# Towards Self-cognitive Exploration: Metacognitive Knowledge Graph Retrieval Augmented Generation

Xujie Yuan
Sun Yat-sen University
Zhuhai, China

Shimin Di
Southeast University
Nanjing, China

Jielong Tang
Sun Yat-sen University
Zhuhai, China

Libin Zheng
Sun Yat-sen University
Zhuhai, China

Jian Yin
Sun Yat-sen University
Zhuhai, China

## ABSTRACT

Knowledge Graph-based Retrieval-Augmented Generation (KG-RAG) significantly enhances the reasoning capabilities of Large Language Models by leveraging structured knowledge. However, existing KG-RAG frameworks typically operate as open-loop systems, suffering from cognitive blindness, an inability to recognize their exploration deficiencies. This leads to relevance drift and incomplete evidence, which existing self-refinement methods, designed for unstructured text-based RAG, cannot effectively resolve due to the path-dependent nature of graph exploration. To address this challenge, We propose Metacognitive Knowledge Graph Retrieval Augmented Generation (**MetaKGRAG**), a novel framework inspired by human metacognition process, which introduces a Perceive-Evaluate-Adjust cycle to enable path-aware, closed-loop refinement. This cycle empowers the system to self-assess exploration quality, identify deficiencies in coverage or relevance, and perform trajectory-connected corrections from precise pivot points. Extensive experiments across five datasets in the medical, legal, and commonsense reasoning domains demonstrate that MetaKGRAG consistently outperforms strong KG-RAG and self-refinement baselines. Our results validates the superiority of our approach and highlights the critical need for path-aware refinement in structured knowledge retrieval.

## KEYWORDS

Large Language Models, Knowledge Graph, Retrieval-Augmented Generation, Metacognition

## 1 INTRODUCTION

Large Language Models (LLMs) have demonstrated remarkable capabilities [2, 28], yet their reliability is often limited by hallucinations and outdated internal knowledge [18, 20, 36]. Retrieval-Augmented Generation (RAG) mitigates these issues by grounding LLMs in external knowledge [6, 8]. While standard RAG uses unstructured text, Knowledge Graph RAG (KG-RAG) [5, 7, 33] leverages explicit, structured relationships. This enables a more verifiable reasoning process, making it particularly powerful for complex queries that require connecting multiple pieces of information (i.e., multi-hop reasoning) to deliver precise answers [11, 14, 16, 21].

Despite these advantages, the effectiveness of KG-RAG is often undermined by a challenge inherent in its exploration process. Current KG-RAG methods typically operate as **open-loop systems**, generating an evidence path in a single linear pass without a feedback mechanism [24, 30]. This design flaw leads to what we term
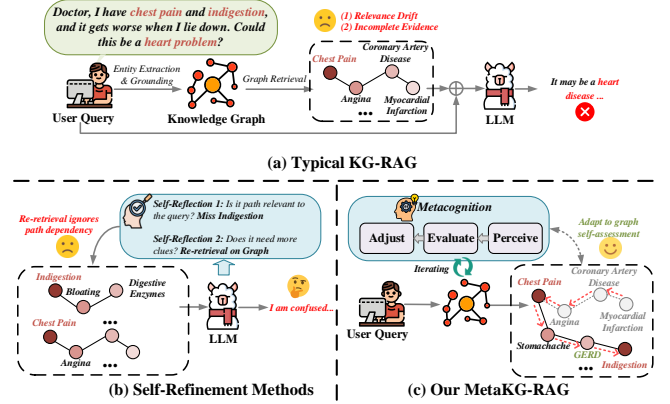


Figure 1: The comparisons of KG-RAG pipeline, Self-Refinement methods and our MetaKGRAG. (a) Typical KG-RAG suffers from cognitive blindness issues, leading to relevance drift and incomplete evidence. (b) Self-Refinement struggles to adapt to KG-RAG due to overlooking path dependency of graph exploration. (c) Our MetaKGRAG achieves graph-based self-cognition through a metacognitive cycle.

**Cognitive Blindness**, a state where the system is unaware of its own exploration deficiencies. As shown in Fig. 1 (a), a real-world medical query: *"Doctor, I have chest pain and indigestion, and it gets worse when I lie down. Could this be a heart problem?"*. Traditional KG-RAG methods typically employ a greedy search based on local similarity. In this case, "chest pain" has a strong semantic link to heart disease, biasing the system towards this path. As the exploration deepens from "chest pain" to "Angina" and then "Coronary Artery Disease", it progressively diverges from other crucial evidence paths, leading to **relevance drift**. As the system delves into the cardiac subgraph, this drift causes it to overlook key concepts such as "indigestion" and "worse when lying down", resulting in **incomplete evidence**. Ultimately, an incorrect answer (heart disease) is generated. Due to the absence of an effective feedback mechanism, current KG-RAG methods struggle to recognize these drifts and omissions, thereby impacting their overall performance.

To address the cognitive blindness issue, a natural inclination is to apply self-refinement mechanism of current text-based RAG methods [10, 32]. They assess discrete units of evidence (e.g., text chunks), if a flaw is found, they can substitute the faulty piece or trigger a second-time searching to fetch a new evidence. However,

this paradigm is inappropriate for adapting to KG-RAG due to the **path-dependent** nature of graph exploration. The core limitation is that existing self-refinement methods treat evidence paths as a set of independent items, failing to grasp the relational trajectories between them. *This is akin to realizing you are on the wrong highway, the solution is not just finding the correct road segment, but also identifying where the wrong turn occurred and re-planning a new route from your current location to the destination.* As illustrated in Fig. 1 (b), self-refinement methods recognize the missing concept "indigestion" and perform a re-retrieval. However, due to the absence of trajectories connecting the original "chest pain" path to the new "indigestion" path, LLMs still struggle to integrate these two separate evidence paths for the final diagnosis: *heart problem* or *digestive issue*? These potential knowledge conflicts across independent evidence paths induce hallucinations in LLMs. Thus, directly employing self-refinement to address the cognitive blindness issue in KG-RAG remains a challenge.

Inspired by human metacognitive processes [12, 22], the ability to "think about thinking", we propose Metacognitive Knowledge Graph Retrieval-Augmented Generation (**MetaKGRAG**). This framework introduces a metacognitive cycle tailored for the path-dependent nature of graph exploration. MetaKGRAG incorporates self-monitoring and self-regulation mechanisms, which are lacking in open-loop KG-RAG systems. It transforms blind path generation into a reflective and closed-loop process through a concrete Perceive-Evaluate-Adjust cycle (as shown in Fig. 1 (c)). Here is how this cycle directly addresses cognitive blindness:

- **Perceive:** First, our method generates an initial candidate path and then holistically assesses it. It systematically checks how well the entire path, as a whole, covers all crucial aspects of the input query, creating a comprehensive initial understanding.
- **Evaluate:** Based on the Perceive, it then diagnoses specific, pre-defined problems. To combat *Incomplete Evidence*, it identifies which key concepts from the query were missed. To correct *Relevance Drift*, it pinpoints the exact node where the exploration began to deviate from the query's overall intent.
- **Adjust:** Different from existing self-refinement methods, the adjustment is not a blind re-retrieval. Based on the specific diagnosis, the system performs a **trajectory-connected correction**. It identifies the optimal pivot point on the flawed path, the last correct step before things went wrong. From there, it initiates a smarter re-exploration, now armed with the knowledge of what to avoid and what to prioritize, effectively re-routing its trajectory to the correct destination.

This cycle equips the system with the self-cognition to master path-dependent retrieval, providing a solution adapted to knowledge graphs. The main contributions of this paper are as follows:

- We identify Cognitive Blindness, along with its manifestations of incomplete coverage and relevance drift, as a core challenge in current KG-RAG methods.
- We propose MetaKGRAG, a novel framework inspired by metacognitive principles, that introduces an evidence-level refinement cycle to overcome the limitations of traditional correction in path-dependent graph exploration.
- We design a concrete three-stage `Perceive-Evaluate-Adjust` cycle that enables iterative assessment, deficiency diagnosis, and strategic re-exploration of candidate evidence paths.

- We conduct comprehensive experiments across medical (ExplainCPE [13], CMB-Exam [26], webMedQA [9]), legal (JEC-QA [35], and commonsense (CommonsenseQA [25]) domains demonstrate that MetaKGRAG achieves substantial improvements over strong LLM baselines, KG-RAG methods, and incorporated self-refinement approaches, validating the superiority of our adaptive control framework.

## 2 INSPIRATION FROM METACOGNITION

Metacognition [12, 22], inspired by human "thinking about thinking" processes, has shown promise in improving system self-awareness and adaptive control. In the context of LLMs, this has led to several promising approaches. Metacognitive prompting [27] guides LLMs through explicit self-reflection steps, asking models to evaluate their own reasoning quality and identify potential errors. A metacognition framework used in text-based RAG is MetaRAG [37], which employs a monitor-evaluate-plan loop to diagnose and rectify specific failures. It assesses an initial response by analyzing the sufficiency and consistency of both its internal knowledge and externally retrieved documents. Based on this diagnosis, it plans a targeted revision, such as generating new search queries to fill knowledge gaps or ignoring distracting evidence.

