# Can We Further Elicit Reasoning in LLMs? Critic-Guided Planning with Retrieval-Augmentation for Solving Challenging Tasks

**Xingxuan Li**[1,2*] **Weiwen Xu**[3,4] **Ruochen Zhao**[2] **Fangkai Jiao**[2]
**Shafiq Joty**[2,5] **Lidong Bing**[1†]

[1]MiroMind  [2]Nanyang Technological University
[3]DAMO Academy, Alibaba Group [4]Hupan Lab, Hangzhou
[5]Salesforce Research

{xingxuan.li, lidong.bing}@miromind.ai xuweiwen.xww@alibaba-inc.com

{ruochen002, fangkai002, srjoty}@ntu.edu.sg

## Abstract

Large language models excel at problem-solving but often struggle with complex reasoning and factual accuracy. While chain-of-thought and retrieval-augmented generation help break down problems and retrieve knowledge, they still falter on challenging tasks like competitive programming due to frequent reasoning errors and irrelevant retrieval. To address this, we introduce **C**ritic-guided planning with **R**etrieval-augmentation, CR-Planner, a novel framework that leverages fine-tuned critic models to guide both reasoning and retrieval processes through planning. CR-Planner iteratively selects and executes sub-goals, guided by critic models. A sub-goal critic identifies promising sub-goals from reasoning, query generation, and retrieval, while an execution critic evaluates outputs of sub-goal executions. We employ Monte Carlo Tree Search to collect data for critic training, allowing systematic exploration of action sequences and effective navigation toward the final answer. We evaluate CR-Planner on challenging domain-knowledge-intensive and reasoning-heavy tasks, including competitive programming, theorem-driven math reasoning, and complex domain retrieval problems. It significantly outperforms baselines, demonstrating effectiveness in both reasoning and retrieval. Our code is available at https://github.com/xingxuanli/CR-Planner.

## 1 Introduction

State-of-the-art large language models (LLMs), while demonstrating remarkable problem-solving capabilities (OpenAI, 2023; Chen et al., 2024; Zhang et al., 2025), still face two key challenges: reasoning for complex tasks (Huang et al., 2024) and domain-specific knowledge (Zhao et al., 2023a). Existing approaches (Yao et al., 2023b;

Zhao et al., 2023b; Li et al., 2024) seek to harness the strengths of both chain-of-thought (CoT) reasoning (Wei et al., 2022) and retrieval-augmented generation (RAG) (Lewis et al., 2020) on knowledge-intensive complex reasoning problems. These methods can potentially apply RAG at each reasoning step, enabling iterative retrieval and reasoning. Insights from reasoning improve retrieval relevance, while retrieved knowledge enhances the factuality of subsequent reasoning. Techniques like Self-RAG (Asai et al., 2024) and its variants (Yan et al., 2024; Islam et al., 2024) further refine this integration by fine-tuning LLMs to decide when and how to retrieve, using special reflection tokens.

While these methods show promise, they are generally limited to simpler reasoning tasks, such as answering two-hop questions like, "What year was the Argentine actor who directed El Tio Disparate born?" They often fail to address **domain-knowledge-intensive** and **reasoning-heavy** problems, such as competitive programming problems (Shi et al., 2024), which require advanced algorithmic knowledge and strong reasoning capability. Two key issues, illustrated in Figure 1 (a), contribute to these failures. The first is **reasoning error**: for instance, when solving "Given a string **s**, find the length of the longest substring without repeating characters in optimal time complexity," a CoT approach may incorrectly generate that "The optimal time complexity is $O(n^2)$," leading to cascading errors and an incorrect final answer. The second is **retrieving error**, where flawed reasoning generates poor search queries or selects incorrect retrieved documents, grounding subsequent reasoning on erroneous information.

To address these errors, we present critic-guided planning with retrieval-augmentation (CR-Planner), a framework designed to tackle reasoning-heavy problems requiring extensive domain knowledge. CR-Planner systematically plans both rea-

**Question**

Given a string **s**, find the length of the longest substring without repeating characters in optimal time complexity.

**(a) Chain of Thought & Retrieval-Augmented Generation**

**Step 1 (Reason)**: The optimal time complexity is $O(n^2)$

**Step 2 (Retrieve)**: A similar task from the textbook solved in $O(n^2)$ complexity: Problem 2707: Extra characters in a string. Given a string **s** and a dictionary of words, our goal is to break **s** into non-overlapping substring, each of which should be present in the dictionary. Any extra characters in **s** that do not belong to any of the substrings should be minimized. We need to determine the minimum number of extra characters left over when breaking **s** optimally.
The problem can be solved in $O(n^2)$ as follows:
```
def minExtraChar(s: str, dictionary: set) -> int:
    n = len(s)
    dp = [float('inf')] * (n + 1)
    dp[n] = 0

    (...)
```

**Step 3 (Reason)**: Here is how I plan to solve the task: Iterate through string characters (outer loop); Start a new substring and track unique characters using a set; Add characters until finding a repeat; Update the max length if current is longer; Repeat for each starting position.

**Step 4 (Reason)**: Below is the code:
```
def lengthOfLongestSubstring(s: str) -> int:
    n, maxLength = len(s), 0

    for i in range(n):
        charSet, currentLength = set(), 0

        for j in range(i, n):
            if s[j] not in charSet:
                charSet.add(s[j])
                currentLength += 1
                maxLength = max(maxLength, currentLength)
            else:
                Break
    return maxLength
```
❌

**(b) Critic-Guided Planning with Retrieval-Augmentation (CR-Planner)**

**Step 1 (Sub-Goal Selection)**: $g^g$(Reason)=11.2; $g^g$(GenQuery)=2.7; $g^g$(Retrieve)=-11.6.

**Step 2 (Execution Selection)**: $g^e_{RATIONALE}$(Rationale 1: The optimal time complexity is $O(n^2)$)=2.1; $g^e_{RATIONALE}$(Rationale 2: The optimal time complexity is $O(n)$)=13.1; $g^e_{RATIONALE}$(Rationale 3: The optimal time complexity is $O(n^3)$)=-4.6.

**Step 3 (Sub-Goal Selection)**: $g^g$(Reason)=0.7; $g^g$(GenQuery)=10.9; $g^g$(Retrieve)=-12.1.

**Step 4 (Execution Selection)**: $g^e_{QUERY}$(Query 1: Given a string **s**, find the length of the longest substring without repeating characters in optimal time complexity)=-5.6; $g^e_{QUERY}$(Query 2: Sliding window technique string problems)=1.1; $g^e_{QUERY}$(Query 3: Max length substring with unique characters with $O(n)$ complexity)=13.4.

**Step 5 (Sub-Goal Selection)**: $g^g$(Reason)=-14.2; $g^g$(GenQuery)=-6.8; $g^g$(Retrieve)=10.1.

**Step 6 (Execution Selection)**: $g^e_{DOC}$(Document 1: The longest strings without substring without repeating characters are …)=0.2; $g^e_{DOC}$(Document 2: The complexity of this is definitely $O(n)$ since they are only moving forward together through the string …)=0.1; $g^e_{DOC}$(Document 3: The intuition behind the solution is to iteratively find the longest substring without repeating characters by maintaining a sliding window approach…)=9.2.

**Step 7 (Sub-Goal Selection)**: $g^g$(Reason)=12.2; $g^g$(GenQuery)=0.4; $g^g$(Retrieve)=-3.3.

**Step 8 (Execution Selection)**: $g^e_{RATIONALE}$(Rationale 1: The retrieved document is not sufficient for problem solving. Therefore, a second-level retrieval is required …)=0.1; $g^e_{RATIONALE}$(Rationale 2: Based on the optimal time complexity and retrieved document, here is how I plan to solve the task …)=15.2; $g^e_{RATIONALE}$(Rationale 3: Here is the code …)=-11.1.

