# 💇🏽‍♀️ Talk to Right Specialists: Iterative Routing in Multi-agent System for Question Answering

**Anonymous ACL submission**

## Abstract

Large language model (LLM)-based applications have increasingly leveraged retrieval-augmented generation (RAG) techniques to provide reliable responses, particularly for queries demanding knowledge of private domains. Practical constraints, such as data sovereignty regulations, can hinder the centralized aggregation of private knowledge. This can create challenges in situations where (1) a user comes with a question but has no idea which applications have the related knowledge to answer, or (2) the question requires cross-domain knowledge to answer.

In this work, we abstract each RAG application with private knowledge as an RAG-based agent. We propose RIRS, a framework with an efficient and accurate routing mechanism and an iterative refining-solving mechanism to orchestrate multiple RAG-based agents with private knowledge bases. The server routes queries to the most relevant agents by identifying the most related knowledge clusters by similarities in a vector space. For complicated questions, the server can iteratively aggregate responses to derive intermediate results and refine the question to bridge the gap toward a comprehensive answer. Extensive experiments demonstrate the effectiveness of RIRS, including how our routing algorithm precisely selects the agents and provides accurate responses to single-hop queries and how an iterative strategy achieves accurate, multi-step resolutions for complex queries.

## 1 Introduction

Large language models (LLMs) have revolutionized natural language processing (NLP) by demonstrating superior performance in question-answering (QA) tasks, often surpassing traditional systems in both accuracy and contextual understanding. Based on LLMs, retrieval-augmented generation (RAG) is a technique to integrate external knowledge sources, extracting the most relevant information for any input query to enable LLMs to answer questions beyond their training data and reduce their hallucination (Wu et al., 2024; Asai et al., 2023a; Lewis et al., 2020; Jiang et al., 2023; Izacard and Grave, 2020; Mallen et al., 2022; Kasai et al., 2024; Xiong et al., 2024). To further improve retrieval and response quality, many RAG applications are built in the form of agents (referred to as RAG-based agents in this paper) (Weng, 2023; Roucher, 2024; Joshi et al., 2024). However, because RAG-based agents' reliability is limited to the domain of their knowledge sources, there are two major inconveniences when serving users. (1) The domain and the boundary of the knowledge source are difficult to clearly define and usually unknown to users, so manual attempts with different agents may be required to obtain reliable answers. (2) Answering some questions may require cross-domain knowledge from different RAG-based agents.

One straightforward solution to overcome the inconveniences is to build a unified, large-scale knowledge repository by collecting all knowledge from worldwide existing RAG-based agents, as illustrated in Figure 1a. A unified RAG-based agent could be built with the repository to handle queries spanning multiple domains by pooling all available information. However, it is confronted with two critical limitations. Firstly, the computational complexity and scalability issues inherent in managing such a vast and diverse knowledge base could undermine the system's overall performance and responsiveness (Fan et al., 2024; Asai et al., 2023a). Secondly, it is impractical or even infeasible to invade knowledge sovereignty and construct a centralized knowledge repository, especially when the knowledge is intellectually protected or sensitive. For example, due to privacy concerns, ophthalmology hospital data derived from internal medical records is unavailable for merging.

An alternative is to deploy a distributed multi-

**(a) Knowledge Base Collection from All Agents**  **(b) Message Collection from All Agents**  **(c) Message Collection from Selected Agents**
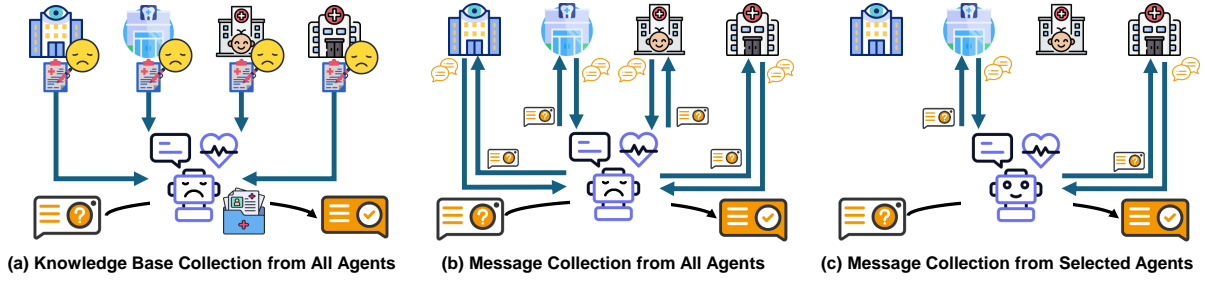
Figure 1: Collaboration Strategies of Multi-agent System for QA.

agent system, as illustrated in Figure 1b, which consists of a central server and multiple RAG-based agents. The server forwards a user's query to *all* available RAG-based agents and subsequently aggregates their responses. In this framework, each agent retains its local, domain-specific knowledge base, ensuring that its knowledge sovereignty is well protected and that its data are not merged into a centralized repository. This approach eliminates the need for a massive, unified knowledge base. However, routing every query to *all* agents incurs redundant computational overhead and cost, as not every agent is equipped to address every query, Consequently, it is worth considering building a distributed multi-agent system with *a routing mechanism* that selects the most appropriate agents in response to a query, as illustrated in Figure 1c.

In this work, we propose RIRS, a framework with a Routing and Iterative Refining-Solving mechanism designed to effectively reconcile multiple RAG-based agents. For the routing mechanism, each agent partitions its local knowledge base into disjoint clusters, and the central server collects these clustered knowledge representations. When a user query arrives, the server computes its similarity to the collected clusters and forwards the query only to those agents whose knowledge clusters are most relevant. This training-free mechanism does not require additional training and is inherently privacy-preserving, and raw knowledge remains confined within the individual agents. Moreover, our method employs an iterative refining-solving strategy to handle complex queries that involve multiple reasoning steps or span across different domains. Specifically, the server sequentially routes a query to the appropriate agents and, after each response, simplifies the query by removing the addressed portion.

**Contributions.**   The major contributions of this work are listed as follows:

- To the best of our knowledge, this is the first work

that considers knowledge sovereignty issues under multi-agent, enabling efficient collaboration across various specialized agents without collecting unnecessary information.

- We introduce RIRS, a training-free iterative routing mechanism that selects the most proper agents and collaborates on complex user queries.

- Our experimental results demonstrate the effectiveness of RIRS in handling various queries: The routing mechanism enables precise agent selections, while the iterative refining-solving mechanism achieves accurate, multi-step resolutions for complex queries.

## 2   Related Works

**Retrieval Augmented Generation (RAG).**   RAG has gained substantial interest in academic research as a robust framework that integrates external knowledge sources into large language models to enhance the quality and reliability of generated responses (Lewis et al., 2020; Jiang et al., 2023; Chen et al., 2017; Guu et al., 2020; Karpukhin et al., 2020; Izacard and Grave, 2020; Borgeaud et al., 2022; Yu, 2022; Shi et al., 2023; Yan et al., 2024; Asai et al., 2023b; Li et al., 2023b; Press et al., 2022; Chan et al., 2024; Su et al., 2024). Notable recent contributions in this domain include RankRAG (Yu et al., 2024), which reranks the selected knowledge pieces and generates a response with genuinely important ones; EfficientRAG (Zhuang et al., 2024), an approach that iteratively generates new queries by sorting out the portion addressed by retrieved knowledge until a multi-hop question can be well-addressed; PlanRAG (Verma et al., 2024), which decomposes complex queries into interrelated atomic sub-queries by formulating a reasoning plan as a directed acyclic graph (DAG). While these methods primarily focus on single-agent retrieval frameworks, our work differentiates itself by addressing the challenge of coordinating knowledge retrieval across multiple

agents, thereby enabling a more comprehensive handling of diverse and cross-domain queries.

**Routing Mechanism in Multi-agent System.** LLM-based multi-agent systems leverage the collective intelligence and specialized capabilities of multiple expert agents to collaboratively tackle complex problems, a research direction that has garnered significant interests (Hong et al., 2023; Li et al., 2023a; Wu et al., 2023; Chen et al., 2023a,b; Zhao et al., 2023; Guo et al., 2024; Chen et al., 2024b). In such systems, a robust routing mechanism is essential to direct each query to the most appropriate agent based on its unique expertise (Shnitzer et al., 2023; Lu et al., 2023; Zhao et al., 2024; Srivatsa et al., 2024; Li et al., 2024; Lu et al., 2024; Chen et al., 2024a; Addison et al., 2024).

Recent works in this area include Chameleon (Lu et al., 2024), which requires a comprehensive textual description of each agent's capabilities such that a well-trained LLM can select the most appropriate agents to address an input query; RouterDC (Chen et al., 2024a), which trains a lightweight model to dynamically route queries to the most suitable agent within a predefined set; and C-FedRAG (Addison et al., 2024), which forwards queries to randomly selected agents to retrieve relevant documents such that a central entity can collect these documents to generate a final answer.

In contrast, our approach accurately characterizes each agent's knowledge capacity based on their own data, thereby eliminating the need for exhaustive textual descriptions and labor-intensive dataset labeling. Moreover, our effective routing mechanism ensures that queries are directed only to agents capable of providing comprehensive answers, with each agent summarizing its response based solely on local knowledge, thereby preserving data privacy and protecting knowledge sovereignty.