While powerful, these methods reveal a critical limitation for graph exploration that their refinement loops operate are not at the path level. For instance, MetaRAG's "plan" step can decide to trigger a completely new search with a revised query, but it lacks the fine-grained control to intervene within a single, ongoing graph traversal. It cannot perform fine-grained interventions, such as backtracking to a specific node on a flawed path and re-exploring from that point. Its mechanism is about replacing or adding entire evidence, not correcting a trajectory.

This highlights an adaptation gap that existing metacognitive frameworks lack the graph-native ability to perform fine-grained diagnosis and correction. They cannot answer the crucial questions, "*Which specific step in my current path was the wrong turn?*" and "*How can I correct my trajectory from that point?*". This gap motivates our development of MetaKGRAG, which implements metacognitive cycles specifically for the unique challenges of structured knowledge retrieval, enabling path-aware refinement.

## 3 METAKGRAG FRAMEWORK

Given a knowledge graph $\mathcal{G} = (\mathcal{E}, \mathcal{R})$ and a natural language question $Q$, the KG-RAG task is to explore $\mathcal{G}$ to retrieve a relevant evidence subgraph $\mathcal{S}$ that provides sufficient context for an LLM to generate an accurate answer. The core of this process is the generation of one or more evidence paths $P = \{(h_1, r_1, t_1), \dots, (h_n, r_n, t_n)\}$, where each triplet $(h, r, t)$ represents a head entity, relation, and tail entity respectively in the KG.

### 3.1 Framework Overview

To simultaneously address the issues of relevance drift and incomplete evidence in a way that respects the core challenge of path dependency that simple self-refinement methods cannot solve, we designed the MetaKGRAG framework. This transforms the path generation process from a blind execution into a reflective one by operationalizing metacognitive principles into an evidence-level
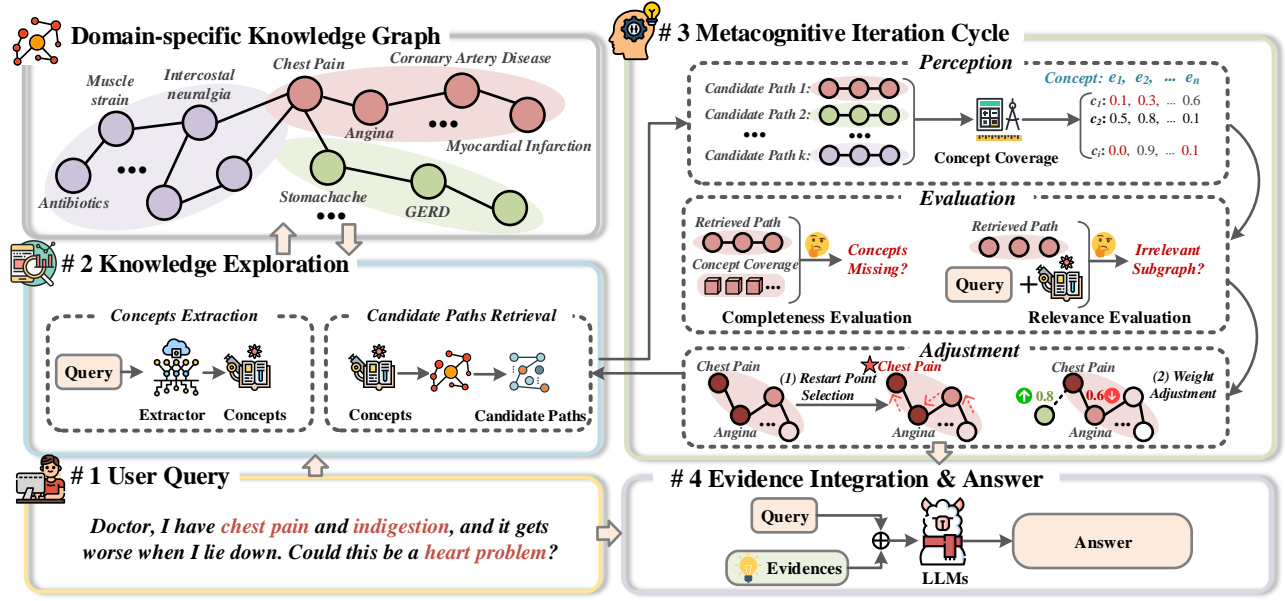
**Figure 2: An overview of our MetaKGRAG framework. It iteratively refines graph exploration via a path-aware Perceive-Evaluate-Adjust cycle to address cognitive blindness in relevance drift and incomplete evidence.**

refinement cycle. At its core, MetaKGRAG is driven by a **Perceive-Evaluate-Adjust** cycle. This cycle is an iterative process that assesses a fully generated candidate path to inform a more strategic subsequent search.

- **Perceive: What is the quality of my current path?** After generating an initial candidate path, the framework first Perceives its overall quality. This involves a holistic assessment of the entire path, checking how well its constituent entities cover the key concepts of the query. This step provides the raw self-awareness that traditional open-loop KG-RAG methods lack.
- **Evaluate: What are the specific problems?** Based on the perceived metrics, the framework Evaluates the path to diagnose specific, pre-defined issues. It explicitly checks if the path suffers from a *Completeness Deficiency* (i.e., key concepts are missed) or a *Relevance Deficiency* (i.e., the path has drifted into an irrelevant subgraph). This turns a vague sense of "low quality" into an actionable diagnosis.
- **Adjust: How can I correct the trajectory?** In stark contrast to the disconnected re-retrievals of other methods, the adjustment is a trajectory-connected correction. Based on the diagnosis, the system performs a Strategic Re-exploration. This involves identifying the optimal pivot point on the flawed path, formulating a new search strategy, and initiating a new, smarter exploration from that point to effectively re-route the trajectory.

This iterative cycle operationalizes the principle of self-correction. By assessing a fully generated candidate path and using that diagnosis to inform a subsequent, more targeted search, the framework systematically improves the quality of evidence before it is passed to the LLM. The overall algorithm is detailed in Algorithm 1.

---

**Algorithm 1** Metacognitive Knowledge Graph RAG

**Input:** Question $Q$, Knowledge Graph $\mathcal{G} = (\mathcal{E}, \mathcal{R})$, Parameters $\Theta$
**Output:** Answer $A$

1: $C \leftarrow$ ExtractConcepts($Q, \Theta$.max_concepts)
2: $E_0 \leftarrow$ MatchEntities($C, \mathcal{E}, \tau_{entity}$)
3: $\mathcal{P}_{all} \leftarrow \emptyset$         ▷ Collection of all refined paths
4: **for** each entity $e \in E_0$ **do**
5:     $P_{candidate} \leftarrow$ InitialPathSearch($e, Q, \mathcal{G}$)
6:     **for** $i = 1$ to $N_{max}$ **do**
7:         $\mathcal{M} \leftarrow$ Perceive($P_{candidate}, C$)   ▷ Assess overall quality
8:         $\mathcal{I} \leftarrow$ Evaluate($\mathcal{M}, Q, C, \Theta$)   ▷ Diagnose deficiencies
9:         **if** $\mathcal{I} = \emptyset$ **then**
10:             **break**               ▷ No issues found
11:         **end if**
12:         $P_{new} \leftarrow$ Adjust($P_{candidate}, \mathcal{I}, Q, \mathcal{G}$)   ▷ Initiate a new search to generate a refined path
13:         **if** PathSimilarity($P_{new}, P_{candidate}$) $> \tau_{similarity}$ **then**
14:             **break**       ▷ No significant improvement
15:         **end if**
16:         $P_{candidate} \leftarrow P_{new}$
17:     **end for**
18:     **if** $P_{candidate} \neq \emptyset$ **then**
19:         $\mathcal{P}_{all} \leftarrow \mathcal{P}_{all} \cup \{P_{candidate}\}$
20:     **end if**
21: **end for**
22: $\mathcal{S} \leftarrow$ IntegrateEvidencePaths($\mathcal{P}_{all}$)
23: $A \leftarrow$ GenerateAnswer($Q, \mathcal{S}$)
24: **return** $A$

## 3.2 Initial Knowledge Exploration

Before entering the metacognitive cycle, MetaKGRAG establishes multiple starting points to ensure comprehensive exploration. Given question $Q$, we use an LLM to extract key concepts $C = \{c_1, c_2, ..., c_k\}$ that represent the essential information needs. This extraction is domain-aware, prioritizing technical terms and named entities relevant to the task. For each concept $c_i$, we identify corresponding entities in $\mathcal{G}$ through semantic matching. Specifically, we compute the semantic similarity:

$$\text{sim}(c_i, e_j) = \frac{\mathbf{v}_{c_i} \cdot \mathbf{v}_{e_j}}{||\mathbf{v}_{c_i}|| \cdot ||\mathbf{v}_{e_j}||}$$

where $\mathbf{v}_{c_i}$ and $\mathbf{v}_{e_j}$ are the embedding vectors of concept $c_i$ and entity $e_j$ respectively. To ensure quality, we only select the most semantically similar entities for each concept, subject to a minimum similarity threshold:

$$E_0 = \{\arg\max_{e \in \mathcal{E}} \text{sim}(c_i, e) : c_i \in C\}$$

where we only retain entities with $\text{sim}(c_i, e) > \tau_{entity}$ to ensure high-quality matches, and $\tau_{entity}$ is the entity matching threshold. This results in a set of high-confidence starting points $E_0 = \{e_1, e_2, ..., e_m\}$. Each starting entity initiates an independent retrieval path. This multi-perspective approach ensures that different aspects of the query are explored, reducing the risk of missing crucial information due to a single suboptimal starting point. For each starting entity $e \in E_0$, we first generate an initial candidate path $P_{candidate}$ through greedy exploration from $e$, following edges with the highest relevance to query $Q$.

