage, it allows users to trace the causes of the incorrect recalls through multi-facet relevance analysis on the graph. Based on this framework, we design *XGraphRAG*, a visual analysis system consisting of four interrelated views. The QA View (Fig 1A) and the Inference Trace View (Fig 1B) serve as the springboard for exploration that helps users compare the discrepancies between the answer and ground truth to identify suspicious recalls. The Topic Explore View (Fig 1C) and the Entity Explore view (Fig 1D) allow users to conduct relevance analysis on the graph from different aspects. The LLM Invocation View (Fig 1E) adaptively displays structured details of the LLM invocation related to the intermediate output the user is interested in during the analysis process. We conduct a usage scenario demonstration and a user study to verify the effectiveness of our approach.

In summary, this paper makes the following contributions:

- We identify challenges in the analysis of the GraphRAG process and distill design requirements.

- We propose a visual analysis framework that allows GraphRAG developers to examine the GraphRAG process efficiently and systematically.

- We present *XGraphRAG*, a prototype system that instantiates our framework, whose effectiveness is verified with a usage scenario demonstration and a user study.

## 2 RELATED WORK

Here, we review works in the field of GraphRAG and visual Analysis for LLM-based Applications, which are closely related to our study.

### 2.1 GraphRAG

GraphRAG is an emerging powerful retrieval-augmented generation paradigm that leverages knowledge graph as an intermediate representation to enable more precise and comprehensive retrieval for LLM generation, which has recently shown great potential in diverse application areas, such as medicine [31], E-commerce [41], intelligence analysis [25], and software engineering [3]. While all GraphRAG

approaches utilize graphs as a crucial intermediate knowledge representation, the underlying design space is vast and has attracted extensive research attention. For the graph construction phase, various indexing strategies (e.g. basic graph indexing [13], text indexing [16], vector indexing [14], and hybrid indexing [27]) have been explored to facilitate downstream retrievals. For the graph retrieval phase, different types of recalls (e.g. nodes [17], relationships [16], subgraphs [6], and paths [21]) can be retrieved. Moreover, the retrieval paradigm can vary from efficient and cost-effective one-time retrieval [13, 14] to more sophisticated and adaptive iterative retrieval [21, 30].

While extensive work has focused on optimizing GraphRAG's process design, none has addressed how to assist developers in analyzing the underlying complex processes—a crucial aspect for practical GraphRAG applications. In this work, we propose a visual analytics framework that enables users to better understand and analyze GraphRAG processes, facilitating its effective deployment in real-world scenarios.

### 2.2 Visual Analysis for LLM-based Applications

Based on the interaction patterns of LLM invocations, existing visualization systems for analyzing LLM behavior can be categorized into two main directions: single-shot LLM invocation analysis and iterative LLM invocation analysis.

For single-shot LLM invocation analysis, researchers have developed various visualization systems to analyze individual LLM responses. PromptIDE [28] enables users to experiment with prompt variations and visualize prompt performance for ad-hoc task adaptation. LLM Comparator [15] facilitates side-by-side evaluation of LLM responses through interactive visual analytics, helping users understand performance differences between models. To analyze specific LLM capabilities, some work focused on specialized analysis tasks. For example, researchers proposed visualization techniques to analyze text style transfer [5], revealing how LLMs apply different stylistic features in their responses. In terms of security analysis, JailbreakLens [8] provides a visual analytics system for analyzing jailbreak attacks against LLMs, enabling users to evaluate model vulnerabilities and defensive capabilities.

For iterative LLM invocation analysis, existing work has explored how to visualize and analyze complex interaction patterns between multiple LLM calls. Sensecape [29] supports multilevel exploration and sensemaking with LLMs by enabling users to manage information complexity through hierarchical organization. ChainForge [4] provides a visual toolkit for prompt engineering and hypothesis testing across multiple LLM responses. Recent work has also focused on visually analyzing LLM-based agent systems. AgentLens [20] visualizes the dynamic evolution of agent behaviors in LLM-based autonomous systems, while AgentCoord [22] helps users explore and design coordination strategies for LLM-based multi-agent collaboration. [19] breaks down the programming tutorial creation process into actionable steps and visualizes the tree-of-thought exploration process.

Extending this research line, we aim to develop an interactive visual analysis framework for analyzing complex GraphRAG processes with diverse and intensive iterative LLM invocations.

## 3 FORMATIVE STUDY

We conduct a formative study to inform the design of our approach. First, to understand the core process of GraphRAG, we surveyed existing papers and high-impact open-source projects involving GraphRAG. Based on the survey, an abstraction of the common pipeline of GraphRAG is summarized in Section 3.1. Second, we conducted a formative interview with 5 RAG experts to identify the challenges encountered during the analysis of the GraphRAG process with existing tools in Section 3.1. Finally, we distill the four design requirements to improve the analysis process.

### 3.1 Common Pipeline of GraphRAG

We survey recent GraphRAG-related papers and open-source projects to abstract the typical components of a GraphRAG pipeline. Given that GraphRAG is an emerging field with many works still awaiting peer
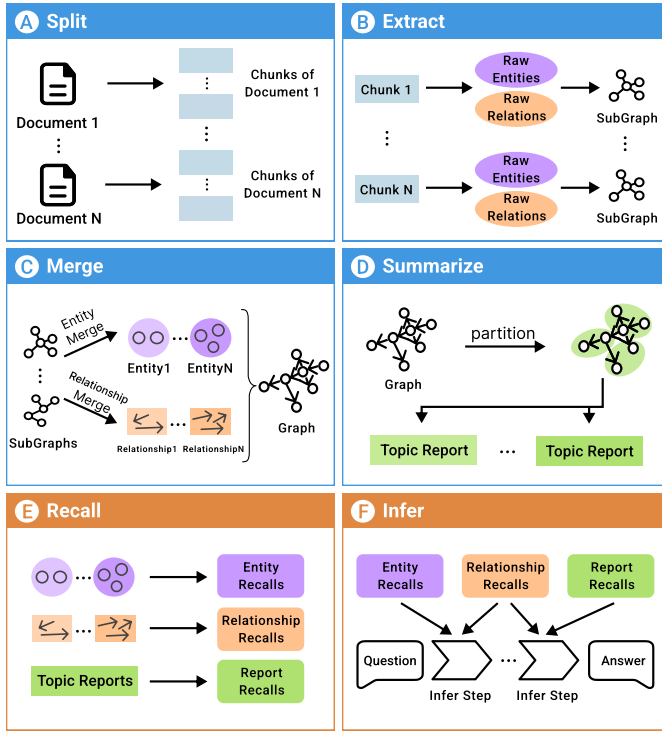
**A** Split

Document 1 → ⋮ → Chunks of Document 1

Document N → ⋮ → Chunks of Document N

**B** Extract

Chunk 1 → Raw Entities / Raw Relations → SubGraph

⋮

Chunk N → Raw Entities / Raw Relations → SubGraph

**C** Merge

SubGraphs — Entity Merge → Entity1 ⋯ EntityN

— Relationship Merge → Relationship1 ⋯ RelationshipN

→ Graph

**D** Summarize

Graph → partition →

Topic Report ⋯ Topic Report

**E** Recall

→ Entity Recalls

→ Relationship Recalls

Topic Reports → Report Recalls

**F** Infer

Entity Recalls | Relationship Recalls | Report Recalls

Question → Infer Step ⋯ Infer Step → Answer

Fig. 2: Common Pipeline of GraphRAG, which is typically divided into an offline construction phase (a-d) and an online retrieval phase (e-f).

review, we identified influential papers/projects in this area based on recommendations from RAG experts and search terms "GraphRAG", "Graph Retrieval Augmented Generation", "Knowledge Graph RAG", "Graph-based RAG", and "Graph-augmented LLM". Starting with several influential seeding works (e.g., GraphRAG [6], KAG [18], LightRAG [12]), we iteratively expanded our corpus by analyzing both the references within these works and the works that cite them. Our final collection comprises 30 GraphRAG-related works (list included in our supplementary materials). Drawing from this corpus, we analyzed and summarized the essential data elements and data processing components shared by them. As shown in Fig 2, the GraphRAG pipeline is typically divided into two main components: an offline construction phase (a-d) and an online retrieval phase (e-f). We elaborate on the essential data process stages for each phase as follows.

