

ReNAct: Augmenting LLM Reasoning with Dynamic Notes Writing for Complex QA

Anonymous ACL submission

Abstract

Retrieval-augmented generation (RAG) for complex multihop QA tasks struggle with long contexts — retrieved information accumulated across steps, leading to redundancy and inefficient queries at later reasoning steps that degrade retrieval and thus the overall response quality. We introduce ReNAct (Reason, iterative *NotesWriting*, and Action) which enhances long context based multi-hop QA reasoning by iteratively retrieving and accumulating only the most relevant information at each reasoning step. *NotesWriting* dynamically maintains a concise list of the relevant information which enables more effective query writing at each reasoning step allowing focus on missing information rather than reprocessing previously retrieved long content. By writing concise notes and guiding query formulation, ReNAct significantly improves both effectiveness and efficiency in multi-hop reasoning. Our approach achieves >20% absolute F1 score gains on long-context benchmarks such as FanOutQA and FRAMES, while reducing the number of reasoning steps by 56% on average compared to the ReAct + BM25 baseline.

1 Introduction

Large Language Models (LLMs) have significantly impacted many real-world applications including complex reasoning based tasks such as multi-hop question answering (QA) (Asai et al., 2023; Tang and Yang, 2024a; Lewis et al., 2020) by enabling frameworks to retrieve relevant documents, reason and generate more accurate and context-aware responses. Iterative RAG methods (Trivedi et al., 2022; Asai et al., 2023; Tang and Yang, 2024a) go further, leveraging the iterative reasoning and planning capabilities of LLMs through multiple reasoning methods, such as ReACT (Yao et al., 2023), Tree-of-Thought (Yao et al., 2024), and STaR (Zelikman et al., 2022), to achieve best-in-class performance on complex multi-hop QA tasks.

One issue with iterative retrieval over multiple steps is the increasingly long accumulated contexts (Kamradt, 2023; Hsieh et al., 2024; Liu et al., 2024a). As a result, query writing at a specific reasoning step requires comprehension of the long accumulated context until the previous step, which results in inefficient and redundant queries (shown in Table 1 and Table 3). Redundant queries lead to longer and ineffective reasoning traces, resulting in lower answer accuracy.

To this end, we propose Reason, iterative *NotesWriting*, and Action (ReNAct) (Figure 1) to address the challenges of stagnant information collection within reasoning steps due to growing context lengths. *NotesWriting* creates concise relevant information (referred to as *Notes*) from documents retrieved in each step. Notes are then used by LLMs for writing more effective queries at each reasoning step, optimizing redundant queries resulting from accumulated unfiltered, longer context. Our work makes the following key contributions:

- We showcase the need for *NotesWriting* - focusing on relevant information and targeted context at each reasoning step (addressing the “needle-in-a-haystack” (Kamradt, 2023) problem), leading to better query writing at each reasoning step, thus substantially improving QA performance. Our iterative approach ReNAct improves performance by over **23% and 20% absolute F1 score** compared to the baseline (ReAct + BM25) for FanoutQA and FRAMES, respectively.
- Our work emphasizes the benefits of concise context that iteratively enables query writing at each step to focus only on retrieving missing information, thus improving overall reasoning efficiency. Using concise context with *NotesWriting*, ReNAct requires **upto 55% fewer reasoning steps** to arrive at the correct answer.
- On mid-sized context based multi-hop QA datasets such as HotPotQA and MultiHop-RAG, ReNAct shows minor to measurable QA perfor-

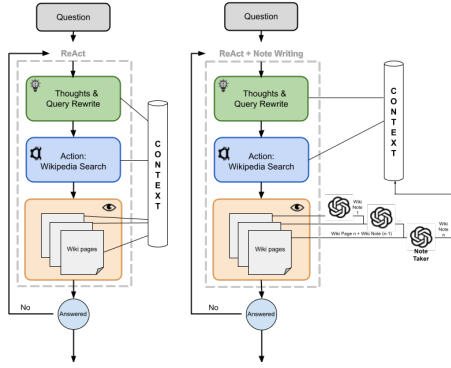


Figure 1: Comparison of 2 pipelines — Left: regular ReAct and Right: ReAct with iterative Note Writing (ReNAct).

mance improvements, but still substantial improvements in reasoning efficiency with 31% lesser reasoning steps.

2 Related Work

Iterative RAG enhances standard RAG by refining the retrieval process through multiple interactions with the knowledge base (Yao et al., 2023). RA-ISF (Liu et al., 2024b) and SELF-RAG (Asai et al., 2023) employ self-feedback, evaluating retrieved documents for relevance. PlanRAG (Lee et al., 2024) incorporates a planning step, explicitly outlining what information is needed before retrieval. FLARE (Jiang et al., 2023) and iMedRAG (Xiong et al., 2024) focus on active query generation, formulating new queries based on previous iterations. While all these algorithms aim to improve retrieval accuracy, they differ in their specific mechanisms, ranging from document evaluation and planning to query refinement and multi-draft verification. Others combine these ideas, Rewrite-Retrieve-Rerank (Cao et al., 2018) integrates self-feedback and active query generation, and IRCOT (Trivedi et al., 2022) interleaves Chain-of-Thought (Wei et al., 2022) updates with retrieval.

Agent-based methods are also employed for reasoning based multi-hop QA tasks. Infogent (Reddy et al., 2024) addresses complex QA tasks using an agent-based approach where its Navigator and Extractor agents handle web search and data formatting, respectively, while an Aggregation module synthesizes information and provides feedback. Due to certain framework similarities with ReNAct, we consider Infogent as one of our primary baselines. ReNAct outperforms all approaches, but more importantly, is significantly more efficient in reasoning with much fewer reasoning steps (Table 1).