## 3.3 Metacognitive Iteration Cycle

This iterative cycle is the core of our framework. It takes a candidate evidence path and progressively refines it by diagnosing and correcting its flaws. The cycle consists of three distinct modules: Perception, Evaluation, and Adjustment.

- **Perception Module.** The Perception module acts as a quality sensor for the initial candidate path. Let $E_P = \{e : e \in \text{head or tail of any triple in } P_{candidate}\}$ be the set of entities in the path. It computes a coverage map $\mathcal{M}$ by assessing how well the path covers each key concept $c_i \in C$:

$$\text{Coverage}(c_i) = \max_{e \in E_P} \text{sim}(c_i, e)$$

This provides the raw data for self-assessment.

- **Evaluation Module.** Based on the output from the Perception module, the Evaluation module acts as a diagnostic engine. It formalizes and diagnoses the two key deficiencies discussed in the introduction. A path $P$ may suffer from:
  - **Completeness Deficiency,** if it fails to cover all key aspects. Formally, a path $P$ is deficient in completeness if there exists a concept $c_i \in C$ for which the coverage, defined as $\max_{e \in \text{entities}(P)} \text{sim}(c_i, e)$, falls below a threshold $\tau_{coverage}$.
  - **Relevance Deficiency,** if it contains entities that are only locally relevant to one aspect but have low relevance to the overall query $Q$. Formally, a path $P$ is deficient in relevance if it contains entities with a low GlobalSupport score.

To implement this, the module analyzes the path entities to identify specific deficiencies:

(1) *To detect Completeness Deficiency,* it flags concepts whose coverage scores are below a threshold $\tau_{coverage}$ as the set of missing concepts $C_{missing}$. This indicates that the current path has failed to explore entities relevant to certain key concepts.

(2) *To detect Relevance Deficiency,* it scrutinizes each entity in the path by calculating its GlobalSupport$(e, Q, C)$. This metric balances concept coverage with overall question relevance:

$$\text{EntityScope}(e, C) = \frac{|\{c \in C : \text{sim}(e, c) > \tau_c\}|}{|C|}$$

$$\text{GlobalSupport}(e, Q, C) = \alpha \cdot \text{EntityScope}(e, C) + (1 - \alpha) \cdot \text{sim}(e, Q)$$

where EntityScope$(e, C)$ measures how many concepts entity $e$ is relevant to, $\alpha$ balances the two components, and $\tau_c$ is the concept relevance threshold. Entities with low global support are flagged as misleading.

The evaluation module produces a diagnosis $\mathcal{I}$ of the issues found in $P_{candidate}$.

- **Adjustment Module.** If deficiencies are diagnosed ($\mathcal{I} \neq \emptyset$), this module executes a corrective action. Instead of triggering new, independent searches, our approach performs a targeted correction of the flawed trajectory itself. This Strategic re-exploration process involves:

(1) *Formulating an Adjustment Strategy.* For a Completeness Deficiency, the system identifies external entities relevant to the missing concepts $C_{missing}$ and assigns them a positive weight adjustment $+\delta$. For a Relevance Deficiency, it assigns a negative weight adjustment $-\delta$ to misleading entities. These adjustments modify the original edge weights during subsequent exploration:

$$w_{adjusted}(e_i, e_j) = w_{original}(e_i, e_j) + \text{adjustment}(e_j)$$

(2) *Executing an Informed Re-search.* Rather than restarting from scratch, the module selects a strategic restart point from the previous path. For Completeness Deficiency, it chooses the entity most relevant to the missing concepts. For Relevance Deficiency, it selects the entity with the highest global support:

$$e_{restart} = \arg\max_{e \in E_P} f(e)$$

where

$$f(e) = \begin{cases} \max_{c \in C_{missing}} \text{sim}(e, c) & \text{Completeness issue} \\ \text{GlobalSupport}(e, Q, C) & \text{Relevance issue} \end{cases}$$

The system then performs greedy search from $e_{restart}$, selecting the next entity with maximum adjusted weight at each step until a new candidate path is generated.

This targeted re-exploration mechanism directly addresses the diagnosed issues, efficiently correcting the path's trajectory without losing all prior progress.

The metacognitive cycle repeats until one of three convergence conditions is met: (1) the Evaluation module detects no deficiencies, (2) the change between the new path and the previous one is minimal, measured by entity overlap similarity exceeding a threshold $\tau_{similarity}$, or (3) a maximum iteration limit $N_{max}$ is reached.

Specifically, two paths are considered similar if:

$$\text{PathSimilarity}(P_1, P_2) = \frac{|E_{P_1} \cap E_{P_2}|}{|E_{P_1} \cup E_{P_2}|} > \tau_{similarity}$$

where $E_{P_1}$ and $E_{P_2}$ are the entity sets of the two paths. In our implementation, we set the default $N_{max} = 3$ and $\tau_{similarity} = 0.8$ based on empirical validation.

## 3.4 Multi-Evidence Integration and Answer Generation

After all initial entities complete their independent metacognitive cycles, MetaKGRAG performs comprehensive multi-path integration to synthesize the collected evidence. The evidence synthesis process merges all refined paths into a comprehensive evidence subgraph $\mathcal{S}$ while removing duplicate information, thus preserving diverse perspectives while eliminating redundancy.

The evidence subgraph is then converted to natural language context through a structured transformation. Each triple $(h_i, r_i, t_i) \in \mathcal{S}$ is converted into natural language statements following the template "Evidence $i$: $h_i$ $r_i$ $t_i$", creating a structured evidence list that maintains the semantic relationships between entities while being interpretable by the LLM. Finally, the LLM generates the answer $A$ using both the original question $Q$ and the curated evidence context through a prompt that instructs the model to analyze the evidence and answer the question. The complete prompt template for answer generation is detailed in Appendix 7.

## 4 EXPERIMENTS

## 4.1 Experimental Setup

*4.1.1 Datasets.* As shown in Tab 1, to evaluate MetaKGRAG's effectiveness and generalization ability across diverse domains, we conduct experiments on five datasets spanning commonsense, medical, and legal domains.

**Table 1: The statistics of datasets.**

| Domain | Dataset | Questions | Language |
|---|---|---|---|
| Commonsense | CommonsenseQA | 700 | English |
| Medical | CMB-Exam | 2,000 | Chinese |
| | ExplainCPE | 507 | Chinese |
| | webMedQA | 500 | Chinese |
| Legal | JEC-QA | 479 | Chinese |

We use CommonsenseQA [25] for commonsense knowledge evaluation, which contains multiple-choice questions requiring broad knowledge integration. For medical domains, we employ three datasets, CMB-Exam [26] with 2,000 sampled questions from Chinese medical professional examinations (including Nursing, Pharmacy, Postgraduate, and Professional), ExplainCPE [13] containing pharmaceutical questions from the National Licensed Pharmacist Examination with both answers and explanations, and webMedQA [9] featuring real-world patient-doctor conversations from online medical platforms. For legal domain, we use JEC-QA [35] from China's National Judicial Examination, which requires logical reasoning to apply legal materials to specific case scenarios. For more detailed information on these datasets, please refer to Appendix A.1.

*4.1.2 Evaluation Metrics.* For evaluation, we adopt a variety of different metrics. *Correct (Accuracy)*, *Wrong*, *Fail* are used for those with ground truth (e.g., CommonsenseQA, CMB-Exam, ExplainCPE), where *Fail* indicates the model fails to generate any answer. For tasks requiring generative answers (ExplainCPE, webMedQA), we use *ROUGE-L* [15] to measure lexical overlap with reference answers and *BERTScore* [34] to assess semantic similarity. To further evaluate the overall quality of the generated responses, we utilize *G-Eval* [17], a framework that leverages LLMs for evaluation, assessing the generated answers based on four key dimensions, *Coherence*, *Consistency*, *Fluency*, and *Relevance*. For all results, the best results are in **bold** and the second best results are underlined.

*4.1.3 Baselines.* To comprehensively evaluate the performance of our proposed MetaKGRAG framework, we select a diverse range of baselines, which can be categorized into three groups: Large Language Models, KG-RAG approaches, and self-refinement methods.