### 3.1.1 Construction Phase

**Split:** The initial step involves dividing the original documents into smaller, manageable **text chunks**. This segmentation is crucial as it sets the foundation for subsequent processing. The text chunks are designed to be coherent units of information, facilitating easier extraction and analysis in later stages.

**Extract:** For each text chunk, we extract **raw entities** and **raw relationships**, forming **subgraphs**. This involves identifying key elements within the text, such as names, dates, and specific terms, and understanding how these elements relate to each other. Each entity and relationship is enriched with additional metadata, such as **descriptions** and **entity types**, which provide context and clarity, enhancing the granularity of the subgraphs.

**Merge:** The next step is to combine these subgraphs into a comprehensive **graph**. This involves resolving conflicts where entities and relationships share names but differ in context. The merging process ensures that information is unified, with **entity** and **relationship** descriptions being consolidated. This step is critical for maintaining data integrity and ensuring that the graph accurately reflects the source material. Each **node** and **edge** in the graph is annotated with **chunk references** to the original text chunks from which they were derived, ensuring traceability.

**Summarize:** Once the graph is complete, it is partitioned based

on the density of relationships among entities, resulting in distinct **topics**. This partitioning helps in organizing the graph into meaningful segments, each representing a coherent topic or theme. Within each topic, the entities and relationships are summarized to create concise **reports**. These summaries are constructed with clear references to the specific entities and relationships they are derived from, ensuring transparency. Depending on the complexity of the original corpus and the graph, this partitioning may occur at multiple hierarchical levels, providing varying degrees of abstraction and detail.

Throughout the construction phase, large language model invocations are extensively employed at different stages for natural language processing tasks. These models assist in tasks such as entity recognition, relationship extraction, and summarization, leveraging their advanced capabilities to enhance the accuracy and efficiency of the process.

### 3.1.2 Retrieval Phase

**Recall:** During the retrieval phase, user **queries** are utilized to extract relevant information from the constructed graph. The system identifies and retrieves various types of **recalls**, such as **entity recalls**, **relationship recalls**, and **report recalls**. This involves searching the graph for nodes and edges that are relevant to the user's query, ensuring that the most pertinent information is brought to the forefront.

**Infer:** The LLM then uses these recalls to perform step-by-step **inference**. Each inference step is explicitly linked to specific recalls, ensuring that the reasoning process is transparent and traceable. This iterative process allows the system to build upon retrieved information, synthesizing it to generate comprehensive answers. The inference culminates in producing a final **response** to the user's query, leveraging the depth and breadth of the graph's information.

This structured approach enhances the GraphRAG system's ability to manage and utilize complex relationships within data, supporting more effective retrieval and reasoning capabilities. By providing clear traceability and leveraging advanced language models, the system ensures both accuracy and transparency in its responses.

### 3.2 Challenges of Analyzing the GraphRAG Process

**Participants and procedure.** To understand the challenges users could encounter during the analysis of the GraphRAG process, we conducted a formative interview with 5 RAG experts. Three of them (E1-E3) are experienced RAG engineers who have first-hand experience with GraphRAG in their projects. The other two (E4-E5) are RAG researchers who have conducted at least one GraphRAG-related research project. In our interview, we first asked the participants about their current workflow for analyzing the GraphRAG process on their datasets and the challenges they encountered during this process. After that, we showed them how to use *Kotaemon* [33], one of the most popular open-source projects that provide interactive analysis support for the GraphRAG process. We then asked the participants to try out *Kotaemon* on their own in a think-aloud way (reporting the analysis target and the challenges they encountered in situ). We also collected participants' overall feedback at the end of the interview. Based on the interview, we summarize the challenges as follows.

**Lack of traceability across the complex information processing pipeline.** All of the participants mentioned the necessity and difficulty of tracing through the information processing pipeline of GraphRAG. To analyze the cause of the final or intermediate output, the participants need to "examine its upstream modules step by step" (E2). Nevertheless, this demand is not well-supported by existing tools. *Kotaemon* [33] provides some support for certain steps (e.g. users can click an entity to trace the chunk it is extracted from), but "lacks complete traceability from the final answer back to the raw document chunks" (E1).

**Lack of revealing of LLM invocation context on the graph.** The GraphRAG process typically involves hundreds or even thousands of LLM invocations to construct and query the graph. The participants mentioned that they often feel "overwhelmed and lost in the sea of LLM calls" (E2) and "hard to connect certain LLM invocation with its specific role on graph" (E5). Existing tools only provide exclusive or separate support for graph analysis for LLM invocation analysis.

However, it is highly desirable to provide support for "analyzing LLM invocations on the graph" (E4).

**Lack of support for multi-facet relevance analysis on the graph.** The key to the success of GraphRAG lies in whether it can utilize the graph structure to capture the relevance between the information in the customized corpus and the user's query. This relevance can be explicit local connectivity on the graph, or semantic connectivity based on global structure and specific techniques applied for graph summarization (as discussed in Section 3.1.1). However, there is a lack of support for systematic analysis of different types of relevance on the graph. As E2 mentioned, "the key lies in finding why some relevance is successfully or unsuccessfully captured and exploited on graphs, which is a very complex and demanding analysis process"

## 3.3 Design Requirement

In response to the challenges identified in our formative study, we aim to design a visual analysis system to facilitate the analysis process for GraphRAG. The design requirements are distilled as follows:

**R1: Facilitate step-by-step trace for the generated answer.** To understand the final or intermediate output, users need to trace back to its upstream modules step by step. This traceability ensures that users can verify the accuracy of each step and identify the root cause of unexpected results. The system should provide flexible interaction and visualization to facilitate this tracing process.

**R2: Reveal LLM invocation context on the graph.** There is a large quantity of LLM invocations during the GraphRAG process. While those LLM invocations play a vital role in information extraction and processing during graph construction and query, it is hard for analyzers to link each invocation with its context on the graph during the analysis process. Therefore, it is highly desirable to provide support to help users understand the structure and context of those invocations.

**R3: Support global relevance analysis on the graph.** A significant advantage of GraphRAG over conventional RAG is its ability to use graphs to capture global semantic information. For example, if certain entities are not directly connected but are frequently mentioned in the context of certain topics, they might be grouped into the same community by GraphRAG, which is essential for answering questions that require a broader, more holistic understanding. Therefore, analysis support for the underlying process is important for understanding why GraphRAG succeeds or fails in different settings.

**R4: Support local relevance analysis on the graph.** While GraphRAG excels at capturing global semantic information, it also greatly enhances recall quality by leveraging local graph connectivity to capture local relevance. For instance, if a user's query pertains to a particular entity, understanding the direct connections and relationships of that entity can provide more precise and contextually relevant answers. Therefore, it is also important to provide support for a detailed, localized analysis of the graph.

## 4 WORKFLOW

Here, we elaborate on the workflow for analyzing the *GraphRAG* process.

### 4.1 Overview

To assist users in exploring graph construction issues, we divide the entire methodological process into two stages:

**Identifying Suspicious Retrievals:** Users conduct a question-and-answer evaluation using test data to preliminarily identify contradictions between actual answers and the ground truth. By analyzing the inference paths of both the actual answers and the ground truth, users can pinpoint problematic inference steps and identify suspicious recalls, understanding what the issues are and where they occur.