3 Approach

Reason, iterative *NotesWriting*, and Action (ReNAct) operates in an iterative reasoning and retrieval paradigm to answer queries by generating structured intermediate steps. Specifically, given a query q , the LLM \mathcal{M} follows a cycle of reasoning, retrieval, and synthesis steps to arrive at a final answer. This process builds upon the ReAct (Yao et al., 2023) framework, with the addition of a *NotesWriting* module which retains the most relevant information from the documents retrieved with each chain, informing the upcoming chain of the focused information retrieved in the step before.

Algorithm 1 ReNAct: Iterative *NotesWriting*

Require: Query q , task prompt \mathcal{P} , summarize prompt \mathcal{P}_s , iteration limit $T = 25$, number of articles $k = 5$

Ensure: Final Answer

- 1: Initialize LLM \mathcal{M} with prompt \mathcal{P}
- 2: $O_0 \leftarrow \emptyset$ $\triangleright O_t$ aggregated observations
- 3: $t \leftarrow 0$
- 4: **while** $t < T$ and not Answer **do**
- 5: **Generate Thought & Action:**
- 6: (Thought, Action) $\leftarrow \mathcal{M}(q, O_t)$
- 7: Let Action = (e, s) $\triangleright e$: entity, s : search query
- 8: **Retrieve Information:**
- 9: Use e and s to fetch top k Wikipedia articles
- 10: $R \leftarrow \emptyset$ \triangleright Not used in non-iterative
- 11: **for** $i \leftarrow 1$ to k **do**
- 12: Convert article w_i to Markdown format
- 13: **Extract Content:**
- 14: $r_i \leftarrow \mathcal{M}_s(w_i, \mathcal{P}_s, O_t \cup R)$ \triangleright Summarize w_i
- 15: $R \leftarrow R \cup r_i$
- 16: **end for**
- 17: **Aggregate Observations:**
- 18: $O_{t+1} \leftarrow O_t \cup R$
- 19: $t \leftarrow t + 1$
- 20: **end while**
- 21: **return** Answer from $\mathcal{M}(q, O_T)$

ReNAct (Algorithm 1) follows these steps:

Step 1 – Initial Thought and Action Generation: \mathcal{M} is initialized with a few-shot prompt \mathcal{P} , which remains consistent across all tasks and models. Given the query q , the model first generates a reasoning step that guides the retrieval process as *Thought*. Then a search *Action* is initiated using a structured tuple (e, s) where e is the search entity, and s is the formulated search query.

Step 2 – Information Retrieval: Based on the entity e , we retrieve the top $k = 5$ relevant Wikipedia articles using the Wikipedia API. Each retrieved article is converted into Markdown format to ensure compatibility with LLM processing.

Step 3 – Writing Notes: For each Wikipedia article $w_{i \in \{1, \dots, k\}}$, an LLM \mathcal{M}_s ($= \mathcal{M}$ in our case) is employed to extract relevant information. This

process is guided by a specialized prompt \mathcal{P}_s (Appendix A.6), being informed by the current context, including the previous notes, query, and the retrieved document, produces a summarized note r_i for each article.

Step 4 – Observation Generation: The extracted content from all retrieved articles is aggregated as observation $O_t = \bigcup_{i=1}^k r_i$. O_t is provided as context to \mathcal{M} to generate the next reasoning step.

Iterative Reasoning and Termination: The process iterates as follows \rightarrow The model generates a new (*Thought*, *Action*) pair based on O_t . \rightarrow Steps 2 - 4 are repeated, updating O_t at each iteration. \rightarrow The cycle continues until the model determines it has sufficient information and executes a *finish action*, at which point it provides the final answer. To ensure computational efficiency, we impose a maximum iteration limit of $T = 25$ steps.

We consider two variations: simple *NotesWriting*, and iterative *NotesWriting* - called ReNAct. For *NotesWriting*, notes are written separately for each document; for ReNAct, each note builds on previous ones, reducing repetition and improving quality even further. However, ReNAct is sequential, while *NotesWriting* can run parallelly.

4 Experiments

Baselines: We compare ReNAct with: 1) **CoT**: relies on internal knowledge with chain-of-thought (closed-book), 2) **ReAct**: concatenates all retrieved documents into the LLM context, 3) **ReAct + BM25**: uses BM25 to assess the impact of focused retrieval by extracting the most relevant segments of the documents based on the search query, 4) **Info-gent** (Reddy et al., 2024): state-of-the-art for agentic retrieval and extraction method (Appendix A.1).

Benchmarks: We evaluated four multi-hop QA datasets: (1) FanOutQA (Zhu et al., 2024) with complex fanout questions, (2) FRAMES (Krishna et al., 2024) requiring 2–15 articles, (3) MultiHopRAG (Tang and Yang, 2024b) involving retrieval and reasoning over news articles, and (4) HotpotQA (Yang et al., 2018) requiring multi-article reasoning. See Appendix A.2 for details.

Models: We experiment with two models, representing closed & open weights, GPT-4o-mini¹ and Llama 3.1-70-Instruct (Dubey et al., 2024). We set the temperature to 0.7 and use the same LLM for for generating reasoning step and

NotesWriting (i.e $\mathcal{M}_s = \mathcal{M}$).

Evaluation Metrics: We report the F1 score between predicted and ground truth answer and the GPT4-as-Judge score (prompt in the appendix 2). Additionally, to measure the efficiency of reasoning steps we report the average number of tokens and reasoning steps.

5 Results

(1) Performance improvements and efficient reasoning - ReNAct outperforms all baselines in all of our evaluations shown in Tables 1 and 2. Especially in long context multi-hop QA dataset (Table 1), ReNAct achieves absolute 23% and 16.2% F1 score improvements over ReAct highlighting importance of concise context provided by *NotesWriting*. ReNAct also has the lowest number of *Avg Steps* highlighting the shortest and most efficient reasoning trace due to effective query writing enabled from concise context (shown by the lowest number of *Avg Tokens* in Tables 1 and 2).

