FROM COMPLEX TO ATOMIC: ENHANCING AUG MENTED GENERATION VIA KNOWLEDGE-AWARE DUAL REWRITING AND REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advancements in Retrieval-Augmented Generation (RAG) systems have significantly enhanced the capabilities of large language models (LLMs) by incorporating external knowledge retrieval. However, the sole reliance on retrieval is often inadequate for mining deep, domain-specific knowledge and for performing logical reasoning from specialized datasets. To tackle these challenges, we present an approach, which is designed to extract, comprehend, and utilize domain knowledge while constructing a coherent rationale. At the heart of our approach lie four pivotal components: a knowledge atomizer that extracts atomic questions from raw data, a query proposer that generates subsequent questions to facilitate the original inquiry, an atomic retriever that locates knowledge based on atomic knowledge alignments, and an atomic selector that determines which follow-up questions to pose guided by the retrieved information. Through this approach, we implement a knowledge-aware task decomposition strategy that adeptly extracts multifaceted knowledge from segmented data and iteratively builds the rationale in alignment with the initial query and the acquired knowledge. We conduct comprehensive experiments to demonstrate the efficacy of our approach across various benchmarks, particularly those requiring multihop reasoning steps. The results indicate a significant enhancement in performance, up to 12.6% over the secondbest method, underscoring the potential of the approach in complex, knowledgeintensive applications.

031 032 033

034

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

Large Language Models (LLMs) have revolutionized the field of natural language processing by 035 demonstrating the capability to generate coherent and contextually relevant text. These advanced models are trained on expansive corpora, equipping them with the versatility to execute a diverse 037 spectrum of linguistic tasks, ranging from text completion to translation and summarization (Achiam et al., 2023; Bahrini et al., 2023; Touvron et al., 2023; Anil et al., 2023). Despite their broad capabilities, LLMs exhibit pronounced limitations when tasked with specialized queries in professional 040 domains (Ling et al., 2024; Wang et al., 2023a), a demand that is particularly acute in practical ap-041 plications. This primarily stem from the scarcity of domain-specific training material and an limited 042 grasp of private knowledge and rationale within these domains. As a result, LLMs may produce 043 responses that are not only potentially erroneous but also lack the detail and precision required for 044 expert-level engagement (Bender et al., 2021). Besides the limitations in the domain-specific tasks, another striking issue with LLMs is the phenomena known as "hallucination", where the model generates information that is not grounded in reality or factual data (Beltagy et al., 2020; Xu et al., 046 2024). Moreover, the knowledge base of LLMs, being static and crystallized at the point of their 047 last update, introduces temporal stasis (Brown et al., 2020). Further compounding these challenges 048 is the issue of long-context comprehension (Li et al., 2024). Existing LLMs struggle to maintain an understanding of task definitions across long context, and their performance tends to deteriorate significantly when confronted with more complex and demanding tasks. 051

To mitigate these issues, Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) has emerged as a promising solution, augmenting LLMs with external knowledge retrieval to anchor generated content in factual data. While RAG systems improve the accuracy and reliability of generated re-

sponses by incorporating relevant external information, they often fall short in handling tasks that
 demand deep domain-specific reasoning. Current RAG methods predominantly rely on text retrieval,
 without explicitly focusing on extracting and understanding the underlying knowledge needed to
 tackle complex, multihop reasoning tasks. In this work, we argue that the key to advancing RAG
 lies in knowledge-aware processing, where domain-specific knowledge is not only retrieved but also
 atomized, contextually-decomposed, and reasoned over in a more structured and adaptive manner.

060 The Importance of Knowledge-aware Processing Addressing complex, logic-driven tasks in spe-061 cialized domains requires more than surface-level retrieval of relevant text passages. It demands 062 knowledge extraction and comprehension to deeply understand both the user's information needs 063 and the underlying context of the retrieved data. For example, specialized queries in fields like 064 medicine, law, or finance often involve domain-specific terminology and logic, which generic LLMs may fail to grasp fully. Traditional RAG systems that retrieve text passages based on keyword match-065 ing (Ram et al., 2023; Jiang et al., 2023) or embedding similarity (Gao et al., 2023) may retrieve 066 contextually relevant but semantically shallow information, insufficient for answering intricate ques-067 tions. The challenge lies in ensuring that the retrieved knowledge aligns with the true intent behind 068 the user's query and the broader context of the problem at hand. 069

Moreover, effective comprehension of knowledge extends beyond extracting factual data. It also involves understanding the semantics and relationships hidden within the data corpus. To this end, some recent efforts (Raina & Gales, 2024) attempted to rephrase the original questions with integrating hypothestical answers (Gao et al., 2022) or transform text passages into lists of simple yet atomic questions (Raina & Gales, 2024). This is especially crucial when multiple sources of information must be integrated to form a coherent response. Without a mechanism to extract and comprehend knowledge from each chunk of domain data, current RAG systems are limited in their ability to handle complex tasks that require deep contextual understanding.

The Need for Iterative Reasoning in multihop Questions In many real-world scenarios, a single retrieval or answer generation step may not be enough to fully resolve a complex query. 079 multihop reasoning tasks, where the answer depends on synthesizing information from multiple 080 sources, demand the decomposition of the original query into a series of simpler, interrelated sub-081 questions (Press et al., 2023). Nonetheless, this approach may face obstacles in domains where the knowledge is not readily accessible to LLMs. We argue that the decomposition in such domains 083 should be contextual, rather than a standalone operation, meaning that decomposed queries can be 084 answered with the retrieved knowledge and context progressively and evolve into refining subse-085 quent queries. This iterative approach allows the system to evolve its understanding of the user's 086 inquiry, ensuring that follow-up questions are informed by the most recent retrieval results.

087 For example, consider the HotpotQA (Yang et al., 2018) dataset question: "Who was born first, 880 Erika Jayne or Marco Da Silva?". A straightforward rewriting strategy might decompose this into 089 two separate queries: "What is the age of Erika Jayne?" and "What is the age of Marco Da Silva?". 090 However, if the available knowledge base only contains information such as "Her father left in 1971 091 when Erika Jayne was 9 months old.", traditional rewriting may not lead to the retrieval of relevant 092 information. Techniques like Think-on-Graph (ToG) (Sun et al., 2024), which apply beam search 093 and iterative reasoning on knowledge graphs (KGs), could help. Yet, the performance of KG-based methods is often hampered by limited richness of the knowledge within their structured triples and 094 the need for a high-quality, domain-specific KG. For instance, GPT-4 might generate triples like (Her 095 father, left, action) and (Erika Jayne, 9 months old, age), which are insufficient for reconstructing 096 the original context, "Her father left in 1971 when Erika Jayne was 9 months old."¹. 097

To overcome these limitations, we introduce a novel framework, KAR³-RAG, that employs a knowledge-aware dual rewriting and reasoning mechanism. Our approach features a dynamic interaction between question rewriting and knowledge retrieval, enabling the system to adaptively refine both the query and the retrieved context at each iteration. The core components of our system include Knowledge atomizer, decomposing raw data into atomic questions for more granular retrieval, Query proposer, generating follow-up questions based on the evolving context, Atomic

 ¹We use GPT-4-1106-preview and the KG construction prompt from https://github.com/rahulnyk/knowledge_graph to extract entity relations from the statement "Her father left in 1971 when Erika Jayne was 9 months old,", it generates six triples as follows: "(Her father, left, action), (Her father, 1971, time), (Her father, Erika Jayne, relation), (Erika Jayne, 9 months old, age), (Erika Jayne, 1971, time), (1971, 9 months old, time relation)"

retriever, identifying and retrieving relevant knowledge based on atomic knowledge alignments, and
 Atomic selector, determining the most relevant follow-up questions based on the retrieved information. By leveraging these components, our system can iteratively refine its understanding of both
 the question and the retrieved knowledge, enabling more accurate and context-aware reasoning over
 multiple hops.

113 Our key contributions are as follows: 1). We present a novel RAG framework that capitalizes 114 on the synergistic interaction between interdependent rewriting and reasoning processes, ensuring 115 full utilization of the available context. 2). We enhance retrieval efficacy through a dual rewriting 116 mechanism that modifies both the original questions and the text passages (chunks). Our reasoning 117 process is context-aware, enabling the adaptive formulation of follow-up questions based on the 118 provided context. 3). We report on comprehensive experimental and ablation studies that validate the superior performance of our approach across multiple benchmark datasets, which is up to 12.6% 119 increase over the second-best method. 120

121 122

2 RELATED WORK

123 124 125

2.1 RAG

Retrieval-Augmented Generation (RAG) has emerged as a promising solution that effectively incor-126 porates external knowledge to enhance response generation. With the booming of LLMs (Bahrini 127 et al., 2023; Touvron et al., 2023), most research in the RAG paradigm has shifted towards a frame-128 work that initially retrieves pertinent information from external data sources and subsequently inte-129 grates it into the context of the query prompt as supplementing knowledge for contextually relevant 130 generation (Ram et al., 2023). To enhance the retrieval quality of the naive RAG, advanced RAG 131 approaches implement specific enhancements across the pre-retrieval, retrieval, and post-retrieval 132 processes, including query optimization (Ma et al., 2023; Zheng et al., 2023), multi-granularity 133 chunking (Chen et al., 2023; Zhong et al., 2024), mixed retrieval (Yang, 2023) and re-ranking (Co-134 here, 2023). On one hand, efforts focus on query rewriting, either explicitly (Zheng et al., 2024) or 135 implicitly (Gao et al., 2022), to enhance retrieval performance. On the other hand, several studies 136 transform raw data sources into structured data, ultimately converting them into valuable knowledge for more effective retrieval and reasoning(Wang et al., 2023b; Zheng et al., 2024; Raina & 137 Gales, 2024). In our system, we introduce atomic rewriting for both queries and chunks, which 138 not only achieves multi-granularity query decomposition but also comprehensively extract inherent 139 knowledge from chunks. It has been recognized that naive RAG systems is insufficient for tack-140 ling complex tasks such as summarization (Hayashi et al., 2021) and multihop reasoning (Ho et al., 141 2020). Consequently, most recent research focuses on developing advanced coordination schemes 142 that leverage existing RAG modules to collaboratively address these challenges. Iter-RetGen (Shao 143 et al., 2023) and DSP (Khattab et al., 2023) employ retrieve-read iteration to leverage generation 144 response as the context for next round retrieval. FLARE (Jiang et al., 2023) propose a confidence-145 based active retrieval mechanism that dynamically adjusts query during iterative retrieval processes. 146 Our approach adopts an iteration-based RAG pipeline that leverages context-aware reasoning pro-147 cess, enabling the adaptive formulation of follow-up questions for each iteration and reducing the difficulty of retrieval and reasoning of complex tasks. 148