**(More Steps …)**

**Step $n$ (Execution Selection)**: $g^e_{RATIONALE}$(Rationale 1: …)=-5.3; $g^e_{RATIONALE}$(Rationale 2: …)=0.6; $g^e_{RATIONALE}$(Rationale 3: Here is the code:
```
def lengthOfLongestSubstring(s: str) -> int:
    n, charSet, left, maxLength = len(s), set(), 0, 0

    for right in range(n):
        while s[right] in charSet:
            charSet.remove(s[left])
            left += 1
        charSet.add(s[right])
        maxLength = max(maxLength, right - left + 1)
    return maxLength   )=10.1.
```
✅

Figure 1: Comparison between (a) chain-of-thought reasoning (Wei et al., 2022) with retrieval-augmented generation (Lewis et al., 2020) and (b) critic-guided planning with retrieval-augmentation or CR-Planner (this work). $g(\cdot)$ indicates the critic model (or value function) that assigns a reward (or value) to an action (see Equation 2). Texts in (b) highlighted in green are actions selected at each step. For succinct presentation, only pivotal steps are shown in the figure.

soning and retrieval processes with specially fine-tuned critic models. An example of CR-Planner in action is illustrated in Figure 1 (b). CR-Planner begins with **Sub-Goal Selection**, where it selects a sub-goal from three options: REASON (generating rationales), GENQUERY (generating search queries), and RETRIEVE (retrieving documents), based on reward scores estimated by a critic model, *the sub-goal critic*. After choosing the sub-goal of REASON in Step 1, CR-Planner proceeds to **Execution Selection**, sampling candidate rationales and using another critic model, *the execution critic* to select the optimal rationale, which in this case is "The optimal time complexity is $O(n)$." The framework alternates between sub-goal and execution selection until the final answer is reached, with each step guided by the appropriate critic model. CR-Planner integrates a large general generator model (*e.g.,* GPT-4) with small critic models (*e.g.,* Llama-3-8B) fine-tuned for domain-specific critiquing. When executing a sub-goal, the generator

model generates multiple candidates (*e.g.,* rationales or search queries), and the execution critic (specially trained for each candidate type) selects the most promising option. This design leverages the generation strengths of large LLMs while keeping the critic models efficient and trainable with domain-specific *(critiquing)* knowledge. To optimize planning performance for sub-goal and execution selection in each domain, critic models are trained separately using reasoning and retrieval trajectories labeled with step-wise rewards. Given the scarcity and cost of human-annotated data (Lightman et al., 2024), we utilize Monte Carlo Tree Search (MCTS) (Browne et al., 2012) to simulate trajectories, estimate long-term rewards, and propagate them back through steps. This approach efficiently trains the critic models to guide reasoning and retrieval at each step.

In summary, our key contributions are: (1) We introduce CR-Planner, a novel framework designed to tackle domain-knowledge-intensive and

reasoning-heavy problems by employing specially fine-tuned critic models that guide both reasoning and retrieval processes through planning; (2) We propose using MCTS to effectively collect training data for the critic models, enhancing their ability to estimate the long-term impact of an action. (3) We perform experiments on challenging tasks that require domain knowledge and complex reasoning, including competitive programming, math reasoning, and complex retrieval. CR-Planner outperforms the baseline by 10.06% on average.

## 2 Related Work

LLMs exhibit strong reasoning capabilities, showing promising performance on most logical reasoning datasets (Liu et al., 2023; Qin et al., 2023). However, they often struggle with complex tasks that require structured thinking or planning (Huang and Chang, 2023), prompting researchers to develop more sophisticated reasoning schemes. CoT (Wei et al., 2022) improves reasoning by prompting LLMs to articulate step-by-step processes. Tree-of-Thought (Yao et al., 2023a) then generalizes further by breaking down a CoT into coherent units of "thoughts", enabling the LLM to consider multiple reasoning paths and self-evaluate. To further improve LLMs in planning-based reasoning, process supervision has shown promise. RAP (Lightman et al., 2024) uses a world model to estimate future rewards of reasoning steps, providing step-wise guidance for reasoning processes. Jiao et al. (2024) learns planning-based reasoning through direct preference optimization (DPO) (Rafailov et al., 2023) on collected trajectories, which are ranked according to synthesized process rewards. As a result, tuned 7B models can surpass GPT-3.5-Turbo. However, this method requires training the base model, which limits its applicability to larger, closed-source models. In comparison, CR-Planner trains external critic models, offering flexibility for use with any base model.

Besides reasoning improvements, RAG effectively reduces hallucinations (Huang et al., 2023) by introducing external knowledge. Specifically, the RAG process involves 3 sub-tasks: pre-retrieval analysis, query generation, and document selection. Currently, most methods attempt to optimize the subtasks separately. Self-ask (Press et al., 2023) optimizes pre-retrieval analysis by breaking down the original problem into sub-problems. Chain-of-Knowledge (Li et al., 2024) rewrites natural-

language questions to database queries for more precise retrieval with structured knowledge. RePlug (Shi et al., 2023) improves document selection with a fine-tuned retriever. As these methods optimize sub-tasks locally, the single-task improvements may not constitute the globally optimal solution. In comparison, CR-Planner trains the critic model by learning the rewards of each individual action for overall performance.

## 3 Critic-Guided Planning with Retrieval-Augmentation

We introduce the critic-guided planning with retrieval-augmentation framework (CR-Planner) to address challenging tasks that are both domain-knowledge-intensive and reasoning-heavy. As shown in Figure 2, CR-Planner operates with two key components during inference: **(1) Sub-Goal Selection**: Given the current state, it employs a sub-goal critic model to determine the sub-goal among REASON, GENQUERY, and RETRIEVE that leads towards the desired answer. **(2) Execution Selection**: Upon selecting a sub-goal, CR-Planner undertakes multiple possible executions to realize the sub-goal (*e.g.,* generating multiple search queries for GENQUERY). Then, an execution critic model specifically designed to assess the executions for the sub-goal is employed to select the optimal execution among these candidates. In this process, a general generator model collaborates with multiple specialized critic models, leveraging the generator's strengths for initial plan generation while relying on the fine-tuned critics to guide optimal routing. To ensure the training data for the critic models is comprehensive and represents global reward information, we employ MCTS to collect the training data.

### 3.1 Problem Formulation

We define the associated planning environment of CR-Planner as a Markov Decision Process (MDP) represented by the tuple $(\mathcal{S}, \mathcal{A}_s, \mathcal{P}, \mathcal{R}, T)$, where:

- $\mathcal{S}$ represents the state space. Specifically, the state at timestamp $t$, denoted by the random variable $s_t$, comprises an action-observation trajectory history $(o_0, a_0, ..., a_{t-1}, o_t)$, where $a_{t-1}$ is the action taken at timestep $t-1$, and $o_t$ is the observation made after. Observation $o$ can be REASON, GENQUERY, or RETRIEVE in sub-goal selection stage, and RATIONALE, QUERY, or DOC in execution selection stage. Additionally, a state

**Question**

Given a string **s**, find the length of the longest substring without repeating characters in optimal time complexity.

