

# MINTQA: A Multi-Hop Question Answering Benchmark for Evaluating LLMs on New and Long-tail Knowledge

Anonymous ACL submission

## Abstract

Retrieval-Augmented Generation (RAG) enhances Large Language Models (LLMs) by integrating external knowledge, enabling them to tackle knowledge-intensive tasks. However, limited research has explored how LLMs effectively leverage RAG techniques for multi-hop Question Answering (QA), particularly when handling knowledge with varying degrees of familiarity. In this paper, we introduce **MINTQA** (Multi-hop Question Answering on New and longTail Knowledge), a benchmark designed to evaluate multi-hop QA performance on questions involving 10,479 question-answer pairs for evaluating old/new knowledge and 17,887 pairs for assessing popular/unpopular knowledge, with each question equipped with its sub-questions and answers. This benchmark primarily evaluates the multi-hop reasoning ability of LLMs and their capacity to handle knowledge with varying levels of familiarity during the reasoning process. We evaluate 22 state-of-the-art LLMs using three distinct QA strategies: LLM-based parameterized knowledge QA, direct RAG-enhanced QA, and multi-hop RAG-enhanced QA. Our experiments reveal key challenges in how LLMs handle knowledge with different familiarity and offer insights into improving their multi-hop reasoning capabilities when combined with RAG techniques.<sup>1</sup>

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in Question Answering (QA) (Kamalloo et al., 2023; Wang and Qin, 2024). However, they face significant challenges when handling questions requiring specific domain knowledge or up-to-date information (Pan et al., 2023). Although Retrieval-Augmented Generation (RAG) provides an effective strategy for generating responses by incorporating external knowledge

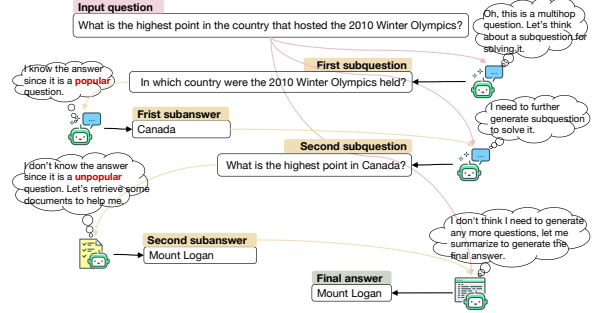


Figure 1: **An example for our benchmark:** Given a complex question, the model must decide whether to decompose it into sub-questions and determine if external knowledge retrieval is required.

(Soudani et al., 2024; Islam et al., 2024), optimizing the combination of LLMs and RAG remains a critical challenge, particularly for multi-hop QA tasks involving knowledge with varying familiarity for models.

Consider a complex question: “What is the highest point in the country that hosted the 2010 Winter Olympics?” As illustrated in Figure 1, answering such questions necessitates decomposing them into sub-questions, such as: (1) In which country were the 2010 Winter Olympics held? and (2) What is the highest point in Canada? Each sub-question may require different knowledge sources. For instance, the first can leverage parametric knowledge, while the second needs related knowledge retrieval.

Current benchmarks for evaluating LLMs on such multi-hop QA scenario have several limitations. First, studies such as (Sun et al., 2023; Maekawa et al., 2024; Zhang et al., 2024) primarily focus on single-hop queries, leaving complex multi-hop questions largely unexplored. Second, while multi-hop benchmarks such as MultiHop-RAG (Tang and Yang, 2024) assess retrieval effectiveness, they do not systematically evaluate the interaction between question decomposition and retrieval, a capability essential for real-world applications. Furthermore, existing works like

<sup>1</sup>The MINTQA benchmark is available at <https://anonymous.open.science/r/multi-hop-D70E/>.

FanoutQA (Zhu et al., 2024a) and HotpotQA (Yang et al., 2018) lack assessment of how models handle queries containing new or unpopular knowledge, which presents unique challenges in both decomposition and retrieval.

To address these gaps and facilitate the effective integration of LLMs and RAG techniques to handle different knowledge during multi-hop QA reasoning, we propose MINTQA, a benchmark for evaluating LLMs on complex multi-hop questions across two critical dimensions: **Unpopular knowledge** (information appearing infrequently in training corpora) and **New Knowledge** (recently emerged entities or relationships). Figure 2 outlines our benchmark construction and the evaluation framework based on the benchmark.

We construct MINTQA by systematically collecting knowledge triplets based on Wikipedia and Wikidata and using GPT-4o to generate multi-hop questions spanning one to four hops. The benchmark comprises two sub-datasets: **MINTQA-POP** (17,887 examples) focusing on unpopular/popular knowledge, and **MINTQA-TI** (10,479 examples) examining new/old knowledge, with each example including sub-questions and answers for fine-grained analysis of models’ reasoning processes. Appendix D presents a comparison of MINTQA with existing benchmarks, highlighting its unique contributions to multi-hop QA evaluation.

Our framework evaluates LLMs across five aspects: using parametric knowledge, retrieval-augmented generation, sub-question generation, and direct or dynamic decomposition-retrieval. Our comprehensive evaluation of 22 state-of-the-art LLMs reveals: **First**, performance differs across knowledge with varying familiarity for the models, and the effectiveness of retrieval and question decomposition strategies also varies for different knowledge. **Second**, all LLMs struggle with multiple reasoning hops. **Third**, generating sub-questions alone doesn’t significantly improve performance and might even degrade results. **Fourth**, retrieval effectiveness decreases with question complexity. **Fifth**, combining decomposition and retrieval benefit models larger than 14B parameters but do not improve smaller models, suggesting smaller models need improved planning and retrieval capabilities. **Sixth**, our implementation of dynamic retrieval methods (Ni et al., 2024) shows that maintaining performance while reducing retrieval frequency remains challenging, with some models showing excessive retrieval dependency.

**Our main contributions are summarized as follows:**

1. To address limitations in existing benchmarks, we introduce MINTQA, a novel benchmark for evaluating LLMs for multi-hop reasoning across knowledge of varying familiarity. The best-performing LLaMA 3.1-70B’s 62.33% accuracy demonstrates both our benchmark’s difficulty and the remaining challenges in complex multi-hop reasoning.
2. We present a systematic evaluation framework examining key aspects of multi-hop QA, enabling comprehensive analysis of models’ reasoning capabilities and the effectiveness of different retrieval and sub-questions generations strategies.
3. Our evaluation of 22 state-of-the-art LLMs reveals their multi-hop reasoning limitations, especially for questions with new or unpopular knowledge. LLMs still struggle to identify the appropriate strategy to help correctly answer questions requiring different knowledge sources. The insights from our evaluation are valuable for improving LLM’s multi-hop QA capabilities.

## 2 Related Work

### 2.1 Multi-hop Question Answering (QA)

Multi-hop QA challenges LLMs by requiring synthesis and reasoning across multiple sources (Huang and Chang, 2023; Feng et al., 2020; Khashabi et al., 2019). While researchers have proposed decomposing complex questions into sub-questions (Min et al., 2019; Wang et al., 2022, 2023; Liu et al., 2024), generating relevant sub-questions and reasoning chains remains challenging. Existing benchmarks (Zhang et al., 2024; Zhu et al., 2024a; Tang and Yang, 2024) assess retrieval and multi-hop reasoning, but overlook when and how to retrieve, interactions between decomposition and retrieval, or queries with new and unpopular knowledge. Our MINTQA fills these gaps by systematically evaluating LLMs’ on multi-hop QA.

### 2.2 Retrieval Augmented Generation (RAG)

RAG enhances LLMs’ performance in multi-question answering by providing access to external documents (Lewis et al., 2020; Xiong et al., 2020), particularly for knowledge-intensive tasks (Yu et al., 2020; Zhu et al., 2023). In sub-question generation, RAG can verify and correct LLMs’ out-

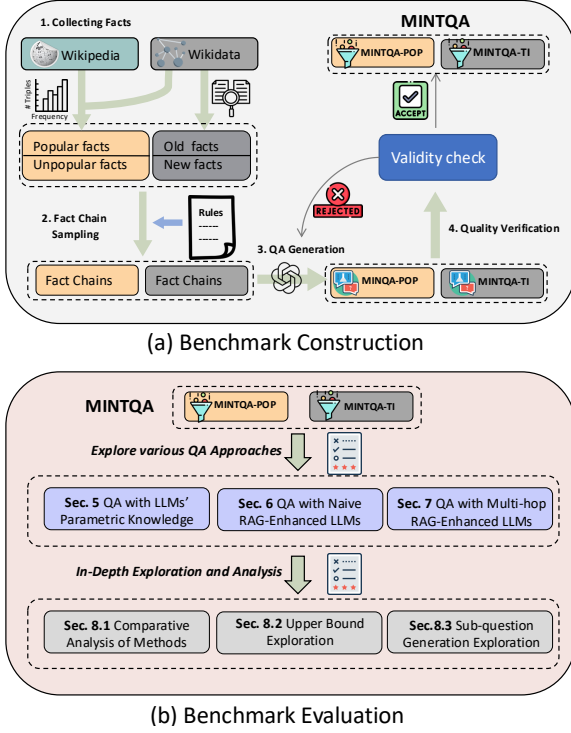


Figure 2: Illustration of the MINTQA construction process (a) and the evaluation framework leveraging MINTQA (b).

puts (Zhao et al., 2023; Shi et al., 2024a). However, irrelevant retrievals can introduce noise (Yoran et al., 2024; Joren et al., 2024), and external knowledge may override model’s inherent knowledge (Xu, 2023; Li et al., 2022), while adding computational overhead (Zhu et al., 2024b). While Jeong et al. (2024) propose using a classifier to determine retrieval necessity, our research investigates LLMs’ inherent ability to recognize when retrieval is needed for sub-questions.

## 2.3 Evaluation of LLMs-based QA

Existing QA datasets for evaluating retrieval-augmented LLMs fall into two categories: (1) Reasoning-focused datasets (Ho et al., 2020; Yang et al., 2018; Sen et al., 2022), such as MuSiQue (Trivedi et al., 2021), FanOutQA (Zhu et al., 2024a), and MultiHop-RAG (Tang and Yang, 2024) that emphasize multi-hop reasoning across multiple documents; (2) Long-tail question datasets (Mallen et al., 2022; Zhang et al., 2024), including WiTQA (Maekawa et al., 2024) focusing on rare single-hop queries and Head-to-Tail (Sun et al., 2023) examining entity and relationship popularity to highlight the value of knowledge graphs. Our work extends these by evaluating both long-tail and new-fact multi-hop QA, while analyzing models’

sub-question generation and retrieval capabilities.

## 3 Benchmark Construction

This section presents our comprehensive methodology for constructing two multi-hop QA benchmarks: MINTQA-POP and MINTQA-TI, designed to evaluate LLM across two critical dimensions: knowledge popularity (popular versus unpopular) and knowledge freshness (new versus old). We first present the data construction methodology for MINTQA-POP (Section 3.1). We then detail the construction process of MINTQA-TI, which follows a similar procedure but focuses on new/old knowledge (Section 3.2). Finally, we describe our QA generation process (Section 3.3) and present comprehensive statistics of our constructed datasets (Section 3.4).

### 3.1 Data Construction of MINTQA-POP

**Collecting Facts** We gather a collection of facts with popularity, denoted as  $\mathcal{G}_{pop} = \{(s, r, o), p | (s, r, o) \in \mathcal{G}, p \in \mathbb{Z}^+\}$ , where  $\mathcal{G}$  refers to Wikidata,  $(s, r, o)$  represents a triple,  $p$  indicates the popularity as the positive integer set  $\mathbb{Z}^+$ . The triples are extracted from Wikipedia (version 2024-05-01). Specifically, we extract raw triples in the format of (Head Span, Relation, Tail Span) from Wikipedia passages using an existing information extraction tool<sup>2</sup>. These raw triples are linked to Wikidata (version 2024-04-22) using WikiMapper<sup>3</sup>, producing structured triples with Wikidata IDs  $(s, r, o)$ . We only keep the triples  $(s, r, o)$  existing in Wikidata. The popularity  $p$  of each triple is calculated as the frequency of its occurrence across the entire Wikipedia corpus.

**Sampling fact chains** We sample facts from  $\mathcal{G}_{pop}$  and concatenate them into a chain  $\mathcal{FC} = \{(s_1, r_1, o_1), \dots, (s_n, r_n, o_n)\}$  as the grounded facts of a multi-hop question. We categorize facts in  $\mathcal{G}_{pop}$  based on their popularity scores into two distinct sets: **unpopular** knowledge ( $\mathcal{P}_{unpop} = [1, 10)$ ) and **popular** knowledge ( $\mathcal{P}_{pop} = [50, \infty)$ ). A fact chain  $\mathcal{FC}$  is constructed as an ordered sequence of connected triples:  $\mathcal{FC} = \{(s_1, r_1, o_1), \dots, (s_n, r_n, o_n)\}$ , where  $n \leq 4$  and each triple can be either popular or unpopular. This construction follows four key constraints:

1. **Connectivity:**  $o_i = s_{i+1}$  for all  $i \in \{1, \dots, n-1\}$ .

<sup>2</sup><https://github.com/Babelscape/rebel>

<sup>3</sup><https://github.com/jcklie/wikimapper>

2. **Acyclicity:**  $o_i \neq s_j$  for all  $i, j \in \{1, \dots, n\}$ .
3. **Uniqueness:** No fact chain  $\mathcal{FC}$  can be a sub-chain of another fact chain.
4. **No Shortcuts:** For each fact chain  $\mathcal{FC}$ , there does not exist a triple  $(s_i, r, o_j)$  in  $\mathcal{G}_{pop}$  such that  $j > i + 1$ , where  $i \in \{1, \dots, n - 1\}$  and  $j \in \{2, \dots, n\}$ .
5. **Single Object:** For each triple  $(s_i, r_i, o_i)$  in the fact chain, there does not exist another triple  $(s_i, r_i, o_j)$  in Wikidata with a different object.

### 3.2 Data Construction of MINTQA-TI

Building on the methodology established for MINTQA-POP, we construct MINTQA-TI, focusing on old and new knowledge. To construct the dataset, we extract two versions of Wikidata: 2021-06-21 and 2024-06-05. We identify triples that are either common to both versions or differ between them. These triples form the knowledge graph  $\mathcal{G}_{ti}$ . We define **old** knowledge as triples present in both Wikidata versions, and **new** knowledge as triples only appearing in the newer version, characterized by a new subject, relation, or object. Following the same chain construction principles outlined in Section 3.1, we create fact chains combining new and old knowledge from  $\mathcal{G}_{ti}$ .

### 3.3 QA Generation and Verification

Following WitQA (Maekawa et al., 2024), we employ GPT-4o to automatically generate questions from extracted triplets, overcoming the diversity and scalability issues of template-based methods like PopQA (Mallen et al., 2022) and the high costs of manual annotation. Given a fact chain  $\mathcal{FC} = \{(s_1, r_1, o_1), \dots, (s_n, r_n, o_n)\}$ , where  $o_{i-1} = s_i, i \in \{2, \dots, n\}$ , we aim to generate a question about  $s_1$  that yields  $o_n$  as the answer. To enhance generation quality, we provide one demonstration example per hop. And to ensure validity, we verify questions by having the model answer them using source contexts; only questions yielding  $o_n$  are retained. Invalid questions are re-generated up to three times, and unsatisfactory examples are discarded. For multi-hop questions (hop count  $\geq 2$ ), sub-questions for each intermediate fact are also generated and validated. Examples are included in the dataset only if the main question and all sub-questions pass validation. This process filtered out 138 and 67 examples from MINTQA-POP and MINTQA-TI, respectively. Prompts and examples are in Appendices C and E.

MINTQA-POP	1-hop	2-hop	3-hop	4-hop	Total
#Samples	5,894	4,428	4,664	2,901	17,887
#Input Tok.	52,488	59,699	74,363	57,579	244,129
#Input Vocab	7,468	5,410	7,094	4,491	18,852
Avg. In. Len.	8.91	13.48	15.94	19.85	13.65
#Output Tok.	11,526	6,249	6,196	5,610	29,581
#Output Vocab	3,398	420	228	210	3,721
Avg. Out. Len.	1.96	1.41	1.33	1.93	1.65
Avg. Ctx. Len.	32.93	375.93	549.18	706.12	361.63
#Relations	124	84	85	98	140
#Entities	7,482	4,357	5,191	3,180	18,501
MINTQA-TI					
#Samples	3,949	2,198	2,057	2,275	10,479
#Input Tok.	34,014	30,471	34,106	44,959	143,550
#Input Vocab	5,586	3,113	2,439	2,695	10,345
Avg. In. Len.	8.61	13.86	16.58	19.76	13.70
#Output Tok.	8,064	5,738	4,365	5,086	23,253
#Output Vocab	1,710	1,004	996	809	3,318
Avg. Out. Len.	2.04	2.61	2.12	2.24	2.22
Avg. Ctx. Len.	56.44	321.32	487.22	645.61	324.47
#Relations	123	147	149	147	189
#Entities	4,096	2,484	2,250	2,346	9,616

Table 1: Data statistics of MINTQA.

### 3.4 Dataset Statistics

Table 1 summarizes the statistics of the MINTQA-POP and MINTQA-TI datasets, which exhibit diverse coverage across multiple dimensions. MINTQA-POP contains 17,887 examples and MINTQA-TI 10,479, with over 2,000 examples per hop category, ensuring robust evaluation. The datasets include 18,501 and 9,616 entities, and 140 and 198 relationships, respectively, demonstrating their diversity. As the number of hops increases, the average context length grows, requiring models to retrieve more documents and face greater challenges. For details, see the App. A.

## 4 Experimental Setup

### 4.1 Language Models and Configurations

**Models** We evaluate state-of-the-art LLMs across various architectures and model sizes: GPT-3.5, GPT-4o, GPT-4o mini, LLaMA-3.1/3.2 (Grattafiori et al., 2024), Gemma-2 (Team et al., 2024), Mistral (Jiang et al., 2023), Phi-3 (Abdin et al., 2024), and Qwen2.5 (Hui et al., 2024). All models are instruct versions. For simplicity, we omitted the “instruct” name in the result presentation. To ensure reproducibility, we set the temperature parameter to 0 across all models and accelerated inference using vLLM (Kwon et al., 2023). For more details, please refer to Appendix F.

### 4.2 Evaluation Metrics

We adopt *Accuracy* (Acc) as our evaluation metric across all experiments by determining whether the ground-truth answer is present in the model’s predicted text across all experiments, consistent with established benchmarks in factual knowledge assessment (Ren et al., 2023; Maekawa et al., 2024; Mallen et al., 2022).