(2) Comparison to baselines – ReNAct substantially outperforms both ReAct + BM25 and Info-gent baselines that have similar iterative reasoning based setup, with the exception of FanOutQA on F1 score. Importantly, top performing agentic framework Info-gent has the highest number of reasoning steps whereas ReNAct has the lowest - highlighting the need of **iterative** *NotesWriting* for complex multi-hop QA. Furthermore, CoT baseline which does not retrieve any context has the lowest number of *Avg Tokens*, but achieves much lower performance compared to our ReNAct method.

(3) NotesWriting vs. ReNAct – ReNAct outperforms *NotesWriting* in the majority of settings in Tables 1 and 2 with nearly the same average number of reasoning steps, highlighting importance of writing notes iteratively. At the same time, *NotesWriting* can be applied on multiple documents in parallel for faster notes writing, and still achieves close performance to ReNAct.

(4) Correlation with context length – ReNAct performs better in long context multi-hop QA datasets (table 1) where *NotesWriting* can be more useful. We also observe that longer context (*Avg Tokens* column) always leads to lower performance and more reasoning steps, emphasizing its impact on query writing and the subsequent retrieval process. On mid-sized context based QA datasets such as HotPotQA and Multi-hopRAG where *Avg Tokens* (Table 2) are relatively

¹<https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

smaller, the performance improvements are not as pronounced as with long context multi-hop QA datasets (Table 1). Even so, ReNAct uses (upto 31.2%) fewer reasoning steps even in mid-size contexts datasets.

Overall, the significant performance gains, coupled with a reduction in total tokens and average reasoning steps, clearly highlight the importance of filtering retrieved context by extracting and presenting only the most relevant information to the reasoning LLM at each step. A few examples of our ReNAct are shown in Appendix A.6, highlighting benefits like 1) better informing the overall reasoning process, for e.g. by better handling of conflicting information (see Table 10 vs Table 8), and 2) informing the next action taken by informing the search query (see Table 9 vs Table 8). We leave further in-depth explorations of these aspects for future work.

M	DS	Method	F1 Score ↑	GPT4 Score ↑	Avg Tokens ↓	Avg Steps ↓
Llama-3.1-70B-Inst	FanoutQA	CoT	38.1	15.2	263.5	1.0
		ReAct	21.5	11.3	9215.8	19.1
		ReAct + BM25	18.0	9.7	10514.5	16.1
		Infogent	47.2	22.9	3938.9	20.0
		<i>NotesWriting</i>	43.0	27.1	1461.7 (-86%)	8.1 (-50%)
		ReNAct	41.4	24.8	1647.4 (-84%)	8.6 (-47%)
	FRAMES	CoT	28.2	32.1	365.1	1.0
		ReAct	34.3	39.7	5770.1	15.9
		ReAct + BM25	33.4	39.3	8180.5	12.1
		Infogent	28.0	29.9	4272.4	19.7
		<i>NotesWriting</i>	49.1	57.6	939.5 (-89%)	7.1 (-41%)
		ReNAct	53.7	65.8	972.3 (-88%)	7.2 (-40%)
GPT-4o-mini	FanoutQA	CoT	42.2	13.9	261.3	1.0
		ReAct	28.6	12.9	8387.8	18.1
		ReAct + BM25	18.8	8.1	13010.8	18.1
		Infogent	47.2	22.9	5255.9	19.5
		<i>NotesWriting</i>	48.0	24.5	1350.4 (-90%)	10.2 (-44%)
		ReNAct	51.4	29.0	1329.2 (-90%)	10.0 (-45%)
	FRAMES	CoT	29.1	28.8	133.5	1.0
		ReAct	28.7	31.1	9498.0	12.4
		ReAct + BM25	28.6	30.0	9599.6	12.2
		Infogent	28.0	29.9	4700.1	18.8
		<i>NotesWriting</i>	44.7	49.2	1190.5 (-88%)	8.5 (-31%)
		ReNAct	44.9	51.4	1237.8 (-87%)	8.6 (-29%)

Table 1: Results on FanoutQA & FRAMES. *NotesWriting* is non-iterative Reason + *NotesWriting* + Action, ReNAct refers to iterative *NotesWriting*. % improvements are in comparison to ReAct + BM25. M = Model, DS = Dataset. CoT results are presented as crude baseline without retrieval, comparisons are made with reasoning frameworks only.

5.1 Redundancy analysis

As shown in Tables 1 and 2, ReAct + BM25 baseline typically has more number of tokens (*Avg Tokens* column in result tables) in context along with higher reasoning steps. To highlight the impact of long context in query writing at different steps, we compute the average number of repeated (or retried) queries. As shown in Table 3, ReAct + BM25 has substantially more repeated queries, highlighting

M	DS	Method	F1 Score ↑	GPT4 Score ↑	Avg Tokens ↓	Avg Steps ↓
Llama-3.1-70B-Inst	MultiHop	ReAct	50.5	61.4	31460.6	5.0
		ReAct + BM25	63.9	70.8	1732.6	6.2
		<i>NotesWriting</i>	64.5	73.4	803.3 (-53.6%)	5.6 (-8.9%)
		ReNAct	63.2	72.6	861.1 (-50.3%)	5.9 (-4.5%)
	HotpotQA	ReAct	56.3	66.6	2858.3	9.5
		ReAct + BM25	50	63	3976.0	7.8
		<i>NotesWriting</i>	56.4	68.2	696.6 (-82.5%)	5.6 (-27.8%)
		ReNAct	57.8	70.0	632.4 (-84.1%)	5.3 (-31.2%)
GPT-4o-mini	MultiHop	ReAct	57.0	70	37614.7	5.6
		ReAct + BM25	56.2	68.4	2703.1	8.5
		<i>NotesWriting</i>	58.0	70.8	711.4 (-73.7%)	6.2 (-26.7%)
		ReNAct	58.0	71.2	693.9 (-74.3%)	6.2 (-27.0%)
	HotpotQA	ReAct	48.4	62.4	2546.5	7.1
		ReAct + BM25	44.6	56	4461.2	7.2
		<i>NotesWriting</i>	54.0	67.0	620.1 (-86.1%)	5.7 (-20.9%)
		ReNAct	53.7	64.2	614.1 (-86.2%)	5.6 (-22.2%)

Table 2: Evaluations on complex multi-hop QA datasets - MultiHop-RAG and HotpotQA. Notations same as Table 1.

the impact of long context on query writing. On the other hand, ReNAct has just 2 or 3 repeated queries on average, suggesting the effective quality of queries in ReNAct’s reasoning process. We leave enhancement of ReNAct to further reduce the number of repeated queries to future work.