**Analyzing Suspicious Retrievals:** Starting with the identified suspicious recalls, users conduct either local or global analyses based on the type of recall. They trace reference chains to investigate the behavior of LLMs, uncover issues in graph construction, and provide evidence for optimizing graph quality, thereby understanding why the issues occur and how to address them.
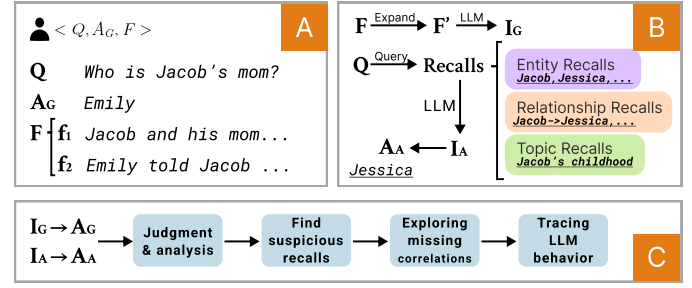


Fig. 3: Overview of the workflow: The user first prepares a testing pair for evaluation (A). Subsequently, two recall-inference-answer chains are constructed for comparison, corresponding to the ground truth answer and actual answer respectively (B). Finally, users conduct a step-by-step analysis of the issue through a linear exploration process (C).

### 4.2 Stage 1: Identifying Suspicious Retrievals

The GraphRAG system utilizes an initial document collection $D$ ($d_i \in D$) to construct a graph $\mathcal{G}$. During the testing phase, users query the constructed graph with test pairs $< Q, A_G, F >$ to retrieve actual answers $A_A$, where $Q$ is the question, $A_G$ is the Ground Truth Answer and $F$ is the set of evidence facts (Fig 3A). We evaluate both the correctness of the generated answers and their corresponding inference steps, guiding users to explore the unexpected results. For **answer evaluation**, we employ an LLM to judge the answer $A_A$ as *correct* or *wrong* based on the Ground Truth $A_G$. For **assessing inference steps**, we develop an LLM-assisted analysis framework [37] to reconstruct the inference process for generated answers and ground truth and compare them to identify suspicious retrieval recalls (Fig 3B). We describe the details of the analysis framework as follows.

#### 4.2.1 Inferences Process Construction

We construct the inference process as a *Question-Recalls-Inferring-Answer* pipeline to facilitate comparative analysis.

**Actual Answer Inferences $I_A$.** When users submit a query, the GraphRAG system first retrieves a series of recalls. These are processed by the LLM to generate detailed inference steps ($I_{A_i} \in I_A$), each step indicating the recalls it relies on.

**Ground Truth Inferences $I_G$.** For the Ground Truth, each fact $F_i \in F$ from the test pair is contextually expanded to $F'_i \in F'$ to facilitate inferring, as the original facts may lack sufficient context. The expanded context test pair $< Q, A_G, F' >$ is then submitted to the LLM for reverse inferring, resulting in detailed inference steps $I_G$ based on $Q$ and $F'$, with each step specifying the facts utilized.

Each fact $F_i$ corresponds to several original document chunks. Entities and relationships related to these chunks are extracted to form a fact subgraph $\mathcal{G}_{F_i}$. For each inference step $I_{G_i}$, the relevant fact subgraphs are merged and filtered by the LLM to select entities and relationships pertinent to that step, forming the inference subgraph $\mathcal{G}_{G_i}$.

Finally, the entities and relationships in the inference subgraph $\mathcal{G}_{G_i}$ are used as recalls for each inference step $I_{G_i}$.

#### 4.2.2 Suspicious Recall Identification

Based on the inference steps of Ground Truth $I_G$ and the Actual Answer $I_A$, users perform a comparative analysis to pinpoint problematic inference steps and identify the corresponding recalls for discrepancies. We flag two types of recalls as suspicious: **Missing Recalls**, which are essential to critical inference steps in the correct answer but absent from the actual inference steps, necessitating analysis to determine why they were not retrieved or used; and **Unexpected Recalls**, which do not belong to the critical inference steps in the correct answer but are erroneously included in the actual inference steps, requiring investigation to understand why they were mistakenly retrieved or utilized.

### 4.3 Stage 2: Analyzing Suspicious Retrievals

In the previous phase, users identified suspicious retrievals by analyzing problematic inference steps. To further investigate the causes of these

issues in the GraphRAG system, users must dive into the internal structure of the system to explore the relationships between entities, relationships, and topic reports. This involves tracing the construction and retrieval processes to uncover the behavior of the LLM, continuing until the root cause of the issue is clarified (Fig 3C).

### 4.3.1 Exploring Missing Correlations

The exploration of missing correlations is categorized into local and global investigations based on recall types, aiming to identify and address missing correlation information. **Local exploration** focuses on entity and relationship recalls, requiring users to examine their relationships with other entities and relationships. Specifically, these relationships may include connections with entities mentioned in the question, entities referenced in the correct answer, or entities involved in critical inference steps.

For **global exploration**, which centers on topic report recalls, users investigate their relationships with other reports. These relationships may include connections to topic reports that reference entities mentioned in the question or entities included in the correct answer. Through this structured approach, users can systematically trace missing correlations and enhance their understanding of the underlying retrieval process.

### 4.3.2 Tracing LLM Behavior

Once a recall with missing correlation information is identified, further tracing of the LLM's behavior is necessary to pinpoint where in the construction process this information was lost.

In the construction and retrieval stages of the GraphRAG system, each stage involves references to preceding content, with the LLM playing a crucial role. This forms the foundation for tracing and analyzing the behavior of the LLMs:

- **Extract**: The LLM extracts multiple entities and relationships from chunks.

- **Merge**: The LLM consolidates same-named entities and relationships, referencing entity-chunks or relationship-chunks.

- **Summarize**: The LLM summarizes entities and relationships within a topic subgraph to create a topic report.

- **Infer**: The LLM uses query retrievals for reasoning, referencing inference step-entities and relationships.

To identify issues in different types of suspicious recalls, users undertake distinct tracing analyses. For entity and relationship recalls, the process involves analyzing the merging behavior and subsequently examining the extraction behavior. In the case of topic report recalls, a comprehensive investigation is conducted, focusing on the summarization, merging, and extraction behaviors.

## 5 SYSTEM

Here, we elucidate the design and usage of each view of *XGraphRAG*.
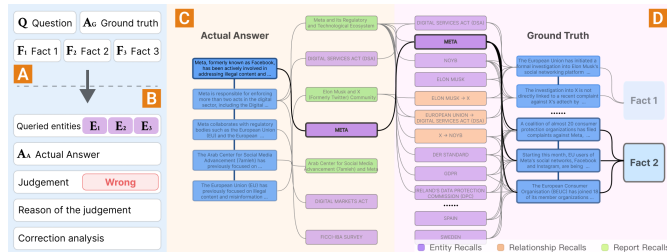
### 5.1 QA View & Inference Trace View



Fig. 4: QA View & Inference Trace View. Once the user inputs the test question in QA View (A), the system returns key information to support discrepancy analysis (B). A detailed comparison of the inference steps for the Actual Answer (C) and the Ground Truth (D) are shown in the Inference Trace View.

The QA (Query-Answer) View (Fig 4A) and the Inference Trace View (Fig 4B) serve as entry points for exploration, helping users identify discrepancies between Actual Answer (Fig 4C) and Ground Truth (Fig 4D). They also help users locate problematic inference steps and identify suspicious recalls (**R1**).

### 5.1.1 Analysis of Question-Answer Discrepancies

Once users input test questions in the form $< Q, A_G, F >$ and perform queries and analyses (Fig 4A), the system initially presents key information to support discrepancy analysis in QA View (Fig 4B). This includes the queried entities explicitly mentioned in the question, the actual answer generated by the GraphRAG system, and the corresponding Ground Truth. Furthermore, the system provides a relevance score that evaluates the alignment between the actual answer and the Ground Truth, offering a quantitative measure of their consistency.

To further support the analysis, justifications for the relevance assessment are provided along with explanations of any observed answer discrepancies. Users can view these justifications and explanations below the Ground Truth, enabling a detailed examination of the underlying causes of discrepancies.

This initial presentation provides users with a clear and structured overview, allowing them to gain a preliminary understanding of the differences between the actual answer and the Ground Truth. It serves as a foundation for a more in-depth analysis of the reasoning and retrieval processes.

### 5.1.2 Analysis of Inference Steps

Beneath the question-answer discrepancy analysis section lies the inference step analysis section, which provides a detailed illustration of the inference steps for both the Actual Answer and the Ground Truth. This section highlights the associated query recalls for each inference step, presenting a functionally symmetric view divided into two areas.