**Legend**

| | |
|---|---|
| `Root` Root node | **S** Simulation score ╱▶ Simulation ⁙ Critic Model ⑤ Generator |
| ☐ SubGoal observation *o*: Reason, GenQuery, Retrieve | |
| ◯ Execution observation *o*: Rationale, Query, Doc | ↗ Action ▤ Rule-based function  ▣ Retriever |

**Training Data Collection via MCTS**

(a) Selection  (b) Expansion  (c) Simulation  (d) Backpropagation

**Value Model Training**

$g$(currentState, Action)

**currentState:**
Previous Thoughts;
Current Action.

→ Reward

**Inference**

Step 1: Sub-Goal Selection

Step 2: Execution Selection

If SubGoal==Reason (Sampling and choosing Rationales)

If SubGoal==GenQuery (Sampling and choosing Queries)

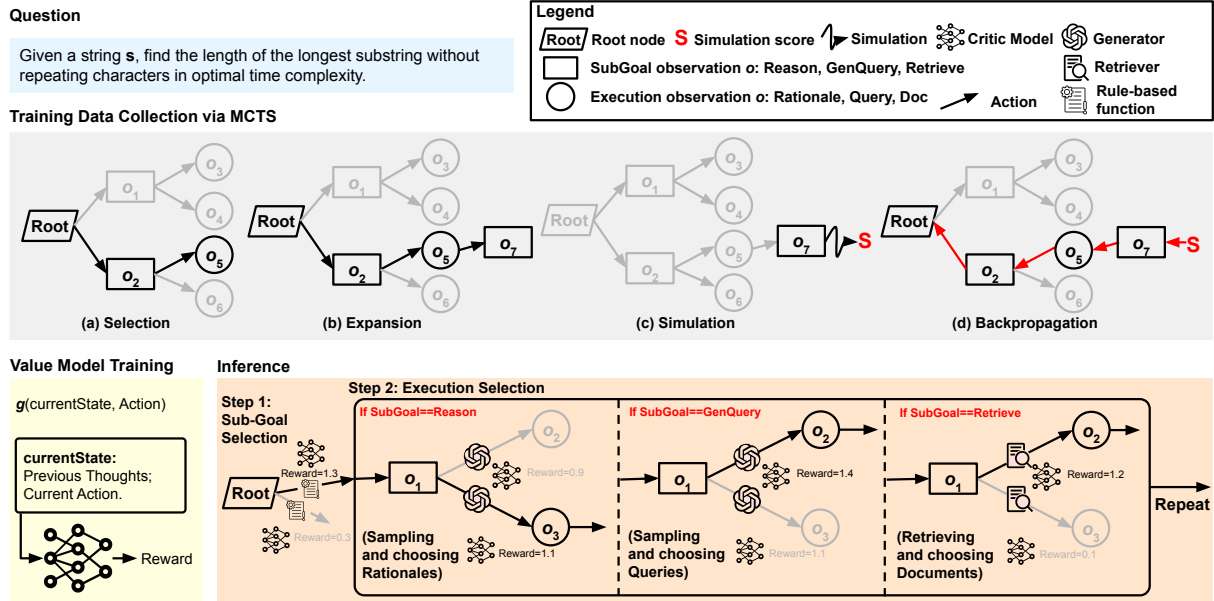If SubGoal==Retrieve (Retrieving and choosing Documents)

Repeat

Figure 2: The retrieval-augmented and critic-guided planning (CR-Planner) framework. The figure illustrates training data collection via MCTS, critic model training, and inference. For succinct presentation, SUBGOAL observations (REASON, GENQUERY, and RETRIEVE) are shown as labeled rectangles and EXECUTION observations (RATIONALE, QUERY, and DOC) as labeled circles. A state $s_t$ includes all preceding nodes (observations) and arrows (actions) up to the last node.

is named after its last observation, *e.g.,* RATIO-NALE state $s_t$ means $o_t$ is a RATIONALE.

- $\mathcal{A}_s$ represents the actions available at each state. For example, the actions available at the sub-goal selection stage, *i.e.,* at the Root state or after observing an outcome of an execution selection are: *reasoning*, *querying*, and *retrieving*. The possible actions available at the execution selection stage arise from the sampling for the corresponding sub-goal (*i.e.,* temperature sampling for REASON and GENQUERY, and top-k candidates for RETRIEVE). For example, Steps 1 and 2 in Figure 1 (b) illustrate the REASON and RATIONALE observations generated following the sub-goal selection and execution selection stages, respectively.

- The state transition $\mathcal{P}$ defines how the states evolve after an action is taken. In our context, state transitions are determined and handled by different functions depending on the current state. During the execution selection stage, a REASON or GENQUERY state transits to the respective RATIONALE or QUERY execution outcomes via the distribution defined by a large general generator model $f_{gen}(\cdot)$. Similarly, a RETRIEVE state transits to a DOC state via a retriever $f_{retr}(\cdot)$. During the sub-goal selection stage, the transition is

more straightforward and done via a rule-based function $f_{rule}(\cdot)$, *e.g.,* selecting *reasoning* action transits to a REASON state.

- The reward function $\mathcal{R}(s_t, a)$ specifies the expected reward received after taking an action $a_t$ at state $s_t$. In our context, fine-tuned critic models estimate the rewards and guide the decision-making process by encouraging actions that contribute the most towards solving the MDP. Details of the critic models are provided in Section 3.2.

- Lastly, $T$ represents the maximum number of steps that can occur within the MDP.

Solving the MDP requires generating an optimal plan in the form of a trajectory: $\tau* = (s_0, a_0, ..., s_t, a_t, ..., s_{T-1}, a_{T-1}, s_T)$ that maximizes the total expected rewards. [1]

## 3.2 Inference of CR-Planner

CR-Planner employs critic models at each time step to guide the decision-making process. Specifically, at time step $t$, given the current state $s_t$, the critic model $g$ assesses the available actions $\mathcal{A}_{s_t}$ and helps select an action $a_t$ that maximizes the expected reward.

---

[1]Details of state types and action spaces are in Appendix E Table 9.

**Action selection using the critic models.** At timestamp $t$, the policy model $\pi$ determines the next action as:

$$a_t = \pi(s_t) = \arg\max_{a \in \mathcal{A}_{s_t}} \mathcal{R}(s_t, a). \tag{1}$$

The action space $\mathcal{A}_{s_t}$ varies depending on $s_t$. As previously discussed in Section 3.1 and outlined in Table 9, for a state in the sub-goal stage, the action space leads to possible executions of that sub-goal, while for a state in the execution stage, the action space leads to the possible subsequent sub-goals. $\mathcal{R}(s_t, a)$ is the expected reward when taking action $a$ in state $s_t$ and estimated by the critic models:

$$\mathcal{R}(s_t, a) = \begin{cases} g^e_{\text{RATIONALE}}(s_t, a), & \text{if } s_t = \text{REASON state} \\ g^e_{\text{QUERY}}(s_t, a), & \text{if } s_t = \text{GENQUERY state} \\ g^e_{\text{DOC}}(s_t, a), & \text{if } s_t = \text{RETRIEVE state} \\ g^g(s_t, a), & \text{otherwise.} \end{cases} \tag{2}$$

Specifically, distinct critic models are utilized for different state types: $g^g(\cdot)$ is for determining the next sub-goal at the current execution state (*i.e.,* the inference Steps 1 in Figure 2), and $g^e(\cdot)$ is for evaluating different execution candidates at the current sub-goal state (*i.e.,* the inference Step 2 in Figure 2). Additionally, according to the sub-goal states, $g^e(\cdot)$ has three variants $g^e_{\text{RATIONALE}}$, $g^e_{\text{QUERY}}$, and $g^e_{\text{DOC}}$, correspondingly evaluating rationales, queries and the retrieved documents.

**State transition with the selected action.** Once $a_t$ is determined and executed, the state is then transited from $s_t$ to $s_{t+1} = (s_t, a_t, o_{t+1})$, where

$$o_{t+1} = \begin{cases} f_{gen}(s_t, a_t), & \text{if } s_t = \text{REASON or GENQUERY state} \\ f_{retr}(s_t, a_t), & \text{if } s_t = \text{RETRIEVE state} \\ f_{rule}(s_t, a_t), & \text{otherwise.} \end{cases} \tag{3}$$

As mentioned in Section 3.1, given the current state $s_t$ and action $a_t$, we employ three specific functions to generate different types of outcomes. The generator $f_{gen}(\cdot)$ generates either a RATIONALE or QUERY. The retriever $f_{retr}(\cdot)$ outputs a DOC. Last but not least, the rule-based function $f_{rule}(\cdot)$ outputs a SUBGOAL. The SUBGOAL is a predefined natural language. For example, a REASON thought is "The next step is to generate a rationale".

**Termination conditions and the final answer.** This process continues until one of two conditions is met. The process ends at step $t$ if the observation $o_t$ includes the complete answer. Otherwise, if $t$ equals $T$ and $o_t$ does not contain the final answer, an extra step occurs to force the model to conclude the answer. In this case, a concluding answer is generated using an LLM.

