

Augmenting LLM Reasoning with Dynamic Notes Writing for Complex MultiHop QA

Rishabh Maheshwary[§] Masoud Hashemi[§] Khyati Mahajan[§]
Shiva Krishna Reddy Malay[§] Sai Rajeswar[§] Sathwik Tejaswi Madhusudhan[§]
Spandana Gella[†] Vikas Yadav[§]
ServiceNow[§], ServiceNow Research[†]
rishabh.maheshwary@servicenow.com, vikas.yadav@servicenow.com

Abstract

Iterative RAG for multi-hop question answering faces challenges with lengthy contexts and the buildup of irrelevant information. This hinders a model's capacity to process and reason over retrieved content and limits performance. While recent methods focus on compressing retrieved information, they are either restricted to single-round RAG, require finetuning or lack scalability in iterative RAG. To address these, we propose *NotesWriting*, a method that generates concise and relevant notes from retrieved documents at each step, thereby reducing noise and retaining only essential information. This increases the *effective context length* of Large Language Models (LLMs), to reason and plan more effectively while processing larger volumes of input text due to the compression in the form of notes. *NotesWriting* is framework agnostic and can be integrated with different iterative RAG methods. We demonstrate its effectiveness with three iterative RAG methods, across two models and four evaluation datasets. *NotesWriting* yields an average improvement of 15.6 percentage points overall, by scaling the amount of ingested information.

Keywords: RAG, MultiHop, Question Answering

1. Introduction

The Retrieval-Augmented Generation (RAG) paradigm has significantly advanced open-domain question answering (Zhang et al., 2022; Kamaloo et al., 2023) by incorporating external knowledge sources (Borgeaud et al., 2022; Shi et al., 2023b; Izacard et al., 2023). This augmentation enables Large Language Models (LLMs) (Hurst et al., 2024; Dubey et al., 2024) to overcome the limitations of static parametric memory (Dhingra et al., 2022; Kasai et al., 2023) and mitigates hallucination issues (Ji et al., 2023) by grounding responses in retrieved evidence.

For complex tasks such as multi-hop question answering (Yang et al., 2018; Zhu et al., 2024; Krishna et al., 2024), which require compositional reasoning across multiple documents, single-round RAG approaches—those retrieving evidence based solely on the initial question—often retrieve only partial evidence, failing to gather all the information necessary to answer the query accurately. To address this limitation, iterative RAG methods such as IRCot (Trivedi et al., 2022), FLARE (Jiang et al., 2023), and ReAct (Yao et al., 2023) have been proposed. These approaches interleave retrieval and reasoning over multiple steps, progressively refining the query and accumulating supporting evidence to improve answer accuracy.

However, iterative retrieval introduces new challenges. Each reasoning step brings in additional retrieved content, much of which can be noisy or irrelevant. Prior studies have shown that excessive or low-quality context can severely degrade

the performance of RAG systems (Petroni et al., 2020; Shi et al., 2023a; Leng et al., 2024; Wu et al., 2024). In multi-hop settings, this problem is exacerbated as iterative methods continuously append new documents, leading to long and cluttered input contexts. Simply concatenating all retrieved documents across steps introduces several key issues:

- **Context Overload:** The aggregated content often exceeds the LLM's context window (Krishna et al., 2024), forcing truncation of potentially relevant information.
- **Computational Cost & Scalability:** Processing long inputs increases inference time and memory consumption, limiting scalability (Yue et al., 2024).
- **Distraction:** Irrelevant or redundant content distracts the model and impairs its ability to reason and plan effectively (Yu et al., 2023; Chen et al., 2024; Xie et al., 2024; Aghzal et al., 2025).
- **Readability:** Long reasoning traces become harder for users to interpret, especially when cluttered with repetitive or tangential content.

To address these challenges, we propose *NotesWriting*, a scalable and effective method that operates at each retrieval step. The core idea is to employ a note-taking LLM to generate concise, query-relevant notes from retrieved documents, thereby retaining only the essential information necessary for the subsequent reasoning step. This focused content selection prevents the primary LLM's context window from being overloaded with irrelevant details, substantially increasing the *effective context length* and enhancing planning

and reasoning capabilities. A key advantage of *NotesWriting* is its modular and easily integrable design, enabling seamless adoption within a wide range of iterative RAG frameworks.

Our key contributions are as follows:

- We introduce *NotesWriting*, a method to enhance the *effective context length* in iterative RAG. It mitigates context overload, reduces redundancy, and enables more efficient reasoning with fewer reasoning steps.
- *NotesWriting* is designed to be modular and easily integrable across different iterative RAG pipelines. By condensing context and filtering out noise, it directly supports improved planning and reasoning.
- We validate the effectiveness of *NotesWriting* through comprehensive experiments across three iterative RAG frameworks (IRCoT, FLARE, and ReAct), four multi-hop QA datasets, and two LLMs. Our results show an average improvement of 15.6 percentage points, demonstrating its ability to scale with larger input volumes.
- We show that *NotesWriting* improves the reasoning process by guiding the model to generate more precise search queries, which leads to the retrieval of more accurate and relevant documents. This effect is especially pronounced when integrated with ReAct (ReNAct), where it achieves the highest performance gains in our evaluation.

2. Background and Related Work

Single-Step vs. Iterative RAG. Traditional RAG often operates in a single step: retrieve relevant documents based on the initial query, then generate the final response conditioned on both the query and the retrieved context. While effective for simpler questions, this retrieve-then-read approach struggles for multi-hop QA, where the information needed evolves throughout the reasoning process. Iterative RAG addresses this limitation by interleaving retrieval and generation. The model can issue multiple queries, gather information incrementally, and refine its reasoning path based on newly retrieved evidence. This dynamic interaction between the LLM and the retriever is better suited for complex, multi-step reasoning.

Advances in Iterative RAG. Several approaches have explored different strategies within this iterative framework. IRCoT (Trivedi et al., 2022) integrates Chain-of-Thought reasoning with retrieval by using each reasoning step as the next query, enabling tighter coupling between thought and evidence. FLARE (Jiang et al., 2023) introduces a confidence-based retrieval mechanism that selectively queries external knowledge only when

uncertainty arises during generation. ReAct (Yao et al., 2023) separates reasoning from retrieval by explicitly generating search actions, allowing more structured control over when and how information is accessed. These approaches highlight the growing emphasis on dynamic, context-aware retrieval strategies in complex language tasks. ReSP (Retrieve, Summarize, Plan) (Jiang et al., 2025) applies query-focused summarization in multi-hop QA, maintaining both global and local summaries to mitigate context overload. Each iteration involves multiple LLM calls for planning sub-questions, summarizing retrieved content, and deciding whether sufficient evidence has been gathered. While modularization helps manage complex tasks, it introduces potential failure points across stages and increases latency due to repeated summarization and planning steps.

Context Management in Iterative RAG. A key challenge in iterative RAG is managing retrieved context across steps. While long-context LLMs help, retrieval remains essential for complex tasks like multi-hop QA (Xu et al., 2023b; Lee et al., 2024a). However, LLMs often struggle with pinpointing relevant content in long inputs (Kamradt, 2023; Hsieh et al., 2024), limiting the effectiveness of retrieval-augmented approaches (Jiang et al., 2024).

To address context overload, recent work has proposed methods to compress, filter, or control retrieved content. RECOMP (Xu et al., 2023a) summarizes documents before generation, while CON (Yu et al., 2023) trains models to assess retrieval reliability through note-taking. Plan-RAG (Lee et al., 2024b) and SmartRAG (Gao et al., 2024) introduce planning and policy modules to optimize when and what to retrieve. Self-RAG (Asai et al., 2023) uses reflection tokens to prune less relevant documents.

Most of these approaches focus on single-step retrieval and do not generalize to multi-step settings. InfoAgent (Reddy et al., 2024) addresses this with dedicated modules (Navigator, Extractor, Aggregator) for iterative context handling, while ReSP (Jiang et al., 2025) maintains global and local memory queues to manage context growth. However, such multi-module designs can reduce modularity and scalability. In contrast, *NotesWriting* is lightweight, modular, and easily integrable into existing iterative RAG workflows.

3. Method

To address the challenges of context overload and information noise in iterative RAG, particularly for multi-hop QA, we introduce *NotesWriting*, a method for generating concise, query-relevant

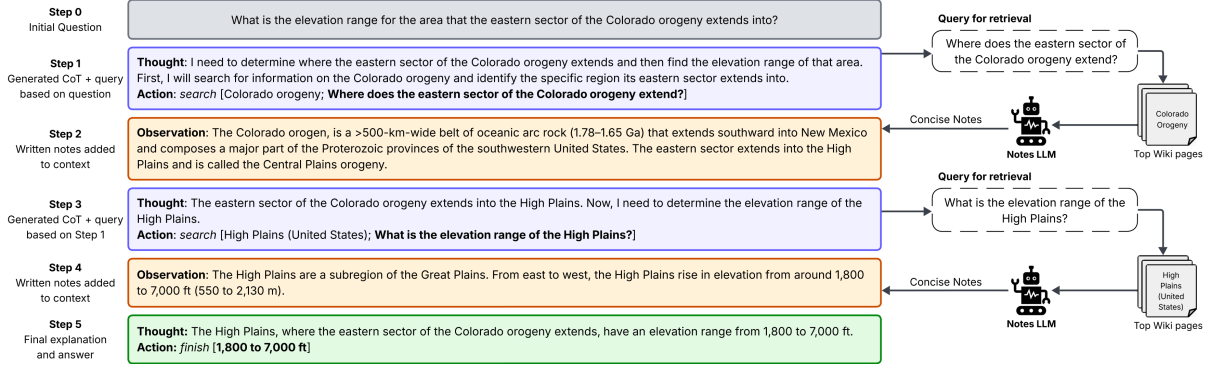


Figure 1: Overview of *NotesWriting* within an iterative RAG framework.

notes from retrieved documents at each step. Instead of feeding raw retrieved documents to the main LM (LM_{main}), *NotesWriting* first processes them to extract key information, thereby reducing context length and filtering irrelevant content.