The left section is dedicated to analyzing the inference chain of the Actual Answer (Fig 4C). It displays the sequence of inference steps and their associated query recalls, allowing users to trace the inferring that led to the Actual Answer. The right section focuses on the Ground Truth (Fig 4D), presenting the reference inference steps alongside their associated query recalls and the factual basis for each step.

Users can systematically compare the inference steps of the Actual Answer on the left with the inference steps of the Ground Truth on the right. This side-by-side view enables users to identify discrepancies, pinpoint problematic inference steps, and analyze the differences in the inferring stage between the two.

To further aid analysis, query recalls associated with each inference step can be traced interactively. When users hover over an inference step, the related recalls and their connections are prominently highlighted, providing immediate context for the inferring stage. In the Ground Truth section, the interface also highlights the specific facts used in each inference step, helping users quickly locate the original basis for correct reasoning. This structured approach facilitates an in-depth examination of the inferring stage and its underlying assumptions.

### 5.1.3 Analysis of Recall Usage

Users can analyze the utilization of specific query recalls within the inferring stage by interacting with the system's interface. Hovering over a query recall provides detailed information on its role in the reasoning pipeline, highlighting the following information.

First, all inference steps that involve the selected query recall are displayed, offering a clear view of how the recall contributes to the reasoning process. Second, related recalls on the opposite side, including related inferences and recalls for both the Actual Answer and the Ground Truth, are presented. For topic recalls, this includes entities or relationships encompassed within the topic recall on the opposite side. For entity-relationship recalls, identical entity or relationship recalls on the opposite side are highlighted, along with included topic recalls on the opposite side.

Additionally, related inference steps for the recalls on the opposite side are displayed, where applicable. Finally, the source of facts for the

recall is revealed, whether it pertains to the Ground Truth or to related recalls on the opposite side in the Actual Answer.

This information enables users to comprehensively analyze the query recall's usage in the inferring stage, compare the similarities and differences between the inferring steps on both sides and efficiently identify suspicious recalls.
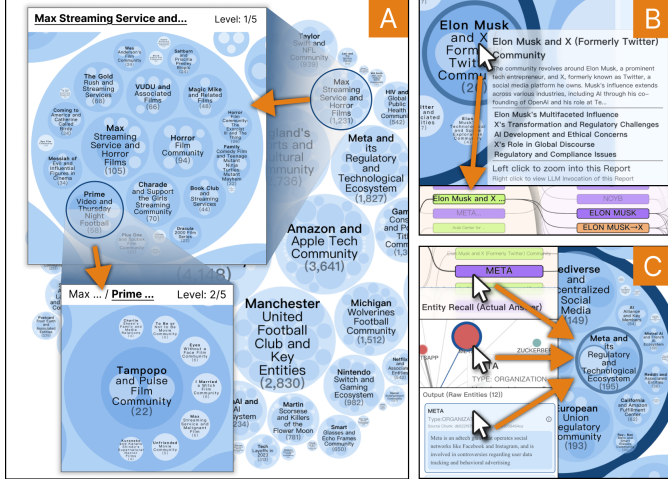
## 5.2 Topic Explore View



Fig. 5: Topic Explore View, which displays all topics with explicit nesting relationships in the graph (A). Whenever the mouse hovers over a topic, the corresponding topic-type recall in the Inference Trace View will be highlighted (B). Conversely, whenever the mouse hovers over any recall in the Inference Trace View, the topic related to it will be highlighted in the Topic Explore View (C).

The Topic Explore View (Fig 1C) utilizes the circle packing hierarchical visualization method to display all topics with explicit nesting relationships in the graph for global relevance exploration (R3). Each topic is represented by a circle at a certain level, showing the topic's title and the number of entities it contains. The number of entities is encoded as the diameter of the circle to intuitively reflect the topic's global size. Although nested circles aren't as space-efficient as treemaps, the extra space helps to better display the hierarchical structure, helping users to understand the relationships between topics more intuitively. Within each topic circle, the next-level topic circles are drawn in a lighter color, providing an intuitive understanding of the composition of subtopics. Users can click on a topic to expand and display the next-level topics it contains or explore other topics at the same level, thereby progressively exploring the composition and structure of topics in the graph (Fig 5A).

The Topic Explore View also provides intuitive interactions to help trace the relationship between a topic and its related recalls: As shown in (Fig 5B), whenever the mouse hovers over a topic, the corresponding topic-type recall in the Inference Trace View will be highlighted. Conversely, as shown in Fig 5C), whenever the mouse hovers over any recall in the Inference Trace View, the topic related to it will be highlighted in the Topic Explore View. This facilitates quick exploration of global semantic relationships among different entities, relationships, and topics.

## 5.3 Entity Explore View

The Entity Explore View (Fig 1D) assists users in investigating the local relevance of suspicious entity recalls with other entity recalls and relationship recalls (R4). Whenever users identify a suspicious entity recall in the Inference Trace View, they can right-click to add it to the list in the Entity Explore View. Subsequently, users can utilize this view to explore other entities and relationships associated with the identified entity. We use a node-link diagram to visualize the local subgraph around this entity. Here the diameter of each node encodes
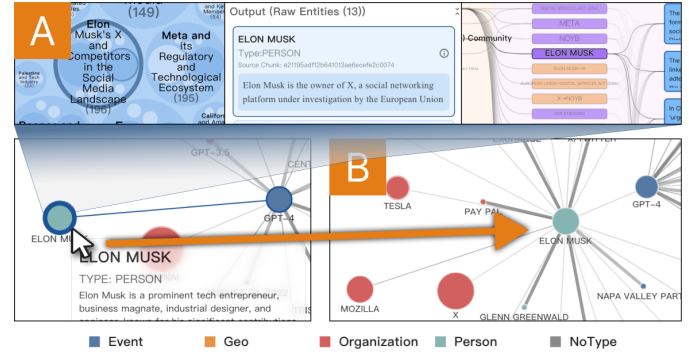


Fig. 6: The Entity Explore View enables investigation of entity recalls through a node-link diagram. Node size encodes entity frequency, edge thickness shows topic distance, and node color indicates entity type. Interactive highlighting helps trace elements across different views (A). Users can expand the local graph by right-clicking nodes/edges to explore related entities and relationships (B).

the number of times it appears in different text chunks, intuitively reflecting the frequency and importance of the entity. The thickness of the edge encodes the distance of nodes in the nested topic tree, which also determines the strength of attraction in the force-directed layout. The node color encodes the types of entities. To ensure that each type has a distinctly different yet soft color, we use the HSL model, evenly sampling the hue values while maintaining moderate saturation and lightness. To help users better trace entities and relationships across views, whenever they hover over an entity point or relationship edge in the local graph, it will also be highlighted in other views (Fig 6A). If the user wants to explore more about a certain node/edge, they can right-click on it to add more entities and relationships near it (Fig 6B).
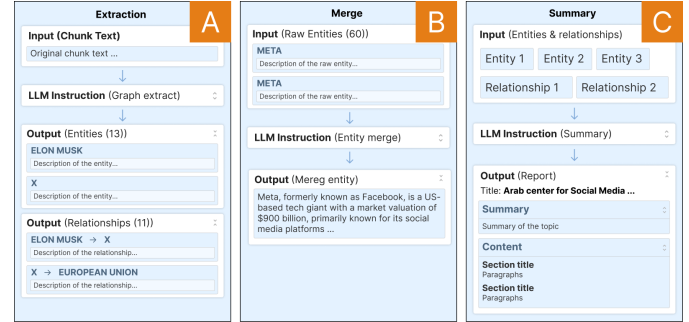
## 5.4 LLM Invocation View



Fig. 7: The LLM Invocation View has three variations, showing LLM behaviors during different stages of graph construction: extraction stage view displays source chunks and extraction details (A), merge stage view presents entity/relationship integration process (B), summary stage view reveals topic report generation with related entities and relationships (C).