Model	Dataset	ReAct + BM25	ReNAct
LLama-3.1-70B-Inst	FanoutQA	7.43	2.01
	FRAMES	5.12	2.20
GPT-4o-mini	FanoutQA	8.03	2.69
	FRAMES	5.06	3.23

Table 3: No of search retries are lower for ReNAct compared to ReAct + BM25 for FanoutQA and FRAMES.

6 Conclusion

We presented ReNAct, a *NotesWriting* technique that accumulates only the most relevant information at each reasoning step iteratively. We show that long context leads to redundant query generation at various reasoning steps, which leads to inefficient reasoning and lower question answering performance. Our approach ReNAct tackles these problems by retaining relevant information concisely, reducing redundancies in contexts, and thus focusing the context on salient information required for the QA task. ReNAct achieves the best performance on long-context multi-hop QA benchmarks FanoutQA and FRAMES, clearly showcasing the effectiveness of ReNAct. Our approach ReNAct also uses the least number of reasoning steps on average highlighting the effectiveness of providing concise context for better query writing at each reasoning step.

Limitations and Societal Impact

Our approach has several limitations. First, our experiments are limited to the two models we experiment with, which could be extended to newer smaller open-source models. Second, we limit on-line searches to the Wikipedia API², which only supports searching for text matching Wiki pages; and third, Wiki pages change often and this could lead to a mismatch with static benchmarks’ ground truth. While these could affect performance, we ensure that the same setup is also followed in all baselines we experiment with, to keep evaluation comparable while reducing the need to utilize paid search APIs. Third, with retrievals based on iterative notes writing, there is a possibility of conflicting information being received (Table 10). In observed examples, the ReNAct based setup is able to resolve discrepancies in the retrieved information, showing that our approach performs well. However, it is possible that the model starts hallucinating facts, and this remains a weakness at large. Lastly, we impose a maximum iteration limit to ensure computational efficiency, which could also impact performance. Further explorations towards improving on weaknesses remain future work.

Potential risks of our work include usage in scenarios where the requested retrieval information is toxic or harmful. While we cannot control how our method is used for prompting, we expect content moderation policies to help with reducing the impact of such queries. Moreover, hallucinations can affect the QA experience, although manual observation of the reasoning traces show that recovery can be better with ReNAct.

We expect our work to significantly enhance the QA user experience, as focused information improves performance and reduced context lengths lower computational costs. We hope our ReNAct method can contribute towards better task handling at large. We will make our code publicly available upon acceptance towards this goal.

References

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*.
- Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. 2018. Retrieve, rerank and rewrite: Soft template based

neural summarization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 152–161, Melbourne, Australia. Association for Computational Linguistics.

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

Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. Ruler: What’s the real context size of your long-context language models? *arXiv preprint arXiv:2404.06654*.

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

Greg Kamradt. 2023. Needle in a haystack-pressure testing llms. *Github Repository*, page 28.

Satyapriya Krishna, Kalpesh Krishna, Anhad Mohananey, Steven Schwarcz, Adam Stambler, Shyam Upadhyay, and Manaal Faruqui. 2024. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation. *arXiv preprint arXiv:2409.12941*.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626.

Myeonghwa Lee, Seonho An, and Min-Soo Kim. 2024. Planrag: A plan-then-retrieval augmented generation for generative large language models as decision makers. *arXiv preprint arXiv:2406.12430*.

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

Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024a. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.

Yanming Liu, Xinyue Peng, Xuhong Zhang, Weihao Liu, Jianwei Yin, Jiannan Cao, and Tianyu Du. 2024b. RA-ISF: Learning to answer and understand from retrieval augmentation via iterative self-feedback. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4730–4749, Bangkok, Thailand. Association for Computational Linguistics.

²<https://www.mediawiki.org/wiki/API:Search>

Revanth Gangi Reddy, Sagnik Mukherjee, Jeonghwan Kim, Zhenhailong Wang, Dilek Hakkani-Tur, and Heng Ji. 2024. Infogent: An agent-based framework for web information aggregation. *arXiv preprint arXiv:2410.19054*.

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

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

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

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

Guangzhi Xiong, Qiao Jin, Xiao Wang, Minjia Zhang, Zhiyong Lu, and Aidong Zhang. 2024. Improving retrieval-augmented generation in medicine with iterative follow-up questions. In *Biocomputing 2025: Proceedings of the Pacific Symposium*, pages 199–214. World Scientific.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. STaR: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488.

Andrew Zhu, Alyssa Hwang, Liam Dugan, and Chris Callison-Burch. 2024. [FanOutQA: A multi-hop, multi-document question answering benchmark for large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 18–37, Bangkok, Thailand. Association for Computational Linguistics.

A Appendix

A.1 Infogent Implementation Details

We use the official code provided by Infogent authors [here](#) (Apache 2.0. License) with the following modifications:

- Due to our limitations in accessing OpenAI, we modify the code to use AzureOpenAI endpoints, and use ChatOpenAI and OpenAI for Llama-70B experiments with vLLM ([Kwon et al., 2023](#)) endpoints³.
- OpenAI embedding is replaced by sentence-transformers’ all-mpnet-base-v2⁴.
- Serper Google Search⁵ is replaced by Wikipedia search API due to credit limitations and to use similar open knowledge tools as those used in our method, reducing the cost needed to conduct RAG experiments.