## 3   Multi-agent Framework

**Preliminary: RAG-based Agent.** An RAG-based agent is an advanced application that integrates retrieval mechanisms with LLMs to deliver accurate and contextually rich responses. Specifically, an RAG-based agent processes a question through three steps: (i) **Knowledge Retrieval:** It extracts relevant knowledge pieces using both sparse and dense retrievers from external sources (Robertson et al., 2004; Izacard et al., 2021; Xu et al., 2023); (ii) **Reranking:** It filters out unhelp-

ful or misleading information, allowing the generation model to focus on the most pertinent content (Yu et al., 2024); (iii) **Response Generation:** It combines the question and the pertinent external knowledge to produce an informed response with its backbone LLM.

**Motivations and Problem Statement.** To protect knowledge sovereignty and harness domain-specialized expertise, we propose a distributed multi-agent framework. In this framework, each RAG-based agent holds its own specialized knowledge base, while a central server coordinates query processing. Upon receiving an input query, the server routes it to all agents and aggregates their responses into a final answer, as described in Figure 1b. However, this standard operating procedure for QA tasks faces two major challenges:
• *(i) Irrelevant Agent Involvement:* When a query falls outside an agent's domain expertise, its participation not only introduces unnecessary computational and communication overhead but may also generate misleading information.
• *(ii) Incomplete Knowledge Fusion:* Queries that span multiple domains require seamless integration of responses from various agents and likely from multiple reasoning steps, and, without effective coordination, the final answer can be fragmented or partial.
These challenges highlight the need for *an intelligent routing mechanism* that strategically directs queries to the most relevant agents and aggregates their responses wisely by filtering out irrelevant knowledge. Such a mechanism is critical to reducing redundant processing, enhancing scalability, ensuring accurate, holistic query resolution, and maintaining data privacy in a decentralized setting.

## 4   RIRS

To address the challenges, we propose an intelligent routing mechanism to deal with both single-hop and multi-hop questions. Section 4.1 outlines the design of a query routing algorithm, where the server routes a query to a subset of agents according to their knowledge coverage represented in a vector space. Section 4.2 extends our design to address more challenging queries that require multi-round or cross-agent knowledge.

### 4.1   Routing Algorithm

The primary objectives of our routing algorithm are twofold: high quality of the final answer and

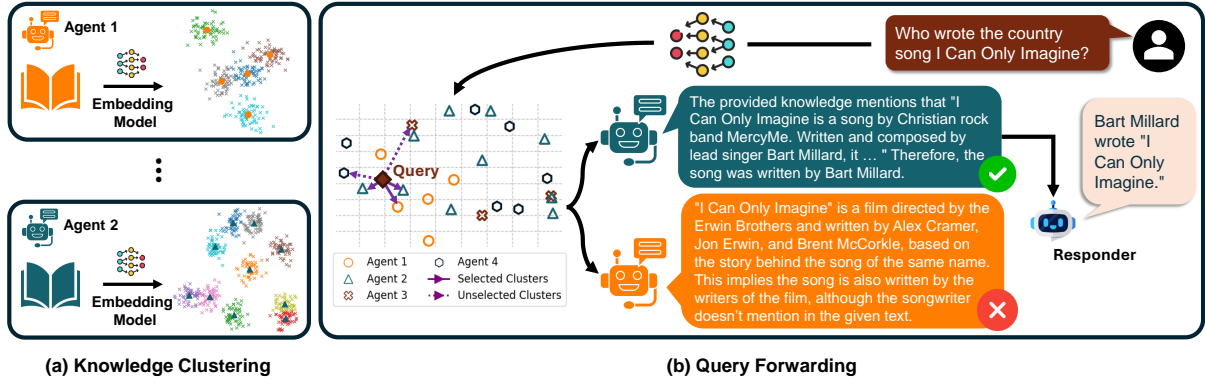**(a) Knowledge Clustering**  **(b) Query Forwarding**

Figure 2: Routing Mechanism

high efficiency in terms of both minimum additional latency of routing and minimum token consumption. To accomplish these objectives, an ideal server must possess three key abilities: (i) accurately assessing each agent's knowledge capabilities, (ii) selecting a necessary and sufficient subset of agents to reduce computational and communication costs, and (iii) verifying the validity of each response, encompassing both the analysis and the final answer. Notably, the second capability is closely linked to the server's understanding of the knowledge boundaries inherent to each agent.

To satisfy the requirements, we design a routing mechanism as shown in Figure 2, which consists of two primary stages: (a) knowledge clustering and (b) the query forwarding workflow. In the knowledge clustering stage (Figure 2a), each agent encodes its local knowledge (e.g., text chunks) into vectors using an identical embedding model, then partitions the knowledge into disjoint clusters and generates a representation for each cluster with the embeddings. The central server then collects these clustered knowledge representations and uses them to determine which agents are most relevant to an incoming query. As shown in Figure 2b, the server forwards the query to those agents whose clusters exhibit the highest similarity to the query. The selected RAG-based agents subsequently process the query based on their own knowledge, and the server aggregates and evaluates their responses. Finally, the server synthesizes the final response to the user. More details are as follows.

**Knowledge Clustering.** Suppose an RAG-based agent contains $m$ distinct knowledge pieces, represented as $e_1, \ldots, e_m$ in a vector space. To evaluate the agent's knowledgeability, we partition the knowledge pieces into $n$ disjoint clusters, $c_1, \ldots, c_n$. The goal is to minimize the maximum intra-cluster similarity, which can be formulated as:

$$\min_{\substack{c_1,\ldots,c_n \\ c_1 \cup \cdots \cup c_n = \{e_1,\ldots,e_m\}}} \max_{\substack{e_a, e_b \in c_k \\ k \in \{1,\ldots,n\}, e_a \neq e_b}} \mathsf{sim}(e_a, e_b) \quad (1)$$

Here, $\mathsf{sim}(e_a, e_b)$ denotes the cosine similarity between two knowledge pieces $e_a$ and $e_b$. By minimizing the maximum similarity within clusters, we ensure that knowledge pieces within the same cluster are as similar as possible, leading to more informative cluster representations.

The RAG-based agent follows a four-step process to solve this clustering objective and report its knowledgeability to the server:

• **Step 1: Compute Embeddings.** An RAG-based agent computes the embeddings for their own knowledge, which can be reused later in the knowledge retrieval of a dense method to an input query (Izacard and Grave, 2020).

• **Step 2: Knowledge Clustering.** Using maximum hierarchical clustering, the agent partitions the $m$ knowledge pieces into $n$ disjoint clusters. The distance between any two knowledge pieces is measured by their embeddings' cosine similarity.

• **Step 3: Calculate Cluster Representations.** For each cluster, the agent calculates a centroid by averaging the embeddings of all knowledge pieces within that cluster. This centroid serves as a representative summary of the cluster.

• **Step 4: Push Representations to Server.** The agent sends the centroids of all clusters to the server who uses this information to make routing decisions.

*Choice on the number of clusters.* Since different RAG-based agents hold varying amounts of knowledge, the number of clusters $n$ should not be constant across agents. Intuitively, agents with more knowledge pieces may have overlapping or redun-

4

dant knowledge, while agents with fewer pieces might specialize in sparse, distinct knowledge domains. To account for this, we set $n = \lfloor\sqrt{m}\rfloor$, aligning with the hypothesis that a larger number of knowledge pieces should correspond to more clusters while maintaining manageable granularity. In fact, the choice of $n$ has proved its effectiveness in the field of inverse file indexing in practice.

**Query Forwarding.** Once the server has gathered the knowledgeability of all RAG-based agents, it must effectively coordinate the agents to handle user queries. This involves selecting the most suitable agents based on the similarity between the query and the centroids provided by each agent.

Let us define $\bar{c}_j^{(i)}$ as the $j$-th centroid of RAG-based agent $i$, and define a function $f(\cdot)$ such that $f\left(\bar{c}_j^{(i)}\right) = i$. This function maps a centroid to the corresponding agent. Suppose there are $M$ RAG-based agents in the multi-agent system. For each agent $i \in \{1, \ldots, M\}$, there are $n_i$ centroids, denoted by the set $\left\{\bar{c}_j^{(i)}\right\}_{j=1}^{n_i}$. Let $x$ be the embedding of a query. The goal is to identify $k$ clusters whose centroids have the highest similarity scores with the query embedding. This can be formulated as:

$$\{\bar{c}_j\}_{j=1}^k \overset{\triangle}{=} \arg \underset{\bar{c} \in \bigcup_{i \in \{1, \ldots, M\}} \left\{\bar{c}_j^{(i)}\right\}_{j=1}^{n_i}}{\mathsf{Top}k} \mathsf{sim}(x, \bar{c}) \quad (2)$$

The agents corresponding to the centroids in the set $\{f(\bar{c}_j)\}_{j=1}^k$ are then invited to answer the query.

Therefore, we define a routing-then-answer function **RTANS(QUERY)**, which the server calls to proceed through the following steps to generate the final response and ensure the response is accurate and well-supported to the input query:

• *Step 1: Agent Selection.* The server selects the most relevant agents based on the similarity between the query and the centroids, as described.