## 3.3 The Critic Models

The CR-Planner framework relies on critic models to evaluate actions and guide sub-goal and execution selection. Accurate assessment of each action's contribution to the problem-solving process is essential, making high-quality training data critical. To generate such data, we use MCTS, which explores long-term impacts of actions while balancing exploration and exploitation. By simulating diverse action-observation trajectories, MCTS creates a rich dataset, helping critic models distinguish effective actions from suboptimal ones.

**Collecting data via MCTS.** As shown in Figure 2, MCTS consists of the four key steps: **(1) Selection.** Starting from the root state $s_0$, the algorithm selects child node (with observation) recursively based on the Upper Confidence Bound (UCB1) that balances exploration and exploitation. The UCB1 value for $o_i$ is computed as $\frac{v_i}{n_i} + c\sqrt{\frac{\ln n_p}{n_i}}$, where $v_i$ is the cumulative rewards of $o_i$, $n_i$ is the number of times $o_i$ has been visited, and $n_p$ denotes the number of visits to the parent thought of $o_i$. This process continues until it reaches a node that is not fully expanded or a terminal node. **(2) Expansion.** If the selected $o_i$ is not terminal and has unexplored child nodes, MCTS expands the tree by adding one or more of these unexplored child nodes. This represents exploring new actions available from the current action space $\mathcal{A}_{s_t}$. **(3) Simulation.** From the newly added observation, MCTS simulates a playthrough to a terminal state by employing a generative model $f_{gen}(\cdot)$ to generate the final answer based on existing observations. This simulation estimates the potential outcome from the observation. **(4) Backpropagation.** The result of the simulation is then propagated back up the tree. Each node along the path to the root updates its statistics, including visit counts and total reward, which informs future selection decisions by reflecting the observed outcomes. For each data point in the training dataset, we run MCTS for $N$ steps and collect pairwise data from the final state for each observation type. In particular, a chosen observation $o_i$ is the one with the highest score, while a rejected observation $o'_i$ is one of the observations sharing the same parent node but a lower score. For critic model $g^e_{\text{RATIONALE}}(\cdot)$, we collect $\mathcal{D}^{\text{RATIONALE}} = \{(O^{\text{RATIONALE}}_i, o_i, o'_i)...\}$, where $O^{\text{RATIONALE}}_i$ represents previous RATIONALEs along the trajectory before the current RATIONALE $o_i$. It is

crucial to evaluate $o_i$ considering all prior rationales. The critic model $g_{\text{QUERY}}^e(\cdot)$ uses $\mathcal{D}^{\text{QUERY}} = \{(o_i^{\text{RATIONALE}}, o_i, o_i')...\}$, where $o_i^{\text{RATIONALE}}$ is one immediately preceding RATIONALE of QUERY $o_i$. For the critic model $g_{\text{DOC}}^e(\cdot)$, we have $\mathcal{D}^{\text{DOC}} = \{(o_i^{\text{RATIONALE}}, o_i^{\text{QUERY}}, o_i, o_i')...\}$, where $o_i^{\text{RATIONALE}}$ and $o_i^{\text{QUERY}}$ are the immediately preceding RATIONALE and QUERY of DOC $o_i$. Lastly, the SUBGOAL critic model $g^g(\cdot)$ uses $\mathcal{D}^{\text{SUBGOAL}} = \{(O_i, o_i, o_i')...\}$, where $O_i$ represents all previous observations of any type along the trajectory.

**Training.** For each of the collected training datasets described above, we train a dedicated critic model as shown in Figure 2. Following Burges et al. (2005) and Ouyang et al. (2022), we employ pairwise ranking loss to optimize the parameters.

# 4 Experiments

## 4.1 Setup

**Models.** GPT-4 (gpt-4o-2024-05-13) is utilized as the black-box LLM for generation during both inference and training data collection. We fine-tune Skywork-Reward-Llama-3.1-8B (Skywork, 2024) with LoRA (Hu et al., 2021) as critic models, leveraging its capability to score complex scenarios like mathematics and coding. [2]

**Baselines.** We compare CR-Planner with both commonly used baselines and state-of-the-art methods to offer a comprehensive evaluation: **(1) Standard prompting (Standard)** (Ouyang et al., 2022), which directly generates the answer. **(2) Chain-of-Thought (CoT)** (Wei et al., 2022), which generates multiple rationales before the final answer to enhance the models' reasoning ability. **(3) Reflexion** (Shinn et al., 2023), a framework uses linguistic feedback to further improve models' reasoning. **(4) Standard retrieval-augmented generation (RAG)** (Lewis et al., 2020), which retrieves relevant knowledge based on the problem itself and then lets the model to generate the final answer using both the problem and the retrieved knowledge. **(5) Chain-of-Knoweldge (CoK)** (Li et al., 2024), a state-of-the-art CoT-based framework designed to enhance prediction accuracy by retrieving and post-editing rationale at each step. All methods are zero-shot by default unless otherwise specified. [3]

---

[2]Details are in Appendix A.1.

[3]We exclude Self-RAG as a baseline because it requires training the base model, which is not feasible in our setup. This further highlights the flexibility of CR-Planner.

| Method | Bronze | Silver | Gold | Platinum | Overall |
|---|---|---|---|---|---|
| Standard | 18.70 | 6.00 | 3.17 | 0.00 | 10.10 |
| CoT | 21.95 | 8.00 | 4.76 | 0.00 | 12.38 |
| RAG | 17.07 | 4.00 | 1.59 | 0.00 | 8.47 |
| CoK | 15.45 | 5.00 | 1.59 | 0.00 | 8.14 |
| CR-Planner | **26.02** | **10.00** | **14.29** | **14.29** | **17.59** |
| Reflexion | 23.58 | 9.00 | 4.76 | 0.00 | 13.36 |
| *Retrieval+Reflexion** | - | - | - | - | *18.05* |
| CR-Planner+Reflexion | **34.15** | **16.00** | **14.29** | **14.29** | **22.80** |

Table 1: Pass@1 performances on USACO. The *Retrieval+Reflexion** result is from Shi et al. (2024).

## 4.2 Competitive Programming

**USACO benchmark.** Computing Olympiads require advanced algorithmic reasoning, problem-solving skills, and efficient code generation, often supplemented by retrieving knowledge from textbooks or similar problems. Following the baseline methods outlined in the USACO benchmark (Shi et al., 2024), we use both textbooks and a problem bank as external knowledge sources. [4]

**Results and observations. (1) CR-Planner outperforms all baselines consistently.** Table 1 presents the results for USACO using various methods. CR-Planner significantly outperforms all baseline methods, achieving a 7.49% improvement in overall performance compared to standard prompting. This highlights the effectiveness of CR-Planner. **(2) Reasoning-driven methods offer limited improvements.** We observe that reasoning-driven methods like CoT and Reflexion do improve the performances of the standard prompting method on bronze, silver, and gold problems, reaffirming that intermediate rationales and critique-based reasoning aid in solving reasoning tasks (Wei et al., 2022; Shinn et al., 2023). However, the improvements are trivial, and these methods fail to improve performance on platinum-level problems. We attribute this to the model's limited knowledge of the tasks or the generation of faulty rationales and critiques. **(3) Faulty retrieval hinders performance.** We observe that both standard RAG and CoK perform worse than the standard prompting method, consistent with the findings of Yao et al. (2023b) and Shi et al. (2024). This decline in performance can be attributed to the quality of retrieval. As demonstrated in Figure 1, if the retrieved example is irrelevant to the original problem, it may mislead the model into generating an

---

[4]Details of USACO benchmark and external knowledge are in Appendix B.1.

| Method | TheoremQA-Math |
|---|---|
| Standard | 39.81 |
| CoT | 41.75 |
| Reflexion | 40.29 |
| RAG | 44.17 |
| CoK | 45.15 |
| CR-Planner | **53.40** |