A.2 Benchmarks

We evaluated four multi-hop QA datasets: (1) FanOutQA ([Zhu et al., 2024](#)), which features complex fanout questions, (2) FRAMES ([Krishna et al., 2024](#)), requiring reasoning over 2–15 articles, (3) MultiHop-RAG ([Tang and Yang, 2024b](#)), which involves retrieval and reasoning over news articles, and (4) HotpotQA ([Yang et al., 2018](#)), which requires multi-article reasoning. For FanOutQA, we evaluated all 310 examples from the development set, while for FRAMES, we used 549 multiple-constraint-tagged questions. For MultiHop-RAG and HotpotQA, we assessed performance on 500 examples from the test and development splits, respectively. FanOutQA, HotpotQA and Wikipedia comes under CC BY-SA 4.0 (Creative Commons Attribution-ShareAlike 4.0 International License), FRAMES under Apache 2.0. license and MultiHop-RAG under ODC-By (Open Data Commons Attribution License).

A.3 Models

Models: We conduct experiments with two models, representing both closed and open weights: GPT-4o-mini⁶ and Llama 3.1-70-Instruct ([Dubey et al., 2024](#)). The temperature is set to 0.7, and the same LLM is

³<https://docs.vllm.ai/en/stable/index.html>

⁴<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

⁵<https://serper.dev/>

⁶<https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

used for generating reasoning steps and *NotesWriting* (i.e., $\mathcal{M}_s = \mathcal{M}$). Llama 3.1-70-Instruct was hosted using vLLM (Kwon et al., 2023) across 8 A100-80GB GPUs, supporting a maximum context length of 64K. With parallelization, evaluation runs took approximately 9–10 hours for MultiHop-RAG, HotpotQA, and FRAMES, and around 15 hours for FanOutQA. GPT-4o-mini, which has a context length of 128K, completed evaluations in approximately 7 hours for FRAMES and FanOutQA, 2 hours for HotpotQA, and 27 minutes for MultiHop-RAG. The reported times include the full end-to-end process, accounting for rate limits, Wikipedia queries, and *NotesWriting*.

A.4 EM and Avg Obs Tokens comparison across benchmark datasets

We report the EM scores and average observations tokens for comparison across all our benchmark dataset in Table 4.

Model	Dataset	Method	EM Score \uparrow	Avg Obs Tokens \downarrow
Llama 3.1 70B Inst	FanoutQA	ReAct	2.3	8531.5
		ReAct + BM25	1.3	9496.8
		<i>NotesWriting</i>	39.0	1013.1 (-89%)
		ReNAct	35.0	1144.5 (-88%)
	FRAMES	ReAct	24.6	5150.2
		ReAct + BM25	21.9	7471.2
		<i>NotesWriting</i>	33.3	566.8 (-92%)
		ReNAct	38.6	559.4 (-93%)
	MultiHop RAG	ReAct	49.2	30655.4
		ReAct + BM25	63.2	1410.6
		<i>NotesWriting</i>	63.6	449.7 (-68.1%)
		ReNAct	61.8	458.6 (-67.5%)
	HotpotQA	ReAct	43.4	2524.5
		ReAct + BM25	37.8	3549.2
		<i>NotesWriting</i>	41.6	420.2 (-88.2%)
		ReNAct	43.6	369.9 (-89.6%)
GPT4o-mini	FanoutQA	ReAct	1.9	7,305.5
		ReAct + BM25	1.0	1,1746.9
		<i>NotesWriting</i>	3.9	648.1 (-94%)
		ReNAct	6.5	640.7 (-95%)
	FRAMES	ReAct	16.9	8604.5
		ReAct + BM25	18.2	8,751.7
		<i>NotesWriting</i>	30.2	603.4 (-93%)
		ReNAct	30.2	643.5 (-93%)
	MultiHop RAG	ReAct	56.0	36843.1
		ReAct + BM25	55.2	2203.6
		<i>NotesWriting</i>	56.6	367.6 (-83.3%)
		ReNAct	56.8	353.1 (-84.0%)
	HotpotQA	ReAct	33.4	2170.8
		ReAct + BM25	31.4	4010.4
		<i>NotesWriting</i>	38.0	277.4 (-93.1%)
		ReNAct	39.8	283.0 (-92.9%)

Table 4: EM and Avg Obs Tokens for all benchmarks. *NotesWriting* refers to non-iterative Reason + *NotesWriting* + Action approach, and ReNAct refers to the iterative *NotesWriting* approach. % improvements are in comparison to the ReAct + BM25 baseline.

A.5 Standard deviation across runs

We ran the *NotesWriting* and ReNAct across all datasets and models three times to see the variance

across different runs. We report the results in Table 5.

A.6 Examples comparing ReNAct with baselines

Model	Dataset	Avg F1	GPT-4 Score	Avg Tokens	Avg Steps
GPT-4o-mini	Fanout	± 1.86	± 2.45	± 28.30	± 0.13
	FRAMES	± 1.10	± 2.35	± 30.07	± 0.13
Llama-3.1 70B	Fanout	± 3.79	± 1.54	± 98.15	± 0.47
	FRAMES	± 4.42	± 5.76	± 37.89	± 0.61

Table 5: Standard deviation across FRAMES & FanoutQA for each metric per model.