- *Large Language Models.* This group serves as a fundamental baseline to evaluate the capabilities of LLMs themselves without any external knowledge retrieval. We select leading commercial models including **GPT-4o**, **Claude 3.5 Sonnet**, **Gemini 1.5 Pro**, and OpenAI **o1-mini**. For open-source models, we use **Qwen2.5-7B** and **Qwen2.5-72B** [31] for Chinese tasks, and **Llama-3-8B** and **Llama-3-70B** [4] for English tasks.

- *KG-RAG Approaches.* We compare against several representative KG-RAG approaches. **Vanilla KGRAG** [23] serves as the fundamental implementation, performing direct entity similarity-based retrieval and presenting facts to LLMs. **ToG (Think-on-Graph)** [24] guides LLMs to explore multiple reasoning paths within knowledge graphs for multi-hop reasoning. **MindMap** [30] constructs structured representations that integrate knowledge from subgraphs to enhance interpretability. **KGGPT** [11] handles complex queries by decomposing them into simpler clauses and constructing evidence graphs through separate retrievals. These baselines providing comprehensive coverage of existing retrieval paradigms from simple similarity matching to sophisticated multi-hop strategies.

- *Self-Refinement Methods.* As discussed in Sec. 5, we compare several Self-Refinement methods. **Chain-of-Thought (CoT)** [29] encourages more thoughtful reasoning by generating intermediate steps. **Metacognitive Prompting** [27] employs a metacognitive prompt to guide the model to self-critique and refine its reasoning. For retrieval-augmented scenarios, we construct baselines by combining a standard KG-RAG retriever with frameworks like **FLARE** [10], which performs active retrieval when generation confidence is low, and **ReAct** [32], which synergizes reasoning and acting to iteratively search for information. We also adapt **Meta RAG** [37], which enhances RAG by implementing a three-step metacognitive process of monitoring, evaluating, and planning to enable the model to introspectively identify and rectify its own knowledge gaps and reasoning errors. These baselines represent a straightforward "stacking" approach that adds a refinement mechanism in KG-RAG.

*4.1.4 Implementation Details.* Our MetaKGRAG framework is implemented using both large and small-scale open-source models as its backbone. Specifically, we utilize the Qwen2.5 series (7B, 72B) for Chinese tasks and the Llama-3 series (8B, 70B) for English tasks.

**Table 2: Performance Comparison on ExplainCPE and JEC-QA using Accuracy.**

| Type | Method | ExplainCPE | JEC-QA |
|------|--------|------------|--------|
| **Without Retrieval** | | | |
| LLM Only | Qwen2.5-7B | 69.76 | 65.06 |
| | Qwen2.5-72B | 81.82 | 80.13 |
| | GPT4o | 79.64 | 78.51 |
| | o1-mini | 75.10 | 70.15 |
| | Claude3.5-Sonnet | 76.88 | 75.33 |
| | Gemini1.5-Pro | 69.37 | 76.18 |
| Self-Refine | Chain-of-Thought | 82.53 | 81.02 |
| | Meta Prompting | 83.11 | <u>81.67</u> |
| **With Retrieval** | | | |
| KG-RAG | KGRAG | 78.53 | 73.88 |
| | ToG | 78.85 | 74.90 |
| | MindMap | 78.41 | 71.55 |
| | KGGPT | 78.86 | 71.83 |
| Self-Refine | FLARE | 80.23 | 75.81 |
| | ReAct | 81.51 | 76.92 |
| | Meta Prompting | 80.88 | 76.25 |
| | Meta RAG | 81.93 | 77.31 |
| Ours | MetaKGRAG (Qwen2.5-7B) | <u>85.97</u> | 77.10 |
| | MetaKGRAG (Qwen2.5-72B) | **91.70** | **88.49** |

**Table 3: Performance Comparison on webMedQA. Prec. and Rec. represent Precision and Recall, respectively.**

| Type | Method | Prec. | Rec. | F1 |
|------|--------|-------|------|-----|
| **Without Retrieval** | | | | |
| LLM Only | Qwen2.5-7B | 66.68 | 71.14 | 68.85 |
| | Qwen2.5-72B | 70.12 | 73.58 | 71.81 |
| | GPT4o | 72.53 | 75.11 | 73.80 |
| | o1-mini | 68.24 | 72.05 | 70.09 |
| | Claude3.5-Sonnet | 71.89 | 74.32 | 73.09 |
| | Gemini1.5-Pro | 72.01 | 74.88 | 73.42 |
| Self-Refine | Chain-of-Thought | 71.21 | 74.63 | 72.88 |
| | Meta Prompting | 71.85 | 75.01 | 73.40 |
| **With Retrieval** | | | | |
| KG-RAG | KGRAG | 73.15 | 76.02 | 74.56 |
| | ToG | 73.89 | 76.55 | 75.19 |
| | MindMap | 72.93 | 75.81 | 74.34 |
| | KGGPT | 73.51 | 76.23 | 74.85 |
| Self-Refine | FLARE | 74.22 | 76.91 | 75.54 |
| | ReAct | 74.98 | 77.53 | 76.23 |
| | Meta Prompting | 74.53 | 77.18 | 75.83 |
| | Meta RAG | 75.11 | 77.82 | 76.44 |
| Ours | MetaKGRAG (Qwen2.5-7B) | <u>76.53</u> | <u>78.91</u> | <u>77.70</u> |
| | MetaKGRAG (Qwen2.5-72B) | **78.02** | **80.15** | **79.07** |

**Table 4: Performance comparison on CommonsenseQA.**

| Type | Method | Correct | Wrong | Fail |
|------|--------|---------|-------|------|
| **Without Retrieval** | | | | |
| LLM Only | Llama-3-8B | 73.82 | 26.04 | 0.14 |
| | Llama-3-70B | 81.76 | 18.24 | 0.00 |
| | GPT-4o | 84.54 | 15.44 | 0.02 |
| | o1-mini | 81.41 | 18.45 | 0.14 |
| | Claude3.5-Sonnet | 82.55 | 17.45 | 0.00 |
| | Gemini1.5-Pro | 83.83 | 16.17 | 0.00 |
| Self-Refine | Chain-of-Thought | 83.52 | 16.48 | 0.00 |
| | Meta Prompting | 84.13 | 15.87 | 0.00 |
| **With Retrieval** | | | | |
| KG-RAG | KGRAG | 85.04 | 14.96 | 0.00 |
| | ToG | 85.81 | 14.19 | 0.00 |
| | MindMap | 85.53 | 14.47 | 0.00 |
| | KGGPT | 86.22 | 13.78 | 0.00 |
| Self-Refine | FLARE | 86.54 | 13.46 | 0.00 |
| | ReAct | 87.31 | 12.69 | 0.00 |
| | Meta Prompting | 86.81 | 13.17 | 0.02 |
| | Meta RAG | 87.52 | 12.48 | 0.00 |
| Ours | MetaKGRAG (Llama-3-8B) | <u>88.54</u> | 11.46 | 0.00 |
| | MetaKGRAG (Llama-3-70B) | **92.11** | 7.89 | 0.00 |

For all semantic similarity calculations, such as matching concepts to entities and assessing path coverage, we employ the distiluse-base-multilingual-cased-v1 [19] embedding model due to its strong multilingual capabilities. The specific knowledge graphs for each dataset were custom-built to ensure relevance and quality. The detailed methodologies for knowledge graph construction, the specific prompt templates used for different stages of the framework (e.g., concept extraction, answer generation), and the experimental environment setup are all detailed in the Appendix A for reproducibility. The core parameters of MetaKGRAG will be analyzed in detail in the subsequent analysis section.

## 4.2 Results and Analysis

*4.2.1 Main Results.* The main experimental results are presented in Tab. 2, 3, 4, and 5. The results across all five datasets validate the effectiveness of our framework. We highlight three key findings.

(1) Our framework consistently outperforms all baselines, achieving 91.70% on ExplainCPE (+9.88% over best LLM), 92.11% on CommonsenseQA (+10.35%), and 88.49% on JEC-QA (+8.36%). These improvements across medical, commonsense, and legal domains demonstrate the effectiveness of our metacognitive approach. Regarding to the ROUGE-L and G-Eval results of webMedQA and ExplainCPE, please refer to the Appendix C. As shown in Tab. 9, our MetaKGRAG method also achieved scores that surpass other baselines. The results demonstrate that our method not only improves the accuracy of multiple-choice question answering, but also excels in explanatory answer generation capabilities.

(2) A critical finding is that simply stacking existing refinement methods on KG-RAG yields only marginal improvements. Methods like FLARE, ReAct, and Meta RAG. When combined with KG-RAG, improve performance by merely 1-3%. In stark contrast, MetaK-GRAG achieves 5-10% improvements over basic KG-RAG methods.

This dramatic difference reveals that document-oriented refinement strategies fail to address the unique challenges of graph exploration. Our perceive-evaluate-adjust cycle, specifically designed to handle path dependencies and structural constraints, proves essential for high-quality evidence retrieval in KG.