Table 2: Results (accuracy) on TheoremQA-Math.

| Method | StackBio nDCG@10 | StackEcon nDCG@10 |
|---|---|---|
| BM25 | 19.20 | 14.90 |
| CoT | 21.06 | 16.33 |
| CoK | 20.82 | 17.45 |
| CR-Planner | **29.51** | **22.80** |

Table 3: Results on complex domain retrieval.

| Method | USACO |
|---|---|
| Claude-3.5 | 9.12 |
| CR-Planner w/ Claude-3.5 | 13.68 |
| Llama-3.1 | 7.49 |
| CR-Planner w/ Llama-3.1 | 10.10 |

Table 4: CR-Planner with various base models.

incorrect answer. Additionally, we notice that CoK performs worse than RAG due to its reliance on multiple retrievals at individual steps, increasing the likelihood of misleading information being introduced and leading to a faulty final answer. **(4) CR-Planner improves harder problems.** CR-Planner notably boosts the performances on gold- and platinum-level problems. As aforementioned, while CoT offers minor improvements, it falls short on more difficult problems, and retrieval can hinder performance due to irrelevant knowledge. In contrast, CR-Planner employs critic models to guide both the reasoning and retrieval through the process, leading to non-trivial improvements at the two highest levels of programming problems. **(5) CR-Planner is orthogonal with other methods.** Reflexion executes the initially generated code and uses the execution results of a few test cases as linguistic feedback to revise the code. CR-Planner works orthogonal with such methods, leading to a significant improvement of 9.44%, further highlighting the effectiveness of critic-guided planning with retrieval-augmentation.

## 4.3 Theorem-Driven Math Problems

**TheoremQA-Maths.** Theorem-driven math problems require both complex reasoning and knowledge of math theorems. To evaluate CR-Planner, we use TheoremQA-Math (Su et al., 2024), a dataset of 206 solvable problems that emphasize reasoning and theorem application, making it highly relevant to this study. Following the BRIGHT benchmark (Su et al., 2024), we employ a collection of processed documents sourced from high-quality STEM datasets as external knowledge sources. [5]

**Results and observation.** Similar to competitive programming, as shown in Table 2, we observe a notable performance improvement from CR-Planner, with 13.59% on TheoremQA-Math

compared to standard prompting method. This further demonstrates the effectiveness of CR-Planner in tasks requiring knowledge retrieval and complex reasoning. Interestingly, Reflexion exhibits inferior performance compared to CoT, which we attribute to Reflexion's tendency to potentially revise initially correct answers into incorrect ones. Furthermore, in contrast to their behavior in the USACO benchmark, retrieval methods, such as standard RAG and CoK, do enhance performance in this task. We attribute this to the shorter context of the retrieved documents in the math domain. With shorter retrieved documents, the base model is easier to determine which information to incorporate. Nevertheless, CR-Planner maximizes the benefits of both retrieval and reasoning, leading to the best performance improvement.

## 4.4 Reasoning-Heavy Domain Retrieval

**StackBio and StackEcon.** To evaluate reasoning-heavy domain retrieval, we use the StackBio and StackEcon datasets from the BRIGHT benchmark (Su et al., 2024), which include 103 biology and 103 economics questions sourced from StackExchange. External sources can include any accessible web content. [6]

**Results and observations.** As shown in Table 3, CR-Planner consistently improves over the standard BM25 method by 10.31% and 7.9% on StackBio and StackEcon, respectively. CoK improves the standard BM25 method, which indicates that reasoning before retrieval is crucial in such reasoning-heavy domain retrieval tasks. However, CoK does not consistently enhance performance; for instance, it performs worse than CoT on StackBio. We attribute this to the potential noise introduced by multiple suboptimal retrieval results. These observations further highlight the effectiveness of the critic models in RC-Planner.

---

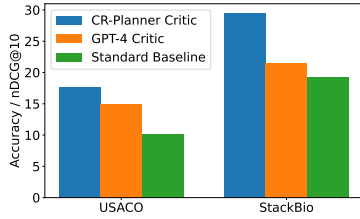[5]Details of TheoremQA-Maths benchmark and external knowledge are in Appendix B.2.

[6]Details of StackBio and StackEcon benchmark and external knowledge are in Appendix B.2.

Figure 3: Performances of different critic models.

| Method | TheoremQA-Math |
|---|---|
| Standard | 31.43 |
| Llama-3.1-Fine-tuned | 34.10 |
| CR-Planner | 38.62 |

Table 5: Fine-tuning critic models vs. base model.

| Method | USACO |
|---|---|
| Standard | 10.10 |
| CR-Planner w/o Retrieval | 14.33 |
| CR-Planner | 17.59 |

Table 6: CR-Planner with and without retrieval.

## 5 Analysis

### 5.1 Domain-Specific Critic Models

Critic models are crucial in CR-Planner, guiding sub-goal and execution selection during inference. Previous works often use proprietary LLMs as critics (*e.g.,* GPT-4), leveraging in-context learning to evaluate actions (Gou et al., 2024; Zhao et al., 2024). This sub-section compares CR-Planner's performance with fine-tuned models versus GPT-4 (gpt-4o-2024-05-13) as critics on USACO and StackBio datasets, with results in Figure 3. While GPT-4 improves over the baseline, CR-Planner performs better with fine-tuned critics, particularly in domain-specific tasks like StackBio. This highlights the importance of domain-specific fine-tuning and validates CR-Planner's choice of fine-tuned critic models.

### 5.2 Flexibility of Critic Models on Various Base Models

Unlike methods like Self-RAG (Asai et al., 2024), CR-Planner does not require fine-tuning the base model, making it adaptable to both open- and closed-source models. This subsection highlights the effectiveness of the critic models on another closed-source model, Claude-3.5 (claude-3-5-sonnet), and an open-source model, Llama-3.1 (Llama-3.1-70B-Instruct). As shown in Table 4, CR-Planner improves Claude-3.5 by 4.56% and Llama-3.1 by 2.61%, though these gains are smaller than the 7.49% boost with GPT-4. This is likely because the critic models were trained on GPT-4 data, making them better suited for it during inference. Nonetheless, the plug-and-play design of CR-Planner's critic models offers a promising way to distill capabilities from powerful LLMs to smaller models, enabling them to benefit from high-quality guidance without generating robust MCTS trajectories themselves.

### 5.3 Fine-tuning Critic Models vs. Base Model

The training of CR-Planner's critic models involves conducting MCTS to gather reasoning and retrieval trajectories for constructing reward datasets. This process may introduce additional computational overhead compared to in-context learning baselines. To enable a more comprehensive comparison, we present results of fine-tuning the base model with the same MCTS-collected data using LoRA. However, due to cost constraints and the inability to fine-tune closed-source LLMs, we cannot include these results in the main experiments. In this subsection, we use Llama-3.1 (Llama-3.1-70B-Instruct) as the base model. As shown in Table 5, fine-tuning Llama-3.1 enhances performance over standard prompting, though it does not surpass CR-Planner. Additionally, CR-Planner requires less computation compared to fine-tuning the entire base model.

### 5.4 Retrieve or Not to Retrieve

Tackling complex domain-specific tasks such as competitive programming requires extensive reasoning as well as advanced algorithmic knowledge, which base models may not inherently possess. In this subsection, we examine the importance of accurately retrieving external knowledge to assist in solving competitive programming problems. We instruct the model to concentrate solely on reasoning, employing the reasoning critic model $g_{\text{REASON}}^g$ to select a rationale for each reasoning step. As shown in Table 6, the performance without retrieval is lower. However, as discussed in Section 4.2 and by Shi et al. (2024), inaccurate retrieval could impair performance. This emphasizes the critical role of accurate retrieval and the overall effectiveness of CR-Planner.

| Method | USACO |
|---|---|
| Standard | 10.10 |
| Vanilla MCTS | 12.42 |
| CR-Planner | 17.59 |

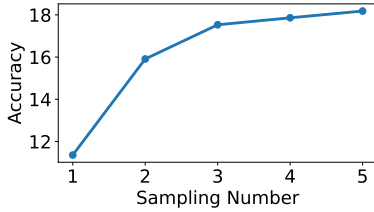Table 7: CR-Planner vs. vanilla MCTS.