The LLM Invocation View (Fig 1E) reveals the behaviors of the LLM during various stages of graph construction in the GraphRAG system (R2). This view has three variations for different stages, as elaborated below.

### 5.4.1 Summary Stage LLM Analysis

Whenever users click on any topic recall in any view, the related information of the topic and the behavioral details of the LLM in generating the topic report during the summary stage are displayed in the LLM Invocation View (see Fig 7C). Users can investigate the entities and relationships used by the LLM in the summary process, the structure of the LLM's prompt words, and the topic report generated by the LLM. Since many entities and relationships may be used in topic summaries, this part can also link with other views for related information highlighting and filtering. For example, when the mouse hovers over an

entity-relationship recall or fact in the Inference Trace View, or an entity point or relationship edge in the Entity Explore View, related entities and relationships are highlighted at the top of the list. This allows for quick confirmation of whether there is a relevant link between the part of the graph used for the topic and the object to be explored.

### 5.4.2 Merge Stage LLM Analysis

Whenever users click on any entity or relationship recall in any view, the related information of the entity/relationship and the behavior details of the LLM in processing the entity or relationship during the merge stage are displayed in the LLM Invocation View (see Fig 7B). Users can investigate the original entities or relationships used by the LLM in the merge process, the structure of the LLM's prompt words, and the merge results generated by the LLM. This information helps users understand how the LLM processes and integrates entity or relationship information from different sources and further traces the source of information for entities or relationships. Similarly, this view links with other views for hover-highlighting, enabling fluid tracing.

### 5.4.3 Extraction Stage LLM Analysis

Whenever users click on any fact or any entity or relationship during Merge Stage LLM Analysis, the source chunk information of the original entity or relationship and the behavior details of the LLM in the extraction stage are displayed in the LLM Invocation View (see Fig 7A). Then users can delve into how the LLM extracts entities and relationships from the original text using this information. This transparency helps users evaluate the reliability and accuracy of the information while also laying the foundation for further data quality improvement. Additionally, this view can help users trace the source of specific information, thereby better understanding and verifying the system's inference process.

## 6 EVALUATION

Here, we verify the effectiveness of our method through a usage scenario demonstration and a formal user study.

### 6.1 Usage Scenario: Analysis on MultiHop-RAG Dataset

This usage scenario illustrates how our system aids users in identifying and investigating issues related to graphs constructed by the GraphRAG system, thereby gaining insights into improving graph quality. Given the inherent complexity of the GraphRAG analysis process, we have also produced a demonstration video for this section to facilitate understanding.[1] The user employs the *MultiHop-RAG* news dataset [32] to construct a graph, intending to explore its quality and potential issues using our system.

The user inputs the query "Which entity previously focused on illegal content and misinformation related to the Israel-Hamas conflict and is enforcing more than two pieces of legislation (act) in the digital domain?" along with the Ground Truth Answer *European Commission* and two key original text excerpts into the inference trace view for query and analysis.

Initially, the user observes that the entities *Israel* and *Hamas* were extracted from the query, ultimately leading to the answer *Meta*. This answer is deemed *wrong* by our system, which provides the reasoning for this judgment, highlighting two contradictions.

1. *Meta* is incorrectly identified as focusing on illegal content and misinformation about the *Israel-Hamas* conflict; this is actually the *European Commission*'s role. (From Fact 1)

2. *Meta* is mistakenly thought to implement more than two digital acts, a responsibility that belongs to the *European Commission*. (From Fact 2)

It is noted that the query consists of two distinct parts: the first part tends to query the local association between explicit entities, while the second part aims to derive an answer based on a global summary topic. In this example, both parts result in contradictions, serving as

an excellent illustration. Subsequently, the user investigates these two parts separately for local and global exploration.

### 6.1.1 Exploration of Local Association Contradictions

First, the reason why the segment "Which entity previously focused on illegal content and misinformation related to the Israel-Hamas conflict" failed to yield the correct answer should be analyzed. Issues with local queries can generally be attributed to a lack of entity or relationship information, so we will focus on exploring these two aspects.

By comparing inference steps, the user discovers that the steps related to *Israel-Hamas* in the Actual Answer inference are steps 1, 3, 4, and 5. None of these steps utilized the *European Commission* entity. However, in step 2 of the Ground Truth inference, entities such as *Hamas* and *European Commission* are used as recalls, indicating the presence of the *European Commission* entity in the graph. Still, there is an omission when attempting to recall the *European Commission* entity through *Israel* and *Hamas* entities.

Thus, the user seeks to explore whether there is a local association between the *European Commission* and *Israel* and *Hamas*. In the Inference Trace View, the user right-clicks to recall the *European Commission* entity and begins exploring its local relationships in the Entity Explore View. By hovering over *Israel* and *Hamas*, the user finds no highlighted edges or nodes in the Entity Explore View, indicating no direct relationship between *Israel* and *Hamas* with the *European Commission*, resulting in a missing relationship on the graph.

Next, the user investigates the cause of the missing relationship. By left-clicking to recall the *European Commission* entity in the Inference Trace View, the LLM Invocation View (for the merge stage) reveals the process of merging this entity from 11 raw entities. Hovering over "Fact 1" (a text segment describing the relationship between the *European Commission* and *Israel-Hamas*) in the Inference Trace View highlights one of the raw entities, indicating it originates from the chunk containing "Fact 1". However, the type and description of this raw entity are empty.

By clicking on "Fact 1", the LLM Invocation View (for the extraction stage) displays the chunks where the raw entity resides. It is found that the chunk does not extract any relationship with *Israel* or *Hamas*. Further reading of chunk text reveals that due to changes in subject and pronoun, the model fails to understand the relationship among the three, resulting in the failure of entity and relationship extraction and ultimately causing the loss of relational information.

### 6.1.2 Exploration of Global Semantic Contradictions

Next, the reason why the segment "is enforcing more than two pieces of legislation (act) in the digital domain" failed to yield the correct answer should be analyzed. Note that this answer originates from "Fact 2", where the *European Commission* is responsible for enforcing the *DMA* and *DSA* digital acts on *Meta*.

In the second step of the Actual Answer inference, information about *Meta* fulfilling the *DMA* and *DSA* acts is indeed included, but the *European Commission* is not mentioned. Further examination reveals that this inference step uses a topic report recall. Clicking on this topic recall, the LLM behavior analysis view shows the process of summarizing the topic into a report. Reading the report reveals descriptions of *Meta*, the *DSA*, and the *DMA*, but no mention of the *European Commission*. Further inspection of all entities and relationships used in this topic, hovering over "Fact 2", highlights entities and relationships related to "Fact 2" in the list, but the *European Commission* is not found among them, indicating that the topic does not include the *European Commission*, and there is no global semantic association between the *European Commission* and the two acts.

Subsequently, the user further investigates the cause of semantic omission. In the previous analysis, the user suspected that the graph lacked the relationship from "Fact 2", which includes the *European Commission* and the two acts. By left-clicking "Fact 2", the LLM Invocation View (for the extraction stage) displays the extraction process of the chunk. The user finds that although the chunk explicitly mentions the *European Commission*'s responsibility for enforcing the two acts, the model fails to accurately capture this critical information during

entity and relationship extraction. This discovery further confirms the user's previous suspicion that the LLM overlooked the extraction of the corresponding relationship.

## 6.2 User Evaluation

This study aims to evaluate the *XGraphRAG* performance in analyzing *GraphRAG* results, focusing on:

- **Efficiency**: The speed and simplification of analysis.

- **Effectiveness**: The comprehensiveness and accuracy of the system's support.

- **Usability**: The system's user-friendliness and ease of use.

It also compares *XGraphRAG* with a baseline system for analysis support and user experience.

### 6.2.1 Participants

We invited 14 RAG engineers or experts as participants. These individuals have a thorough understanding of RAG principles but were not involved in the design requirements phase of this study, thereby enhancing the assessment's validity and the results' generalizability. Some have experience with naive RAG systems, while others are familiar with GraphRAG systems.