149

150 2.2 MULTIHOP QA

151 multihop Question Answering (MHQA) (Yang et al., 2018) involves answering questions that re-152 quire reasoning over multiple pieces of information, often scattered across different documents or 153 paragraphs. This task presents unique challenges as it necessitates not only retrieving relevant infor-154 mation but also effectively combining and reasoning over the retrieved pieces to arrive at a correct 155 answer. The traditional graph-based methods in MHQA solves the problem by building graphs and 156 inferring on graph neural networks(GNN) to predict answers (Qiu & other authors, 2019; Fang & 157 other authors, 2020). With the advent of LLMs, recent graph-based methods (Li & Du, 2023; Panda 158 et al., 2024) have evolved to construct knowledge graphs for retrieval and generate response through 159 LLMs. Another branch of methods dynamically convert multihop questions into a series of subqueries by generating subsequent questions based on the answers to previous ones (Trivedi et al., 160 2023; Khattab et al., 2023; Feng et al., 2023). The subqueries guides the sequential retrieval and the 161 retrieved results in turn are used to improve reasoning. Treating MHQA as a supervised problem,



Figure 1: Overview of the KAR³-RAG workflow, illustrating knowledge atomizing by the atomizer, 173 174 and knowledge-aware task decomposition using the query proposer, atomic retrieval and atomic selector. The query proposer generates atomic query proposals based on the original query and 175 reference context. These proposals are used to retrieve the relevant atomic questions, producing 176 retrieved atomic pairs. The atomic selector chooses the most relevant pair and the corresponding 177 chunk, which is added to the reference context for task decomposition in the subsequent iteration. 178 Once the atomic selector determines that no further information is required and no atomic pair are 179 selected, the query and reference context are passed to the generator to produce the final answer. 180

Self-RAG (Zhang et al., 2024) trains an LM to learn to retrieve, generate, and critique text passages, and beam-retrieval (Asai et al., 2023) models the multihop retrieval process in an end-to-end manner by jointly optimizing an encoder and classification heads across all hops. Self-Ask (Press et al., 2023) improves CoT by explicitly asking itself follow-up questions before answering the initial question. This method enables the automatic decomposition of questions and can be seamlessly integrated with retrieval mechanisms to tackle multihop QA.

3 1

162 163 164

165 166

167 168

169

170

171 172

181

182

183

184

185

186

187 188

189

197

212

3 METHODOLOGY

190 3.1 Preliminary

In a RAG system, the textual corpus is divided into a collection of document chunks, denoted as $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$, where d_i represents the *i*-th document chunk. The original query is denoted as q, and its corresponding ground truth answer is represented by a. The retrieval phase involves evaluating the similarity between the query q and each document chunk d_i , after which the top-kmost relevant chunks are selected as retrieval results, forming the basis for subsequent generation.

$$\mathcal{R}: \underset{d_i \in \mathcal{D}}{\text{topk Sim}(q, d_i)} \to D^q \tag{1}$$

Here, the retriever \mathcal{R} selects the top-k most relevant chunks D^q based on the similarity function Sim(·). Finally, the original query and retrieved chunks are fed into the large language model to generate the answer, denoted as $\hat{a} = \mathcal{LLM}(q, D^q)$. In the advanced RAG systems, query rewriting is employed to bridge the semantic gap between the query and the chunks to be retrieved. The rewritten query is represented as $\hat{q} = f_{re}(q)$. The workflow of the advanced RAG system is further improved as follows,

$$\hat{a} = \mathcal{LLM}(q, D^{\hat{q}}), \text{ where } D^{\hat{q}} = \mathcal{R}(\hat{q}, \mathcal{D})$$
(2)

This enhancement allows the system to better align queries with relevant document chunks, enhancing retrieval accuracy and answer generation. However, addressing complex multihop questions remains challenging. These questions often require reasoning across multiple chunks and integrating information through several retrieval and generation steps-a process that a single pass may not fully capture.

211 3.2 FRAMEWORK

To address complex multihop questions, we introduce an enhanced RAG system with Knowledge-Aware **dual R**ewriting and **R**easoning, termed as KAR³. This system employs an iterative retrievalreasoning-generation mechanism that facilitates gradual collection of relevant information and progressive reasoning over incremental context. An overview of the proposed workflow is depicted



(b) Illustrative example of KAR³-RAG case

Figure 2: Illustrative examples of KAR³-RAG cases: (a) Example of knowledge atomizing, (b) RAG case with knowledge-aware task decomposition. As iterations progress, the reference context is enriched by accumulating relevant chunks through atomic retrieval and selection. With the expansion of context, the number of atomic query proposals generated decreases until no further proposals are produced. Subsequently, the iteration process terminates, and the combined query and context are harnessed to produce the final response.

248

241

in Figure 1. In our framework, raw data chunks are broken down into atomic questions using a 249 knowledge atomizer to construct atomic knowledge base for the subsequent retrieval. Oueries are 250 similarly atomized by a query proposer to generate atomic query proposals, which are utilized to 251 retrieve the relevant atomic questions from the knowledge base. Both chunks and queries are rewritten to bridge the semantic gap and improve the alignment of knowledge. An atomic retriever then 253 selects the top-k atomic pairs for each atomic query proposal. Based on these retrieved atomic pairs, 254 an atomic selector, as a reasoner, identifies the most useful atomic pair for problem-solving and adds the corresponding raw chunk to the context. This context is then aggregated with query for the task 256 decomposition in next iteration. The iteration process may terminate earlier if it fails to retrieve suitable atomic questions, either due to the generation of low-quality question proposals or the lack of 257 relevant atomic question candidates. At this point, the query and context are passed to the generator 258 to produce the final answer. 259

260 261

262

3.3 KNOWLEDGE ATOMIZING

Chunked text often contains multifaceted information, and typically, only a subset is needed to address a specific task. Traditional information retrieval methods, which consolidate all information
within a single chunk may not facilitate the efficient retrieval of the precise information required.
Recent research have explored the extraction of triple knowledge units from chunked text and constructing knowledge graphs to facilitate efficient information retrieval (Edge et al., 2024; Panda et al.,
2024). However, the construction of these knowledge graphs is costly, and the inherent knowledge
may not always be fully explored. To better present the knowledge embedded in documents, we propose atomizing the original documents for knowledge extraction, a process we refer as *Knowledge*

Atomizing. This approach leverage the context understanding and content generation capabilities of
 LLMs to automatically tag atomic knowledge pieces within each document chunk.

The presentation of the atomic knowledge can be various. Instead of utilizing declarative sentences or subject-relationship-object tuples, we propose using questions as knowledge indexes to further bridge the gap between stored knowledge and query. In knowledge atomizing process, we input the document chunk to an LLM as context, ask it to generate relevant questions that can be answered by the given chunk as many as possible. These generated atomic questions are stored together with the given chunks. The knowledge atomizer applies atomizing operation on each chunk.

$$f_a(d_k) = \{q_{k1}, q_{k2}, \cdots, q_{km}\}$$
(3)

The atomic questions are generated by atomizer for every chunk, forming an atomic knowledge base, denoted as $\mathcal{KB} = \{f_a(d_k), d_k\}$. An example of knowledge atomizing is illustrated in Figure 2(a), where the atomic questions encapsulate various aspects of the knowledge contained within the chunk. Since each chunk is tagged with multiple atomic questions, an atomic query can be used to locate relevant atomic questions, which then leads to the associated reference chunks.

285

3.4 KNOWLEDGE-AWARE TASK DECOMPOSITION

287 Addressing complex multihop questions often requires integrating multiple pieces of knowledge, 288 which implicitly demands the ability to break down the original question into several sequential or 289 parallel atomic questions for retrieval. We refer to this operation as Task Decomposition. By com-290 bining the extracted atomic knowledge with the original chunks, we construct an atomic knowledge 291 base. Each time a task is decomposed, the atomic knowledge base provides insights into the available 292 knowledge, enabling knowledge-aware task decomposition. We design the Knowledge-Aware Task 293 Decomposition workflow, and the complete algorithm for solving task is detailed in Algorithm 1, and an example is illustrated in Figure 2(b). 294

Initially, the reference context C_0 is initialized as an empty set. In the first iteration, task decomposition relies solely on the query to generate atomic query proposals. As iterations progress, the accumulated context at *t*-th iteration denoted as C_{t-1} , consists of chunks retrieved from previous iterations. During the *t*-th iteration, the query proposer generates atomic query proposals based on the original query and the accumulated context.

300 301

302

311 312

$$f_p(q, \mathcal{C}_{t-1}) = \{\hat{q}_1^t, \hat{q}_2^t, \cdots, \hat{q}_n^t\}$$
(4)

The query proposer $f_p(\cdot)$ can be implemented as either an LLM or a learnable component. we 303 leverage an LLM to generate query proposals that are potentially beneficial for task completion, 304 represented as $\hat{q}^t = \{\hat{q}_i^t\}$. During this process, the selected reference chunks C_{t-1} are provided as 305 context to avoid generating proposals linked to already known knowledge. Consequently, the query 306 proposals evolve with each iteration, adapting to the updated context and aiming to explore addi-307 tional knowledge beyond chunks in the context. For each atomic question proposal, we retrieve its 308 top-k relevant atomic question candidates along with their source chunks from the knowledge base. 309 The atomic retrieval process is: 310

$$\mathcal{R}_{atom} : \underset{q_{kl} \in f_{a}(\mathcal{D})}{\operatorname{topk}} \operatorname{Sim}(\hat{q}_{i}^{t}, q_{kl}) \xrightarrow{\mathcal{KB}} P^{\hat{q}_{i}^{t}}$$
(5)

where the atomic retriever, denoted as \mathcal{R}_{atom} , produces a set of retrieved atomic pairs for each atomic query proposal, represented as $P^{\hat{q}_i^t} = \{(\hat{q}_i^t, q_{k_i l_i}, d_{k_i})\}$. All the retrieved atomic pairs from each atomic query proposal are aggregated to generate an overall set $P^{\hat{q}^t}$. We employ cosine similarity of the corresponding embeddings to retrieve the top-k atomic questions, provided their similarity to a proposed atomic question meets or exceeds a specified threshold δ . With the original question, the accumulated context, and the list of retrieved atomic pairs, the atomic selector employ an LLM to select the most useful atomic pair for problem-solving.

320 321

322

 $\mathcal{LLM}(q, \mathcal{C}_{t-1}, P^{\hat{q}^t}) = (\hat{q}_s^t, q_{k_s l_s}, d_{k_s})$ (6)

The atomic selector, denoted as S_{atom} , further retrieve the relevant raw chunk of the atomic pair selected as the new context added in the *t*-th iteration, denoted as c_t . This chunk corresponds to d_{k_e} .