• *Step 2: Response Generation.* The selected RAG-based agents generate responses that include both evidence and an answer. The evidence provides evidence, such as a supporting text passage, to justify the answer.

• *Step 3: Evaluation of Responses.* The server collects the responses and evaluates them based on the quality of the analysis provided. It categorizes the answers as "Addressed" or "Not Addressed."

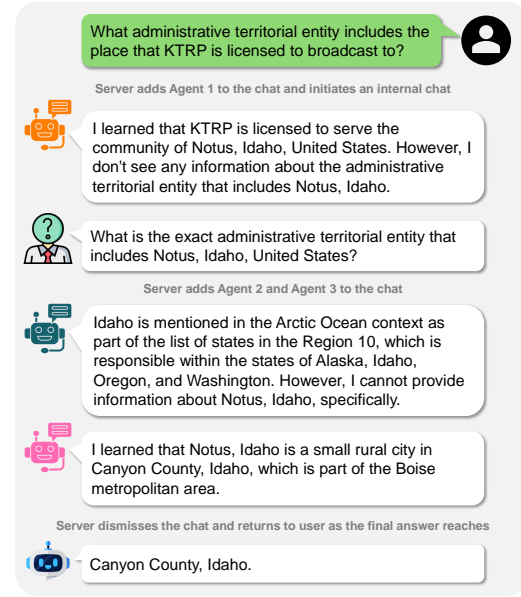• *Step 4: Final Answer Curation.* The server utilizes the "Addressed" answers and finalizes the response to the user.



Figure 3: An example of Answering a Multi-hop Query

## 4.2 Iterative Refining-solving

A multihop question is one that inherently requires multiple reasoning steps, with each step drawing on distinct pieces of supporting knowledge (Kwiatkowski et al., 2019; Ho et al., 2020; Trivedi et al., 2022b; Tang and Yang, 2024; Welbl et al., 2018). In our setting, such queries are especially challenging because the required information can be distributed among different agents. The single-step routing process, which relies solely on an initial similarity match to select agents, may be inadequate for handling multihop questions because it cannot effectively identify or integrate the sequential pieces of evidence required for a complete answer. Thus, a more robust solution is needed to address complex and often cross-domain queries.

**Iterative Routing.** Previous works (Zhuang et al., 2024; Yang et al., 2024b; Press et al., 2022; Trivedi et al., 2022a; Ma et al., 2023; Shao et al., 2023; Jiang et al., 2023) introduce a simple yet effective solution to handle a multi-hop question by repeatedly generating new queries based on the extracted knowledge. Inspired by that, we propose an iterative routing process to manage multi-hop questions, as depicted in Figure 3. In this approach, the server iteratively invokes a routing-then-answer function, RTANS(QUERY), which progressively refines the query and synthesizes supporting evidence until a comprehensive answer is reached.

**Advantages.** This proposed method appears to have twofold advantages from the perspective of *efficiency* and *adaptiveness*. For multi-hop questions

that require sequential reasoning, our approach dynamically refines the query based on the acquired knowledge from agents. Compared to those question decomposition methods (Zhou et al., 2022; Verma et al., 2024; Chan et al., 2024), the proposed RIRS reduces unnecessary query rounds and allows an agent to address multiple reasoning steps within its expertise. For those questions requiring parallel reasoning, particularly across diverse domains, the server enables multiple specialized agents to work concurrently, each leveraging its smaller, more efficient knowledge base, which significantly accelerates the overall inference process.

## 5 Experiments

### 5.1 Experimental Setup

**Datasets.** Our experiments cover both single-hop and multi-hop open-ended QA tasks within a unified evaluation framework. We employ the Natural Questions (NQ) dataset (Kwiatkowski et al., 2019) for single-hop QA to emphasize the effectiveness of the routing mechanism. In the multi-hop QA setting, we assess performance on four benchmark datasets: HotpotQA (Yang et al., 2018), 2WikiMQA (Ho et al., 2020), MusiQue (Trivedi et al., 2022b), and Multi-Hop RAG (MHR) Benchmark (Tang and Yang, 2024) to show the performance of the complete RIRS.

**Models.** Our experimental setup employs a suite of models to support both retrieval and generation tasks. Specifically, we use text-embedding-v2 (Zhang et al., 2024) as our embedding model to generate dense representations for effective text retrieval. The agents are powered by two large language models, i.e., llama-3.1-8b-instruct (Touvron et al., 2023a,b; Dubey et al., 2024) and qwen-plus-2024-12-20 (Bai et al., 2023; Yang et al., 2024a). To evaluate the open-ended QA tasks, we make use of both qwen-max-0125 and gpt-4o-2024-08-06 (Achiam et al., 2023) by comparing the generated responses against designated groundtruth answers.

**Baselines.** Throughout the experiments, we compare the proposed RIRS with the following baselines. RankRAG and EfficientRAG operate in a single RAG-based agent scenario, where all knowledge is managed within one agent. In contrast, Chameleon, RouterDC, and GoldRouter focus on the multi-agent setting, serving as alternative routing strategies. The detailed descriptions of each

baseline are deferred to Appendix B.2.

- *RankRAG (Yu et al., 2024):* This method retrieves documents using dense and sparse retrievers, then ranks them based on helpfulness before generating a response. If no relevant documents are found, a LLM is used to answer independently.
- *EfficientRAG (Zhuang et al., 2024):* This method iteratively simplifies the query by removing resolved components, enabling a more targeted retrieval process.
- *Chameleon (Lu et al., 2024):* This method acts as a routing mechanism by means of an LLM to select relevant agents based on their specialized topics and the given query.
- *RouterDC (Chen et al., 2024a):* This method selects agents by computing similarity between query embeddings and precomputed agent representations. We adapt it using historical queries to approximate each agent's knowledge capacity.
- *GoldRouter:* This method serves as an upper bound by eliminating routing uncertainty. The router has prior knowledge of the optimal agent(s) for each query, ensuring the most appropriate selection without error.

**Multi-Agent Settings.** Chameleon, RouterDC, and RIRS are in multi-agent setting. To simulate the practical knowledge domain segmentation, we construct the two groups of RAG-base agents for different datasets as follows.
- *WikiAgents for NQ, HotpotQA, 2WikiMQA and MusiQue:* WikiAgents group is built upon a corpus of over 121K Wikipedia pages, dumped as of November 1, 2023, and made publicly available via the HuggingFace dataset. The system comprises exactly 64 RAG-based agents. This specific number is derived from the inherent limitation of the ORES legacy service, which can only classify a Wikipedia page into 64 predefined categories. Consequently, each RAG-based agent is designated to handle one of these 64 categories, ensuring that the categorization of pages is consistent and aligned with the predefined taxonomy established by the ORES service. However, these agents cannot cover all required documents, and once the question cannot be answered, a knowledgeable agent will be called to answer the question because the existing LLMs have been pretrained with Wikipedia corpus.
- *NewsAgents for MHR:* NewsAgents group is constructed on a corpus of 609 news articles from 49 distinct news media and spanning six domains (Tang and Yang, 2024). To analyze the effects of