(3) MetaKGRAG demonstrates effectiveness with both small and large backbone models. With smaller models (7B/8B), MetaKGRAG already achieves competitive performance, for example, reaching 85.97% on ExplainCPE. When scaled to larger models (72B/70B), the performance gains are even more pronounced, suggesting that

**Table 5: Performance Comparison on CMB-Exam with six different types.**

| Type | Method | Nursing | | | Pharmacy | | | Postgraduate | | | Professional | | |
|------|--------|---------|---|---|----------|---|---|--------------|---|---|--------------|---|---|
| | | Correct | Wrong | Fail | Correct | Wrong | Fail | Correct | Wrong | Fail | Correct | Wrong | Fail |
| **Without Retrieval** | | | | | | | | | | | | | |
| LLM Only | Qwen2.5-7B | 80.96 | 18.84 | 0.20 | 77.56 | 22.24 | 0.20 | 80.36 | 19.64 | 0.00 | 74.15 | 25.85 | 0.00 |
| | Qwen2.5-72B | 89.80 | 10.08 | 0.12 | 90.08 | 9.92 | 0.00 | 88.18 | 12.62 | 0.20 | 83.98 | 16.02 | 0.00 |
| | GPT4o | 83.13 | 16.87 | 0.00 | 72.89 | 26.91 | 0.20 | 76.95 | 22.44 | 0.60 | 78.96 | 21.04 | 0.00 |
| | o1-mini | 74.50 | 25.50 | 0.00 | 60.44 | 39.56 | 0.00 | 63.13 | 36.27 | 0.60 | 73.55 | 26.45 | 0.00 |
| | Claude3.5-Sonnet | 75.90 | 24.10 | 0.00 | 65.86 | 34.14 | 0.00 | 69.54 | 30.46 | 0.00 | 73.75 | 26.25 | 0.00 |
| | Gemini1.5-Pro | 80.72 | 19.28 | 0.00 | 70.68 | 29.32 | 0.00 | 75.95 | 24.05 | 0.00 | 77.56 | 22.44 | 0.00 |
| Self-Refine | CoT | 90.15 | 9.73 | 0.12 | 90.31 | 9.69 | 0.00 | 88.54 | 11.26 | 0.20 | 84.22 | 15.78 | 0.00 |
| | Meta Prompting | 90.55 | 9.33 | 0.12 | 90.72 | 9.28 | 0.00 | 88.91 | 10.89 | 0.20 | 84.67 | 15.33 | 0.00 |
| **With Retrieval** | | | | | | | | | | | | | |
| KG-RAG | KGRAG | 88.15 | 11.85 | 0.00 | 85.33 | 14.47 | 0.20 | 85.12 | 14.88 | 0.00 | 82.71 | 17.29 | 0.00 |
| | ToG | 89.18 | 10.62 | 0.20 | 86.77 | 13.23 | 0.00 | 85.17 | 14.83 | 0.00 | 83.37 | 16.63 | 0.00 |
| | MindMap | 85.77 | 14.02 | 0.20 | 81.95 | 17.65 | 0.41 | 80.97 | 18.83 | 0.00 | 81.19 | 18.81 | 0.00 |
| | KGGPT | 86.74 | 13.26 | 0.00 | 86.13 | 13.87 | 0.00 | 86.15 | 13.85 | 0.00 | 83.62 | 16.38 | 0.00 |
| Self-Refine | FLARE | 89.11 | 10.89 | 0.00 | 87.15 | 12.65 | 0.20 | 87.52 | 12.48 | 0.00 | 85.18 | 14.82 | 0.00 |
| | ReAct | 90.12 | 9.88 | 0.00 | 88.23 | 11.57 | 0.20 | 88.43 | 11.57 | 0.00 | 86.27 | 13.73 | 0.00 |
| | Meta Prompting | 89.55 | 10.45 | 0.00 | 87.64 | 12.16 | 0.20 | 87.91 | 12.09 | 0.00 | 85.73 | 14.27 | 0.00 |
| | Meta RAG | 90.53 | 9.47 | 0.00 | 88.71 | 11.09 | 0.20 | 88.88 | 11.12 | 0.00 | 86.74 | 13.26 | 0.00 |
| Ours | MetaKGRAG (Qwen2.5-7B) | <u>91.78</u> | 8.22 | 0.00 | <u>91.58</u> | 8.42 | 0.00 | <u>89.78</u> | 10.22 | 0.00 | <u>88.78</u> | 11.22 | 0.00 |
| | MetaKGRAG (Qwen2.5-72B) | **96.54** | 3.46 | 0.00 | **98.79** | 1.21 | 0.00 | **92.99** | 7.01 | 0.00 | **94.59** | 5.41 | 0.00 |

stronger base models can better leverage the high-quality evidence paths produced by our metacognitive cycle.

*4.2.2 Ablation Studies.* To gain deeper insights into the specific mechanisms driving MetaKGRAG's performance, we conduct fine-grained ablation studies that examine individual functional components within our framework. Rather than removing entire **Metacognitive Cycle**, we selectively disable specific mechanisms to understand their individual contributions. In Evaluation module, **Completeness Check** identifies missing key concepts, **Relevance Check** detects entities to prevent misleading information. In Adjustment module, **Strategic Restart** intelligently selects optimal starting points for path re-exploration when issues are detected, rather than naively restarting from the original entity.

Table 6 shows that Completeness Check provides the largest contribution (5.04% on ExplainCPE, 5.48% on JEC-QA and 2.39% on CommonsenseQA), indicating that ensuring comprehensive information coverage is more critical than filtering irrelevant content. Relevance Check contributes substantially (4.32%, 4.22%, and 1.99% respectively), demonstrating the importance of avoiding misleading entities that satisfy keyword matching without providing meaningful information. Strategic Restart offers consistent but modest gains (3.16%, 3.56%, and 1.29%), proving that intelligent restart point selection outperforms naive re-exploration strategies.

**Table 6: Ablation study on ExplainCPE, JEC-QA, and CommonsenseQA with backbone model Qwen2.5-7B.**

| Configuration | ExplainCPE | JEC-QA | CommonsenseQA |
|---------------|-----------|--------|---------------|
| MetaKGRAG (Original) | 85.97 | 77.10 | 88.54 |
| w/o Metacognitive Cycle | 78.51 | 73.84 | 85.04 |
| *Ablating Specific Mechanisms* | | | |
| w/o Completeness Check | 80.93 | 71.62 | 86.15 |
| w/o Relevance Check | 81.65 | 72.88 | 86.55 |
| w/o Strategic Restart | 82.81 | 73.54 | 87.25 |

*4.2.3 Detailed Analysis.* To gain deeper insights into the internal mechanisms and effectiveness of our MetaKGRAG framework, we conduct a series of detailed analyses across multiple dimensions.

**Effectiveness Analysis.** We first verify whether MetaKGRAG's metacognitive cycle genuinely improves evidence quality through effective path refinement. We introduce the Path Refinement Rate (PRR) metric to quantify the degree of change in final evidence paths relative to initial paths. As shown in Fig. 4, our MetaKGRAG method demonstrates a high path refinement rate of 38.5% while achieving 85.97% accuracy. In contrast, baselines that simply incorporate self-refinement methods (Meta RAG and ReAct) exhibit refinement rates below 15% with correspondingly lower accuracy, while ToG shows zero refinement. This strong correlation between a high refinement rate and high accuracy proves that our adjustments are targeted and beneficial; our method makes more meaningful changes to improve answer quality, whereas other methods either do not refine the path or make fewer, less effective adjustments.

**Parameter Sensitivity Analysis.** We examine the framework's sensitivity to two core hyperparameters: the coverage threshold $\tau_{coverage}$ (used to judge path entity coverage of question key concepts) and the maximum iteration number $N_{max}$. As shown in Fig. 3 (a), $\tau_{coverage}$ performs optimally around 0.6 across both datasets. This is because a threshold set too low risks triggering unnecessary adjustments for irrelevant concepts, while one set too high may fail to identify and correct important evidence gaps. For $N_{max}$, as shown in Fig. 3 (b), its performance peaks at 3 iterations and then stabilizes, suggesting that most path deficiencies can be effectively resolved within three cycles, with further iterations offering diminishing returns. The consistency of these optimal values across different domains validates our default parameter choices and demonstrates the framework's robustness.