3.1. Formulation of Iterative RAG.

Let x be the user input question, and $\mathcal{D} = \{d_i\}_{i=1}^{|\mathcal{D}|}$ represent the external knowledge corpus (e.g., Wikipedia). An iterative RAG process aims to generate a sequence of reasoning steps or partial outputs $s = [s_1, s_2, \dots, s_n]$. We denote the language model as $LM(\cdot)$ and the retrieval function, which returns the top- k documents for a query q , as $\text{ret}(q)$.

At each step $t \geq 1$, the typical process involves:

1. **Query Formulation:** A query q_t is generated based on the initial input x and the preceding steps $s_{<t} = [s_1, \dots, s_{t-1}]$. This is governed by a query formulation function $\mathcal{Q}(\cdot)$:

$$q_t = \mathcal{Q}(x, s_{<t}) \quad (1)$$

For the first step, $s_{<1} = \emptyset$, and often $q_1 = x$.

2. **Retrieval:** The retriever fetches the top- k relevant documents: $\mathcal{D}_{q_t} = \text{ret}(q_t)$.
3. **Generation:** The LM generates the next reasoning step s_t using the original input, previous steps, and the newly retrieved documents:

$$s_t = LM([\mathcal{D}_{q_t}, x, s_{<t}]) \quad (2)$$

This process continues until a final answer is generated, or a maximum number of steps is reached. Several approaches instantiate this iterative RAG framework with different strategies for query formulation and control flow:

IRCoT uses the last generated reasoning step as the next retrieval query, tightly coupling generation and retrieval:

$$q_t = \mathcal{Q}_{\text{IRCoT}}(x, s_{<t}) = s_{t-1} \quad (3)$$

FLARE (Jiang et al., 2023) proactively triggers retrieval based on token-level confidence during generation. If low-confidence tokens appear, retrieval is initiated using the prior sentence; otherwise, it is

skipped:

$$q_t = \begin{cases} s_{t-1}, & \text{if } \text{conf}(s_{t-1}) > \theta \\ \emptyset, & \text{otherwise} \end{cases} \quad (4)$$

ReAct (Reason and Act) (Yao et al., 2023) explicitly separates reasoning (Thought) from retrieval (Action). The query at each step is formed from the generated action:

$$q_t = \mathcal{Q}_{\text{ReAct}}(x, s_{<t}) = \text{Action}_t \quad (5)$$

This decoupling enables explicit control over when and how retrieval is invoked during reasoning.

3.2. *NotesWriting*: Iterative Note Extraction

The core idea is to use a dedicated, smaller language model (LM_{notes}) to act as a note-taker. At each iteration t , after retrieving the top- k documents $\mathcal{D}_{q_t} = \{d_1, d_2, \dots, d_k\}$ based on the query q_t , *NotesWriting* performs the following:

1. **Note Extraction:** For each retrieved document d_i , LM_{notes} is prompted (using prompt $\mathcal{P}_{\text{notes}}$, see Section 3.2) to extract concise notes r_i relevant to the current query q_t :

$$r_i = LM_{\text{notes}}(q_t, d_i) \quad (6)$$

2. **Note Aggregation:** The extracted notes from all k documents are aggregated as \mathbf{O}_t :

$$\mathbf{O}_t = \bigcup_{i=1}^k r_i \quad (7)$$

This process replaces the direct feeding of potentially long and noisy documents \mathcal{D}_{q_t} with the much shorter and focused notes \mathbf{O}_t .

3.3. ReNAcT: ReAct with *NotesWriting*

While *NotesWriting* is a generic module that can be integrated with different iterative RAG methods, results in Section 5 demonstrates that it works best with the ReAct framework (Yao et al., 2023). Therefore, we propose leveraging the ReAct framework as a suitable base for our approach. ReAct's structure explicitly separates reasoning (Thought) from

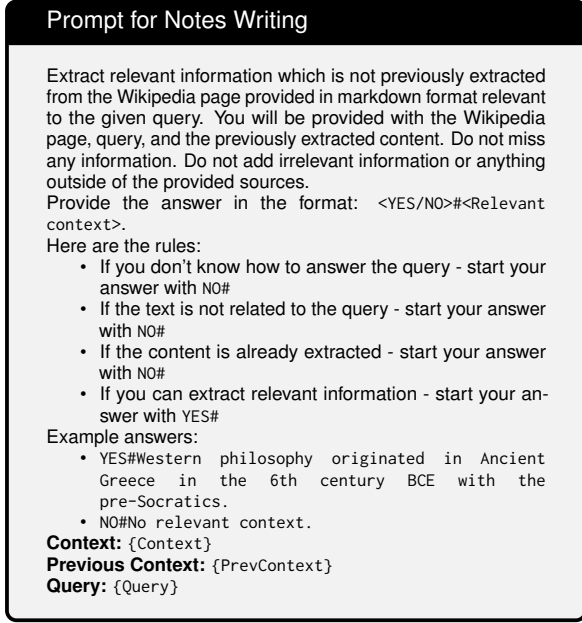


Figure 2: Notes writing prompt for extracting the relevant information.

information gathering (Action), where the action often involves generating a specific search query. This explicit query generation aligns naturally with the goal of targeted retrieval followed by focused note-taking. We combine ReAct with *NotesWriting* with the process at step t as follows:

1. LM_{main} generates Thought step outlining reasoning, along with an Action step, typically containing a search query. This query becomes q_t :

$$q_t = Q_{ReAct}(x, s_{<t}) = \text{SearchQueryFrom}(\text{Action}_t) \quad (8)$$

Retrieval is performed using q_t to get D_{q_t} .

2. *NotesWriting* processes D_{q_t} using LM_{notes} to generate aggregated notes, presented as the observation O_t .
3. The LM_{main} receives O_t and uses it along with x and $s_{<t}$ to generate the next Thought and Action pair:

$$s_t(\text{next Thought+Action}) = LM([\text{Observation: } O_t, x, s_{<t}]) \quad (9)$$

Iterations continue until the model generates a final answer within s_t or reaches a maximum number of iterations T , after which a final answer is synthesized based on the full history s and the collected notes $\{O_t\}_{t=1}^T$. This approach (illustrated in Figure 1) aims to combine the structured reasoning of ReAct with the context management benefits of *NotesWriting*, leading to a more robust and efficient iterative RAG system for complex QA.

Model	Dataset	Over Limit (n / %)	F1 (%)
GPT-4o-mini	Frames	84.3%	5.2
	FanoutQA	78.7%	4.7
LLaMA 3.1 70B	Frames	88.8 %	4.6
	FanoutQA	82.2%	3.8

Table 1: Full context RAG results. Over Limit shows the percentage of questions not answered due to context length getting exceeded.

3.4. NotesWriting: A Plug-and-Play Module for Iterative RAG

NotesWriting is designed as a complementary module that can be integrated into various iterative RAG frameworks. It modifies the generation step (Eq. 5) while keeping the specific query formulation Q of the base framework. We demonstrate this integration with two SOTA iterative RAG frameworks: IRCoT and FLARE.

IRCoT: Query remains the last generated sentence ($q_t = s_{t-1}$). Generation step becomes:

$$s_t = LM([O_t, x, s_{<t}]) \text{ where } O_t \text{ is derived from } D_{q_t} = \text{ret}(s_{t-1}) \quad (10)$$

FLARE: Query formulation remains conditional based on confidence θ ($q_t = Q_{FLARE}(x, s_{<t})$). If retrieval occurs ($q_t \neq \emptyset$), the generation step uses the extracted notes:

$$s_t = LM([O_t, x, s_{<t}]) \text{ where } O_t \text{ is derived from } D_{q_t} \quad (11)$$

If retrieval is skipped ($q_t = \emptyset$), $O_t = \emptyset$ and generation proceeds without new retrieved context.

4. Experiments

4.1. Datasets

(1) **FanoutQA** (Zhu et al., 2024) focuses on "fanout" multi-hop, multi-document complex questions that require gathering information about a large set of entities. We report results on the dev set containing 310 questions.

(2) **FRAMES** (Krishna et al., 2024) a challenging multi-hop QA dataset requiring 2–15 hops to answer questions. We exclude questions requiring tabular reasoning and evaluate on 549 examples.

(3) **HotpotQA** (Yang et al., 2018) a popular multi-hop QA dataset that requires reasoning over 2-3 Wikipedia articles. reasoning. We report results on 500 examples from the dev set.

(4) **MultiHop-RAG** (Tang and Yang, 2024) a non-wikipedia based benchmark that involves retrieval over recent news articles. It has ~ 600 news articles. For each question, we used BM25 to get the top five news articles in each iteration.