324 Algorithm 1 Task Solving with Knowledge-Aware Decomposition 325 1: Initialize context $C_0 \leftarrow \phi$ 326 2: for $t = 1, 2, \ldots, N$ do 327 Generate atomic question proposals $\hat{q}^t \leftarrow f_p(q, C_{t-1})$ 3: 328 Retrieve top-k atomic pairs for each atomic query proposal from knowledge base 4: $P^{\hat{q}^t} \xleftarrow{\mathcal{KB}} \mathcal{R}_{atom}(\hat{q}^t, f_a(\mathcal{D}))$ 330 331 5: Select the most useful atomic question or None when additional information is unnecessary 332 333 $q_{k_s l_s} \leftarrow \mathcal{LLM}(q, \mathcal{C}_{t-1}, P^{\hat{q}^t})$ 334 if $q_{k_s l_s}$ is None then 6: 335 7: $\mathcal{C}_t \leftarrow \mathcal{C}_{t-1}$ 336 8: break 337 9: else 338 10: Fetch the relevant chunk c^t corresponding to $q_{k_s l_s}$ 339 Update context $C_t \leftarrow C_{t-1} \cup c^t$ 11: 340 12: end if 341 13: end for 342 14: Generate answer $\hat{a} \leftarrow \mathcal{LLM}(q, \mathcal{C}_t)$ 343 344 345 in equation 6. The chunk retrieval process can be represented by the following formula, 346

$$c_t = \mathcal{S}_{atom}(\mathcal{R}_{atom}(f_p(q, \mathcal{C}_{t-1}), f_a(\mathcal{D})))) \tag{7}$$

This retrieved chunk is aggregated into the reference context for the next round of decomposition, 348 expressed as $C_t = c_t \cup C_{t-1}$. Knowledge-aware decomposition can iterate up to N times, where N 349 is a hyperparameter set to control computational cost. The iteration process may conclude earlier 350 if it fails to retrieve suitable atomic questions, either due to the generation of low-quality question 351 proposals or the absence of relevant atomic question candidates. Alternatively, the process can be 352 halted if the \mathcal{LLM} deems the accumulated knowledge adequate for task completion. This early 353 termination mechanism allows the process to conclude before completing all iterations, reducing 354 computational costs without compromising accuracy. Finally, the accumulated context C_t is utilized 355 to generate answer \hat{a} for the given query q in line 14.

356 It is worth mentioning that the knowledge-aware decomposition can be a learnable component. 357 For each private knowledge base, we can utilize the data collected in each decomposition iter-358 ation—specifically $(q, a, \hat{a}, \{\hat{q}_s^t, c^t, \hat{q}^t, P^{\hat{q}^t}, C_t\})$. This trained proposer can then directly suggest 359 atomic queries q^t during inference, which means lines 3 to 5 in Algorithm 1 can be replaced by a 360 single call to this learned proposer, thereby reducing both inference time and computational cost. 361 We leave the exploration of training an efficient query proposer as future work.

362 364

365

367

347

4 **EVALUATION AND METRICS**

The experimental setup is detailed in Section 4.1, while the primary experimental results are outlined 366 in Section 4.2. Ablation studies are discussed in Section 4.3. Additionally, evaluations on two legal domain-specific benchmarks and three real case studies are included in Appendix A.4 and 368 Appendix A.5, respectively, due to content constraints. 369

370 4.1 EXPERIMENTAL SETUP 371

372 Methods To thoroughly evaluate the performance of our proposed knowledge-aware decomposi-373 tion approach, we have selected a variety of baseline methods that represent different strategies for 374 task-solving with LLMs. We include **Zero-Shot CoT**(Kojima et al., 2022) to assess the inherent rea-375 soning capabilities and built-in knowledge of the underlying LLM without any additional context. **Naive RAG**(Lewis et al., 2020), which introduces external knowledge through retrieval, serves as a 376 benchmark for evaluating the incremental benefits of augmented knowledge. The Self-Ask frame-377 work(Press et al., 2023) is employed to investigate the impact of an iterative question decomposition 378 and answering strategy on task performance. The IRCoT(Trivedi et al., 2023), which iteratively 379 generates the rationale to process the multihop questions, and the Iter-RetGen(Shao et al., 2023), 380 which iteratively uses the recent response as a retrieval query to improve the response quality, are 381 also conducted for performance comparison. Detailed descriptions of these experimental methods 382 are provided in Appendix A.3, and below are the brief summaries: Zero-Shot CoT: Questions are addressed using solely the Chain-Of-Thought (CoT) technique without any example demonstrations 383 or supplemental context. Naive RAG: This approach employs dense retrieval from a flat knowl-384 edge base to procure relevant information for each question as the context in question answering. 385 Self-Ask w/ Retrieval: This method employs a task decomposition strategy wherein the LLMs is 386 prompted to iteratively generate and answer follow-up questions. Furthermore, naive RAG is ap-387 plied for answering each follow-up question. **IRCoT**: This approach iteratively prompts LLMs to 388 generate one more sentence of rationale with retrieved passages, and retrieves new passages with the 389 newly generated reason. Iter-RetGen: This method iteratively answers questions with retrieved pas-390 sages, and uses the newly generated rationale and answer for the next-round retrieval. KAR³: The 391 proposed approach that iteratively decomposes complex questions into sub-questions and retrieves 392 relevant knowledge.

In our experiments, we employ GPT-4 (1106-Preview version) across all the methods outlined previously. For the experiments presented in Section 4.2, the iteration number N is set to 5 for Self-Ask with Retrieval, IRCoT, Iter-RetGen and KAR³. Additionally, the atomic retriever is initialized with k = 4 and $\delta = 0.5$. A comprehensive list of hyper-parameters for the retrieval and LLM can be found in Appendix A.2.

Metrics To ensure consistency with established benchmarks, we adopt F1 as a conventional metric in our experimental evaluation. To more accurately assess the the alignment of responses with the intended answers—beyond mere lexical matching—we introduce a novel evaluation metric employing *GPT-4*. In this process, *GPT-4* acts as an evaluator, assessing the correctness of a response in relation to the question and the correct answer labels. We refer to this metric as Accuracy (Acc). Upon manual inspection of a sample set, the judgments rendered by *GPT-4* demonstrate complete agreement with human evaluators, affirming the reliability of this metric.

Specifically, in cases where multiple correct answer labels are available, we employ a conservative scoring approach for F1 by retaining the highest score achieved. While in the context of computing
 Accuracy (Acc), all admissible answer labels are furnished concurrently to the evaluation process, resulting in a singular accuracy score. Furthermore, a full evaluation results with Exact Match (EM),
 Recall and Precision can be found in Appendix A.3.

411

Datasets Since we are targeting at solving multihop reasoning tasks, three widely-recognized multihop datasets: HotpotQA(Yang et al., 2018), 2WikiMultiHopQA(Ho et al., 2020), and MuSiQue(Trivedi et al., 2022) are used in our evaluation. A brief introduction to these datasets can be found in Appendix A.1. For each dataset, we randomly sample 500 QA data from the *dev* set, disregarding the question type and the number of hops to ensure randomness. We compile the context paragraphs from all sampled QA data into a single knowledge base for each benchmark, creating a more complex retrieval scenario. This design choice aims to rigorously assess the task decomposition and relevant context retrieval capabilities of our model.

- 419
- 420 4.2 MAIN RESULTS

422 As demonstrated in Table 1, our approach achieves superior performance across all datasets, yielding 423 approximately +1.4(1.6%), +7.2(9.6%) and +7.0(12.6%) increases in accuracy over the second 424 best results for HotpotQA, 2WikiMultiHopQA and MuSiQue, respectively. For brevity, 2WikiMul-425 tiHopQA is abbreviated as 2Wiki in the result tables.

When comparing the performance of Zero-Shot CoT and Naive RAG, the inclusion of retrieved context significantly boosts accuracy, with gains ranging from +10.03(42.7%) to +29.0(54.1%).
Furthermore, by incorporating decomposition mechanisms-either by asking follow-up questions as in Self-Ask w/ Retrieval method or by reasoning step-by-step as in IRCoT approach-performance can be further improved. Although both IRCoT and Iter-RetGen utilize rationale as the retrieval query to minimize the semantic gap, Iter-RetGen outperforms IRCoT, as shown in Table 1. A key distinction between these two methods is that IRCoT uses only the newly generated rationale sen-

Method	HotpotQA		2Wiki		MuSiQue		
Method	F1	Acc	F1	Acc	F1	Acc	
Zero-Shot CoT	43.94	53.60	41.40	43.87	22.90	23.47	
Naive RAG	72.67	82.60	59.74	62.80	43.31	44.40	
Self-Ask w/ Retrieval	71.40	80.00	<u>69.06</u>	<u>75.00</u>	46.76	51.40	
IRCoT	67.30	81.00	63.83	70.40	47.57	49.20	
Iter-RetGen	<u>75.27</u>	86.60	67.21	73.60	<u>52.48</u>	55.60	
KAR ³ (Ours)	76.48	88.00	75.00	82.20	57.86	62.60	

Table 1: Performance comparison on multihop QA datasets. Best in bold, second-best underlined.

441 442 443

tence as the query, without revisiting previously generated rationales, while Iter-RetGen uses the
entire rationale generated in last round as the retrieval query, allowing for the reevaluation or correction of past rationales. This suggest that incorporating a mechanism for rethinking or correcting
historical generations may be critical for enhancing performance.

448 Our proposed approach, KAR³, emphasizes knowledge-aware task decomposition and differs from 449 the spontaneous decomposition mechanism reliant on given demonstrations, as employed by Self-450 Ask. It performs decomposition with an awareness of available knowledge and effectively uses 451 atomic questions as an intermediate medium to bridge the semantic gap. The "proposal first, then select" framework, detailed in Algorithm 1, provides an opportunity to validate the intent of the 452 question and rectify potential errors in the historical rationale generation process. A practical ap-453 plication of this point can be seen in Case(a) of Appendix A.5. Consequently, the experimental 454 results demonstrate that KAR³ consistently outperforms other methods, validating its effectiveness 455 in complex reasoning scenarios. 456

457 458

4.3 Ablation Study

The selection of N. We first conducted experiments with the iteration upper bound N set to 1, 2, ... 10, and the results are presented in Figure 3. Detailed performance metrics are available in Table 6 of Appendix A.3. Across all three datasets, there is a consistent uptrend in both Supporting Fact Recall and Answer Accuracy. This pattern underscores the approach's capability to incrementally enhance its outputs through additional iterations, particularly when more detailed and contextually relevant information is required to address problem.