### 6.2.2 Baseline Systems and Experimental Setup

An additional baseline system was set up for comparative study [9, 10] alongside our system. Both systems are based on Microsoft's *GraphRAG* [6], constructed using the *MultiHop-RAG* dataset to build identical graphs, and utilize GPT-4 as the default LLM. The test dataset used in the experiments also comes from the *MultiHop-RAG* dataset, including questions, answers, and graphs.

Baseline: We used *Kotaemon*, a popular GraphRAG visualization interface system with high ratings on GitHub, as the baseline system. This system provides a visual question-and-answer interface for *GraphRAG* [6], *Nano-GraphRAG* [43], *LightRAG* [12], and other systems. Users integrate GraphRAG for Q&A through this system and display recalled entities, relationships, and topics referenced in the answers in the form of graph diagrams. It also shows various detailed information about recalls, such as explanations of entities. Users can click on entities and relationships in the graph to view the original chunk from which they originated.

Test Dataset: To test the ability of the GraphRAG system to handle both local and global questions, we improved the test questions of the test dataset. Firstly, we selected questions whose answers were names of people or organizations and ensured that these names had corresponding entities in the graph, with sufficient related text in the dataset. Then, we input these texts into the GPT-4 model to generate summary reports related to the entities. Subsequently, we verified the accuracy of the generated reports and posed global questions to them in a manner similar to the GraphRAG approach. Combining the original local questions, we formed the final test questions. Finally, to ensure the validity of the questions, we submitted these questions and related texts to three human experts for review. If the three experts could independently provide clear answers based on the texts and their answers were consistent, the question was deemed valid and could be used as a test case.

### 6.2.3 Procedure and Tasks

**Introduction (10 min):** We began by introducing the system to the participants, explaining the research motivation and methodology, and collecting basic information such as age and gender. With the participants' consent, we recorded and analyzed their usage behavior during subsequent tasks. Finally, we provided a detailed description of the features and characteristics of both the baseline system and the *XGraphRAG* [6] system, demonstrating their use through specific example scenarios.

**Task-based Analysis (60 min):** In this phase, participants interacted with both the baseline system and the *XGraphRAG* [6] system in a randomized order to eliminate any prior knowledge effects. The tasks

were designed to assess the effectiveness and usability of each system, ensuring a fair level of complexity. We recorded the completion time and accuracy for each task to facilitate comparative analysis.

**Semi-structured Interview (30 min):** To further evaluate the interface design and usability of the systems, participants completed a questionnaire consisting of 6 items, using a five-point Likert scale to capture their perceptions and satisfaction. Participants rated each item from 1 (strongly disagree) to 5 (strongly agree). Concurrently, an open feedback session was conducted, allowing participants to explain the reasons behind their ratings and provide detailed insights into their user experience.

### 6.2.4 Task Completion Analysis

| No. | Task | Accuracy | | Time Consuming (s) | |
|---|---|---|---|---|---|
| T1 | What are the deficiencies of the Actual Answer compared to the Ground Truth? | 14/14 | 14/14 | 23.00 | 52.79 |
| T2 | What Suspicious Recalls occurred during the Recall and Infer stage? | 14/14 | 14/11 | 54.36 | 180.14 |
| T3 | What Entity is missing in the graph that leads to the incorrect answer? | 14/11 | 14/7 | 177.57 | 408.43 |
| T4 | What Relationship is missing in the graph that leads to the incorrect answer? | 14/12 | 14/0 ✗ | 213.71 | |
| T5 | What Topic is missing in the graph that leads to the incorrect answer? | 14/8 | 14/0 ✗ | 381.07 | |
| T6 | At which stage does the LLM exhibit anomalies that result in missing information on the graph? | 14/10 | 14/0 ✗ | 308.29 | |

■ **XGraphRAG**   ■ **Baseline**

Table 1: Accuracy and time consumption comparison result between XGraphRAG and the baseline system.

For each task (T1-T6, in Table 1) and each participant, we randomly assigned a test case to explore. Before starting any task, users are provided with all necessary prerequisite information and acknowledge that they understand it. For example, "It is already known that the question was answered incorrectly because the Google entity was not recalled; please try to find information related to the missing Google entity." (Example of T3) and "It is known that the relationship between Taylor Swift and Travis Kelce is missing in the chunk where Fact2 is located; please explore the reason." (Example of T6). Participants need to arrive at an answer within the stipulated time to consider the task complete, such as "The relationship of Epic filing a monopoly lawsuit against Google in Fact1 is missing." (Answer to the example of T3).

Although participants varied in accuracy and time consumption in tasks, using *XGraphRAG* significantly improved task accuracy while reducing time consumption. Below, we analyze each type of task and provide the setup of task problems and objectives.

**Analysis of Suspicious Recalls (T1 - T2):** This set of tasks focuses on the system's ability to identify contradictions and suspicious recalls in the results and their locations. Compared to the baseline system, *XGraphRAG* performed more efficiently and accurately. Its graphical interface allows participants to quickly highlight contradictions in retrieval results in T1, while enhanced visualization tools help locate suspicious recalls in T2, thus reducing the need for extensive manual review and improving diagnostic accuracy ($p_1 = 0.00012$, $p_2 = 0.00012$).

**Analysis of Missing Information (T3 - T5):** This set of tasks emphasizes the system's ability to help participants identify missing information in the graph. Compared to the baseline system, *XGraphRAG* significantly improved time efficiency and success rate, especially in T4 and T5, where the baseline system's inability to analyze relationships and topic reports posed a considerable challenge for participants in analyzing global relevance. The visual encoding in *XGraphRAG* helps participants effectively locate areas of missing nodes or relationships. This design provides clear visual cues, significantly reducing the effort required to locate and analyze incomplete parts of the graph ($p_3 = 0.0033$, $p_4 = 0.0022$, $p_5 = 0.012$).

**Analysis of LLM Behavior (T6)**: In this task, the baseline system lacks the functionality to present and analyze the behavior of the GraphRAG framework model, unable to assist participants in tracking factors leading to missing information in the graph. Using *XGraphRAG*, participants can effectively identify and analyze LLM behavior at various stages of the GraphRAG process and accurately pinpoint issues related to missing information ($p_6 = 0.0051$).

Overall, the *XGraphRAG* system excels in both accuracy and time efficiency in task completion. These findings indicate that the enhanced visualization tools and graphical reasoning capabilities of *XGraphRAG* significantly improve participant efficiency and accuracy in complex analytical scenarios.

## 6.3 Semi-structured Interview Analysis

| Question | Average | Score Distribution |
|---|---|---|
| Q1. QA View can help me to quickly understand the flaw of the answer. | 4.42 | 2 — 6 — 7 |
| Q2. Inference Trace View helps me quickly identify Suspicious Recall and Inference issues. | 4.29 | 2 — 5 — 7 |
| Q3. Global Explore View assists me in exploring the global relevance and topic structure. | 3.79 | 4 — 7 — 3 |
| Q4. Local Explore View aids me in investigating the local correlation of Suspicious Entity Recall. | 4.36 | 2 — 6 — 6 |
| Q5. LLM Invocation View allows me to analyze the detailed behavior of the LLM at each stage. | 4.29 | 1 — 4 — 9 |
| Q6. The system is useful. | 4.58 | 4 — 8 |
| Q7. The system is easy to learn. | 3.53 | 2 — 3 — 6 — 3 |
| Q8. I will use this system again. | 4.21 | 1 — 1 — 5 — 7 |

Legend: 1 (Strongly Disagree) · 2 · 3 · 4 · 5 (Strongly Agree)

Table 2: The evaluation result for our system's effectiveness and usability.

We conducted a qualitative evaluation through semi-structured discussions to assess the system's usability and effectiveness. Participants were encouraged to provide detailed feedback on specific features, including their ability to identify reasoning process flaws, the clarity of visualizations, and the ease of interface interaction.

