# How Significant Are the Real Performance Gains?
# An Unbiased Evaluation Framework for GraphRAG

Qiming Zeng[†]    Xiao Yan[#]    Hao Luo[†]    Yuhao Lin[†]    Yuxiang Wang[†]
Fangcheng Fu[§]    Bo Du[†]    Quanqing Xu[‡]    Jiawei Jiang[†,*]

[†]School of Computer Science, Wuhan University    [§]School of Computer Science, Peking University
[#]Centre for Perceptual and Interactive Intelligence (CPII)    [‡]OceanBase

[†] {kirin_z,lohozz,yuhao_lin,nai.yxwang,dubo,jiawei.jiang}@whu.edu.cn
[#]yanxiaosunny@gmail.com    [§]ccchengff@pku.edu.cn    [‡]xuquanqing.xqq@oceanbase.com

arXiv:2506.06331v1 [cs.CL] 31 May 2025

## Abstract

By retrieving contexts from knowledge graphs, graph-based retrieval-augmented generation (GraphRAG) enhances large language models (LLMs) to generate quality answers for user questions. Many GraphRAG methods have been proposed and reported inspiring performance in answer quality. However, we observe that the current answer evaluation framework for GraphRAG has two critical flaws, i.e., *unrelated questions* and *evaluation biases*, which may lead to biased or even wrong conclusions on performance. To tackle the two flaws, we propose an unbiased evaluation framework that uses *graph-text-grounded question generation* to produce questions that are more related to the underlying dataset and an *unbiased evaluation procedure* to eliminate the biases in LLM-based answer assessment. We apply our unbiased framework to evaluate 3 representative GraphRAG methods and find that their performance gains are much more moderate than reported previously. Although our evaluation framework may still have flaws, it calls for scientific evaluations to lay solid foundations for GraphRAG research.

## Keywords
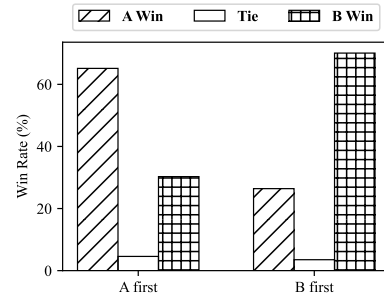
LLMs, RAG, GraphRAG, performance evaluation

## 1 Introduction

Large language models (LLMs) have achieved tremendous success in generating answers to user questions [11, 12, 38]. However, they may produce wrong or fictitious responses when lacking specific knowledge to answer the questions [18, 28, 30], e.g., due to outdated training data or missing domain-specific knowledge. In these cases, *retrieval-augmented generation* (RAG) [21] improves answer quality by retrieving contexts related to the user question from an external database and feeding these contexts along with the user question to the LLMs for answer generation. Among RAG methods, *graph-based RAG* (a.k.a., GraphRAG) is particularly popular, which uses knowledge graphs as the external databases [15, 16, 25, 36]. This is because knowledge graph is a powerful format of knowledge expression, where the graph nodes represent real-world entities and edges represent the complex relations between the entities [4, 8].

Many GraphRAG methods have been proposed, and they mainly differ in how to traverse the knowledge graph and retrieve the related contexts. For instance, Microsoft GraphRAG (called *MGRAG*

---

∗ Corresponding author.



Figure 1: Questions produced for Harry Potter novel by the summary-based method in current evaluation framework.



Figure 2: The position bias for LLM-based answer assessment. We compare the answers generated by LightRAG in 2 runs for the same questions but change the order of the two answers (i.e., 'AB' and 'BA') when feeding them to the LLM.

hereafter to avoid confusion with GraphRAG) [7] performs community detection on the knowledge graph and generates text summaries for the communities offline, and during online question answering, it selects the community summaries that have large relevance scores for the user question. LightRAG [10] first extracts keywords referring to entities and relations from the user question and then retrieves these entities and relations along with their 1-hop neighbors from the knowledge graph. FastGraphRAG (called *FGRAG* hereafter) [6] first extracts entities from the user question, then retrieves the relevant contents from the knowledge graph according to semantic similarity w.r.t. the extract entities (measured by their vector embeddings), and finally ranks the contents according to their PageRank scores to select the top ones. GraphCot [19] adopts chain of thought [32] and involves an agent (which is an LLM) to determine the node or edge to check in each graph traversal step. Please refer to [13, 25] for a survey of GraphRAG methods.
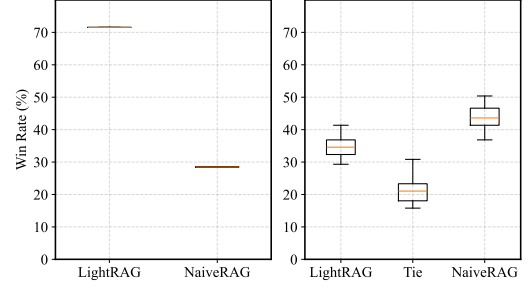
Since human inspection is expensive, GraphRAG methods are usually evaluated with LLMs for their answer quality. In particular, an evaluation framework consists of two parts, i.e., *question generation* and *answer assessment*. For question generation, MGRAG [7] and LightRAG [10] extract the head and tail of each passage in the dataset as the passage summary, and the summaries of all passages are aggregated as dataset summaries, which are fed to an LLM for question generation. For answer assessment, an LLM is requested to compare the answers produced by two RAG methods for the same question, in aspects such as *Comprehensiveness*, *Diversity*, *Empowerment*, and *Directness*, e.g., using a prompt like '*Please choose the better answer from A and B for question Q in term of ⋯*', where *A* and *B* are the answers produced by two RAG methods. An overall winner is determined based on the separate winners of each aspect, and the win rate of a method is the percentage of questions it wins.

We observe that current method for GraphRAG evaluation has two critical flaws that hinder scientific performance assessment.

- *Unrelated questions*. Figure 1 lists some example questions generated using the *summary-based method* in existing evaluations for the Harry Potter novel. It is obvious that these questions are overly broad and vague. Moreover, they are not (closely) related to the Harry Potter story. This is because the LLM is only provided with vague summaries of the passages in the dataset for question generation, and as a result, the questions do not involve the fine-grained details of the passages. However, GraphRAG is expected to retrieve fine-grained details from knowledge graphs to supplement the LLMs for answer generation, and such ability can not be evaluated using the questions in Figure 1.

- *Evaluation biases*. We find that the current evaluation protocol has three biases that hinder scientific performance assessment, i.e., *position bias*, *length bias*, and *trial bias*. In particular, position bias means that simply changing the positions of the two compared answers in the evaluation prompt has a significant impact on the evaluation result. We show such an example in Figure 2, where the positions of two answers '*A*' and '*B*' are changed from '*AB*' to '*BA*'. Since we compare the same GraphRAG method (i.e., LightRAG) with itself, the result should be a tie. However, LLMs favor the answer that comes in the front, and the win rates can differ by over 30%. Similarly, length bias means that LLMs favor longer answers, and trial bias means that an LLM can produce different results when evaluating two answers multiple times.

To tackle the flaws of the current GraphRAG evaluation, we design an unbiased evaluation framework with two innovations.

- *Graph-text-ground question generation*. To ensure that the questions correspond to fine-grained details of the dataset, we use the knowledge graph derived from the dataset to guide question generation. In particular, to generate a question, we first sample a structure from the knowledge graph (i.e., node, edge, or subgraph) and then feed the structure and its related contexts to an LLM for question generation. This graph-text-grounded method allows us to configure the type of the generated question, i.e., related to a node, a relation, or a subgraph. We show with concrete examples that our questions are much more related to the dataset than those generated using the summary-based method.



**Figure 3: The evaluation results using the existing method (left) and our unbiased framework (right). Note that existing method does not allow ties. In the box plot (right), the red line is the median, and the lower and upper box boundaries are the 25th and 75th percentiles, respectively.**

- *Unbiased evaluation procedure*. We design a comprehensive evaluation procedure to eliminate the aforementioned biases. In particular, the position bias is handled by evaluating both '*AB*' and '*BA*', and a tie is declared if the overall scores obtained by the two answers are the same. For the length bias, we first let two methods generate their answers and then perform answer alignment to align the shorter answer with the longer one in length. For the trial bias, we conduct multiple rounds of evaluation with each round covering all the questions and report the statistics of the rounds (e.g., median, and the percentiles).