**Multi-start Retrieval Analysis.** To quantify the importance of the multi-start retrieval strategy, we analyze the impact of the number of starting entities on final performance. As illustrated in Fig. 3
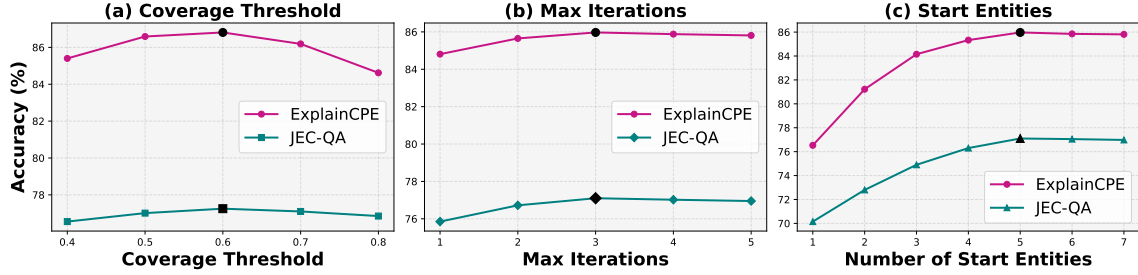
**Figure 3: Evaluation of different hyperparameters on ExplainCPE and JEC-QA.**

(c), accuracy on both datasets improves significantly as the number of starting entities increases, reaching an optimal point at 5 before plateauing. This is because complex questions often require multiple perspectives to construct a complete evidence subgraph, and starting from just one entity frequently misses critical information. Notably, medical questions (ExplainCPE) benefit more from additional starting points than legal questions (JEC-QA), likely because medical problems often involve more interconnected concepts that necessitate exploration from diverse angles. This result echoes our ablation study conclusions, further demonstrating that exploring knowledge graphs from multiple perspectives is crucial for ensuring retrieval comprehensiveness. **Please refer to Appendix C for more experimental results analysis and case study.**
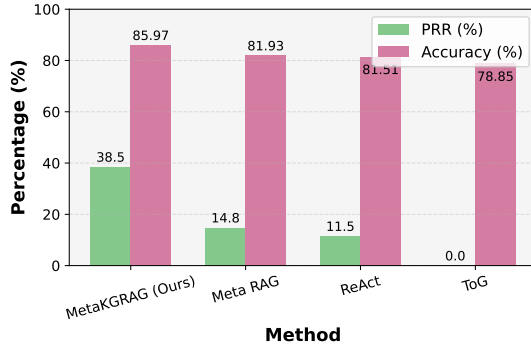


**Figure 4: Path refinement rate (PPR) and accuracy of different methods on ExplainCPE.**

## 5 RELATED WORK

**Knowledge Graph RAG.** Recent advances in KG-RAG have demonstrated the value of structured knowledge for complex reasoning tasks. ToG (Think-on-Graph) [24] guides LLMs to explore multiple reasoning paths within knowledge graphs through a beam search mechanism, where the model iteratively selects the most promising entities and relations to expand the current path. MindMap [30] constructs hierarchical structured representations by first identifying key entities, then building subgraph clusters around them, and finally integrating these clusters into a unified knowledge structure for enhanced interpretability. KGGPT [11] handles complex queries through a decomposition strategy, breaking down

multi-faceted questions into simpler sub-queries and constructing evidence graphs through separate retrieval processes for each component. Despite their sophisticated exploration strategies, these methods operate as open-loop systems that generate evidence paths in a single forward pass without mechanisms to assess or refine path quality during the retrieval process.

**Self-Refinement Methods.** The importance of quality control in retrieval has been recognized in the broader RAG domain, leading to several self-refinement approaches. ReAct [32] introduces a reasoning-acting loop where the system can perform new searches based on intermediate reasoning states, using action sequences like "Search[query]" and "Finish[answer]" to iteratively gather evidence. FLARE [10] employs a forward-looking approach that generates answers incrementally, triggering retrieval when confidence drops below a threshold, thus ensuring continuous evidence support throughout generation. Self-RAG [1] implements a more comprehensive framework with reflection tokens that enable the model to critique its own outputs and decide whether to retrieve additional information, revise responses, or continue generation. These methods excel in document-based retrieval, where individual pieces of evidence can be independently assessed and substituted. However, this approach mismatches the path-dependent nature of graph exploration, where each step constrains subsequent choices and simple replacement invalidates entire reasoning chains.

## CONCLUSION

In this paper, we proposed MetaKGRAG, a novel framework that enhances Knowledge Graph RAG by integrating a human-inspired metacognitive process to solve Cognitive Blindness. Through its Perceive-Evaluate-Adjust cycle, MetaKGRAG empowers the system to identify path-level deficiencies like incomplete evidence and relevance drift, and to perform targeted, trajectory-connected corrections. Experimental results across diverse medical, legal, and commonsense datasets demonstrated the superior performance of MetaKGRAG over strong baselines. For future work, we aim to explore more advanced, learnable adjustment strategies to further enhance the efficiency and adaptability of the metacognitive cycle.

## REFERENCES
[1] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum?id=hSyW5go0v8

[2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda

Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) *(NIPS '20)*. Curran Associates Inc., Red Hook, NY, USA, Article 159, 25 pages.

[3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. https://doi.org/10.18653/v1/N19-1423

[4] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).

[5] 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* abs/2404.16130 (2024). https://api.semanticscholar.org/CorpusID:269363075

[6] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A Survey on RAG Meeting LLMs: Towards Retrieval-Augmented Large Language Models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Barcelona, Spain) *(KDD '24)*. Association for Computing Machinery, New York, NY, USA, 6491–6501. https://doi.org/10.1145/3637528.3671470

[7] Xinyan Guan, Yanjiang Liu, Hongyu Lin, Yaojie Lu, Ben He, Xianpei Han, and Le Sun. 2024. Mitigating large language model hallucinations via autonomous knowledge graph-based retrofitting. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence (AAAI'24/IAAI'24/EAAI'24)*. AAAI Press, Article 2022, 9 pages. https://doi.org/10.1609/aaai.v38i16.29770

[8] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: retrieval-augmented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*. JMLR.org, Article 368, 10 pages.

[9] Junqing He, Mingming Fu, and Manshu Tu. 2019. Applying deep matching networks to Chinese medical question answering: A study and a dataset. *BMC Medical Informatics and Decision Making* 19, 2 (2019), 52. https://doi.org/10.1186/s12911-019-0761-8

[10] Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active Retrieval Augmented Generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 7969–7992. https://doi.org/10.18653/v1/2023.emnlp-main.495

[11] Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi. 2023. KG-GPT: A General Framework for Reasoning on Knowledge Graphs Using Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 9410–9421. https://doi.org/10.18653/v1/2023.findings-emnlp.631

[12] Emily R Lai. 2011. Metacognition: A Literature Review. https://api.semanticscholar.org/CorpusID:146606759

[13] Dongfang Li, Jindi Yu, Baotian Hu, Zhenran Xu, and Min Zhang. 2023. ExplainCPE: A Free-text Explanation Benchmark of Chinese Pharmacist Examination. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 1922–1940. https://doi.org/10.18653/v1/2023.findings-emnlp.129

[14] Zhuoqun Li, Xuanang Chen, Haiyang Yu, Hongyu Lin, Yaojie Lu, Qiaoyu Tang, Fei Huang, Xianpei Han, Le Sun, and Yongbin Li. 2025. StructRAG: Boosting Knowledge Intensive Reasoning of LLMs via Inference-time Hybrid Information Structurization. In *The Thirteenth International Conference on Learning Representations*. https://openreview.net/forum?id=GhexuBLxbO

[15] Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*. Association for Computational Linguistics, Barcelona, Spain, 74–81. https://www.aclweb.org/anthology/W04-1013

[16] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics* 12 (2024), 157–173. https://doi.org/10.1162/tacl_a_00638

[17] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG Evaluation using Gpt-4 with Better Human Alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 2511–2522. https://doi.org/10.18653/v1/2023.emnlp-main.153

[18] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 9802–9822. https://doi.org/10.18653/v1/2023.acl-long.546

[19] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. http://arxiv.org/abs/1908.10084

[20] Pranab Sahoo, Prabhash Meharia, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. 2024. A Comprehensive Survey of Hallucination in Large Language, Image, Video and Audio Foundation Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 11709–11724. https://doi.org/10.18653/v1/2024.findings-emnlp.685

[21] Ahmmad O. M. Saleh, Gokhan Tur, and Yucel Saygin. 2024. SG-RAG: Multi-Hop Question Answering With Large Language Models Through Knowledge Graphs. In *Proceedings of the 7th International Conference on Natural Language and Speech Processing (ICNLSP 2024)*, Mourad Abbas and Abed Alhakim Freihat (Eds.). Association for Computational Linguistics, Trento, 439–448. https://aclanthology.org/2024.icnlsp-1.45/

[22] Gregory Schraw and David Moshman. 1995. Metacognitive Theories. *Educational Psychology Review* 7 (12 1995), 351–371. https://doi.org/10.1007/BF02212307

[23] Karthik Soman, Peter W Rose, John H Morris, Rabia E Akbas, Brett Smith, Braian Peetoom, Catalina Villouta-Reyes, Gabriel Cerono, Yongmei Shi, Angela Rizk-Jackson, Sharat Israni, Charlotte A. Nelson, Sui Huang, and Sergio Baranzini. 2023. Biomedical knowledge graph-optimized prompt generation for large language models. *Bioinformatics* 40 (2023). https://api.semanticscholar.org/CorpusID:265498312

[24] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Sai Wang, Chen Lin, Yeyun Gong, Lionel M. Ni, Heung yeung Shum, and Jian Guo. 2023. Think-on-Graph: Deep and Responsible Reasoning of Large Language Model on Knowledge Graph. In *International Conference on Learning Representations*. https://api.semanticscholar.org/CorpusID:263333907

[25] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, Minneapolis, Minnesota, 4149–4158. https://doi.org/10.18653/v1/N19-1421

[26] Xidong Wang, Guiming Chen, Song Dingjie, Zhang Zhiyi, Zhihong Chen, Qingying Xiao, Junying Chen, Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, and Haizhou Li. 2024. CMB: A Comprehensive Medical Benchmark in Chinese. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Kevin Duh, Helena Gomez, and Steven Bethard (Eds.). Association for Computational Linguistics, Mexico City, Mexico, 6184–6205. https://doi.org/10.18653/v1/2024.naacl-long.343

[27] Yuqing Wang and Yun Zhao. 2024. Metacognitive Prompting Improves Understanding in Large Language Models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Kevin Duh, Helena Gomez, and Steven Bethard (Eds.). Association for Computational Linguistics, Mexico City, Mexico, 1914–1926. https://doi.org/10.18653/v1/2024.naacl-long.106

[28] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent Abilities of Large Language Models. *Transactions on Machine Learning Research* (2022). https://openreview.net/forum?id=yzkSU5zdwD Survey Certification.