Figure 4: Performances of various execution sampling.

## 5.5 Vanilla MCTS

In this subsection, we evaluate the performance of directly applying MCTS during inference. Since scoring necessitates an answer, we experiment with the USACO benchmark, using test cases to score the simulation results. As shown in Table 7, MCTS performs worse than CR-Planner. This is likely because execution sampling for a single data point can overfit to local optima based on the test cases, potentially leading to suboptimal performance on hidden test cases. In contrast, CR-Planner benefits from leveraging learned patterns from the training data, resulting in improved outcomes.

## 5.6 Execution Sampling

Throughout both training and inference of CR-Planner, executing sub-goals involves sampling several candidates. Increasing the number of candidates can improve the likelihood of selecting a better option. In this subsection, we study how varying the number of candidates sampled for sub-goal execution during inference impacts performance. Due to cost concerns, we do not conduct ablation studies for the training phase. As shown in Figure 4, the improvements on USACO are substantial when increasing from one to two, but converge around three. This limitation likely stems from the

generator model's reasoning capabilities and the retriever's accuracy. Without fine-tuning both generator and retriever models, further performance gains are difficult to achieve. Therefore, to balance performance and cost, we select three as the sampling number for the main experiments.

## 5.7 Single Unified vs. Task-Specific Critics

While training a unified critic is possible, we find that smaller reward models, such as `Skywork-Reward-Llama-3.1-8B` used in our experiments, struggle to generalize across all critic tasks effectively. As shown in Table 8, a single unified critic underperforms compared to task-specific critics.

## 6 Conclusions

In this paper, we present critic-guided planning with retrieval-augmentation (CR-Planner), a novel framework for handling domain-knowledge-specific and reasoning-heavy tasks by leveraging fine-tuned critic models to guide both the reasoning and retrieval processes. We further employ the Monte Carlo Tree Search for systematic data collection to enhance the training of the critic models. Our approach, validated across challenging domains like competitive programming, math reasoning, and complex domain retrieval tasks, has shown substantial performance improvements over existing methods. By combining the strengths of large generalist models with domain-specific fine-tuned critics, CR-Planner offers a flexible and scalable solution for solving problems that require both intricate reasoning and accurate knowledge retrieval.

## Limitations

A limitation of CR-Planner is that both its training and inference processes involve sampling multiple candidates when executing sub-goals. While increasing the number of candidates can enhance the chances of selecting a better option, the framework itself does not enhance the capabilities of the base model. Consequently, for particularly challenging problems, it is likely that none of the sampled candidates will be correct, which limits overall performance. Additionally, running MCTS for data collection is computationally expensive. Developing more efficient methods for MCTS data collection is a potential direction for future research.

| Method | Bronze | Silver | Gold | Platinum | Overall |
|---|---|---|---|---|---|
| Standard | 18.70 | 6.00 | 3.17 | 0.00 | 10.10 |
| Single Unified | 21.14 | 8.00 | 4.76 | 4.76 | 12.38 |
| Task-Specific | **26.02** | **10.00** | **14.29** | **14.29** | **17.59** |

Table 8: Single Unified vs. Task-Specific Critics.

## Ethical Impact

We do not foresee any potential ethical issues with our proposed method.

## Acknowledgements

## References

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *Proceedings of ICLR*.

Cameron Browne, Edward Powley, Daniel Whitehouse, Simon Lucas, Peter Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. 2012. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in Games*.

Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. 2005. Learning to rank using gradient descent. *Proceedings of ICML*.

Hailin Chen, Fangkai Jiao, Xingxuan Li, Chengwei Qin, Mathieu Ravaut, Ruochen Zhao, Caiming Xiong, and Shafiq Joty. 2024. Chatgpt's one-year anniversary: Are open-source large language models catching up? *arXiv preprint arXiv:2311.16989*.

Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023. Theoremqa: A theorem-driven question answering dataset. *Proceedings of EMNLP*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2024. Critic: Large language models can self-correct with tool-interactive critiquing. *Proceedings of ICLR*.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards reasoning in large language models: A survey. *Proceedings of ACL*.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*.

Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, Yiwei Qin, Yan Ma, Jiadi Su, Yixiu Liu, Yuxiang Zheng, Shaoting Zhang, Dahua Lin, Yu Qiao, and Pengfei Liu. 2024. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent ai. *arXiv preprint arXiv:2406.12753*.

Shayekh Islam, Md Rahman, K Hossain, Enamul Hoque, Shafiq Joty, and Md Parvez. 2024. Open-RAG: Enhanced Retrieval Augmented Reasoning with Open-Source Large Language Models. *Findings of EMNLP*.

Fangkai Jiao, Chengwei Qin, Zhengyuan Liu, Nancy F. Chen, and Shafiq Joty. 2024. Learning planning-based reasoning by trajectories collection and process reward synthesizing. *Proceedings of EMNLP*.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Proceedings of NIPS*.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for "mind" exploration of large language model society. *arXiv preprint arXiv:2303.17760*.

Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing. 2024. Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting over heterogeneous sources. *Proceedings of ICLR*.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let's verify step by step. *Proceedings of ICLR*.

Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *Proceedings of ACL*.

Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*.

OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774.*

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155.*

Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350.*

Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is ChatGPT a general-purpose natural language processing task solver? *Proceedings of EMNLP.*

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Proceedings of NIPS.*

Quan Shi, Michael Tang, Karthik Narasimhan, and Shunyu Yao. 2024. Can language models solve olympiad programming? *arXiv preprint arXiv:2404.10952.*

Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652.*

Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366.*

Skywork. 2024. Skywork reward model series. *HuggingFace.*

Hongjin Su, Howard Yen, Mengzhou Xia, Weijia Shi, Niklas Muennighoff, Han yu Wang, Haisu Liu, Quan Shi, Zachary S. Siegel, Michael Tang, Ruoxi Sun, Jinsung Yoon, Sercan O. Arik, Danqi Chen, and Tao Yu. 2024. Bright: A realistic and challenging benchmark for reasoning-intensive retrieval. *arXiv preprint arXiv:2407.12883.*

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Proceedings of NIPS.*

Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective Retrieval Augmented Generation. *arXiv preprint arXiv:2401.15884.*

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. *Proceedings of NIPS.*

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models. In *Proceedings of ICLR.*

Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and Jingren Zhou. 2023. Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825.*

Chong Zhang, Yue Deng, Xiang Lin, Bin Wang, Dianwen Ng, Hai Ye, Xingxuan Li, Yao Xiao, Zhanfeng Mo, Qi Zhang, and Lidong Bing. 2025. 100 days after deepseek-r1: A survey on replication studies and more directions for reasoning language models. *arXiv preprint arXiv:2505.00551.*

Ruochen Zhao, Xingxuan Li, Yew Ken Chia, Bosheng Ding, and Lidong Bing. 2023a. Can chatgpt-like generative models guarantee factual accuracy? on the mistakes of new generation search engines. *arXiv preprint arXiv:2304.11076.*

Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023b. Verify-and-edit: A knowledge-enhanced chain-of-thought framework. *Proceedings of ACL.*

Ruochen Zhao, Wenxuan Zhang, Yew Ken Chia, Deli Zhao, and Lidong Bing. 2024. Auto arena of llms: Automating llm evaluations with agent peer-battles and committee discussions. *arXiv preprint arXiv:2405.20267.*

## A  Experiment Details

### A.1  Models

In our experiments, we employ GPT-4 (gpt-4o-2024-05-13) as the black-box LLM for generation during both inference and training data collection. Since CR-Planner requires the sampling of diverse RATIONALE and QUERY, we set the decoding temperature to 0.7. To ensure training and inference efficiency, we limit the sampling to three instances. For the critic models, we fine-tune Skywork-Reward-Llama-3.1-8B (Skywork, 2024) with LoRA (Hu et al., 2021), which was trained as a sequence classifier with the Skywork Reward Data Collection and excels at scoring in complex scenarios, such as mathematics and coding. The first logit value of the model output is used as the reward score of our critic models.