We applied our unbiased evaluation framework to compare 3 representative GraphRAG methods, i.e., MGRAG, LightRAG, and FGRAG, and directly retrieving the text chunks without knowledge graph (called NaiveRAG) [9] is included as a baseline. The results show that the significant performance gains reported by existing research may be caused by the biases, and after eliminating the biases, the performance gains become much more moderate or even vanish. Figure 3 provides such an example, where LightRAG originally outperforms NaiveRAG with a win rate of 72% vs. 28% but NaiveRAG slightly outperforms LightRAG after eliminating the biases. This is because LightRAG usually generates longer answers than NaiveRAG, and the original evaluation always places its answer before NaiveRAG. We do observe that FGRAG outperforms the others but again the performance gains are moderate, i.e., below 10% in win rate expect when compared with LightRAG.

To summarize, we make the following contributions:

- We observe that the current method for GraphRAG performance evaluation has critical flaws, i.e., *unrelated questions* and *evaluation biases*, that may lead to biased or wrong conclusions.

- We design a new evaluation framework for GraphRAG, which tackles the flaws of the existing method with *graph-text-ground question generation* and *unbiased evaluation procedure*.

- We compare representative GraphRAG methods using our evaluation framework. The results suggest that existing GraphRAG researches are overly optimistic about their performance gains.

We note that our evaluation framework may still be far from perfect. However, given the rapid advances of GraphRAG methods, our work does show that a scientific evaluation framework for their

performance is lacking, and without it, we cannot assess whether GraphRAG researches make real progress. We hope that our work can inspire more work on performance evaluation.

## 2  Current Evaluation Method for GraphRAG

GraphRAG evaluation consists of two main components, i.e., *question generation* and *answer evaluation* [35]. Specifically, question generation produces a set of user queries to test the GraphRAG methods, while answer evaluation asses the quality of the answers generated by the GraphRAG methods [27].

For question generation, MGRAG [7] and LightRAG [10] employ a summary-based method with three steps. The first step extracts the beginnings and ends of the articles in the dataset. For an article, the beginning part usually outlines the topic and provides background information, while the ending part typically contains the conclusion or summary. In the second step, the extracted segments are concatenated to form a summary of each article and a dataset description is obtained by stacking the summaries of all articles. The third step feeds the dataset description to an LLM for question generation, using prompts like '*Given the following description of a dataset, please generate questions for potential user · · ·*'.

Answer evaluation is usually conducted by using an LLM as automated judge [1, 2, 5, 20, 23], and two answers for the same question are compared each time [3, 14, 17, 24, 27, 31]. In particular, the GraphRAG methods are first instructed to generate answers for a set of questions. Then, the answers produced by two methods for the same questions are paired for comparison. The LLM is tasked with selecting a superior answer from the pair in some predefined aspects (e.g., *Comprehensiveness*, *Diversity*, *Empowerment*, and *Directness*), and each aspect comes with specific criteria as text instructions. For each aspect, the LLM selects a winner and provides an explanation justifying its decision. Subsequently, the LLM integrates the evaluations across all the aspects to determine an overall winner, e.g., using a prompt like '*select an overall winner based on aspects including · · ·*'. For two GraphRAG methods, win rates are computed as the percentages of questions they win in the direct comparison. We note that existing researches place the answers of two GraphRAG methods in a fixed position (e.g., always 'AB') in the evaluation prompt.

Some researches [6, 19] adopt datasets with predefined question-answer pairs. In these cases, answer quality can be assessed using the similarity between the generated answers and the ground truth answers. However, this approach is susceptible to training data contamination [37], wherein the questions and answers from public datasets may have been included in the LLM's training data, and thus it is difficult to discern whether the generated answers are attributed to the GraphRAG method or the LLM. Human experts can provide highly reliable evaluations for answer quality [29] but are very expensive to employ. As such, an automated evaluation framework is favorable due to its good generality and low cost.

## 3  Graph-text-grounded Question Generation

As shown in Figure 1, the questions generated by the current summary-based method are too macroscopic to be answered using
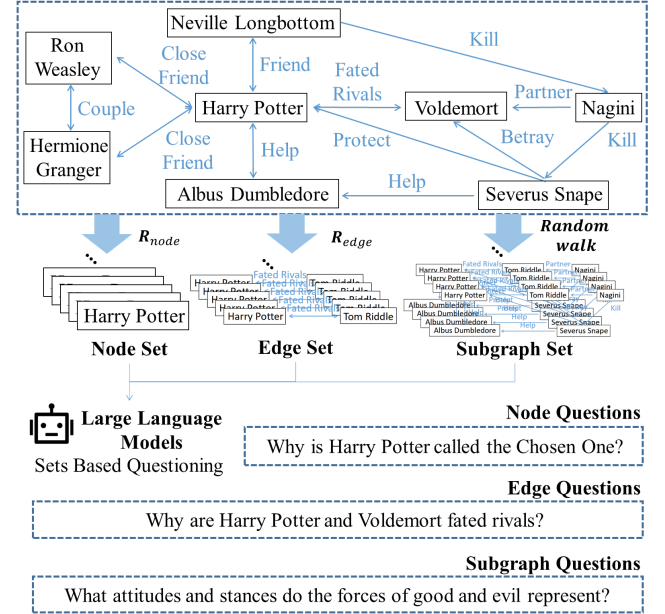


**Figure 4: Workflow of graph-text-grounded question generation.**

the dataset's knowledge or remain relevant. To tackle its defects, we propose a *graph-text-grounded question-generation method*.

**Question Generation.** Our graph-text-grounded question generation leverages the knowledge graph to produce questions that are closely related to the dataset. This process consists of two steps: *knowledge graph construction* and *graph-text-grounded question generation*.
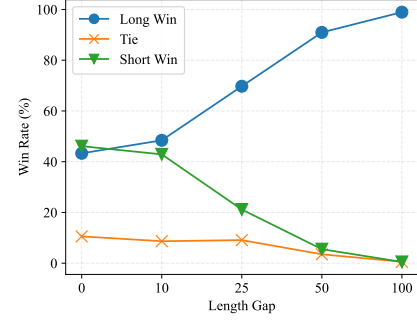
- *Knowledge graph construction.* We first segment the dataset articles into fine-grained chunks and employ an LLM to extract the entities and relations. We prompt the LLM to conduct deeper mining, e.g., using prompts like '*identify the missing entities and add them*'. The extracted entities and relations are enriched with the descriptions derived from their source articles. These descriptions are stored in a knowledge graph, with the duplicate entities merged to maintain consistency and completeness. After establishing the knowledge graph, we will associate each entity (node and relationship) in the knowledge graph with its corresponding original text paragraph to ensure that the source of the text paragraph is not lost.

- *Graph-text-grounded question generation.* Figure 4 shows our procedure for question generation. We leverage the knowledge graph to generate questions at three levels, i.e., *node*, *edge*, and *subgraph*. At the node and edge levels, we conduct random sampling for the entities and relations of the knowledge graph. These sampled entities and relations form context structures, which are then used to prompt the LLM to generate questions. In particular, at the node level, we instruct the LLM to focus on the contents of the entities to generate relevant questions. At the edge level,

| Node Level | • Why is Harry Potter called the Chosen One? |
|---|---|
| | • How does the Invisibility Cloak aid Harry and impact his adventures? |
| Edge Level | • Why are Harry Potter and Voldemort fated rivals? |
| | • How crucial is Prof. McGonagall's mentorship for the Gryffindors? |
| Subgraph Level | • What attitudes and stances do the force of good and evil represent? |
| | • How can Quidditch symbolize teamwork and strategy? |

**Figure 5: Examples of the questions generated by our graph-grounded method for the Harry Potter novel.**

we prompt the LLM to focus on the descriptions of the relations to propose inferential questions. To prevent the quality of the knowledge graph itself from being strongly correlated with the quality of the problem, we will simultaneously take the original text segment corresponding to the entity as supplementary input. In this way, even if the quality of the knowledge graph itself has problems, it can still ensure to the greatest extent that the problem is faithful to the original text. To generate subgraph-level questions, we first sample a node in the knowledge graph and perform random walk algorithm [22] from the seed node to extract a subgraph with multiple walk steps. We enforce a minimum size threshold (which is 50 non-repeating hops in default) for the subgraph such that the subgraph encompasses rich entities and relations. If a sampled subgraph falls below the threshold, we resample to ensure adequate size. Given a subgraph, we use it to form context structures that guide the LLM in generating questions that require an overall understanding and reasoning about the subgraph. For subgraph scenarios, we will first generate summaries of the text segments corresponding to all entities within the subgraph. Since the text segments involved in the subgraph at this time are far less than those directly generated for the full text, the summaries will be more accurate. Next, we use the generated summary as supplementary input to assist in generating subgraph-level questions.