[29] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems*, Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (Eds.). https://openreview.net/forum?id=_VjQlMeSB_J

[30] Yilin Wen, Zifeng Wang, and Jimeng Sun. 2024. MindMap: Knowledge Graph Prompting Sparks Graph of Thoughts in Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 10370–10388. https://doi.org/10.18653/v1/2024.acl-long.558

[31] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115* (2024).

[32] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *International Conference on Learning Representations (ICLR)*.

[33] Miao Zhang, Rufeng Dai, Ming Dong, and Tingting He. 2022. DRLK: Dynamic Hierarchical Reasoning with Language Model and Knowledge Graph for Question Answering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 5123–5133. https://doi.org/10.18653/v1/2022.emnlp-main.342

[34] Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. In *International Conference on Learning Representations*. https://openreview.net/forum?id=SkeHuCVFDr

[35] Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020. JEC-QA: A Legal-Domain Question Answering Dataset. In *Proceedings of AAAI*.

[36] Lexin Zhou, Wout Schellaert, Fernando Martínez-Plumed, Yael Moros-Daval, Cèsar Ferri, and José Hernández-Orallo. 2024. Larger and more instructable language models become less reliable. *Nature* (2024), 1–8. https://doi.org/10.1038/s41586-024-07930-y

[37] Yujia Zhou, Zheng Liu, Jiajie Jin, Jian-Yun Nie, and Zhicheng Dou. 2024. Metacognitive Retrieval-Augmented Large Language Models. In *The Web Conference 2024*. https://openreview.net/forum?id=TW2gJyR6Mj

# A COMPLEMENTARY EXPERIMENTAL SETTINGS

## A.1 Datasets

The detailed descriptions of the adopted datasets are summarized as follows:

- CommonsenseQA [25] is a multiple-choice QA dataset specifically designed to evaluate commonsense reasoning capabilities. Each question is accompanied by five candidate answers, only one of which is correct.
- CMB-Exam [26] covers 280,839 questions from six major medical professional qualification examinations, including physicians, nurses, medical technologists and pharmacists, as well as Undergraduate Disciplines Examinations and Graduate Entrance Examination in the medical field at China. Given the extensive scale of CMB-Exam, we sample a subset of CMB-Exam that comprises 2,000 questions, where 500 questions are randomly sampled from nurses, pharmacists, Undergraduate Disciplines Examination, and Graduate Entrance Examination categories.
- ExplainCPE [13] is a Chinese medical benchmark dataset containing over 7K instances from the National Licensed Pharmacist Examination. This dataset is distinctive in providing both multiple-choice answers and their corresponding explanations.
- webMedQA [9] is a large-scale Chinese medical QA dataset constructed from professional health consultation websites (such as Baidu Docto). The dataset contains 63,284 questions covering various clinical departments including internal medicine, surgery, gynecology, and pediatrics, with answers provided by doctors and experienced users. The dataset has been preprocessed to remove web tags, links, and garbled characters, retaining only Chinese and English characters, numbers, and punctuation.
- JEC-QA [35] is a legal domain dataset collected from the National Judicial Examination of China. It serves as a comprehensive evaluation of professional skills required for legal practitioners. The

dataset is particularly challenging as it requires logical reasoning abilities to retrieve relevant materials and answer questions correctly.

## A.2 Implementation details

Our framework is built on LangChain[1]. The local open-source LLMs are deployed based on the llama.cpp[2] project. Except for the context window size, which is adjusted according to the dataset, all other parameters use default configurations, such as temperature is 0.8. Both LangChain and llama.cpp are open-source projects, providing good transparency and reproducibility. For computational resources, all experiments were conducted on a cluster of 8 NVIDIA RTX 3090 GPUs. Due to computational constraints and to ensure fair comparison across different model scales, we applied 4-bit quantization to locally deployed LLMs. For the evaluation, we employed Bert Score metrics using "bert-base-chinese [3]" model, while ROUGE Score version 0.1.2 was utilized. Due to resource constraints, G-Eval assessments were conducted using locally deployed Qwen2.5-72B.

# B KNOWLEDGE GRAPH CONSTRUCTION

We employed a consistent KG construction method for all datasets, utilizing LLMs to extract knowledge triples from the datasets to build specialized KGs. The prompt example is shown in Tab. 8. All KGs were deployed using Neo4j[3].

# C COMPLEMENTARY EXPERIMENTAL RESULTS

## C.1 Analysis of Generative Quality

In addition to evaluating answer accuracy, we assessed the quality of the generated explanations for tasks requiring free-form responses (ExplainCPE and webMedQA). We used ROUGE-L to measure lexical overlap with reference answers and G-Eval to evaluate the Coherence, Consistency, Fluency, and Relevance of the generated text. The results, presented in Table 9, demonstrate that MetaKGRAG's advantages extend beyond correctness to significantly enhance generative quality. Across both datasets, MetaK-GRAG consistently achieves the highest scores in all evaluated dimensions. For instance, on ExplainCPE, MetaKGRAG (Qwen2.5-72B) achieves a ROUGE-L score of 28.45, a substantial improvement of over 4.4 points compared to the best-performing baseline (Meta RAG). This significant improvement can be attributed to the high-quality evidence subgraphs produced by our metacognitive cycle. The path-aware refinement process does not just retrieve relevant facts; it constructs a coherent, logically connected evidence narrative. This well-structured context enables the LLM to generate answers that are not only factually accurate but also more fluent, consistent, and directly relevant to the user's query. In contrast, while other retrieval methods provide a performance lift, the potentially fragmented or incomplete evidence they retrieve limits the ultimate quality of the generated text. This analysis confirms that the path-level self-correction of MetaKGRAG is crucial for both accuracy and the quality of explanatory generation.

---

[1]https://www.langchain.com/

[2]https://github.com/ggml-org/llama.cpp

[3]https://neo4j.com/

## C.2 Analysis of Concept Relevance Threshold

The concept relevance threshold, $\tau_c$, is a hyperparameter within the Evaluation module used to determine if an entity provides meaningful support for a query concept (i.e., if $\text{sim}(e, c) > \tau_c$). A threshold set too low might accept noisy entities, while one set too high could prematurely discard useful ones. To assess the framework's sensitivity to this parameter, we conducted a tuning experiment on both ExplainCPE and JEC-QA datasets, varying $\tau_c$ from 0.1 to 0.5. The results are illustrated in Figure 5 (a). As shown in the figure, the model's performance remains highly stable across the entire range of $\tau_c$ values. For both datasets, the accuracy fluctuates by less than 0.4%, with a slight peak around $\tau_c = 0.3$. This indicates that while a reasonably calibrated threshold is beneficial, MetaKGRAG is not overly sensitive to its precise value. The framework's robustness in this regard simplifies its deployment in new domains, as it does not require extensive, dataset-specific tuning of this parameter.

## C.3 Analysis of Weight Adjustment

The weight adjustment parameter, $\delta$, controls the magnitude of the positive or negative incentive applied to entities during the Strategic Re-exploration phase of the Adjustment module. A larger $\delta$ more aggressively steers the search away from misleading entities and towards missing concepts. We validated the impact of $\delta$ by testing values from 0.1 to 0.5. The results, depicted in Figure 5 (b), demonstrate that the framework exhibits low sensitivity to the specific value of $\delta$. Performance on both datasets forms a plateau, with optimal results achieved for $\delta$ values between 0.2 and 0.3.

The minimal variation in accuracy suggests that as long as the adjustment provides a clear directional signal, its exact magnitude is not a critical factor.
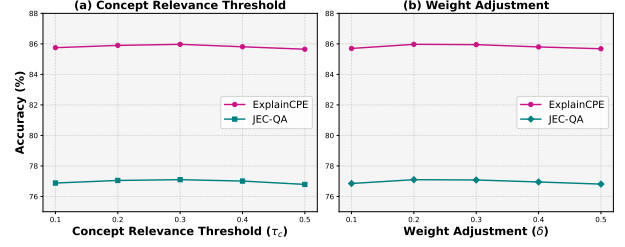


**Figure 5: Evaluation of other hyperparameters on ExplainCPE and JEC-QA.**

## C.4 Case Study

We conduct a case study on the ExplainCPE, . To visually demonstrate the operational difference between MetaKGRAG and baseline methods when handling complex queries, we present a representative case from our medical dataset. As shown in Tab. 10, this case illustrates the "Cognitive Blindness" issue in traditional KG-RAG, the path-dependency failure of self-refinement methods, and showcases how our metacognitive cycle effectively corrects the evidence-gathering path.