6

| Methods | Models | Natural Questions | | | HotpotQA | | | | 2WikiMultiHopQA | | | | MuSiQue | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | EM | Acc. | Time | EM | Acc. | Rounds | Time | EM | Acc. | Rounds | Time | EM | Acc. | Rounds | Time |
| *Without RAG* | | | | | | | | | | | | | | | | |
| CoT | Qwen-Plus | 65.90 | 80.51 | 3.17 | 34.41 | 59.12 | 1.0 | 2.99 | 23.09 | 49.85 | 1.0 | 3.10 | 22.13 | 35.12 | 1.0 | 3.00 |
| | LLaMA-3.1-8B | 61.17 | 74.87 | 0.89 | 43.18 | 53.17 | 1.0 | 0.88 | 47.51 | 35.71 | 1.0 | 0.86 | 17.23 | 30.70 | 1.0 | 0.87 |
| *Single RAG-based Agent* | | | | | | | | | | | | | | | | |
| RankRAG | Qwen-Plus | 67.52 | 88.01 | 12.96 | 43.23 | 61.69 | 1.0 | 13.28 | 47.51 | 53.61 | 1.0 | 13.31 | 24.08 | 43.07 | 1.0 | 13.19 |
| | LLaMA-3.1-8B | 61.80 | 78.62 | 7.00 | 43.12 | 61.27 | 1.0 | 9.31 | 37.51 | 44.17 | 1.0 | 9.77 | 19.30 | 35.65 | 1.0 | 10.44 |
| EfficientRAG | Qwen-Plus | 67.52 | 88.01 | 12.96 | 55.72 | 72.79 | 1.45 | 19.48 | 58.93 | 64.79 | 1.42 | 21.30 | 27.09 | 45.89 | 1.82 | 23.89 |
| | LLaMA-3.1-8B | 61.80 | 78.62 | 7.00 | 44.98 | 63.32 | 1.68 | 17.01 | 46.64 | 56.52 | 1.66 | 16.64 | 21.68 | 37.67 | 2.01 | 17.39 |
| *Multiple RAG-based Agents* | | | | | | | | | | | | | | | | |
| GoldRouter[1] | Qwen-Plus | 51.29 | 86.78 | 17.46 | 57.78 | 77.48 | 1.24 | 25.53 | 70.71 | 76.33 | 1.36 | 25.16 | 29.76 | 46.96 | 1.48 | 31.89 |
| | LLaMA-3.1-8B | 55.04 | 76.57 | 5.50 | 48.39 | 62.23 | 1.54 | 10.57 | 50.05 | 52.49 | 1.58 | 12.24 | 20.12 | 43.01 | 1.84 | 14.42 |
| Chameleon | Qwen-Plus | 50.47 | 81.26 | 17.37 | 50.70 | 70.78 | 1.47 | 27.95 | 60.69 | 65.93 | 1.70 | 31.57 | 23.90 | 39.38 | 1.72 | 36.69 |
| | LLaMA-3.1-8B | 38.01 | 75.00 | 7.38 | 39.96 | 54.34 | 1.81 | 23.68 | 41.61 | 46.45 | 1.90 | 24.58 | 14.24 | 33.75 | 2.01 | 27.45 |
| RouterDC | Qwen-Plus | 49.62 | 80.53 | 12.91 | 44.02 | 62.38 | 1.73 | 33.62 | 49.64 | 57.20 | 1.69 | 35.58 | 21.06 | 35.31 | 1.69 | 33.59 |
| | LLaMA-3.1-8B | 36.10 | 75.32 | 8.79 | 36.37 | 44.31 | 1.78 | 17.38 | 35.60 | 34.49 | 1.70 | 18.13 | 14.63 | 20.05 | 1.94 | 20.24 |
| RIRS | Qwen-Plus | 54.61 | 81.56 | 11.88 | 53.46 | 74.62 | 1.46 | 28.93 | 62.42 | 66.80 | 1.70 | 31.25 | 24.86 | 43.48 | 1.71 | 28.03 |
| | LLaMA-3.1-8B | 44.62 | 75.57 | 8.29 | 47.83 | 60.75 | 2.14 | 18.40 | 46.20 | 47.06 | 2.26 | 20.45 | 21.65 | 40.46 | 2.36 | 22.20 |
| | Mixed[2] | 53.10 | 80.15 | 10.89 | 52.04 | 74.14 | 1.84 | 24.74 | 53.92 | 65.01 | 1.92 | 30.91 | 23.46 | 43.20 | 2.06 | 24.50 |

[1] The inference is completely based on the selected agents.
[2] The RAG-based agents use LLaMA-3.1-8B, while other modules in the server use Qwen-Plus.

Table 1: Performance comparison of different methods under various datasets and the knowledge of WikiAgents.

source and domain characteristics, we further split NewsAgent into two variants: *NewsAgent-Source* and *NewsAgent-Domain*, corresponding to partitions based on news media and domain categories, respectively.

## 5.2 Quantative Analysis with WikiAgents

Table 1 presents the performance of our proposed RIRS alongside various baselines on four Wikipedia-related QA tasks using the knowledge from the WikiAgents. Below, we highlight key comparisons and insights based on the acquired experimental results.

**Comparison with GoldRouter.** In the GoldRouter, the system benefits from comprehensive prior knowledge about which agents are best suited for different aspects of a question, effectively serving as an upper bound for multi-agent methods. While GoldRouter can accurately route queries to the optimal agents, our approach (RIRS) occasionally encounters routing errors, which lead to additional query rounds and minor performance degradation compared to this ideal scenario. These routing errors highlight the challenges in dynamically estimating each agent's expertise on-the-fly, yet the overall performance remains competitive even with these extra iterations.

**Comparison with the scenario of Single RAG-based agent.** Single RAG-based methods, such as RankRAG and EfficientRAG, consolidate all knowledge into one unified base, allowing them to review a comprehensive document set for each query. In expectation, RankRAG should perform best under a single-hop QA task (i.e., Natural Questions), while EfficientRAG can generate the most accurate responses to multihop queries. However, we observe that both methods sometimes underperform compared to multi-agent approaches like RIRS. The reason is that partitioning the knowledge into domain-specific agents can limit irrelevant or distracting content during retrieval, enabling each agent to focus on a smaller, more relevant subset of documents. Additionally, the decentralized nature of our approach allows for the review of more knowledge chunks without being restricted by input token limits, thereby improving inference focus and overall performance.

**Comparison with other routing methods.** When compared to other routing strategies such as RouterDC and Chameleon, RIRS achieves higher accuracy by providing a more reliable outline of each agent's knowledge capacity. RouterDC, which relies on caching 100 historical questions per agent, often falls short in representing the full spectrum of an agent's expertise, while Chameleon's reliance on static textual descriptions can lead to misrouting. In contrast, our iterative routing mechanism dynamically refines query assignments based on actual performance feedback, resulting in more precise and robust routing that better aligns with the agents' strengths.

**Comparison with various LLMs.** Focusing on our methods, we observe that larger LLM models like Qwen-Plus consistently achieve better performance than smaller models such as LLaMA-

| Methods | Models | NewsAgent-Source | | | NewsAgent-Domain | | |
|---|---|---|---|---|---|---|---|
| | | Acc. | Hall.[1] | Rounds | Acc. | Hall.[1] | Rounds |
| *Single RAG-based Agent* [2] | | | | | | | |
| RankRAG | Qwen-Plus | 78.52 | 3.37 | 1.0 | - | - | - |
| | LLaMA-3.1-8B | 65.77 | 8.71 | 1.0 | - | - | - |
| EfficientRAG | Qwen-Plus | 81.80 | 1.19 | 1.57 | - | - | - |
| | LLaMA-3.1-8B | 66.80 | 6.05 | 1.32 | - | - | - |
| *Multiple RAG-based Agents* | | | | | | | |
| GoldRouter | Qwen-Plus | 89.32 | - | 1.06 | 90.31 | - | 1.12 |
| | LLaMA-3.1-8B | 79.53 | - | 1.08 | 80.16 | - | 1.18 |
| RIRS | Qwen-Plus | 84.75 | 0.93 | 1.65 | 88.78 | 1.07 | 1.66 |
| | LLaMA-3.1-8B | 71.91 | 4.91 | 1.37 | 75.12 | 4.73 | 1.39 |

[1] Hall. means a hallucination rate that misclassifies a null query as answerable, which should be the lower, the better.

[2] Both NewsAgent-Source and NewsAgent-Domain share the same knowledge repository for the single RAG-based agent, thus producing identical results. We therefore report these single-agent outcomes under NewsAgent-Source only, leaving the corresponding entries for NewsAgent-Domain blank to avoid duplication.

Table 2: Performance comparison of different methods using NewsAgent-Source and NewsAgent-Domain.

3.1-8B. Notably, when the server is equipped with Qwen-Plus while the agents operate with lightweight models, the overall performance is nearly on par with a scenario where all agents use Qwen-Plus. This gap underscores the limitations of LLaMA-3.1-8B in following complex instructions, but also highlights a practical advantage: agents can run lightweight models locally to preserve efficiency and knowledge sovereignty, while a more powerful central model can manage complex reasoning tasks without exposing sensitive information to third-party providers.

## 5.3 Quantitative Analysis for Multi-hop RAG

In this section, we analyze the performance of our proposed method, RIRS, on the MHR benchmark under two NewsAgents settings: **NewsAgent-Source** (with 49 agents) and **NewsAgent-Domain** (with six agents). As shown in Table 2 and Figure 4, we focus on two key perspectives: (i) comparing RIRS with single RAG-based agents, and (ii) discussing the effect of varying the number of RAG-based agents.

**Comparison with Single RAG-based Agent.** In the single-agent scenario, all knowledge is consolidated into a single agent (e.g., RankRAG and EfficientRAG). While this setup can sometimes simplify retrieval by reviewing a comprehensive document pool, it also risks introducing irrelevant or distracting information. By contrast, RIRS partitions knowledge across multiple specialized agents and iteratively refines query routing, thereby reducing the likelihood of retrieving spurious content. As Table 2 indicates, RIRS achieves not only higher accuracy but also a notably lower hallucination rate compared to the single-agent methods. In particular, when dealing with multi-hop queries,
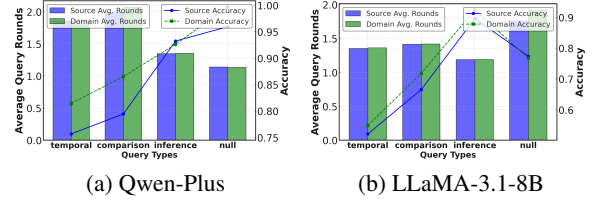


Figure 4: Performance Comparison of Different Models for MHR Benchmarks under Different Query Types. (Zoom in for the best view)

the smaller domain-specific knowledge bases (or source-specific segments) mitigate confusion and enhance the reliability of retrieved evidence, helping RIRS avoid erroneous or fabricated answers.

**Discussion of Different Numbers of RAG-based Agents.** From Table 2 and Figure 4, we observe that having a large number of agents (49) can offer very fine-grained coverage but may incur additional query rounds due to routing overhead; in contrast, having fewer agents reduces the routing complexity, sometimes resulting in fewer rounds while maintaining high accuracy. Notably, the domain-based approach (6 agents) appears to strike a more balanced trade-off between specialized coverage and routing overhead, often leading to efficient query resolution for multi-hop questions.