**Question examples.** We utilize the Harry Potter novel as the dataset to compare the questions generated by the existing summary-based method and our graph-text-grounded method. In particular, Figure 5 shows the questions generated by our method, while Figure 1 shows the questions generated by the existing method. It can be observed that the summary-based method produces questions that are too dispersed and sometimes unanswerable or irrelevant given the underlying dataset. Take the last question in Figure 1 for example, it is evident that GraphRAG methods cannot forecast the box office receipts using knowledge from the Harry Potter novel. This is because the summary-based method relies on vague summaries of the articles and dataset, which lose the fine-grained details. This problem becomes more significant for larger datasets due to more articles and increased ambiguity. In contrast, our graph-text-grounded method produces questions that are closely related to multiple levels of the underlying dataset. At the node level, it generates questions focused on the attributes and characteristics of individual entities, testing the ability to retrieve detailed knowledge of specific nodes for GraphRAG methods. At the edge level, it creates



**Figure 6: The effect of length bias. We compare the answers generated by LightRAG from the same questions in two different runs, and the length gap is the token difference between the longer answer and the shorter answer. The average length of the answers is around 200.**

inferential questions that explore relations between entities, such as causality or correlation, assessing relational reasoning abilities. At the subgraph level, it formulates macroscopic questions requiring a holistic understanding of multiple interconnected entities, evaluating the ability to synthesize information and comprehend broader structural and semantic contexts.
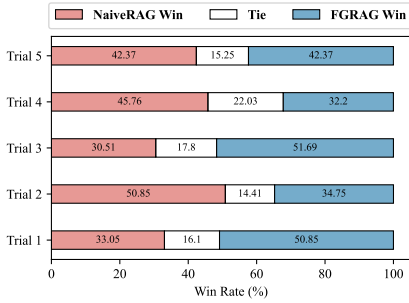
## 4 Unbiased Evaluation Procedure

In this part, we demonstrate the biases in the current evaluation framework and introduce our methods to tackle them. We also present some additional modifications to the evaluation.

### 4.1 Biases in Existing Evaluation

Our examination reveals three types of biases when using an LLM as the judge to asses answer quality, i.e., *position bias*, *length bias*, and *trial bias*. These biases severely hinder the fairness of assessment and can lead to absurd conclusions.

- *Position bias.* As shown in Figure 3, even when comparing the answers generated by the same GraphRAG method, the method whose answers come in front of the evaluation prompt is much more likely to win. This may be because the attention mechanism of LLMs tends to focus on the initial tokens [34].

- *Length bias.* LLMs tend to perceive the longer answer as having a higher quality, and Figure 6 provides an illustration of this phenomenon. For this experiment, we compare answers of different lengths generated by LightRAG. Note that the position bias has already been eliminated using a method that will be presented later. The results show that even when the core contents of the answers are similar, a length gap of 25 tokens, which is relatively small compared with the average answer length of 200 tokens, can lead to a win rate gap of over 50%. This may be because the longer answer has more tokens and thus receives more attention. We also observe that adding meaningless tokens to an answer may improve its chances of winning.

- *Trial bias.* LLMs can produce different win rate results when comparing two GraphRAG methods multiple times. We show this phenomenon in Figure 7, where NaiveRAG and FGRAG are

**Figure 7: The effect of trial bias. We compare NaiveRAG and FGRAG on the Agriculture dataset five times with the same set of questions. The win rates differ across the trials.**

compared five times using the same set of questions. Note that both position bias and length bias have already been eliminated in this experiment. The results show that different trials lead to contradictory conclusions. In Trial 5, NaiveRAG and FGRAG have a tie; FGRAG wins for Trials 1 and 3, but NaiveRAG wins for Trials 2 and 4. With a single trial, we may arrive at any of the 3 possible conclusions, i.e., tie, NaiveRAG wins, and FGRAG wins. The trial bias is caused by the inherent randomness of LLM generation, which involves sampling the possible tokens.

## 4.2 Unbiased Evaluation Framework

To tackle the evaluation biases discussed earlier, we design an unbiased evaluation framework with three main elements, i.e., *length alignment*, *position exchange*, and *trial statistics*. The workflow of our evaluation framework is illustrated in Figure 8.

- *Length alignment.* To tackle length bias, we align their lengths when comparing two answers for the same question. Given that the GraphRAG methods produce answers of different lengths for the same question, imposing a uniform length constraint may impair their answer quality. As such, we adopt a *generate-adjust* strategy. In particular, both methods first generate answers without restrictions to preserve their logic, and then we adjust the shorter answer to match the length of the longer answer. Instead of directly expanding the shorter answer with the LLM, we adjust the response-generating prompt to include the target length. Compared with expanding the answer, this gives the LLM more room to re-organize the retrieved contexts.

  We observe that it is difficult to exactly match the lengths of two answers. Therefore, we set a tolerance threshold for the length difference, which is 10 words by default (observed to have a negligible impact on the evaluation results). If two answers still exceed the threshold after several length adjustments, we apply a forced alignment method. Specifically, we calculate the remaining length discrepancy after the adjustments and instruct the LLM to append tokens to the shorter answer without altering its meaning. Our length aliment procedure can succeed for 85% of the answer pairs, and we discard the failing answer pairs whose length difference still exceeds the threshold after all adjustments.

- *Position exchange.* To address the position bias, we impose position exchange when comparing two answers for the same question. In particular, denote the two answers as $A$ and $B$, position

exchange evaluates both $AB$ and $BA$ using two separate prompts. For each position (i.e., $AB$), the LLM is instructed to score the two answers, and the final score of each answer (e.g., $A$) is its average score over the two positions. Position exchange ensures fairness because both answers enumerate the advantageous position (i.e., front) and the disadvantageous position (i.e., back), and thus the average score faithfully indicates answer quality.

- *Trial statistics.* To address the trial bias, we compare two answers for the same question multiple times. We adopt two kinds of such repetitions. First, when comparing each position of two answers (i.e., $AB$), we issue the same prompt $N$ times and collect the average scores of these prompts as the scores for a certain position. This is intended to account for the randomness in the LLM's evaluations. Second, when comparing two GraphRAG methods, we call it *a trial* when all question-answer pairs are evaluated and obtain win rates. We conduct $M$ such trials to account for the LLM's randomness in generating responses, and thus the evaluation process is re-executed for each trial. By default, we set $N = 2$ and $M = 25$.

For the final evaluation results, we report the statistics of win rates across the trials using the box plot [33], and one example can be found in the right plot of Figure 3. In particular, for a box plot, the middle line is the median, and the box boundaries are the 25% and 75% percentiles. The size of the box (called interquartile range, IQR) reflects the variability of the win rates, and the outliers are depicted as individual points to show extreme deviations. Compared with a single win rate, the box plot provides more information about the results.

## 4.3 Enhancement to the Evaluation Protocol

Besides tackling the evaluation biases, we also enhance the existing evaluation framework with three modifications, i.e., *evaluation perspectives*, *scoring mechanism*, and *tie cases*.

- *Evaluation perspectives.* We initially adopt the perspectives of MGRAG (i.e.,*Comprehensiveness*, *Diversity*, *Empowerment*, and *Directness*) to evaluate answer quality. However, we observe that diversity is not effective when evaluating a single answer; instead, the relevance of an answer to the question is more indicative for the perceived quality. As such, we supplant the *Diversity* aspect with *Relevance*.

- *Scoring mechanism.* Existing evaluation instructs the LLM to directly determine a winner according to the text comments generated for each aspect, which are coarse-grained and lack quantifiable metrics. To address this limitation, we introduce a scoring mechanism that enhances the granularity and objectivity of evaluation. In particular, we ask the LLM to give a score to an answer in each evaluation aspect and provide detailed text instructions to the LLM for scoring. An example of the score grading instructions is provided in Table 1. After scoring all the aspects, we compute the overall score for each answer to determine the overall winner.