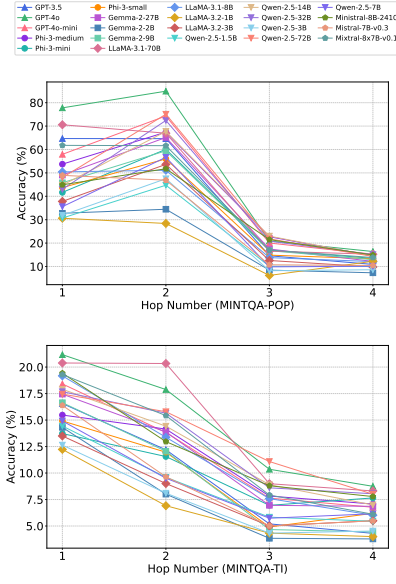


Figure 3: Zero-shot accuracy of different LLMs across various hops.

### 4.3 Research Questions

Based on the two proposed multi-hop QA sub-datasets, MINTQA-POP, which involves popular/unpopular knowledge, and MINTQA-TI, involving old/new knowledge, we investigate the following research questions:

$RQ_1$ : How do various LLMs perform in the two multi-hop QA scenarios relying solely on their internal knowledge? (Section 5)  $RQ_2$ : How do LLMs perform, when enhanced with a direct retrieval approach, in the multi-hop QA scenarios, and how do different retrievers perform? (Section 6)  $RQ_3$ : How do LLMs enhanced with a multi-hop RAG strategy (decomposition-then-retrieval approach) perform in the multi-hop QA scenarios, particularly when handling popular/unpopular and old/new knowledge during the decomposition and retrieval processes? (Section 7.1)  $RQ_4$ : Whether decomposition-dynamic retrieval can help achieve an optimal balance between performance and efficiency, given the different knowledge types (popular/unpopular, old/new knowledge) involved in multi-hop QA? (Section 7.2)

## 5 LLMs' Performance on MINTQA with Parametric Knowledge

We evaluate LLMs on MINTQA using their parametric knowledge to understand intrinsic model capabilities and dataset challenges. Specifically, we prompt LLM using its own knowledge to directly answer questions, as shown in Table 14. The

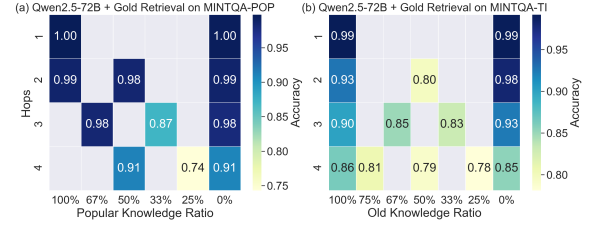


Figure 4: The performance of Qwen2.5-72B with gold retrieval across two datasets. The X-axis represents the proportion of popular knowledge required in the question, and the Y-axis indicates question hops.

results are shown in Figure 3 and further elaborated in App. G.1. Our findings reveal **significant performance gaps between the MINTQA-POP and MINTQA-TI**. Models perform reasonably on MINTQA-POP (e.g., GPT-4o: 77.79%, LLaMA3.1-8B: 50.42% for single-hop questions) but struggle on MINTQA-TI, with GPT-4o's accuracy dropping to 21.17% for single-hop questions. This confirms MINTQA-TI's effectiveness in evaluating knowledge beyond training data, and low performance across models from LLaMA-3.2-1B (7.78%) to GPT-4o (21.17%) demonstrates scaling model size alone doesn't address this. Moreover, **increased reasoning complexity further highlights these limitations**. On MINTQA-POP, performance drops sharply for three-hop (20.03%) and four-hop (16.41%) questions, while on MINTQA-TI, accuracy consistently declines with complexity.

## 6 Effectiveness of Direct Retrieval

After analyzing the performance of LLMs using only their parametric knowledge on MINTQA in Section 5, we find that LLMs struggle to answer MINTQA questions independently. In this section, we follow prior work (Mallen et al., 2022; Maekawa et al., 2024) to assess the effectiveness of the naive RAG approach, direct retrieval, when applied to LLMs for handling our complex multi-hop questions. For each retrieval, we select the top-5 passages that are relevant to the question and input them as context. The prompt is shown in Table 15.

### 6.1 Retrieval Model Setup

**Retrieval Models** We evaluate seven retrieval approaches across three categories: 1) Sparse retriever: **BM25** (Robertson and Zaragoza, 2009). 2) Vector retrievers pre-trained on large unlabeled corpora: **Contriever** (Izacard et al., 2021), **GTR-LARGE/XL** (Ni et al., 2021) and **BGE** (Xiao et al., 2023). 3) Instruction-tuned text embedding

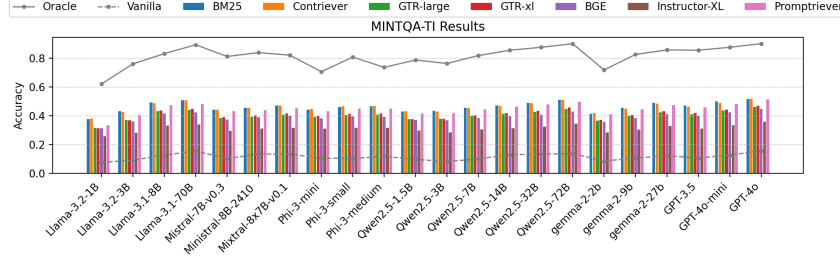


Figure 5: Performance comparison of LLMs on MINTQA-TI using different retrieval methods: “Oracle” uses gold-standard retrieval passages, while “Vanilla” involves models answering without retrieval content.

retrievers: **Instructor-XL** (Su et al., 2022) and **Promptriever** (Weller et al., 2024)..

**Configuration** We follow the approach of Yu et al. (2023) to construct the retrieval corpus by linearizing the knowledge graph  $\mathcal{G}$  into text.  $\mathcal{G}$  consists of  $\mathcal{G}_{pop}$  (Section 3.1) and  $\mathcal{G}_{ti}$  (Section 3.2). See Appendix F.2 for details.

## 6.2 Performance Analysis

Figure 5 demonstrates that **retrieval significantly enhances performance**, especially on MINTQA-TI, with an average 30% accuracy gain over the Vanilla setting (no retrieval). Similar trends are observed on MINTQA-POP (refer to Figure 14). Notably, in the Oracle setting, where gold-standard passages are used, even small models like Llama-3.2-1B achieve a 25% accuracy improvement compared to the average performance of all retrievers we used, emphasizing the potential for better retrievers. Appendix G.3 provides more analysis of the retriever.

We analyze the impact of knowledge popularity and newness on QA performance. **Models with different retrievers show inconsistent patterns when varying proportions of new and popular knowledge.** To isolate retrieval quality, we pair models with gold retrieval. Figures 4(a) and (b) show that with Qwen2.5-72B with gold retrieval, performance initially declines and then improves as the proportion of popular or old knowledge decreases. This likely occurs because the model effectively determine whether using parametric knowledge and retrieval for fully familiar (100% popular/old) or unfamiliar (100% unpopular/new) questions but struggles with mixed knowledge, leading to errors. Further analyses are in Appendix G.3.

## 7 Enhancing Multi-hop QA through Integrating Decomposition and Retrieval

Li and Peng (2023) and Shi et al. (2024b) highlight the importance of effectively combining question decomposition and retrieval for solving multi-hop questions. In this section, we explore two advanced RAG strategies (Decomposition-then-Retrieval and Decomposition-Dynamic Retrieval).

### 7.1 Decomposition-then-Retrieval

Building on prior work (Li and Peng, 2023; Shi et al., 2024b), we implement an iterative **decomposition-then-retrieval (DTR)** approach for multi-hop QA. At each step, the LLM determines one of two options: (1) further decomposing the question into sub-questions, or (2) summarizing a final answer based on the results of previously resolved sub-questions. If option (1) is selected, the LLM uses the full history of sub-questions and their corresponding answers as context to generate a new sub-question. Subsequently, five relevant documents are retrieved to assist in answering the new sub-question. If option (2) is chosen, the LLM summarizes the entire history, including all previously resolved sub-questions and their results, to produce a final answer. This iterative process concludes when either option (2) is selected or a maximum of five iterations is reached. The prompts used for both decision-making and summarization are shown in Tables 19 and 20. We evaluate this approach using BM25, Contriever, and Promptriever<sup>4</sup>.

Figure 6 shows patial results and Figure 13 shows the full results. On MINTQA-POP, larger models (>14B) benefit from decomposition and retrieval compared to direct retrieval, while smaller

<sup>4</sup>GPT models were excluded due to high cost and limited performance advantages over open-source LLMs (70B+).

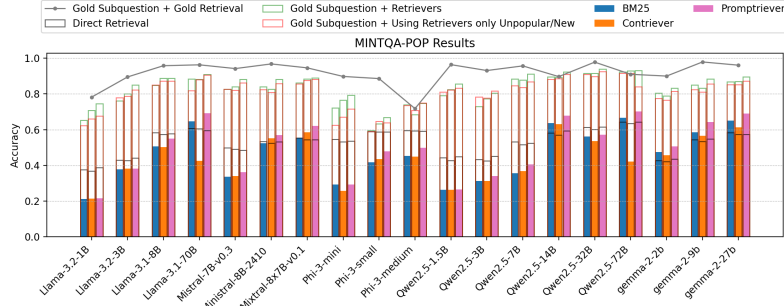


Figure 6: Performance comparison of all models using three retrievers under the decomposition-then-retrieval (DTR) approach on the MINTQA-POP dataset (represented by bars with three colors). **Gold Subquestion + Gold Retriever** indicates that the model utilizes gold subquestions and gold retrieval results. **Gold Subquestion + Retrievers** indicates that the model uses gold subquestions and employs different retrievers for retrieval. **Gold Subquestion + Using Retrievers Only Unpopular/New** denotes that the model uses gold subquestions and retrieves only for subquestions involving unpopular or new knowledge, while relying on the model itself to directly answer subquestions involving popular or old knowledge. **Direct Retrieval** refers to the use of different retrievers for direct retrieval (see Section 6) instead of adopting the DTR approach. See full results in Figure 13.

Model	BM25		
	Acc (%)	Avg. Sub	Avg. Ret
MINTQA-POP			
Qwen2.5-1.5B	25.86 (-0.5)	1.13 (1.0)	0.32 (0.19)
Qwen2.5-3B	31.16 (-0.13)	1.78 (0.95)	1.54 (0.71)
Qwen2.5-7B	32.98 (-2.68)	2.18 (0.97)	1.14 (-0.07)
Qwen2.5-14B	53.77 (-10.02)	3.44 (1.1)	1.22 (-1.12)
Qwen2.5-32B	50.33 (-5.86)	2.79 (1.02)	1.18 (-0.59)
Qwen2.5-72B	58.63 (-7.93)	3.01 (1.1)	1.34 (-0.57)
LLaMA-3.2-1B	20.53 (-0.75)	1.79 (1.0)	1.36 (0.57)
LLaMA-3.2-3B	37.23 (-0.64)	3.48 (0.72)	3.26 (0.5)
LLaMA-3.1-8B	50.01 (-0.72)	3.88 (0.56)	3.79 (0.47)
LLaMA-3.1-70B	64.80 (0.17)	3.41 (0.89)	3.39 (0.87)
Mistral-7B-v0.3	29.47 (-4.23)	3.13 (0.76)	1.84 (-0.53)
Mistral-8B-2410	35.91 (-16.94)	2.95 (0.96)	0.02 (-1.97)
Mistral-8x7B-v0.1	48.28 (-7.08)	3.61 (0.77)	1.29 (-1.55)
Phi-3-mini	26.81 (-2.54)	4.75 (0.04)	2.23 (-2.58)
Phi-3-small	37.09 (-4.58)	2.92 (0.44)	0.78 (-1.7)
Phi-3-medium	40.16 (-5.24)	2.98 (0.87)	1.08 (-1.03)
Gemma-2-2B	34.20 (-13.31)	4.96 (0.13)	0.98 (-3.85)
Gemma-2-9B	39.99 (-18.52)	3.93 (0.18)	0.32 (-3.43)
Gemma-2-27B	64.64 (-0.37)	4.05 (0.74)	4.01 (0.7)

Table 2: The results for Decomposition-Dynamic Retrieval approach. **Acc** represents the accuracy of the model in answering questions, **Avg. Sub** indicates the average number of sub-questions generated by the model, **Avg. Ret** refers to the average number of sub-questions that are deemed necessary for retrieval by the model. The value in brackets indicates the value of DDR minus that of DTR.

models (<8B) perform worse due to decomposition errors. On MINTQA-TI, direct retrieval outperforms decomposition-then-retrieval for most models, suggesting new knowledge poses greater challenges than question decomposition.

## 7.2 Decomposition-Dynamic Retrieval

The iterative DTR strategy in Section 7.1 faces two key challenges: high computational overhead from repeated retrievals (Zhuang et al., 2024) and performance degradation from unnecessary retrievals (Mallen et al., 2022; Maekawa et al., 2024). To ad-

Model	MINTQA-POP		MINTQA-TI	
	(1)	(2)	(1)	(2)
LLaMA-3.2-1B	19.62	41.89	6.88	14.23
LLaMA-3.2-3B	27.02	59.93	7.42	18.26
LLaMA-3.1-8B	37.28	70.69	9.82	23.75
LLaMA-3.1-70B	54.08	72.97	16.73	23.15
Mistral-7B-v0.3	31.41	58.80	9.52	16.00
Mistral-8B-2410	35.76	59.61	10.54	17.44
Mistral-8x7B-v0.1	45.66	68.90	12.08	20.93
Phi-3-mini	21.41	37.61	5.84	11.46
Phi-3-small	26.98	38.53	8.09	13.52
Phi-3-medium	38.91	64.29	10.18	19.45
Qwen2.5-1.5B	25.74	59.29	9.07	18.33
Qwen2.5-3B	23.86	55.90	6.99	15.87
Qwen2.5-7B	26.32	55.83	8.93	17.85
Qwen2.5-14B	41.34	65.67	12.96	20.01
Qwen2.5-32B	39.16	63.44	12.51	19.79
Qwen2.5-72B	44.99	54.62	14.10	20.51
Gemma-2-2B	25.67	49.97	8.04	14.32
Gemma-2-9B	38.18	59.65	11.17	18.09
Gemma-2-27B	44.54	66.29	12.36	18.90

Table 3: The accuracy of LLMs evaluated under query decomposition settings: (1) the model generates and answers sub-questions itself, and (2) the model answers given gold sub-questions.

dress this, we explore the **decomposition-dynamic retrieval (DDR)** approach that whether LLMs can dynamically determine retrieval necessity. Following (Ni et al., 2024), we implement a confidence-guided retrieval mechanism, where LLM determines whether to directly answer a sub-question or adopt the retrieval action when it has low confidence for answering the sub-question (details in App. F). Table 2 shows some results, with complete results in App. G.4. The prompt used for this experiment is shown in Table 21.

Our analysis reveals two key findings. **First**, reducing retrievals while maintaining performance proves challenging, with only the largest models (LLaMA-3.1-70B and Gemma-2-27B) maintain-

ing accuracy despite high retrieval rates (>98%). Other models show significant performance drops, reflecting our datasets’ emphasis on rare and new information. **Second**, models exhibit varying retrieval dependencies. Mistral and Phi models show high self-confidence (55% retrieval rate), LLaMA variants consistently trigger retrieval (>90%), while Gemma models exhibit size-dependent behavior, with retrieval rates ranging from <10% (2-9B) to >98% (2-27B) on MINTQA-POP.

## 8 More Investigations

### 8.1 Comparison of DTR and DDR

In Tables 2 and 11, we can compare the performance of the DTR strategy and the DDR strategy in terms of accuracy (Acc), average number of sub-questions (Avg.Sub), and average number of retrieved sub-questions (Avg.Ret). In terms of accuracy, most models perform worse with the DDR strategy compared to the DTR strategy. Under DDR, models generate more sub-questions than with DTR, but retrieve for fewer sub-questions on average. This suggests that while DDR improves the retrieval efficiency of the model, it comes at the cost of accuracy. The observed accuracy drop indicates that, for current models on the MINTQA dataset, there is a trade-off between accuracy and efficiency that needs further improvement.

### 8.2 Oracle Analysis with Gold Component

We evaluate system limitations using gold-standard sub-questions and their retrieved documents. Figure 6 shows notable gains across all models and retrievers when using gold sub-questions (i.e. Gold Subquestion + Retrievers), especially for smaller LLMs, highlighting their difficulties in generating accurate sub-questions independently. Additionally, previous work (Mallen et al., 2022; Maekawa et al., 2024) has shown that retrieval introduces noise that hurts QA performance when answering popular single-hop questions, and we verify this conclusion in the complex multi-hop question scenario by having LLM perform retrieval only on unpopular or new questions. The results (i.e. Gold Subquestion + Using Retrievers only Unpop/New v.s. Gold Subquestion + Retrievers) show that answering complex multi-hop questions generally yields better QA performance when using a common retriever that always performs the retrieval operation.

Notably, even with perfect decomposition and

relevant documents (i.e. Gold Subquestion + Gold Retrieval), the accuracy of various LLM remains below 100%. This reveals two challenges: extracting relevant information from documents containing multiple facts and synthesizing information across sub-questions, suggesting areas for future improvement beyond retrieval and decomposition.

### 8.3 Effectiveness of Only Decomposition

In this section, we investigate whether only generating and answering sub-questions or providing sub-questions for answering improves the accuracy on our benchmark. Results can be seen in Table 3.

**Self-Generated Sub-Questions:** On MINTQA-POP, self-generated sub-questions improve performance slightly (e.g., LLaMA-3.1-8B: 34.83% to 37.28%), but they degrade accuracy on MINTQA-TI (12.83% to 9.28%). This contrast reflects the reliance on models’ knowledge bases: for known but unpopular facts in MINTQA-POP, decomposition organizes existing knowledge, while on MINTQA-TI, knowledge gaps might lead to flawed decomposition and additional errors.

**Providing Sub-Questions:** Gold sub-questions significantly boost performance on MINTQA-POP (e.g., LLaMA-3.1-8B sees a 33.41% increase) by clarifying reasoning paths and allowing models to focus on synthesis. On MINTQA-TI, improvements are modest, with the best accuracy (23.75%) still from LLaMA-3.1-8B. This differences can be expected. While decomposition can help models better utilize their existing knowledge, it cannot compensate for the fundamental lack of information when handling questions about new facts.

## 9 Conclusion

In this work, we introduce **MINTQA**, a multi-hop QA benchmark reasoning across popular/unpopular and old/new knowledge. MINTQA spans reasoning chains from one to four hops, enabling systematic assessment of LLMs’ complex reasoning abilities. We also propose a comprehensive evaluation framework to assess key aspects of multi-hop QA, including the effectiveness of different question decomposition and retrieval strategies, which allows detailed analysis of models’ reasoning capabilities. Evaluations on state-of-the-art LLMs reveal that even the best LLM with retrieval still struggle on our benchmark. We give detailed discussions on the future work based on this paper in App.I.



