

Search-o1: Agentic Search-Enhanced Large Reasoning Models

Xiaoxi Li¹, Guanting Dong¹, Jiajie Jin¹, Yuyao Zhang¹, Yujia Zhou²,
Yutao Zhu¹, Peitian Zhang¹, Zhicheng Dou^{1*}
¹Renmin University of China ²Tsinghua University
{xiaoxi_li, dou}@ruc.edu.cn

Project Page: <https://search-o1.github.io/>

Abstract

Large reasoning models (LRMs) like OpenAI-o1 have demonstrated impressive long stepwise reasoning capabilities through large-scale reinforcement learning. However, their extended reasoning processes often suffer from knowledge insufficiency, leading to frequent uncertainties and potential errors. To address this limitation, we introduce **Search-o1**, a framework that enhances LRMs with an agentic retrieval-augmented generation (RAG) mechanism and a Reason-in-Documents module for refining retrieved documents. Search-o1 integrates an agentic search workflow into the reasoning process, enabling dynamic retrieval of external knowledge when LRMs encounter uncertain knowledge points. Additionally, due to the verbose nature of retrieved documents, we design a separate Reason-in-Documents module to deeply analyze the retrieved information before injecting it into the reasoning chain, minimizing noise and preserving coherent reasoning flow. Extensive experiments on complex reasoning tasks in science, mathematics, and coding, as well as six open-domain QA benchmarks, demonstrate the strong performance of Search-o1. This approach enhances the trustworthiness and applicability of LRMs in complex reasoning tasks, paving the way for more reliable and versatile intelligent systems. The code is available at <https://github.com/sunnynexus/Search-o1>.

1 Introduction

Recently emerged large reasoning models (LRMs), exemplified by OpenAI’s o1 [22], Qwen-QwQ [54] and DeepSeek-R1 [7], employ large-scale reinforcement learning foster impressive long-sequence stepwise reasoning capabilities, offering promising solutions to complex reasoning problems [46, 31, 59, 84, 73, 74, 67]. This advancement has inspired a series of foundational efforts aimed at exploring and reproducing o1-like reasoning patterns, to broaden their application to a wider range of foundational models [49, 19, 77, 80, 71, 25, 45].

It is noteworthy that o1-like reasoning patterns guide LRMs to engage in a slower thinking process [6, 61] by implicitly breaking down complex problems, generating a long internal reasoning chain and then discovering suitable solutions step by step. While this characteristic enhances logical coherence and interpretability of reasoning, an extended chain of thought may cause overthinking [4] and increased risks of knowledge insufficiency [60, 51, 2], where any knowledge gap can propagate errors and disrupt the entire reasoning chain [79, 40, 44, 41].

*Corresponding author.

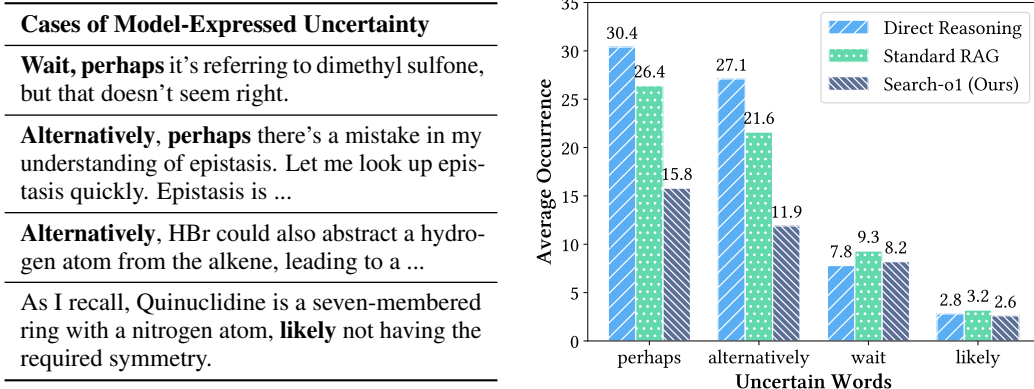


Figure 1: Analysis of reasoning uncertainty with QwQ-32B-Preview. **Left:** Examples of uncertain words identified during the reasoning process. **Right:** Average occurrence of high-frequency uncertain words per output in the GPQA diamond set.

To address this limitation, we conduct preliminary experiments to assess the frequency of uncertain words decoded by the LRMs due to knowledge gaps. As shown in Figure 1, the extended thinking process leads LRM to frequently decode numerous uncertain terms in challenging reasoning problems, with “*perhaps*” averaging over 30 occurrences in each reasoning process. Notably, the high specialization of these problems also complicates manual reasoning verification, often incurring significant costs [63]. Consequently, automating the supplementation of knowledge required for the o1-like reasoning process has become a significant challenge, limiting the progress of LRMs in achieving universally trustworthy reasoning.

To shed light on this topic, our core motivation is to enhance the LRMs with o1-like reasoning pattern through autonomous retrieval. We propose **Search-o1**, which integrates the reasoning process of LRMs with two core components: an agentic retrieval-augmented generation (RAG) mechanism and a knowledge refinement module. This design aims to enable LRMs to incorporate the agentic search workflow into the reasoning process, retrieving external knowledge on demand to support step-wise reasoning while preserving coherence throughout.

Specifically, our results in Figure 1 reveal that traditional problem-oriented RAG techniques do not effectively address the knowledge gaps compared to direct reasoning (Standard RAG vs. Direct Reasoning). This finding aligns with human intuition, as standard RAG retrieves relevant knowledge only once in a problem-oriented manner, while the knowledge required for each step in complex reasoning scenarios is often varied and diverse [83, 41, 11]. Unlike them, **Search-o1** employs an agentic RAG technique that guides the model to actively decode search queries when facing knowledge shortages, thereby triggering the retrieval mechanism to obtain relevant external knowledge. Owing to the benefits of this design, our retrieval mechanism can be triggered and iterated multiple times within a single reasoning session to fulfill the knowledge needs of various reasoning steps.

To effectively integrate retrieved knowledge into the LRM’s reasoning process, we further identify two key challenges when directly incorporating retrieved documents into the reasoning chain during practical experiments: (1) **Redundant Information in Retrieved Documents**. Retrieved documents are often lengthy and contain redundant information, directly inputting them into LRMs may disrupt the original coherence of reasoning and even introduce noise [62, 72, 26]. (2) **Limited Ability to Understand Long Documents**. Most LRMs have been specifically aligned for complex reasoning tasks during the pre-training and fine-tuning stages. This focus has resulted in a degree of catastrophic forgetting in their general capabilities [39, 10], ultimately limiting their long-context understanding of retrieved documents.

To address these challenges, we introduce the Reason-in-Documents module, which operates independently from the main reasoning chain. This module first conducts a thorough analysis of retrieved documents based on both the current search query and previous reasoning steps, and then produces refined information that seamlessly integrates with the prior reasoning chain.

In summary, our contributions are as follows:

- We propose Search-o1, the first framework that integrates the agentic search workflow into the o1-like reasoning process of LRM for achieving autonomous knowledge supplementation.
- To effectively integrate external knowledge during reasoning, Search-o1 combines the reasoning process with an agentic RAG mechanism and a knowledge refinement module. This design enables the LRM to retrieve external knowledge on demand, seamlessly incorporating it into the reasoning chain while maintaining the original logical flow.
- With five complex reasoning domains and six open-domain QA benchmarks, we demonstrate that Search-o1 achieves remarkable performance in the reasoning field while maintaining substantial improvements in the general knowledge. Further quantitative analysis confirms its efficiency and scalability, offering practical guidance for trustworthy reasoning in LRMs.

2 Related Work

Large Reasoning Models. Large reasoning models focus on enhancing performance at test time by utilizing extended reasoning steps, contrasting with traditional large pre-trained models that achieve scalability during training by increasing model size or expanding training data [17, 66, 50, 85, 76]. Studies have shown that test-time scaling can improve the reasoning abilities of smaller models on complex tasks [15, 75]. Recently, models like OpenAI-o1 [22], Qwen-QwQ [54] and DeepSeek-R1 [7] explicitly demonstrate chain-of-thought reasoning [59], mimicking human problem-solving approaches in domains such as mathematics, coding, and so on.

Various approaches have been explored to achieve o1-like reasoning capabilities. Some methods combine policy and reward models with Monte Carlo Tree Search (MCTS) [25], though this does not internalize reasoning within the model. Other studies incorporate deliberate errors in reasoning paths during training to partially internalize these abilities [49, 71]. Additionally, distilling training data has been shown to enhance models’ o1-like reasoning skills [45]. The o1-like reasoning paradigm has demonstrated strong performance across diverse domains, including vision-language reasoning [65, 11, 48, 69], code generation [81, 32], healthcare [3], and machine translation [57]. However, these approaches are limited by their reliance on static, parameterized models, which cannot leverage external world knowledge when internal knowledge is insufficient.

Retrieval-Augmented Generation. Retrieval-augmented generation (RAG) introduces retrieval mechanisms to address the limitations of static parameters in generative models, allowing access to external knowledge to solve more complex problems [30, 82, 35, 86]. Advanced research in this field enhances the RAG system from multiple aspects, including the necessity of retrieval [53], pre-processing of queries [43, 58], retrieved documents compressing [64], denoising [42, 12], refining [24, 27, 88], instruction following [9, 8, 87] and so on. Furthermore, some studies have explored end-to-end model training to implement RAG systems [1, 36, 33, 34] and knowledge-graph-based RAG systems [14, 37].