- *Tie cases.* Existing evaluation mandates selecting a winner from the two compared answers, which can be restrictive when the two answers have comparable quality. To address this limitation, we introduce tie cases. Specifically, when the total scores of the
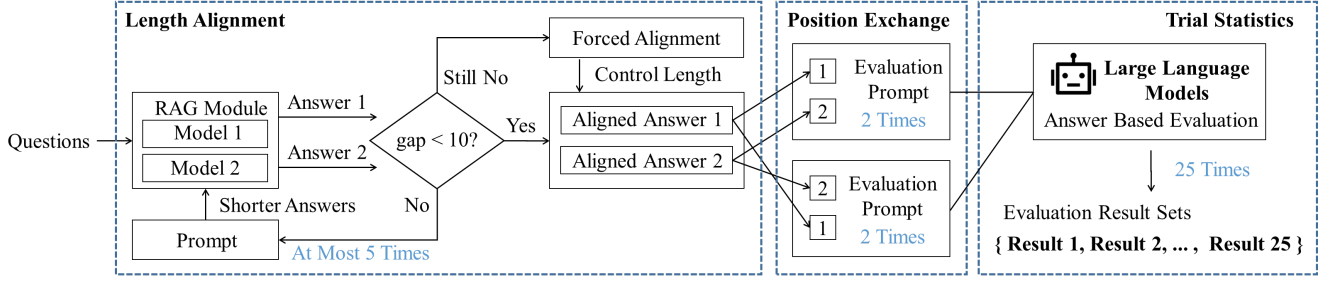
Figure 8: An overview of our unbiased evaluation framework.

Table 1: The grading criteria for the Directness aspect.

| Aspect | Score | Description |
|---|---|---|
| | 0 point: | The answer is extremely indirect, failing to address the question specifically and clearly. |
| | 1 point: | The answer is indirect and deviates significantly from the question, making it hard to discern the intended response. |
| Directness | 2 point: | The answer is somewhat indirect, occasionally straying from the question, but still touching on relevant points. |
| | 3 point: | The answer is moderately direct, addressing the question with some clarity but could be more specific and focused. |
| | 4 point: | The answer is clear and direct, effectively addressing the question with specificity and clarity. |
| | 5 point: | The answer is exceptionally direct, precisely and specifically addressing the question without any ambiguity. |

two answers are the same, we declare a tie, acknowledging that the answers have similar quality.
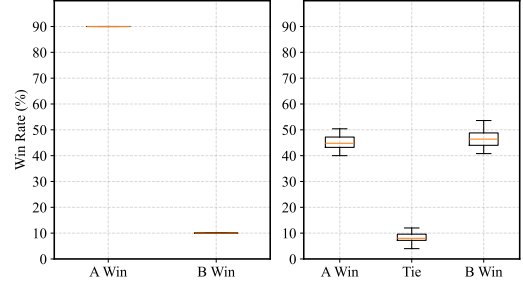
## 5 Experiments

In this part, we adopt our evaluation framework to examine 3 representative GraphRAG methods and show that it leads to conclusions that are radically different from those reported previously.

### 5.1 Experiment Settings

**GraphRAG methods.** We evaluate 3 representative GraphRAG methods, i.e., MGRAG [7], LightRAG [10], and FGRAG [6], which are popular recently. We include NaiveRAG as a baseline, which directly retrieves text chunks according to embedding similarity w.r.t. the question and does not utilize knowledge graph. Details of the GraphRAG methods are provided in Appendix A. We use their publicly accessible codebases to conduct answer generation.

**Datasets.** We use the Ultra Domain Benchmark (UDB) datasets [26] as the knowledge base. In particular, UDB is a diverse and representative collection compiled from 428 university textbooks and spanning 18 subjects. Each sub-corpus contains between 600,000 and 5 million tokens. We use the following 3 sub-corpus.

- *Mix:* A compact and diverse sub-corpus of UDB. It is the smallest among our three datasets with 5.74 MB.
- *Agriculture:* It focuses on agricultural contents and covers topics like beekeeping and urban farming. Its size is 70 MB.
- *Music:* It encompasses a variety of subjects related to music (such as musical styles, biographies, etc). It is the largest among the three datasets, totaling 143 MB.



Figure 9: Sanity check. We compare the answers generated by the same GraphRAG method (i.e., LightRAG). The current evaluation (left plot) leads to wrong conclusion, i.e., the method outperforms itself, while our evaluation (right plot) correctly determines that the method matches itself.

These datasets are also used in the experiments of LightRAG [10]. Additionally, we utilize the *GPT-4o-mini* model as the LLM to generate answers as well as serve as the evaluation judge. We generate 50 distinct questions for each of the 3 levels, i.e., node, edge, and subgraph, using our graph-text-grounded question generation method. As such, the total number of questions is 150.

**Sanity check.** Before commencing the experiments, we check whether our evaluation framework removes all the biases. In Figure 9, we create two instances of LightRAG (denoted as *A* and *B*) and compare them as different GraphRAG methods. Intuitively, a reasonable evaluation should declare that the two match each other in answer quality as they are essentially the same GraphRAG method. However, the current evaluation shows that *A* has a win rate of 90%, while *B* has a win rate of only 10%, translating into an advantage of 80% and suggesting significant performance gains.
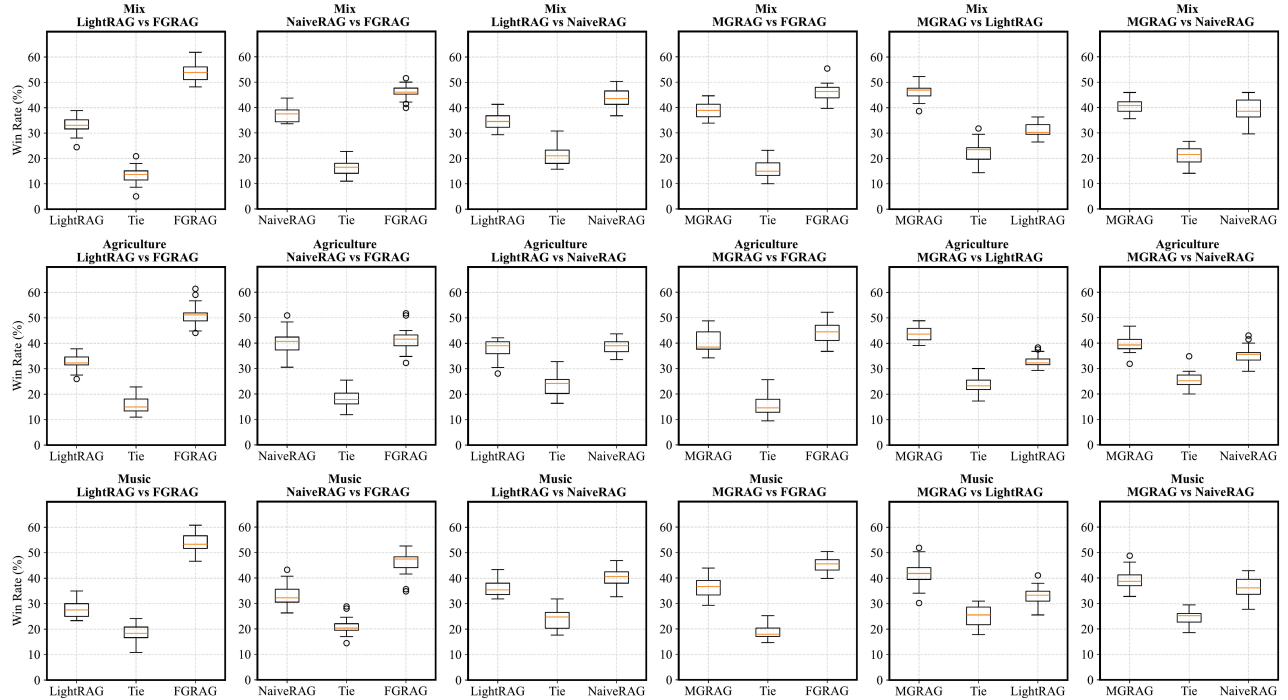
**Figure 10: One-to-one comparison between the 3 GraphRAG methods and NaiveRAG, each row corresponds to a dataset.**

This is because the answers of *A* are placed in the front (i.e., position bias) and only a single trial is conducted (i.e., trail bias). In contrast, our evaluation correctly determines that the two match each other. In Appendix D, we provide a case study to show how our evaluation framework assesses answer quality.