Model	Setting	Benchmark	F1 (%)	GPT-4	Main Tokens		Notes Tokens	
				Score (%)	Input	Output	Input	Output
GPT-4o-mini	ReAct	Fanout-QA	28.6	12.9	116K	916	-	-
		Frames	8.7	31.1	67K	707	-	-
		Hotpot-QA	42.2	56.4	26K	319	-	-
		MultiHop-RAG	58.0	64.2	188K	278	-	-
	ReAct + <i>NotesWriting</i> (ReNAcT)	Fanout-QA	50.0	28.0	17K	598	359K	675
		Frames	46.8	52.3	16K	543	277K	607
		Hotpot-QA	51.0	64.0	9K	326	130K	321
		MultiHop-RAG	58.0	70.6	46K	368	68K	390
LLaMA-3.1-70B	ReAct	Fanout-QA	13.5	8.7	113K	506	-	-
		Frames	21.7	26.8	85K	433	-	-
		Hotpot-QA	43.7	52.6	49K	289	-	-
		MultiHop-RAG	53.6	61.4	180K	295	-	-
	ReAct+ <i>NotesWriting</i> (ReNAcT)	Fanout-QA	43.0	26.1	15K	485	265K	1116
		Frames	49.0	57.6	13K	412	193K	717
		Hotpot-QA	55.5	67.4	8K	274	109K	391
		MultiHop-RAG	63.5	73.0	47K	262	76K	425

Table 2: ReAct and *NotesWriting* results for GPT-4o-mini and LLaMA. Main tokens represent the total number of input & output tokens for the main LLM across all steps (average on all questions). Similarly, notes tokens represent the total number of input & output tokens across all steps by the notes writing LLM (averaged on all questions). Token counts are rounded to the nearest thousand.

Evaluation metrics. Following (Krishna et al., 2024) we report the **F1 score** between predicted and ground truth answer and **GPT4-as-Judge score**. We also measure the **effective context length** by reporting the average number of input & output tokens processed by LM_{main} and LM_{notes} across all steps/iterations. We report the token counts only for *NotesWriting* and ReAct baseline due to space constraints and similar observed trends across other baselines. Finally, we conduct quality evaluation and the number of reasoning step comparison

4.2. Models & Implementation Details

We experiment with two LLMs, 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 generating reasoning step and *NotesWriting* (i.e $LM = LM_{notes}$). 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. GPT-4o-mini, which has a context length of 128K. 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*.

For *NotesWriting*, we utilize the Wikipedia API to fetch the top 5 relevant pages based on the query q_t . Each retrieved Wikipedia page is converted to Markdown format using `markdownify`² before being processed by LM_{notes} .

4.3. Baselines

Full-context RAG — In this baseline, we simply pass the top retrieved documents at each step to the LM_{main} . We report its results separately in Table 1, as over 80% of the questions are unanswerable due to the LLM’s context length limitations, making this setting non-competitive compared to *NotesWriting* and other baselines.

ReAct — We re-implement the original baseline (Yao et al., 2023) which allows the LLM to search that retrieves the first paragraph of the top 5 Wikipedia pages, `select` that allows ReAct to select relevant page for first 10 passages and `lookup` that returns paragraphs containing that specific string in the selected page.

IRCoT & FLARE — These were originally designed for older completion-based models such as `text-davinci-003` which relied heavily on internal parametric knowledge to generate CoTs. However, such a design is not directly applicable to questions requiring step by step planning and up-to-date knowledge. To address this, we modify the existing few-shot prompts (Trivedi et al., 2022; Jiang et al., 2023) to be compatible with chat mod-

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

²<https://pypi.org/project/markdownify/>

els, enabling them not only to generate CoTs but also to explicitly plan ahead. Moreover, the original baselines used BM25 from an older Wikipedia dump. However, in initial experiments we observed that the older dump is outdated for latest datasets. Therefore, we used a recent dump 20231101.en³ and dense passage retrieval with ef-base-v2 embeddings (Wang et al., 2022). We set the selective retrieval parameter θ to 0.8 for all our experiments. We would like to highlight that IRCOT & FLARE adds the entire content retrieved to the context of LLM similar to the full context RAG baseline which frequently results in LLMs’ context limit being exceeded (see Table 1). Therefore, to enable fair comparison, we use a chunked document setup to fit retrieved information in the context window instead of full Wiki pages for the two baselines as described in the respective papers. While ReSP represents another possible baseline, direct comparisons are difficult owing to their framework not having publicly released code, and missing key implementations in the paper leading to replication being difficult. Moreover, the ReSP methodology is less scalable owing to the need for global and local memory creation as well as the limited maximum number of steps.

InfoAgent (Reddy et al., 2024) — employs a navigator, extractor, and aggregator to retrieve and manage complex information. We use the official code provided here (Apache 2.0. License) with some modifications. We replaced the OpenAI embedding by sentence-transformers’ all-mpnet-base-v2⁴, and the Serper Google Search⁵ by Wikipedia search API due to costs and to use similar open knowledge tools as used in our method.

5. Results

Full context RAG performs poorly Table 1 reveals that feeding the full content of the top-5 retrieved Wikipedia pages directly into the main LM (as standard iterative RAG) caused context length errors, especially on challenging benchmarks like Fanout-QA and FRAMES. Table 1 shows that approximately 80% questions are unanswerable due to context build-up causing the context window to be exceeded.

Enhanced performance with NotesWriting. From Table 2, 3 and 4 in comparison to the respective baselines, *NotesWriting* shows significant improvements across all models and benchmarks. Specifically from Table 2, on complex long-form

³<https://huggingface.co/datasets/wikimedia/wikipedia>

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

⁵<https://serper.dev/>

Model	Setting	Benchmark	F1 (%)	GPT-4 (%)
GPT-4o-mini	IRCoT	FanoutQA	33.6	15.2
		Frames	24.0	22.0
		HotpotQA	36.4	42.8
		M-RAG	23.7	48.0
	+NotesWriting	FanoutQA	41.9	21.3
		Frames	43.9	42.3
		HotpotQA	46.2	53.8
		M-RAG	36.0	65.6
LLaMA-3.1-70B	IRCoT	FanoutQA	21.0	8.4
		Frames	19.4	21.1
		HotpotQA	31.5	38.2
		M-RAG	35.4	64.8
	+NotesWriting	FanoutQA	36.0	22.9
		Frames	26.9	33.3
		HotpotQA	36.5	53.0
		M-RAG	38.0	64.8

Table 3: IRCOT scores for GPT-4o-mini and LLaMA-3.1-70B M-RAG represents MultiHop-RAG.

multihop-QA datasets like *FRAMES* and *Fanout-QA*, on average ReNAct achieves an absolute improvement of 29.1 points in F1 score and 21.1 points in GPT-4 score. On relatively easier datasets such as *Hotpot-QA* and *MultiHop-RAG*, ReNAct yields absolute improvements of 10.3 and 5.0 points, respectively. The strong results compared to the baseline demonstrate that the LLM is receiving correct and relevant information at each step with *NotesWriting*.

From Tables 3 and 4 on challenging datasets *NotesWriting* coupled with each of IRCOT and FLARE leads to 14.4 and 10.5 points improvement on F1 and GPT-4 score. Similarly on Hotpot-QA and MultiHop-RAG we find 7.0 and 10.8 points improvement on F1 and GPT-4 score respectively. Finally Table 5 shows that *NotesWriting* outperforms InfoAgent by significant margin.

Increased effective context length. Table 2 show the average number of input and output tokens across all steps for the baseline and *NotesWriting*. The total number of tokens processed by the system (sum of input tokens across LM_{main} and LM_{notes}) increases by 100K on average for *NotesWriting*, allowing the model to reason over more retrieved content. However, it is important to note that this information cannot be naïvely appended to the LM_{main} context as full context RAG (Table 1) performs poorly and summing columns 7 and 9 exceeds the context window. This shows how *NotesWriting* enables scalable use of the large retrieval context by delegating information management to a specialized LLM.

Reduced context overload for the main LLM. With ReAct (Table 2), the number of tokens for the LM_{main} reduce significantly from baseline to ReNAct

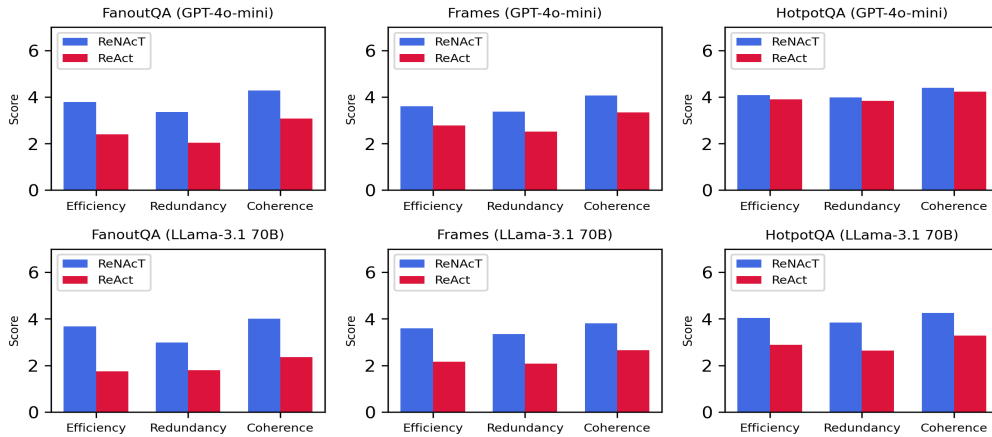


Figure 3: Quality evaluation of ReAct and ReNAct reasoning chain.

across all benchmarks — by 77K tokens for GPT-4o-mini and 86K tokens for LLaMA-3.1-70B on average. The same trend is observed with IRCot and where the LM_{main} input tokens reduces by at least $4\times$ and $1.5\times$ for GPT-4o-mini and LLaMA-3.1-70B. This shows that concise notes are added at each retrieval step and help the LM_{main} plan better as seen in reasoning quality analysis in Section 6 (Figure 3).

NotesWriting trades off performance with cost.

Table 2 show that with *NotesWriting*, the combined output tokens from the main and note-taking LLM are on average $2\times-3\times$ higher than the baseline. However, we emphasize that the increase in output tokens (approximately 500 tokens on average) is still significantly smaller than the increase in input tokens (around 100K on average). Since output tokens are the primary contributors to both latency and cost⁶, *NotesWriting* offers a much more cost, compute, and performance efficient approach.