## Limitation

This work has several key limitations. **First**, our definition of long-tail and new facts relies solely on Wikidata distribution patterns, which may not very accurately reflect knowledge representation in LLMs’ diverse pre-training corpora. **Second**, our simplified approach to constructing the retrieval corpus by concatenating entity-related facts into sequential sentences—differs from the complexity of real-world documents and might potentially overestimate the performance of retrieval-augmented methods. **Third**, budget constraints limited our evaluation of powerful closed-source models like GPT-4, though preliminary results suggest our benchmark remains challenging even for these advanced systems. **Fourth**, using only GPT-4o for eval data generation may cause bias. However, we counter that this bias is negligible. We chose GPT-4o for its superior quality. The data-generation task is simple, and our focus is on diverse knowledge and retrieval. GPT-4o mainly adds auxiliary words, not affecting factual QA. Also, single-LLM dataset construction is common in related work (Maekawa et al., 2024; Zhong et al., 2023). Regarding methodology, while our prompting strategy proved effective on the sampled data, we did not explore advanced techniques such as iterative prompt optimization. However, we hypothesize that such optimizations would yield limited improvements, as the core challenge lies in models’ knowledge gaps.

## References

- Marah Abidin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, and Martin Cai etc. 2024. *Phi-3 technical report: A highly capable language model locally on your phone*. *Preprint*, arXiv:2404.14219.
- Yufei Feng, Mo Yu, Wenhan Xiong, Xiaoxiao Guo, Junjie Huang, Shiyu Chang, Murray Campbell, Michael Greenspan, and Xiaodan Zhu. 2020. *Learning to recover reasoning chains for multi-hop question answering via cooperative games*. *Preprint*, arXiv:2004.02393.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, and . Christian Keller etc. 2024. *The llama 3 herd of models*. *Preprint*, arXiv:2407.21783.
- Xanh Ho, A. Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. *Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps*. *ArXiv*, abs/2011.01060.
- Jie Huang and Kevin Chen-Chuan Chang. 2023. *Towards reasoning in large language models: A survey*. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. 2024. *Qwen2.5-coder technical report*. *Preprint*, arXiv:2409.12186.
- Shayekh Bin Islam, Md Asib Rahman, K S M Tozammel Hossain, Enamul Hoque, Shafiq R. Joty, and Md. Rizwan Parvez. 2024. *Open-rag: Enhanced retrieval-augmented reasoning with open-source large language models*. *ArXiv*, abs/2410.01782.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. *Unsupervised dense information retrieval with contrastive learning*. *Trans. Mach. Learn. Res.*, 2022.
- Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C. Park. 2024. *Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity*. In *North American Chapter of the Association for Computational Linguistics*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. *Mistral 7b*. *Preprint*, arXiv:2310.06825.
- Hailey Joren, Jianyi Zhang, Chun-Sung Ferng, Da-Cheng Juan, Ankur Taly, and Cyrus Rashtchian. 2024. *Sufficient context: A new lens on retrieval augmented generation systems*. *CoRR*, abs/2411.06037.
- Ehsan Kamalloo, Nouha Dziri, Charles L. A. Clarke, and Davood Rafiei. 2023. *Evaluating open-domain question answering in the era of large language models*. *ArXiv*, abs/2305.06984.

696	Daniel Khashabi, Erfan Sadeqi Azer, Tushar Khot,	Sewon Min, Eric Wallace, Sameer Singh, Matt Gardner,	753
697	Ashish Sabharwal, and Dan Roth. 2019. <a href="#">On the</a>	Hannaneh Hajishirzi, and Luke Zettlemoyer. 2019.	754
698	<a href="#">possibilities and limitations of multi-hop reasoning</a>	<a href="#">Compositional questions do not necessitate multi-hop</a>	755
699	<a href="#">under linguistic imperfections.</a> <i>arXiv: Computation</i>	<a href="#">reasoning.</a> In <i>Annual Meeting of the Association for</i>	756
700	<a href="#">and Language.</a>	<i>Computational Linguistics.</i>	757
701	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gus-	758
702	field, Michael Collins, Ankur P. Parikh, Chris Alberti,	tavo Hernández Abrego, Ji Ma, Vincent Zhao,	759
703	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei	760
704	ton Lee, Kristina Toutanova, Llion Jones, Matthew	Yang. 2021. <a href="#">Large dual encoders are generalizable</a>	761
705	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	<a href="#">retrievers.</a> <i>ArXiv</i> , abs/2112.07899.	762
706	Uszkoreit, Quoc V. Le, and Slav Petrov. 2019. <a href="#">Natu-</a>	Shiyu Ni, Keping Bi, J. Guo, and Xueqi Cheng. 2024.	763
707	<a href="#">ral questions: A benchmark for question answering</a>	<a href="#">When do llms need retrieval augmentation? miti-</a>	764
708	<a href="#">research.</a> <i>Transactions of the Association for Compu-</i>	<a href="#">gating llms' overconfidence helps retrieval augmen-</a>	765
709	<i>tational Linguistics</i> , 7:453–466.	<a href="#">tation.</a> In <i>Annual Meeting of the Association for</i>	766
710	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	<i>Computational Linguistics.</i>	767
711	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gon-	Jeff Pan, Simon Razniewski, Jan-Christoph Kalo, Sneha	768
712	zalez, Haotong Zhang, and Ion Stoica. 2023. <a href="#">Effi-</a>	Singhania, Jiaoyan Chen, Stefan Dietze, Hajira	769
713	<a href="#">cient memory management for large language model</a>	Jabeen, Janna Omeliyanenko, Wen Zhang, Matteo	770
714	<a href="#">serving with pagedattention.</a> <i>Proceedings of the 29th</i>	Lissandrini, et al. 2023. Large language models	771
715	<i>Symposium on Operating Systems Principles.</i>	and knowledge graphs: Opportunities and challenges.	772
716	Patrick Lewis, Ethan Perez, Aleksandara Piktus, Fabio	<i>Transactions on Graph Data and Knowledge.</i>	773
717	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin	774
718	rich Kuttler, Mike Lewis, Wen tau Yih, Tim Rock-	Zhao, J. Liu, Hao Tian, Huaqin Wu, Ji rong Wen,	775
719	täschel, Sebastian Riedel, and Douwe Kiela. 2020.	and Haifeng Wang. 2023. <a href="#">Investigating the factual</a>	776
720	<a href="#">Retrieval-augmented generation for knowledge-</a>	<a href="#">knowledge boundary of large language models with</a>	777
721	<a href="#">intensive nlp tasks.</a> <i>ArXiv</i> , abs/2005.11401.	<a href="#">retrieval augmentation.</a> <i>ArXiv</i> , abs/2307.11019.	778
722	Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin	Stephen E. Robertson and Hugo Zaragoza. 2009. <a href="#">The</a>	779
723	Wang, Michal Lukasik, Andreas Veit, Felix X. Yu,	<a href="#">probabilistic relevance framework: Bm25 and be-</a>	780
724	and Surinder Kumar. 2022. <a href="#">Large language mod-</a>	<a href="#">yond.</a> <i>Found. Trends Inf. Retr.</i> , 3:333–389.	781
725	<a href="#">els with controllable working memory.</a> <i>ArXiv</i> ,	Priyanka Sen, Alham Fikri Aji, and Amir Saffari.	782
726	abs/2211.05110.	2022. <a href="#">Mintaka: A complex, natural, and multilingual</a>	783
727	Zekai Li and Wei Peng. 2023. <a href="#">Self-adaptive reason-</a>	<a href="#">dataset for end-to-end question answering.</a> <i>ArXiv</i> ,	784
728	<a href="#">ing on sub-questions for multi-hop question answer-</a>	abs/2210.01613.	785
729	<a href="#">ing.</a> <i>ICASSP 2023 - 2023 IEEE International Con-</i>	Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen	786
730	<i>ference on Acoustics, Speech and Signal Processing</i>	Zhong, Kaixiong Zhou, and Ninghao Liu. 2024a.	787
731	<i>(ICASSP)</i> , pages 1–5.	<a href="#">Retrieval-enhanced knowledge editing in language</a>	788
732	Yanming Liu, Xinyue Peng, Xuhong Zhang, Weihao	<a href="#">models for multi-hop question answering.</a> In <i>Inter-</i>	789
733	Liu, Jianwei Yin, Jiannan Cao, and Tianyu Du. 2024.	<i>national Conference on Information and Knowledge</i>	790
734	<a href="#">Ra-isf: Learning to answer and understand from re-</a>	<i>Management.</i>	791
735	<a href="#">trieval augmentation via iterative self-feedback.</a> In	Zhengliang Shi, Shuo Zhang, Weiwei Sun, Shen Gao,	792
736	<i>Annual Meeting of the Association for Computational</i>	Pengjie Ren, Zhumin Chen, and Zhaochun Ren.	793
737	<i>Linguistics.</i>	2024b. <a href="#">Generate-then-retrieve in retrieval-augmented</a>	794
738	Seiji Maekawa, Hayate Iso, Sairam Gurajada, and Nikita	<a href="#">generation for multi-hop question answering.</a> In <i>An-</i>	795
739	Bhutani. 2024. <a href="#">Retrieval helps or hurts? a deeper</a>	<i>annual Meeting of the Association for Computational</i>	796
740	<a href="#">dive into the efficacy of retrievals augmentation to</a>	<i>Linguistics.</i>	797
741	<a href="#">language models.</a> In <i>Proceedings of the 2024 Con-</i>	Heydar Soudani, Evangelos Kanoulas, and Faegheh Ha-	798
742	<i>ference of the North American Chapter of the Asso-</i>	sibi. 2024. <a href="#">Fine tuning vs. retrieval augmented gener-</a>	799
743	<i>ciation for Computational Linguistics: Human Lan-</i>	<a href="#">ation for less popular knowledge.</a> In <i>Proceedings of</i>	800
744	<i>guage Technologies (Volume 1: Long Papers)</i> , pages	<i>the 2024 Annual International ACM SIGIR Confer-</i>	801
745	5506–5521, Mexico City, Mexico. Association for	<i>ence on Research and Development in Information</i>	802
746	Computational Linguistics.	<i>Retrieval in the Asia Pacific Region, SIGIR-AP 2024,</i>	803
747	Alex Troy Mallen, Akari Asai, Victor Zhong, Rajarshi	page 12–22, New York, NY, USA. Association for	804
748	Das, Hannaneh Hajishirzi, and Daniel Khashabi.	Computing Machinery.	805
749	2022. <a href="#">When not to trust language models: Investigat-</a>	Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang,	806
750	<a href="#">ing effectiveness of parametric and non-parametric</a>	Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A.	807
751	<a href="#">memories.</a> In <i>Annual Meeting of the Association for</i>	Smith, Luke Zettlemoyer, and Tao Yu. 2022. <a href="#">One</a>	808
752	<i>Computational Linguistics.</i>		

809	embedder, any task: Instruction-finetuned text embeddings.	865	on Empirical Methods in Natural Language Processing, pages 17716–17736.	866
810				
811	Kai Sun, Y. Xu, Hanwen Zha, Yue Liu, and Xinhsuai Dong. 2023. <a href="#">Head-to-tail: How knowledgeable are large language models (llms)? a.k.a. will llms replace knowledge graphs?</a> <i>ArXiv</i> , abs/2308.10168.	867	Orion Weller, Benjamin Van Durme, Dawn Lawrie, Ashwin Paranjape, Yuhao Zhang, and Jack Hes-	868
812		869	sel. 2024. <a href="#">Promptriever: Instruction-trained retrievers can be prompted like language models.</a> <i>ArXiv</i> , abs/2409.11136.	870
813		871		
814				
815	Yixuan Tang and Yi Yang. 2024. <a href="#">Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries.</a> <i>ArXiv</i> , abs/2401.15391.	872	Shitao Xiao, Zheng Liu, Peitian Zhang, and Xing-	873
816		874	grun Xing. 2023. <a href="#">Lm-cocktail: Resilient tuning of language models via model merging.</a> <i>ArXiv</i> , abs/2311.13534.	875
817				
818	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupati-	876	Wenhan Xiong, Xiang Lorraine Li, Srini Iyer, Jingfei	877
819	raju, Léonard Hussenot, Thomas Mesnard, Bobak	878	Du, Patrick Lewis, William Yang Wang, Yashar	879
820	Shahriari, Alexandre Ramé, Johan Ferret, Peter	880	Mehdad, Wen tau Yih, Sebastian Riedel, Douwe	881
821	Liu, Pouya Tafti, Abe Friesen, Michelle Casbon,		Kiela, and Barlas Oğuz. 2020. <a href="#">Answering complex open-domain questions with multi-hop dense</a>	
822	Sabela Ramos, Ravin Kumar, Charline Le Lan, and		retrieval. <i>ArXiv</i> , abs/2009.12756.	
823	Sammy Jerome etc. 2024. <a href="#">Gemma 2: Improving open language models at a practical size.</a> <i>Preprint</i> , arXiv:2408.00118.			
824				
825				
826				
827	H. Trivedi, Niranjan Balasubramanian, Tushar Khot,	882	Shicheng Xu. 2023. <a href="#">Search-in-the-chain: Towards accurate, credible and traceable large language models</a>	883
828	and Ashish Sabharwal. 2021. <a href="#">Musique: Multi-hop questions via single-hop question composition.</a> <i>Transactions of the Association for Computational Linguistics</i> , 10:539–554.	884	for knowledge-intensive tasks.	
829				
830				
831				
832	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry	885	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-	886
833	Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny	887	gio, William W. Cohen, Ruslan Salakhutdinov, and	888
834	Zhou, Quoc Le, and Thang Luong. 2023. <a href="#">Freshllms: Refreshing large language models with search engine</a>	889	Christopher D. Manning. 2018. <a href="#">Hotpotqa: A dataset for diverse, explainable multi-hop question answer-</a>	890
835	<a href="#">augmentation.</a> In <i>Annual Meeting of the Association for Computational Linguistics</i> .		ing. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	
836				
837				
838	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry	891	Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Be-	892
839	Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny	893	rant. 2024. <a href="#">Making retrieval-augmented language models robust to irrelevant context.</a> In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	894
840	Zhou, Quoc Le, and Thang Luong. 2024. <a href="#">Fresh-LLMs: Refreshing large language models with search</a>	895		896
841	<a href="#">engine augmentation.</a> In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 13697–13720, Bangkok, Thailand. Association for Computational Linguistics.			
842				
843				
844				
845				
846	Jinyuan Wang, Junlong Li, and Hai Zhao. 2023. <a href="#">Self-prompted chain-of-thought on large language models for open-domain multi-hop reasoning.</a> <i>ArXiv</i> , abs/2310.13552.	897	Donghan Yu, Sheng Zhang, Patrick Ng, Henghui	898
847		899	Zhu, Alexander Hanbo Li, Jun Wang, Yiqun Hu,	900
848		901	William Yang Wang, Zhiguo Wang, and Bing Xiang.	902
849		903	2023. <a href="#">Decaf: Joint decoding of answers and logical forms for question answering over knowledge bases.</a> In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	904
		905	W. Yu, Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhit-	906
850	Shouhui Wang and Biao Qin. 2024. <a href="#">No need for large-scale search: Exploring large language models in complex knowledge base question answering.</a> In <i>International Conference on Language Resources and Evaluation</i> .	907	ing Hu, Qingyun Wang, Heng Ji, and Meng Jiang.	908
851			2020. <a href="#">A survey of knowledge-enhanced text generation.</a> <i>ACM Computing Surveys</i> , 54:1 – 38.	
852				
853				
854				
855	Siyuan Wang, Zhongyu Wei, Zhihao Fan, Qi Zhang,	909	Zihan Zhang, Meng Fang, and Ling Chen. 2024. <a href="#">Retrievalqa: Assessing adaptive retrieval-augmented</a>	910
856	and Xuanjing Huang. 2022. <a href="#">Locate then ask: Interpretable stepwise reasoning for multi-hop question</a>	911	generation for short-form open-domain question	912
857	<a href="#">answering.</a> In <i>International Conference on Computational Linguistics</i> .	912	answering. <i>ArXiv</i> , abs/2402.16457.	
858				
859				
860	Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran	913	Ruochen Zhao, Xingxuan Li, Shafiq R. Joty, Chengwei	914
861	Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi,	915	Qin, and Lidong Bing. 2023. <a href="#">Verify-and-edit: A knowledge-enhanced chain-of-thought framework.</a> <i>ArXiv</i> , abs/2305.03268.	916
862	Zhengyuan Wang, Shizheng Li, Qi Qian, et al. 2024. Searching for best practices in retrieval-augmented	917	Zexuan Zhong, Zhengxuan Wu, Christopher Manning,	918
863	generation. In <i>Proceedings of the 2024 Conference</i>	918	Christopher Potts, and Danqi Chen. 2023. <a href="#">MQuAKE: Assessing knowledge editing in language models via</a>	919
864		919		

multi-hop questions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15686–15702, Singapore. Association for Computational Linguistics.

Andrew Zhu, Alyssa Hwang, Liam Dugan, and Chris Callison-Burch. 2024a. [Fanoutqa: A multi-hop, multi-document question answering benchmark for large language models](#). In *Annual Meeting of the Association for Computational Linguistics*.

Yun Zhu, Jia-Chen Gu, Caitlin Sikora, Ho Ko, Yinxiao Liu, Chu-Cheng Lin, Lei Shu, Liangchen Luo, Lei Meng, Bang Liu, and Jindong Chen. 2024b. [Accelerating inference of retrieval-augmented generation via sparse context selection](#). *ArXiv*, abs/2405.16178.

Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou, and Ji rong Wen. 2023. [Large language models for information retrieval: A survey](#). *ArXiv*, abs/2308.07107.

Ziyuan Zhuang, Zhiyang Zhang, Sitao Cheng, Fangkai Yang, Jia Liu, Shujian Huang, Qingwei Lin, S. Rajmohan, Dongmei Zhang, and Qi Zhang. 2024. [Efficientrag: Efficient retriever for multi-hop question answering](#). *ArXiv*, abs/2408.04259.