Recently, agentic RAG systems empower models to autonomously determine when and what knowledge to retrieve as needed, showcasing enhanced planning and problem-solving capabilities [5, 56, 70]. There is also research combining agent-based systems with MCTS to optimize complex workflows, leveraging retrievers and other tools to accomplish tasks [78]. However, existing RAG approaches have not combined the strong reasoning capabilities of o1-like models, limiting the potential to further enhance system performance in solving complex tasks.

3 Methodology

3.1 Problem Formulation

We consider a complex reasoning task that necessitates multi-step reasoning and the retrieval of external knowledge to derive solutions. The objective is to generate a comprehensive solution for each question q , consisting of both a logical reasoning chain \mathcal{R} and the final answer a . In this work, we enable the reasoning model to utilize external knowledge sources during the reasoning process. Specifically, we consider three primary inputs in the problem-solving process: the task instruction I , the question q , and externally retrieved documents \mathcal{D} . Here, I provides an overarching description of

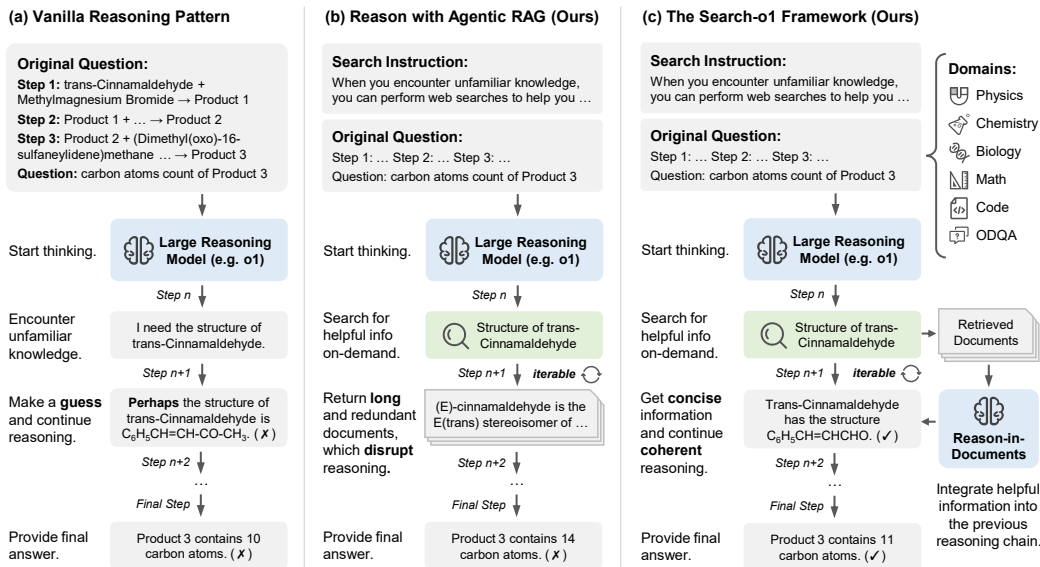


Figure 2: Comparison of reasoning approaches: (a) Direct reasoning without retrieval often results in inaccuracies due to missing knowledge. (b) Our agentic retrieval-augmented reasoning approach improves knowledge access but usually returns lengthy, redundant documents, disrupting coherent reasoning. (c) Our Search-o1 integrates concise and accurate retrieved knowledge seamlessly into the reasoning process, enabling precise and coherent problem-solving.

the reasoning task, q is the specific complex question to be answered, and \mathcal{D} comprises background knowledge dynamically retrieved from a relevant knowledge base.

The goal is to design a reasoning mechanism that effectively integrates I , q , and \mathcal{D} to produce a coherent reasoning chain \mathcal{R} and a final answer a . This can be formalized as the mapping $(I, q, \mathcal{D}) \rightarrow (\mathcal{R}, a)$. The generation of the reasoning sequence and the final answer can be expressed as:

$$P(\mathcal{R}, a \mid I, q, \mathcal{D}) = \underbrace{\prod_{t=1}^{T_r} P(\mathcal{R}_t \mid \mathcal{R}_{<t}, I, q, \mathcal{D}_{<t})}_{\text{Reasoning Process}} \cdot \underbrace{\prod_{t=1}^{T_a} P(a_t \mid a_{<t}, \mathcal{R}, I, q)}_{\text{Answer Generation}}, \quad (1)$$

where T_r is the number of tokens in the reasoning sequence \mathcal{R} . The token at the position t is \mathcal{R}_t , and $\mathcal{R}_{<t}$ represents all tokens generated before position t . $\mathcal{D}_{<t}$ represents all documents retrieved up to token t in the reasoning chain. Similarly, T_a is the length of the answer sequence a , with a_t being the token at the position t and $a_{<t}$ indicating all generated answer tokens before the position t .

3.2 Overview of the Search-o1 Framework

The Search-o1 framework addresses knowledge insufficiency in large reasoning models (LRMs) by seamlessly integrating external knowledge retrieval into their reasoning process while maintaining chain-of-thought coherence. As illustrated in Figure 2, we present a comparative analysis of three approaches: vanilla reasoning, agentic retrieval-augmented generation (RAG), and our proposed Search-o1 framework.

- **Vanilla Reasoning Pattern:** Consider the example in Figure 2(a), where the task involves determining the carbon atom count in the final product of a three-step chemical reaction. The vanilla reasoning approach falters when encountering knowledge gaps (e.g., the “structure of trans-Cinnamaldehyde”). Without access to accurate information, the model must rely on assumptions, potentially leading to cascading errors throughout subsequent reasoning steps.
- **Agentic RAG:** To bridge the knowledge gaps during reasoning, we build the agentic RAG mechanism (Figure 2(b)) to enable the model to autonomously retrieve external knowledge when needed. When uncertainty arises—such as regarding the compound’s structure—the model generates targeted search queries (e.g., “structure of trans-Cinnamaldehyde”). However, the direct insertion

of retrieved documents, which often contain lengthy and irrelevant information, may disrupt the reasoning flow and hurt coherence.

- **Search-o1:** Our Search-o1 framework (Figure 2(c)) extends the agentic RAG mechanism by incorporating a Reason-in-Documents module. This module condenses retrieved documents into focused reasoning steps that integrate external knowledge while maintaining the logical flow of the reasoning chain. It considers the current search query, retrieved documents, and the existing reasoning chain to generate coherent steps. This iterative process continues until the final answer is reached. The following sections provide detailed explanations of agentic RAG, Reason-in-Documents, and the Search-o1 inference process.

3.3 Agentic Retrieval-Augmented Generation Mechanism

The agentic RAG mechanism is a pivotal component of the Search-o1 framework, empowering the reasoning model to autonomously determine when to retrieve external knowledge during the reasoning process. This mechanism allows the model itself to decide whether to continue generating reasoning steps or to initiate a retrieval step. Detailed model instructions can be found in Appendix A.1.

During the generation of the reasoning chain \mathcal{R} , the model may intermittently generate search queries $q_{\text{search}}^{(i)}$ encapsulated between special symbols `<begin_search_query>` and `<end_search_query>`, where i indexes the i -th search step. Each search query is generated based on the current state of the reasoning process and the previously retrieved knowledge. The generation of each search query is expressed as:

$$P(q_{\text{search}}^{(i)} \mid I, q, \mathcal{R}^{(i-1)}) = \prod_{t=1}^{T_q^{(i)}} P\left(q_{\text{search},t}^{(i)} \mid q_{\text{search},<t}^{(i)}, I, q, \mathcal{R}^{(i-1)}\right), \quad (2)$$

where $T_q^{(i)}$ is the length of the i -th search query, $q_{\text{search},t}^{(i)}$ denotes the token generated at step t of the i -th search query, and $\mathcal{R}^{(i-1)}$ represents all the reasoning steps prior to the i -th search step, including both search queries and search results.

Once a new pair of special symbols for the search query is detected in the reasoning sequence, we pause the reasoning process, and the search query $q_{\text{search}}^{(i)}$ is extracted. The retrieval function `Search` is invoked to obtain relevant documents:

$$\mathcal{D}^{(i)} = \text{Search}(q_{\text{search}}^{(i)}), \quad (3)$$

where $\mathcal{D}^{(i)} = d_1^{(i)}, d_2^{(i)}, \dots, d_{k_i}^{(i)}$ represents the set of top- k_i relevant documents retrieved for the i -th search query. The retrieved documents $\mathcal{D}^{(i)}$ are subsequently injected into the reasoning chain $\mathcal{R}^{(i-1)}$ between the special symbols `<begin_search_result>` and `<end_search_result>`, allowing the reasoning model to utilize the external knowledge to continue the reasoning process.

This agentic mechanism enables the model to dynamically and efficiently incorporate external knowledge, maintaining the coherence and relevance of the reasoning process while avoiding information overload from excessive or irrelevant retrieval results.

3.4 Knowledge Refinement via Reason-in-Documents

While the agentic RAG mechanism addresses knowledge gaps in reasoning, directly inserting full documents can disrupt coherence due to their length and redundancy. To overcome this, the Search-o1 framework includes the knowledge refinement module, which selectively integrates only relevant and concise information into the reasoning chain through a separate generation process using the original reasoning model. This module processes retrieved documents to align with the model’s specific reasoning needs, transforming raw information into refined, pertinent knowledge while maintaining coherence and logical consistency of the main reasoning chain.