6. Analysis

Reasoning Quality Analysis. We evaluate the reasoning chains generated by ReAct and ReNAct using GPT-4o as a judge across three axes, (1) Efficiency — to measure redundant searches and how well each step contributes to the final answer, (2) Redundancy — to assess repeated search queries, or unnecessary repetition or duplication of steps (3) Coherence — to check if the chain is comprehensible, logically connected, and free from unnecessary complexity or ambiguity. Figure 3 shows the results. ReNAct is better across all three axes than ReAct across on all models and datasets. Specifically, on Frames and FanoutQA across both models efficiency, redundancy and coherence improve

Model	Setting	Benchmark	F1 (%)	GPT-4 (%)
GPT-4o-mini	FLARE	Fanout-QA	35.1	14.2
		Frames	26.3	23.7
		HotpotQA	34.8	39.0
		M-RAG	28.9	65.7
	+NotesWriting	Fanout-QA	42.3	22.2
		Frames	27.7	29.8
		HotpotQA	34.5	45.8
		M-RAG	30.2	66.6
LLaMA-3.1-70B	FLARE	FanoutQA	23.0	11.4
		Frames	16.4	18.6
		Hotpot-QA	24.7	31.2
		M-RAG	36.1	67.0
	+NotesWriting	FanoutQA	35.8	24.2
		Frames	20.0	25.3
		HotpotQA	34.0	47.0
		M-RAG	30.5	66.4

Table 4: FLARE scores for GPT-4o-mini and LLaMA-3.1-70B. M-RAG is MultiHop-RAG.

Setting	Benchmark	F1 (%)	GPT-4 (%)
InfoAgent	FanoutQA	47.2	22.9
	Frames	28.0	29.9
NotesWriting	FanoutQA	50.0	28.0
	Frames	46.8	52.3

Table 5: Infoagent vs *NotesWriting* performance comparison on GPT-4o-mini.

by at least $1.5\times$. On HotpotQA, the improvement is $1.2\times$.

Figure 4 shows the comparison of the number of ground truth steps, ReAct and ReNAct search steps for each question in each dataset across both models. The dashed lines for each method represents the in-correct answers and the solid line represents correct ones. The x-axis is the index of the question in the dataset sorted by the number of ground-truth search steps. From the Figure 4,

⁶<https://openai.com/api/pricing/>

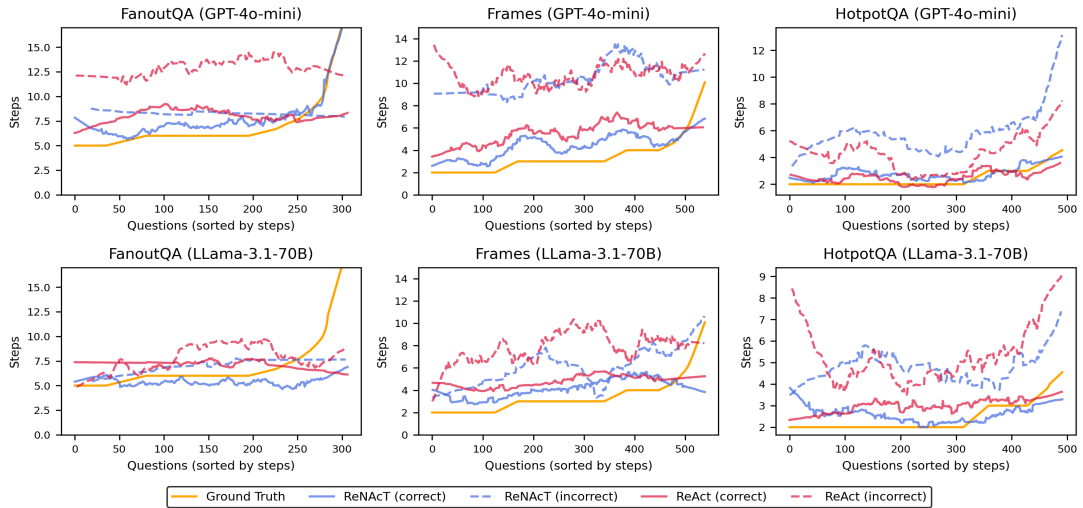


Figure 4: Steps (smoothed) by ReNAct, ReAct vs the ground truth steps for GPT-4o-mini and LLaMA-3.1-70B.

Model	Dataset	ReNAct	ReAct
GPT-4o-mini	Fanout	22.83	36.84
	Frames	23.18	29.89
	HotpotQA	30.56	46.02
LLaMA	Fanout	49.44	35.71
	Frames	28.22	30.94
	HotpotQA	31.33	32.82

Table 6: % of correct questions having search steps less than number of ground truth Wikipedia pages.

it can be observed that ReNAct (solid blue line) is much closer to the ground truth steps with ReAct (solid red line) being relatively far demonstrating the effectiveness of *NotesWriting* in coming up with correct stepwise plan and search query for retrieval. Figure 4 also shows that the in-correct questions (dashed red & blue line) have a higher number of steps that shows that it fails after many re-tries. The difference among ReNAct ReAct and ground truth steps is more significant in challenging datasets like Frames and Fanout-QA as opposed to HotpotQA.

Correct Answers with Fewer Searches than Ground Truth. Figure 4 shows cases where ReNAct and ReAct (solid blue & red lines) are below ground truth steps. Table 6 reports percentages of cases for the number of correctly answered questions which took less searches than the number of ground truth Wikipedia pages required to answer the question correctly.

Search Steps Comparison.

Standard deviation across runs We ran the *NotesWriting* and across all datasets and models three times to see the variance across different runs. We report the results in Table 7.

Model	Dataset	Avg F1	GPT-4 Score
GPT-4o-mini	Fanout	± 1.86	± 2.45
	Frames	± 1.10	± 2.35
LLaMA-3.1 70B	Fanout	± 3.79	± 1.54
	Frames	± 4.42	± 5.76

Table 7: Standard deviation across Frames & FanoutQA.

7. Conclusion

We present *NotesWriting*, a plug-and-play module that improves *effective context length* in Iterative RAG by accumulating only the most relevant information at each reasoning step. Experiments on *three* RAG baselines (IRCoT, FLARE, and ReAct), *four* multi-hop QA datasets, and *two* LLMs show that *NotesWriting* improves performance by up to 15.6 points, while also reducing context overload, the number of reasoning steps, and redundancy. In the ReAct setting, *NotesWriting* enables better planning by guiding the model to generate more accurate search queries and retrieve the correct documents. Moreover, *NotesWriting* consistently improves coherence and efficiency of planning and search across models in ReAct. Therefore, we suggest ReNAct as an effective iterative RAG framework. Our results show that ReNAct (ReAct + *NotesWriting*) makes iterative RAG more scalable and precise.

8. 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 online searches to the Wikipedia API⁷, 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. 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 *NotesWriting*.

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 *NotesWriting* method can contribute towards better task handling at large. We will make our code, generated reasoning traces, evaluation/quality assessment prompts and additional experimental details publicly available upon acceptance towards this goal.

9. Bibliographical References

Mohamed Aghzal, Erion Plaku, Gregory J Stein, and Ziyu Yao. 2025. A survey on large language models for automated planning. *arXiv preprint arXiv:2502.12435*.

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*.

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

Yanan Chen, Ali Pesaranghader, Tanmana Sadhu, and Dong Hoon Yi. 2024. Can we rely on llm agents to draft long-horizon plans? let's take travelplanner as an example. *arXiv preprint arXiv:2408.06318*.

Bhuvan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. 2022. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273.

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*.

Jingsheng Gao, Linxu Li, Weiyuan Li, Yuzhuo Fu, and Bin Dai. 2024. Smartrag: Jointly learn rag-related tasks from the environment feedback. *arXiv preprint arXiv:2410.18141*.

Cheng-Ping Hsieh, Simeng Sun, Samuel Krieman, 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*.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, 24(251):1–43.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12):1–38.

Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,

⁷<https://www.mediawiki.org/wiki/API:Search>

- Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*.
- Zhouyu Jiang, Mengshu Sun, Lei Liang, and Zhiqiang Zhang. 2025. [Retrieve, summarize, plan: Advancing multi-hop question answering with an iterative approach](#).
- Ziyang Jiang, Xueguang Ma, and Wenhui Chen. 2024. Longrag: Enhancing retrieval-augmented generation with long-context llms. *arXiv preprint arXiv:2406.15319*.
- Ehsan Kamaloo, Nouha Dziri, Charles LA Clarke, and Davood Rafiei. 2023. Evaluating open-domain question answering in the era of large language models. *arXiv preprint arXiv:2305.06984*.
- Greg Kamradt. 2023. Needle in a haystack-pressure testing llms. *GitHub Repository*, page 28.
- Jungo Kasai, Keisuke Sakaguchi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, Kentaro Inui, et al. 2023. Real-time qa: What’s the answer right now? *Advances in neural information processing systems*, 36:49025–49043.
- Satyapriya Krishna, Kalpesh Krishna, Anhad Mohananeey, 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.
- Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua, Devendra Singh Sachan, Michael Boratko, Yi Luan, Sébastien MR Arnold, Vincent Perot, Siddharth Dalmia, et al. 2024a. Can long-context language models subsume retrieval, rag, sql, and more? *arXiv preprint arXiv:2406.13121*.
- Myeonghwa Lee, Seonho An, and Min-Soo Kim. 2024b. Planrag: A plan-then-retrieval augmented generation for generative large language models as decision makers. *arXiv preprint arXiv:2406.12430*.
- Quinn Leng, Jacob Portes, Sam Havens, Matei Zaharia, and Michael Carbin. 2024. Long context rag performance of large language models. *arXiv preprint arXiv:2411.03538*.
- Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2020. How context affects language models’ factual predictions. *arXiv preprint arXiv:2005.04611*.
- Revant 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*.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023a. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023b. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652*.
- Yixuan Tang and Yi Yang. 2024. Multihoprag: 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*.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.
- Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai Zhang, and Yanghua Xiao. 2024. How easily do irrelevant inputs skew the responses of large language models? *arXiv preprint arXiv:2404.03302*.
- Jian Xie, Kexun Zhang, Jiangjie Chen, Siyu Yuan, Kai Zhang, Yikai Zhang, Lei Li, and Yanghua Xiao. 2024. Revealing the barriers of language agents in planning. *arXiv preprint arXiv:2410.12409*.
- Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2023a. Recomp: Improving retrieval-augmented llms with compression and selective augmentation. *arXiv preprint arXiv:2310.04408*.
- Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee, Chen Zhu, Zihan Liu, Sandeep Subramanian, Evelina Bakhturina, Mohammad

10. Appendix

10.1. Examples comparing ReNAct with baselines

Shoeybi, and Bryan Catanzaro. 2023b. Retrieval meets long context large language models. In *The Twelfth International Conference on Learning Representations*.