465 Additionally, upon examining the relationship between the number of iterations and the observed 466 growth in supporting fact recall, we note that for HotPotQA and 2WikiMultiHopQA datasets, the 467 recall curves exhibit a pronounced increase up to the fourth iteration. Conversely, the recall for the 468 MuSiQue dataset continues to rise sharply beyond this point, even though the maximum number 469 of hops per question is capped at four, as mentioned in Appendix A.1. This discrepancy implies 470 that while KAR³ is adept at retrieving relevant and useful information within a limited number of 471 iterations, it still has certain limitation: KAR³ relis on the reasoning capability of the used LLM, 472 and further iterations may be required to fully capture the necessary information, especially as the complexity of the questions increases. 473

474 Although the algorithm, as outlined in Algorithm 1, does incorporate early-stopping mechanisms 475 that prevent every question from reaching the maximum iteration limit, a higher N invariably leads 476 to increased computational demands. Therefore, the selection of an appropriate N calls for a del-477 icate balance between computational resources and the expected enhancement in performance. To 478 this end, we choose N = 5 for the experiments in Section 4.2. This value is slightly above the 479 maximum number of hops in the datasets and is justified by the plateauing of performance gains beyond this point as evidenced in Figure 3. It reflects a pragmatic trade-off that accounts for both 480 the computational cost and the retrieval efficacy of our approach. 481

482

The variants of approach components. KAR³ is comprised of four key components: (a) a knowl edge atomizer, (b) a query proposer, (c) an atomic retriever, and (d) an atomic selector. We conduct
 ablation studies to ascertain the individual and collective contributions of these components by in troducing several method variants. For the query proposer, the Single Proposer variant assesses



Figure 3: Supporting fact recall and answer accuracy over iterations. Supporting fact recall is depicted in blue, while answer accuracy in orange.

Variant Mathad	Hotp	otQA	2W	<i>v</i> iki	MuS	iQue
	F1	Acc	F1	Acc	F1	Acc
w/ Single Proposer	75.06	85.60	70.19	76.40	49.67	52.20
w/ Chunk Selector	72.80	83.20	61.65	65.80	49.31	53.40
w/ Chunk Retriever	76.31	86.60	67.14	72.40	49.05	53.00
KAR^3 (Ours)	76.48	88.00	75.00	82.20	57.86	62.60

Table 2: Ablation study on components of KAR³.

the impact of generating only a single query as opposed to multiple queries. In the case of the 507 atomic selector, the Chunk Selector variant examines the implications of selecting information in 508 larger segments, or chunks, rather than focusing on atomic questions. Finally, the Chunk Retriever 509 variant combines the modifications of the previous two variants: it generates a single query and 510 retrieves information in the form of chunks from a knowledge base that has not been pre-processed 511 into atomic units. The atomic selection phase, corresponding to Algorithm 1 line 5, is then replaced 512 by directly selecting the most useful chunk since no available atomic question in this variant. These 513 variants enable us to isolate the impact of each component and understand how they interact to 514 produce the system's output. The experimental results of these variants are presented in Table 2. 515

As evidenced by the results in Table 2, the individual contributions of the components were evaluated. We observed that limiting the approach to propose only a single atomic query led to accuracy reductions of 2.8%, 7.0% and 16.6% for the respective datasets. Similarly, opting for chunk-based selections over atomic questions resulted in Accuracy declines of 5.5%, 16.2% and 14.7%. The substitution of the atomic retriever with a general chunk retriever caused the Accuracy to drop by approximately 1.6%, 11.9%, 15.3%, respectively. These ablation studies imply that each designed component is crucial for achieving optimal retrieval performance and coherent reasoning traces.

522

496

497

504 505 506

Limitation Discussion. Beyond the need for additional iterations to extract crucial information
 for complex questions, our experiments with GPT-3.5 (details in Table7 in AppendixA.3), indicate
 a limitation in relying on LLMs' reasoning capabilities. The performance of KAR³ does not significantly surpass that of methods like IRCoT and Self-Ask w/ Retrieval and occasionally falls short
 compared to Self-Ask w/ Retrieval. This highlights that KAR³'s success hinges on its advanced
 reasoning skills and its ability to robustly follow complex instructions.

529 530

531

5 CONCLUSION

532 We present an advanced RAG system, enhanced with knowledge-aware dual rewriting, and reason-533 ing capabilities, designed to improve knowledge extraction and rationale formulation within special-534 ized datasets. The comprehensive results from extensive experiments underscore the efficacy of our approach, particularly in scenarios involving benchmarks with multihop questions. For future work, 535 we aim to refine the system's proficiency through the integration of in-context learning (Wei et al., 536 2022), by adaptively selecting demonstrations for the query proposer. This will further enhance its 537 ability to perform knowledge-aware question rewriting. Additionally, we are interested in develop-538 ing a knowledge-aware atomizer capable of incorporating feedback from sample questions, thereby improving its understanding of the most beneficial types of atomic knowledge.

540 REFERENCES 541

554

559

560

564

565

566 567

569

570

576

581

587

- J Achiam, S Adler, S Agarwal, L Ahmad, I Akkaya, FL Aleman, D Almeida, J Altenschmidt, 542 S Altman, S Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. 543
- 544 R Anil, S Borgeaud, Y Wu, J-B Alayrac, J Yu, R Soricut, J Schalkwyk, AM Dai, A Hauth, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023. 546
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning 547 to retrieve, generate, and critique through self-reflection, 2023. URL https://arxiv.org/ 548 abs/2310.11511. 549
- 550 Aram Bahrini, Mohammadsadra Khamoshifar, Hossein Abbasimehr, Robert J Riggs, Maryam Es-551 maeili, Rastin Mastali Majdabadkohne, and Morteza Pasehvar. Chatgpt: Applications, opportuni-552 ties, and threats. In 2023 Systems and Information Engineering Design Symposium (SIEDS), pp. 553 274–279. IEEE, 2023.
- Iz Beltagy, Arman Cohan, and Kyle Lo. Fact or fiction: Verifying scientific claims. In Proceedings 555 of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 556 7534-7550, 2020.
- 558 Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Mitchell. On the dangers of stochastic parrots: Can language models be too big? In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, pp. 610–623. ACM, 2021.
- 561 T Brown, B Mann, N Ryder, M Subbiah, JD Kaplan, P Dhariwal, A Neelakantan, P Shyam, G Sas-562 try, A Askell, et al. Language models are few-shot learners. Advances in neural information 563 processing systems, 33:1877-1901, 2020.
 - Umar Butler. Open australian legal qa, 2023. URL https://huggingface.co/datasets/ umarbutler/open-australian-legal-ga.
- Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang, 568 and Dong Yu. Dense x retrieval: What retrieval granularity should we use? arXiv preprint arXiv:2312.06648, 2023. URL https://arxiv.org/pdf/2312.06648.pdf.
- Cohere. Say goodbye to irrelevant search results: Cohere rerank is here. https://txt. 571 cohere.com/rerank/, 2023. Accessed: 2023-08-28. 572
- 573 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, 574 and Jonathan Larson. From local to global: A graph rag approach to query-focused summariza-575 tion, 2024. URL https://arxiv.org/abs/2404.16130.
- Yu Fang and other authors. Hierarchical graph network for multi-hop question answering. In Pro-577 ceedings of the Annual Meeting of the Association for Computational Linguistics (ACL). Associ-578 ation for Computational Linguistics, 2020. 579
- 580 Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. Lawbench: Benchmarking legal knowledge of large language models. arXiv preprint arXiv:2309.16289, 2023. 582
- 583 Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. Retrieval-generation 584 synergy augmented large language models, 2023. URL https://arxiv.org/abs/2310. 585 05149. 586
 - Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. Precise zero-shot dense retrieval without relevance labels, 2022. URL https://arxiv.org/abs/2212.10496.
- 589 Y Gao, Y Xiong, X Gao, K Jia, J Pan, Y Bi, Y Dai, J Sun, and H Wang. Retrieval-augmented 590 generation for large language models: A survey. arXiv preprint arXiv:2312.10997, 2023. 591
- Hiroaki Hayashi, Prashant Budania, Peng Wang, Chris Ackerson, Raj Neervannan, and Graham 592 Neubig. WikiAsp: A Dataset for Multi-domain Aspect-based Summarization. Transactions of 593 the Association for Computational Linguistics, 9:211–225, 2021.

622

628

630

631

635

636

594	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop
595	ga dataset for comprehensive evaluation of reasoning steps. arXiv preprint arXiv:2011.01060.
596	2020.
597	

- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, 598 Jamie Callan, and Graham Neubig. Active retrieval augmented generation, 2023. URL https: //arxiv.org/abs/2305.06983. 600
- 601 Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, 602 and Matei Zaharia. Demonstrate-search-predict: Composing retrieval and language models for 603 knowledge-intensive nlp, 2023. URL https://arxiv.org/abs/2212.14024.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large 605 language models are zero-shot reasoners. Advances in neural information processing systems, 606 35:22199-22213, 2022. 607
- 608 P Lewis, E Perez, A Piktus, F Petroni, V Karpukhin, N Goyal, H Kuttler, M Lewis, WT Yih, T Rocktaschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. Advances in 609 Neural Information Processing Systems, 33:9459–9474, 2020. 610
- 611 Ruosen Li and Xinya Du. Leveraging structured information for explainable multi-hop question 612 answering and reasoning, 2023. URL https://arxiv.org/abs/2311.03734. 613
- 614 Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with long in-context learning, 2024. URL https://arxiv.org/abs/2404.02060. 615
- 616 Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy 617 Chowdhury, Yun Li, Hejie Cui, Xuchao Zhang, Tianjiao Zhao, Amit Panalkar, Dhagash Mehta, 618 Stefano Pasquali, Wei Cheng, Haoyu Wang, Yanchi Liu, Zhengzhang Chen, Haifeng Chen, Chris 619 White, Quanquan Gu, Jian Pei, Carl Yang, and Liang Zhao. Domain specialization as the key 620 to make large language models disruptive: A comprehensive survey, 2024. URL https:// 621 arxiv.org/abs/2305.18703.
- X. Ma, Y. Gong, P. He, H. Zhao, and N. Duan. Query rewriting for retrieval-augmented large 623 language models. arXiv preprint arXiv:2305.14283, 2023. 624
- 625 Pranoy Panda, Ankush Agarwal, Chaitanya Devaguptapu, Manohar Kaul, and Prathosh A P. 626 Holmes: Hyper-relational knowledge graphs for multi-hop question answering using llms, 2024. 627 URL https://arxiv.org/abs/2406.06027.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. Measuring 629 and narrowing the compositionality gap in language models, 2023. URL https://arxiv. org/abs/2210.03350.
- 632 Minghui Qiu and other authors. Dynamically fusing recurrent neural networks for multi-hop ques-633 tion answering. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP). Association for Computational Linguistics, 2019. 634
 - Vatsal Raina and Mark Gales. Question-based retrieval using atomic units for enterprise rag, 2024. URL https://arxiv.org/abs/2405.12363.
- 638 O Ram, Y Levine, I Dalmedigos, D Muhlgay, A Shashua, K Levton-Brown, and Y Shoham. In-639 context retrieval-augmented language models. Transactions of the Association for Computational Linguistics, 11:1316–1331, 2023. 640
- 641 Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. Enhancing 642 retrieval-augmented large language models with iterative retrieval-generation synergy, 2023. URL 643 https://arxiv.org/abs/2305.15294. 644
- 645 Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel M. Ni, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of 646 large language model on knowledge graph, 2024. URL https://arxiv.org/abs/2307. 647 07697.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
 language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving re trieval with chain-of-thought reasoning for knowledge-intensive multi-step questions, 2023. URL
 https://arxiv.org/abs/2212.10509.
- ⁶⁵⁹
 ⁶⁶⁰
 ⁶⁶¹
 ⁶⁶²
 ⁶⁶²
 ⁶⁶³
 ⁶⁶³
 ⁶⁶³
 ⁶⁶⁴
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁷
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶¹
 ⁶⁶²
 ⁶⁶²
 ⁶⁶³
 ⁶⁶³
 ⁶⁶³
 ⁶⁶⁴
 ⁶⁶⁴
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁷
 ⁶⁶⁷
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶¹
 ⁶⁶²
 ⁶⁶²
 ⁶⁶²
 ⁶⁶³
 ⁶⁶³
 ⁶⁶⁴
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁷
 ⁶⁶⁷
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶¹
 ⁶⁶²
 ⁶⁶²
 ⁶⁶²
 ⁶⁶³
 ⁶⁶³
 ⁶⁶⁴
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁷
 ⁶⁶⁷
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁶¹
 ⁶⁶²
 ⁶⁶²
 ⁶⁶²
 ⁶⁶³
 ⁶⁶³
 ⁶⁶³
 ⁶⁶⁴
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁵
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁷
 ⁶⁶⁶
 ⁶⁶⁸
 ⁶⁶⁶
 ⁶⁶⁶
 ⁶⁶⁷
 ⁶⁶⁶
 ⁶⁶⁸
 ⁶⁶⁷
 ⁶⁶⁸
 ⁶⁶⁸
 ⁶⁶⁹
- Yu Wang, Nedim Lipka, Ryan A. Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. Knowledge
 graph prompting for multi-document question answering, 2023b. URL https://arxiv.
 org/abs/2308.11730.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V
 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models.
 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in *Neural Information Processing Systems*, volume 35, pp. 24824–24837. Curran Associates, Inc.,
 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/
 file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf.
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models, 2024. URL https://arxiv.org/abs/2401.11817.
- S. Yang. Advanced rag 01: Small-to-big retrieval. https://towardsdatascience.com/
 advanced-rag-01-small-to-big-retrieval-172181b396d4, 2023. Accessed:
 2023-08-28.
- ⁶⁷⁹
 ⁶⁷⁹ Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
 - Jiahao Zhang, Haiyang Zhang, Dongmei Zhang, Yong Liu, and Shen Huang. End-to-end beam retrieval for multi-hop question answering, 2024. URL https://arxiv.org/abs/2308.08973.
- H. S. Zheng, S. Mishra, X. Chen, H.-T. Cheng, E. H. Chi, Q. V. Le, and D. Zhou. Take a step back: Evoking reasoning via abstraction in large language models. *arXiv preprint arXiv:2310.06117*, 2023.
- Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H. Chi, Quoc V Le, and Denny Zhou. Take a step back: Evoking reasoning via abstraction in large language models, 2024. URL https://arxiv.org/abs/2310.06117.
- Zijie Zhong, Hanwen Liu, Xiaoya Cui, Xiaofan Zhang, and Zengchang Qin. Mix-of-granularity:
 Optimize the chunking granularity for retrieval-augmented generation, 2024. URL https://
 arxiv.org/abs/2406.00456.
- 696

683

684

685

- 697
- 030

699

702 A APPENDIX

A.1 INTRODUCTION TO OPEN-DOMAIN BENCHMARKS

We provide a brief overview of the multihop QA datasets used in our experiments, noting that our method does not leverage the question type information nor the number of hops information during the solving process, as our approach is designed to be agnostic to such classifications. Table 3 outlines the distribution of question types within our sampled sets, offering insight into the variety of reasoning challenges presented in our evaluation, though this does not directly impact our method.

711

720

704

705

HotpotQA The HotpotQA dataset is a well-known multihop QA benchmark primarily consist-712 ing of 2-hop questions, each associated with 10 Wikipedia paragraphs. Among these, some para-713 graphs contain supporting facts essential to answering the question, while the rest serve as dis-714 tractors. The dataset also includes a question type field, which delineates the logical reasoning 715 required-comparison questions involve contrasting two entities, and bridge questions require in-716 ferring the bridge entity, or inferring the property of an entity through an intermediary entity, or 717 locating the answer entity (Yang et al., 2018). Although our method operates independently of 718 these types, their description here is to exemplify the nature of questions within the dataset and to 719 contextualize the expected performance variance across different benchmarks.

2WikiMultiHopQA Inspired by HotpotQA, 2WikiMultiHopQA expands the diversity of question 721 types. It retains the *comparison* type from HotpotQA and introduces *inference* and *compositional* 722 questions that evolve from the *bridge* type by focusing on entity attribute deduction and entity lo-723 cation, respectively. Additionally, the bridge comparison type is a novel category that requires a 724 synthesis of *bridge* and *comparison* reasoning. This dataset typically presents 2-hop to 4-hop ques-725 tions, each accompanied by 10 Wikipedia paragraphs containing supporting facts and distractors. 726 While these types inform the dataset's structure, they are not utilized by our method, which treats all 727 questions uniformly regardless of their categorization. For the sake of brevity, 2WikiMultiHopQA 728 is abbreviated to 2Wiki in the result tables throughout this paper. 729

730 MuSiQue Addressing the issue that many multihop questions can be solved via short-731 cuts—arriving at correct answers without proper reasoning—MuSiQue implements stringent filters 732 and additional mechanisms specifically designed to encourage connected reasoning, as reported by 733 Trivedi et al., (Trivedi et al., 2022). Unlike the other datasets, MuSiQue does not categorize questions by type, but it does provide explicit information on the number of hops required for each question, 734 ranging from 2 to 4 hops. Each question is associated with 20 context paragraphs, which introduce 735 a mix of relevant and irrelevant information, further complicating the task of discerning the correct 736 reasoning path. This explicit hop information, while not used by our method, underscores the com-737 plexity of the dataset and the robustness required by models to handle such challenges effectively. 738

Table 3: Distribution of question types across three distinct multihop QA datasets.

			Туре	Count	Ratio			
			comparison	132	26.4%	#Hops	Count	Ratio
Туре	Count	Ratio	inference	64	12.8%	2	263	52.6%
comparison	107	21.4%	compositional	196	39.2%	3	169	33.8%
bridge	393	78.6%	bridge_comparison	108	21.6%	4	68	13.6%
(a) HotPotQA			(b) 2WikiM	ultiHopQ	DA A	(c) MuSiQ	ue

746 747

739

748 749

A.2 HYPER-PARAMETERS

During the knowledge extraction phase, we utilize a *temperature* setting of 0.7 specifically for the *Knowledge Atomizing* process, promoting a balance between diversity and determinism in the generated atomic knowledge. Conversely, for all question-answering (QA) steps in each method, we implement a *temperature* of 0, ensuring consistent responses from the model.

Regarding the retrieval component, we engage the *text-embedding-ada-002* (version 2) as our embedding model for both the general knowledge bases and the atomic knowledge bases. For the

756 general knowledge bases used in Naive RAG and Iter-RetGen, the retriever is configured to fetch 757 up to 16 knowledge chunks, applying a retrieval score threshold of 0.2. For the general knowledge 758 bases used in Self-Ask w/ Retrieval and IRCoT, where the retrieval chunks are used for a single 759 follow-up question answering or the generation of single continuous rationale sentence, the refer-760 ence chunks for whole rationale or final question answering are accumulated. The system retrieves 761 4 relevant chunks per request, maintaining the same score threshold of 0.2. In the case of atomic 762 knowledge bases, the retriever is set to retrieve 4 relevant atomic questions for each atomic query 763 but with a higher threshold 0.5 due to the shorter content length.

- 764
- 765

776

777

778

779

780

781

782

783

794

797

798 799

800

801

802

804

805

A.3 DETAILED EXPERIMENTAL RESULTS ON OPEN-DOMAIN DATASETS

766 In addition to the methods outlined in Table 1, we also conduct experiments with a Knowledge 767 Graph-based method, GraphRAG (Edge et al., 2024), to determine its effectiveness in multihop 768 reasoning tasks. GraphRAG was inferred in both local and global modes. While we acknowledge the 769 existence of other knowledge graph-based methods designed to tackle multihop questions-such as 770 KGP (Wang et al., 2023b), which necessitates further fine-tuning, and ToG (Sun et al., 2024), which 771 depends on external, predefined knowledge graphs-we contend that the additional requirements of 772 these approaches render them not directly comparable to our approach, which is both training-free 773 and relying only on a specific knowledge base.

- The methods evaluated in this study are listed as follows:
 - Zero-Shot CoT: Questions are addressed using solely the Chain-Of-Thought (CoT) technique, which prompts the LLMs to articulate its reasoning process step-by-step without the aid of example demonstrations or supplemental context. This method assesses the LLMs' intrinsic knowledge and reasoning capabilities in a zero-shot setting.
 - Naive RAG: This approach employs dense retrieval from a flat knowledge base to procure relevant information for each question. The knowledge base consists of pre-embedded chunks are matched to the original question based on semantic similarity. The retrieval process is direct, without any intermediate task decomposition.
- 784 • Self-Ask w/ Retrieval: This method employs a task decomposition strategy wherein the 785 LLMs is prompted to iteratively generate and answer follow-up questions, thereby breaking 786 down complex problems into more manageable sub-tasks. General demonstrations illustrating the logic and methodology of task decomposition are provided for all benchmarks 788 to guide the LLMs' reasoning process. Different to the original setting(Press et al., 2023), where the framework relies solely on LLM's own knowledge to answer each follow-up question, in this setting, we introduces an additional retrieval component. Relevant chunks 791 are retrieved with the follow-up question as the query from a flat knowledge base to provide 792 a reference context. What's more, we also limit the decomposition process to raise up to Nfollow-up questions to align with other methods. 793
 - **IRCoT**: This approach iteratively prompts LLMs to generate one more sentence of rationale with retrieved passages, and retrieves new passages with the newly generated reason. The original setting limit the process with a maximum token number(Trivedi et al., 2023). In our experiments, we limit the total number of iterations to the constant N we used for our methods.
 - Iter-RetGen: This method iteratively answers questions with retrieved passages, and uses the newly generated rationale and answer for the next-round retrieval. In this setting, we also limit the total number of iterations to the same N.
 - **GraphRAG Local**: The knowledge base is pre-processed to construct a knowledge graph in accordance with the public guidance. The evaluation is inferred in local mode.
 - **GraphRAG Global**: The knowledge base is pre-processed to construct a knowledge graph in accordance with the public guidance. The evaluation is inferred in global mode.
- KAR³ (Ours): The proposed knowledge-aware decomposition method iteratively decomposes complex questions into sub-questions and retrieves relevant knowledge up to a maximum of N iterations. This process limits the context for the final answer to the five most useful knowledge chunks.