The refinement guidelines for Reason-in-Documents are detailed in Appendix A.1. These guidelines instruct the model to analyze the retrieved web pages based on the previous reasoning steps, current search query, and the content of the searched web pages. The objective is to extract relevant and accurate information that directly contributes to advancing the reasoning process for the original question, ensuring seamless integration into the existing reasoning chain.

Algorithm 1 Search-o1 Inference

Require: Reasoning Model \mathcal{M} , Search function Search

```

1: Input: Questions  $\mathcal{Q}$ , Task instruction  $I$ , Reason-in-documents instruction  $I_{\text{docs}}$ 
2: Initialize set of unfinished sequences  $\mathcal{S} \leftarrow \{I \oplus q \mid q \in \mathcal{Q}\}$ 
3: Initialize set of finished sequences  $\mathcal{F} \leftarrow \{\}$ 
4: while  $\mathcal{S} \neq \emptyset$  do
5:   Generate all sequences in  $\mathcal{S}$  until EOS or <end_search_query>:  $\mathcal{T} \leftarrow \mathcal{M}(\mathcal{S})$   $\triangleright$  Batch Generate
6:   Initialize empty set  $\mathcal{S}_r \leftarrow \{\}$   $\triangleright$  Reason-in-documents Inputs
7:   for each sequence  $\text{Seq} \in \mathcal{T}$  do
8:     if  $\text{Seq}$  ends with <end_search_query> then
9:       Extract search query:  $q_{\text{search}} \leftarrow \text{Extract}(\text{Seq}, \text{<begin\_search\_query>}, \text{<end\_search\_query>})$ 
10:      Retrieve documents:  $\mathcal{D} \leftarrow \text{Search}(q_{\text{search}})$   $\triangleright$  Retrieval
11:      Construct input for Reason-in-documents:  $I_{\mathcal{D}} \leftarrow I_{\text{docs}} \oplus q_{\text{search}} \oplus \text{Seq}$ 
12:      Append the tuple  $(I_{\mathcal{D}}, \text{Seq})$  to  $\mathcal{S}_r$ 
13:     else if  $\text{Seq}$  ends with EOS then
14:       Remove  $\text{Seq}$  from  $\mathcal{S}$ , add  $\text{Seq}$  to  $\mathcal{F}$   $\triangleright$  Sequence Finished
15:   if  $\mathcal{S}_r \neq \emptyset$  then
16:     Prepare batch inputs:  $\mathcal{I}_r \leftarrow \{I_{\mathcal{D}} \mid (I_{\mathcal{D}}, \text{Seq}) \in \mathcal{S}_r\}$ 
17:     Reason-in-documents:  $\mathcal{T}_r \leftarrow \mathcal{M}(\mathcal{I}_r)$   $\triangleright$  Batch Generate
18:     for  $i \leftarrow \{1, \dots, |\mathcal{T}_r|\}$  do
19:       Let  $r \leftarrow \mathcal{T}_r[i]$ ,  $\text{Seq} \leftarrow \mathcal{S}_r[i].\text{Seq}$ 
20:       Extract knowledge-injected reasoning step:  $r_{\text{final}} \leftarrow \text{Extract}(r)$ 
21:       Update sequence in  $\mathcal{S}$ :  $\text{Seq} \leftarrow \text{Insert}(\text{<begin\_search\_result>}, r_{\text{final}}, \text{<end\_search\_result>})$ 
22: Output: Finished Sequences  $\mathcal{F}$ 
  
```

For each search step i , let $\mathcal{R}^{(<i)}$ denote the reasoning chain accumulated up to just before the i -th search query. Given $\mathcal{R}^{(<i)}$, the current search query $q_{\text{search}}^{(i)}$, and the retrieved documents $\mathcal{D}^{(i)}$, the knowledge refinement process operates in two stages: **first generating an intermediate reasoning sequence $r_{\text{docs}}^{(i)}$ to analyze the retrieved documents, then producing refined knowledge $r_{\text{final}}^{(i)}$ based on this analysis.** The generation of the intermediate reasoning sequence $r_{\text{docs}}^{(i)}$ is expressed as:

$$P(r_{\text{docs}}^{(i)} \mid \mathcal{R}^{(<i)}, q_{\text{search}}^{(i)}, \mathcal{D}^{(i)}) = \prod_{t=1}^{T_d^{(i)}} P(r_{\text{docs},t}^{(i)} \mid r_{\text{docs},<t}^{(i)}, \mathcal{R}^{(<i)}, q_{\text{search}}^{(i)}, \mathcal{D}^{(i)}), \quad (4)$$

where $T_d^{(i)}$ is the length of the intermediate reasoning sequence, and $r_{\text{docs},t}^{(i)}$ denotes the token at step t . The refined knowledge $r_{\text{final}}^{(i)}$ is then generated based on this analysis:

$$P(r_{\text{final}}^{(i)} \mid r_{\text{docs}}^{(i)}, \mathcal{R}^{(<i)}, q_{\text{search}}^{(i)}) = \prod_{t=1}^{T_r^{(i)}} P(r_{\text{final},t}^{(i)} \mid r_{\text{final},<t}^{(i)}, r_{\text{docs}}^{(i)}, \mathcal{R}^{(<i)}, q_{\text{search}}^{(i)}), \quad (5)$$

where $T_r^{(i)}$ is the length of the refined knowledge sequence, and $r_{\text{final},t}^{(i)}$ denotes the token at step t . The refined knowledge $r_{\text{final}}^{(i)}$ is then incorporated into the reasoning chain $\mathcal{R}^{(i)}$, enabling the model to continue generating coherent reasoning steps with access to the external knowledge.

$$P(\mathcal{R}, a \mid I, q) = \prod_{t=1}^{T_r} P(\mathcal{R}_t \mid \mathcal{R}_{<t}, I, q, \{r_{\text{final}}^{(j)}\}_{j \leq i(t)}) \cdot \prod_{t=1}^{T_a} P(a_t \mid a_{<t}, \mathcal{R}, I, q), \quad (6)$$

where $\{r_{\text{final}}^{(j)}\}_{j \leq i(t)}$ denotes all previously refined knowledge up to the $i(t)$ -th search step. Here, $i(t)$ represents the index of the search step corresponding to the current reasoning step t . **This refined knowledge integration ensures that each reasoning step can access relevant external information while maintaining the conciseness and focus of the reasoning process.**

3.5 Search-o1 Inference Process

Inference Logic for a Single Question. For each question, the Search-o1 inference begins by initializing the reasoning sequence with the task instruction I concatenated with the specific question

q . As the reasoning model \mathcal{M} generates the reasoning chain \mathcal{R} , it may produce search queries encapsulated within the special symbols `<begin_search_query>` and `<end_search_query>`. Upon detecting the `<end_search_query>` symbol, the corresponding search query q_{search} is extracted, triggering the retrieval function `Search` to obtain relevant external documents \mathcal{D} . These retrieved documents, along with the reason-in-documents instruction I_{docs} and the current reasoning sequence \mathcal{R} , are then processed by the Reason-in-Documents module. This module refines the raw documents into concise, pertinent information r_{final} , which is seamlessly integrated back into the reasoning chain \mathcal{R} within symbols `<begin_search_result>` and `<end_search_result>`. This iterative process ensures that the reasoning model incorporates necessary external knowledge while maintaining coherence and logical consistency, leading to the generation of a comprehensive reasoning chain \mathcal{R} and the final answer a .

Batch Inference Mechanism. To efficiently handle multiple questions simultaneously, the Search-ol framework employs a batch inference mechanism that optimizes both token generation and knowledge refinement. Initially, a set of unfinished reasoning sequences \mathcal{S} is created by concatenating the task instruction I with each question q in the batch \mathcal{Q} . The reasoning model \mathcal{M} then generates tokens for all sequences in \mathcal{S} in parallel, advancing each reasoning chain until it either completes or requires external knowledge retrieval. When a search query is identified within any sequence, the corresponding queries are extracted and processed in batches through the `Search` function to retrieve relevant documents \mathcal{D} . These documents are then collectively refined by the Reason-in-Documents module, which generates the refined knowledge r_{final} for each sequence. The refined knowledge is subsequently inserted back into the respective reasoning chains. Completed sequences are moved to the finished set \mathcal{F} , while ongoing sequences remain in \mathcal{S} for further processing. By leveraging parallel processing for both generation and refinement steps, the batch inference mechanism enhances system throughput associated with handling multiple inputs concurrently.

4 Experiments

4.1 Tasks and Datasets

The evaluations used in this experiment include the following two categories:

Challenging reasoning tasks: (1) **GPQA** [52] is a PhD-level science multiple-choice QA dataset. The questions are authored by domain experts in physics, chemistry, and biology. In our main experiments, we use the highest quality diamond set containing 198 questions, and in Table 2, we use a more comprehensive extended set containing 546 questions to compare with the performance of human experts. (2) **Math benchmarks** include **MATH500** [38], **AMC2023**², and **AIME2024**³. **MATH500** consists of 500 questions from the MATH test set [16]. **AMC2023** and **AIME2024** are middle school math competitions covering arithmetic, algebra, geometry, etc., containing 40 and 30 questions respectively. Among these three datasets, **MATH500** and **AMC** are relatively simple, while **AIME** is more difficult. (3) **LiveCodeBench** [23] is a benchmark for evaluating LLMs’ coding capabilities, consisting of easy, medium, and hard difficulty problems. It collects recently published programming problems from competitive platforms to avoid data contamination. We utilize problems from August to November 2024, comprising 112 problems.