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, 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)*.

Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. 2023. Chain-of-note: Enhancing robustness in retrieval-augmented language models. *arXiv preprint arXiv:2311.09210*.

Zhenrui Yue, Honglei Zhuang, Aijun Bai, Kai Hui, Rolf Jagerman, Hansi Zeng, Zhen Qin, Dong Wang, Xuanhui Wang, and Michael Bendersky. 2024. Inference scaling for long-context retrieval augmented generation. *arXiv preprint arXiv:2410.04343*.

Qin Zhang, Shangsi Chen, Dongkuan Xu, Qingqing Cao, Xiaojun Chen, Trevor Cohn, and Meng Fang. 2022. A survey for efficient open domain question answering. *arXiv preprint arXiv:2211.07886*.

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.

Model	Setting	Benchmark	F1 (%)	GPT-4 Score (%)	Main Tokens		Notes Tokens	
					Input	Output	Input	Output
GPT-4o-mini	Baseline	Fanout-QA	33.6	15.2	273K	385	-	-
		Frames	24.0	22.0	183K	312	-	-
		Hotpot-QA	36.4	42.8	99K	205	-	-
		MultiHop-RAG	23.7	48.0	909K	322	-	-
	NotesWriting	Fanout-QA	41.9	21.3	68K	444	902K	1.8K
		Frames	43.9	42.3	37K	280	658K	862
		Hotpot-QA	46.2	53.8	26K	193	433K	650
		MultiHop-RAG	36.0	65.6	189K	412	40K	324
LLaMA-3.1-70B	Baseline	Fanout-QA	21.0	8.4	197K	511	-	-
		Frames	19.4	21.1	192K	732	-	-
		Hotpot-QA	31.5	38.2	167K	557	-	-
		MultiHop-RAG	35.4	64.8	429K	313	-	-
	NotesWriting	Fanout-QA	36.0	22.9	154K	552	1.4M	7.6K
		Frames	26.9	33.3	116K	735	1.2M	3.9K
		Hotpot-QA	36.5	53.0	100K	561	1.1M	3.3K
		MultiHop-RAG	38.0	64.8	99K	680	338K	2.9K

Table 8: IRCOT results for GPT-4o-mini and LLaMA-3.1-70B across Baseline and *NotesWriting*. Token counts are rounded to the nearest thousand (input/output).

Model	Setting	Benchmark	F1 (%)	GPT-4 Score (%)	Main Tokens		Notes Tokens	
					Input	Output	Input	Output
GPT-4o-mini	Baseline	Fanout-QA	35.1	14.2	268K	384	-	-
		Frames	26.3	23.7	50K	314	-	-
		Hotpot-QA	34.8	39.0	36K	218	-	-
		MultiHop-RAG	28.9	65.7	204K	319	-	-
	NotesWriting	Fanout-QA	42.3	22.2	65K	409	968K	2.2K
		Frames	27.7	29.8	43K	358	47K	56
		Hotpot-QA	34.5	45.8	33K	258	45K	49
		MultiHop-RAG	30.2	66.6	205K	312	14K	49
LLaMA-3.1-70B	Baseline	Fanout-QA	23.0	11.4	195K	453	-	-
		Frames	16.4	18.6	86K	640	-	-
		Hotpot-QA	24.7	31.2	83K	573	-	-
		MultiHop-RAG	36.1	67.0	301K	558	-	-
	NotesWriting	Fanout-QA	35.8	24.2	154K	555	1.3M	7.5K
		Frames	20.0	25.3	67K	592	47K	66
		Hotpot-QA	34.0	47.0	65K	563	45K	77
		MultiHop-RAG	30.5	66.4	67K	710	15K	38

Table 9: FLARE results for GPT-4o-mini and LLaMA-3.1-70B across Baseline and *NotesWriting*. Token counts are rounded to the nearest thousand (input/output).

<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 10: An example of *NotesWriting* on the FanoutQA dataset using GPT-4o-mini.

<p>Question: According to the 2000 United States census, what was the 2000 population of the birth city of the only 21st-century mayor of Austin, Texas who also served as mayor in the 1990s? Round your answer to the nearest thousand.</p>
<p>Thought: I need to determine the birth city of the only 21st-century mayor of Austin, Texas who also served as mayor in the 1990s and then find its population in the 2000 United States census. First, I will identify the mayor who meets this criterion.</p> <p>Action: search[Austin, Texas mayors; Who was the 21st-century mayor of Austin, Texas who also served as mayor in the 1990s?]</p> <p>Observation: (Result 1) Mayor of Austin - Kirk Watson took office as mayor on January 6, 2023, for a second term, having served as mayor from 1997 to 2001.</p> <p>(Result 2) 2024 Austin mayoral election - Kirk Watson</p> <p>(Result 3) Austin, Texas - Kirk Watson</p>
<p>Thought: The only 21st-century mayor of Austin, Texas who also served as mayor in the 1990s is Kirk Watson. I now need to find his birth city.</p> <p>Action: search[Kirk Watson; Where was Kirk Watson born?]</p> <p>Observation: (Result 1) Kirk Watson - Oklahoma City, Oklahoma, U.S.</p> <p>(Result 2) Austin, Texas - Kirk Watson was born in 1958 in Oak Park, Illinois</p>
<p>Thought: The information provided contains conflicting birth cities for Kirk Watson. To confirm his birth city, I will perform another search.</p> <p>Action: search[Kirk Watson (American politician); What is the confirmed birth city of Kirk Watson?]</p> <p>Observation: (Result 1) Kirk Watson - Oklahoma City</p>
<p>Thought: Based on the confirmed information, Kirk Watson's birth city is Oklahoma City. Now, I need to find the population of Oklahoma City in the 2000 United States census.</p> <p>Action: search[Oklahoma City, Oklahoma; What was the population of Oklahoma City in the 2000 United States census?]</p> <p>Observation: (Result 1) Oklahoma City - 506,132</p>
<p>Thought: The population of Oklahoma City in the 2000 United States census was 506,132. To answer the question, I need to round this number to the nearest thousand.</p> <p>Action: finish[506,000]</p>

Table 11: An example of *NotesWriting* on the Frames dataset using LLama-3.1-70B.

<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 12: An Example of *NotesWriting* on Frames dataset using GPT-4o-mini.

<p>Question: If my future wife has the same first name as the 15th first lady of the United States' mother and her surname is the same as the second assassinated president's mother's maiden name, what is my future wife's name?</p>
<p>Thought: I need to determine the first name of the 15th first lady of the United States and the surname of the second assassinated president's mother's maiden name.</p> <p>Action: search[15th first lady of the United States; Who is the 15th first lady and what is her mother's name?]</p> <p>Observation: No relevant information, try a different search term.</p>
<p>Thought: I need to rephrase my search to find the correct information about the 15th first lady of the United States and her mother's name. I will search for the name of the 15th first lady directly.</p> <p>Action: search[15th first lady of the United States; Who is the 15th first lady?]</p> <p>Observation: No relevant information, try a different search term.</p>
<p>Thought: It seems that I'm having trouble finding information about the 15th first lady of the United States. I will now search for a list of all first ladies to identify her name and then find her mother's name.</p> <p>Action: search[List of first ladies of the United States; What is the name of the 15th first lady?]</p> <p>Observation: (Result 1) List of first ladies of the United States - The 15th first lady was Harriet Lane, who served from March 4, 1857, to March 4, 1861.</p>
<p>Thought: I have identified the 15th first lady of the United States as Harriet Lane. Now, I need to find information about her mother's name.</p> <p>Action: search[Harriet Lane; What is the name of Harriet Lane's mother?]</p> <p>Observation: (Result 1) Harriet Lane - Harriet Lane's mother was Jane Ann Buchanan Lane.</p>
<p>Thought: I have found that Harriet Lane's mother was named Jane Ann Buchanan Lane. Now, I need to identify the second assassinated president of the United States and find his mother's maiden name.</p> <p>Action: search[second assassinated president of the United States; Who is the second assassinated president?]</p> <p>Observation: (Result 1) List of United States presidential assassination attempts and plots - James A. Garfield was the second U.S. president to be assassinated, occurring on July 2, 1881, less than four months after he took office.</p>
<p>Thought: I need to rephrase my search to find the correct information about James A. Garfield's mother's maiden name. I will search for information about James A. Garfield and his family.</p> <p>Action: search[James A. Garfield's mother; What was the maiden name of James A. Garfield's mother?]</p> <p>Observation: No relevant information, try a different search term.</p>

Table 13: An example of *NotesWriting* on the Frames dataset using GPT-4o-mini (1/2).