812						
813	Method	EM	F1	Acc	Precision	Recall
814	Zero-Shot CoT	32.60	43.94	53.60	46.56	43.97
815	Naive RAG	56.80	72.67	82.60	74.52	74.86
816	Self-Ask w/ Retrieval	57.00	71.40	80.00	73.25	73.95
817	IRCoT	51.40	67.30	81.00	69.32	72.15
818	Iter-RetGen	<u>59.60</u>	75.27	86.60	77.18	77.62
819	GraphRAG Local	0.00	10.66	89.00	5.90	83.07
820	GraphRAG Global	0.00	7.42	64.80	4.08	63.16
821	KAR ³ (Ours)	61.40	76.48	88.00	78.53	<u>78.96</u>
822		(a)]	II. 4D. 4C			
823		(a)	HOIPOIQ	ĮA		
824	Method	FM	F1	Acc	Precision	Recall
825	Zero-Shot CoT	35.67	41.40	13.87	41.43	/3 11
826	Naiva PAG	51.20	50.74	43.07 62.80	50.06	43.11 62.30
827	Salf Ask w/ Detrieval	60.60	59.74 60.06	75.00	59.00 67.88	02.30 73.15
828	BCoT	55.00	<u>09.00</u>	70.40	62.47	13.13
829	IRC01	57.00	67.01	70.40	02.47	71.00
830	Grank DAC Logal	37.80	07.21	75.00	674	71.09
831	GraphRAG Local	0.00	11.85	/1.20	0.74	<u>/5.1/</u> 55.42
832	GraphRAG Global	0.00	7.35	45.00	4.09	55.4 <i>3</i>
833	KAR [°] (Ours)	65.80	75.00	82.20	73.63	79.08
834		(b) 2Wil	kiMultiH	IopQA		
835		E14	F1	•	D · ·	
836	Method	EM	FI	Acc	Precision	Recall
837	Zero-Shot CoT	12.93	22.90	23.47	24.40	24.10
838	Naive RAG	32.00	43.31	44.40	44.42	47.29
839	Self-Ask w/ Retrieval	38.20	46.76	51.40	46.75	51.00
840	IRCoT	36.00	47.57	49.20	48.70	50.30
841	Iter-RetGen	40.20	<u>52.48</u>	<u>55.60</u>	<u>53.51</u>	<u>56.45</u>
842	GraphRAG Local	0.60	9.62	49.80	5.73	55.82
843	GraphRAG Global	0.00	5.16	44.60	2.82	52.19
844	KAR^3 (Ours)	47.40	57.86	62.60	58.52	61.37
845		(c)	MuSiO	10		
846		(0)	musiQi	ic.		
847						

810 Table 4: Detailed performance comparison on multihop QA datasets. Best in bold, second-best 811 underlined.

848 What's more, three more metrics are employed in Appendix. Exact Match (EM), which assesses whether the response is identical to a predefined correct answer is applied as the community usually 849 did. Furthermore, we encounter situations where a method achieves high accuracy (Acc) scores 850 yet registers low F1 scores. To elucidate the underlying factors of such discrepancies, we also 851 report on the Recall and Precision of the generated responses. Recall measures the proportion of 852 relevant tokens from the answer labels that are captured in the response, while precision evaluates 853 the relevance of the tokens in the generated answer with respect to the correct labels. 854

The detailed evaluation results across HotpotQA, 2WikiMultiHopQA, and MuSiQue are presented 855 in Table 4. Notably, knowledge graph-based method, GraphRAG Local, excels in HotpotQA-a 856 dataset predominantly comprised of 2-hop questions. However, in the other two datasets, which con-857 tain questions involving more hops, GraphRAG Local is merely on par with IRCoT. This highlights 858 the challenge that knowledge graph-based methods face in addressing complex multihop questions. 859

860 Regarding GraphRAG, originally designed for the query-focused summarization (QFS) task as out-861 lined by (Edge et al., 2024), we observe its suboptimal performance in both local and global modes compared to our method. GraphRAG exhibits a curious trend: it achieves higher accuracy and recall 862 scores while performing lower on EM, F1, and Precision metrics. A closer analysis of GraphRAG's 863 outputs reveals a tendency to echo the query and include meta-information about the answer within

Question	Which country is home to Alsa Mall and Spencer Plaza?
Answer Labels	India
Answer of GraphRAG	Alsa Mall and Spencer Plaza are both located in Chennai, India [Data: India and Chennai Community (2391): Entities (4901, 4904): Relation-
	ships (9479, 1687, 5215, 5217)].

Table 5: An Example of GraphRAG Local output on a HotpotQA question. The table showcases the tendency to repeat the question and include meta-information in its response.

Table 6: Ablation study on hyper-parameter N. Recall^{*} indicates the recall of supporting facts.

N	H	HotpotQA			2Wiki			MuSiQue			
11	Recall [*]	F1	Acc	Recall [*]	F1	Acc	Recall [*]	F1	Acc		
1	42.96	59.46	70.20	40.41	41.08	43.00	31.20	32.55	32.80		
2	82.04	74.27	84.80	78.83	70.22	77.20	56.43	48.46	50.00		
3	90.16	76.90	87.20	87.71	72.84	79.40	64.82	53.50	57.20		
4	92.46	76.49	87.80	92.86	74.68	81.80	69.87	55.73	59.40		
5	92.83	76.48	88.00	94.06	75.00	82.20	73.08	57.86	62.60		
6	93.35	77.67	89.00	94.76	75.12	81.80	74.88	57.03	61.20		
7	93.68	77.32	88.80	94.91	75.44	82.40	76.07	56.66	61.40		
8	93.78	76.88	88.40	95.06	75.16	82.00	76.72	57.65	62.40		
9	93.78	76.99	88.60	95.11	74.89	81.80	76.90	57.17	61.40		
10	93.78	77.52	89.00	95.16	75.09	82.00	77.20	57.69	62.40		

its graph structure. Despite attempts to refine its QA prompt, this behavior persists. An illustrative example is presented in Table 5, which shows GraphRAG Local's response to a question from HotpotQA.

Table 6 lists the granular performance metrics according to those we shown in Figure 3 for the ablation study on the iteration upper bound N. Different to the **Recall** we reported in Table 4, which indicates the recall tokens of the answer labels, the **Recall**^{*} here represents the recall of the supporting facts provided by these datasets.

As introduced in the limitation discussion section, we have carried out a series of experiments uti-lizing GPT-3.5. The outcomes of these experiments are delineated in Table 7. For these specific trials, we substituted GPT-4 (1106-Preview) with GPT-3.5 (1106-Preview) as the language model, while maintaining all other experimental settings identical to those employed in the experiments summarized in Table 1.

A.4 EVALUATION ON LEGAL BENCHMARKS

In this subsection, we present the performance of our approach on two legal benchmarks: Law-Bench Fei et al. (2023) and Open Australian Legal QA Butler (2023). Before doing so, we provide a brief description of each benchmark.

LawBench LawBench is a comprehensive legal benchmark for Chinese laws. It comprises 20 meticulously designed tasks aimed at accurately assessing the legal capabilities of LLMs. Unlike some existing benchmarks that rely solely on multiple-choice questions, LawBench includes a va-riety of task types that are closely related to real-world applications. These tasks encompass legal entity recognition, reading comprehension, crime amount calculation, and legal consulting, among others. Since not all tasks are RAG-oriented (e.g., reading comprehension), we have selected 6 specific tasks, which are detailed in Table 8. The number of questions of each task is 500.

We also provide example questions of these tasks for the readers reference (translated using GPT-4).

1-1: Answer the following question by directly providing the content of \hookrightarrow the article:What is the content of Article 76 of the Securities Law \rightarrow ?

939

940

941

942

943

944

945

946

947

948

949

950

971

918	Table 7: Performance comparison of implementations with GPT-3.5. Best Accuracy in bold, second-
919	best Accuracy underlined.

Mathad	HotpotQA		2W	/iki	MuSiQue		
Method	F1	Acc	F1	Acc	F1	Acc	
Self-Ask w/ Retrieval	49.52	61.40	53.83	60.00	31.05	35.20	
IRCoT	56.39	68.40	40.31	46.00	33.93	34.40	
Iter-RetGen	48.63	66.80	44.32	55.20	25.77	37.80	
KAR^3 (Ours)	46.37	68.80	41.95	58.20	26.80	39.60	

Table 8: Overview of LawBench tasks

Task No.	Task	Туре	Metric
1-1	Statute Recitation	Generation	F1
1-2	Legal Knowledge Q&A	Single Choice	EM
3-1	Statute Prediction (Fact-based)	Multiple Choices	EM
3-2	Statute Prediction (Scenario-based)	Generation	F1
3-6	Case Analysis	Single Choice	EM
3-8	Consultation	Generation	F1

1-2: According to the 'Securities Law', which of the following statements \Rightarrow about stock exchanges is incorrect? A: Without the permission of \Rightarrow the stock exchange, no entity or individual may publish real-time \Rightarrow securities trading information; B: The stock exchange may restrict \Rightarrow trading on securities accounts that exhibit major abnormal trading \Rightarrow conditions as needed, and report to the securities regulatory \Rightarrow authority under the State Council for record; C: The accumulated \Rightarrow property of a member-based stock exchange belongs to the members, \Rightarrow and their rights are jointly enjoyed by the members; during its \Rightarrow existence, the accumulated property may not be distributed to the \Rightarrow members; D: The stock exchange formulates listing rules, trading \Rightarrow rules, member management rules, and other relevant rules in \Rightarrow accordance with securities laws and administrative regulations, and \Rightarrow reports to the securities regulatory authority under the State

 \hookrightarrow Council for record.