<p>Question: What law schools did the most recent four Supreme Court justices attend?</p> <p>Thought: I need to gather information on the law schools attended by the most recent four Supreme Court justices. First, I will search for the names of the most recent four justices to find out where they studied law.</p> <p>Action: search[Supreme Court justices; Who are the most recent four Supreme Court justices?]</p> <p>Observation: (Result 1) List of justices of the Supreme Court of the United States - The most recent four Supreme Court justices are: Amy Coney Barrett, Ketanji Brown Jackson, Neil Gorsuch, and Brett Kavanaugh.</p> <p>(Result 2) Ideological leanings of United States Supreme Court justices - The most recent four Supreme Court justices are: Amy Coney Barrett (appointed in 2020), Ketanji Brown Jackson (appointed in 2022), Brett Kavanaugh (appointed in 2018), and Neil Gorsuch (appointed in 2017).</p> <p>(Result 3) Demographics of the Supreme Court of the United States - The most recent four Supreme Court justices are: Ketanji Brown Jackson (appointed in 2022), Amy Coney Barrett (appointed in 2020), Brett Kavanaugh (appointed in 2018), and Neil Gorsuch (appointed in 2017).</p>
<p>Thought: The most recent four Supreme Court justices are Amy Coney Barrett, Ketanji Brown Jackson, Neil Gorsuch, and Brett Kavanaugh. Next, I need to find out which law schools they attended. I will start by searching for the law school attended by Amy Coney Barrett.</p> <p>Action: search[Amy Coney Barrett; Which law school did Amy Coney Barrett attend?]</p> <p>Observation: (Result 1) Amy Coney Barrett attended Rhodes College (BA) and University of Notre Dame (JD). She graduated from Notre Dame Law School in 1997 with a Juris Doctor, <i>summa cum laude</i>, ranked first in her class.</p>
<p>Thought: Amy Coney Barrett attended Notre Dame Law School. Next, I need to find out which law school Ketanji Brown Jackson attended.</p> <p>Action: search[Ketanji Brown Jackson; Which law school did Ketanji Brown Jackson attend?]</p> <p>Observation: (Result 1) Ketanji Brown Jackson attended Harvard Law School, where she graduated in 1996 with a Juris Doctor, <i>cum laude</i>.</p> <p>(Result 2) Joe Biden Supreme Court candidates - Ketanji Brown Jackson attended Harvard University for her law degree.</p> <p>(Result 3) Demographics of the Supreme Court of the United States - Ketanji Brown Jackson attended Harvard Law School.</p>
<p>Thought: Ketanji Brown Jackson attended Harvard Law School. Now, I need to find out which law school Neil Gorsuch attended.</p> <p>Action: search[Neil Gorsuch; Which law school did Neil Gorsuch attend?]</p> <p>Observation: (Result 1) Neil Gorsuch received his legal education at Harvard Law School and graduated in 1991 with a Juris Doctor, <i>cum laude</i>.</p>
<p>Thought: Neil Gorsuch also attended Harvard Law School. Finally, I need to find out which law school Brett Kavanaugh attended.</p> <p>Action: search[Brett Kavanaugh; Which law school did Brett Kavanaugh attend?]</p> <p>Observation: (Result 1) Brett Kavanaugh - Kavanaugh graduated from Yale Law School with a Juris Doctor degree in 1990.</p>
<p>Thought: Brett Kavanaugh graduated from Yale Law School. Now I have gathered the law schools attended by the most recent four Supreme Court justices: Amy Coney Barrett attended Notre Dame Law School, Ketanji Brown Jackson and Neil Gorsuch both attended Harvard Law School, and Brett Kavanaugh attended Yale Law School.</p> <p>Action: finish[Notre Dame Law School, Harvard Law School, Harvard Law School, Yale Law School]</p>

Table 6: An Example of ReNAct on the FanOUT dataset.

<p>Question: What law schools did the most recent four Supreme Court justices attend?</p>
<p>Thought: I need to find out which law schools the most recent four Supreme Court justices attended. First, I'll identify the current Supreme Court justices and then gather information about their law school backgrounds.</p> <p>Action: search[Supreme Court justices; Who are the current Supreme Court justices?]</p> <p>Observation: (Result 1) The U.S. Supreme Court Building, current home of the Supreme Court, which opened in 1935..... (Result 2) The graph below, based on data from the Supreme Court Database, shows this partisan split.[6] In all the non-unanimous decisions.... (Result 3) This graphical timeline depicts the length of each current Supreme Court justice's tenure (not seniority, as the chief justice has.... (Result 4) There are currently nine justices on the Supreme Court: Chief Justice John Roberts and eight associate justices. Among the current.... (Result 5) From this data, scholars have inferred the ideological leanings of each justice and how the justices are likely to vote.....</p>
<p>Thought: The search results did not provide specific names of the current Supreme Court justices. I need to refine my search to get the names of the most recent four justices so I can find out which law schools they attended.</p> <p>Action: search[most recent Supreme Court justices; What are the names of the most recent four Supreme Court justices?]</p> <p>Observation: (Result 1) As of March 20, 2024[update], there are fourteen associate justices on the Supreme Court, with the most recent appointment being.... (Result 2) Once the committee reports out the nomination, it is put before the full Senate for final consideration. A simple majority..... (Result 3) Potter Stewart, the most recent Supreme Court justice initially appointed through a recess appointment..... (Result 4) William Rehnquist, the most recent incumbent associate justice to assume the position of chief justice..... (Result 5) Until 1973, only men were appointed as Associate Justices to the Court. Cecilia Muñoz-Palma, an appointee of President Ferdinand Marcos....</p>
<p>Thought: The search results did not provide the names of the most recent four U.S. Supreme Court justices. I need to further refine my search to get the specific names and then gather their law school information.</p> <p>Action: search[U.S. Supreme Court justices; Who are the current justices and when were they appointed?]</p> <p>Observation: (Result 1) The graph below, based on data from the Supreme Court Database, shows this partisan split.[6] In all the non-unanimous decisions.... (Result 2) Five of the current justices are assigned to circuits on which they previously sat as circuit judges: Chief Justice Roberts.... (Result 3) From this data, scholars have inferred the ideological leanings of each justice and how the justices are likely to vote..... (Result 4) These two graphs differ because of the choices of data sources, data coverage, coding of complicated cases, smoothing parameters, and.... (Result 5) Under Article Three of the United States Constitution, the composition and procedures of the Supreme Court were originally established by....</p>
<p>Thought: The search results are still not yielding the names of the most recent Supreme Court justices. I will modify my search approach to focus specifically on the recent appointments to the Supreme Court.</p> <p>Action: search[U.S. Supreme Court justices names and law schools of the last four justices]</p> <p>Observation: (Result 1) Associate justices James F. Byrnes, whose short tenure lasted from June 1941 to October 1942, and Robert H. Jackson, who.... (Result 2) The Court returned to the issue of abortion in the 2000 case of <i>Stenberg v. Carhart</i>*, in which it invalidated.... (Result 3) In modern discourse, the justices of the Court are often categorized as having conservative, moderate, or liberal philosophies of law.... (Result 4) <i>* Ideological leanings of United States Supreme Court justices *</i> List of U.S. Supreme Court justices who also served in.... (Result 5) The nine Supreme Court justices base their decisions on their interpretation of both legal doctrine and the precedential application of....</p>
<p>Thought: I have not yet obtained the specific names of the most recent four U.S. Supreme Court justices or their law schools. I will make one final attempt to gather this information.</p> <p>Action: search[most recent four U.S. Supreme Court justices and their law schools]</p> <p>Observation: (Result 1) Article III, Section I of the Constitution of the United States (1776) establishes that the justices hold their office during..... (Result 2) Law clerks have assisted the justices of the United States Supreme Court in various capacities since the first one was..... (Result 3) In modern discourse, the justices of the Court are often categorized as having conservative, moderate, or liberal philosophies of law..... (Result 4) At least six justices are Roman Catholics, one is Jewish, and one is Protestant. It is unclear whether Neil Gorsuch..... (Result 5) Each justice is permitted to have three or four law clerks per Court term. Most clerks are recent law school.....</p>
<p>Thought: The searches have consistently failed to yield the specific names of the most recent four U.S. Supreme Court justices or their law schools. I will conclude my task as I have exhausted the options for retrieving this information.</p> <p>Action: finish[Unable to retrieve names and law schools of the most recent four U.S. Supreme Court justices.]</p>