<p>Thought: It seems I'm having difficulty finding the maiden name of James A. Garfield's mother. I will search for biographies of James A. Garfield to find more details about his family.</p> <p>Action: search[James A. Garfield biography; What is known about James A. Garfield's family?]</p> <p>Observation: (Result 1) James A. Garfield - James Abram Garfield was born on November 19, 1831, in a log cabin in Moreland Hills, Ohio.</p>
<p>Thought: I still need to find the maiden name of James A. Garfield's mother. I will focus my search specifically on his mother's name.</p> <p>Action: search[James A. Garfield's family; What is the name of James A. Garfield's mother?]</p> <p>Observation: (Result 1) James A. Garfield - James A. Garfield's mother was named Eliza.</p>
<p>Thought: I am struggling to find the maiden name of James A. Garfield's mother. I will search for historical records or references that detail Eliza Garfield's family background.</p> <p>Action: search[Eliza Garfield biography; What is known about Eliza Garfield's family and maiden name?]</p> <p>Observation: No relevant information, try a different search term.</p>
<p>Thought: It seems I am having persistent difficulties finding the maiden name of James A. Garfield's mother, Eliza. I will try searching for historical documentation or genealogical resources that might provide her maiden name.</p> <p>Action: search[Eliza Garfield maiden name; What is the maiden name of James A. Garfield's mother?]</p> <p>Observation: (Result 1) Abbie C. B. Robinson - The maiden name of President James A. Garfield's mother is Eliza Ballou.</p>
<p>Thought: I have found that the maiden name of James A. Garfield's mother is Eliza Ballou. Now, I have the necessary information to form the name of my future wife. Her first name is Jane (from Jane Ann Buchanan Lane) and her last name is Ballou (from Eliza Ballou).</p> <p>Action: finish[Jane Ballou]</p>

Table 14: An example of *NotesWriting* on the Frames dataset using GPT-4o-mini (2/2).

<p>Question: What 2016 film stars actor Xavier Samuel and an actress who portrayed Vasilisa Dragomir in the film "Vampire Academy"?</p>
<p>Thought: I need to search for the film that stars Xavier Samuel, then find out who played Vasilisa Dragomir in "Vampire Academy" to identify the actress. Finally, I will check if both appeared in the same film in 2016.</p>
<p>Action: search[Xavier Samuel filmography]</p>
<p>Observation: Wikipedia Title: Xavier Samuel</p> <p>Xavier Samuel (born 10 December 1983) is an Australian film and theatre actor. He has appeared in leading roles in the feature films Adore, September, Further We Search, Newcastle, The Loved Ones, Frankenstein, A Few Best Men, and played Riley Biers...</p> <p>Wikipedia Title: Mohanlal filmography</p> <p>Mohanlal is an Indian actor, producer, distributor, host, and playback singer who has starred in mainstream blockbuster and art-house films. He has made over 355 appearances in feature films during the span of his career, primarily in Malayalam...</p> <p>Wikipedia Title: Benedict Samuel</p> <p>Benedict Samuel is an Australian actor, writer, and director best known for playing Jervis Tetch / Mad Hatter in the Fox crime series Gotham and as Owen in The Walking...</p> <p>Wikipedia Title: Gérard Depardieu</p> <p>Gérard Xavier Marcel Depardieu is a multinational French actor, known to be one of the most prolific in film history. An icon of French cinema, considered a world star in...</p> <p>Wikipedia Title: Jagadish filmography</p> <p>P. V. Jagadish Kumar, known as Jagadish, is an Indian actor, screenwriter, and television presenter. He has starred in over 400 Malayalam...</p>
<p>Thought: The search still did not yield the specific filmography of Xavier Samuel. I will try searching for his name along with the year 2016 to see if I can find a film he starred in that year.</p>
<p>Action: search[Xavier Samuel 2016 film]</p>
<p>Observation: Wikipedia Title: Xavier Samuel</p> <p>Xavier Samuel (born 10 December 1983) is an Australian film and theatre actor. He has appeared in leading roles in the feature films Adore, September, Further We Search, Newcastle, The Loved Ones, Frankenstein, A Few Best Men, and played Riley Biers...</p> <p>Wikipedia Title: Mr. Church</p> <p>Mr. Church is a 2016 American drama film directed by Bruce Beresford and written by Susan McMartin. The film stars Eddie Murphy as the title character with Britt Robertson, Xavier Samuel, Lucy Fry, Christian Madsen and Natascha McElhone also...</p> <p>Wikipedia Title: Benedict Samuel</p> <p>Benedict Samuel is an Australian actor, writer, and director best known for playing Jervis Tetch / Mad Hatter in the Fox crime series Gotham and as Owen in The Walking...</p> <p>Wikipedia Title: Samuel Hunt (actor)</p> <p>Samuel Caleb Hunt (born October 10, 1986) is an American actor best known for his reoccurring role as Greg "Mouse" Gerwitz in the television series Chicago Fire and Chicago P.D., playing Xavier Rosen in the third season of Empire, and for starring...</p> <p>Wikipedia Title: Spin Out (film)</p> <p>Spin Out is a 2016 Australian romantic comedy film directed by Tim Ferguson and Marc Gracie and starring Xavier Samuel and Morgan...</p>

Table 15: An example of ReAcT on the HotpotQA dataset using GPT-4o-mini. It starts by repeatedly searching similar query.

<p>Thought: I have found that Xavier Samuel stars in the film “Spin Out” in 2016. Now I need to find out who portrayed Vasilisa Dragomir in “Vampire Academy.”</p>
<p>Action: search[Vasilisa Dragomir Vampire Academy actress]</p>
<p>Observation: Wikipedia Title: Vampire Academy (film)</p> <p>Vampire Academy is a 2014 American fantasy comedy horror film directed by Mark Waters from a screenplay by Daniel Waters, based on the first book of the novel series of the same name by Richelle Mead. The film stars an ensemble cast led by Zoey...</p> <p>Wikipedia Title: Bloodlines (book series)</p> <p>Bloodlines is a book series by Richelle Mead. It is a spin-off of Mead’s Vampire Academy series. The books follow alchemist Sydney Sage.</p> <p>After the incident with Rose and Dimitri, Sydney is being closely watched by the Alchemists (humans who help...</p> <p>Wikipedia Title: List of fictional princesses</p> <p>This is a list of fictional princesses that have appeared in various works of fiction. This list is organized by medium and limited to well-referenced, notable examples of fictional...</p>
<p>Thought: The search did not provide the specific actress who portrayed Vasilisa Dragomir. I will select the “Vampire Academy” film page to find the cast details.</p>
<p>Action: select[Vampire Academy (film)]</p>
<p>Observation: 2014 fantasy comedy horror film by Mark Waters</p> <p> Vampire Academy </p> <p> Theatrical release poster </p> <p> Directed by Mark Waters </p> <p> Screenplay by Daniel Waters </p> <p> Based on <i>Vampire Academy</i> by Richelle Mead </p> <p> Produced by Don Murphy, Susan Montford, Michael Preger, Deepak Nayar </p> <p> Cinematography Tony Pierce-Roberts </p> <p> Edited by Chris Gill </p> <p> Music by Rolfe Kent </p> <p> Production companies Kintop Pictures, Reliance Entertainment, Angry Films, Preger Entertainment </p> <p> Distributed by The Weinstein Company </p> <p> Release dates February 7, 2014 (US), July 14, 2014 (UK straight-to-DVD) </p> <p><i>Vampire Academy</i> is a 2014 American fantasy comedy horror film directed by Mark Waters from a screenplay by Daniel Waters, based on the first book of the novel series by Richelle Mead. The film stars an ensemble cast led by Zoey Deutch as Rose Hathaway, a Dhampir, and Lucy Fry as Lissa Dragomir, her Moroi best friend, and follows their return to the St Vladimir’s Academy boarding school, after being on the run for one year.</p>
<p>Thought: I found that Lucy Fry portrays Vasilisa Dragomir in “Vampire Academy.” Now, I will check if both Xavier Samuel and Lucy Fry starred in the same film in 2016. Since Xavier Samuel is in “Spin Out,” I need to find out if Lucy Fry is also in this film.</p>
<p>Action: search[Spin Out film cast]</p>

Table 16: Continued example of ReAcT on the HotpotQA dataset using GPT-4o-mini. It focuses on "Spin Out" but misses the potential answer "Mr. Church".