## 5.2 Results and Analysis

**Overall performance.** We perform a one-by-one comparison for the 3 GraphRAG methods and NaiveRAG and report the results in Figure 10. We can make the following observations.

First, the win rates are generally smaller than reported by previous work, and the tie rates are notable, making the performance gaps between the methods small. For instance, LightRAG [10] reports that it achieves win rates of 66.70% (vs. NaiveRAG) and 56.38% (vs. MGRAG) on the Agriculture dataset, while in Figure 10, the win rates of LightRAG become 39.06% (vs. NaiveRAG) and 32.33% (vs. MGRAG) on the same dataset. Moreover, the tie rates are over 20% in most of the cases. These results are because our evaluation framework resolves the biases and allows ties. Moreover, our evaluation also changes the relative performance of the methods. For example, LightRAG is reported to outperform NaiveRAG [10] but our results show that NaiveRAG performs better than LightRAG.

Second, FGRAG performs the best among the methods, followed by MGRAG, then NaiveRAG, and finally LightRAG. This is because FGRAG assigns more concise and informative contexts to the entities and relations in the knowledge graphs and adopts the PageRank scores to rank the retrieved entities and relations. Microsoft also shows that FGRAG can accurately retrieve information relevant
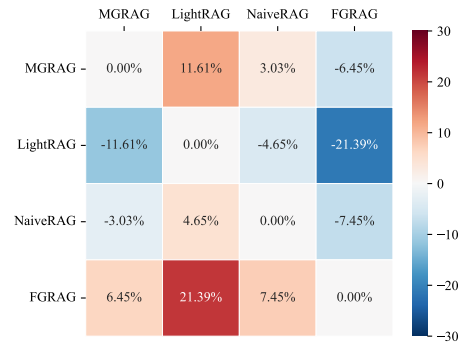


**Figure 11: Relative win rates (defined as row minus column) between the 3 GraphRAG methods and NaiveRAG. A positive value means that the Row method outperforms the Column method. The results are averaged over 3 the datasets.**

to the question [6]. The results suggest that a high-quality knowledge graph and proper re-ranking for the retrieved contents are beneficial for performance.

Third, the performance advantages of the GraphRAG methods over NaiveRAG generally improve with dataset size. This is because GraphRAG methods use knowledge graphs for context retrieval while NaiveRAG directly retrieves text chunks. The knowledge graphs help to accurately retrieve related information, and such ability becomes more important for large datasets.
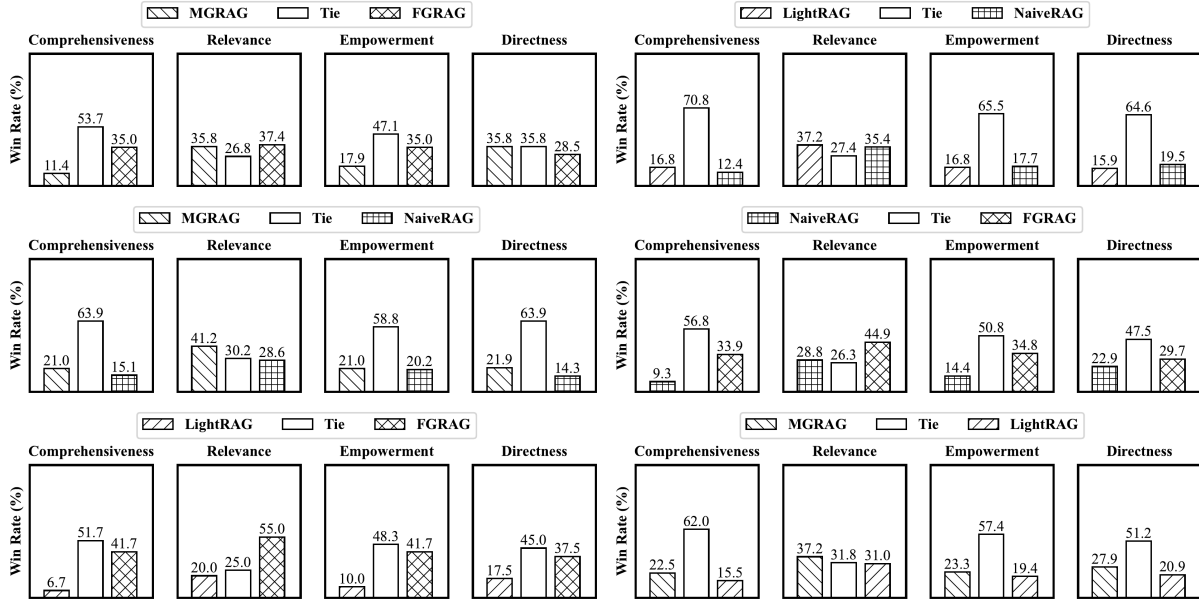
**Figure 12: Detailed win rate comparison for the 4 aspects of answer quality, the dataset is Music.**

**Relative win rate.** To better visualize the performance differences between the methods, we introduce *relative win rate*, which is defined as follows:

$$\text{Relative Win Rate} = \frac{A_{\text{win}} - B_{\text{win}}}{A_{\text{win}} + B_{\text{win}} + \text{Tie}},$$

where $A_{\text{win}}$ is the number of questions method $A$ wins, and similarly for $B_{\text{win}}$ and Tie. Intuitively, the relative win rate measures the advantage of one method (i.e., $A$) over another (i.e., $B$).

Figure 11 plots the relative win rate between the methods. We observe again that the performance gaps between the methods are moderate. Specifically, with the exception of FGRAG and MGRAG, which improve LightRAG by 21% and 11%, respectively, the relative win rates between the other methods are all below 8%. We provide the relative win rate for each individual dataset in Appendix B, and the observations are similar.

**Detailed performance.** In Figure 12, we look into how the methods perform in the detailed aspects of answer quality (i.e., comprehensiveness, relevance, empowerment, and directness). We only report the results on the Music dataset due to the page limit, and the results on the other datasets can be found in Appendix C. We make the following observations.

First, the tie rates are much higher than the overall answer quality results in Figure 10. In particular, the tie rates are above 50% for many cases and seldom below 25%. This suggests that the answers generated by different methods are likely to have similar quality in individual evaluation aspects, and as a consequence, even if our evaluation decides that a method wins another for a question, the actual quality gap between their answers can be small.

Second, the performance of the methods for the detailed aspects reflects their designs. For instance, when comparing FGRAG and MGRAG, FGRAG excels in comprehensiveness, while MGRAG leads

**Table 2: Average token consumption and query time**

| Method name | Token consumption | Query time (s) |
|---|---|---|
| MGRAG | 326,200 | 25.84 |
| LightRAG | 19,650 | 18.27 |
| NaiveRAG | 4,213 | 11.55 |
| **FGRAG** | **3,310** | **8.77** |

in directness. This is because FGRAG retrieves broader knowledge, which enhances comprehensiveness, while MGRAG directly leverages the relevant summaries for more targeted answers. Similar observations can also be made for LightRAG and NaiveRAG.

In Appendix E, we also explore how the methods perform for different question types, i.e., node, edge, and subgraph. The results show that MGRAG excels for node questions as it generates highly high-quality community summaries around nodes, while FGRAG performs the best for subgraph questions due to its comprehensive retrieval and re-ranking.

**Execution costs.** We also test the methods for their average number of consumed tokens and query time for answering a question, and Table 2 reports the results. The results are averaged over 100 questions based on the Mix dataset, and the query time is the duration from question submission to answer generation, which includes API request time. The results show that MGRAG consumes the most tokens and the longest time. This is because it retrieves and sorts the community summaries, which require repeated LLM calls. FGRAG is the most efficient, with minimal token consumption and the shortest query time. This benefits from its concise knowledge format. LightRAG ranks second, while NaiveRAG performs nearly as well as FGRAG.

## 6 Conclusions and Future Directions

We find that the current evaluation framework for the answer quality of GraphRAG methods asks unrelated questions and suffers from evaluation biases. To tackle the two critical limitations, we propose graph-text-grounded question generation to produce questions that are closely related to the fine-grained details of the dataset, and an unbiased evaluation procedure to resolve the biases in answer position, length, and evaluation trials. We apply our new evaluation framework to compare 3 representative GraphRAG methods and find that their performance gains are much more moderate than reported previously. Besides overall performance, we also look into detailed quality aspects, question types, and execution costs.