**Effectiveness:** All participants agreed that the Q&A View effectively helped users quickly understand deficiencies in generated answers (Q1). P10 commented, "The Q&A View can automatically make judgments and provide reasons, along with three-point analysis suggestions, which greatly simplifies my thought process." P3 noted, "In previous systems, I had to spend time reading the large outputs provided by the model, but now I can immediately see the key issues." Participants found the Inference Trace View particularly useful for identifying suspicious recalls and reasoning issues (Q2). P2 suggested providing an automatic translation feature for non-native speakers to accelerate understanding and analysis. P9 recommended integrating more context into the Inference Trace View, such as displaying additional information on mouse hover. P1 stated, "Previously, I needed to use tools like neo4j to cross-system trace links for analysis, but now I can accomplish everything in one view, which gives me a sense of security." The Topic Explore View was highly praised for its assistance in exploring global correlations and corpus thematic structures (Q3). P8 emphasized, "I've always wondered how to quickly understand what exactly a corpus is about, and this view has helped me; it's as important as a nautical chart." P11 stated, "Even though the view is unnecessary for local problem analysis, I can quickly understand the distance between two unconnected entities by hovering over them, which is impressive and something I couldn't achieve with traditional graph tools." Participants also appreciated the Entity Explore View for its role in investigating local correlations of suspicious entity recalls (Q4). P13 pointed out, "The view is simple and practical, allowing right-click expansion of entities for continuous exploration, and combined with the Topic Explore View, it enables targeted exploration." P1 suggested adding search filtering functionality to this view to handle large graphs. The LLM Invocation View received positive feedback for its ability to analyze the detailed behavior of large language models at each stage (Q5). P5 stated, "This

is a missing link in GraphRAG and even RAG systems, and we need more interpretability work in LLM systems." P9 remarked, "I want to understand the root cause of issues, and this view clearly shows each step, which is friendly to RAG developers; we no longer need to check lengthy run logs."

**Usability:** Participants unanimously agreed that the system was useful for their tasks (Q6). P6 expressed, "I previously tried tools like Kotaemon, but they only presented graph retrieval results, leaving a gap in presenting the complete chain from construction to retrieval reasoning, which made me reluctant to trust the model's answers. *XGraphRAG* restored my confidence." Participants generally acknowledged the system's ease of use (Q7). We focused on feedback from participants who rated it 2 or 3. P2, P5, and others suggested adding a guided tutorial to help new users get started quickly. P7 mentioned, "I hope there will be a mock example to guide users step-by-step through an analysis during their first use, which would enhance product usability." Most participants indicated they were likely to use the system again (Q8). Two low-scoring participants (Q8) acknowledged the system's functionality but had little need for diagnosing GraphRAG currently, so they were unsure about future use. P14 stated, "XGraphRAG provides inspiration for designing subsequent Q&A products in the industry and accelerates our development and testing processes."

## 7 DISCUSSION

In this section, we discuss the system's generality, limitations, and directions for future improvements.

**System Generalizability:** Although we adopted Microsoft's *GraphRAG* [6] as an integration tool in our study, all designs are based on our proposed general architecture. This architecture is compatible with other mainstream GraphRAG systems. For instance, the six stages of our architecture are reflected in all systems, where recall, entity, and relationship are present and serve as the core functionalities. While topic reports may have different meanings and expressions in different systems, we regard them as a type of global summary recall. Other systems may have additional types of recall, leading to extra stages, but this does not conflict with the GraphRAG architecture.

**Enhancing temporal analysis.** While our system demonstrates satisfactory performance, dynamic changes in entity relationships or topics within a graph, such as the relationship between A and B evolving from friends to enemies, require further investigation. Therefore, it is essential to develop advanced temporal analysis tools to better guide users in exploring the temporal evolution of graphs.

**Expanding recall capabilities.** Although we currently support entity, relationship, and topic views in mainstream GraphRAG systems with promising results, the development of future GraphRAG systems will necessitate supporting more types of recall to accommodate a wider range of scenarios. This expansion represents another critical direction for improving system capabilities.

## 8 CONCLUSION

This paper presents *XGraphRAG*, a visual analysis system designed to enhance the understanding and optimization of GraphRAG processes. By addressing key challenges such as traceability, LLM invocation context, and relevance analysis, *XGraphRAG* offers a comprehensive framework that empowers developers to systematically examine and refine GraphRAG systems. Through our formative study and user evaluations, we demonstrated the system's effectiveness in identifying and resolving graph quality issues, thus improving the precision and comprehensiveness of generated results. The integration of interactive visualizations and structured exploration tools not only facilitates more informed decision-making but also advances the field of visual analytics in complex data environments. Future work will focus on expanding the system's capabilities to support dynamic, real-time datasets and enhance user guidance for open-ended exploration.

## REFERENCES

[1] A. Abd-Alrazaq, R. AlSaad, D. Alhuwail, A. Ahmed, P. M. Healy, S. Latifi, S. Aziz, R. Damseh, S. A. Alrazak, J. Sheikh, et al. Large language models in medical education: opportunities, challenges, and future directions. *JMIR Medical Education*, 9(1):e48291, 2023. doi: 10.2196/48291 1

[2] M. T. Agler, B. A. Wrenn, S. H. Zinder, and L. T. Angenent. Waste to bioproduct conversion with undefined mixed cultures: the carboxylate platform. *Trends in biotechnology*, 29(2):70–78, 2011. doi: 10.1016/j.tibtech.2010.11.006 1

[3] M. Alhanahnah and Y. Boshmaf. Depsrag: Towards agentic reasoning and planning for software dependency management. *CoRR*, abs/2405.20455, 2024. doi: 10.48550/arXiv.2405.20455 2

[4] I. Arawjo, C. Swoopes, P. Vaithilingam, M. Wattenberg, and E. L. Glassman. Chainforge: A visual toolkit for prompt engineering and llm hypothesis testing. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, article no. 304, 18 pages. Association for Computing Machinery, New York, NY, USA, 2024. doi: 10.1145/3613904.3642016 2

[5] R. Brath, A. Bradley, and D. Jonker. Visualizing llm text style transfer: visually dissecting how to talk like a pirate. In *Proceedings of IEEE VIS 2023*, 2023. 2

[6] 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. doi: 10.48550/ARXIV.2404.16130 1, 2, 3, 8, 9

[7] W. Fan, Y. Ding, L. Ning, S. Wang, H. Li, D. Yin, T. Chua, and Q. Li. A survey on RAG meeting llms: Towards retrieval-augmented large language models. In R. Baeza-Yates and F. Bonchi, eds., *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024*, pp. 6491–6501. ACM, 2024. doi: 10.1145/3637528.3671470 1

[8] Y. Feng, Z. Chen, Z. Kang, S. Wang, M. Zhu, W. Zhang, and W. Chen. Jailbreaklens: Visual analysis of jailbreak attacks against large language models. *arXiv preprint arXiv:2404.08793*, 2024. doi: 10.48550/arXiv.2404.08793 2

[9] Y. Feng, X. Wang, B. Pan, K. K. Wong, Y. Ren, S. Liu, Z. Yan, Y. Ma, H. Qu, and W. Chen. Xnli: Explaining and diagnosing nli-based visual data analysis. *IEEE Trans. Vis. Comput. Graph.*, pp. 1–14, 2023. doi: 10.1109/TVCG.2023.3240003 8

[10] Y. Feng, X. Wang, K. Wong, S. Wang, Y. Lu, M. Zhu, B. Wang, and W. Chen. Promptmagician: Interactive prompt engineering for text-to-image creation. *IEEE Trans. Vis. Comput. Graph.*, 30(1):295–305, 2024. doi: 10.1109/TVCG.2023.3327168 8

[11] 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. *CoRR*, abs/2312.10997, 2023. doi: 10.48550/ARXIV.2312.10997 1

[12] Z. Guo, L. Xia, Y. Yu, T. Ao, and C. Huang. Lightrag: Simple and fast retrieval-augmented generation. *arXiv preprint arXiv:2410.05779*, 2024. doi: 10.48550/arXiv.2410.05779 1, 3, 8