These findings suggest that there exists an *optimal number* of agents for a given knowledge base: too many agents risk increased routing complexity and query overhead, whereas too few agents risk merging domains too broadly, which can reintroduce the problem of irrelevant knowledge retrieval. Thus, system designers must weigh the benefits of granular specialization against the costs of additional query rounds when determining the appropriate level of knowledge partitioning.

## 6 Conclusion

In this work, we introduce RIRS, a novel framework that coordinates multiple RAG-based agents in response to a query while preserving their knowledge sovereignty. Specifically, this framework consists of two mechanisms: the routing mechanism directs a user query to the most appropriate agents, and the iterative refining-solving mechanism enhances the system's ability to tackle complex, multi-hop queries by progressively synthesizing intermediate responses into a comprehensive final answer. Extensive experiments using Wikipedia-related and News-related corpus and datasets demonstrate the effectiveness of the proposed method, regardless of the complexities of the given questions.

8

## Limitations

Despite the promising results, RIRS has two limitations that warrant further investigation. First, the effectiveness of the routing mechanism heavily relies on the quality of the knowledge boundary representations derived from embedding clusters; in cases where knowledge domains overlap significantly, or embeddings are less distinct, the router may misidentify relevant agents. Second, this work is unable to handle a multimodality scenario, where the provided corpus contains a modality other than texts.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774.*

Parker Addison, Minh-Tuan H Nguyen, Tomislav Medan, Mohammad T Manzari, Brendan McElrone, Laksh Lalwani, Aboli More, Smita Sharma, Holger R Roth, Isaac Yang, et al. 2024. C-fedrag: A confidential federated retrieval-augmented generation system. *arXiv preprint arXiv:2412.13163.*

Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. 2023a. Retrieval-based language models and applications. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 6: Tutorial Abstracts)*, pages 41–46.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023b. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511.*

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609.*

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.

Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. 2024. Rq-rag: Learning to refine queries for retrieval augmented generation. *arXiv preprint arXiv:2404.00610.*

Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879.

Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F Karlsson, Jie Fu, and Yemin Shi. 2023a. Autoagents: A framework for automatic agent generation. *arXiv preprint arXiv:2309.17288.*

Shuhao Chen, Weisen Jiang, Baijiong Lin, James T Kwok, and Yu Zhang. 2024a. Routerdc: Query-based router by dual contrastive learning for assembling large language models. *arXiv preprint arXiv:2409.19886.*

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. 2023b. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*.

Zehui Chen, Kuikun Liu, Qiuchen Wang, Jiangning Liu, Wenwei Zhang, Kai Chen, and Feng Zhao. 2024b. Mindsearch: Mimicking human minds elicits deep ai searcher. *arXiv preprint arXiv:2407.20183.*

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.*

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*, pages 6491–6501.

Dawei Gao, Zitao Li, Xuchen Pan, Weirui Kuang, Zhijian Ma, Bingchen Qian, Fei Wei, Wenhao Zhang, Yuexiang Xie, Daoyuan Chen, et al. 2024. Agentscope: A flexible yet robust multi-agent platform. *arXiv preprint arXiv:2402.14034.*

Michael Glass, Gaetano Rossiello, Md Faisal Mahbub Chowdhury, Ankita Rajaram Naik, Pengshan Cai, and Alfio Gliozzo. 2022. Re2g: Retrieve, rerank, generate. *arXiv preprint arXiv:2207.06300.*

Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680.*

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060.*

Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*.

Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*.

Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*.

Ashutosh Joshi, Sheikh Muhammad Sarwar, Samarth Varshney, Sreyashi Nag, Shrivats Agrawal, and Juhi Naik. 2024. Reaper: Reasoning based retrieval planning for complex rag systems. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 4621–4628.

Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.

Jungo Kasai, Keisuke Sakaguchi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, Kentaro Inui, et al. 2024. Realtime qa: what's the answer right now? *Advances in Neural Information Processing Systems*, 36.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.

Andrew K Lampinen, Ishita Dasgupta, Stephanie CY Chan, Kory Matthewson, Michael Henry Tessler, Antonia Creswell, James L McClelland, Jane X Wang, and Felix Hill. 2022. Can language models learn from explanations in context? *arXiv preprint arXiv:2204.02329*.

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*, 33:9459–9474.

Ao Li, Yuexiang Xie, Songze Li, Fugee Tsung, Bolin Ding, and Yaliang Li. 2024. Agent-oriented planning in multi-agent systems. *arXiv preprint arXiv:2410.02189*.

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. Camel: Communicative agents for" mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008.

Xiaonan Li, Changtai Zhu, Linyang Li, Zhangyue Yin, Tianxiang Sun, and Xipeng Qiu. 2023b. Llatrieval: Llm-verified retrieval for verifiable generation. *arXiv preprint arXiv:2311.07838*.

Jerry Liu. 2022. LlamaIndex.

Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2023. Routing to the expert: Efficient reward-guided ensemble of large language models. *arXiv preprint arXiv:2311.08692*.

Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2024. Chameleon: Plug-and-play compositional reasoning with large language models. *Advances in Neural Information Processing Systems*, 36.

Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting for retrieval-augmented large language models. *arXiv preprint arXiv:2305.14283*.

Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511*.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.

Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2020. Document ranking with a pretrained sequence-to-sequence model. *arXiv preprint arXiv:2003.06713*.

Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*.

Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*, 11:1316–1331.

Stephen Robertson, Hugo Zaragoza, and Michael Taylor. 2004. Simple bm25 extension to multiple weighted fields. In *Proceedings of the thirteenth ACM international conference on Information and knowledge management*, pages 42–49.

Aymeric Roucher. 2024. Agentic RAG: tur-bocharge your RAG with query reformulation and self-query. https://huggingface.co/learn/cookbook/en/agent_rag. Accessed: 2025-02-15.

Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. *arXiv preprint arXiv:2305.15294*.

Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652*.

Tal Shnitzer, Anthony Ou, Mírian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson, and Mikhail Yurochkin. 2023. Large language model routing with benchmark datasets. *arXiv preprint arXiv:2309.15789*.

EuiYul Song, Sangryul Kim, Haeju Lee, Joonkee Kim, and James Thorne. 2024. Re3val: Reinforced and reranked generative retrieval. *arXiv preprint arXiv:2401.16979*.

KV Srivatsa, Kaushal Kumar Maurya, and Ekaterina Kochmar. 2024. Harnessing the power of multiple minds: Lessons learned from llm routing. *arXiv preprint arXiv:2405.00467*.

Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. Dragin: Dynamic retrieval augmented generation based on the real-time information needs of large language models. *arXiv preprint arXiv:2403.10081*.

Yixuan Tang and Yi Yang. 2024. Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries. *arXiv preprint arXiv:2401.15391*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022a. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022b. Musique: Multi-hop questions via single-hop question composition.

*Transactions of the Association for Computational Linguistics*, 10:539–554.

Prakhar Verma, Sukruta Prakash Midigeshi, Gaurav Sinha, Arno Solin, Nagarajan Natarajan, and Amit Sharma. 2024. Plan × rag: Planning-guided retrieval augmented generation. *arXiv preprint arXiv:2410.20753*.

Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Constructing datasets for multi-hop reading comprehension across documents. *Transactions of the Association for Computational Linguistics*, 6:287–302.

Lilian Weng. 2023. LLM powered autonomous agents. https://lilianweng.github.io/posts/2023-06-23-agent/. Accessed: 2025-01-09.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.

Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can Chen, Ye Yuan, Lianming Huang, Xue Liu, Tei-Wei Kuo, Nan Guan, et al. 2024. Retrieval-augmented generation for natural language processing: A survey. *arXiv preprint arXiv:2407.13193*.

Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. *arXiv preprint arXiv:2402.13178*.

Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee, Chen Zhu, Zihan Liu, Sandeep Subramanian, Evelina Bakhturina, Mohammad Shoeybi, and Bryan Catanzaro. 2023. Retrieval meets long context large language models. *arXiv preprint arXiv:2310.03025*.

Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective retrieval augmented generation. *arXiv preprint arXiv:2401.15884*.

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

Diji Yang, Jinmeng Rao, Kezhen Chen, Xiaoyuan Guo, Yawen Zhang, Jie Yang, and Yi Zhang. 2024b. Im-rag: Multi-round retrieval-augmented generation through learning inner monologues. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 730–740.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

11

Wenhao Yu. 2022. Retrieval-augmented generation across heterogeneous knowledge. In *Proceedings of the 2022 conference of the North American chapter of the association for computational linguistics: human language technologies: student research workshop*, pages 52–58.

Yue Yu, Wei Ping, Zihan Liu, Boxin Wang, Jiaxuan You, Chao Zhang, Mohammad Shoeybi, and Bryan Catanzaro. 2024. Rankrag: Unifying context ranking with retrieval-augmented generation in llms. *arXiv preprint arXiv:2407.02485*.