Table 7: An Example of ReACT with BM25 on the Fanout dataset. The agent keeps on searching but is unable to extract relevant information.

<p>Question: Diago Costa played for which club when he was awarded the first FIFA World Cup Goal based on a VAR Decision?</p> <p>Thought: I need to determine which club Diago Costa was playing for when he was awarded the first FIFA World Cup Goal based on a VAR Decision. First, I will gather information on the first FIFA World Cup Goal awarded based on a VAR Decision.</p> <p>Action: search[FIFA World Cup Goal; What was the first FIFA World Cup Goal awarded based on a VAR Decision, and who scored it?]</p> <p>Observation: (Result 1) 2022 FIFA World Cup - Ao Tanaka</p>
<p>Thought: It seems like the search result did not provide the relevant information about the first FIFA World Cup Goal awarded based on a VAR Decision.</p> <p>Action: search[FIFA World Cup VAR Decision; What was the first FIFA World Cup Goal awarded based on a VAR Decision, and who scored it?]</p> <p>Observation: (Result 1) 2018 FIFA World Cup - Diego Costa's first goal against Portugal became the first World Cup goal based on a VAR decision.</p> <p>(Result 2) 2022 FIFA World Cup - The first FIFA World Cup Goal awarded based on a VAR Decision was scored by Ao Tanaka of Japan against Spain.</p> <p>(Result 3) 2019 FIFA Women's World Cup - The first FIFA World Cup goal awarded based on a VAR decision was scored by French midfielder Antoine Griezmann in the 58th minute of a group stage match against Australia on June 16, 2018.</p> <p>(Result 4) Video assistant referee - The first VAR decision at the World Cup came on 16 June 2018 in a group stage match between France and Australia, where referee Andres Cunha awarded a penalty to France after consulting with the VAR. The resulting penalty kick was successfully converted by Antoine Griezmann of France, but the goal was initially disallowed for offside. However, the VAR team correctly advised the referee that the offside infraction was committed by an Australian player, not a French player, and the goal was awarded.</p>
<p>Thought: It seems like the search results provided conflicting information about the first FIFA World Cup Goal awarded based on a VAR Decision. However, according to Result 1 of the second search, Diego Costa's first goal against Portugal became the first World Cup goal based on a VAR decision.</p> <p>Action: search[Diego Costa; Which club did Diego Costa play for in the 2018 FIFA World Cup?]</p> <p>Observation: (Result 1) Diego Costa - Atlético Madrid</p>
<p>Thought: Diego Costa played for Atlético Madrid when he was awarded the first FIFA World Cup Goal based on a VAR Decision.</p> <p>Action: finish[Atlético Madrid]</p> <p>Observation: Atlético Madrid</p>

Table 8: An Example of ReNAct from the FRAMES dataset.

GPT-4 Judge Prompt
<p>===Task===</p> <p>I need your help in evaluating an answer provided by an LLM against a ground truth answer. Your task is to determine if the ground truth answer is present in the LLM's response. Please analyze the provided data and make a decision.</p> <p>===Instructions===</p> <ol style="list-style-type: none"> 1. Carefully compare the "Predicted Answer" with the "Ground Truth Answer". 2. Consider the substance of the answers – look for equivalent information or correct answers. Do not focus on exact wording unless the exact wording is crucial to the meaning. 3. Your final decision should be based on whether the meaning and the vital facts of the "Ground Truth Answer" are present in the "Predicted Answer." <p>===Input Data===</p> <ul style="list-style-type: none"> • Question: «question» • Predicted Answer: «LLM_response» • Ground Truth Answer: «ground_truth_answer» <p>===Output Format===</p> <p>Provide your final evaluation in the following format:</p> <p>Explanation: (How you made the decision?)</p> <p>Decision: ("TRUE" or "FALSE")</p> <p>Please proceed with the evaluation.</p>

Figure 2: GPT-4 prompt for evaluating the correctness of predicted answer.