Dataset	# Questions	Multi-hop	Time	Popularity	Sub-questions	Original of Generated Questions
PopQA (Mallen et al., 2022)	14,268	×	×	✓	×	Templates
WiTQA (Mackawa et al., 2024)	14,837	×	×	✓	×	Machine
Head-to-tail (Sun et al., 2023)	18,171	×	×	✓	×	Template
RetrievalQA (Zhang et al., 2024)	1,271	×	✓	✓	×	Mixed
FreshQA (Vu et al., 2024)	600	✓	✓	×	×	Human
MultihopQA-RAG (Tang and Yang, 2024)	2,556	✓	×	×	×	Machine
HotpotQA (Yang et al., 2018)	112,779	✓	×	×	×	Human
MuSiQue (Trivedi et al., 2021)	24,814	✓	×	×	×	Human
2WikiMultiHopQA (Ho et al., 2020)	192,606	✓	×	×	×	Templates
Mintaka (Sen et al., 2022)	20,000	✓	×	×	×	Human
FanoutQA (Zhu et al., 2024a)	1,034	✓	×	×	✓	Human
MINTQA (Ours)	28,366	✓	✓	✓	✓	Machine

Table 4: Comparison between MINTQA and other datasets.

## C Details of Benchmark

### C.1 Details of Benchmark Curation

In Section 3, we present a comprehensive description of our benchmark construction methodology. Our approach includes carefully designed prompts for both question generation and validation processes. The complete specifications of these prompts are illustrated in figures 9 through 15.

### C.2 License

Our benchmark data are released under the MIT License, which is detailed in <https://open-source.org/licenses/MIT>.

## D Comparison with Existing Benchmarks

In this section, we provide a comprehensive comparison with question answering benchmarks closely related to our own in Table 4.

Compared to previous benchmarks, ours encompasses both old/new knowledge and unpopular/popular knowledge, presenting new challenges for retrieval-augmented large language model systems. Furthermore, unlike RetrievalQA (Zhang et al., 2024), which covers old/new or unpopular/popular knowledge but relies on integrating existing QA datasets, our benchmark generates questions using language models, enabling scalable data construction. RetrievalQA, on the other hand, is constrained by the limited availability of existing datasets and focuses solely on short-form open-domain question answering.

Additionally, while multi-hop datasets exist, only FreshQA (Vu et al., 2023) involves new knowledge in questions. However, FreshQA’s data is manually created, limited to just 600 samples, and lacks scalability. Our dataset, by contrast, provides sub-questions that assist in evaluating or training models on intermediate reasoning steps in multi-hop processes, enabling a more comprehensive assessment of LLMs’ capabilities on similar tasks.

This more integrated benchmark can help the research community gain deeper insights into the weaknesses of large models in question answering, improve training methods, and address the limitations of current benchmark practices.

## E Qualitative Analysis

Table 22 to Table 30 present representative examples of multi-hop questions and their corresponding sub-questions generated by GPT-4o for both MINTQA-POP and MINTQA-TI datasets. We have selected three representative instances for each hop level, ranging from single-hop to four-hop questions. As demonstrated in the table, GPT-4o effectively converted the triplets into well-structured, coherent questions. The high quality of these generated questions makes them suitable for evaluating retrieval-augmented LLMs’ capabilities in handling multi-hop questions that involve rare and new knowledge.

## F Additional Experimental Details

### F.1 Implementation Details

In our experiments, we utilized the following state-of-the-art LLMs, with detailed version specifications: GPT-3.5 (gpt-3.5-turbo-1106), GPT4o-mini (gpt-4o-mini-2024-07-18), GPT4o (gpt-4o-2024-08-06), LLaMA-3.1-8B (LLaMA-3.1-8B-instruct), LLaMA-3.1-70B (LLaMA-3.1-70B-instruct), LLaMA-3.2-1B (LLaMA-3.2-1B-instruct), LLaMA-3.2-3B (LLaMA-3.2-3B-Instruct), Qwen-2.5-1.5B (Qwen-2.5-1.5B-Instruct), Qwen-2.5-3B (Qwen-2.5-3B-Instruct), Qwen-2.5-7B (Qwen-2.5-7B-Instruct), Qwen-2.5-14B (Qwen-2.5-14B-Instruct), Qwen-2.5-32B (Qwen-2.5-32B-Instruct), Qwen-2.5-72B (Qwen-2.5-72B-Instruct), Gemma-2-2B (Gemma-2-2B-it), Gemma-2-9B (Gemma-2-9B-it), Gemma-2-27B (Gemma-2-27B-it), Phi-3-mini (Phi-3-mini-4k), Phi-3-small (Phi-3-small-8k), Phi-3-medium (Phi-3-medium-

Model	POP						TI						
	1	0.67	0.5	0.33	0.25	0	1	0.75	0.67	0.5	0.33	0.25	0
GPT													
GPT-3.5	83.94	2.74	69.59	18.46	9.40	49.34	9.68	4.34	8.79	8.88	2.02	2.30	17.97
GPT-4o	89.11	3.23	87.83	29.95	14.65	63.42	14.73	7.85	15.80	14.03	3.85	6.91	22.89
GPT-4o-mini	84.32	2.88	78.55	26.24	12.30	47.96	12.31	6.20	12.09	11.32	3.49	4.93	18.90
Llama													
Llama-3.2-1B	64.51	0.49	29.30	4.24	5.45	23.59	6.35	2.69	5.36	5.33	2.57	2.96	14.29
Llama-3.2-3B	75.87	1.41	58.91	14.22	7.05	29.17	7.40	5.37	6.18	6.84	3.12	4.44	15.91
Llama-3.1-8B	75.87	2.18	54.91	19.88	7.70	38.29	11.29	5.37	8.65	9.59	5.50	5.43	21.73
Llama-3.1-70B	90.42	2.32	68.99	22.84	12.15	55.23	13.47	7.44	11.95	15.80	4.22	7.57	23.06
Qwen													
Qwen-1.5B	62.44	3.79	49.13	14.84	7.00	22.68	8.49	5.17	8.10	7.32	2.39	4.44	16.41
Qwen-3B	62.82	3.65	51.57	12.06	6.40	23.40	7.89	4.96	5.22	5.68	3.30	4.28	13.62
Qwen-7B	73.62	5.34	61.49	12.81	6.35	26.87	9.40	4.55	7.83	7.81	3.12	4.93	16.15
Qwen-14B	78.59	5.20	73.66	37.68	11.75	35.18	13.29	6.82	12.36	10.39	4.77	5.43	18.97
Qwen-32B	81.13	6.04	76.87	30.17	12.25	35.69	12.56	7.23	13.46	11.45	3.49	7.89	19.47
Qwen-72B	80.56	5.41	78.74	31.85	10.60	40.55	12.63	8.47	14.29	11.90	3.67	5.26	20.63
Gemma													
Gemma-2-2B	57.18	1.05	36.83	10.73	4.40	24.39	7.30	4.13	5.08	5.73	1.47	2.80	16.18
Gemma-2-9B	80.28	1.62	65.33	26.63	11.75	33.99	10.73	4.55	6.04	7.19	2.39	4.28	18.14
Gemma-2-27B	80.56	2.74	70.07	27.83	9.85	37.74	12.17	6.40	9.20	9.50	3.85	7.40	18.87
Phi													
Phi-3-mini	79.53	2.67	64.41	23.06	10.70	33.17	9.61	5.58	11.13	9.19	2.39	7.57	15.45
Phi-3-small	74.74	2.18	60.06	21.33	9.60	34.34	9.93	5.17	7.42	9.41	2.02	5.26	15.85
Phi-3-medium	84.79	2.60	70.35	25.57	11.40	45.97	10.87	6.82	12.50	11.10	2.20	4.28	17.34
Mistral													
Mistral-7B-v0.3	81.31	1.48	48.78	12.81	8.00	36.20	9.12	3.93	7.14	7.55	2.39	4.11	18.21
Ministral-8B-2410	76.06	4.50	56.34	28.75	11.75	35.57	10.84	5.58	11.26	9.28	5.32	8.72	23.29
Mixtral-8x7B-v0.1	84.69	2.74	66.03	24.16	10.95	47.30	13.43	5.58	10.85	10.92	4.22	5.10	20.40

Table 5: The model’s accuracy in the zero-shot setting is analyzed within MINTQA-POP and MINTQA-TI, categorized based on the proportion of popular facts and old facts. A value of 0 indicates that the questions are entirely composed of unpopular facts or new facts, with other numbers increasing proportionally.

4k), Mistral-7B (mistral-7B-instruct-v0.3), Mixtral-8x7B (Mixtral-8x7B-instruct-v0.1), and Ministral-8B (Ministral-8B-instruct-2410). All experiments were conducted using 4 A100 (80GB) GPUs. From Table 14 to 21, we provide the prompts used to instruct these models in completing their respective tasks.

## F.2 Retrievers and KG Linearization Details

We evaluate seven retrieval approaches across three categories: 1) Sparse retriever: **BM25** (Robertson and Zaragoza, 2009). 2) Vector retrievers pre-trained on large unlabeled corpora: **Contriever** (Izacard et al., 2021): Fine-tuned on MS-MARCO, **GTR-LARGE/XL** (Ni et al., 2021) and **BGE** (Xiao et al., 2023): Further fine-tuned on NQ (Kwiatkowski et al., 2019) and HotpotQA (Yang et al., 2018). 3) Instruction-tuned text embedding retrievers: **Instructor-XL** (Su et al., 2022): Multi-task trained on 330 tasks for instruction robustness. **Promptriever** (Weller et al., 2024): Uses LLaMA backbone, trained on curated instance-level instruction sets from MS-MARCO, demonstrating superior retrieval performance compared to Instructor-XL.

We linearise the knowledge graph (KG)  $\mathcal{G}$  as a source of text retrieval in the corpus, with reference to the work in Yu et al. (2023). Specifically, for each entity in  $\mathcal{G}$ , we extract a 1-hop subgraph centered on the entity and convert it into linearized text, treating it as a passage. Since  $\mathcal{G}$  includes

both old and new versions of the Wikidata dump, knowledge conflicts may arise due to updates. Conflicting triples are separated into different passages. Each passage is split into chunks of 512 tokens, a size shown to be effective for practical applications (Wang et al., 2024).

## G Additional Experiments and Result Analysis

### G.1 Zero-shot: Performance Across Retrieval Categories

In Table 5, we present the performance of LLMs in a zero-shot evaluation setting across different proportions of unpopular/popular and old/new facts. As observed, the accuracy is highest when questions are composed solely of popular or old facts. For example, LLaMA-3.1-70B achieves an accuracy of 90.42% on MINTQA-POP and 13.47% on MINTQA-TI.

However, as the proportion of unpopular or new facts increases, the accuracy of the models shows a declining trend. Interestingly, when this proportion reaches 1, the accuracy tends to rise compared to lower ratios. This is likely because the proportion of 1 often includes many 1-hop questions, which are comparatively easier for the models to resolve.

### G.2 Sub-question Generation Analysis

From Figures 15, we illustrate the relationship between the number of sub-questions generated by

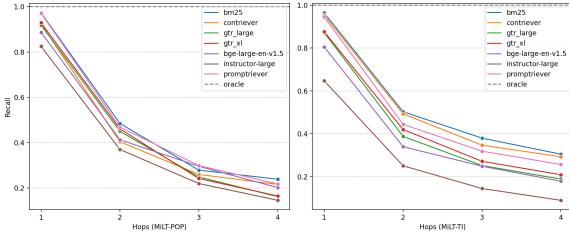


Figure 11: Recall performance of retrieval methods across two datasets for varying question hops.

models and the corresponding gold sub-question counts. This analysis considers scenarios where models are required to independently generate and answer sub-questions.

We observe substantial differences among models of similar sizes. For instance, the Qwen2.5-7B model tends to generate fewer sub-questions, with most counts falling in the range of 1 or 2. In contrast, the Mistral-7B model produces sub-questions with a more uniform distribution, primarily ranging from 2 to 5. Despite these differences, smaller models, such as Qwen2.5-1.5B and LLaMA-3.2-1B, exhibit similar trends. Both predominantly generate only 1 sub-question, reflecting the limited capability of these smaller LLMs to generate sub-questions as part of their answering process. Examining the distributions of larger models on the MINTQA-POP and MINTQA-TI datasets reveals that, despite differences in the datasets, large models exhibit similar distributions in terms of actual step counts and the number of sub-questions generated by the models.

### G.3 More Analysis of Direct Retrieval

Direct retrieval strategy have limitations when handling multi-hop questions. Figure 11 reveals significant limitations in current direct retrieval approaches when handling multi-hop questions. First, the retrieval effectiveness decreases markedly with increased hop count. We observe a consistent decline in recall rates across all retrieval methods as question complexity increases, indicating fundamental limitations in the direct retrieval approach. Second, among the retrieval methods evaluated, BM25 demonstrated the best performance. This can be explained by the highly structured nature of our KG-linearized corpus. While dense retrieval methods excel at capturing semantic similarities in natural text, BM25’s lexical matching approach is well-suited for knowledge graph-derived text.

We demonstrate the influence of knowledge new-

Model	MINTQA-POP		MINTQA-TI	
	Acc	F1	Acc	F1
GPT-3.5	54.82	41.21	52.04	26.67
GPT-4o-mini	65.61	48.67	59.56	28.72
GPT-4o	68.84	51.35	65.46	38.18
LLaMA-3.2-1B	67.05	26.76	62.32	25.59
LLaMA-3.2-3B	56.40	32.91	47.37	33.65
LLaMA-3.1-8B	59.93	25.83	55.55	26.80
LLaMA-3.1-70B	70.03	44.88	63.26	29.54
Qwen2.5-1.5B	64.01	26.47	58.42	24.92
Qwen2.5-3B	47.08	38.89	48.09	41.86
Qwen2.5-7B	70.90	48.23	62.22	34.47
Qwen2.5-14B	69.31	42.08	62.39	26.28
Qwen2.5-32B	28.74	29.80	18.69	11.04
Qwen2.5-72B	65.39	55.33	51.11	37.00
Gemma-2-2B	5.39	3.84	21.07	14.53
Gemma-2-9B	70.90	46.18	63.05	30.39
Gemma-2-27B	73.22	55.67	63.64	36.15
Phi-3-mini	64.47	26.44	56.60	24.72
Phi-3-small	75.79	55.56	62.53	40.50
Phi-3-medium	28.43	21.69	18.50	10.64
Mistral-7B-v0.3	37.58	30.94	26.28	18.88
Minstral-8B-2401	67.07	27.35	62.27	25.66
Mixtral-8x7B-v0.1	28.00	15.14	18.48	10.51

Table 6: The model’s accuracy and F1 score for the task of determining question retrieval, sub-question generation, or direct answering.

ness and popularity on direct retrieval scenarios, using Qwen2.5-72B paired with BM25 as a representative example. As shown in Figure 12 (a) and (c), QA performance declines with an increasing proportion of unpopular or new knowledge in questions. However, performance improves when the proportion of new knowledge reaches 100% (i.e., no old knowledge), as higher new knowledge presence boosts recall rates (Figure 12 (d)), ultimately enhancing QA accuracy on MINTQA-TI. This highlights the retriever’s effectiveness in handling new knowledge.

### G.4 Complete Results for Decomposition-Dynamic Retrieval

Table 11 presents the complete results on MINTQA-POP and MINTQA-TI using large models to output confidence scores for sub-questions and determine whether retrieval is needed based on the confidence values. We conducted experiments with three retrievers: BM25, Contriever, and PromptRetrieve.

## H Evaluating LLMs’ Decision-Making Capabilities in Multi-hop QA

When evaluating the ability of LLMs to answer questions using their parametric knowledge, we frequently observed “I don’t know” responses or no answers at all. This clearly indicates the difficulties LLMs encounter in solving complex multi-hop



Model	MINTQA-POP		MINTQA-TI	
	Acc	F1	Acc	F1
GPT-3.5	49.32	49.15	47.32	47.27
GPT-4o-mini	47.77	46.98	47.48	47.14
GPT-4o	37.13	33.19	44.67	43.83
LLaMA-3.2-1B	31.40	24.71	43.07	30.74
LLaMA-3.2-3B	30.82	23.56	43.05	30.10
LLaMA-3.1-8B	30.88	23.66	43.05	30.15
LLaMA-3.1-70B	54.60	54.53	48.46	47.70
Qwen2.5-1.5B	41.76	41.55	51.03	51.03
Qwen2.5-3B	33.42	28.22	43.62	33.43
Qwen2.5-7B	32.11	25.63	43.12	31.24
Qwen2.5-14B	65.33	63.42	53.25	42.62
Qwen2.5-32B	68.13	63.53	53.87	43.39
Qwen2.5-72B	32.25	26.34	43.57	34.20
Gemma-2-2B	30.82	23.56	43.06	30.11
Gemma-2-9B	54.37	54.37	50.99	50.65
Gemma-2-27B	69.39	63.95	55.58	41.72
Phi-3-mini	32.28	26.05	43.23	32.96
Phi-3-small	35.28	30.68	44.35	41.48
Phi-3-medium	40.77	38.12	44.19	42.75
Mistral-7B-v0.3	38.13	36.23	46.37	45.53
Ministral-8B-2410	30.86	23.63	43.06	30.13
Mixtral-8x7B-v0.1	68.29	43.32	56.96	36.97

Table 7: The accuracy and F1 scores of different models in determining whether sub-questions should be retrieved or directly answered.

questions relying solely on their internal knowledge or limited single-step reasoning capabilities. To overcome these challenges, LLMs often employ sub-question decomposition and retrieval strategies. Nevertheless, the effectiveness of these strategies highly depends on the model’s capacity to determine when to utilize them. We analyze this from three crucial aspects. Our objective here is to assess the degree to which the model’s decisions are consistent with the heuristic labels derived from Wikidata. We regard this as an evaluation of **preference**.

## H.1 Direct Answer vs. Decompositions vs. Retrieval

When encountering multi-hop questions, models must choose between direct answering, sub-question generation, or retrieval. This decision significantly impacts system efficiency and accuracy. Specifically, simple factual questions are often answered directly, while multi-hop or rare fact queries benefit from decomposition or retrieval. The evaluation is conducted on the main question only and involves a three - class classification: decomposition, retrieval, and direct answer. The labels are assigned as follows:

- All queries with  $hop\_num \geq 2$  are labeled as decomposition.
- Queries with  $hop\_num == 1$  constructed from unpopular or new knowledge are labeled

Model	MINTQA-POP		MINTQA-TI	
	Acc	F1	Acc	F1
GPT-3.5	37.74	30.78	43.74	38.20
GPT-4o-mini	59.31	58.96	56.25	56.02
GPT-4o	71.30	71.22	59.30	57.91
LLaMA-3.2-1B	23.14	20.09	33.86	23.47
LLaMA-3.2-3B	28.62	22.63	37.52	26.98
LLaMA-3.1-8B	34.81	25.82	40.83	28.99
LLaMA-3.1-70B	58.19	57.67	52.43	52.43
Qwen2.5-1.5B	34.81	25.82	40.83	28.99
Qwen2.5-3B	77.49	72.19	59.52	49.39
Qwen2.5-7B	62.16	62.13	53.25	52.85
Qwen2.5-14B	79.05	78.53	60.02	57.83
Qwen2.5-32B	95.94	95.52	62.53	58.74
Qwen2.5-72B	83.68	83.06	62.03	59.54
Gemma-2-2B	34.81	25.82	40.83	28.99
Gemma-2-9B	40.05	34.31	43.79	38.25
Gemma-2-27B	65.83	65.82	55.32	54.93
Phi-3-mini	34.87	25.92	41.11	30.07
Phi-3-small	47.19	44.44	46.71	44.78
Phi-3-medium	35.36	26.77	42.65	35.16
Mistral-7B-v0.3	34.91	25.99	40.87	29.31
Ministral-8B-2410	34.81	25.82	40.88	29.14
Mixtral-8x7B-v0.1	35.96	28.18	41.73	32.52

Table 8: The accuracy and F1 scores of the model in determining whether the main question has been answered based on the given sub-question-answer pair.

as retrieval.