Xin Zhang, Yanzhao Zhang, Dingkun Long, Wen Xie, Ziqi Dai, Jialong Tang, Huan Lin, Baosong Yang, Pengjun Xie, Fei Huang, et al. 2024. mgte: Generalized long-context text representation and reranking models for multilingual text retrieval. *arXiv preprint arXiv:2407.19669*.

Qinlin Zhao, Jindong Wang, Yixuan Zhang, Yiqiao Jin, Kaijie Zhu, Hao Chen, and Xing Xie. 2023. Competeai: Understanding the competition behaviors in large language model-based agents. *arXiv preprint arXiv:2310.17512*.

Ziyu Zhao, Leilei Gan, Guoyin Wang, Wangchunshu Zhou, Hongxia Yang, Kun Kuang, and Fei Wu. 2024. Loraretriever: Input-aware lora retrieval and composition for mixed tasks in the wild. *arXiv preprint arXiv:2402.09997*.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.

Ziyuan Zhuang, Zhiyang Zhang, Sitao Cheng, Fangkai Yang, Jia Liu, Shujian Huang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. 2024. Efficientrag: Efficient retriever for multi-hop question answering. *arXiv preprint arXiv:2408.04259*.

12

# A Key Modules on the Server

The proposed RIRS is comprised of two key components: multiple RAG-based agents, each specialized in a domain of expertise, and a central server that coordinates their collaborative efforts. The server serves as the nexus of our framework by orchestrating the interaction between agents through a series of specialized roles designed to ensure that user queries are addressed with both logical rigor and relevant supporting evidence. Below, we detail the primary agents managed by the server and their corresponding responsibilities.

## A.1 Question Evaluator

The Question Evaluator is the first checkpoint in the server's processing pipeline. Given that the server itself lacks domain-specific background knowledge, the evaluator assesses agent responses using a common-sense, logic-based approach. Its primary tasks are:

- **Logical Assessment:** Evaluating whether the provided response is coherent and free from logical fallacies.

- **Evidence Verification:** Confirming that necessary citations and supporting evidence are included in the response.

- **Response Classification:** Determining if the response is "addressed" or "not addressed." A response is marked "not addressed" if it fails to provide sufficient evidence, contains logical deficiencies, or is entirely irrelevant to the user query.

- **Granular Evaluation for Multihop Questions:** For multistep queries, even a response classified as "addressed" is further divided into "fully addressed" or "partially addressed." When a response is rated as "partially addressed," the unaddressed part is expected to be answered by other agents and/or in the further reasoning steps.

This preliminary evaluation ensures that only logically sound and evidence-backed responses are propagated in the subsequent stages.

## A.2 Response Summarizer

Following the evaluation stage, the Response Summarizer plays a crucial role in consolidating agent responses. Depending on the evaluator's results, two variants of the summarizer are deployed:

- **Fully Addressed Summarizer:** This variant consolidates one or more responses that have been deemed to fully address the query. It integrates the responses into a comprehensive answer which is then forwarded directly to the user.

- **Partially Addressed Summarizer:** When none of the responses can fully address the query, the partially addressed summarizer steps in to compile a more complete solution based on the available partial responses. Once consolidated, it re-evaluates the answer. If the unified response is deemed fully addressed, it is returned to the user; otherwise, further action is initiated.

Notably, although the partially addressed summarizer can directly consolidate the information from all responses, the overall performance of response summarizing benefits from the initial filtering conducted by the Question Evaluator and fine-granularity of the standard collaboration flow.

## A.3 Question Simplifier

When the query remains insufficiently addressed, as rated by the partially addressed summarizer, the Question Simplifier intervenes to decompose the problem into more manageable sub-questions. This agent performs the following tasks:

- **Identification of Addressed Components:** It examines the existing responses to isolate the aspects of the query that have already been effectively addressed.

13

- **Generation of a New Question:** By removing the resolved parts, the simplifier formulates a new query targeting the unresolved components.

- **On-the-Fly Decomposition:** This dynamic simplification enables the server to continue resolving the query in an iterative, step-by-step manner without the need for prior, rigid question planning.

This adaptive approach leverages already acquired knowledge, ensuring that subsequent agent interactions are focused on the remaining aspects of the problem, thereby reducing the need for redundant multistep reasoning in later stages.

### A.4 Discussion: Unhandleable Queries

In some instances, the server may fail to obtain any reliable ("addressed") responses. This scenario may arise due to several factors:

- **Query Outside the System's Domain:** The user may pose a question that falls outside the scope of all RAG-based agents' expertise. For instance, in a medical suggestion QA system, a travel recommendation query would not align with the agents' specialized knowledge, resulting in no suitable answer.

- **Ambiguous or Incomplete Queries:** A query that is vague, under-specified, or contains numerous typographical errors can hinder the server's ability to correctly map the question to the appropriate agents. For example, in a medical context, ambiguous terminology or poorly structured queries may impede the identification of a clear problem statement, leading to an inability to retrieve a fully addressed response.

- **Rapidly Evolving Information Domains:** In areas where information is rapidly changing, some RAG-based agents may not have the most current data or guidelines. This lag can result in responses that are either outdated or insufficient, prompting the system to classify the query as out-of-scope.

In such cases, the system will inform the user that the question cannot be answered based on the current scope of the available agents. However, if the server integrates a knowledgeable agent with broader capabilities, this agent may be employed to attempt an answer. For multihop questions, the knowledgeable agent is provided with a simplified version of the query, since evidence suggests that large language models perform better when fewer reasoning steps are required (Zhuang et al., 2024).

## B Implementations and Baselines

### B.1 Implementation Details

Our implementation builds upon the open-source AgentScope project (Gao et al., 2024), and each agent's knowledge base is processed by LLaMAIndex (Liu, 2022) and stored in ElasticSearch. We leverage ElasticSearch's Mixed Retrieval for knowledge retrieval, where each agent extracts 20 chunks and selects the best five for answer generation. Without special annotations, five agents are chosen per query round. The entire system is deployed within a pseudo-distributed environment that utilizes multi-threading to enhance scalability and efficiency. Our code and the setup of multiagent systems (including corpus) will be released upon acceptance. Due to the limited space, more experimental setups (e.g., prompts, metrics, and baseline implementations) and some experiential results are deferred to the appendix.

### B.2 Baselines

In order to comprehensively evaluate our proposed system, we have reproduced several baselines inspired by existing works. These baselines are implemented manually to fit within our experimental framework. Notably, our study focuses on an off-the-shelf scenario, which does not fine-tune or train any models and instead leverages the capabilities of existing LLMs to achieve desired effects. Toward the goal, the implementation details of the baselines are given as follows:

**Single-agent Scenarios.** In the single-agent scenario, the knowledge contained within all RAG-based agents is merged into a single agent. This is a unified setup for conventional RAG-based methods. In this setting, baseline methods can review all documents and retrieve the most relevant ones within the system. Although the single-agent setup does not fully align with the scenario we aim to examine, we include it to demonstrate the challenges associated with managing a large knowledge base, i.e., longer retrieval times and potential distractions from plausible yet irrelevant information. Ideally, if the retrieved knowledge were perfectly clean, the single-agent setup could serve as an upper bound of the multi-agent setup in terms of accuracy when comparing the generated results against the ground truth.

- **RankRAG:** This method retrieves some documents from the knowledge base using both dense and sparse retrievers. Next, a pretrained model is introduced to evaluate the helpfulness of each retrieved document and select the most appropriate document(s) to answer the given query. Based on the retrieved document(s), the LLM is asked to generate a response. If the provided document(s) are irrelevant to the question, the LLM is supposed to generate the answer on its own ability. Therefore, this method maintains a single query round for all types of questions. This baseline method covers a number of the existing works (Yu et al., 2024; Glass et al., 2022; Song et al., 2024; Ram et al., 2023; Ma et al., 2023; Nogueira et al., 2020), which focuses on using reranking to enhance LLM content generation, while they use different ways to train the ranking model.

- **EfficientRAG:** This approach iteratively simplifies the query by retrieving query-related documents from the knowledge base until it can be fully addressed. The simplification process depends entirely on the LLM's ability to remove portions of the query that have already been resolved. In cases where the remaining question cannot be further simplified or adequately answered using the provided context, a knowledgeable agent is used to generate the final answer. This baseline follows the EfficientRAG framework (Zhuang et al., 2024), which generates training data from an LLM to train a compact model for question simplification, and is further inspired by related works such as Self-ASK (Press et al., 2022), SelfRAG (Asai et al., 2023b), and IM-RAG (Yang et al., 2024b). This approach has proven effective for multi-hop questions due to its adoption of multiple reasoning steps to converge on a final answer.

**Multi-agent Scenarios.** In this scenario, we implement two routing strategies for comparison with our proposed routing mechanism, RIRS, while keeping the other server modules unchanged and still employing iterative routing. These two methods are derived from settings that involve multiple LLMs and use a router to identify the best LLM for handling a given task. Inspired by these approaches, we extend their ideas to our scenario to construct an effective routing mechanism. In our experiments, these two routing strategy are used to handle Wikipedia-related QA tasks, i.e., Natural Questions, HotpotQA, 2WikiMultiHopQA, and MuSiQue.