951 3-1: Based on the following facts and charges, provide the relevant ← articles of the Criminal Law. Facts: The Yushu City, Jilin Province 952 \hookrightarrow , accused that on November 15, 2015, the defendant He signed a car 953 \hookrightarrow rental agreement with Guo, the owner of a taxi with license plate 954 \hookrightarrow number . The agreement stipulated a monthly rent of RMB 3,900.00, 955 \hookrightarrow payable monthly. On January 19, 2016, without the knowledge of Guo, 956 \hookrightarrow the defendant He concealed the truth and falsely claimed to be the $\,\hookrightarrow\,$ owner of the taxi. He signed a car rental agreement with the 957 \hookrightarrow victim Ma, with a monthly rent of RMB 3,800.00 and a rental period 958 \rightarrow of one year, collecting a total of RMB 50,600.00 from Ma for one 959 \hookrightarrow year's rent and vehicle deposit. On February 26, 2016, the taxi was 960 \hookrightarrow retrieved by its owner Guo from the victim Ma. The victim Ma 961 \hookrightarrow repeatedly asked the defendant He to return the rent and deposit, \hookrightarrow but the defendant He refused to return them. The prosecution 962 \hookrightarrow provided evidence including the defendants confession, the victims 963 \hookrightarrow statement, witness testimonies, and documentary evidence, and 964 \hookrightarrow believed that the defendant He, with the purpose of illegal 965 \hookrightarrow possession, defrauded others of their property by fabricating facts 966 \rightarrow and concealing the truth during the signing and performance of the \hookrightarrow contract. The amount was relatively large, and his actions 967 \hookrightarrow violated the provisions of Article of the Criminal Law of the 968 \hookrightarrow Peoples Republic of China, and he should be held criminally 969 \hookrightarrow responsible for . Charge: Contract Fraud. 970 3-2: Please provide the legal basis according to the specific scenario

⇒ and question, only the content of the specific legal provision is → needed, each scenario involves only one legal provision. Scenario:

73							
74		Т	ask	Zero-Shot CoT	GraphRAG Local	Ours (N=5)	
75			1-1	21.31	23.27	78.58	
76			1-2	54.24	62.60	70.60	
77			3-1	53.32	74.60	83.16	
78		LawBench	3-2	27.51	25.98	46.05	
79			3-6	51.16	47.64	61.91	
30			3-8	17.44	18.43	23.58	
81		Open Austra	lian Legal QA	25.10	34.35	63.34	
2							
33 34		Table 10): Evaluation F	Results on Legal 1	Benchmarks (Metri	c is Acc)	
35		Т	ask	Zero-Shot CoT	GraphRAG Local	Ours (N=5)	
86			1-1	1 23	16.60	90.12	
7			1-2	54.00	<u>63 40</u>	70.60	
8			3-1	49.90	75.40	88.82	
9		LawBench	3-1	49.90	$\frac{73.40}{27.60}$	67.54	
0			3-2	51.12	<u>27.00</u> 57.00	62 73	
1			3-0 2.0	51.12	<u>57.00</u>	02.73	
2			3-8	49.70	<u>38.80</u>	01.72	
3		Open Austra	lian Legal QA	16.48	88.27	98.59	
<u>л</u>							
5							
6	$\hookrightarrow A$	cargo ship	o arrives a	t the port of	the goods Upd	t the cons	ıgnee
7	\rightarrow Id \rightarrow pr	ovision ca	n the cant.	e to correct ain unload th	re goods at and	ther appro	zya⊥ nriate
202	rq ↔ ⇔ pl	ace?	in ene cape	ain anioaa ci	le goodb de dife	cher appro	prince
0	3-6: One	year afte	r the bar o	ppened, the b	usiness enviro	nment chang	jed
00	\hookrightarrow dr	astically,	and all p	artners held	a meeting to d	liscuss	
01	\rightarrow co	untermeasu	ires. Accor	ding to the '	Partnership Er	terprise L	aw,' the
02	\rightarrow 1 \rightarrow he	lieves that	t the name	'Tongcheng'	is not attract	ive and pr	ang onoses
02	, se ⇔ to	change it	to 'Tongs	heng Bar.' Wa	ing and Zhao ac	ree, but L	i
0.4	↔ op	poses; B:	In view of	the sluggish	n business, War	ng proposes	to
05	↔ su	spend oper	ations for	one month fo	or renovation a	and reorgan	ization
05	\hookrightarrow Z	hang and Z	Zhao agree,	but Li oppos	ses; C: Due to	the urgent	needs
07	$\rightarrow $ oI $\rightarrow $ h	the bar,	and Wang a	ses to sell a gree but Li	opposes: D: Gi	ven the for	s to the ur
07	a ↔	rtners' la	and wang a ack of expe	rience in bar	management, I	i proposes	to
00	\hookrightarrow ap	point his	friend Wan	g as the mana	ging partner.	Zhang and	Wang
10	↔ ag	ree, but Z	Zhao oppose	s.			
10	3-8: Res	ident A re	nted out th	ne house to B	. With A's con	sent, B rer	lovated
10	\hookrightarrow th	e rented f	nouse and s	ublet it to (Why can a me	ly altered	the
12	\rightarrow 10 \rightarrow 1i	ability fo	r breach o	f contract?	. Wily Call A Ie	equest b to	Deal
13							
14							
015	Open Aust	ralian Legal	QA The ber	hchmark consists	of 2,124 questions	and answers s	ynthesized
10	by GPT-4 f	rom the Aust	ralian legal co	rpus. All questic	ons are of the gener	ation type. Or	ie example
117	is: "What is	s the landlord	l's general obli	igation under sec	tion 63 of the Act i	n the case of A	Anderson v
10	Armitage [2	2014] NSWC	CATCD 157 in	New South Wale	s?"		
19	Evaluation	results are lis	sted in Table 9	where we only	compare to "Graph	RAG Local"	as it gener
20	ally perform	ns better than	"GraphRAG	Global" on these	tasks.	<u>-</u> ,	
21	E d C		1		1		
<u>12</u>	For the afor	rementioned	reasons, we a	iso use GPI-4 to	evaluate all experi-	mental results	, reporting
123	ine accurac	y (Acc) in la	ule IU. when a	comparing the res	suits in Table 9 and	Table 10, we o	userve that

Table 9: Evaluation Results on Legal Benchmarks (Metric is F1 / EM as indicated in Table 8)

the order of the results is preserved, even though some metrics change significantly. In the following
 section, we aim to identify the reasons behind these changes, which may provide valuable insights
 for designing better metrics to evaluate RAG frameworks in the future.

1026 1. The accuracy of our approach increases significantly for generation tasks (1-1, 3-2, Open 1027 Australian Legal QA). For these tasks, our answers are often semantically equivalent but 1028 syntactically different from the golden answers. This explains the improved metric perfor-1029 mance, as GPT-4 can compare the semantic content of the answers. This also applies to the 1030 "GraphRAG Local" results for the "Open Australian Legal QA" task. 1031 2. The accuracy of "GraphRAG Local" decreases for generation tasks 1-1 and 3-2. These 1032 tasks involve statute recitation and prediction, requiring the retrieval of specific articles. 1033 Upon detailed examination, We find that "GraphRAG Local" often fails to retrieve the 1034 correct articles or references the wrong ones, but it tends to repeat the legal information. 1035 Therefore, token-level recall can be improved by simply rephrasing legal names and com-1036 mon prefixes, such as "According to XX law, XX articles...". 1037 3. Both our approach and "GraphRAG Local" show significant accuracy improvements on task 3-8. Besides the reason mentioned in the first point, the quality of the golden answers 1039 may also contribute to this difference. The questions and golden answers in task 3-8 are 1040 sourced from a consulting website, resulting in varying quality. For example, one question 1041 asks "Do the children from the original marriage have an obligation to support the father?" 1042 However, the provided golden answer includes an irrelevant article, "Article 1067," which 1043 pertains to parents' obligations to support minor children. 1044 Question: In the case where both parents are divorced and have 1045 \hookrightarrow formed their own families with new children, and according 1046 \rightarrow to the court's judgment, the father is required to pay 1047 \hookrightarrow monthly child support to the mother until the child is 18 \hookrightarrow years old. Do the children from the original marriage have 1048 \hookrightarrow an obligation to support the father? 1049 Reference Answer: In our country, biological children have an 1050 \hookrightarrow obligation to support their parents who are divorced. The 1051 \hookrightarrow relationship between children and parents does not dissolve 1052 \hookrightarrow because of the divorce of the children or parents. \hookrightarrow Supporting parents is a legal obligation of the children. If 1053 \hookrightarrow the children do not support the elderly, the parents can 1054 \hookrightarrow directly file a lawsuit in the People's Court requesting the 1055 \hookrightarrow court to order the daughter to pay support. Legal basis: 1056 \hookrightarrow Article 1067 of the Civil Code of the People's Republic of 1057 \hookrightarrow China states that if parents do not fulfill their obligation \hookrightarrow to support, minor children or adult children who cannot 1058 \hookrightarrow live independently have the right to request support from \hookrightarrow their parents. If adult children do not fulfill their \hookrightarrow obligation to support, parents who are incapacitated or have 1061 difficulties in living have the right to request support \rightarrow 1062 \hookrightarrow from adult children. Article 1084 states that the \hookrightarrow relationship between parents and children does not dissolve 1063 \hookrightarrow because of the parents' divorce. After the divorce, whether 1064 \hookrightarrow the children are directly raised by the father or the mother 1065 \hookrightarrow , they are still the children of both parents. Lawyer's \hookrightarrow explanation: After the parents are divorced, adult children 1067 \hookrightarrow still have the obligation to support both parents. The 1068 \hookrightarrow obligation of adult children to support their parents does \rightarrow not change because of changes in the parents' relationship. 1069 \hookrightarrow If adult children do not fulfill their obligation to support 1070 \hookrightarrow , parents who are incapacitated or have difficulties in 1071 \hookrightarrow living have the right to request support from adult children \hookrightarrow . Our Answer: Yes, there is an obligation to support. Legal basis: \hookrightarrow According to Article 1069 of the Civil Code of the People's 1074 \hookrightarrow Republic of China, it is stipulated that: 'The obligation of 1075 \rightarrow children to support their parents does not terminate \hookrightarrow because of changes in the parents' marital relationship. 1077 \hookrightarrow Therefore, even if the parents are divorced and remarried \hookrightarrow with new children, the original children still have an 1078 \hookrightarrow obligation to support their parents. 1079



Figure 4: Case (a): Given the lesser-known film "What Women Love" as opposed to the more popular "What Women Want," single-path methods like Self-Ask on the left are predisposed to generating follow-up questions about the latter, leading to an incorrect final answer. Conversely, KAR³ can effectively discern the intended meaning of the original question by positing several atomic queries and postpone the task understanding to atomic selection phase with relevant atomic questions provided, and subsequently arriving at an accurate conclusion.