Our work calls for exploring automatic answer quality evaluation, which is crucial for laying solid foundations for RAG and GraphRAG research. We think two promising directions are: (i) comparing multiple methods in one pass, as the results of the current 1-to-1 comparison can be huge and difficult to comprehend; (ii) instead of simply declaring a winner, quantifying how quantitatively how far two answers differ in quality. We suggest that further exploration along these lines may be beneficial for the community.

# References

[1] Bernd Bohnet, Vinh Q Tran, Pat Verga, Roee Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, et al. 2022. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037* (2022).

[2] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology* 15, 3 (2024), 1–45.

[3] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 17754–17762.

[4] Xiaojun Chen, Shengbin Jia, and Yang Xiang. 2020. A review: Knowledge reasoning over knowledge graph. *Expert systems with applications* 141 (2020), 112948.

[5] Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937* (2023).

[6] CircleMind. 2024. Streamlined and promptable fast graphrag framework designed for interpretable, high-precision, agent-driven retrieval workflows. https://github.com/circlemind-ai/fast-graphrag Accessed: 2025-02-07.

[7] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130* (2024).

[8] Dieter Fensel, Umutcan Şimşek, Kevin Angele, Elwin Huaman, Elias Kärle, Oleksandra Panasiuk, Ioan Toma, Jürgen Umbrich, Alexander Wahler, Dieter Fensel, et al. 2020. Introduction: what is a knowledge graph? *Knowledge graphs: Methodology, tools and selected use cases* (2020), 1–10.

[9] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* (2023).

[10] Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. 2024. LightRAG: Simple and Fast Retrieval-Augmented Generation. (2024). arXiv:2410.05779 [cs.IR]

[11] Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. 2023. Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. *Authorea Preprints* (2023).

[12] Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. 2023. A survey on large language models: Applications, challenges, limitations, and practical usage. *Authorea Preprints* (2023).

[13] Haoyu Han, Yu Wang, Harry Shomer, Kai Guo, Jiayuan Ding, Yongjia Lei, Mahantesh Halappanavar, Ryan A Rossi, Subhabrata Mukherjee, Xianfeng Tang, et al. 2024. Retrieval-augmented generation with graphs (graphrag). *arXiv preprint arXiv:2501.00309* (2024).

[14] Rujun Han, Yuhao Zhang, Peng Qi, Yumo Xu, Jenyuan Wang, Lan Liu, William Yang Wang, Bonan Min, and Vittorio Castelli. 2024. Rag-qa arena: Evaluating domain robustness for long-form retrieval augmented question answering. *arXiv preprint arXiv:2407.13998* (2024).

[15] Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *arXiv preprint arXiv:2402.07630* (2024).

[16] Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. 2024. GRAG: Graph Retrieval-Augmented Generation. *arXiv preprint arXiv:2405.16506* (2024).

[17] Yucheng Hu and Yuxing Lu. 2024. Rag and rau: A survey on retrieval-augmented language model in natural language processing. *arXiv preprint arXiv:2404.19543* (2024).

[18] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232* (2023).

[19] Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng Tang, Suhang Wang, Yu Meng, et al. 2024. Graph chain-of-thought: Augmenting large language models by reasoning on graphs. *arXiv preprint arXiv:2404.07103* (2024).

[20] Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. *arXiv preprint arXiv:2302.14520* (2023).

[21] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.

[22] Rong-Hua Li, Jeffrey Xu Yu, Lu Qin, Rui Mao, and Tan Jin. 2015. On random walk based graph sampling. In *2015 IEEE 31st international conference on data engineering*. IEEE, 927–938.

[23] Yen-Ting Lin and Yun-Nung Chen. 2023. Llm-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models. *arXiv preprint arXiv:2305.13711* (2023).

[24] Yuanjie Lyu, Zhiyu Li, Simin Niu, Feiyu Xiong, Bo Tang, Wenjin Wang, Hao Wu, Huanyong Liu, Tong Xu, and Enhong Chen. 2024. Crud-rag: A comprehensive chinese benchmark for retrieval-augmented generation of large language models. *ACM Transactions on Information Systems* (2024).

[25] Tyler Thomas Procko and Omar Ochoa. 2024. Graph retrieval-augmented generation for large language models: A survey. In *2024 Conference on AI, Science, Engineering, and Technology (AIxSET)*. IEEE, 166–169.

[26] Hongjin Qian, Peitian Zhang, Zheng Liu, Kelong Mao, and Zhicheng Dou. 2024. MemoRAG: Moving towards Next-Gen RAG Via Memory-Inspired Knowledge Discovery. (2024). arXiv:2409.05591 https://arxiv.org/abs/2409.05591

[27] Alireza Salemi and Hamed Zamani. 2024. Evaluating retrieval quality in retrieval-augmented generation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2395–2400.

[28] SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv preprint arXiv:2401.01313* (2024).

[29] Yu-Min Tseng, Wei-Lin Chen, Chung-Chi Chen, and Hsin-Hsi Chen. 2024. Are Expert-Level Language Models Expert-Level Annotators? *arXiv preprint arXiv:2410.03254* (2024).

[30] Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al. 2023. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. *arXiv preprint arXiv:2310.07521* (2023).

[31] Yang Wang, Alberto Garcia Hernandez, Roman Kyslyi, and Nicholas Kersting. 2024. Evaluating quality of answers for Retrieval-Augmented Generation: A strong LLM is all you need. *arXiv preprint arXiv:2406.18064* (2024).

[32] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.

[33] David F Williamson, Robert A Parker, and Juliette S Kendrick. 1989. The box plot: a simple visual method to interpret data. *Annals of internal medicine* 110, 11 (1989), 916–921.

[34] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2023. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453* (2023).

[35] Hao Yu, Aoran Gan, Kai Zhang, Shiwei Tong, Qi Liu, and Zhaofeng Liu. 2024. Evaluation of retrieval-augmented generation: A survey. In *CCF Conference on Big Data*. Springer, 102–120.

[36] Haozhen Zhang, Tao Feng, and Jiaxuan You. 2024. Graph of records: Boosting retrieval augmented generation for long-context summarization with graphs. *arXiv preprint arXiv:2410.11001* (2024).

[37] Yujia Zhou, Yan Liu, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Zheng Liu, Chaozhuo Li, Zhicheng Dou, Tsung-Yi Ho, and Philip S Yu. 2024. Trustworthiness in retrieval-augmented generation systems: A survey. *arXiv preprint arXiv:2409.10102* (2024).

[38] Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zheng Liu, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107* (2023).

## A Graph RAG Introduction

The Retrieval-Augmented Generation (RAG) system enhances large language models (LLMs) by integrating external knowledge bases, enabling real-time access to up-to-date information without re-traininging. This approach addresses LLMs' limitations in handling new or domain-specific knowledge while reducing training costs and mitigating the "illusion" problem, where models generate factually incorrect content. Traditional Naïve RAG achieves this by converting external data into text chunks, vectorizing them, and retrieving the most relevant segments for question answering. However, real-world knowledge often exists in graph structures, where entities (nodes) and their relationships (edges) are key. This limitation has driven the evolution from text-based RAG to Graph RAG, which leverages graph-structured knowledge for more accurate and contextually rich retrieval and generation.
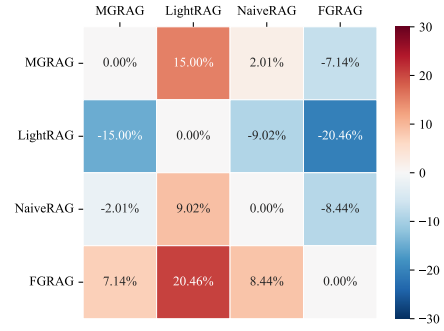
Graph RAG then proposes its Graph RAG solution, which further performs entity extraction and relationship extraction operations on the text chunk partitioned by the traditional RAG, and generates entity- and relationship-specific summaries as attributes after extraction. After constructing a graph knowledge structure consisting of nodes and edges, Graph RAG further runs community discovery algorithms on the graph structure form summaries for the communities, and reorders the summaries of the communities as auxiliary contexts when answering questions. Such an approach, when applied in practice, creates the problem of long data preparation time as well as limited scalability because of the need for community discovery when constructing the graph structure.