- **Chameleon:** This method (Lu et al., 2024) leverages a collection of tools, including LLMs and off-the-shelf vision models, to accomplish complex reasoning tasks step by step, selecting the best tool for each step. The router, which is based on a well-trained LLM, decomposes a complex task into multiple steps and identifies the most suitable tool for each. In our adaptation, each WikiAgent is dedicated to a specific topic from Wikipedia. By providing the router with a description of each agent's specialized topic, it can select up to five agents whose expertise best aligns with the given query.

- **RouterDC:** This method (Chen et al., 2024a) utilizes several LLMs by training a representation vector for each model based on a collection of questions and the optimal candidate from a pool of models. During inference, the router determines the best agents by computing the similarity between the query embedding and the stored representation vectors.

  Inspired by this approach, we calculate the similarity between the query embedding and the representation vectors of our RAG-based agents. However, obtaining these vectors typically requires additional training, which is not compatible with our training-free setting. Instead, we leverage insights from recent works (Lampinen et al., 2022; Mishra et al., 2021) and cache 100 historical questions for each

15

agent to represent their knowledge capacity. For a new question, the router computes the average similarity between the query and the cached questions for each agent, then selects the five agents with the highest similarity scores to generate an answer.

### B.3 Evaluation Metrics

In this section, we provide the details of the most common evaluation metrics as follows:

- **Lexical Match:** This metric measures the percentage of questions where the groundtruth answer appears within the generated response. Since the LLM-generated answers may be longer than the groundtruth, we focus on whether the groundtruth is included in the predicted answer rather than requiring an exact match.

- **GPT Evaluation:** We adopt *gpt-4o-2024-08-06* (Achiam et al., 2023) to evaluate the correctness of the generated responses. This metric captures cases where the generated response conveys the same meaning as the groundtruth, even if the wording is different.

- **Cost:** We calculate the total token consumption for each query across all agents, measuring the computational cost associated with each query.

- **Time:** We compute the wall-clock time from the question that appears to be a valid response. It is noted that the time may not be accurately measured because of the existence of network or threading congestion, especially if an API call is required.

### B.4 Data Distributions for Multi-hop Question

In this section, we show the data distribution for MuSiQue, 2WikiMultiHopQA (or 2WikiMQA), and HotpotQA across the minimum number of required agents. In our experiments, we sample our dataset to speed up our inference progress, while ensuring at most 3% error within 95% confidence interval.
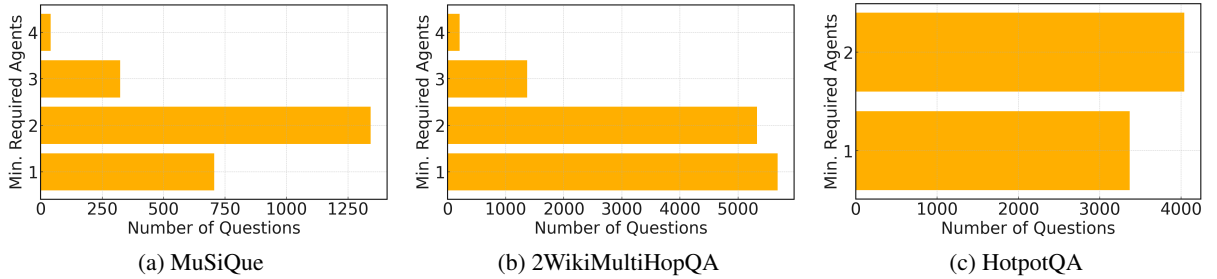


Figure 5: Data distribution v.s. Minimum required agents under different datasets
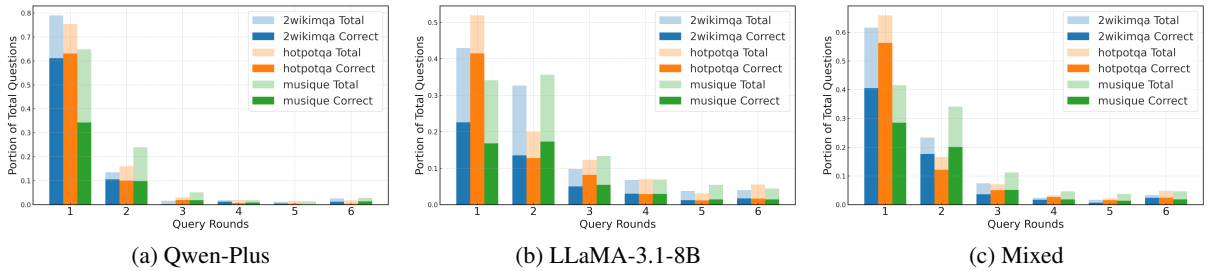
### B.5 More Experimental Results



Figure 6: Evaluation for Wikipedia-related QA (2WikiMQA, HotpotQA, and MuSiQue)

## C  Prompts

**System Prompts for Chain-of-Thought (CoT) Agent**

You're a knowledgeable assistant. You are provided with a question, and you should answer the question in the following two steps. FIRST, you should utilize your knowledge and analyze the question step by step. SECOND, you should finalize an answer based on your analysis with no more than 30 words.

Your output should be in the json format:
```json
{
    "analysis": "<a paragraph with no more than six sentences>",
    "answer": "<a response within 30 words>"
}
```

Figure 7: Prompt Templates for the CoT Agents.

**System Prompts of Wikipedia Agents for A Single Question**

You are provided with one question and a collection of knowledge. Based on the given knowledge, you should try to analyze and tackle the question as thoroughly as possible, even if you cannot fully answer the question. Unless the given context is unrelated to the question, you must directly quote the evidence (i.e., sentences) without being altered to support your analysis, enclosing it in double asterisks (**). You should not state any arguments that are not explicitly mentioned or implied from the pieces of evidence or without quoting them. The analysis should be in one paragraph with no more than ten sentences. Moreover, the analysis should start with "I" and not mention that the analysis is generated based on the given knowledge, documents, or information.

Your output should be in the json format:
```json
{
    "analysis": "<one paragraph of up to ten sentences, directly quoting supporting evidence from
        the provided knowledge>"
}
```

Figure 8: Prompt Templates for the Wikipedia Agents.

17

You are given one question and a collection of news articles. Each article contains content along with its source information, including the title, news source, author, and published time. Your task is to analyze and address the question as thoroughly as possible based on the provided news articles, even if you cannot fully answer the question. Unless the given context is unrelated to the question, you must directly quote the evidence (i.e., sentences) without being altered to support your analysis. You should enclose quoted evidence (sentences) in double asterisks (**), followed by the source in brackets, including the title, news source, author, and published time, separated by semicolons (;). You should not state any arguments that are not explicitly mentioned or implied from the pieces of evidence or without quoting them. Your analysis to each question should be concise, limited to one paragraph per question, with no more than ten sentences. The analysis must begin with "I" and should not mention that the analysis is based on provided news articles, knowledge, or information.

Your output should be in the json format:
```json

{
    "analysis": "<one paragraph of up to ten sentences, directly quoting supporting evidence from
        the provided knowledge>"
}

```

Figure 9: Prompt Templates for the News Agents.

You are provided with a question and a response. Your task is to evaluate the response according to the following steps:
1. Assess the response against three criteria:
- Relevance: Does the response help answer the question, even if the response does not fully resolve it?
- Evidence-Based Support: Are statements supported by explicitly mentioned evidence enclosed in double asterisks (**)?
- Logical Coherence: Is it well-structured, logically reasoned, and free from logical fallacies or contradictions?
2. Assign one of the following ratings:
- Fully addressed: The response meets all criteria and completely answer the question.
- Partially addressed: The response meets all criteria but not fully resolve the question.
- Not addressed: The response fails to meet one or more of the criteria.

Your output should be in the json format:
```json

{
    "evaluation": {
        "relevance": <a sentence within 30 words>,
        "evidence_support": <a sentence within 30 words>,
        "logical_coherence": <a sentence within 30 words>
    },
    "rating": <"Fully addressed" or "Partially addressed" or "Not addressed">
}

```

Figure 10: Prompt Templates for the evaluator.

## System Prompts of Summarizing Fully-addressed Responses

You are given a question, and one or more responses that fully resolve the question. Your task is to produce a final answer by following these steps:
1. Analysis: Incorporate all relevant information from the given responses, quoting any supporting evidence word-for-word in double asterisks (**).
2. Answer: Provide a concise conclusion in no more than 30 words that summarizes the analysis.

Your output should be in the json format:
```json

{
    "analysis": "<a paragraph that directly quotes relevant evidence in **double asterisks**>",
    "answer": "<a concise final answer within 30 words>"
}

```

## System Prompts of Summarizing Partially-addressed Responses

You are given a question and several partially addressed responses. Your task is to combine these responses to create a comprehensive solution, then evaluate its completeness. Follow these steps:
1. Solution Synthesis:
- Incorporate all relevant information from the provided responses to form a solution that addresses the question as thoroughly as possible.
- Directly quote supporting evidence (word-for-word) using double asterisks (**).
2. Evaluation and Justification:
- Determine whether this combined solution fully addresses the question.
- Provide a clear explanation of why it does or does not fully address the question.
3. Answerability Determination:
- Based on your evaluation, decide if the solution makes the question answerable.
- Respond with either "yes" (if the solution fully addresses the question) or "no" (if it does not).
- Most importantly, if the solution explicitly states that the question cannot be fully addressed or identifies missing aspects or necessary additional information, you must answer "no".
4. Final Answer:
- If the solution is deemed answerable ("yes"), provide a concise conclusion in no more than 30 words that summarizes the solution.
- If the solution is not answerable ("no"), return "None" as the final answer.