**Table 7: Prompt Example for Answer Generation**

```
prompt = f"""Your task is to accurately understand the question requirements and provide
          reasonable answers and explanations based on the provided reference content.
          Input Question:
          {question}
          You have the following medical evidence knowledge:
          {evidence_text}
          What is the answer to this multiple-choice question? Answer the question by
          referring to the provided medical evidence knowledge. First, choose the answer
          from (A\B\C\D\E), output the answer option, then explain the reasoning.
          """
```

**Table 8: Prompt Example for Knowledge Graph Construction**

```
prompt = f"""As a professional knowledge extraction assistant, your task is to extract knowledge triples from the given question.
1. Carefully read the question description, all options, and the correct answer.
2. Focus on the core concept "{question_concept}" in the question.
3. Extract commonsense knowledge triples related to the question.
4. Each triple should be in the format: subject\predicate\object
5. Focus on the following types of relationships:
 - Conceptual relations
 - Object properties
 - Object functions
 - Spatial relations
 - Temporal relations
 - Causal relations
6. Each triple must be concrete and valuable commonsense knowledge.
7. Avoid subjective or controversial knowledge.
8. Ensure triples are logically sound and align with common sense.
Please extract knowledge triples from this multiple-choice question:
Question: {question}
Core Concept: {question_concept}
Correct Answer: {correct_answer}
Please output knowledge triples directly, one per line, in the format: subject\predicate\object. """
```

**Table 9: Comparison of ROUGE-L and G-Eval scores on ExplainCPE and webMedQA. Coh., Cons., Flu., and Rel. indicate Coherence, Consistency, Fluency, and Relevance, respectively.**

| Method | | ExplainCPE (Medical) | | | | | webMedQA (Medical) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | ROUGE-L | Coh. | Cons. | Flu. | Rel. | ROUGE-L | Coh. | Cons. | Flu. | Rel. |
| **Without Retrieval** | | | | | | | | | | | |
| LLM Only | Qwen2.5-7B | 18.55 | 95.11 | 90.23 | 91.88 | 85.74 | 20.15 | 94.88 | 89.95 | 90.15 | 84.69 |
| | Qwen2.5-72B | 21.34 | 96.53 | 92.88 | 93.15 | 88.91 | 23.41 | 96.15 | 91.83 | 92.55 | 87.81 |
| | GPT4o | 20.89 | 96.81 | 93.05 | 93.55 | 89.15 | 22.95 | 96.53 | 92.11 | 92.98 | 88.05 |
| | o1-mini | 17.98 | 94.88 | 89.75 | 91.03 | 85.01 | 19.88 | 94.13 | 89.15 | 89.87 | 84.11 |
| | Claude3.5-Sonnet | 19.52 | 95.83 | 91.54 | 92.51 | 87.13 | 21.73 | 95.77 | 90.89 | 91.83 | 86.58 |
| | Gemini1.5-Pro | 20.15 | 96.11 | 92.18 | 92.88 | 88.05 | 22.51 | 96.01 | 91.55 | 92.71 | 87.92 |
| Self-Refine | CoT | 21.88 | 96.95 | 93.51 | 93.88 | 89.53 | 23.98 | 96.83 | 92.55 | 93.18 | 88.51 |
| | Meta Prompting | 22.15 | 97.03 | 93.88 | 94.01 | 89.98 | 24.33 | 97.01 | 92.93 | 93.55 | 88.99 |
| **With Retrieval** | | | | | | | | | | | |
| KG-RAG | KGRAG | 22.58 | 97.15 | 94.01 | 94.18 | 90.15 | 24.88 | 97.18 | 93.11 | 93.88 | 89.53 |
| | ToG | 22.91 | 97.33 | 94.52 | 94.55 | 90.83 | 25.13 | 97.51 | 93.82 | 94.13 | 90.11 |
| | MindMap | 22.43 | 97.01 | 93.85 | 94.03 | 89.95 | 24.71 | 97.05 | 92.95 | 93.72 | 89.25 |
| | KGGPT | 22.83 | 97.28 | 94.41 | 94.48 | 90.71 | 25.05 | 97.43 | 93.71 | 94.01 | 89.98 |
| Self-Refine | FLARE | 23.15 | 97.51 | 94.88 | 94.83 | 91.15 | 25.53 | 97.83 | 94.18 | 94.51 | 90.83 |
| | ReAct | 23.88 | 97.98 | 95.53 | 95.11 | 92.01 | 26.15 | 98.15 | 94.98 | 95.03 | 91.95 |
| | Meta Prompting | 23.41 | 97.72 | 95.01 | 94.95 | 91.53 | 25.81 | 97.99 | 94.51 | 94.82 | 91.21 |
| | Meta RAG | 24.03 | 98.05 | 95.81 | 95.33 | 92.35 | 26.33 | 98.33 | 95.21 | 95.38 | 92.18 |
| Ours | MetaKGRAG (Qwen2.5-7B) | 25.17 | 98.53 | 96.58 | 96.01 | 93.51 | 27.55 | 98.81 | 96.03 | 96.15 | 93.11 |
| | MetaKGRAG (Qwen2.5-72B) | **28.45** | **99.11** | **98.03** | **97.88** | **96.15** | **30.18** | **99.25** | **98.15** | **98.01** | **96.53** |

**Table 10: A comparative analysis of a drug interaction case, showing how MetaKGRAG overcomes the limitations of KG-RAG baseline and self-refinement method.**

---

**Input Question**: Which of the following statements about medications for Alzheimer's disease is **incorrect**?
A. Cholinesterase inhibitors can increase the risk of stomach bleeding.
B. Donepezil is suitable for co-administration with NSAIDs.
C. Memantine's clearance can be affected by urinary pH.
D. Galantamine's metabolism can be inhibited by ketoconazole.

---

**Baseline Method (Vanilla KG-RAG)**:

*Analysis*: Employs a greedy search. It strongly associates `Donepezil` with its primary function, `treating Alzheimer's`, and explores this path, overlooking the interaction with `NSAIDs`.

*Initial Path*:

1. `Donepezil` $\xrightarrow{\text{treats}}$ `Alzheimer's Disease`

2. `Donepezil` $\xrightarrow{\text{is\_a}}$ `Cholinesterase Inhibitor`

*Problem Diagnosis*:

**Incomplete Evidence**: The path completely fails to retrieve the critical information about the adverse interaction between Cholinesterase Inhibitors and NSAIDs.

*Generate Answer*: Based on the retrieved evidence, the model cannot identify the incorrect statement. It might wrongly conclude that statement (B) is plausible.

---

**Self-Refinement Method**:

*Analysis*: After generating an initial path like the baseline, the refinement mechanism recognizes the answer is insufficient as it doesn't address `NSAIDs`.

*Corrective Action (Flawed)*: It triggers a **separate and isolated search** for the interaction between `Donepezil` and `NSAIDs`. This returns a new, isolated fact: `Cholinesterase Inhibitors + NSAIDs → increased risk of GI bleeding`.

*Problem Diagnosis (Path Dependency Failure)*: The system now has two **disconnected evidence fragments**. It cannot integrate the new interaction fact with the original path. It fails to build a coherent reasoning chain explaining *why* the interaction occurs (e.g., via increased gastric acid secretion), and thus struggles to synthesize a confident and well-grounded response.

*Generate Answer*: The model might mention a risk but cannot provide a clear explanation, or it may be confused by the fragmented evidence and fail to definitively identify (B) as incorrect.

---

**Our Method (MetaKGRAG)**:

*1. Initial Path Generation*: Generates a similar preliminary path: `Donepezil → treats → Alzheimer's Disease`.

*2. Perceive-Evaluate-Adjust Cycle*:

**Perceive**: The framework assesses the path's coverage and detects that the relationship with `NSAIDs` is missing.

**Evaluate**: It diagnoses a **Completeness Deficiency**, flagging `co_administration_with_NSAIDs` as the missing concept ($C_{\text{missing}}$).

**Adjust**: Instead of a disconnected search, it performs a **trajectory-connected correction**. It selects `Cholinesterase Inhibitor` as a strategic restart point to explore its class-level properties.

*Refined Path (Coherent and Connected)*:

1. `Donepezil` $\xrightarrow{\text{is\_a}}$ `Cholinesterase Inhibitor`

2. `Cholinesterase Inhibitor` $\xrightarrow{\text{increases}}$ `Gastric Acid Secretion`

3. `Gastric Acid Secretion` $\xrightarrow{\text{heightens\_risk\_with}}$ `NSAIDs`

4. `NSAIDs` $\xrightarrow{\text{can\_cause}}$ `Gastrointestinal Bleeding`

*Generate Answer*: With the complete and connected evidence path, the model understands the full causal chain: Donepezil increases gastric acid, which heightens the risk of bleeding when combined with NSAIDs. It thus confidently and correctly identifies statement (B) as **incorrect**.

---