<p>Observation: Wikipedia Title: Spin Out (film)</p> <p>Spin Out is a 2016 Australian romantic comedy film directed by Tim Ferguson and Marc Gracie and starring Xavier Samuel and Morgan...</p> <p>Wikipedia Title: Spinning Out</p> <p>Spinning Out is an American drama television series, created by Samantha Stratton, that premiered on Netflix on January 1, 2020. In February 2020, the series was canceled after one season.</p> <p>The series stars Kaya Scodelario as Kat Baker, a young ice...</p> <p>Wikipedia Title: Spin (2021 film)</p> <p>Spin is a 2021 American comedy film released as a Disney Channel Original Movie (DCOM). It is directed by Manjari Makijany from a screenplay by Carley Steiner and Josh A. Gagan. Produced by Disney Channel Original Productions, the film stars...</p> <p>Wikipedia Title: Spinning Man</p> <p>Spinning Man is a 2018 thriller film directed by Simon Kaijser based on a novel by George Harrar and adapted to screenplay by Matthew Aldrich. It stars Guy Pearce, Pierce Brosnan, Minnie Driver, Alexandra Shipp, Clark Gregg and Odeya Rush. The film...</p> <p>Wikipedia Title: Dog Man (film)</p> <p>Dog Man is a 2025 American animated superhero comedy film based on Dav Pilkey's children's graphic novel series Dog Man, produced by DreamWorks Animation and Scholastic Entertainment, and distributed by Universal Pictures. It is a spin-off and story...</p>
<p>Thought: The search did not provide the cast details for "Spin Out." I will select the "Spin Out (film)" page to find the cast.</p> <p>Action: select[Spin Out (film)]</p>
<p>Observation: 2016 Australian film</p> <p> Spin Out </p> <p> — — </p> <p> Theatrical film poster </p> <p> Directed by Tim Ferguson, Marc Gracie </p> <p> Written by Edwina Exton, Tim Ferguson </p> <p> Produced by Marc Gracie, David Redman </p> <p> Starring Xavier Samuel, Morgan Griffin, Lincoln Lewis </p> <p> Cinematography Justin Brickle </p> <p> Edited by Ken Sallows </p> <p> Music by Jon Hume </p> <p><i>Spin Out</i> is a 2016 Australian romantic comedy film directed by Tim Ferguson and Marc Gracie and starring Xavier Samuel and Morgan Griffin.</p>
<p>Thought: The search results confirm that "Spin Out" stars Xavier Samuel, but there is no mention of Lucy Fry in the cast list. Thus, they did not appear in the same film. I will now summarize the findings to answer the question.</p>
<p>Action: finish[There is no common film between Xavier Samuel and Lucy Fry in 2016.]</p>
<p>Observation: There is no common film between Xavier Samuel and Lucy Fry in 2016.</p>

Table 17: Continued example of ReAcT on the HotpotQA dataset using GPT-4o-mini. It fails by focusing exclusively on Xavier Samuel's film "Spin Out" while completely overlooking "Mr. Church," despite it being mentioned in the search results alongside both actors' names. This prevents it from discovering that "Mr. Church" is the 2016 film starring both Xavier Samuel and Lucy Fry (who played Vasilisa Dragomir in "Vampire Academy").

<p>Question: What 2016 film stars actor Xavier Samuel and an actress who portrayed Vasilisa Dragomir in the film "Vampire Academy"?</p>
<p>Thought: I need to identify the actress who portrayed Vasilisa Dragomir in "Vampire Academy" and then find a 2016 film that stars Xavier Samuel alongside her. First, I will search for the cast of "Vampire Academy" to determine who played Vasilisa Dragomir.</p> <p>Action: search[Vampire Academy cast; Who portrayed Vasilisa Dragomir in the film "Vampire Academy"?]</p> <p>Observation: (Result 1) Vampire Academy (film) - Lucy Fry portrayed Vasilisa "Lissa" Dragomir in the film "Vampire Academy".</p>
<p>Thought: Lucy Fry portrayed Vasilisa Dragomir in "Vampire Academy." Now, I need to find a 2016 film that stars both Xavier Samuel and Lucy Fry.</p> <p>Action: search[Xavier Samuel; What films did Xavier Samuel star in during 2016?]</p> <p>Observation: (Result 1) Xavier Samuel - In 2016, Xavier Samuel played the titular character in *The Death and Life of Otto Bloom*, and performed alongside Kate Beckinsale, Chloe Sevigny, Stephen Fry, and James Fleet in the critically-acclaimed film *Love and Friendship*, adapted from Jane Austen novel *Lady Susan*.</p>
<p>Thought: I found two films that Xavier Samuel starred in during 2016: *The Death and Life of Otto Bloom* and *Love and Friendship*. Now, I need to check if Lucy Fry appeared in either of these films.</p> <p>Action: search[Lucy Fry; What films did Lucy Fry star in during 2016?]</p> <p>Observation: (Result 1) Lucy Fry - In 2016, Lucy Fry starred in the films *Mr. Church* as Poppy and *The Darkness* as Stephanie Taylor.</p>
<p>Thought: Lucy Fry starred in *Mr. Church* and *The Darkness* in 2016. I need to check if either of these films also features Xavier Samuel.</p> <p>Action: search[Mr. Church; Did Xavier Samuel star in Mr. Church?]</p> <p>Observation: (Result 1) Mr. Church - Xavier Samuel is listed as one of the stars in the film *Mr. Church*.</p>
<p>Thought: I have confirmed that both Xavier Samuel and Lucy Fry starred in the film *Mr. Church* in 2016.</p> <p>Action: finish[*Mr. Church*]</p> <p>Observation: *Mr. Church*</p>

Table 18: The previous example with *NotesWriting* on the HotpotQA dataset using GPT-4o-mini.

<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 19: 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.

<p>Question: A general motors vehicle is named after the largest ward in the country of Monaco. How many people had walked on the moon as of the first model year of the vehicle? Note: the model year is not the same as the year the model was first produced.</p>
<p>Thought: I need to determine which General Motors vehicle is named after the largest ward in Monaco and find out the first model year of that vehicle. After that, I will find out how many people had walked on the moon by that year. Let's start by searching for the General Motors vehicle named after the largest ward in Monaco.</p> <p>Action: search[General Motors vehicle named after the largest ward in Monaco; What is the largest ward in Monaco and which GM vehicle is named after it?]</p> <p>Observation: No relevant information, try a different search term.</p>
<p>Thought: I need to find the largest ward in Monaco first and then see if there is a General Motors vehicle named after it. I will search for the largest ward in Monaco.</p> <p>Action: search[largest ward in Monaco; What is the largest ward in Monaco?]</p> <p>Observation: (Result 1) Monaco - The largest quarter in Monaco is Monte Carlo.</p>
<p>Thought: The largest ward in Monaco is Monte Carlo. Now, I need to search for any General Motors vehicle named after Monte Carlo and find out its first model year.</p> <p>Action: search[Monte Carlo; What is the General Motors vehicle named after Monte Carlo and what is its first model year?]</p> <p>Observation: (Result 1) Chevrolet Monte Carlo - The Chevrolet Monte Carlo is the General Motors vehicle named after Monte Carlo, and it was introduced for the 1970 model year.</p>
<p>Thought: The Chevrolet Monte Carlo is named after the largest ward in Monaco and was introduced for the 1970 model year. Now, I need to find out how many people had walked on the moon by the end of 1970.</p> <p>Action: search[moon landings; How many people had walked on the moon by the end of 1970?]</p> <p>Observation: (Result 1) Moon landing - A total of twelve astronauts have landed on the Moon. (Result 2) Apollo program - By the end of 1970, twelve people had walked on the Moon.</p>
<p>Thought: By the end of 1970, a total of twelve astronauts had walked on the Moon. Now I can provide the final answer regarding how many people had walked on the Moon as of the first model year of the Chevrolet Monte Carlo.</p> <p>Action: finish[12]</p> <p>Observation: 12</p>

Table 20: An example of hallucination by the note taking LLM in *NotesWriting* on the Frames dataset by GPT-4o-mini. The in-correct note taken is in bold. Instead the note writing LLM should have calculated the number of people walked before 1970 which is 4.

<p>Question: How old would the founder of the publishing company of the magazine that serialized the manga series Raw Hero have been the year the magazine ended publication?</p>
<p>Thought: I need to first identify the founder of the publishing company that serialized the manga series Raw Hero. Then I will determine the year the magazine ended publication and calculate the founder's age at that time.</p> <p>Action: search[Raw Hero; Which publishing company serialized the manga series Raw Hero?]</p> <p>Observation: (Result 1) Raw Hero - *Raw Hero* was serialized in Kodansha's *Evening*.</p>
<p>Thought: The manga series Raw Hero was serialized in Kodansha's *Evening* magazine. Now, I need to find out who the founder of Kodansha is.</p> <p>Action: search[Kodansha; Who is the founder of Kodansha?]</p> <p>Observation: (Result 1) Kodansha - Kodansha was founded by Seiji Noma in 1909.</p>
<p>Thought: Seiji Noma is the founder of Kodansha, which serialized *Raw Hero*. Next, I need to determine the year *Evening* magazine ended publication to calculate his age at that time.</p> <p>Action: search[Evening magazine; When did Evening magazine cease publication?]</p> <p>Observation: (Result 1) *Evening* (magazine) - The magazine ended publication on February 28, 2023. (Result 2) PM Magazine - *Evening Magazine* ceased publication on August 30, 1991.</p>
<p>Thought: The *Evening* magazine ceased publication on February 28, 2023. I will use the founding year of Kodansha (1909) to calculate the age of Seiji Noma at the time of the magazine's end in 2023.</p> <p>Action: finish[113]</p> <p>Observation: 113</p>

Table 21: An example of hallucination by the main LLM in *NotesWriting* on the Frames dataset by GPT-4o-mini. The LLM uses the founding year of the magazine as the birth year of the founder (highlighted in bold) instead of first identifying the birth year.

IRCoT and FLARE fewshot prompt (1/2)

Task: Given a question and relevant documents, generate a step-by-step reasoning process to solve it.

Instructions:

- Enclose each reasoning step within `<s>` `</s>` tags.
- Reason in short sentences as the sentences will be used as queries to retrieve relevant documents for subsequent steps.
- If provided with a partial reasoning chain, complete it by generating the next set of steps.
- If you do not have up-to-date knowledge, try to think of a plan to solve the question.
- Provide the final answer inside `<answer>` `</answer>` tags.

DOCUMENTS:

Nobody Loves You (When You're Down and Out) - "Nobody Loves You (When You're Down and Out)" is a song written by John Lennon released on his 1974 album "Walls and Bridges". The song is included on the 1986 compilation "Menlove Ave.", the 1990 boxset "Lennon", the 1998 boxset "John Lennon Anthology", the 2005 two-disc compilation "", and the 2010 boxset "Gimme Some Truth".