Question: Diego Costa played for which club when he was awarded the first FIFA World Cup Goal based on a VAR Decision?
Think "I need to gather information about Diago Costa and the FIFA World Cup Goal awarded based on a VAR decision." Search "Diego Costa club history" (x3)
Think "I need to find information about Diego Costa and the FIFA World Cup goal awarded based on a VAR decision." Search "Diego Costa biography" (x1)
Think "The first search result is about Diego Costa, which is relevant to my query about his club and the FIFA World Cup Goal based on a VAR decision. I will extract information from the Diego Costa Wikipedia page." Extract https://en.wikipedia.org/wiki/Diego_Costa (x2)
Think "I need to perform a search specifically targeting Diego Costa's participation in the FIFA World Cup where he was awarded a goal based on a VAR decision, including the club he was with at that time." Search "Diego Costa FIFA World Cup goal VAR decision club" (x2)
Think "I need to find out which club Diego Costa was playing for when he was awarded the first FIFA World Cup Goal based on a VAR decision." Search "Diego Costa club history" (x3)
Think "The initial search did not yield relevant information regarding Diego Costa and the FIFA World Cup Goal awarded based on a VAR decision." Search "Diego Costa club history FIFA World Cup VAR decision goal" (x1)
Think "I need to gather information about Diago Costa and the FIFA World Cup." Search "Diago Costa career history" (x1)
Think "I need to find information about Diego Costa and the FIFA World Cup Goal based on a VAR decision." Search "Diego Costa FIFA World Cup Goal VAR Decision" (x1)
Think "The search results include relevant information about the FIFA World Cup, but I need specifics about Diego Costa and any goals awarded based on VAR decisions." Search "Diego Costa World Cup goals VAR decisions" (x1)
Think "I need to find out which club Diego Costa was playing for when he was awarded the first FIFA World Cup Goal based on a VAR Decision." Search "Diego Costa World Cup appearances" (x1)
Think "I found relevant information about Diego Costa, but I need to check his details to find out when he was awarded the first FIFA World Cup goal based on a VAR decision." Extract https://en.wikipedia.org/wiki/Diego_Costa (x2)
Think "I need to search for specific information regarding the year Diego Costa was awarded a FIFA World Cup goal based on a VAR decision, as well as the club he was playing for at that time." Search "Diego Costa FIFA World Cup goal VAR decision year" (x1)
Think "The next step is to find information specifically about Diego Costa and the FIFA World Cup goal awarded based on a VAR decision." Search "Diego Costa FIFA World Cup goal VAR decision" (x2)
Think "The search results provided general information about the FIFA World Cup events, but I need to focus on Diego Costa's specific performance and goals in the World Cup." Search "Diego Costa World Cup goals VAR decision club" (x2)
Think "I need to find information about Diego Costa and the FIFA World Cup Goal based on a VAR decision." Search "Diego Costa career achievements FIFA World Cup" (x1)
Think "I need to find out which club Diego Costa was playing for when he was awarded the first FIFA World Cup Goal based on a VAR Decision." Search "Diego Costa career history" (x1)

Table 9: Example from FRAMES and Infogent: the agent searches for repeated queries without any modifications.

<p>Question: As of August 1, 2024, what is the largest city of the 9th largest country by land area in Europe?</p>
<p>Thought: I need to determine the 9th largest country by land area in Europe and then find its largest city. First, I will identify the 9th largest country by land area in Europe.</p> <p>Action: search[Largest countries in Europe by land area; Which is the 9th largest country by land area in Europe?]</p> <p>Observation: (Result 1) List of countries and dependencies by area - The 9th largest country by land area in Europe is Kazakhstan (2,699,700 km²).</p> <p>(Result 2) List of European countries by area - Poland</p>
<p>Thought: It seems like there are multiple conflicting sources for the 9th largest country by land area in Europe, with one source saying it's Kazakhstan, but Kazakhstan is not in Europe, and another source saying it's Poland but without providing more information. In this case, I'll try to find the correct 9th largest country by land area in Europe.</p> <p>Action: search[Largest countries in Europe by land area; What are the top 10 largest countries by land area in Europe?]</p> <p>Observation: (Result 1) List of countries and dependencies by area -</p> <ol style="list-style-type: none"> 1. Russia - 16,376,870 km² 2. Ukraine - 579,330 km² 3. France - 640,427 km² 4. Spain - 498,980 km² 5. Sweden - 407,284 km² 6. Norway - 366,704 km² 7. Germany - 349,390 km² 8. Finland - 303,815 km² 9. Poland - 304,255 km² 10. Italy - 295,717 km²
<p>Thought: The 9th largest country by land area in Europe is indeed Poland, with a land area of 312,679 km² or 304,255 km². Now, I need to determine the largest city of Poland.</p> <p>Action: search[Poland; What is the largest city in Poland as of August 1, 2024?]</p> <p>Observation: (Result 1) Poland - Warsaw</p> <p>(Result 2) Invasion of Poland - The largest city in Poland as of August 1, 2024, is Warsaw (Polish: Warszawa), which is also the capital of Poland.</p>
<p>Thought: The largest city of the 9th largest country by land area in Europe as of August 1, 2024, is Warsaw.</p> <p>Action: finish[The largest city of the 9th largest country by land area in Europe as of August 1, 2024, is Warsaw.]</p>

Table 10: An example from the FRAMES dataset where note taking leads to conflicting information (in bold). To clarify the confusion the agent re-tries the search to find the correct information.

Prompt for Notes Writing
<p>Extract relevant information which is not previously extracted from the Wikipedia page provided in markdown format relevant to the given query. You will be provided with the Wikipedia page, query, and the previously extracted content. Do not miss any information. Do not add irrelevant information or anything outside of the provided sources. Provide the answer in the format: <YES/NO>#<Relevant context>.</p> <p>Here are the rules:</p> <ul style="list-style-type: none"> • If you don't know how to answer the query - start your answer with NO# • If the text is not related to the query - start your answer with NO# • If the content is already extracted - start your answer with NO# • If you can extract relevant information - start your answer with YES# <p>Example answers:</p> <ul style="list-style-type: none"> • YES#Western philosophy originated in Ancient Greece in the 6th century BCE with the pre-Socratics. • NO#No relevant context. <p>Context: {Context}</p> <p>Previous Context: {PrevContext}</p> <p>Query: {Query}</p>

Figure 3: Notes writing prompt for extracting the relevant information.