- Queries with  $hop\_num == 1$  constructed from popular or old knowledge are labeled as direct answer.

The purpose of this evaluation is to analyze the model’s initial decision - making when presented with a question.

As shown in Table 6, Phi-3-small-8k performs best on MINTQA-POP (Accuracy: 75.79%, F1: 55.56%), while GPT-4o leads on MINTQA-TI (Accuracy: 65.46%, F1: 38.18%). However, model size doesn’t always predict performance; Qwen2.5-32B underperforms its 14B variant. Lower-performing models, like Gemma-2-2B, favor direct answering (92.59% on MINTQA-TI), likely due to their limited ability to assess question complexity.

## H.2 Direct Answer vs. Retrieval for Sub-questions

When handling sub-questions, models must decide between direct answering and retrieval based on the required knowledge. Popular facts might be answered directly, while tail knowledge or recent information often requires retrieval. This evaluation focuses on queries with  $hop\_num \geq 2$ . Given the main question and sub-questions, the model must decide whether to retrieve or directly answer each sub-question. The labels are assigned as follows:

- Sub-questions constructed from unpopular or new knowledge are labeled as retrieval.
- Sub-questions constructed from popular or old knowledge are labeled as direct answer.

Our experiments results in Table 7 reveal a general correlation between model size and decision quality, with some exceptions. LLaMA-3.1-70B outperforms other LLaMA variants, achieving 54.60% and 48.46% accuracy on MINTQA-POP and MINTQA-TI, respectively. However, GPT-4o underperforms GPT-3.5, likely due to overconfidence in its parametric knowledge, as it selects direct answering on 93.48% of MINTQA-POP and 69.20% of MINTQA-TI questions. Additionally, models perform better on MINTQA-TI, indicating new knowledge provides a clearer signal for retrieval compared to knowledge of varying popularity, where the decision boundary is less distinct.

### H.3 Decomposition vs. Synthesis

For multi-hop questions (hop count  $\geq 2$ ), we evaluate models' ability to decide whether to decompose further or synthesize the final answer from intermediate results. This evaluation is also conducted on queries with *hop\_num*  $\geq 2$ . Given the main question, sub-questions, and corresponding sub-answers, the model must decide whether to continue decomposition or synthesize an answer. The labels are assigned as follows:

- If the number of sub-questions and sub-answers is less than the main question's hop number, the model should continue decomposition.
- If the number matches the hop number, the model should synthesize an answer.

As shown in Table 8, performance generally correlates with model size. Qwen2.5-32B achieves 95% accuracy on MINTQA-POP but drops to 62.53% on MINTQA-TI, reflecting new knowledge poses challenges for synthesizing. Some models like Mistral-7B, show extreme biases, predicting the main answer always within sub-answers for 99.90% cases of MINTQA-POP.

### H.4 Accuracy and F1 Across Categories

Table 9 and 10 reports the accuracy and F1 scores for each category under the evaluation setup described in Section H.2 and H.3. From the table, we can observe that most models demonstrate high accuracy, often exceeding 90% or even reaching 100% in identifying sub-questions that can directly generate answers. However, the F1 scores are sig-

Model	Class A		Class B	
	Acc.	F1	Acc.	F1
<b>GPT</b>				
GPT-3.5	89.4	52.1	31.5	46.2
GPT-4o-mini	97.3	53.5	25.7	40.5
GPT-4o	99.7	49.4	9.3	16.9
<b>Llama</b>				
Llama-3.2-1B	99.3	47.1	1.2	2.3
Llama-3.2-3B	100.0	47.1	0.0	0.0
Llama-3.1-8B	100.0	47.1	0.1	0.2
Llama-3.1-70B	94.9	56.3	36.6	52.8
<b>Qwen</b>				
Qwen-1.5B	57.9	38.0	34.6	45.1
Qwen-3B	97.9	47.5	4.7	8.9
Qwen-7B	100.0	47.6	1.9	3.7
Qwen-14B	68.9	55.1	63.7	71.8
Qwen-32B	52.9	50.6	74.9	76.5
Qwen-72B	98.2	47.2	2.8	5.5
<b>Gemma</b>				
Gemma-2-2B	100.0	47.1	0.0	0.0
Gemma-2-9B	88.7	54.5	39.1	54.2
Gemma-2-27B	49.6	50.0	78.2	78.0
<b>Phi</b>				
Phi-3-mini	99.5	47.5	2.4	4.6
Phi-3-small	99.0	48.5	6.9	12.8
Phi-3-medium	99.7	50.9	14.5	25.3
<b>Mistral</b>				
Mistral-7B-v0.3	89.8	47.2	15.1	25.2
Ministral-8B-2410	100.0	47.1	0.1	0.1
Mixtral-8x7B-v0.1	3.1	5.7	97.3	80.9

Table 9: The per-label accuracy and F1 scores for the tasks of sub-question judgment, retrieval, or direct answer generation.

nificantly lower. This discrepancy indicates that models tend to predict that all examples are solvable, revealing an overconfidence in their ability to answer our constructed benchmarks.

The table also highlights similar phenomena across models, particularly for LLaMA-3.1-8B, LLaMA-3.2-1B, LLaMA-3.2-3B, Qwen2.5-1.5B, Gemma-2-2B, and Ministral-8B-2410. These models consistently predict that the main question can be derived from existing sub-question answers. On the other hand, models in the same series, such as Qwen2.5 variants, exhibit more balanced accuracy and F1 scores across categories. This reflects significant inconsistencies among large models in determining whether sub-question answers suffice to answer the main question.

Such findings indicate the challenges of relying on large models for complex reasoning tasks and highlight the need for more robust evaluation metrics and methodologies.

## I Discussion on Future Directions

Based on the findings and limitations highlighted in this study, we propose several promising directions

Model	Class A		Class B	
	Acc.	F1	Acc.	F1
<b>GPT</b>				
GPT-3.5	99.7	52.7	4.6	8.8
GPT-4o-mini	98.5	62.8	38.4	55.1
GPT-4o	95.2	69.8	58.5	72.7
<b>Llama</b>				
Llama-3.2-1B	100.0	51.6	0.0	0.0
Llama-3.2-3B	100.0	51.6	0.0	0.0
Llama-3.1-8B	100.0	51.6	0.0	0.0
Llama-3.1-70B	99.6	62.4	36.1	53.0
<b>Qwen</b>				
Qwen-1.5B	100.0	51.6	0.0	0.0
Qwen-3B	48.6	60.1	92.9	84.3
Qwen-7B	93.4	63.2	45.5	61.0
Qwen-14B	91.1	75.2	72.6	81.9
Qwen-32B	93.7	94.1	97.1	96.9
Qwen-72B	92.7	79.8	78.9	86.3
<b>Gemma</b>				
Gemma-2-2B	100.0	51.6	0.0	0.0
Gemma-2-9B	100.0	53.7	8.1	14.9
Gemma-2-27B	97.1	66.4	49.1	65.2
<b>Phi</b>				
Phi-3-mini	100.0	51.7	0.1	0.2
Phi-3-small	99.7	56.8	19.1	32.1
Phi-3-medium	100.0	51.8	0.9	1.7
<b>Mistral</b>				
Mistral-7B-v0.3	100.0	51.7	0.2	0.3
Mixtral-8x7B-v0.1	98.9	51.8	2.3	4.6
Ministral-8B-2410	100.0	51.6	0.0	0.0

Table 10: The per-label accuracy and F1 scores for the task where the model is required to determine whether the answer to the main question has been found, given the sub-questions and their answers.

for future research:

▷ *Enhanced Sub-Question Generation and Planning* Smaller models (<14B) struggle with sub-question generation, suggesting a need for specialized training or architectural improvements. Future work could explore fine-tuning on decomposed reasoning chains or integrating reinforcement learning to optimize decomposition strategies. Hybrid approaches combining rule-based methods with LLM-driven planning may also help.

▷ *Adaptive Retrieval-Augmented Strategies* Models often over-rely on retrieval, leading to inefficiencies. Future research could focus on confidence calibration, uncertainty-aware retrieval triggers, or lightweight classifiers to reduce unnecessary retrievals while maintaining accuracy. Integrating retrieval necessity prediction into training could further optimize performance.

▷ *Knowledge-Type-Aware Reasoning* Models struggle to distinguish between popular/unpopular and old/new knowledge, leading to suboptimal strategy selection. Future frameworks could incorporate explicit knowledge-type classifiers or metadata to

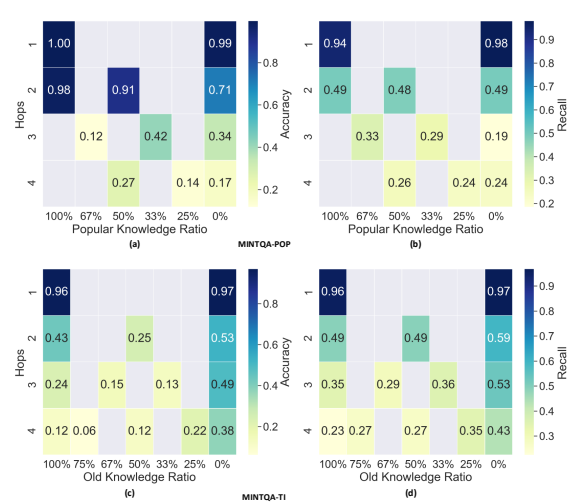


Figure 12: Heatmaps (a) and (c) show Qwen2.5-72B with BM25 performance on two datasets, while heatmaps (b) and (d) shows BM25 recall. The X-axis represents the proportion of popular knowledge required in the question, and the Y-axis indicates question hops.

guide retrieval and parametric knowledge usage, improving decision-making during inference.

▷ *Cross-Hop Information Synthesis* Even with perfect retrieval and decomposition, models fail to fully synthesize information across sub-questions. This calls for the exploration of novel techniques and architectures that can effectively integrate knowledge from multiple sub-questions, leading to more comprehensive and accurate answers.

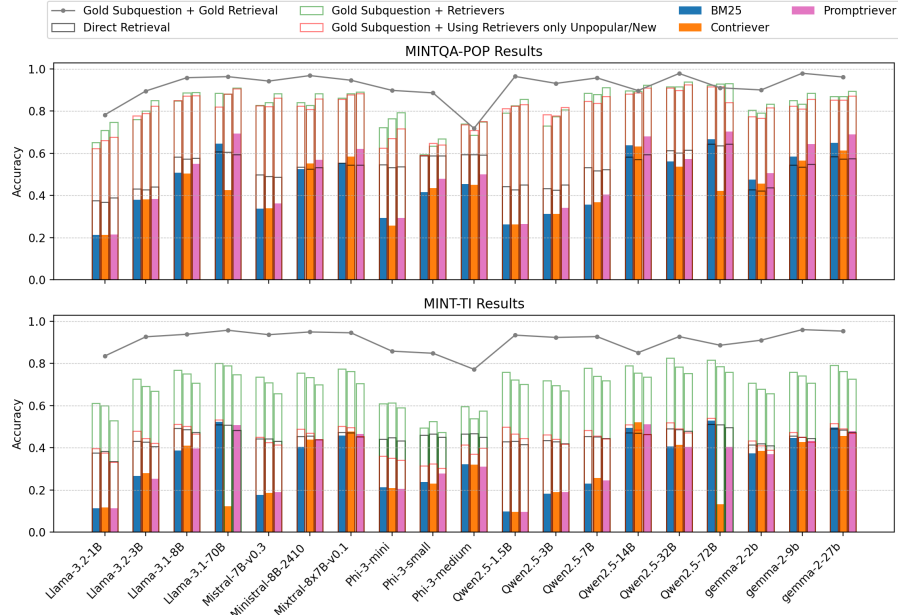


Figure 13: Full datasets evaluation results of performance of all models with three retrievers (i.e. BM25, Contriever and Promptriever) using decomposition-retrieval approach on two datasets. Gold Subquestion + Gold Retrieval means that the model uses gold subquestion and gold retrieval results. Direct Retrieval means that the model uses different retrievers for direct retrieval instead of the decomposition-retrieval approach. Gold Subquestion + Retrievers indicates that the model uses gold subquestion and uses different retrievers to retrieve. Gold Subquestion + Using Retrievers only Unpopular/New indicates that the model uses gold subquestion and uses different retrievers to only retrieve the question involving unpopular knowledge or new knowledge, while models rely on their own to direct answer the popular knowledge or old knowledge.

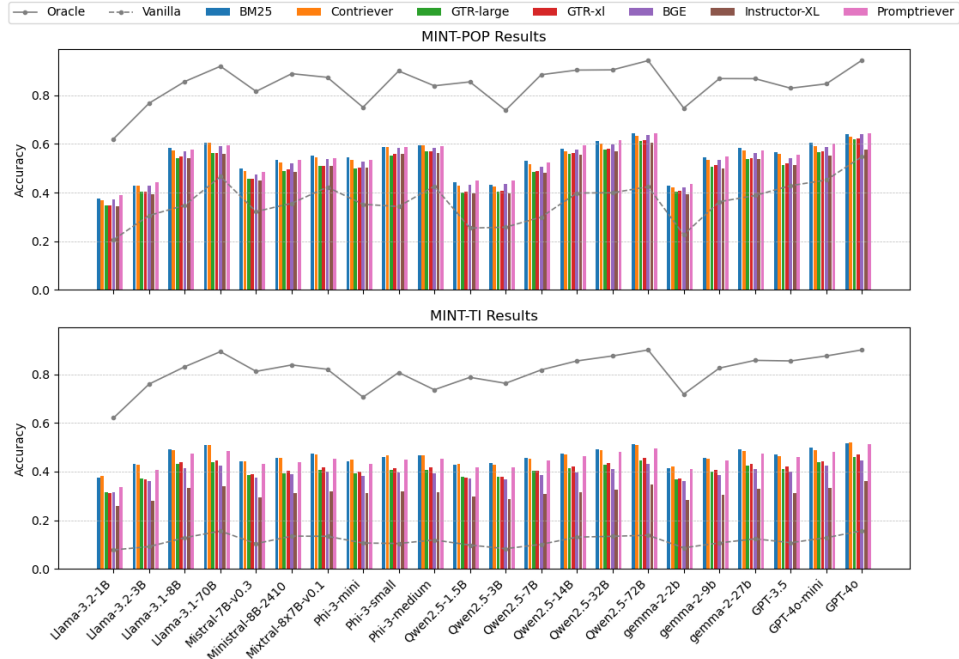


Figure 14: Performance comparison of LLMs on MINTQA-POP and MINTQA-TI using different retrieval methods. “Oracle” uses gold-standard retrieval passages, while “Vanilla” involves models answering without retrieval content.



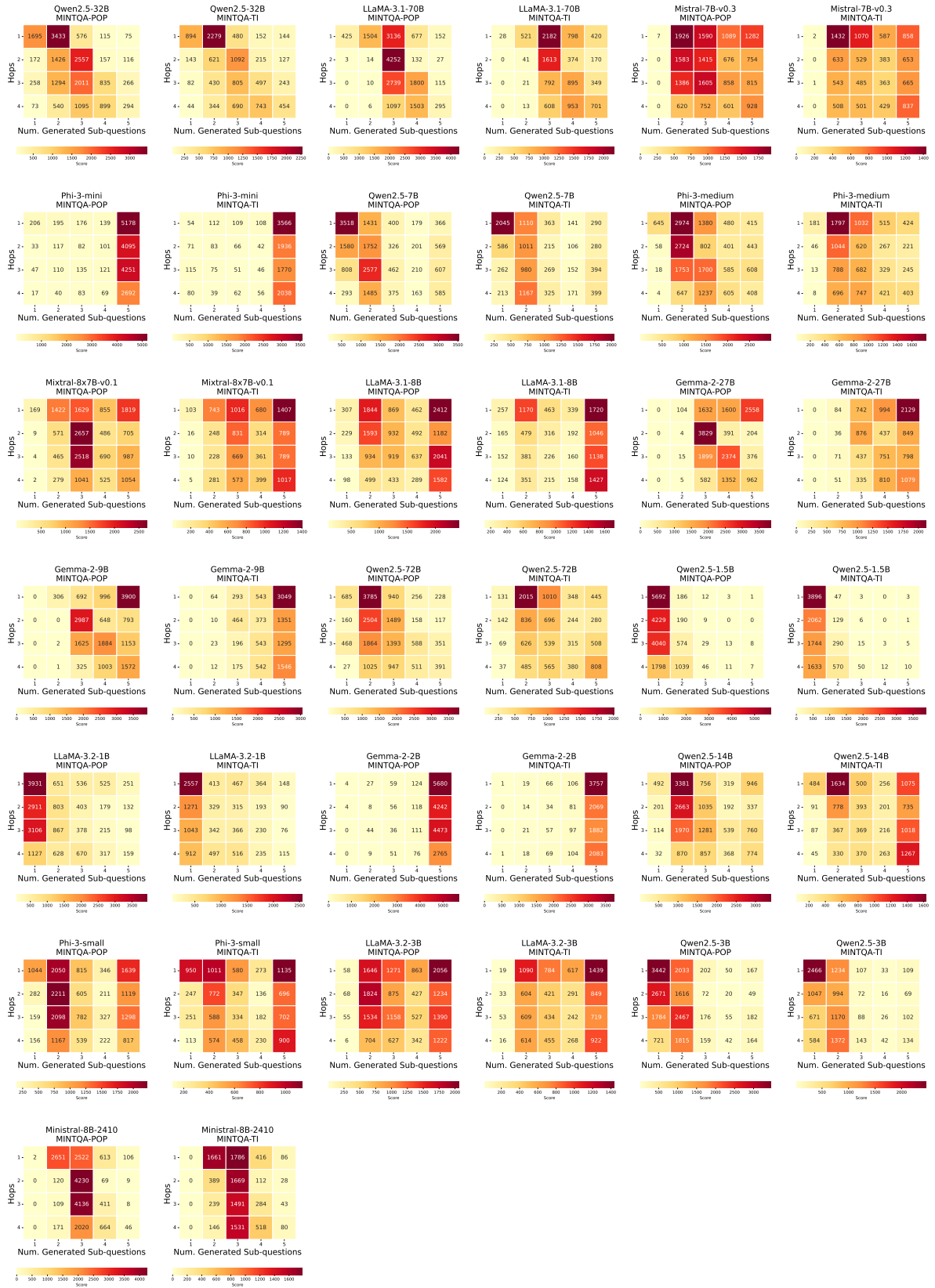


Figure 15: The confusion matrix of the number of sub-questions generated by the LLMs for main questions categorized by hops in the setting of purely generating sub-questions.