Light RAG proposes a hierarchical context retrieval approach, specifically, after constructing the knowledge of the graph structure, Light RAG does not further operate against the communities, but rather, from the user's question instead, Light RAG introduces a dual-level retrieval paradigm that extracts entity-related keywords and relationship-related keywords from the user's problem, and directly retrieves the keywords. Entity keywords can be used to solve more specific problems, while relationship keywords can be used to solve more macroscopic conceptual problems, and Light RAG integrates and rearranges the node and edge information retrieved from the two hierarchical levels to form the final context.

Fast Graph RAG is an efficient graph retrieval framework optimized for large-scale application scenarios, which is designed to improve the efficiency and accuracy of knowledge retrieval through efficient graph construction operations and optimized retrieval strategies. The framework introduces the PageRank algorithm to reorder the retrieved nodes and edges to further enhance the relevance and reliability of the retrieval results. In addition, Fast Graph RAG provides a graphical knowledge view that allows users to observe and navigate the retrieval process intuitively, which enhances the interpretability and debuggability of the system.
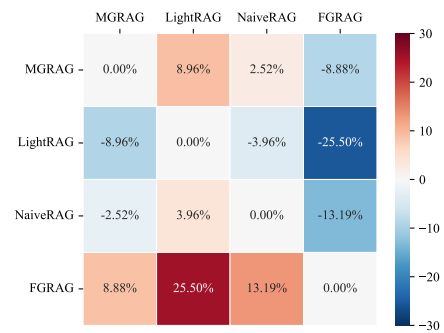
## B Heat Maps

This section shows the specific relative win rates on three datasets, and it can be seen that the distribution of relative win rates for the three plots is generally consistent with the mean plots shown in the paper, i.e., the relative win rates for the vast majority of comparisons are not particularly high, but there are a couple of exceptional values (e.g., FGRAG vs. LightRAG)



Figure 13: The performance comparison heatmap shows the relative win rate between model pairs over Mix dataset. A positive value indicates that the Row model outperforms the Column model and vice versa.



Figure 14: The performance comparison heatmap shows the relative win rate between model pairs over Agriculture dataset. A positive value indicates that the Row model outperforms the Column model and vice versa.



Figure 15: The performance comparison heatmap shows the relative win rate between model pairs over Music dataset. A positive value indicates that the Row model outperforms the Column model and vice versa.

## C Detailed Performance on Other Datasets

This section presents comparisons of specific aspects of the two datasets not included in the main text. The observations align with
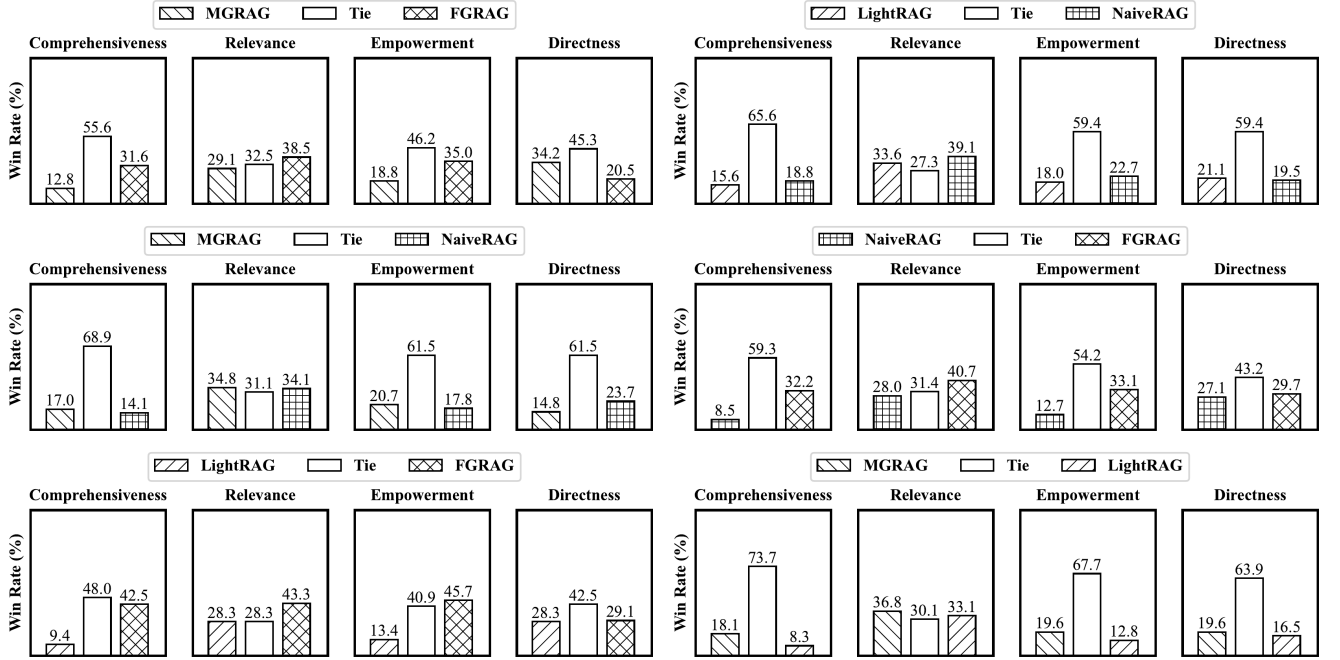
**Figure 16: Detailed win rate comparison across the four introduced aspects for model pairs under Agriculture dataset.**
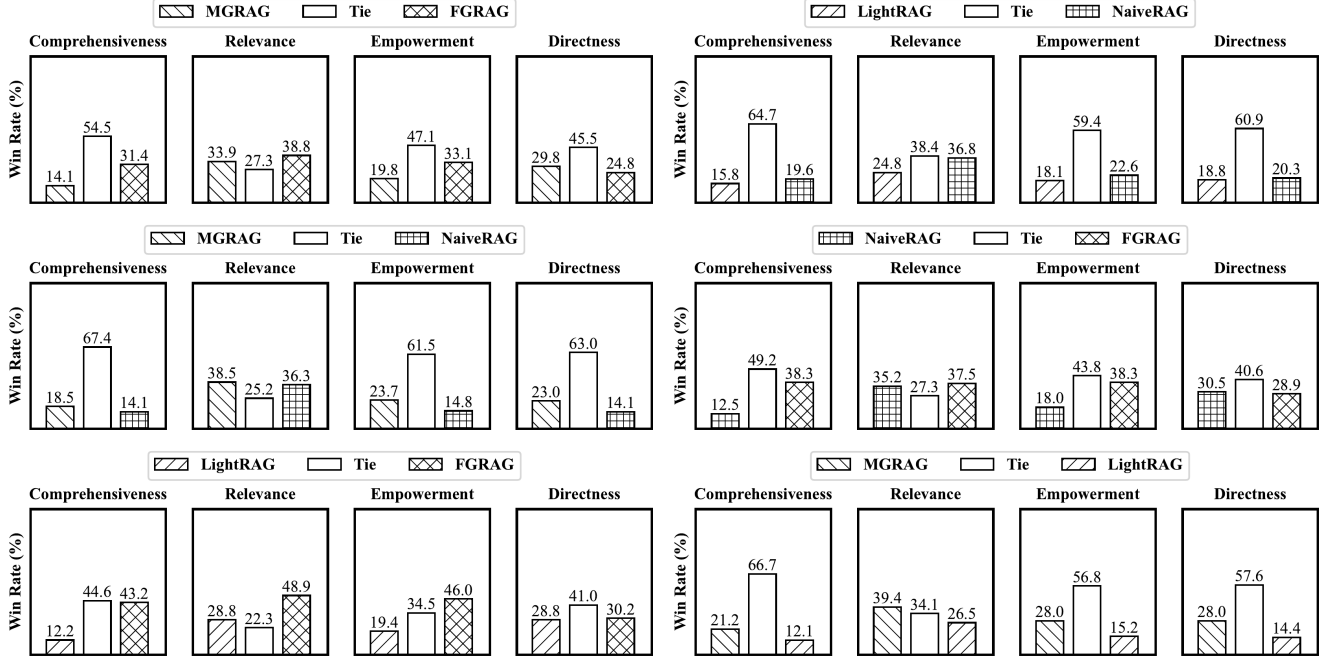


**Figure 17: Detailed win rate comparison across the four introduced aspects for model pairs under Mix dataset.**

those in the main text, but the comparisons between LightRAG and NaiveRAG do not exhibit the same level of consistency as those between MGRAG and FGRAG. The comparative results between MGRAG and FGRAG remain valid, and many instances of ties can still be observed. We speculate that this discrepancy is primarily influenced by the size of the dataset. When the dataset is significantly larger, the superior comprehensiveness of LightRAG's answers becomes more apparent.