Open-domain QA tasks: (1) **Single-hop QA datasets:** **Natural Questions (NQ)** [29] contains questions from real Google search queries with answers from Wikipedia articles. **TriviaQA** [28] is a large-scale dataset with questions from trivia websites and competitions, featuring complex entity relationships. (2) **Multi-hop QA datasets:** **HotpotQA** [68] is the first large-scale dataset requiring reasoning across multiple Wikipedia paragraphs. **2WikiMultiHopQA (2WIKI)** [18] provides explicit reasoning paths for multi-hop questions. **MuSiQue** [55] features 2-4 hop questions built from five existing single-hop datasets. **Bamboogle** [47] collects complex questions that Google answers incorrectly to evaluate models’ compositional reasoning across various domains.

4.2 Baselines

We evaluate our approach against the following baseline methods:

²<https://huggingface.co/datasets/AI-MO/aimo-validation-amc>

³<https://huggingface.co/datasets/AI-MO/aimo-validation-aime>

Table 1: Main results on challenging reasoning tasks, including PhD-level science QA, math, and code benchmarks. We report Pass@1 metric for all tasks. For models with 32B parameters, the best results are in **bold** and the second-best are underlined. Results from larger or non-proprietary models are in gray color for reference. Symbol “†” indicates results from their official releases.

Method	GPQA (PhD-Level Science QA)				Math Benchmarks			LiveCodeBench			
	Physics	Chemistry	Biology	Overall	MATH500	AMC23	AIME24	Easy	Medium	Hard	Overall
<i>Direct Reasoning (w/o Retrieval)</i>											
Qwen2.5-32B	57.0	33.3	52.6	45.5	75.8	57.5	23.3	42.3	18.9	<u>14.3</u>	22.3
Qwen2.5-Coder-32B	37.2	25.8	57.9	33.8	71.2	67.5	20.0	<u>61.5</u>	16.2	12.2	25.0
QwQ-32B	75.6	39.8	68.4	58.1	83.2	<u>82.5</u>	<u>53.3</u>	<u>61.5</u>	<u>29.7</u>	20.4	33.0
Qwen2.5-72B	57.0	37.6	68.4	49.0	79.4	67.5	20.0	53.8	29.7	24.5	33.0
Llama3.3-70B	54.7	31.2	52.6	43.4	70.8	47.5	36.7	57.7	32.4	24.5	34.8
DeepSeek-R1-Lite†	-	-	-	58.5	91.6	-	52.5	-	-	-	51.6
GPT-4o†	59.5	40.2	61.6	50.6	60.3	-	9.3	-	-	-	33.4
o1-preview†	89.4	59.9	65.9	73.3	85.5	-	44.6	-	-	-	53.6
<i>Retrieval-augmented Reasoning</i>											
RAG-Qwen2.5-32B	57.0	37.6	52.6	47.5	82.6	72.5	30.0	<u>61.5</u>	24.3	8.2	25.9
RAG-QwQ-32B	<u>76.7</u>	38.7	<u>73.7</u>	58.6	84.8	<u>82.5</u>	50.0	57.7	16.2	12.2	24.1
RAgent-Qwen2.5-32B	58.1	33.3	63.2	47.0	74.8	65.0	20.0	57.7	24.3	6.1	24.1
RAgent-QwQ-32B	<u>76.7</u>	<u>46.2</u>	68.4	<u>61.6</u>	<u>85.0</u>	85.0	56.7	65.4	18.9	12.2	<u>26.8</u>
<i>Retrieval-augmented Reasoning with Reason-in-Documents</i>											
Search-o1 (Ours)	77.9	47.3	78.9	63.6	86.4	85.0	56.7	57.7	32.4	20.4	33.0

Direct Reasoning: These methods utilize the model’s internal knowledge without retrieval. The open-source models include Qwen2.5-32B-Instruct [50], Qwen2.5-Coder-32B-Instruct [20], QwQ-32B-Preview [54], Qwen2.5-72B-Instruct [50], and Llama3.3-70B-Instruct [13]. Closed-source non-proprietary models include DeepSeek-R1-Lite-Preview [7], OpenAI GPT-4o [21], and o1-preview [22]. Results for open-source models are based on our implementations, while closed-source model results are sourced from their official releases.

Retrieval-augmented Reasoning: These methods retrieve external information to enhance the reasoning process. We consider two retrieval augmentation approaches: (1) **Standard RAG:** Retrieves the top-10 documents for the original question and inputs them alongside the question into the model for reasoning and answer generation. (2) **RAG Agent (RAgent):** Allows the model to decide when to generate queries for retrieval, as detailed in Section 3.3. To manage the length of retrieved documents, inspired by ReAct [70], we first retrieve the top-10 snippets during reasoning. The model then decides which URLs to obtain for the full documents when necessary.

4.3 Implementation Details

For the backbone large reasoning model in Search-o1, we utilize the open-sourced QwQ-32B-Preview [54]. For generation settings, we use a maximum of 32,768 tokens, temperature of 0.7, top_p of 0.8, top_k of 20, and a repetition penalty of 1.05 across all models. For retrieval, we employ the Bing Web Search API, setting the region to US-EN and the top-k retrieved documents to 10. We use Jina Reader API to fetch the content of web pages for given URLs. For all retrieval-based methods, following [52], we apply a back-off strategy where, when a final answer is not provided, we use the result from direct reasoning. For baseline models not specifically trained for o1-like reasoning, we apply Chain-of-Thought (CoT) [59] prompting to perform reasoning before generating answers. Detailed instructions for all models are provided in Appendix A. All experiments are conducted on eight NVIDIA A800-80GB GPUs.

4.4 Results on Challenging Reasoning Tasks

Main Results. Table 1 presents Search-o1’s performance on complex reasoning tasks, with the main results outlined below:

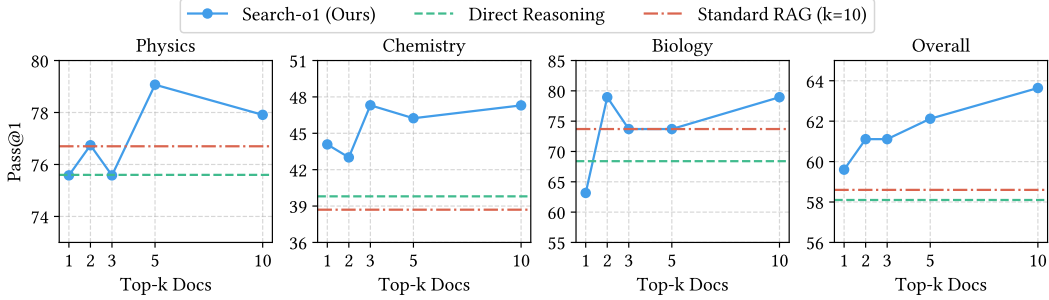


Figure 3: Scaling analysis of top-k retrieved documents utilized in reasoning. All results are based on QwQ-32B-Preview model.

1. For both settings without retrieval and with retrieval augmentation, **the large reasoning model QwQ-32B-Preview consistently shows superior performance** compared to traditional instruction-tuned LLMs. The QwQ model with 32B parameters even outperforms larger LLMs such as Qwen2.5-72B and Llama3.3-70B in the direct reasoning setting, demonstrating the effectiveness of the o1-like long CoT approach in complex reasoning.
2. **RAgent-QwQ-32B surpasses both standard RAG-based models and direct reasoning QwQ-32B in most tasks**, thanks to its agentic search mechanism, which autonomously retrieves information to supplement the knowledge required for reasoning at each step. Additionally, we find that the non-reasoning model Qwen2.5-32B using agentic RAG performs similarly to standard RAG on GPQA and even shows decreased performance on math and code tasks. This indicates that ordinary LLMs cannot effectively utilize search as a tool to solve complex reasoning tasks.
3. **Our Search-o1 further outperforms RAgent-QwQ-32B in most tasks**, demonstrating the effectiveness of our Reason-in-Documents strategy by integrating external knowledge while ensuring that it does not affect the coherence of the original reasoning. Specifically, on average across all five datasets, Search-o1 exceeds RAgent-QwQ-32B by 4.7% and 3.1%, respectively, and significantly outperforms non-reasoning models Qwen2.5-32B and Llama3.3-70B by 44.7% and 39.3%.

Scaling Analysis on Number of Retrieved Documents. In this experiment, we analyze the performance variation with respect to the number of retrieved documents, as shown in Figure 3. Our results demonstrate that **Search-o1 can effectively leverage an increasing number of retrieved documents, leading to improvements in handling complex reasoning tasks.** We also observe that for overall performance, retrieving even one document can surpass Direct Reasoning and standard RAG models that use ten retrieved documents, showcasing the effectiveness of the agentic search and Reason-in-Documents strategies.

Table 2: Performance comparison with human experts on the GPQA extended set [52].