Model	BM25			Contriever			PromptRetrieval		
	Acc (%)	Avg. Sub	Avg. Ret	Acc (%)	Avg. Sub	Avg. Ret	Acc (%)	Avg. Sub	Avg. Ret
MINTQA-POP									
<b>Qwen Models</b>									
Qwen2.5-1.5B	25.86 (-0.50)	1.13 (+1.00)	0.32 (+0.19)	26.02 (-0.36)	1.15 (+1.02)	0.31 (+0.18)	26.13 (-0.35)	1.15 (+1.01)	0.32 (+0.18)
Qwen2.5-3B	31.16 (-0.13)	1.78 (+0.95)	1.54 (+0.71)	29.56 (-1.64)	1.70 (+0.89)	1.48 (+0.67)	31.55 (-2.50)	1.69 (+0.89)	1.47 (+0.67)
Qwen2.5-7B	32.98 (-2.68)	2.18 (+0.97)	1.14 (-0.07)	34.89 (-1.90)	2.11 (+0.93)	1.09 (-0.09)	36.58 (-3.93)	2.09 (+0.92)	1.10 (-0.07)
Qwen2.5-14B	53.77 (-10.02)	3.44 (+1.10)	1.22 (-1.12)	53.35 (-9.80)	3.45 (+0.99)	1.23 (-1.23)	55.53 (-12.39)	3.41 (+1.10)	1.21 (-1.10)
Qwen2.5-32B	50.33 (-5.86)	2.79 (+1.02)	1.18 (-0.59)	48.56 (-5.03)	2.77 (+0.92)	1.19 (-0.66)	50.84 (-6.40)	2.81 (+0.91)	1.18 (-0.72)
Qwen2.5-72B	58.63 (-7.93)	3.01 (+1.10)	1.35 (-0.56)	57.42 (+15.28)	3.05 (+3.05)	1.36 (+1.36)	60.89 (-9.33)	3.02 (+1.09)	1.32 (-0.61)
<b>Llama Models</b>									
LLaMA-3.2-1B	20.53 (-0.75)	1.79 (+1.00)	1.36 (+0.57)	20.86 (-0.47)	1.79 (+1.00)	1.32 (+0.53)	21.13 (-0.46)	1.80 (+1.00)	1.32 (+0.52)
LLaMA-3.2-3B	37.23 (-0.64)	3.48 (+0.72)	3.26 (+0.50)	37.70 (-0.53)	3.49 (+0.68)	3.28 (+0.47)	38.04 (-0.27)	3.50 (+0.70)	3.29 (+0.49)
LLaMA-3.1-8B	50.01 (-0.72)	3.88 (+0.56)	3.79 (+0.47)	50.19 (-0.13)	4.05 (+0.49)	3.96 (+0.40)	54.38 (-0.54)	4.01 (+0.56)	3.91 (+0.46)
LLaMA-3.1-70B	64.80 (+0.17)	3.41 (+0.89)	3.40 (+0.88)	62.55 (+19.98)	3.38 (+3.38)	3.37 (+3.37)	69.21 (-0.07)	3.31 (+0.96)	3.30 (+0.95)
<b>Mistral Models</b>									
Mistral-7B-v0.3	29.47 (-4.23)	3.13 (+0.76)	1.84 (-0.53)	28.80 (-5.18)	3.24 (+0.67)	1.77 (-0.80)	30.73 (-5.60)	3.20 (+0.72)	1.84 (-0.64)
Ministral-8B-2410	35.91 (-16.49)	2.95 (+0.96)	0.02 (-1.97)	36.04 (-19.19)	2.94 (+0.95)	0.02 (-1.97)	36.04 (-20.94)	2.94 (+0.97)	0.01 (-1.96)
Mixtral-8x7B-v0.1	48.28 (-7.08)	3.61 (+0.77)	1.29 (-1.55)	48.05 (-10.47)	3.60 (+0.86)	1.30 (-1.44)	49.00 (-13.10)	3.59 (+0.83)	1.23 (-1.53)
<b>Phi Models</b>									
Phi-3-mini	26.81 (-2.54)	4.75 (+0.04)	2.23 (-2.48)	25.64 (-0.12)	4.71 (+0.03)	2.38 (-2.30)	27.63 (-1.79)	4.71 (+0.03)	2.28 (-2.40)
Phi-3-small	37.09 (-4.58)	2.92 (+0.44)	0.78 (-1.70)	37.19 (-6.28)	2.88 (+0.65)	0.79 (-1.44)	39.14 (-8.85)	2.87 (+0.62)	0.77 (-1.48)
Phi-3-medium	40.16 (-5.24)	2.98 (+0.87)	1.08 (-1.03)	39.65 (-5.31)	2.94 (+0.80)	1.07 (-1.07)	41.90 (-8.05)	2.96 (+0.81)	1.07 (-1.08)
<b>Gemma Models</b>									
Gemma-2-2B	34.20 (-13.31)	4.96 (+0.13)	0.99 (-3.84)	34.62 (-11.07)	4.96 (+0.07)	1.00 (-3.89)	35.62 (-14.95)	4.96 (+0.07)	0.99 (-3.90)
Gemma-2-9B	39.99 (-18.52)	3.93 (+0.18)	0.32 (-3.43)	40.20 (-16.35)	3.92 (+0.11)	0.34 (-3.47)	40.56 (-23.73)	3.92 (-0.07)	0.33 (-3.66)
Gemma-2-27B	64.64 (-0.37)	4.05 (+0.74)	4.01 (+0.70)	60.82 (-0.50)	4.09 (+0.71)	4.04 (+0.66)	68.64 (-0.40)	4.32 (+0.57)	4.28 (+0.53)
MINTQA-TI									
<b>Qwen Models</b>									
Qwen2.5-1.5B	9.49 (-0.29)	1.13 (+1.00)	0.35 (+0.22)	9.39 (-0.20)	1.14 (+1.01)	0.36 (+0.23)	9.47 (-0.11)	1.14 (+1.01)	0.36 (+0.23)
Qwen2.5-3B	18.18 (+0.00)	1.83 (+0.95)	1.65 (+0.77)	17.55 (-1.41)	1.78 (+0.91)	1.62 (+0.75)	17.76 (-1.27)	1.75 (+0.89)	1.59 (+0.73)
Qwen2.5-7B	20.36 (-2.60)	2.35 (+0.92)	1.58 (+0.15)	21.83 (-3.93)	2.21 (+0.92)	1.54 (+0.25)	21.34 (-3.26)	2.21 (+0.91)	1.56 (+0.26)
Qwen2.5-14B	39.40 (-10.01)	3.65 (+0.84)	2.01 (-0.80)	41.53 (-10.61)	3.63 (+0.89)	2.00 (-0.74)	40.91 (-10.25)	3.62 (+0.85)	1.98 (-0.79)
Qwen2.5-32B	33.87 (-6.81)	2.86 (+1.02)	1.87 (+0.03)	34.91 (-6.53)	2.87 (+0.98)	1.89 (+0.00)	33.75 (-6.75)	2.89 (+0.97)	1.89 (-0.03)
Qwen2.5-72B	44.79 (-8.13)	3.42 (+0.95)	2.04 (-0.43)	45.99 (+32.78)	3.45 (+3.45)	2.06 (+2.06)	44.44 (-9.64)	3.47 (+0.92)	2.07 (-0.48)
<b>Llama Models</b>									
LLaMA-3.2-1B	9.25 (-2.17)	1.94 (+0.99)	1.41 (+0.46)	9.47 (-2.36)	1.98 (+1.00)	1.34 (+0.36)	8.99 (-2.29)	1.99 (+0.99)	1.36 (+0.36)
LLaMA-3.2-3B	25.89 (-0.69)	3.73 (+0.65)	3.47 (+0.39)	27.39 (-0.68)	3.75 (+0.64)	3.49 (+0.38)	24.82 (-0.57)	3.74 (+0.62)	3.48 (+0.36)
LLaMA-3.1-8B	38.24 (-0.51)	3.99 (+0.49)	3.90 (+0.40)	40.40 (-0.64)	4.10 (+0.46)	4.01 (+0.37)	39.02 (-0.61)	4.11 (+0.48)	4.02 (+0.39)
LLaMA-3.1-70B	51.99 (-0.31)	3.70 (+0.91)	3.68 (+0.89)	52.16 (+39.78)	3.69 (+3.69)	3.67 (+3.67)	50.59 (-0.27)	3.71 (+0.90)	3.69 (+0.88)
<b>Mistral Models</b>									
Mistral-7B-v0.3	9.86 (-7.75)	3.19 (+0.76)	1.72 (-0.71)	9.77 (-8.87)	3.41 (+0.65)	1.76 (-1.00)	10.16 (-8.92)	3.36 (+0.70)	1.81 (-0.85)
Ministral-8B-2410	10.69 (-29.81)	2.99 (+0.89)	0.10 (-2.00)	10.75 (-33.21)	2.99 (+0.88)	0.10 (-2.01)	10.71 (-33.46)	2.98 (+0.90)	0.10 (-1.98)
Mixtral-8x7B-v0.1	24.91 (-20.91)	3.97 (+0.78)	1.33 (-1.86)	26.19 (-21.47)	3.95 (+0.77)	1.32 (-1.86)	25.72 (-20.78)	3.96 (+0.80)	1.29 (-1.87)
<b>Phi Models</b>									
Phi-3-mini	17.38 (-3.95)	4.68 (+0.05)	2.44 (-2.19)	16.94 (-3.91)	4.65 (+0.06)	2.64 (-1.95)	17.21 (-3.43)	4.65 (+0.05)	2.51 (-2.09)
Phi-3-small	15.38 (-8.37)	3.15 (+0.48)	1.23 (-1.44)	15.71 (+15.71)	3.11 (+3.11)	1.22 (+1.22)	16.20 (-11.60)	3.11 (+0.61)	1.21 (-1.29)
Phi-3-medium	19.58 (-12.64)	3.19 (+0.80)	1.35 (-1.04)	19.12 (-12.92)	3.21 (+0.77)	1.34 (-1.10)	19.18 (-11.87)	3.19 (+0.77)	1.33 (-1.09)
<b>Gemma Models</b>									
Gemma-2-2B	24.67 (-12.69)	4.93 (+0.06)	1.78 (-3.09)	25.34 (-13.12)	4.94 (+0.05)	1.82 (-3.07)	24.48 (-12.59)	4.93 (+0.07)	1.80 (-3.06)
Gemma-2-9B	15.29 (-29.38)	4.59 (+0.29)	0.91 (-3.39)	16.08 (-26.57)	4.58 (+0.27)	0.91 (-3.40)	16.50 (-26.89)	4.58 (+0.19)	0.92 (-3.47)
Gemma-2-27B	48.39 (-1.21)	4.49 (+0.51)	4.43 (+0.45)	45.35 (-0.30)	4.53 (+0.47)	4.47 (+0.41)	47.45 (-0.47)	4.53 (+0.48)	4.47 (+0.42)

Table 11: The full results for Decomposition-Dynamic Retrieve. **Acc** represents the model’s accuracy (%), **Avg. Sub** is the average number of sub-questions generated, **Avg. Ret** is the average number of sub-questions actually used for retrieval. The difference between **DDR** and **DTR** is shown in brackets ( $\geq 0$  is **green**,  $< 0$  is **red**).

---

You are a powerful multi-hop question generator. Users will provide a chain of Wikidata triplets, and you will help write questions to ask the tail entity from the head entity. The format of a wikidata triple is (subject, relation, object). You shouldn't include bridge entities in generated questions. The questions should only include the head entity. **All involved relations must be reflected in the question.**

**#Example 1**

**Wikidata triplets:** (Four Peaks, mountain range, x1), (x1, located in the administrative territorial entity, x2), (x2, located in the administrative territorial entity, x3), (x3, office held by head of government, x4)

**Generated question:** Who holds the office of the head of government for the administrative entity where the mountain range Four Peaks is located??

**#Example 2**

**Wikidata triplets:** (Alena Vostrá, place of birth, x1)

**Generated question:** Where was Alena Vostrá born?

**#Example 3**

**Wikidata triplets:** (Anguilla, country, x1), (x1, capital, x2)

**Generated question:** what is the capital of the country of the Anguilla?

**#Example 4**

**Wikidata triplets:** (Nazko River, mouth of the watercourse, x1), (x1, mouth of the watercourse, x2), (x2, country, x3)

**Generated question:** In which country does the Nazko River ultimately discharge its waters?

**#Example 5**

**Wikidata triplets:** {Sampled facts}

**Generated question:**

---

Table 12: The prompt used to generate questions is based on sampled facts. Additionally, we include 4 demonstrations showcasing examples ranging from 1-hop to 4-hop reasoning.

---

You are a powerful question answering system. Users will provide a question and useful context. The provided context are some wikidata triplets which format is (subject, relation, object). You should answer the question based on the context. The answer should be a single entity or a list of entities. If the answer is a list of entities, you should return the most relevant one.

Context: {related documents}

Question: {question}

---

Table 13: The prompt used for question quality inspection provides a given question and its corresponding facts. We aim for the GPT-4o to correctly answer the question based on this information.

---

Below is a question, please answer it directly  
and keep your answer as short as possible.  
**Question:** {question }  
**Answer:**

---

Table 14: The prompt designed to guide the model in providing a concise answer directly to the question.

---

Given some related documents: {retrieved\_documents}. This is a question: {question}. Please answer the question directly. Please keep your answer as short as possible.  
**Answer:**

---

Table 15: The prompt instructs the model to provide a concise answer to the question based on the retrieved documents.



---

Here is a question: {question}

To answer this question. You have to three choices now:

⟨**choice A**⟩ Generate a sub-question.

⟨**choice B**⟩ Answer the question directly if you are confident to answer it.

⟨**choice C**⟩ retrieve some document to help you answer the question.

If you choose ⟨**choice A**⟩, please output:

```
{{"choice A": {{ "sub-question": "your_sub_question_here" }}}}
```

If you choose ⟨**choice B**⟩, please output:

```
{{"choice B": {{ "answer": "your_answer_here" }}}}
```

If you choose ⟨**choice C**⟩, please output:

```
{{"choice C": retrieval}}
```

The final output should be in the form of a JSON string, without any additional content. Please keep your answer as short as possible.

Output:

---

Table 16: The prompt is used for retrieval tasks, directly generating answers or creating sub-questions for judgment purposes.

---

Given a question: {question}

The subsequent sub-questions: {sub\_questions}

You have two choices now:

⟨**choice A**⟩ answer the final sub-question directly.

⟨**choice B**⟩ retrieve some document to help you answer the question. Just output retrieval as a placeholder.

If you choose ⟨**choice A**⟩, please output:

```
{{"choice A": {{ "answer": "your_answer_here" }}}}
```

If you choose ⟨**choice B**⟩, please output:

```
{{"choice B": retrieval}}
```

The final output should be in the form of a JSON string, without any additional content. Please keep your answer as short as possible.

Output:

---

Table 17: The prompt is used for evaluating sub-questions, performing retrieval, or directly generating answers.

---

Given a main question: {question}  
And sub-question-answer pairs: {sub\_question\_answer\_pairs}

Please judge if the main question has been finished. You have two choices now:

⟨**choice A**⟩ The answer can be found in the sub-question-answer pairs. If you choose this choice, please output the final answer.

⟨**choice B**⟩ The answer cannot be found and a new sub-question needs to be generated.

If you choose ⟨**choice A**⟩, please output:  
{{"choice A": {"answer": "final\_answer\_here"}}}}

If you choose ⟨**choice B**⟩, please output:  
{{"choice B": {"sub-question": "new\_sub-question\_here"}}}}

The final output should be in the form of a JSON string, without any additional content. Please keep your answer as short as possible.

Output:

---

Table 18: The prompt provides sub-questions and their answers, requiring the model to determine whether the answer to the main question has been found.

---

To answer this question, you may need to generate subquestions following these guidelines:

Given a main question and optional previous subquestion-answer pairs, you may need to generate subquestions to help answer this main question. Please ensure to only generate subquestions that are relevant to answering the main question. When there are no more subquestions needed, output "finish".

**Input Format**

liRequired:

- Main Question: [question]

Optional:

- Previous Subquestion: [subquestion]

- Previous Answer: [subanswer]

**Output Format**

One of:

- Next Subquestion: [new subquestion]

- "finish" (when no further subquestions are needed)

**Generation Guidelines**

1. Subquestions should:

- Break down complex aspects of the main question

- Follow a logical progression

- Be specific and focused

- Build upon previous answers when available

2. Output "finish" when:

- All relevant aspects have been covered

- Further breakdown would not add value

- The question has been fully addressed

**Examples**

Example 1:

Input:

- Main Question: "What is the location of the headquarters of the institution where Percival Lowell was educated?"

- Previous Subquestion: "Where did Percival Lowell receive his education?"

- Previous Answer: "Harvard University."

Output:

- Next Subquestion: "Where is the headquarters of Harvard University?"

Example 2:

Input:

- Main Question: "What is the capital of France?"

Output:

- "finish"

Main Question: {question}

{previous\_subquestion\_answer\_pairs}

Output:

---

Table 19: This prompt indicates that the model determines whether to generate sub-questions based on the previous answer history, to further break down the main question, or to finish the continued generation of sub-questions.

---

Based on the main question and all subquestion-answer pairs, please provide a comprehensive final answer. Please keep your answer as short as possible.

Main Question: {main\_question}

Previous Subquestions and Answers:

{history\_str}

Final Answer:

---

Table 20: The prompt instructs the model to summarize and generate the answer to the main question based on the sub-questions and their answers.

---

Answer the following question based on your internal knowledge with one or few words.

Add a confidence indicator after your answer: - "certain" if you are completely confident in the accuracy - "uncertain" if you have any doubts

**Input Format**

Input:

- Question: [question]

**Output Format**

Output:

- Answer: [brief answer]

- Confidence: [certain/uncertain]

Question: {question}

Output:

---

Table 21: The prompt requires the model to output a confidence score for the generated sub-questions, which will be used to determine whether retrieval is necessary.