## B  Benchmark Details

### B.1  Competitive Programming

**USACO benchmark.**  USACO problems are categorized into four difficulty levels (*i.e.,* 123 bronze, 100 silver, 63 gold, and 21 platinum problems) and test various core skills, including complete search, binary search, and segment tree implementation. Typically, solving a USACO problem involves several steps: restating the problem in simple terms since many are framed within real-world contexts; retrieving relevant knowledge from textbooks or similar problems from a problem bank; conceptualizing the solution in plain English; drafting a pseudocode solution; and finally, producing the complete Python solution with comments. This multi-step process highlights the benchmark's suitability for evaluating complex reasoning and retrieval. This dataset is under the CC 4.0 license which is free to share and adapt.

**External knowledge.**  Following the baseline methods outlined in the USACO benchmark (Shi et al., 2024), we use both textbooks and a problem bank as external knowledge sources. The textbooks consist of 30 human-written chapters covering algorithmic concepts tailored specifically for the USA Computing Olympiad. The problem bank includes all other USACO problems except for the one currently being solved. Following Shi et al. (2024), we employ both textbooks and the problem bank as external sources for all methods. Additionally, we employ a BM25 retriever to execute the retrieval

process, obtaining relevant information from external knowledge sources.

### B.2  Theorem-Driven Math Problems

**TheoremQA-Maths.**  When tackling a new theorem-driven math problem, people often reference solved problems with similar reasoning logic. However, finding such problems can be challenging because even if two problems share similar reasoning logic, they might appear very different on the surface. Moreover, in theorem-driven math problems, the reasoning process is critical. A single flawed step in the logic can lead to wrong subsequent rationales and finally an incorrect final answer. In this task, we use the rewritten Math set from TheoremQA (Chen et al., 2023), named TheoremQA-Math, as introduced in the BRIGHT dataset (Su et al., 2024). TheoremQA-Math consists of 206 solvable questions that have been improved for fluency and coherence, with all questions requiring the application of math theorems (*e.g.,* the binomial theorems). To solve a problem in the TheoremQA-Math dataset, the process typically involves the following steps: understanding and restating the problem in simple terms; retrieving relevant knowledge from solved problems; conceptualizing the solution in plain English; and finally, generating the solution. Solving problems from TheoremQA-Math requires both complex reasoning and knowledge of Math theorems, making it pertinent to this paper. This dataset is under the MIT license which is free to share and adapt.

**External knowledge.**  Following the BRIGHT benchmark (Su et al., 2024), we employ a collection of processed documents sourced from high-quality STEM datasets, including GSM8K (Cobbe et al., 2021), GSM8K-RFT (Yuan et al., 2023), MATH (Hendrycks et al., 2021), AQuA-RAT (Ling et al., 2017), TheoremQA (Chen et al., 2023) and CAMEL-MATH (Li et al., 2023). To ensure efficient retrieval during both the training data collection and inference stages, we opt for the term-based retrieval method BM25, similar to what is used in competitive programming.

### B.3  Reasoning-Heavy Domain Retrieval

**StackBio and StackEcon.**  Complex domain queries often demand in-depth reasoning to identify relevant documents that go beyond simple surface-level matching. To evaluate models' ability in reasoning-heavy domain retrieval, we use biology-

and economics-related queries from the BRIGHT benchmark (Su et al., 2024), specifically Stack-Bio and StackEcon. Both StackBio and StackEcon contain 103 questions sourced from StackExchange, with the gold labels being the documents cited in the answers. As the evaluation metric is nDCG@10, which requires the top 10 documents, we set the number of retrieved documents to 10 when PC-Planner performs the final retrieval. This dataset is under the CC 4.0 license which is free to share and adapt.

**External knowledge.** In line with the BRIGHT benchmark (Su et al., 2024), external sources can include any accessible web content such as articles, tutorials, news, blogs, and reports. Since this information has already been gathered and incorporated into the benchmark, we employ BM25 for document retrieval to ensure efficiency.

## C  Prompts Used in Different Methods

### C.1  RC-Planner (Competitive Programming)

#### C.1.1  Instruction

Reason through the problem and think step by step. Specifically:
1. Restate the problem in plain English.
2. Conceptualize a solution first in plain English.
3. Write a pseudocode solution
4. Output the Python 3 solution to the problem. Make sure to wrap your code in "'python and '" Markdown delimiters, and include exactly one block of code with the entire solution.
No outside libraries are allowed.
    [BEGIN PROBLEM]

    [END PROBLEM]

#### C.1.2  SubGoal Selection

To proceed, below are the available actions:
    [REASON] - Provide a reasoning step.
    [GENQUERY] - Generate a query to retrieve information from external knowledge sources.
    [RETRIEVE] - Retrieve documents using the query.
    The next step is [].

#### C.1.3  Execution Selection - Rationale Sampling

Reason through the problem and think step by step. Specifically:
1. Restate the problem in plain English.
2. Conceptualize a solution first in plain English.

3. Write a pseudocode solution
4. Output the Python 3 solution to the problem. Make sure to wrap your code in "'python and '" Markdown delimiters, and include exactly one block of code with the entire solution.
No outside libraries are allowed.
    [BEGIN PROBLEM]

    [END PROBLEM]
    Generate one next reasoning step (*e.g.,* [BEGIN REASON] Restate the problem: ... [END REASON]). It starts with [BEGIN REASON] and ends with [END REASON]. Do not include the subsequent reasoning steps.

#### C.1.4  Execution Selection - Query Sampling

To verify or solve the reasoning step, I need additional information from external knowledge sources (*e.g.,* textbook). And I need to generate a query to get that information. The query needs to be conceptual but relevant to the reasoning step. The query should not contain any specific numbers or entities of the reasoning step. The query starts with [BEGIN QUERY] and ends with [END QUERY]. Stop the generation when the query is completed.
    [BEGIN REASON]

    [END REASON]

#### C.1.5  Force Termination

Based on the above rationales and information, generate python code directly to solve the problem. Make sure to wrap your code in "'python and '" Markdown delimiters, and include exactly one block of code with the entire solution. No outside libraries are allowed.

### C.2  CoT

Reason through the problem and think step by step. Specifically:
1. Restate the problem in plain English.
2. Conceptualize a solution first in plain English.
3. Write a pseudocode solution
4. Output the Python 3 solution to the problem. Make sure to wrap your code in "'python and '" Markdown delimiters, and include exactly one block of code with the entire solution.
No outside libraries are allowed.
    [BEGIN PROBLEM]

    [END PROBLEM]

## C.3 Chain-of-Knowledge

### C.3.1 Reasoning Generation

Reason through the problem and think step by step. Specifically:
1. Restate the problem in plain English.
2. Conceptualize a solution first in plain English.
3. Write a pseudocode solution
4. Output the Python 3 solution to the problem. Make sure to wrap your code in "'python and "' Markdown delimiters, and include exactly one block of code with the entire solution.
No outside libraries are allowed.
  [BEGIN PROBLEM]

  [END PROBLEM]

### C.3.2 Rationale Correction

The given sentence may have errors, please correct them based on the given external knowledge.
  Sentence: [Rationale]
Knowledge: [Knowledge]
Edited sentence:

### C.3.3 Next Rationale Generation

Reason through the problem and think step by step. Specifically:
1. Restate the problem in plain English.
2. Conceptualize a solution first in plain English.
3. Write a pseudocode solution
4. Output the Python 3 solution to the problem. Make sure to wrap your code in "'python and "' Markdown delimiters, and include exactly one block of code with the entire solution.
No outside libraries are allowed.
  [BEGIN PROBLEM]

  [END PROBLEM]

  [START PRECEDING RATIONALES]

  [END PRECEDING RATIONALES]

## C.4 Reflexion

### C.4.1 Actor

You are a Python writing assistant. You will be given your previous implementation of a function, a series of unit tests results, and your self-reflection on your previous implementation. Apply the necessary changes below by responding only with the improved body of the function. Do not include the signature in your response. The first line of your response should have 4 spaces of indentation so that it fits syntactically with the user provided signature.