Method	GPQA Extended Set			
	Physics	Chemistry	Biology	Overall
Human Experts				
Physicists	57.9	31.6	42.0	39.9
Chemists	34.5	72.6	45.6	48.9
Biologists	30.4	28.8	<u>68.9</u>	37.2
Reasoning Models				
QwQ-32B	61.7	36.9	61.0	51.8
RAG-QwQ-32B	<u>64.3</u>	38.3	66.7	<u>54.6</u>
Search-o1 (Ours)	68.7	<u>40.7</u>	69.5	57.9

Comparison with Human Experts. We compare the performance of our Search-o1 with human experts across various domains in the GPQA extended set. Table 2 presents the evaluation of human experts from various disciplines, including physics, chemistry, and biology. Our Search-o1 model outperforms human experts in overall performance (57.9), as well as in both physics (68.7) and biology (69.5), demonstrating superior handling of complex reasoning tasks. While Search-o1 slightly trails chemists in the chemistry subdomain (40.7 vs. 72.6), it still provides a competitive edge overall, particularly in terms of general performance across multiple domains. **This highlights the effectiveness of Search-o1 in leveraging document retrieval and reasoning to achieve cross-domain performance that rivals or exceeds expert-level capabilities.**

Table 3: Performance comparison on open-domain QA tasks, including single-hop QA and multi-hop QA datasets. For models with 32B parameters, the best results are in **bold** and the second-best are underlined. Results from larger models are in gray color for reference.

Method	Single-hop QA				Multi-hop QA							
	NQ		TriviaQA		HotpotQA		2WIKI		MuSiQue		Bamboogle	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
<i>Direct Reasoning (w/o Retrieval)</i>												
Qwen2.5-32B	22.8	33.9	52.0	60.3	25.4	34.7	29.8	36.3	8.4	18.0	49.6	63.2
QwQ-32B	23.0	33.1	53.8	60.7	25.4	33.3	34.4	40.9	9.0	18.9	38.4	53.7
Qwen2.5-72B	27.6	41.2	56.8	65.8	29.2	38.8	34.4	42.7	11.4	20.4	47.2	61.7
Llama3.3-70B	36.0	48.7	68.8	76.8	37.8	49.1	46.0	54.2	14.8	23.6	54.4	67.8
<i>Retrieval-augmented Reasoning</i>												
RAG-Qwen2.5-32B	33.4	<u>49.3</u>	65.8	79.2	38.6	50.4	31.6	40.6	10.4	19.8	52.0	66.0
RAG-QwQ-32B	29.6	44.4	<u>65.6</u>	<u>77.6</u>	34.2	46.4	35.6	46.2	10.6	20.2	<u>55.2</u>	<u>67.4</u>
RAgent-Qwen2.5-32B	32.4	47.8	63.0	72.6	<u>44.6</u>	56.8	55.4	69.7	13.0	25.4	54.4	66.4
RAgent-QwQ-32B	<u>33.6</u>	48.4	62.0	74.0	43.0	55.2	58.4	<u>71.2</u>	<u>13.6</u>	<u>25.5</u>	52.0	64.7
<i>Retrieval-augmented Reasoning with Reason-in-Documents</i>												
Search-o1 (Ours)	34.0	49.7	63.4	74.1	45.2	57.3	<u>58.0</u>	71.4	16.6	28.2	56.0	67.8

4.5 Results on Open-Domain QA Tasks

In addition to the reasoning tasks where LRMs excel, we also explore the performance of our Search-o1 on open-domain QA tasks. Table 3 presents the overall results. The key observations are:

1. For direct reasoning without retrieval, the performance of the LRM QwQ-32B is overall similar to the non-reasoning LLM Qwen2.5-32B, with a slight decrease in average EM across all QA datasets (31.3 vs. 30.7). **This indicates that LRMs do not perform as strongly on open-domain QA tasks as they do on reasoning tasks.**
2. When employing retrieval-augmented reasoning, retrieval significantly improves performance for both reasoning and non-reasoning models across all tasks, suggesting that **models have knowledge gaps in these tasks**. Additionally, for the QwQ-32B model, agentic RAG achieves an average EM improvement of 23.2% over standard RAG on multi-hop QA tasks, demonstrating the **effectiveness of our agentic RAG strategy in knowledge-based multi-hop QA**. However, we also observe that there is no significant performance change for single-hop tasks (47.8 vs. 47.6 on average EM), as these questions only require information from a single knowledge point without the need for multiple retrievals. **This also verifies that the agentic search mechanism can better unleash the potential of LRMs in more complex and challenging reasoning tasks.**
3. **For our Search-o1, we find that it generally outperforms all baselines on multi-hop tasks.** Specifically, in terms of the average EM metric, our Search-o1 exceeds RAG-QwQ-32B and RAgent-QwQ-32B by 29.6% and 5.3%, respectively, demonstrating the effectiveness of our Reason-in-Documents strategy in complex QA tasks. **This further emphasizes the importance of maintaining consistency between external knowledge and the logical chain of reasoning.**

5 Conclusion

In this work, we present Search-o1, a framework that addresses the knowledge insufficiency inherent in large reasoning models (LRMs) by integrating an agentic retrieval-augmented generation mechanism alongside a Reason-in-Documents module. Our approach enables LRMs to autonomously retrieve and seamlessly incorporate external knowledge during the reasoning process, thereby enhancing both the accuracy and coherence of their long-step reasoning capabilities. Comprehensive experiments across diverse complex reasoning tasks in science, mathematics, and coding, as well as multiple open-domain QA benchmarks, demonstrate that Search-o1 consistently outperforms existing retrieval-augmented and direct reasoning methods. Notably, Search-o1 not only surpasses baseline models in handling intricate reasoning challenges but also achieves performance levels comparable to

or exceeding human experts in specific domains. These findings underscore the potential of Search-o1 to significantly improve the reliability and versatility of LRMs, paving the way for more trustworthy and effective intelligent systems in complex problem-solving scenarios.

References

- [1] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- [2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’ Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [3] Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. Huatuoqpt-o1, towards medical complex reasoning with llms. *arXiv preprint arXiv:2412.18925*, 2024.
- [4] Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Do not think that much for $2+3=?$ on the overthinking of o1-like llms, 2024.
- [5] Zehui Chen, Kuikun Liu, Qiuchen Wang, Jiangning Liu, Wenwei Zhang, Kai Chen, and Feng Zhao. Mindsearch: Mimicking human minds elicits deep ai searcher. *arXiv preprint arXiv:2407.20183*, 2024.
- [6] Kahneman Daniel. *Thinking, fast and slow*. 2017.
- [7] DeepSeek-AI. Deepseek-r1-lite-preview is now live: unleashing supercharged reasoning power!, November 2024.
- [8] Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu, Chang Zhou, and Jingren Zhou. Self-play with execution feedback: Improving instruction-following capabilities of large language models. *CoRR*, abs/2406.13542, 2024.
- [9] Guanting Dong, Xiaoshuai Song, Yutao Zhu, Runqi Qiao, Zhicheng Dou, and Ji-Rong Wen. Toward general instruction-following alignment for retrieval-augmented generation. *CoRR*, abs/2410.09584, 2024.
- [10] Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected by supervised fine-tuning data composition. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 177–198. Association for Computational Linguistics, 2024.
- [11] Guanting Dong, Chenghao Zhang, Mengjie Deng, Yutao Zhu, Zhicheng Dou, and Ji-Rong Wen. Progressive multimodal reasoning via active retrieval. *arXiv preprint arXiv:2412.14835*, 2024.
- [12] Guanting Dong, Yutao Zhu, Chenghao Zhang, Zechen Wang, Zhicheng Dou, and Ji-Rong Wen. Understand what LLM needs: Dual preference alignment for retrieval-augmented generation. *CoRR*, abs/2406.18676, 2024.
- [13] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- [14] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization, 2024.
- [15] Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information Processing Systems*, 36, 2024.
- [16] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In Joaquin Vanschoren and Sai-Kit Yeung, editors, *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021.
- [17] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- [18] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing A multi-hop QA dataset for comprehensive evaluation of reasoning steps. In Donia Scott, Núria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 6609–6625. International Committee on Computational Linguistics, 2020.
- [19] Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei Qin, Weizhe Yuan, and Pengfei Liu. O1 replication journey—part 2: Surpassing o1-preview through simple distillation, big progress or bitter lesson? *arXiv preprint arXiv:2411.16489*, 2024.
- [20] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report. *CoRR*, abs/2409.12186, 2024.
- [21] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [22] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [23] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *CoRR*, abs/2403.07974, 2024.
- [24] Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression. *arXiv preprint arXiv:2310.06839*, 2023.
- [25] Jinhao Jiang, Zhipeng Chen, Yingqian Min, Jie Chen, Xiaoxue Cheng, Jiapeng Wang, Yiru Tang, Haoxiang Sun, Jia Deng, Wayne Xin Zhao, et al. Technical report: Enhancing llm reasoning with reward-guided tree search. *arXiv preprint arXiv:2411.11694*, 2024.
- [26] Bowen Jin, Jinsung Yoon, Jiawei Han, and Sercan Ö. Arik. Long-context llms meet RAG: overcoming challenges for long inputs in RAG. *CoRR*, abs/2410.05983, 2024.
- [27] Jiajie Jin, Yutao Zhu, Yujia Zhou, and Zhicheng Dou. Bider: Bridging knowledge inconsistency for efficient retrieval-augmented llms via key supporting evidence. *arXiv preprint arXiv:2402.12174*, 2024.
- [28] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In *ACL*, pages 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics.