<b>Triplets:</b> [[Pigeon Bay Domain, country, New Zealand]] <b>Main Question:</b> In which country is Pigeon Bay Domain located? <b>Main Answer:</b> New Zealand <b>Type:</b> New
<b>Triplets:</b> [[Eveline Hoffmann, place of detention, Theresienstadt Ghetto]] <b>Main Question:</b> Where was Eveline Hoffmann detained? <b>Main Answer:</b> Theresienstadt Ghetto <b>Type:</b> Old

Table 22: One-hop question-answer pairs and their corresponding types in MINTQA-TI.

<b>Triplets:</b> [[Scram Kitty and his Buddy on Rails, publisher, Dakko Dakko], [Dakko Dakko, industry, video game industry]] <b>Main Question:</b> In which industry does the publisher of Scram Kitty and his Buddy on Rails operate? <b>Main Answer:</b> video game industry <b>Subquestion pairs:</b> <b>Sub-question 0:</b> Who is the publisher of Scram Kitty and his Buddy on Rails? <b>Sub-answer 0:</b> Dakko Dakko. <b>Type:</b> New <b>Sub-question 1:</b> In which industry does Dakko Dakko operate? <b>Sub-answer 1:</b> video game industry. <b>Type:</b> New
<b>Triplets:</b> [[CineKink NYC, location, New York City], [New York City, capital of, United States of America]] <b>Main Question:</b> CineKink NYC is located in the city that is the capital of which entity? <b>Main Answer:</b> United States of America <b>Subquestion pairs:</b> <b>Sub-question 0:</b> Where is CineKink NYC located? <b>Sub-answer 0:</b> New York City. <b>Type:</b> New <b>Sub-question 1:</b> What entity has New York City as its capital? <b>Sub-answer 1:</b> United States of America. <b>Type:</b> Old
<b>Triplets:</b> [[Sanna Aunesluoma, residence, Espoo], [Espoo, member of, Union of the Baltic Cities]] <b>Main Question:</b> Which organization or group is the residence of Sanna Aunesluoma a member of? <b>Main Answer:</b> Union of the Baltic Cities <b>Subquestion pairs:</b> <b>Sub-question 0:</b> Where does Sanna Aunesluoma reside? <b>Sub-answer 0:</b> Espoo. <b>Type:</b> Old <b>Sub-question 1:</b> Of which entity is Espoo a member? <b>Sub-answer 1:</b> Union of the Baltic Cities. <b>Type:</b> New
<b>Triplets:</b> [[Horst Hoffmann, country of citizenship, German Democratic Republic], [German Democratic Republic, legislative body, Volkskammer]] <b>Main Question:</b> What is the legislative body of the country where Horst Hoffmann holds citizenship? <b>Main Answer:</b> Volkskammer <b>Subquestion pairs:</b> <b>Sub-question 0:</b> What is the country of citizenship of Horst Hoffmann? <b>Sub-answer 0:</b> German Democratic Republic. <b>Type:</b> Old <b>Sub-question 1:</b> What is the legislative body of the German Democratic Republic? <b>Sub-answer 1:</b> Volkskammer. <b>Type:</b> Old

Table 23: Two-hop question-answer pairs and their corresponding types in MINTQA-TI.

<p><b>Triplets:</b> [[Systems and methods for mesh augmentation and prevention of incisional hernia, owned by, The Trustees of the University of Pennsylvania], [The Trustees of the University of Pennsylvania, headquarters location, Philadelphia], [Philadelphia, member of, Organization of World Heritage Cities]]</p> <p><b>Main Question:</b> Of which entity is the headquarters location of the owner of the "Systems and methods for mesh augmentation and prevention of incisional hernia" a member?</p> <p><b>Main Answer:</b> Organization of World Heritage Cities</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who owns the patent for Systems and methods for mesh augmentation and prevention of incisional hernia? <b>Sub-answer 0:</b> The Trustees of the University of Pennsylvania. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Where is the headquarters of The Trustees of the University of Pennsylvania located? <b>Sub-answer 1:</b> Philadelphia. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> What is Philadelphia a member of? <b>Sub-answer 2:</b> Organization of World Heritage Cities. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[De grote Gwen en Geraldine show, nominated for, Dutch Podcast Award for Chatcast Vermaak], [Dutch Podcast Award for Chatcast Vermaak, country, Netherlands], [Netherlands, language used, Dutch]]</p> <p><b>Main Question:</b> What is the language used in the country for which "De grote Gwen en Geraldine show" was nominated?</p> <p><b>Main Answer:</b> Dutch</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> For what award was "De grote Gwen en Geraldine show" nominated? <b>Sub-answer 0:</b> Dutch Podcast Award for Chatcast Vermaak. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> In which country is the Dutch Podcast Award for Chatcast Vermaak given? <b>Sub-answer 1:</b> Netherlands. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> What language is used in the Netherlands? <b>Sub-answer 2:</b> Dutch. <b>Type:</b> Old</p>
<p><b>Triplets:</b> [[Gathering to Celebrate Old Age, creator, Tomioka Tessai], [Tomioka Tessai, location, Tokyo National Museum], [Tokyo National Museum, member of, Japan Consortium for Open Access Repository]]</p> <p><b>Main Question:</b> Which organization or group is the location associated with the creator of "Gathering to Celebrate Old Age" a member of?</p> <p><b>Main Answer:</b> Japan Consortium for Open Access Repository</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the creator of Gathering to Celebrate Old Age? <b>Sub-answer 0:</b> Tomioka Tessai. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Where is Tomioka Tessai located? <b>Sub-answer 1:</b> Tokyo National Museum. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> What organization or association is the Tokyo National Museum a member of? <b>Sub-answer 2:</b> Japan Consortium for Open Access Repository. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[The Woman Who Cooked Her Husband, author, Debbie Isitt], [Debbie Isitt, country of citizenship, United Kingdom], [United Kingdom, continent, Europe]]</p> <p><b>Main Question:</b> On which continent does the author of "The Woman Who Cooked Her Husband" hold citizenship?</p> <p><b>Main Answer:</b> Europe</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the author of "The Woman Who Cooked Her Husband"? <b>Sub-answer 0:</b> Debbie Isitt. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> What is the country of citizenship of Debbie Isitt? <b>Sub-answer 1:</b> United Kingdom. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> On which continent is the United Kingdom located? <b>Sub-answer 2:</b> Europe. <b>Type:</b> Old</p>
<p><b>Triplets:</b> [[Mubarak Shah, religion or worldview, Islam], [Islam, item operated, Qalab], [Qalab, cause of death, Ajal]]</p> <p><b>Main Question:</b> What was the cause of death for the operator of the religion or worldview followed by Mubarak Shah?</p> <p><b>Main Answer:</b> Ajal</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> What is the religion or worldview of Mubarak Shah? <b>Sub-answer 0:</b> Islam. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> What item is operated by Islam? <b>Sub-answer 1:</b> Qalab. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> What was the cause of death for Qalab? <b>Sub-answer 2:</b> Ajal. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[Felipe Borrego Estrada, place of birth, Zacatecas], [Zacatecas, member of, Organization of World Heritage Cities], [Organization of World Heritage Cities, headquarters location, Quebec City]]</p> <p><b>Main Question:</b> Where is the headquarters of the entity that the birthplace of Felipe Borrego Estrada is a member of?</p> <p><b>Main Answer:</b> Quebec City</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where was Felipe Borrego Estrada born? <b>Sub-answer 0:</b> Zacatecas. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Of which organization is Zacatecas a member? <b>Sub-answer 1:</b> Organization of World Heritage Cities. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Where is the headquarters of the Organization of World Heritage Cities located? <b>Sub-answer 2:</b> Quebec City. <b>Type:</b> Old</p>
<p><b>Triplets:</b> [[Hykjeberget, operator, Dalarna County Administrative Board], [Dalarna County Administrative Board, headquarters location, Falun], [Falun, twinned administrative body, Hamina]]</p> <p><b>Main Question:</b> What administrative body is twinned with the location of the headquarters of the operator of Hykjeberget?</p> <p><b>Main Answer:</b> Hamina</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who operates Hykjeberget? <b>Sub-answer 0:</b> Dalarna County Administrative Board. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Where is the headquarters of the Dalarna County Administrative Board located? <b>Sub-answer 1:</b> Falun. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Which administrative body is twinned with Falun? <b>Sub-answer 2:</b> Hamina. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[University of California Italian Studies Multicampus Research Group, country, United States of America], [United States of America, highest point, Denali], [Denali, mountain range, Alaska Range]]</p> <p><b>Main Question:</b> What is the mountain range that contains the highest point in the country where the University of California Italian Studies Multicampus Research Group is located?</p> <p><b>Main Answer:</b> Alaska Range</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> In which country is the University of California Italian Studies Multicampus Research Group located? <b>Sub-answer 0:</b> United States of America. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> What is the highest point in the United States of America? <b>Sub-answer 1:</b> Denali. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> In which mountain range is Denali located? <b>Sub-answer 2:</b> Alaska Range. <b>Type:</b> Old</p>

Table 24: Three-hop question-answer pairs and their corresponding types in MINTQA-T1.

<p><b>Triples:</b> [[Patricia Florence Suthers, sibling, Elaine Suthers], [Elaine Suthers, mother, Elsie Suthers], [Elsie Suthers, country of citizenship, United Kingdom], [United Kingdom, highest point, Ben Nevis]]</p> <p><b>Main Question:</b> What is the highest point in the country where the mother of Patricia Florence Suthers' sibling is a citizen?</p> <p><b>Main Answer:</b> Ben Nevis</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the sibling of Patricia Florence Suthers? <b>Sub-answer 0:</b> Elaine Suthers. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Who is the mother of Elaine Suthers? <b>Sub-answer 1:</b> Elsie Suthers. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Which country is Elsie Suthers a citizen of? <b>Sub-answer 2:</b> United Kingdom. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> What is the highest point in the United Kingdom? <b>Sub-answer 3:</b> Ben Nevis. <b>Type:</b> Old</p>
<p><b>Triples:</b> [[Patricia Florence Suthers, mother, Elsie Suthers], [Elsie Suthers, spouse, Robert Suthers], [Robert Suthers, relative, Miriam Farid], [Miriam Farid, country of citizenship, United Kingdom]]</p> <p><b>Main Question:</b> What is the country of citizenship of the relative of Patricia Florence Suthers' mother's spouse?</p> <p><b>Main Answer:</b> United Kingdom</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the mother of Patricia Florence Suthers? <b>Sub-answer 0:</b> Elsie Suthers. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Who is the spouse of Elsie Suthers? <b>Sub-answer 1:</b> Robert Suthers. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Who is a relative of Robert Suthers? <b>Sub-answer 2:</b> Miriam Farid. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> Which country is Miriam Farid a citizen of? <b>Sub-answer 3:</b> United Kingdom. <b>Type:</b> New</p>
<p><b>Triples:</b> [[May Hnin Aw Kanya, mother, May Hnin Htapi], [May Hnin Htapi, father, Loethai], [Loethai, child, Lithai], [Lithai, notable work, Traibhumikatha]]</p> <p><b>Main Question:</b> What is the notable work of the child of the father of the mother of May Hnin Aw Kanya?</p> <p><b>Main Answer:</b> Traibhumikatha</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the mother of May Hnin Aw Kanya? <b>Sub-answer 0:</b> May Hnin Htapi. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Who is the spouse of Elsie Suthers? <b>Sub-answer 1:</b> Loethai. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Who is the child of Loethai? <b>Sub-answer 2:</b> Lithai. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> What is a notable work created by Lithai? <b>Sub-answer 3:</b> Traibhumikatha. <b>Type:</b> New</p>
<p><b>Triples:</b> [[SEOlytics, parent organization, Sistrix], [Sistrix, country, Germany], [Germany, continent, Europe], [Europe, shares border with, Asia]]</p> <p><b>Main Question:</b> Which continent shares a border with the continent where the country of SEOlytics' parent organization is located?</p> <p><b>Main Answer:</b> Asia</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> What is the parent organization of SEOlytics? <b>Sub-answer 0:</b> Sistrix. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> In which country is Sistrix located? <b>Sub-answer 1:</b> Germany. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> On which continent is Germany located? <b>Sub-answer 2:</b> Europe. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> Which continent shares a border with Europe? <b>Sub-answer 3:</b> Asia. <b>Type:</b> Old</p>
<p><b>Triples:</b> [[Sri Dhamasokaraj, relative, Saileuthai], [Saileuthai, father, Lithai], [Lithai, sibling, May Hnin Htapi], [May Hnin Htapi, place of death, Mottama]]</p> <p><b>Main Question:</b> Where did the sibling of the father of Sri Dhamasokaraj pass away?</p> <p><b>Main Answer:</b> Mottama</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is a relative of Sri Dhamasokaraj? <b>Sub-answer 0:</b> Saileuthai. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Who was the father of Saileuthai? <b>Sub-answer 1:</b> Lithai. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Who is Lithai's sibling? <b>Sub-answer 2:</b> May Hnin Htapi. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> Where did May Hnin Htapi die? <b>Sub-answer 3:</b> Mottama. <b>Type:</b> New</p>
<p><b>Triples:</b> [[Frank Gailor, educated at, New College], [New College, founded by, William of Wykeham], [William of Wykeham, country of citizenship, Kingdom of England], [Kingdom of England, replaced by, Kingdom of Great Britain]]</p> <p><b>Main Question:</b> Which entity replaced the country of citizenship of the founder of the institution where Frank Gailor was educated?</p> <p><b>Main Answer:</b> Kingdom of Great Britain</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where was Frank Gailor educated? <b>Sub-answer 0:</b> New College. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> Who founded New College? <b>Sub-answer 1:</b> William of Wykeham. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Which country was William of Wykeham a citizen of? <b>Sub-answer 2:</b> Kingdom of England. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> What entity replaced the Kingdom of England? <b>Sub-answer 3:</b> Kingdom of Great Britain. <b>Type:</b> Old</p>
<p><b>Triples:</b> [[The Life You Can Save, author, Peter Singer], [Peter Singer, mother, Cora Singer], [Cora Singer, father, David Ernst Oppenheim], [David Ernst Oppenheim, academic degree, doctorate]]</p> <p><b>Main Question:</b> What academic degree does the father of the author of "The Life You Can Save" hold?</p> <p><b>Main Answer:</b> doctorate</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the author of "The Life You Can Save"? <b>Sub-answer 0:</b> Peter Singer. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Who is Peter Singer's mother? <b>Sub-answer 1:</b> Cora Singer. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Who is the father of Cora Singer? <b>Sub-answer 2:</b> David Ernst Oppenheim. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> What academic degree does David Ernst Oppenheim hold? <b>Sub-answer 3:</b> doctorate. <b>Type:</b> Old</p>
<p><b>Triples:</b> [[Geoffrey Howe, creator, June Mendoza], [June Mendoza, place of birth, Melbourne], [Melbourne, located in or next to body of water, Yarra River], [Yarra River, continent, Australian continent]]</p> <p><b>Main Question:</b> On which continent is the body of water located next to the place where the creator Geoffrey Howe was born?</p> <p><b>Main Answer:</b> Australian continent</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> What did Geoffrey Howe create? <b>Sub-answer 0:</b> June Mendoza. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Where was June Mendoza born? <b>Sub-answer 1:</b> Melbourne. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Which body of water is Melbourne located near? <b>Sub-answer 2:</b> Yarra River. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> On which continent is the Yarra River located? <b>Sub-answer 3:</b> Australian continent. <b>Type:</b> New</p>

Table 25: Four-hop question-answer pairs and their corresponding types in MINTQA-T1 (part 1).