Reflexion Actor generations follow the form:
[Instruction]
[Function implementation]
[Unit test feedback]
[Self-reflection]
[Instruction for next function implementation]

### C.4.2 Self-Reflection

You are a Python writing assistant. You will be given your previous implementation of a function, a series of unit tests results, and your self-reflection on your previous implementation. Apply the necessary changes below by responding only with the improved body of the function. Do not include the signature in your response. The first line of your response should have 4 spaces of indentation so that it fits syntactically with the user provided signature.

Reflexion Self-Reflection generations follow the form:
[Instruction]
[Function implementation]
[Unit test feedback]

## D  A Running Example

Below is a running example of CR-Planner. Selected action for each step is highlighted in green :

**Problem:** Given a string $s$, find the length of the longest substring without repeating characters in optimal time complexity.

**Instruction:** Reason through the problem and think step by step. Specifically:
1. Restate the problem in plain English.
2. Conceptualize a solution first in plain English.
3. Write a pseudocode solution
4. Output the Python 3 solution to the problem. Make sure to wrap your code in "'python and "' Markdown delimiters, and include exactly one block of code with the entire solution.
No outside libraries are allowed.

**Step 1 (Sub-Goal Selection):** REASON : The next step is to generate a rationale; GENQUERY: The next step is to generate a query; RETRIEVE: The next step is to retrieve a document.

**Step 2 (Execution Selection):** RATIONALE 1: The optimal time complexity is $O(n^2)$; RATIONALE 2 : The optimal time complexity is $O(n)$; RATIONALE 3: The optimal time complexity is $O(n^3)$.

**Step 3 (Sub-Goal Selection):** REASON: The next step is to generate a rationale; GENQUERY : The next step is to generate a query; RETRIEVE: The next step is to retrieve a document.

**Step 4 (Execution Selection):** QUERY 1: Given a string $s$, find the length of the longest substring without repeating characters in optimal time complexity; QUERY 2: Sliding window technique string problems; QUERY 3 : Max length substring with unique characters with $O(n)$ complexity.

**Step 5 (Sub-Goal Selection):** REASON: The next step is to generate a rationale; GENQUERY: The next step is to generate a query; RETRIEVE : The next step is to retrieve a document.

**Step 6 (Execution Selection):** DOC 1: Given a string s, find the length of the longest substring without repeating characters. Examples: Input: "ABCBC" Output: 3 Explanation: The longest substring without repeating characters is "ABC" Input: "AAA" Output: 1 Explanation: The longest substring without repeating characters is "A" Input: "GEEKSFORGEEKS" Output: 7 Explanation: The longest substrings without repeating characters are "EKSFORG" and "KSFORGE" with lengths of 7; DOC 2: To deal with time complexity problems, it always helps to scale the problem up and think of a massive case. If your string was thousands of characters long, we still only have one start pointer and one end pointer. The key thing is that both the pointers will only ever move forward (along the string). Therefore, the complexity of this is definitely $O(n)$ since they are only moving forward together through the string - so the time this process would take is proportional to the length of the string (the time it takes to get to the end); DOC 3 : The intuition behind the solution is to iteratively find the longest substring without repeating characters by maintaining a sliding window approach. We use two pointers (left and right) to represent the boundaries of the current substring. As we iterate through the string, we update the pointers and adjust the window to accommodate new unique characters and eliminate repeating characters.

**Step 7 (Sub-Goal Selection):** REASON : The next step is to generate a rationale; GENQUERY: The next step is to generate a query; RETRIEVE: The next step is to retrieve a document.

**Step 8 (Execution Selection):** RATIONALE 1: The retrieved document is not sufficient for solving the problem. Therefore, a second-level retrieval is required; RATIONALE 2 : To solve the problem efficiently, use a sliding window technique: 1.

Initialize a Window: Start with a window at the beginning of the string, which represents the current substring without repeating characters. 2. Expand the Window: Move through the string one character at a time, adding each character to a data structure (like a set or dictionary) that keeps track of characters in the current window. 3. Check for Repeats: If you encounter a character that is already in the data structure, it means there's a repetition within the current window. 4. Adjust the Window: Move the start of the window forward, removing characters until the repeated character is excluded from the window. This ensures the window contains only unique characters. 5. Update Maximum Length: Keep track of the maximum size of the window throughout the process. This represents the length of the longest substring without repeating characters. 6. Continue Until End of String: Repeat the expand and adjust steps until you have traversed the entire string. This approach ensures that you examine each character at most twice (once when added and once when removed), resulting in optimal linear time complexity; RATIONALE 3: Here is the code:

```python
def length_of_longest_substring(s):
    char_index = {}
    max_length = 0
    start = 0

    for idx, char in enumerate(s):
        if char in char_index:
            start = char_index[char]
        char_index[char] = idx
        current_length = idx - start
        if current_length > max_length:
            max_length = current_length

    return max_length
```

**Step 9 (Sub-Goal Selection):** REASON : The next step is to generate a rationale; GENQUERY: The next step is to generate a query; RETRIEVE: The next step is to retrieve a document.

**Step 10 (Execution Selection):** RATIONALE 1: Here's a step-by-step plan: Initialize Pointers: Start with two pointers, left and right, both at the beginning of the string. These pointers define the current window of unique characters. Create a Character Map: Use a hash map (or an array if dealing with a limited character set like ASCII) to keep track of characters and their latest indices in the string. Iterate Through the String: Move the right pointer through the string one character at a time. a. Check for Repeats: For each character s[right], check if it exists in the character map and if its last recorded index is on or after the left pointer. This

| Current State $s_t$ | Action Space $\mathcal{A}_{s_t}$ | Next State $s_{t+1}$ |
|---|---|---|
| REASON | temperature sampling | RATIONALE |
| GENQUERY | temperature sampling | QUERY |
| RETRIEVE | top-k candidates | DOC |
| RATIONALE | reasoning, querying, retrieving | REASON / GENQUERY / RETRIEVE |
| QUERY | reasoning, querying, retrieving | REASON / GENQUERY / RETRIEVE |
| DOC | reasoning, querying, retrieving | REASON / GENQUERY / RETRIEVE |

Table 9: Details of state types and action spaces.

indicates a repeating character within the current window. b. Adjust the Left Pointer: If a repeat is found, move the left pointer to one position after the last occurrence of the current character. This effectively skips the repeated character and shortens the window from the left side. c. Update the Character Map: Record or update the index of the current character s[right] in the character map with the current right index. d. Calculate the Maximum Length: Compute the length of the current window as current_length = right - left + 1. Update a variable max_length if current_length is greater than the previously recorded maximum. Continue Until End of String: Repeat step 3 until the right pointer reaches the end of the string. Return the Result: After the loop ends, max_length will contain the length of the longest substring without repeating characters; RATIONALE 2 : Here is the code:

```python
def length_of_longest_substring(s):
    char_index = {}
    max_length = 0
    start = 0

    for idx, char in enumerate(s):
        if char in char_index and
            char_index[char] >= start:
            start = char_index[char] + 1
        char_index[char] = idx
        max_length = max(max_length, idx
            - start + 1)

    return max_length
```

RATIONALE 3: Here is the code:

```python
def length_of_longest_substring(s):
    char_set = set()
    left = 0
    max_length = 0

    for right in range(len(s)):
        if s[right] in char_set:
            char_set.clear()
            left = right + 1
        char_set.add(s[right])
        max_length = max(max_length,
            right - left + 1)

    return max_length
```

# E   CR-Planner State Types and Action Spaces

We provide detailed information on state types and action spaces for CR-Planner in Table 9.