- [29] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- [30] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [31] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay V. Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- [32] Chengpeng Li, Guanting Dong, Mingfeng Xue, Ru Peng, Xiang Wang, and Dayiheng Liu. Dotamath: Decomposition of thought with code assistance and self-correction for mathematical reasoning. *CoRR*, abs/2407.04078, 2024.
- [33] Xiaoxi Li, Zhicheng Dou, Yujia Zhou, and Fangchao Liu. Corpuslm: Towards a unified language model on corpus for knowledge-intensive tasks. In Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang, editors, *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, pages 26–37. ACM, 2024.
- [34] Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yongkang Wu, Zhonghua Li, Qi Ye, and Zhicheng Dou. Retrollm: Empowering large language models to retrieve fine-grained evidence within generation. *arXiv preprint arXiv:2412.11919*, 2024.
- [35] Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian Zhang, Yutao Zhu, and Zhicheng Dou. From matching to generation: A survey on generative information retrieval. *CoRR*, abs/2404.14851, 2024.
- [36] Xiaoxi Li, Yujia Zhou, and Zhicheng Dou. Unigen: A unified generative framework for retrieval and question answering with large language models. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan, editors, *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada*, pages 8688–8696. AAAI Press, 2024.
- [37] Lei Liang, Mengshu Sun, Zhengke Gui, Zhongshu Zhu, Zhouyu Jiang, Ling Zhong, Yuan Qu, Peilong Zhao, Zhongpu Bo, Jin Yang, Huaidong Xiong, Lin Yuan, Jun Xu, Zaoyang Wang, Zhiqiang Zhang, Wen Zhang, Huajun Chen, Wenguang Chen, and Jun Zhou. Kag: Boosting llms in professional domains via knowledge augmented generation, 2024.
- [38] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [39] Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Raghavi Chandu, Chandra Bhagavatula, and Yejin Choi. The unlocking spell on base llms: Rethinking alignment via in-context learning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [40] Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. Deductive verification of chain-of-thought reasoning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.

- [41] Jingyu Liu, Jiaen Lin, and Yong Liu. How much can RAG help the reasoning of llm? *CoRR*, abs/2410.02338, 2024.
- [42] Jingyu Liu, Jiaen Lin, and Yong Liu. How much can RAG help the reasoning of llm? *CoRR*, abs/2410.02338, 2024.
- [43] Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. Query rewriting for retrieval-augmented large language models. *arXiv preprint arXiv:2305.14283*, 2023.
- [44] Ning Miao, Yee Whye Teh, and Tom Rainforth. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [45] Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwu Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, et al. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems. *arXiv preprint arXiv:2412.09413*, 2024.
- [46] OpenAI. Learning to reason with llms, 2024.
- [47] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 5687–5711. Association for Computational Linguistics, 2023.
- [48] Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma Gongque, Shanglin Lei, Zhe Wei, Miaoxuan Zhang, Runfeng Qiao, Yifan Zhang, Xiao Zong, Yida Xu, Muxi Diao, Zhimin Bao, Chen Li, and Honggang Zhang. We-math: Does your large multimodal model achieve human-like mathematical reasoning? *CoRR*, abs/2407.01284, 2024.
- [49] Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Hector Liu, Yuanzhi Li, et al. O1 replication journey: A strategic progress report—part 1. *arXiv preprint arXiv:2410.18982*, 2024.
- [50] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2024.
- [51] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020.
- [52] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. *CoRR*, abs/2311.12022, 2023.
- [53] Jiejun Tan, Zhicheng Dou, Yutao Zhu, Peidong Guo, Kun Fang, and Ji-Rong Wen. Small models, big insights: Leveraging slim proxy models to decide when and what to retrieve for llms. *arXiv preprint arXiv:2402.12052*, 2024.
- [54] Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown, November 2024.
- [55] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022.
- [56] Prakhar Verma, Sukruta Prakash Midigeshi, Gaurav Sinha, Arno Solin, Nagarajan Natarajan, and Amit Sharma. Planxrag: Planning-guided retrieval augmented generation. *arXiv preprint arXiv:2410.20753*, 2024.

- [57] Jiaan Wang, Fandong Meng, Yunlong Liang, and Jie Zhou. Drt-o1: Optimized deep reasoning translation via long chain-of-thought. *arXiv preprint arXiv:2412.17498*, 2024.
- [58] Liang Wang, Nan Yang, and Furu Wei. Query2doc: Query expansion with large language models. *arXiv preprint arXiv:2303.07678*, 2023.
- [59] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [60] Ken C. L. Wong, Hongzhi Wang, Etienne E. Vos, Bianca Zadrozny, Campbell D. Watson, and Tanveer F. Syeda-Mahmood. Addressing deep learning model uncertainty in long-range climate forecasting with late fusion. *CoRR*, abs/2112.05254, 2021.
- [61] Siwei Wu, Zhongyuan Peng, Xinrun Du, Tuney Zheng, Minghao Liu, Jialong Wu, Jiachen Ma, Yizhi Li, Jian Yang, Wangchunshu Zhou, Qunshu Lin, Junbo Zhao, Zhaoxiang Zhang, Wenhao Huang, Ge Zhang, Chenghua Lin, and Jiaheng Liu. A comparative study on reasoning patterns of openai’s o1 model. *CoRR*, abs/2410.13639, 2024.
- [62] Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai Zhang, and Yanghua Xiao. How easily do irrelevant inputs skew the responses of large language models?, 2024.
- [63] Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. Evaluating mathematical reasoning beyond accuracy. *CoRR*, abs/2404.05692, 2024.
- [64] Fangyuan Xu, Weijia Shi, and Eunsol Choi. Recomp: Improving retrieval-augmented lms with compression and selective augmentation. *arXiv preprint arXiv:2310.04408*, 2023.
- [65] Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. Llava-o1: Let vision language models reason step-by-step. *arXiv preprint arXiv:2411.10440*, 2024.
- [66] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report. *CoRR*, abs/2407.10671, 2024.
- [67] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *CoRR*, abs/2409.12122, 2024.
- [68] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *EMNLP*, pages 2369–2380, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.
- [69] Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, et al. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. *arXiv preprint arXiv:2412.18319*, 2024.
- [70] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [71] Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.2, how to learn from mistakes on grade-school math problems. *arXiv preprint arXiv:2408.16293*, 2024.

- [72] Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. Making retrieval-augmented language models robust to irrelevant context. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [73] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [74] Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. Scaling relationship on learning mathematical reasoning with large language models. *CoRR*, abs/2308.01825, 2023.
- [75] Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.
- [76] Zhiyuan Zeng, Qinyuan Cheng, Zhangyue Yin, Bo Wang, Shimin Li, Yunhua Zhou, Qipeng Guo, Xuanjing Huang, and Xipeng Qiu. Scaling of search and learning: A roadmap to reproduce o1 from reinforcement learning perspective. *arXiv preprint arXiv:2412.14135*, 2024.
- [77] Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jiatong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco Pavone, Yuqiang Li, Wanli Ouyang, and Dongzhan Zhou. Llama-berry: Pairwise optimization for o1-like olympiad-level mathematical reasoning. *CoRR*, abs/2410.02884, 2024.
- [78] Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, et al. Aflow: Automating agentic workflow generation. *arXiv preprint arXiv:2410.10762*, 2024.
- [79] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren’s song in the AI ocean: A survey on hallucination in large language models. *CoRR*, abs/2309.01219, 2023.
- [80] Yuxiang Zhang, Shangxi Wu, Yuqi Yang, Jiangming Shu, Jinlin Xiao, Chao Kong, and Jitao Sang. o1-coder: an o1 replication for coding. *CoRR*, abs/2412.00154, 2024.
- [81] Yuxiang Zhang, Shangxi Wu, Yuqi Yang, Jiangming Shu, Jinlin Xiao, Chao Kong, and Jitao Sang. o1-coder: an o1 replication for coding. *arXiv preprint arXiv:2412.00154*, 2024.
- [82] Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, and Bin Cui. Retrieval-augmented generation for ai-generated content: A survey. *arXiv preprint arXiv:2402.19473*, 2024.
- [83] Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. Processbench: Identifying process errors in mathematical reasoning. *arXiv preprint arXiv:2412.06559*, 2024.
- [84] Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, Chao Cao, Hanqi Jiang, Hanxu Chen, Yiwei Li, Junhao Chen, Huawen Hu, Yihen Liu, Huaqin Zhao, Shaochen Xu, Haixing Dai, Lin Zhao, Ruidong Zhang, Wei Zhao, Zhenyuan Yang, Jingyuan Chen, Peilong Wang, Wei Ruan, Hui Wang, Huan Zhao, Jing Zhang, Yiming Ren, Shihuan Qin, Tong Chen, Jiaxi Li, Arif Hassan Zidan, Afrar Jahin, Minheng Chen, Sichen Xia, Jason Holmes, Yan Zhuang, Jiaqi Wang, Bochen Xu, Weiran Xia, Jichao Yu, Kaibo Tang, Yaxuan Yang, Bolun Sun, Tao Yang, Guoyu Lu, Xianqiao Wang, Lilong Chai, He Li, Jin Lu, Lichao Sun, Xin Zhang, Bao Ge, Xintao Hu, Lian Zhang, Hua Zhou, Lu Zhang, Shu Zhang, Ninghao Liu, Bei Jiang, Linglong Kong, Zhen Xiang, Yudan Ren, Jun Liu, Xi Jiang, Yu Bao, Wei Zhang, Xiang Li, Gang Li, Wei Liu, Dinggang Shen, Andrea Sikora, Xiaoming Zhai, Dajiang Zhu, and Tianming Liu. Evaluation of openai o1: Opportunities and challenges of AGI. *CoRR*, abs/2409.18486, 2024.