1115 1116

1117

1118

1119

4. The accuracy of all methods on choice tasks 1-2, 3-1, and 3-6 almost coincides with the F1 score, as expected. An exception is task 3-1, where the difference is mainly due to GPT-4's capacity to understand Chinese, particularly in distinguishing numbers in Arabic and Chinese. In Chinese law, all numbers are written in Chinese, while in the golden answers, all numbers are given in Arabic.

1120 1121

1122 A.5 REAL CASE STUDIES

This section presents three case studies from our evaluation benchmark to illustrate the underlying principles of our proposed decomposition pipeline, as detailed in Algorithm 1. Through these realworld examples, we aim to highlight the benefits of our systematic approach. These cases will shed light on how each step of the pipeline contributes to improved performance and the insights gained from their implementation.

Our task decomposition strategy involves generating multiple atomic queries rather than producing
a single deterministic follow-up question, as demonstrated in the Self-Ask approach. Contemporary decomposition methods typically employ a generative model to formulate a singular follow-up
question. However, this approach carries an intrinsic risk of generating erroneous questions, potentially leading to an incorrect decomposition pathway and, ultimately, an erroneous answer. Consider
the Case (a) depicted in Figure 4, where the original question pertains to a film titled "What Women



Figure 5: Case (b): By proposing multiple atomic queries, KAR³ effectively retrieves the relevant knowledge chunk, whereas the single deterministic follow-up question approach employed by Self-Ask fails to align with the knowledge base's schema, resulting in a retrieval failure.

1160

Love." Due to the existence of a more prominent film, "What Women Want," the employed language 1161 1162 model tends to 'correct' the original question. Consequently, methods like Self-Ask (as shown on the left side of Figure 4) generate only one follow-up question related to this erroneously assumed 1163 object. In the illustrated instance, although the target chunk has been retrieved due to the similarity 1164 in embeddings, a 'false' intermediate answer is produced for the 'false' follow-up question, culmi-1165 nating in an incorrect final response. In contrast, our methodology posits atomic queries concerning 1166 both "What Women Love" and "What Women Want," thereby seeking to clarify the true intent of 1167 the initial question. With both films existing and relevant atomic questions being retrieved, our ap-1168 proach subsequently gains the advantage of verifying the question's intent and selecting the correct 1169 and most pertinent chunk during the atomic selection phase. 1170

Furthermore, the discrepancy between the formulation of the corpus and the query, is another criti-1171 cal factor advocating for a multi-query approach over a singular deterministic one. The presentation 1172 gap can impede the retrieval process even when the generated follow-up question is semantically 1173 accurate. For instance, as illustrated in Case (b) in Figure 5, a single-path method such as Self-Ask 1174 on the left side might directly inquire 'Who is the mother of Oskar Roehler?' However, the knowl-1175 edge base articulates familial relationships using a different schema, 'A is the son of B and C' in this 1176 case, thus the retrieval process falters despite the correctness of the question. Even when we applied 1177 the hierarchical retrieval to Self-Ask, the Self-Ask with Hierarchical Retrieval did not succeed in 1178 bridging this gap. In contrast, our approach, which generates multiple atomic queries, encompasses 1179 a broader range of phrasings that correspond to the diverse representations in the knowledge base. In the depicted case, while the atomic query specifically asking for Oskar Roehler's mother encounters 1180 the same retrieval issue, an alternative query seeking information about his parents successfully re-1181 trieves the target chunk. This exemplifies how our method's flexibility in query generation enhances 1182 the likelihood of aligning with the knowledge base's structure and obtaining accurate information. 1183

Our methodology emphasizes the retrieval of atomic questions rather than directly retrieving chunks. This design choice is exemplified in Case (b) depicted in Figure 5. The knowledge chunk in the corpus is structured using the pattern 'A ... as the son of B and C', which poses challenges for direct retrieval by queries such as 'Who is the mother of ...'. In our specialized knowledge base, such direct queries tend to retrieve chunks conforming to the patterns 'A is the mother of B' or 'A is the father

1222 1223



Figure 6: Case (c): KAR³ has the advantage of leveraging a concise list of atomic questions for targeted selection and retaining full chunks for rich contextual support. Conversely, Self-Ask's approach, although successful in retrieving relevant chunks, is compromised by its dependency on intermediate answers for context, which ultimately results in the generation of incorrect final answers.

of B'. By utilizing atomic questions as intermediaries for retrieval, our approach effectively narrows
 the gap between a single query and the multiple sentence structures found in the knowledge base. It
 facilitates bridging the expression pattern differences exemplified by 'the mother of' versus 'the son of' in this scenario.

In contrast to methods like Self-Ask, which only retains intermediate answers for subsequent pro-1228 cessing, our method preserves the entire chunk as contextual information. During the atomic se-1229 lection phase, we present a list of atomic questions as candidate summaries of the relevant content 1230 from the original chunk. This strategy significantly reduces token usage and simplifies the process 1231 of selecting the pertinent information. Case (c) in Figure 6 demonstrates the dual benefits of our ap-1232 proach: first, by selecting from a curated list of atomic questions, we streamline the identification of 1233 relevant information; second, by retaining the entire selected chunk rather than just the intermediate 1234 answer, we ensure a rich context is maintained for accurate and comprehensive subsequent processing. While the Self-Ask method on the left retrieves the target chunk, it fails to correctly identify 1236 the pertinent 'Ernie Watts' due to the excessive contextual information. Since retrieved chunks in 1237 Self-Ask are discarded after generating an intermediate answer, the method potentially follows an incorrect pathway, leading to an inaccurate conclusion. In contrast, our approach can efficiently filter and select the appropriate atomic question from a concise list. Although the atomic question in 1239 this round pertains to the role of Ernie Watts, there is no need to inquire further about his birthplace, 1240 as this information is encapsulated within the selected chunk, which remains available for context 1241 in subsequent rounds.

Method	HotpotQA		2Wiki		MuSiQue		
Method	F1	Acc	F1	Acc	F1	Acc	
Zero-Shot Self-Ask w/ Retrieval	55.76	76.20	54.98	76.20	40.97	50.40	
Self-Ask w/ Retrieval	71.40	80.00	69.06	75.00	46.76	51.40	
Zero-Shot IRCoT	58.22	75.80	49.69	60.20	37.17	43.00	
IRCoT	67.30	81.00	63.83	70.40	47.57	49.20	

Table 11: Performance comparison: Zero-Shot vs. Few-Shot.

1252 A.6 PROMPT DESIGN

Our approach employs four distinct prompts, detailed at the end of the appendix. (1) Atomic ques-tion tagging prompt: the one used to pre-processing the source paragraphs that linking each para-graphs with several atomic questions; (2) Atomic query proposer prompt: the one used when gen-erating multiple atomic query proposals, referring to line 3 in Algorithm 1; (3) Atomic question selection prompt: the one used when selecting the most useful atomic question from the given ques-tion list, referring to line 5 in Algorithm 1; (4) Question answering prompt: the one applied upon exiting the decomposition loop to generate the final answer to the given question, as described in line 14 of Algorithm 1.

Demonstration Discussion In our current experiments, all prompts are zero-shot, meaning no demonstrations are provided to illustrate the expected reasoning logic. To explore whether demonstrations could enhance performance, we designed an ablation study. We adapted the Self-Ask w/ Retrieval and IRCoT methodologies previously employed, modifying the prompts and task descriptions to create zero-shot, demonstration-free variants of these methods. These were denoted as Zero-Shot Self-Ask w/ Retrieval and Zero-Shot IRCoT. The results of the experiment are presented in Table 11.

The experimental results reveal that the Zero-Shot Self-Ask w/ Retrieval method experiences a marginal decline in accuracy for the 2WikiMultiHopQA and MuSiQue datasets, potentially due to the inherent randomness in generation. However, the inclusion of demonstrations significantly improves all F1 scores and enhances the overall performance of the IRCoT method. This suggests that demonstrations could be particularly beneficial for methods that rely on a step-by-step decomposition approach. Consequently, integrating demonstrations is identified as a promising direction for future work within the KAR³ framework.

```
1296
        Atomic Question Tagging Prompt
1297
        # Task
1298
        Your task is to extract as many questions as possible that are relevant
1299
        and can be answered by the given content. Please try to be diverse
1300
        and avoid extracting duplicated or similar questions. Make sure your
        question contain necessary entity names and avoid to use pronouns like
1301
        it, he, she, they, the company, the person etc.
1302
1303
        # Output Format
1304
        Output your answers line by line, with each question on a new line,
1305
        without itemized symbols or numbers.
1306
        # Content
1307
        {content}
1308
1309
        # Output
1310
1311
        Atomic Query Proposer Prompt
1312
        # Task
1313
        Your task is to analyse the providing context then raise atomic
1314
        sub-questions for the knowledge that can help you answer the question
        better. Think in different ways and raise as many diverse questions as
1315
        possible.
1316
1317
        # Output Format
1318
        Please output in following JSON format:
1319
        {{
            "thinking": <your thinking for this task, including analysis to
1320
        the question and the given context>,
1321
            "sub_questions": <a list of sub-questions indicating what you
1322
        need>
1323
        }}
1324
        # Context
1325
        The context we already have:
1326
        {chosen_content}
1327
1328
        # Question
1329
        {content}
1330
        # Your Output
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
```

```
1350
        Atomic Question Selection Prompt
1351
        # Task
1352
        Your task is to analyse the providing context then decide which
1353
        sub-questions may be useful to be answered before you can answer
1354
        the given question. Select a most relevant sub-question from the
        given question list, avoid selecting sub-question that can already
1355
        be answered with the given context or with your own knowledge.
1356
1357
        # Output Format
1358
        Please output in following JSON format:
1359
        {{
            "thinking": <your thinking for this selection task>,
1360
            "question_idx": <a sub-question index, an integer from 1 to
1361
        {num_atom_questions}>
1362
        }}
1363
1364
        # Context
        The context we already have:
1365
        {chosen_content}
1366
1367
        # Sub-Questions You Can Choose From
1368
        {atom_question_list_str}
1369
        # Question
1370
        {content}
1371
1372
        # Your Output
1373
1374
        Question Answering Prompt
1375
        # Task
1376
        Your task is to answer a question referring to a given context, if
1377
        any. For answering the Question at the end, you need to first read the
1378
        articles, reports, or context provided, then give your final answer.
1379
        # Output format
1380
        Your output should strictly follow the format below. Make sure your
1381
        output parsable by json in Python.
1382
        {{
1383
            "answer": <Your Answer, format it as a string.>,
            "rationale": <rationale behind your choice>
1384
        }}
1385
1386
        # Context, if any
1387
        {context_if_any}
1388
        # Question
1389
        {content}{yes_or_no_limit}
1390
1391
        Let's think step by step.
1392
1393
1394
1395
1396
```

1399

1400

1401

1402