Your output should be in the json format:
```json

{
    "solution": "<a paragraph that integrates responses with direct quotes in **double asterisks
        **>",
    "evaluation and justification": "<a paragraph explaining whether and why the solution fully
        addresses the question>",
    "answerable": <"yes" or "no">,
    "answer": "<a final answer within 30 words if 'answerable' is 'yes', or 'None' if 'answerable'
        is 'no'>"
}

```

Figure 11: Prompt Templates for the Summarizer.

## System Prompts of Question Simplifier

You are given a question and a piece of knowledge that partially addresses the question. Your task is to simplify or refactor the original question so that answering the simplified question will yield the same final answer as answering the original one. Follow the steps below:

1. Map Known Information:
- Identify Answered Parts: Examine the original question and determine which parts have been explicitly answered by the provided solution. Extract the relevant words or phrases from the question.
- Map to Solution: For each identified part, find the corresponding words or phrases in the provided solution and include supporting evidence by directly quoting the exact sentences, enclosed in double asterisks.
- Note: If the provided knowledge does not address any aspect of the original question (i.e., you cannot identify the answered parts from the question), return "None" for the mapping.

2. Identify Unresolved Aspects:
- Determine Gaps: Identify the parts of the question that remain unanswered or unclear after considering the provided solution.
- List Unresolved Points: Clearly list these unresolved aspects. They must be retained in the simplified question to ensure that no essential detail is lost.

3. Generate a Simplified/Refactored Question:
- Remove Redundancies: Remove from the original question any words or phrases that have been fully addressed by the solution (as determined in Step 1).
- Incorporate Mapped Knowledge: Replace the removed parts with their corresponding mapped words or phrases from the solution, ensuring to preserve the supporting evidence.
- Retain Unresolved Aspects: Ensure that any unresolved aspects identified in Step 2 are still present in the final simplified question.
- Ensure Consistency: The new question must be clear, concise, and structured so that answering it will yield the same final answer as if the original question were fully addressed.
- Note: If the mapping is "None" (as indicated in Step 1), keep the simplified/refactored question identical to the original question.

Your output should be in the json format:
```json

{
    "question_knowledge_mapping": <a list of tuples in the format (words from given question,
        words from provided knowledge, supporting evidence) or "None">,
    "unresolved_aspects": <a list of aspects that remain unanswered>,
    "new_question": <the simplified/refactored question; remove all "words from given question"
        found in the mapping and include the corresponding "words from provided knowledge" to
        replace them, ensuring all unresolved aspects are retained. If mapping is "None", use the
        original question>
}
```

Figure 12: System Prompt for Question Simplifier.

**User Input**

who is darrell brother in the walking dead

**Expected Answer**

Merle Dixon

**Assistant: Culture.Visual arts.Comics and Anime**

The provided documents do not contain any information about Darrell's brother in The Walking Dead. The documents mainly discuss the comic book series, its publication history, and various adaptations, but do not provide details about specific characters like Darrell and his brother.

**Assistant: Culture.Media.Media***

The provided documents contain information about Daryl Dixon, a character from the television series The Walking Dead. According to the document, Daryl has an older brother named Merle Dixon, who is described as "The racist and volatile older brother of Daryl Dixon." The document also mentions that "After Merle's disappearance, Daryl lets his guard down and starts to bond with the group."

**Assistant: History and Society.Politics and government**

The provided documents do not contain any information about a character named Darrell or his brother in the context of "The Walking Dead." The documents discuss various topics unrelated to the TV show or comic series, including Australian rules football families, Charles Darwin, the phrase "Dead Man Walking," and the TV show "Neighbours."

**Assistant: Culture.Media.Television**

The document mentions Daryl's brother, Merle, and their relationship: "Daryl had significant periods of time alone and, throughout these lonely periods, learned to fend for himself and adopted a hard-boiled survivalist mindset. When the outbreak occurs he and Merle fend for themselves and drift around, avoiding walkers." It also states, "The character was introduced in the first season as a southerner, expert tracker, living in the shadow of his older brother, Merle." These quotes indicate that Darrell's brother in The Walking Dead is Merle Dixon.

**Assistant: Culture.Biography.Biography***

The provided documents do not contain any information about a character named Darrell or Darrell's brother in the context of "The Walking Dead."

**Server: Evaluation**

Evaluation:
- Response 1 and 3 are unreliable because they claim there is no information, which is incorrect.
- Response 2 and 4 are reliable as they correctly identify Merle Dixon as Daryl's brother.
- Response 5 is unreliable because it incorrectly claims there is no information.

Analysis: Responses 2 and 4 provide accurate information about Daryl Dixon having an older brother named Merle Dixon in The Walking Dead series.

Answer: Daryl's brother in The Walking Dead is Merle Dixon.

**Response**

Daryl's brother in The Walking Dead is Merle Dixon.

Figure 13: A example of single-hop QA task

Where was the director of film Rough Going born?

New Hyde Park, New York

**Assistant: Culture.Media.Films**

I found that **Rough Going is a 1925 American silent Western film directed by Wally Van**. However, the provided information does not specify where the director, Wally Van, was born.

**Assistant: Culture.Media.Media***

I learned that **Richard Rush (April 15, 1929 – April 8, 2021) was an American film director, scriptwriter, and producer.** This indicates that Richard Rush, the director of Rough Going, was born in the United States. However, the specific place of his birth is not mentioned in the provided context.

**Server: Evaluation**

```
{
    'response_1': {
        'relevance': 'Identifies the director of Rough Going but does not provide his place of
            birth.',
        'evidence_support': 'Provides relevant information about the film and its director, Wally
            Van, with evidence.',
        'logical_coherence': 'The response is logically structured and coherent, though it
            acknowledges incomplete information.',
        'rating': 'Partially addressed'
    },
    'response_2': {
        'relevance': 'Incorrectly identifies Richard Rush as the director of Rough Going, making it
            irrelevant to the question asked.',
        'evidence_support': "Offers evidence about Richard Rush's lifespan and profession, but this
            does not relate to the query.",
        'logical_coherence': 'Logically coherent for statements about Richard Rush, but fails to
            address the actual question.',
        'rating': 'Not addressed'
    }
}
```

**Server: PartiallySummarizer**

```
{
    'solution': 'The film **Rough Going is a 1925 American silent Western film directed by Wally
        Van**. However, the provided information does not specify where the director, Wally Van,
        was born. This indicates a gap in the available data regarding the birthplace of Wally Van
        .',
    'evaluation and justification': 'The solution identifies the director of the film Rough Going
        as Wally Van but acknowledges that there is no information provided about his place of
        birth. Since the question specifically asks for the birthplace of the director and this
        detail is missing, the solution does not fully address the question.',
    'answerable': 'no',
    'answer': 'None'
}
```

**Server: QuestionSimplifier**

```
{
    'question_knowledge_mapping': [
        ['the director of film Rough Going', 'Wally Van', '**Rough Going is a 1925 American silent
            Western film directed by Wally Van**.']
    ],
    'unresolved_aspects': ['The birthplace of Wally Van'],
    'new_question': "Where was Wally Van born?"
}
```

**Assistant: Culture.Biography.Biography***

I learned that **Van was born Charles Wallace Van Nostrand on September 27, 1880, in New Hyde Park, New York.** This explicitly answers the question about where Wally Van was born, which is New Hyde Park, New York.

**Assistant: Culture.Linguistics**

I found that **Wally Williams (water polo) (1921–2009), New Zealand water polo player** is listed, indicating a Wally Williams who was born in 1921 and was from New Zealand. However, there is no specific mention of where exactly he was born within New Zealand.

**Server: Evaluation**

```
{
    'response_1': {
        'relevance': "Directly answers the question by specifying Wally Van's birthplace as New
            Hyde Park, New York.",
        'evidence_support': 'Provides explicit evidence enclosed in double asterisks about Wally
            Van's birth details.',
        'logical_coherence': 'The response is logically coherent and clearly structured without
            contradictions.',
        'rating': 'Fully addressed'
    },
    'response_2': {
        'relevance': 'Partially relevant as it discusses a person named Wally but not the specific
            Wally Van in question.',
        'evidence_support': 'Provides evidence about Wally Williams but does not link this
            information to Wally Van.',
        'logical_coherence': "Coherent statement but fails to address the specific query about
            Wally Van's birthplace.",
        'rating': 'Not addressed'
    }
}
```

**Server: FullySummarizer**

```
{
    'analysis': "The response provides explicit information about Wally Van's birthplace: **Van
        was born Charles Wallace Van Nostrand on September 27, 1880, in New Hyde Park, New York
        .**",
    'answer': 'Wally Van was born in New Hyde Park, New York.'
}
```

**Response**

Wally Van was born in New Hyde Park, New York.

Figure 14: A example of a multihop QA task using the proposed RIRS