<p><b>Triplets:</b> [[Descenso a los fascismos, place of publication, Barcelona], [Barcelona, member of, Creative Cities Network], [Creative Cities Network, operator, UNESCO], [UNESCO, operating area, worldwide]]</p> <p><b>Main Question:</b> In what area does the operator of the organization that includes the place where "Descenso a los fascismos" was published operate?</p> <p><b>Main Answer:</b> worldwide</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where was "Descenso a los fascismos" published? <b>Sub-answer 0:</b> Barcelona. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> What organization or group is Barcelona a member of? <b>Sub-answer 1:</b> Creative Cities Network. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Who operates the Creative Cities Network? <b>Sub-answer 2:</b> UNESCO. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> What is the operating area of UNESCO? <b>Sub-answer 3:</b> worldwide. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[Monument to Terenzio Mamiani, commemorates, Terenzio, Count Mamiani della Rovere], [Terenzio, Count Mamiani della Rovere, award received, Order of the Redeemer], [Order of the Redeemer, founded by, Otto of Greece], [Otto of Greece, spouse, Amalia of Oldenburg]]</p> <p><b>Main Question:</b> Who is the spouse of the founder of the award received by the person commemorated by the Monument to Terenzio Mamiani?</p> <p><b>Main Answer:</b> Amalia of Oldenburg</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is commemorated by the Monument to Terenzio Mamiani? <b>Sub-answer 0:</b> Terenzio, Count Mamiani della Rovere. <b>Type:</b> New</p> <p><b>Sub-question 1:</b> What award did Terenzio, Count Mamiani della Rovere receive? <b>Sub-answer 1:</b> Order of the Redeemer. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Who founded the Order of the Redeemer? <b>Sub-answer 2:</b> Otto of Greece. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> Who was the spouse of Otto of Greece? <b>Sub-answer 3:</b> Amalia of Oldenburg. <b>Type:</b> Old</p>
<p><b>Triplets:</b> [[Tansen, religion or worldview, Islam], [Islam, item operated, Qalab], [Qalab, cause of death, Ajal], [Ajal, location, treasures of God in Islam]]</p> <p><b>Main Question:</b> Where did the cause of death of the religious figure associated with Tansen occur?</p> <p><b>Main Answer:</b> treasures of God in Islam</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> What is the religion or worldview associated with Tansen? <b>Sub-answer 0:</b> Islam. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> What item is operated by Islam? <b>Sub-answer 1:</b> Qalab. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> What was the cause of death for Qalab? <b>Sub-answer 2:</b> Ajal. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> Where is Ajal located? <b>Sub-answer 3:</b> treasures of God in Islam. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[Irma Stern, place of birth, Bratislava], [Bratislava, member of, League of Historical Cities], [League of Historical Cities, headquarters location, Kyoto], [Kyoto, highest point, Mount Minako]]</p> <p><b>Main Question:</b> What is the highest point of the location where the headquarters of the entity that includes the birthplace of Irma Stern is situated?</p> <p><b>Main Answer:</b> Mount Minako</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where was Irma Stern born? <b>Sub-answer 0:</b> Bratislava. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Of which organization is Bratislava a member? <b>Sub-answer 1:</b> League of Historical Cities. <b>Type:</b> New</p> <p><b>Sub-question 2:</b> Where is the headquarters of the League of Historical Cities located? <b>Sub-answer 2:</b> Kyoto. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> What is the highest point in Kyoto? <b>Sub-answer 3:</b> Mount Minako. <b>Type:</b> Old</p>
<p><b>Triplets:</b> [[Andrew Cogglesby, present in work, Evan Harrington], [Evan Harrington, author, George Meredith], [George Meredith, spouse, Mary Meredith], [Mary Meredith, cause of death, kidney failure]]</p> <p><b>Main Question:</b> What was the cause of death of the spouse of the author who created the work featuring Andrew Cogglesby?</p> <p><b>Main Answer:</b> kidney failure</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> In which work does Andrew Cogglesby appear? <b>Sub-answer 0:</b> Evan Harrington. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Who is the author of "Evan Harrington"? <b>Sub-answer 1:</b> George Meredith. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Who is the spouse of George Meredith? <b>Sub-answer 2:</b> Mary Meredith. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> What was the cause of death of Mary Meredith? <b>Sub-answer 3:</b> kidney failure. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[Federico Coccozza, employer, Curie Institute], [Curie Institute, founded by, Marie Curie], [Marie Curie, ethnic group, Poles], [Poles, language used, Church Slavonic]]</p> <p><b>Main Question:</b> What language is used by the ethnic group of the founder of Federico Coccozza's employer?</p> <p><b>Main Answer:</b> Church Slavonic</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who employs Federico Coccozza? <b>Sub-answer 0:</b> Curie Institute. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Who founded the Curie Institute? <b>Sub-answer 1:</b> Marie Curie. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> What is the ethnic group of Marie Curie? <b>Sub-answer 2:</b> Poles. <b>Type:</b> New</p> <p><b>Sub-question 3:</b> Which language is used by Poles? <b>Sub-answer 3:</b> Church Slavonic. <b>Type:</b> Old</p>
<p><b>Triplets:</b> [[Devespresso Games, headquarters location, Seoul], [Seoul, member of, Creative Cities Network], [Creative Cities Network, operator, UNESCO], [UNESCO, operating area, worldwide]]</p> <p><b>Main Question:</b> What is the operating area of the operator of the member organization where Devespresso Games' headquarters is located?</p> <p><b>Main Answer:</b> worldwide</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where is the headquarters of Devespresso Games located? <b>Sub-answer 0:</b> Seoul. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Of which organization is Seoul a member? <b>Sub-answer 1:</b> Creative Cities Network. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> Who operates the Creative Cities Network? <b>Sub-answer 2:</b> UNESCO. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> What is the operating area of UNESCO? <b>Sub-answer 3:</b> worldwide. <b>Type:</b> New</p>
<p><b>Triplets:</b> [[Sonetto I, author, Vittorio Alfieri], [Vittorio Alfieri, place of death, Florence], [Florence, present in work, Civilization V], [Civilization V, developer, Firaxis Games]]</p> <p><b>Main Question:</b> Who is the developer of the work where the place of death of the author of Sonetto I is present?</p> <p><b>Main Answer:</b> Firaxis Games</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the author of Sonetto I? <b>Sub-answer 0:</b> Vittorio Alfieri. <b>Type:</b> Old</p> <p><b>Sub-question 1:</b> Where did Vittorio Alfieri die? <b>Sub-answer 1:</b> Florence. <b>Type:</b> Old</p> <p><b>Sub-question 2:</b> In which work is Florence present? <b>Sub-answer 2:</b> Civilization V. <b>Type:</b> Old</p> <p><b>Sub-question 3:</b> Who developed Civilization V? <b>Sub-answer 3:</b> Firaxis Games. <b>Type:</b> Old</p>

Table 26: Four-hop question-answer pairs and their corresponding types in MINTQA-TI (part 2)..

<p><b>Triplets:</b> [[Papanasam taluk, country, India]]</p> <p><b>Main Question:</b> In which country is Papanasam taluk located?</p> <p><b>Main Answer:</b> India</p> <p><b>Type:</b> Popular</p>
<p><b>Triplets:</b> [[Jerod Swallow, sports discipline competed in, ice dance]]</p> <p><b>Main Question:</b> In which sports discipline does Jerod Swallow compete?</p> <p><b>Main Answer:</b> ice dance</p> <p><b>Type:</b> Unpopular</p>

Table 27: One-hop question-answer pairs and their corresponding types in MINTQA-POP.



<b>Triplets:</b> [[Gmina Szypliszki, country, Poland], [Poland, capital, Warsaw]] <b>Main Question:</b> What is the capital of the country where Gmina Szypliszki is located? <b>Main Answer:</b> Warsaw <b>Subquestion pairs:</b> <b>Sub-question 0:</b> In which country is Gmina Szypliszki located? <b>Sub-answer 0:</b> Poland. <b>Type:</b> Popular <b>Sub-question 1:</b> What is the capital of Poland? <b>Sub-answer 1:</b> Warsaw. <b>Type:</b> Popular
<b>Triplets:</b> [[Canary Islands, country, Spain], [Spain, legislative body, Cortes Generales]] <b>Main Question:</b> What is the legislative body of the country to which the Canary Islands belong? <b>Main Answer:</b> Cortes Generales <b>Subquestion pairs:</b> <b>Sub-question 0:</b> Which country are the Canary Islands part of? <b>Sub-answer 0:</b> Spain. <b>Type:</b> Popular <b>Sub-question 1:</b> What is the legislative body of Spain? <b>Sub-answer 1:</b> Cortes Generales. <b>Type:</b> Unpopular
<b>Triplets:</b> [[Pabna Cadet College, country, Bangladesh], [Bangladesh, capital, Dhaka]] <b>Main Question:</b> What is the capital of the country where Pabna Cadet College is located? <b>Main Answer:</b> Dhaka <b>Subquestion pairs:</b> <b>Sub-question 0:</b> In which country is Pabna Cadet College located? <b>Sub-answer 0:</b> Bangladesh. <b>Type:</b> Unpopular <b>Sub-question 1:</b> What is the capital of Bangladesh? <b>Sub-answer 1:</b> Dhaka. <b>Type:</b> Popular
<b>Triplets:</b> [[Brackendale Eagles Provincial Park, country, Canada], [Canada, highest point, Mount Logan]] <b>Main Question:</b> What is the highest point in the country where Brackendale Eagles Provincial Park is located? <b>Main Answer:</b> Mount Logan <b>Subquestion pairs:</b> <b>Sub-question 0:</b> In which country is Brackendale Eagles Provincial Park located? <b>Sub-answer 0:</b> Canada. <b>Type:</b> Unpopular <b>Sub-question 1:</b> What is the highest point in Canada? <b>Sub-answer 1:</b> Mount Logan. <b>Type:</b> Unpopular

Table 28: Two-hop question-answer pairs and their corresponding types in MINTQA-POP.

<b>Triplets:</b> [[Cuzco Department, country, Peru], [Peru, capital, Lima], [Lima, located in or next to body of water, Rímac River]] <b>Main Question:</b> Which body of water is located in or next to the capital of the country where the Cuzco Department is found? <b>Main Answer:</b> Rímac River <b>Subquestion pairs:</b> <b>Sub-question 0:</b> In which country is the Cuzco Department located? <b>Sub-answer 0:</b> Peru. <b>Type:</b> Popular <b>Sub-question 1:</b> What is the capital of Peru? <b>Sub-answer 1:</b> Lima. <b>Type:</b> Popular <b>Sub-question 2:</b> Which body of water is Lima located next to? <b>Sub-answer 2:</b> Rímac River. <b>Type:</b> Unpopular
<b>Triplets:</b> [[Kirkovo Municipality, country, Bulgaria], [Bulgaria, highest point, Musala], [Musala, mountain range, Rila]] <b>Main Question:</b> Which mountain range includes the highest point in the country of Kirkovo Municipality? <b>Main Answer:</b> Rila <b>Subquestion pairs:</b> <b>Sub-question 0:</b> Which country is Kirkovo Municipality located in? <b>Sub-answer 0:</b> Bulgaria. <b>Type:</b> Popular <b>Sub-question 1:</b> What is the highest point in Bulgaria? <b>Sub-answer 1:</b> Musala. <b>Type:</b> Unpopular <b>Sub-question 2:</b> In which mountain range is Musala located? <b>Sub-answer 2:</b> Rila. <b>Type:</b> Unpopular
<b>Triplets:</b> [[Nicu Stroia, participant in, 1992 Summer Olympics], [1992 Summer Olympics, country, Spain], [Spain, capital, Madrid]] <b>Main Question:</b> What is the capital of the country where Nicu Stroia participated in an event? <b>Main Answer:</b> Madrid <b>Subquestion pairs:</b> <b>Sub-question 0:</b> In which events or activities did Nicu Stroia participate? <b>Sub-answer 0:</b> 1992 Summer Olympics. <b>Type:</b> Unpopular <b>Sub-question 1:</b> In which country were the 1992 Summer Olympics held? <b>Sub-answer 1:</b> Spain. <b>Type:</b> Popular <b>Sub-question 2:</b> What is the capital of Spain? <b>Sub-answer 2:</b> Madrid. <b>Type:</b> Popular
<b>Triplets:</b> [[Bunk Moreland, present in work, The Wire], [The Wire, original broadcaster, HBO], [HBO, parent organization, WarnerMedia]] <b>Main Question:</b> What is the parent organization of the original broadcaster of the work featuring Bunk Moreland? <b>Main Answer:</b> WarnerMedia <b>Subquestion pairs:</b> <b>Sub-question 0:</b> In which work does the character Bunk Moreland appear? <b>Sub-answer 0:</b> The Wire. <b>Type:</b> Unpopular <b>Sub-question 1:</b> What is the original broadcaster of The Wire? <b>Sub-answer 1:</b> HBO. <b>Type:</b> Popular <b>Sub-question 2:</b> What is the parent organization of HBO? <b>Sub-answer 2:</b> WarnerMedia. <b>Type:</b> Unpopular
<b>Triplets:</b> [[Ewout van Asbeck, sport, field hockey], [field hockey, country of origin, England], [England, capital, London]] <b>Main Question:</b> What is the capital of the country of origin of the sport in which Ewout van Asbeck participates? <b>Main Answer:</b> London <b>Subquestion pairs:</b> <b>Sub-question 0:</b> What sport does Ewout van Asbeck participate in? <b>Sub-answer 0:</b> field hockey. <b>Type:</b> Unpopular <b>Sub-question 1:</b> Which country is the origin of field hockey? <b>Sub-answer 1:</b> England. <b>Type:</b> Unpopular <b>Sub-question 2:</b> What is the capital of England? <b>Sub-answer 2:</b> London. <b>Type:</b> Popular
<b>Triplets:</b> [[College Hockey in the D, sport, ice hockey], [ice hockey, authority, International Ice Hockey Federation], [International Ice Hockey Federation, headquarters location, Zürich]] <b>Main Question:</b> Where is the headquarters of the authority governing the sport of College Hockey in the D located? <b>Main Answer:</b> Zürich <b>Subquestion pairs:</b> <b>Sub-question 0:</b> What sport is associated with College Hockey in the D? <b>Sub-answer 0:</b> ice hockey. <b>Type:</b> Unpopular <b>Sub-question 1:</b> Which organization is the governing authority for ice hockey? <b>Sub-answer 1:</b> International Ice Hockey Federation. <b>Type:</b> Unpopular <b>Sub-question 2:</b> Where are the headquarters of the International Ice Hockey Federation located? <b>Sub-answer 2:</b> Zürich. <b>Type:</b> Unpopular

Table 29: Three-hop question-answer pairs and their corresponding types in MINTQA-POP.

<p><b>Triplets:</b> [[National Hockey League, sport, ice hockey], [ice hockey, authority, International Ice Hockey Federation], [International Ice Hockey Federation, country, Switzerland], [Switzerland, continent, Europe]]</p> <p><b>Main Question:</b> On which continent is the country that has authority over the sport played in the National Hockey League located?</p> <p><b>Main Answer:</b> Europe</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> What sport is played in the National Hockey League? <b>Sub-answer 0:</b> ice hockey. <b>Type:</b> Popular</p> <p><b>Sub-question 1:</b> Which organization is the governing authority for ice hockey? <b>Sub-answer 1:</b> International Ice Hockey Federation. <b>Type:</b> Unpopular</p> <p><b>Sub-question 2:</b> Which country is the International Ice Hockey Federation based in? <b>Sub-answer 2:</b> Switzerland. <b>Type:</b> Unpopular</p> <p><b>Sub-question 3:</b> On which continent is Switzerland located? <b>Sub-answer 3:</b> Europe. <b>Type:</b> Unpopular</p>
<p><b>Triplets:</b> [[Rafael Bejarano, place of birth, Arequipa], [Arequipa, country, Peru], [Peru, capital, Lima], [Lima, located in or next to body of water, Rímac River]]</p> <p><b>Main Question:</b> Which body of water is the capital of the country where Rafael Bejarano was born located next to?</p> <p><b>Main Answer:</b> Rímac River</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where was Rafael Bejarano born? <b>Sub-answer 0:</b> Arequipa. <b>Type:</b> Unpopular</p> <p><b>Sub-question 1:</b> In which country is Arequipa located? <b>Sub-answer 1:</b> Peru. <b>Type:</b> Popular</p> <p><b>Sub-question 2:</b> What is the capital of Peru? <b>Sub-answer 2:</b> Lima. <b>Type:</b> Popular</p> <p><b>Sub-question 3:</b> Which body of water is Lima located next to? <b>Sub-answer 3:</b> Rímac River. <b>Type:</b> Unpopular</p>
<p><b>Triplets:</b> [[The Perfect Cocktail, part of the series, How I Met Your Mother], [How I Met Your Mother, original broadcaster, CBS], [CBS, owned by, Paramount Global], [Paramount Global, industry, mass media]]</p> <p><b>Main Question:</b> In which industry does the owner of the original broadcaster of the series that includes "The Perfect Cocktail" operate?</p> <p><b>Main Answer:</b> mass media</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Of which series is "The Perfect Cocktail" a part? <b>Sub-answer 0:</b> How I Met Your Mother. <b>Type:</b> Unpopular</p> <p><b>Sub-question 1:</b> Which network originally broadcasted "How I Met Your Mother"? <b>Sub-answer 1:</b> CBS. <b>Type:</b> Popular</p> <p><b>Sub-question 2:</b> Who owns CBS? <b>Sub-answer 2:</b> Paramount Global. <b>Type:</b> Unpopular</p> <p><b>Sub-question 3:</b> In which industry does Paramount Global operate? <b>Sub-answer 3:</b> mass media. <b>Type:</b> Unpopular</p>
<p><b>Triplets:</b> [[Saint George Killing the Dragon, creator, Bernat Martorell], [Bernat Martorell, place of death, Barcelona], [Barcelona, country, Spain], [Spain, capital, Madrid]]</p> <p><b>Main Question:</b> What is the capital of the country where the creator of Saint George Killing the Dragon died?</p> <p><b>Main Answer:</b> Madrid</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who is the creator of Saint George Killing the Dragon? <b>Sub-answer 0:</b> Bernat Martorell. <b>Type:</b> Unpopular</p> <p><b>Sub-question 1:</b> Where did Bernat Martorell die? <b>Sub-answer 1:</b> Barcelona. <b>Type:</b> Unpopular</p> <p><b>Sub-question 2:</b> In which country is Barcelona located? <b>Sub-answer 2:</b> Spain. <b>Type:</b> Popular</p> <p><b>Sub-question 3:</b> What is the capital of Spain? <b>Sub-answer 3:</b> Madrid. <b>Type:</b> Popular</p>
<p><b>Triplets:</b> [[DWNX-FM, owned by, Radio Mindanao Network], [Radio Mindanao Network, headquarters location, Makati], [Makati, country, Philippines], [Philippines, continent, Asia]]</p> <p><b>Main Question:</b> On which continent is the country located where the headquarters of the owner of DWNX-FM is situated?</p> <p><b>Main Answer:</b> Asia</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who owns DWNX-FM? <b>Sub-answer 0:</b> Radio Mindanao Network. <b>Type:</b> Unpopular</p> <p><b>Sub-question 1:</b> Where is the headquarters of Radio Mindanao Network located? <b>Sub-answer 1:</b> Makati. <b>Type:</b> Unpopular</p> <p><b>Sub-question 2:</b> In which country is Makati located? <b>Sub-answer 2:</b> Philippines. <b>Type:</b> Popular</p> <p><b>Sub-question 3:</b> On which continent is the Philippines located? <b>Sub-answer 3:</b> Asia. <b>Type:</b> Unpopular</p>
<p><b>Triplets:</b> [[2008 FA Trophy Final, location, Wembley Stadium], [Wembley Stadium, owned by, The Football Association], [The Football Association, applies to jurisdiction, England], [England, capital, London]]</p> <p><b>Main Question:</b> What is the capital of the jurisdiction that owns the location of the 2008 FA Trophy Final?</p> <p><b>Main Answer:</b> London</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Where was the 2008 FA Trophy Final held? <b>Sub-answer 0:</b> Wembley Stadium. <b>Type:</b> Unpopular</p> <p><b>Sub-question 1:</b> Who owns Wembley Stadium? <b>Sub-answer 1:</b> The Football Association. <b>Type:</b> Unpopular</p> <p><b>Sub-question 2:</b> Which jurisdiction does The Football Association apply to? <b>Sub-answer 2:</b> England. <b>Type:</b> Unpopular</p> <p><b>Sub-question 3:</b> What is the capital of England? <b>Sub-answer 3:</b> London. <b>Type:</b> Popular</p>
<p><b>Triplets:</b> [[Rothschild banking family of France, founded by, James Mayer de Rothschild], [James Mayer de Rothschild, place of birth, Frankfurt], [Frankfurt, located in or next to body of water, Main], [Main, mouth of the watercourse, Rhine]]</p> <p><b>Main Question:</b> Into which body of water does the river located next to the birthplace of the founder of the Rothschild banking family of France ultimately flow?</p> <p><b>Main Answer:</b> Rhine</p> <p><b>Subquestion pairs:</b></p> <p><b>Sub-question 0:</b> Who founded the Rothschild banking family of France? <b>Sub-answer 0:</b> James Mayer de Rothschild. <b>Type:</b> Unpopular</p> <p><b>Sub-question 1:</b> Where was James Mayer de Rothschild born? <b>Sub-answer 1:</b> Frankfurt. <b>Type:</b> Unpopular</p> <p><b>Sub-question 2:</b> Which body of water is Frankfurt located next to? <b>Sub-answer 2:</b> Main. <b>Type:</b> Unpopular</p> <p><b>Sub-question 3:</b> Into which watercourse does the Main River flow? <b>Sub-answer 3:</b> Rhine. <b>Type:</b> Unpopular</p>

Table 30: Four-hop question-answer pairs and their corresponding types in MINTQA-POP.