## D  Case Study

In this section, we provide an evaluation case to prove that our evaluation metric is efficient. The results are present in Table 3

**Table 3: Case Study**

| Question | What school did Harry Potter attend? |
|---|---|
| **Answer 1** | # The School Harry Potter Attends:Harry Potter attends **Hogwarts School of Witchcraft and Wizardry**. This institution, often simply called **Hogwarts**, serves as the central educational setting within the Harry Potter series. ## Significance of Hogwarts Hogwarts is not just a school; it is a vital part of the narrative, providing a rich backdrop for magical education and diverse experiences for its students. The institution fosters learning in a variety of magical subjects, which play a crucial role in the development of Harry and his peers throughout the series. In summary, Hogwarts School of Witchcraft and Wizardry is integral to the story, shaping the characters and their journeys within the realm of magic. Harry Potter attends Hogwarts School of Witchcraft and Wizardry, a prestigious magical institution. It is the primary setting for his adventures and self-discovery. |
| **Answer 2** | ## Significance of Hogwarts Hogwarts is known for its strong educational framework. Students are sorted into four houses: Gryffindor, Hufflepuff, Ravenclaw, and Slytherin. The school provides comprehensive magical education and serves as a safe haven in troubled times.## Harry's Experience At Hogwarts, Harry faces challenges, forms lasting friendships with Hermione Granger and Ron Weasley, and participates in magical adventures, including playing Quidditch as a Seeker for Gryffindor. These experiences shape his identity and give him a sense of belonging he lacks with the Dursleys. Overall, Hogwarts is crucial to Harry's growth as a young wizard. |
| **Random Select Judge** | Aspect 1:<br>Explanation: "Both answers provide detailed information about Hogwarts School of Witchcraft and Wizardry, covering its significance in the Harry Potter narrative. However, Answer 2 goes deeper into Harry's experiences at Hogwarts, including details about house sorting and personal challenges, while also addressing the educational framework and community aspects. Answer 1, though comprehensive, lacks the same level of detail concerning Harry's experiences."<br>Answer 1: 4, **Answer 2: 5**<br><br>Aspect 2:<br>Explanation: Both answers directly address the question, identifying Hogwarts as the school Harry attends. They both remain relevant throughout, discussing its significance and how it impacts Harry. However, Answer 2 provides more context and richer details about Hogwarts and its relevance to the character, making it more aligned overall.<br>Answer 1: 4, **Answer 2: 5**<br><br>Aspect 3:<br>Explanation: Both answers help the reader understand the role of Hogwarts in Harry's life, but Answer 2 more effectively empowers the reader by detailing Harry's friendships and challenges, thereby enabling a deeper understanding of how the school influences his journey. Answer 1 offers a good overview but does not provide the same degree of insight into personal growth.<br>Answer 1: 3, **Answer 2: 5**<br><br>Aspect 4:<br>Explanation: Answer 1 is clear and direct in addressing the question but uses slightly more generalized language when discussing the significance of Hogwarts, whereas Answer 2 is exceptionally direct, offering specific details about Harry's school life, thus addressing the question with greater clarity.<br>Answer 1: 4, **Answer 2: 5** |
| **Overall** | Answer 1 Aspect 1: 3.75<br>Answer 1 Aspect 2: 4.25<br>Answer 1 Aspect 3: 3.25<br>Answer 1 Aspect 4: 4.0<br>Answer 1 Total: 15.25<br><br>Answer 2 Aspect 1: 5.0<br>Answer 2 Aspect 2: 4.75<br>Answer 2 Aspect 3: 5.0<br>Answer 2 Aspect 4: 5.0<br>**Answer 2 Total: 19.75** |

As illustrated in the example provided, the evaluation system we proposed assigned scores to both answers to the same question from four distinct aspects. Upon analysis, it is evident that Answer 2 achieved higher scores across all dimensions. This superior performance can be attributed to the fact that Answer 2 not only encompassed the content of Answer 1 but also expanded upon it by incorporating a detailed description of Harry Potter's experiences at Hogwarts School of Witchcraft and Wizardry. This additional information rendered Answer 2 more comprehensive and contextually relevant to the question posed. Consequently, Answer 2 demonstrates greater utility in addressing the query, providing enhanced clarity and depth for the inquirer. This outcome underscores the importance of incorporating supplementary details that align with the thematic focus of the question, thereby elevating the overall quality of the response.
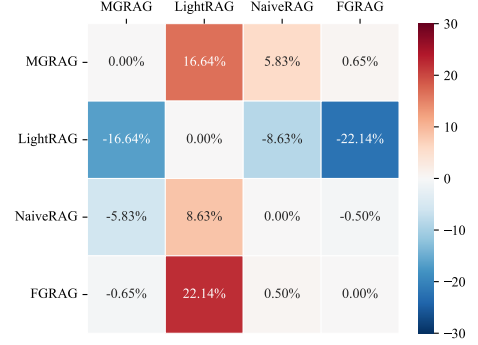
## E  Performance Across Different Question Types

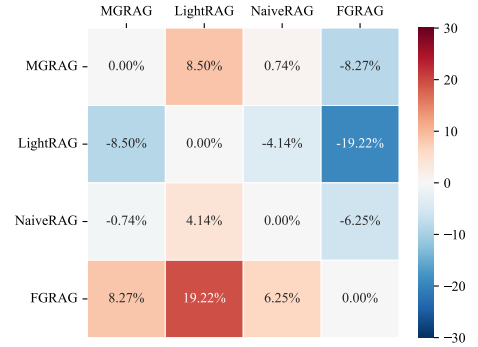In this section, we show the average win rates of different GraphRAG methods under different types of questions.

Our experimental results indicate that MGRAG achieves superior performance on node-level questions, which can be attributed to its ability to generate highly fine-grained community reports with minimal redundancy while maintaining rich information for a limited number of entities. In contrast, the overall performance of FGRAG follows a hierarchical ranking: subgraph level > edge level > node level. This aligns with its original design objective of addressing complex problems that require a comprehensive understanding of the dataset. Upon analyzing the information retrieved by FGRAG, we observed that its node-level summaries are comparatively less detailed, which may explain its relatively weaker performance at this level.

Additionally, Light RAG, a similarly designed Graph RAG approach, demonstrates limited advantages over other methods, despite showing marginal improvements on more complex questions. We believe that this is due to two factors: (1) the inclusion of extensive information from one-hop expansions, which introduces noise into the retrieved data, and (2) the lack of an effective ranking mechanism, such as the PageRank algorithm utilized by FGRAG, to prioritize relevant information. These observations provide insights into the strengths and limitations of different Graph RAG methodologies.
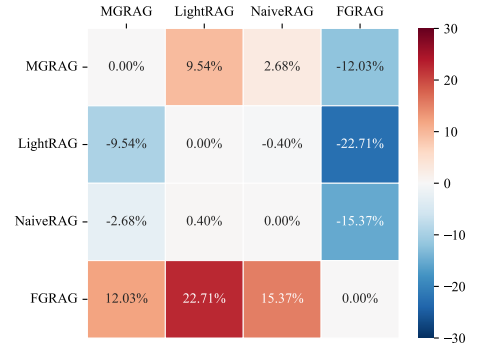
The overall performance of NaiveRAG exhibits an inverse trend compared to Graph RAG approaches, with significantly lower effectiveness on subgraph-level questions relative to node-level and edge-level questions. This aligns with its design, which relies on direct retrieval of fixed-length text segments for augmentation-a strategy that proves inadequate for addressing the complexity of subgraph-level questions. This limitation underscores the challenges of using simplistic retrieval methods for questions requiring broader contextual understanding.



Figure 18: The average win rate comparison across the four baseline model pairs evaluated on node-level questions under three datasets (Mix, Agriculture, and Music).



Figure 19: The average win rate comparison across the four baseline model pairs evaluated on edge-level questions under three datasets (Mix, Agriculture, and Music).



Figure 20: The average win rate comparison across the four baseline model pairs evaluated on subgraph-level questions under three datasets (Mix, Agriculture, and Music).