

# Multi-LLM Collaborative Search for Complex Problem Solving

Sen Yang<sup>1</sup> Yafu Li<sup>2</sup> Wai Lam<sup>1</sup> Yu Cheng<sup>1,2</sup>

## Abstract

Large language models (LLMs) often struggle with complex reasoning tasks due to their limitations in addressing the vast reasoning space and inherent ambiguities of natural language. We propose the Mixture-of-Search-Agents (MOSA) paradigm, a novel approach leveraging the collective expertise of multiple LLMs to enhance search-based reasoning. MOSA integrates diverse reasoning pathways by combining independent exploration with iterative refinement among LLMs, mitigating the limitations of single-model approaches. Using Monte Carlo Tree Search (MCTS) as a backbone, MOSA enables multiple agents to propose and aggregate reasoning steps, resulting in improved accuracy. Our comprehensive evaluation across four reasoning benchmarks demonstrates MOSA’s consistent performance improvements over single-agent and other multi-agent baselines, particularly in complex mathematical and commonsense reasoning tasks.

## 1. Introduction

Large language models (LLMs) face challenges with complex reasoning, even when augmented with linearized reasoning chains (e.g., Chain-of-Thought), due to the vast reasoning space inherent in the complexity and ambiguity of natural languages. A promising approach is step-wise search-based reasoning, which decomposes the reasoning problem into a traversal over a directed graph, where nodes and edges represent individual reasoning sub-steps distributed across the expansive reasoning space. Related methods have applied various search algorithms to LLMs, such as breadth-first search (BFS), depth-first search (DFS) (Yao et al., 2024; Besta et al., 2024), and best-first search (Hao et al., 2023; Zhang et al., 2024a; Qi et al., 2024).

A successful search trial is featured with diverse yet effective explorations (Hao et al., 2023; Yao et al., 2024). A straight-

<sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>Shanghai AI Laboratory. Correspondence to: Yu Cheng <chengyu@cse.cuhk.edu.hk>.

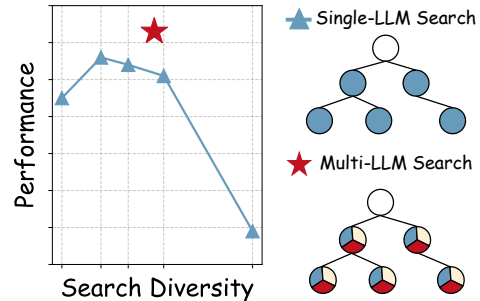


Figure 1: Reasoning performance on MATH-500 against search trajectory diversity. While the diversity of single-LLM search varies with different sampling temperatures, the multi-LLM search consistently achieves superior performance. More details are provided in § 4.1.

forward method to enhance diversity involves increasing the temperature, thereby making the probability distribution more uniform. This is typically combined with top- $k$  and top- $p$  sampling to balance diversity and quality. However, as shown in Figure 1, despite these sampling techniques, achieving a balance between diversity and quality remains challenging and necessitates careful tuning. Besides, even with near-optimal sampling parameters, a single LLM might still get trapped in local optima due to constraints inherent in its training data and architectural design.

To mitigate this limitation, an alternative solution is to aggregate the specialized strengths of multiple LLMs. Recent work (Wang et al., 2024b) has demonstrated that multiple LLMs can collaboratively enhance their instruction-following capabilities by post-editing each other’s responses to the same instruction. Motivated by this progress, we explore leveraging the collective expertise of multiple LLMs for search-based reasoning, which, to the best of our knowledge, has not been previously tested. Figure 1 illustrates the reasoning accuracy on the MATH-500 dataset as a function of search diversity. The performance of search using a single LLM initially improves with increased temperature but subsequently degrades, remaining consistently lower than that of multiple-LLM search.

In this work, we propose Mixture-of-Search-Agents (MOSA), an advanced paradigm for step-wise search-based reasoning that aggregates the complementary strengths of multiple LLMs, leveraging both independent and collabo-

rative contributions to search for reasoning sub-steps more effectively. As illustrated in Figure 3, multiple LLMs propose diverse potential search directions at each reasoning step, either independently or through iterative refinement of each other’s outputs. This hybrid approach ensures that the reasoning process is not constrained by the limitations or biases of any single model. For instance, one model may excel at identifying a promising initial direction, while another might build on it to refine or extend the reasoning path. By combining independence and collaboration, the framework avoids local optima while enhancing reasoning accuracy in the search process.

We performed a comprehensive evaluation of MOSA across four reasoning benchmarks. The findings indicate that MOSA consistently outperforms its single-LLM counterpart in reasoning accuracy with an average improvement of 1.71%. Additionally, our results indicate a synergistic interaction between multi-agent collaboration and search-based reasoning. Further analysis and ablation studies reveal a key challenge for single-agent search-based reasoning: balancing diversity and quality varies across different benchmarks. We also confirm a positive correlation between reasoning performance and the number of distinct search agents, validating the efficacy of multi-agent search. Finally, experiments with an extended action set demonstrate the robustness of MOSA across diverse types of search actions.

## 2. Method

Search-based methods have been extensively used to tackle complex reasoning tasks, such as coding and mathematics, by breaking these problems into multiple search steps (Zhou et al., 2023b; Yao et al., 2024; Hao et al., 2023). Our proposed paradigm is readily applicable to various search algorithms, with the Monte Carlo Tree Search (MCTS) algorithm (Kocsis & Szepesvári, 2006; Coulom, 2007) adopted as the search backbone in this work. This section first introduces the baseline MCTS-based reasoning method with a single search agent (Hao et al., 2023; Qi et al., 2024) in § 2.1, followed by our method, which leverages the expertise of multiple LLMs as search agents in § 2.2.