[13] B. J. Gutiérrez, Y. Shu, Y. Gu, M. Yasunaga, and Y. Su. Hipporag: Neurobiologically inspired long-term memory for large language models. *CoRR*, abs/2405.14831, 2025. doi: 10.48550/arXiv.2405.14831 2

[14] X. He, Y. Tian, Y. Sun, N. V. Chawla, T. Laurent, Y. LeCun, X. Bresson, and B. Hooi. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *CoRR*, abs/2402.07630, 2024. doi: 10.48550/ARXIV.2402.07630 2

[15] M. Kahng, I. Tenney, M. Pushkarna, M. X. Liu, J. Wexler, E. Reif, K. Kallarackal, M. Chang, M. Terry, and L. Dixon. Llm comparator: Interactive analysis of side-by-side evaluation of large language models. *IEEE Transactions on Visualization and Computer Graphics*, 31(1):503–513, 2025. doi: 10.1109/TVCG.2024.3456354 2

[16] Z. Li, L. Deng, H. Liu, Q. Liu, and J. Du. Unioqa: A unified framework for knowledge graph question answering with large language models, 2024. 2

[17] Z. Li, Q. Guo, J. Shao, L. Song, J. Bian, J. Zhang, and R. Wang. Graph neural network enhanced retrieval for question answering of llms. *CoRR*, abs/2406.06572, 2024. doi: 10.48550/arXiv.2406.06572 2

[18] L. Liang, M. Sun, Z. Gui, Z. Zhu, Z. Jiang, L. Zhong, Y. Qu, P. Zhao, Z. Bo, J. Yang, H. Xiong, L. Yuan, J. Xu, Z. Wang, Z. Zhang, W. Zhang, H. Chen, W. Chen, and J. Zhou. KAG: boosting llms in professional domains via knowledge augmented generation. *CoRR*, abs/2409.13731, 2024. doi: 10.48550/ARXIV.2409.13731 3

[19] Y. Liu, Z. Wen, L. Weng, O. Woodman, Y. Yang, and W. Chen. SPROUT: an interactive authoring tool for generating programming tutorials with the visualization of large language models. *IEEE Transactions on Visualization and Computer Graphics*, pp. 1–15, 2024. doi: 10.1109/TVCG.2024.3410523 2

[20] J. Lu, B. Pan, J. Chen, Y. Feng, J. Hu, Y. Peng, and W. Chen. AgentLens: Visual analysis for agent behaviors in llm-based autonomous systems. *CoRR*, abs/2402.08995, Feb. 2024. doi: 10.48550/arXiv.2402.08995 2

[21] S. Ma, C. Xu, X. Jiang, M. Li, H. Qu, C. Yang, J. Mao, and J. Guo. Think-on-graph 2.0: Deep and faithful large language model reasoning with knowledge-guided retrieval augmented generation. *arXiv preprint arXiv:2407.10805*, 2024. doi: 10.48550/arXiv.2407.10805 2

[22] B. Pan, J. Lu, K. Wang, L. Zheng, Z. Wen, Y. Feng, M. Zhu, and W. Chen. Agentcoord: Visually exploring coordination strategy for llm-based multi-agent collaboration. *CoRR*, abs/2404.11943, 2024. doi: 10.48550/arXiv.2404.11943 2

[23] B. Peng, Y. Zhu, Y. Liu, X. Bo, H. Shi, C. Hong, Y. Zhang, and S. Tang. Graph retrieval-augmented generation: A survey. *arXiv preprint arXiv:2408.08921*, 2024. doi: 10.48550/arXiv.2408.08921 1

[24] T. T. Procko and O. Ochoa. Graph retrieval-augmented generation for large language models: A survey. *Available at SSRN*, pp. 166–169, 2024. doi: 10.1109/AIxSET62544.2024.00030 1

[25] P. Ranade and A. Joshi. Fabula: Intelligence report generation using retrieval-augmented narrative construction. In *Proceedings of the International Conference on Advances in Social Networks Analysis and Mining*, ASONAM '23, p. 603–610. ACM, Nov. 2023. doi: 10.1145/3625007.3627505 2

[26] M. Sallam. Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In *Healthcare*, vol. 11, p. 887. MDPI, 2023. doi: 10.3390/healthcare11060887 1

[27] B. Sarmah, D. Mehta, B. Hall, R. Rao, S. Patel, and S. 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, ICAIF 2024, Brooklyn, NY, USA, November 14-17, 2024*, pp. 608–616. ACM, 2024. doi: 10.1145/3677052.3698671 2

[28] H. Strobelt, A. Webson, V. Sanh, B. Hoover, J. Beyer, H. Pfister, and A. M. Rush. Interactive and visual prompt engineering for ad-hoc task adaptation with large language models. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):1146–1156, 2023. doi: 10.1109/TVCG.2022.3209479 2

[29] S. Suh, B. Min, S. Palani, and H. Xia. Sensecape: Enabling multilevel exploration and sensemaking with large language models. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, article no. 1, 18 pages. Association for Computing Machinery, New York, NY, USA, 2023. doi: 10.1145/3586183.3606756 2

[30] J. Sun, C. Xu, L. Tang, S. Wang, C. Lin, Y. Gong, H.-Y. Shum, and J. Guo. Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph. *CoRR*, abs/2307.07697, 2023. doi: 10.48550/arXiv.2307.07697 2

[31] Y. Tang and Y. Yang. Mdalk: Dynamic co-augmentation of llms and kg to answer alzheimer's disease questions with scientific literature. *CoRR*, abs/2405.04819, 2024. doi: 10.48550/arXiv.2405.04819 2

[32] Y. Tang and Y. Yang. Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries. *CoRR*, abs/2401.15391, 2024. doi: 10.48550/ARXIV.2401.15391 7

[33] K. Team. Kotaemon team, 2024. 3

[34] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023. doi: 10.1038/s41591-023-02448-8 1

[35] D. Tiro. The possibility of applying chatgpt (ai) for calculations in mechanical engineering. In *International Conference "New Technologies, Development and Applications"*, pp. 313–320. Springer, 2023. 1

[36] X. Wang, N. Anwer, Y. Dai, and A. Liu. Chatgpt for design, manufacturing, and education. *Procedia CIRP*, 119:7–14, 2023. doi: 10.1016/j.procir.2023.04.001 1

[37] X. Wang, Z. Wu, W. Huang, Y. Wei, Z. Huang, M. Xu, and W. Chen. Vis+ai: integrating visualization with artificial intelligence for efficient data analysis. *Front. Comput. Sci.*, 17(6), 12 pages, June 2023. doi: 10.1007/s11704-023-2691-y 4

[38] Y. Wardat, M. A. Tashtoush, R. AlAli, and A. M. Jarrah. Chatgpt: A revolutionary tool for teaching and learning mathematics. *Eurasia Journal of Mathematics, Science and Technology Education*, 19(7):em2286, 2023. doi: 10.29333/ejmste/13272 1

[39] J. Wu, J. Zhu, and Y. Qi. Medical graph RAG: towards safe medical large language model via graph retrieval-augmented generation. *CoRR*, abs/2408.04187, 2024. doi: 10.48550/ARXIV.2408.04187 1

[40] S. Wu, O. Irsoy, S. Lu, V. Dabravolski, M. Dredze, S. Gehrmann, P. Kambadur, D. S. Rosenberg, and G. Mann. Bloomberggpt: A large language model for finance. *CoRR*, abs/2303.17564, 2023. doi: 10.48550/ARXIV.2303.17564 1

[41] Z. Xu, M. J. Cruz, T. W. Matthew Guevara, M. Deshpande, X. Wang, and Z. Li. Retrieval-augmented generation with knowledge graphs for customer service question answering. *CoRR*, abs/2404.17723, 2024. doi: 10.1145/3626772.3661370 2

[42] H. Yang, X. Liu, and C. D. Wang. Fingpt: Open-source financial large language models. *CoRR*, abs/2306.06031, 2023. doi: 10.48550/ARXIV.2306.06031 1

[43] G. Ye. nano-graphrag: A simple, easy-to-hack GraphRAG implementation. https://github.com/gusye1234/nano-graphrag, 2024. 8