- [85] Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, et al. Evaluation of openai o1: Opportunities and challenges of agi. *arXiv preprint arXiv:2409.18486*, 2024.
- [86] Yujia Zhou, Yan Liu, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Zheng Liu, Chaozhuo Li, Zhicheng Dou, Tsung-Yi Ho, and Philip S. Yu. Trustworthiness in retrieval-augmented generation systems: A survey. *CoRR*, abs/2409.10102, 2024.
- [87] Yujia Zhou, Zheng Liu, and Zhicheng Dou. Assistrag: Boosting the potential of large language models with an intelligent information assistant. *CoRR*, abs/2411.06805, 2024.
- [88] Yujia Zhou, Zheng Liu, Jiajie Jin, Jian-Yun Nie, and Zhicheng Dou. Metacognitive retrieval-augmented large language models. In Tat-Seng Chua, Chong-Wah Ngo, Ravi Kumar, Hady W. Lauw, and Roy Ka-Wei Lee, editors, *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*, pages 1453–1463. ACM, 2024.

Appendix

A Instruction Templates

A.1 Instructions for Search-o1

Instruction for Search-o1

You are a reasoning assistant with the ability to perform web searches to help you answer the user's question accurately. You have special tools:

To perform a search: write `<begin_search_query>` your query here `<end_search_query>`.

Then, the system will search and analyze relevant web pages, then provide you with helpful information in the format `<begin_search_result>` ...search results... `<end_search_result>`.

You can repeat the search process multiple times if necessary. The maximum number of search attempts is limited to {MAX_SEARCH_LIMIT}.

Once you have all the information you need, continue your reasoning.

Example:

Question: "..."

Assistant thinking steps:

- I might need to look up details about ...

Assistant:

`<begin_search_query>`...`<end_search_query>`

(System returns processed information from relevant web pages)

Assistant continues reasoning with the new information...

Remember:

- Use `<begin_search_query>` to request a web search and end with `<end_search_query>`.
- When done searching, continue your reasoning.

Instruction for Reason-in-Documents

Task Instruction:

You are tasked with reading and analyzing web pages based on the following inputs: Previous Reasoning Steps, Current Search Query, and Searched Web Pages. Your objective is to extract relevant and helpful information for Current Search Query from the Searched Web Pages and seamlessly integrate this information into the Previous Reasoning Steps to continue reasoning for the original question.

Guidelines:

1. Analyze the Searched Web Pages:

- Carefully review the content of each searched web page.
- Identify factual information that is relevant to the Current Search Query and can aid in the reasoning process for the original question.

2. Extract Relevant Information:

- Select the information from the Searched Web Pages that directly contributes to advancing the Previous Reasoning Steps.
- Ensure that the extracted information is accurate and relevant.

3. Output Format:

- If the web pages provide helpful information for current search query: Present the information beginning with 'Final Information' as shown below.

Final Information

[Helpful information]

- If the web pages do not provide any helpful information for current search query: Output the following text.

Final Information

No helpful information found.

Inputs:

- Previous Reasoning Steps:

{prev_reasoning}

- Current Search Query:

{search_query}

- Searched Web Pages:

{document}

Now you should analyze each web page and find helpful information based on the current search query "{search_query}" and previous reasoning steps.

A.2 Instructions for Standard RAG

Instruction for Standard RAG

You are a knowledgeable assistant that utilizes the provided documents to answer the user's question accurately.

Question: {question}

Documents: {documents}

Guidelines:

- Analyze the provided documents to extract relevant information. Synthesize the information to formulate a coherent and accurate answer.
- Ensure that your response directly addresses the user's question using the information from the documents.

A.3 Instructions for RAG Agent

Instruction for RAG Agent

You are a reasoning assistant with the ability to perform web searches and retrieve webpage content to help you answer the user's question accurately. You have special tools:

- To perform a search: Write '<begin_search_query>' your query here '<end_search_query>'.

The system will call the web search API with your query and return the search results in the format '<begin_search_result> ...search results... <end_search_result>'.

The search results will include a list of webpages with titles, URLs, and snippets (but not full content).

- To retrieve full page content: After receiving the search results, if you need more detailed information from specific URLs, write '<begin_url> url1, url2, ... <end_url>'.

The system will fetch the full page content of those URLs and return it as '<begin_full_page> ...full page content... <end_full_page>'.

You can repeat the search process multiple times if necessary. The maximum number of search attempts is limited to {MAX_SEARCH_LIMIT}. You can fetch up to {MAX_URL_FETCH} URLs for detailed information.

Once you have all the information you need, continue your reasoning.

Example:

Question: "..."

Assistant thinking steps: - I need to find out ...

Assistant: '<begin_search_query>...<end_search_query>'

(System returns search results)

Assistant: '<begin_search_result> ...search results without full page... <end_search_result>'

Assistant thinks: The search results mention several URLs. I want full details from one of them.

Assistant: '<begin_url>http://...<end_url>'

(System returns full page content)

Assistant: '<begin_full_page> ...full page content... <end_full_page>'

Now the assistant has enough information and can continue reasoning.

Remember:

- Use '<begin_search_query>' to request a web search and end with '<end_search_query>'.
- Use '<begin_url>' to request full page content and end with '<end_url>'.
- When done retrieving information, continue your reasoning.

A.4 Task-Specific Instructions

A.4.1 Open-Domain QA Tasks Instruction

Instruction for Open-Domain QA Tasks

Please answer the following question.

You should provide your final answer in the format \boxed{YOUR_ANSWER}.

Question:

{question}

A.4.2 Math Tasks Instruction

Instruction for Math Tasks

Please answer the following math question.
You should provide your final answer in the format `\boxed{YOUR_ANSWER}`.
Question:
{question}

A.4.3 Multi-choice Tasks Instruction

Instruction for Multi-choice Tasks

You are to answer the following multiple-choice question by selecting the correct option.
Your final choice should be one of the letters A, B, C, or D. Do not include any answer content beyond the choice letter.
You should provide your final choice in the format `\boxed{YOUR_CHOICE}`.
Question: {question}

A.4.4 Code Tasks Instruction

Instruction for Code Tasks

Generate a correct Python program that passes all tests for the given problem. You should provide your final code within a Python code block using triple backticks.

```
```python
YOUR CODE HERE
```
```

Problem Title: {question_title}
Problem Statement:
{question}

A.5 Additional Notes

For all the instructions above, we input them as user prompts, not system prompts. The task-specific instructions in A.4 are used for the QwQ-32B-Preview model. For non-reasoning models like Qwen2.5-32B-Instruct, Qwen2.5-72B-Instruct, and Llama3.3-70B-Instruct, etc., we add a Chain-of-Thought prompt "You should think step by step to solve it." before the question to explicitly make these models reason before giving the final answer.

B Case Study

Tables 4, 5, and 6 present examples of outputs from our Search-o1 model on the GPQA, AMC2023, and HotpotQA datasets, respectively. The model-generated search queries are enclosed within `<|begin_search_query|>` and `<|end_search_query|>`, while the refined search results are enclosed within `<|begin_search_result|>` and `<|end_search_result|>`. We observe that our Reason-in-Documents mechanism provides coherent information that effectively meets the information needs of the model's current reasoning step and seamlessly integrates with the preceding reasoning process.

Table 4: An example from Search-o1 on GPQA dataset, with special symbols used in the search queries and search results highlighted in green box and blue box, respectively.

Example #1 from GPQA Dataset

Question:

trans-cinnamaldehyde was treated with methylmagnesium bromide, forming product 1.
 1 was treated with pyridinium chlorochromate, forming product 2.
 3 was treated with (dimethyl(oxo)-l6-sulfaneylidene)methane in DMSO at elevated temperature, forming product 3.
 how many carbon atoms are there in product 3? Choices:
 (A) 10
 (B) 11
 (C) 14
 (D) 12

Labeled Answer:

B

Model Output by Search-o1

Alright, so I have this problem here about some chemical reactions, and I need to figure out how many carbon atoms are in the final product, which is product 3. Let's break this down step by step.

First, there's trans-cinnamaldehyde, and it's being treated with methylmagnesium bromide to form product 1. I know that methylmagnesium bromide is a Grignard reagent, which is nucleophilic and usually adds to carbonyl compounds. So, trans-cinnamaldehyde has a carbonyl group, right? It's an aldehyde, so it has a C=O bond. The Grignard reagent should attack the carbonyl carbon, adding the methyl group and turning the aldehyde into an alcohol, I think. But let me make sure.

Wait, maybe I should look up what exactly happens when a Grignard reagent reacts with an aldehyde. Let me try to search for that.

<lbegin_search_query> reaction of grignard reagent with aldehyde <lend_search_query>

<lbegin_search_result> Reacting a Grignard reagent with an aldehyde results in the formation of a secondary alcohol, incorporating one additional carbon atom from the Grignard reagent. Therefore, if trans-cinnamaldehyde (9 carbon atoms) reacts with methylmagnesium bromide (1 carbon atom), the resulting product (product 1) will have 10 carbon atoms. Subsequent protonation to form product 2 and further reactions to form product 3 do not change the carbon atom count, assuming no carbon is added or removed in those steps. Hence, product 3 should have 10 carbon atoms. <lend_search_result>

Hmm, according to this, product 1 should have 10 carbon atoms, starting from trans-cinnamaldehyde, which has 9 carbon atoms, plus the methyl group from the Grignard reagent, adding one more carbon atom. That makes sense.