John Lennon/Plastic Ono Band - John Lennon/Plastic Ono Band is the debut studio album by English rock musician John Lennon. It was released in 1970, after Lennon had issued three experimental albums with Yoko Ono and "Live Peace in Toronto 1969", a live performance in Toronto credited to the Plastic Ono Band. The album was recorded simultaneously with Ono's debut avant garde solo album, "Yoko Ono/Plastic Ono Band", at Ascot Sound Studios and Abbey Road Studios using the same musicians and production team and nearly identical cover artwork.

Walls and Bridges - Walls and Bridges is the fifth studio album by English musician John Lennon. It was issued by Apple Records on 26 September 1974 in the United States and on 4 October in the United Kingdom. Written, recorded and released during his 18-month separation from Yoko Ono, the album captured Lennon in the midst of his "Lost Weekend". "Walls and Bridges" was an American "Billboard" number-one album and featured two hit singles, "Whatever Gets You thru the Night" and "#9 Dream". The first of these was Lennon's first number-one hit in the United States as a solo artist, and his only chart-topping single in either the US or Britain during his lifetime.

Question: Nobody Loves You was written by John Lennon and released on what album that was issued by Apple Records, and was written, recorded, and released during his 18 month separation from Yoko Ono?

Step-by-step reasoning:

`<s>`Identify album issued by Apple Records and recorded during John Lennon's 18-month separation from Yoko Ono.`</s>`

`<s>`The album "Walls and Bridges" was issued by Apple Records and recorded during this period.`</s>`

`<s>`Determine the album on which "Nobody Loves You (When You're Down and Out)" was released.`</s>`

`<s>`"Nobody Loves You (When You're Down and Out)" was released on "Walls and Bridges".`</s>`

`<answer>`Walls and Bridges`</answer>`

DOCUMENTS:

Route 13 (Laos) - Route 13 is the most important highway in the country of Laos. It begins at Boten in the North of Laos at the Chinese border. It connects the city of Vientiane to Luang Prabang in the north and roughly follows the line of the Mekong River down to the border with Cambodia. The road then continues at National Highway 7 in Cambodia. Also, this national highway leads past all three international airports in Laos: Vientiane Airport, Luang Prabang Airport, and Pakse Airport. Between Boten and Nateuy, Route 13 is at a length of 20 km as part of the Kunming-Bangkok Expressway, which is also part of the Asian Highway AH3. From Nateuy to Vientiane, Route 13 is part of AH12, and from Vientiane to the border of Cambodia, it is part of AH11.

National Route 13 (Vietnam) - National Route 13 is a highway in southern Vietnam stretching from the northeastern outskirts of Ho Chi Minh City, the commercial center of the country, towards the border to Cambodia. The highway starts around on the northern outskirts of Ho Chi Minh City, once the site of the military academy of the Army of the Republic of Vietnam, and travels north through the provinces.

Cambodia - Cambodia (Khmer: or Kampuchea), officially known as the Kingdom of Cambodia, is a country located in the southern portion of the Indochina Peninsula in Southeast Asia. It is 181,035 km² in area, bordered by Thailand to the northwest, Laos to the northeast, Vietnam to the east, and the Gulf of Thailand to the southwest.

Figure 5: Fewshot prompt used for the evaluation of IRCoT and FLARE methods (1/2).

IRCoT and FLARE fewshot prompt (2/2)

Question: What is known as the Kingdom and has National Route 13 stretching towards its border?

Step-by-step reasoning:

<s>Cambodia is officially known as the Kingdom of Cambodia.</s>

<s>National Route 13 in Vietnam stretches towards the border of Cambodia.</s>

<s>Since Cambodia is referred to as the Kingdom and National Route 13 leads to its border, Cambodia is the correct answer.</s>

<answer>Cambodia</answer>

DOCUMENTS:

Vandenbrink GTO - The Vandenbrink GTO is a limited re-bodied version of the Ferrari 599 GTB Fiorano. It features an entirely new coachwork designed by Michiel van den Brink of Vandenbrink Design, inspired by the 1962 Ferrari 250 GTO as a tribute.

1969 World Sportscar Championship - The 1969 World Sportscar Championship was the 17th season of FIA World Sportscar racing, featuring Group 6 Prototype Sports Cars, Group 4 Sports Cars, and Group 3 Grand Touring Cars.

Ferrari 250 GTO - The Ferrari 250 GTO was a GT car produced by Ferrari from 1962 to 1964 for homologation into the FIA's Group 3 Grand Touring Car category. It was powered by Ferrari's Tipo 168/62 V12 engine.

Question: Which car, produced by Ferrari from 1962 to 1964 for homologation into the FIA's Group 3 Grand Touring Car category, inspired the Vandenbrink GTO?

Step-by-step reasoning:

<s>The Ferrari 250 GTO was produced by Ferrari from 1962 to 1964 for homologation into the FIA's Group 3 Grand Touring Car category.</s>

<s>The Vandenbrink GTO was inspired by the Ferrari 250 GTO.</s>

<s>Since the question asks for the car that inspired the Vandenbrink GTO, the correct answer is the Ferrari 250 GTO.</s>

<answer>Ferrari 250 GTO</answer>

DOCUMENTS:

Talula - "Talula" is a song by Tori Amos, released as the second single from her 1996 album "Boys For Pele". It reached #22 on the UK Singles Chart and appears in the Jan de Bont film "Twister".

Nema aviona za Zagreb - "Nema aviona za Zagreb" ("There Is No Plane to Zagreb") is a 2012 Dutch film by Louis van Gasteren. The film is a retrospective of events in the director's life from 1964 to 1969, filmed by him in that period and reflected on from his vantage point over 40 years later at the age of 90. It is also Jan de Bont's first film credit as cinematographer since 1992's "Basic Instinct".

Jan de Bont - Jan de Bont (born 22 October 1943) is a Dutch cinematographer, director, and film producer. He is widely known for directing the 1994 action film "Speed," starring Keanu Reeves and Sandra Bullock, and the disaster film "Twister." As a director of photography, he has shot numerous blockbusters and genre films, including "Cujo, Flesh and Blood, Die Hard, The Hunt for Red October," and "Basic Instinct."

Question: Who was born first? Jan de Bont or Raoul Walsh?

Step-by-step reasoning:

<s>Jan de Bont was born on October 22, 1943.</s>

<s>Raoul Walsh was born on March 11, 1887.</s>

<s>Since 1887 is earlier than 1943, Raoul Walsh was born first.</s>

<answer>Raoul Walsh</answer>

Figure 6: Fewshot prompt used for the evaluation of IRCoT and FLARE methods (2/2).

GPT-4 Judge Prompt

===Task===

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.

===Instructions===

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."

===Input Data===

- **Question:** «question»
- **Predicted Answer:** «LLM_response»
- **Ground Truth Answer:** «ground_truth_answer»

===Output Format===

Provide your final evaluation in the following format:

Explanation: (How you made the decision?)

Decision: ("TRUE" or "FALSE")

Please proceed with the evaluation.

Figure 7: GPT-4 prompt for evaluating the correctness of predicted answer.

Quality evaluation prompt

You are asked to evaluate the reasoning chain produced in response to a question, particularly focusing on how effectively tools were used throughout the process. The evaluation should be based on the following clearly defined criteria. For each criterion, provide a numerical rating on a scale from 0 to 5, where 5 represents excellent performance and 0 indicates poor or entirely absent performance.

Criterion 1: Efficiency of the Steps Taken

Definition:

Evaluate the overall efficiency of each step in the reasoning chain, with specific focus on whether the tool calls and reasoning steps helped progress toward the final correct answer. Efficient steps reduce uncertainty, narrow the solution space, or directly contribute to solving the problem.

Rating Guide:

- 5 – Extremely efficient: Every step clearly advances the reasoning; no wasted effort.
- 4 – Highly efficient: Most steps are purposeful, with only minor inefficiencies.
- 3 – Moderately efficient: Some steps are valuable, others contribute little.
- 2 – Minimally efficient: Several steps are misdirected or low-impact.
- 1 – Poorly efficient: Most steps offer minimal or no progress toward the answer.
- 0 – Not efficient at all: Steps are irrelevant, aimless, or distracting.

Criterion 2: Redundancy of Steps

Definition:

Assess the reasoning chain for unnecessary repetition or duplication of steps, including redundant tool calls or rephrasing of the same logic without new insight. A low-redundancy chain avoids rework and keeps the progression streamlined.

Rating Guide:

- 5 – No redundancy: Each step is unique and adds distinct value.
- 4 – Very low redundancy: Only minor repetition, quickly resolved.
- 3 – Moderate redundancy: Some ideas or tool uses are repeated without added benefit.
- 2 – Noticeable redundancy: Multiple steps repeat similar content or actions unnecessarily.
- 1 – High redundancy: Repetition significantly detracts from conciseness.
- 0 – Extremely redundant: Most of the chain rehashes prior reasoning with no new value.

Criterion 3: Clarity and Coherence of the Reasoning Chain

Definition:

Examine how clearly and logically the reasoning chain progresses from the question to the final answer. This includes whether steps are easy to follow, logically connected, and free of ambiguity or excessive complexity.

Rating Guide:

- 5 – Exceptionally clear and coherent: The reasoning is logical, concise, and easy to follow.
- 4 – Mostly clear: The chain is understandable with minor clarity issues.
- 3 – Moderately clear: Some transitions or justifications are unclear or weak.
- 2 – Confusing in parts: Multiple unclear, inconsistent, or disjointed steps.
- 1 – Difficult to follow: Lacks logical flow or clear structure.
- 0 – Incomprehensible: The chain cannot be understood or followed logically.

First provide your reasoning of your evaluation then structure your responses as a json with the keys "Criterion 1", "Criterion 2", "Criterion 3" and the values as the ratings you provided.

Chain: {}

Figure 8: Prompt for quality evaluation of reasoning chain.