Next, product 1 is treated with pyridinium chlorochromate to form product 2. Pyridinium chlorochromate (PCC) is used for oxidations, particularly for oxidizing primary alcohols to aldehydes or secondary alcohols to ketones. Since product 1 is a secondary alcohol (from the previous step), PCC should oxidize it to a ketone. Again, this shouldn't change the number of carbon atoms; it's just changing the functional group.

...

Alternatively, perhaps the initial trans-cinnamaldehyde has more or fewer carbon atoms than I thought. Let me confirm its structure.

<lbegin_search_query> structure of trans-cinnamaldehyde <lend_search_query>

<lbegin_search_result> Trans-cinnamaldehyde, also known as (E)-cinnamaldehyde, has the molecular formula C₉H₈O. It consists of a phenyl ring connected to a propenal group via a double bond, forming a conjugated system that makes the molecule planar. The structure comprises 9 carbon atoms, 8 hydrogen atoms, and 1 oxygen atom, with the functional groups indicating its reactivity in chemical transformations. <lend_search_result>
 Yes, confirmed: trans-cinnamaldehyde has 9 carbon atoms.

Now, reacting with methylmagnesium bromide, which provides a methyl group (-CH₃), adds one carbon atom, making product 1 have 10 carbon atoms.

Then, oxidation with PCC to convert the alcohol to a ketone doesn't change the carbon skeleton, so product 2 still has 10 carbon atoms.

...

Therefore, starting from 9 carbon atoms in trans-cinnamaldehyde, adding one carbon atom from the Grignard reagent to make 10, and then adding another carbon atom in the cyclopropanation to make 11.

So, the final answer is B

Table 5: An example from Search-o1 on AMC2023 dataset, with special symbols used in the search queries and search results highlighted in green box and blue box, respectively.

| Example #2 from AMC2023 Dataset |
|--|
| <p>Question:
There is a unique sequence of integers $a_1, a_2, \dots, a_{2023}$ such that</p> $\tan 2023x = \frac{a_1 \tan x + a_3 \tan^3 x + a_5 \tan^5 x + \dots + a_{2023} \tan^{2023} x}{1 + a_2 \tan^2 x + a_4 \tan^4 x + \dots + a_{2022} \tan^{2022} x}$ <p>whenever $\tan 2023x$ is defined. What is a_{2023}?</p> <p>Labeled Answer:
-1</p> |
| Model Output by Search-o1 |
| <p>Alright, so I've got this problem here involving trigonometric identities, specifically with the tangent function. The equation given is:</p> <p>...</p> <p>And I need to find the value of a_{2023}.</p> <p>First off, I know that there are multiple-angle identities for tangent, but I'm not sure if there's a direct formula for such a high multiple like 2023. Maybe there's a general formula for $\tan nx$ in terms of $\tan x$, and perhaps this problem is asking for the coefficients in that expansion.</p> <p>...</p> <p>Wait, I remember that tangent of multiple angles can be expressed using the Bernoulli polynomials or something similar, but I'm not sure. Maybe I need to search for a general formula for $\tan nx$.</p> <p>Let me try to search for tangent multiple angle formula to see if there's a general expression.</p> <p><begin_search_query>tangent multiple angle formula</end_search_query></p> <p><begin_search_result>The recursive formula for $\tan(n\theta)$ is given by:</p> $\tan(n\theta) = \frac{\tan[(n-1)\theta] + \tan \theta}{1 - \tan[(n-1)\theta] \tan \theta}$ <p>This formula can be used iteratively to compute $\tan(n\theta)$ for any integer n, making it particularly useful for large values of n, such as $n=2023$, as encountered in the problem. This approach avoids the complexity of deriving high-degree polynomial expressions directly and leverages the tangent addition formula in a step-by-step manner.</end_search_result></p> <p>Okay, so there's a recursive formula for $\tan(n\theta)$:</p> $\tan(n\theta) = \frac{\tan[(n-1)\theta] + \tan \theta}{1 - \tan[(n-1)\theta] \tan \theta}$ <p>That seems useful, especially for large n, like 2023. But I'm not sure how to apply this recursively up to 2023 without going through each step, which isn't practical.</p> <p>Maybe there's a general expression for $\tan(nx)$ in terms of $\tan(x)$ as a rational function.</p> <p>...</p> <p>Looking at $n = 1$:</p> $\tan x = \frac{a_1 \tan x}{1} \implies a_1 = 1$ <p>For $n = 3$, $a_3 = -1$</p> <p>For $n = 5$, $a_5 = 1$</p> <p>So it seems that $a_n = (-1)^{(n-1)/2}$ for odd n.</p> <p>Wait, let's check:</p> <p>For $n = 1$, $(-1)^{(1-1)/2} = (-1)^0 = 1$, which matches $a_1 = 1$</p> <p>For $n = 3$, $(-1)^{(3-1)/2} = (-1)^1 = -1$, which matches $a_3 = -1$</p> <p>For $n = 5$, $(-1)^{(5-1)/2} = (-1)^2 = 1$, which matches $a_5 = 1$</p> <p>So, generalizing, for odd n, $a_n = (-1)^{(n-1)/2}$</p> <p>Now, 2023 is odd, so $a_{2023} = (-1)^{(2023-1)/2} = (-1)^{1011} = -1$</p> <p>Therefore, $a_{2023} = -1$</p> <p>I think that's the answer.</p> <p>*Final Answer**</p> <div style="border: 1px solid black; width: 30px; margin: 0 auto; text-align: center; padding: 2px;">-1</div> |

Table 6: An example from Search-o1 on HotpotQA dataset, with special symbols used in the search queries and search results highlighted in green box and blue box, respectively.

| Example #3 from HotpotQA Dataset |
|--|
| <p>Question:
Were Goo Goo Dolls and Echosmith formed in the same city?</p> <p>Labeled Answer:
no</p> |
| Model Output by Search-o1 |
| <p>Alright, I've got this question: Were Goo Goo Dolls and Echosmith formed in the same city? Hmm, I don't know much about these bands, so I'll need to look that up.</p> <p>First, I should find out where each band was formed. Maybe start with Goo Goo Dolls. I think they're an older band, but I'm not sure. Let me search for their formation place.</p> <p>< begin_search_query >where was Goo Goo Dolls formed?< end_search_query ></p> <p>< begin_search_result >The Goo Goo Dolls were formed in Buffalo, New York, in 1986. They initially performed under the name Sex Maggots but changed to Goo Goo Dolls in 1986 for a gig, choosing the name from an ad in True Detective magazine.< end_search_result ></p> <p>Okay, so Goo Goo Dolls were formed in Buffalo, New York. Now, I need to find out where Echosmith was formed.</p> <p>< begin_search_query >where was Echosmith formed?< end_search_query ></p> <p>< begin_search_result >Echosmith was formed in Chino, California, in February 2009.< end_search_result ></p> <p>Chino, California, huh? That sounds different from Buffalo, New York. But maybe I should double-check to make sure I have the right information.</p> <p>Let me look up Echosmith again to confirm their formation city.</p> <p>< begin_search_query >Echosmith formation location< end_search_query ></p> <p>< begin_search_result >Echosmith was formed in Toluca Lake, Los Angeles, California, in February 2009. The initial confusion about the formation location being Chino, California, or Buffalo, New York, is resolved by multiple recent sources confirming their formation in Los Angeles. Chino might be their hometown, but the band was established in Toluca Lake, Los Angeles.< end_search_result ></p> <p>Wait, now it says Toluca Lake, Los Angeles, California. Earlier, I had Chino, California. Maybe there's some confusion here.</p> <p>I should check more sources to make sure.</p> <p>< begin_search_query >Echosmith origin city< end_search_query ></p> <p>< begin_search_result >Echosmith was born in Chino, California, but formed the band in Toluca Lake, Los Angeles, California, in February 2009. Chino is their hometown, while Toluca Lake is where they officially formed as a band.< end_search_result ></p> <p>So, it seems like Echosmith originates from Chino, California, but was formed in Toluca Lake, Los Angeles, California.</p> <p>Now, comparing that to Goo Goo Dolls, who were formed in Buffalo, New York.</p> <p>Clearly, Buffalo, New York, is different from Toluca Lake, Los Angeles, California.</p> <p>But just to be thorough, maybe there's a connection I'm missing.</p> <p>Let me see if Echosmith has any ties to Buffalo, New York.</p> <p>< begin_search_query >Echosmith Buffalo, New York< end_search_query ></p> <p>< begin_search_result >Echosmith has performed in Buffalo, New York, in the past, as evidenced by a setlist from a concert at Waiting Room on November 1, 2013. However, there are no current tour dates listed for Buffalo in the searched web pages.< end_search_result ></p> <p>So, they've played in Buffalo, but that doesn't mean they were formed there.</p> <p>Let me see if there's any other connection between the two bands.</p> <p>< begin_search_query >Goo Goo Dolls and Echosmith connection< end_search_query >< begin_search_result >The maximum search limit is exceeded. You are not allowed to search.< end_search_result ></p> <p>Hmm, it seems I've reached my search limit. Based on the information I have so far, Goo Goo Dolls were formed in Buffalo, New York, and Echosmith was formed in Toluca Lake, Los Angeles, California. These are different cities, so the answer should be no, they were not formed in the same city.</p> <p>Therefore, I can confidently say that Goo Goo Dolls and Echosmith were not formed in the same city.</p> <p>*Final Answer**</p> <div style="border: 1px solid black; padding: 2px; width: fit-content; margin: 0 auto;">no</div> |