### 2.1. Baseline Framework

**Overview** Given a problem  $x$  and a generator  $\pi^*$ , MCTS involves iteratively building a search tree starting from the root node  $x$ . We first define the state space  $\mathcal{S}$  and the action space  $\mathcal{A}$ . In our case, each state  $s_j \in \mathcal{S}$  captures the actions (i.e., reasoning steps) generated so far alongside a specific trajectory in the search tree, while each action  $a_j \in \mathcal{A}$  represents the next reasoning step based on the current state and the type of action chosen. As shown in the upper part of Figure 2, given the selected node  $s_i$  (i.e., the reasoning steps generated so far), a step of *Expansion* essentially creates a

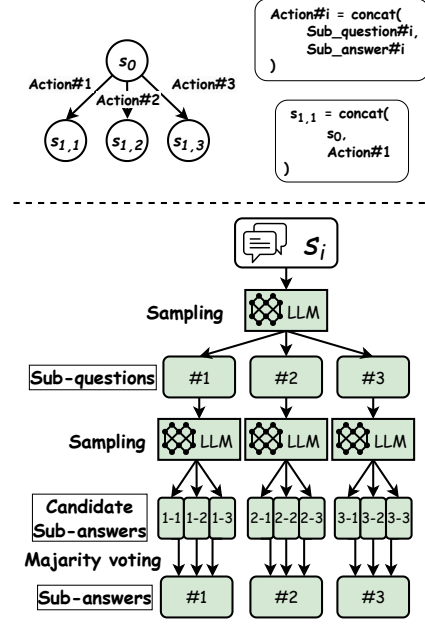


Figure 2: *Top*: An overview of the root node  $s_0$  and its expanded child nodes. *Bottom*: The detailed framework for generating new actions (i.e., sampling sub-questions and sub-answers).

set of child nodes. A child node is created by concatenating  $s_i$  with the new action, with that action being the next reasoning step generated by a search agent (e.g., an LLM) given  $s_i$ .

**Action Space** We follow rStar (Qi et al., 2024) to define a comprehensive set of actions into MCTS-based LLM reasoning. The set of actions,  $\mathcal{A} = \{A1, A2, A3, A4, A5\}$ , includes:

- A1: Propose a one-step thought;
- A2: Propose the remaining thought steps;
- A3: Propose the next sub-question along with its answer;
- A4: Answer the sub-question again;
- A5: Rephrase the question.

Among these actions, we designate A3 as the primary action, comprising a *sub-question* and its corresponding *sub-answer*, i.e.,  $\text{action}_i \equiv \text{concat}(\text{sub\_question}_i, \text{sub\_answer}_i)$ . For instance, an action can be “### Sub-question 3: Does the sum of the previous two digits equal 8? ### Sub-answer 3: The two digits are 3 and 5. We have  $3 + 5 = 8$ , so the answer is yes.”. We consider the other actions along with their effects in an ablation analysis (§ 4.3). We present a detailed illustration of generating new actions, i.e., combinations of sub-question & sub-answer in Algorithm 1. For a given state  $s_i$ , the algorithm traverses all possible actions, where

**Algorithm 1** `GenerateActions`: A function for generating actions, i.e., a sub-question along with a sub-answer, given the current state. The implementation of this function using conventional MCTS and MoSA mainly differs in two aspects: (1) In conventional MCTS methods, the number of search agents  $m = 1$ , while for MoSA,  $m > 1$ . (2) The `FinalizeSubAnswer` function employs heuristic majority voting for single-model search, while employing an additional neural aggregation function for MoSA (see the right section of Figure 3).

**Require:** Selected node  $s_i$ ; Number of sub-questions  $n_q$ ; Number of candidate sub-answers per sub-question  $n_a$ ; A set of LLMs  $\pi^{\text{mix}} = \{\pi_1, \pi_2, \dots, \pi_m\}$

**Ensure:** A set of new actions  $\text{new\_actions} = \{\text{action}_1, \text{action}_2, \dots, \text{action}_{n_q}\}$

Initialize  $\text{new\_actions} \leftarrow \emptyset$  % Prepare the set of new actions

**for**  $i = 1$  **to**  $n_q$  **do**

$\pi^{\text{sub-q}} \leftarrow \text{SelectLLM}(\pi^{\text{mix}})$  % Select an LLM for generating sub-question

$\text{sub\_question}_i \leftarrow \text{GenerateSubQuestion}(\pi^{\text{sub-q}}, s_i)$  % Generate the  $i$ -th sub-question

Initialize  $\text{candidate\_sub\_answers} \leftarrow \emptyset$  % Store candidate sub-answers for sub-question  $i$

**for**  $j = 1$  **to**  $n_a$  **do**

$\pi^{\text{sub-a}} \leftarrow \text{SelectLLM}(\pi^{\text{mix}})$  % Select an LLM for generating a sub-answer

$\text{candidate\_sub\_answer}_j \leftarrow \text{GenerateSubAnswer}(\pi^{\text{sub-a}}, s_i, \text{sub\_question}_i)$  % Generate the  $j$ -th candidate sub-answer

$\text{candidate\_sub\_answers.add}(\text{candidate\_sub\_answer}_j)$  % Store the candidate sub-answer

**end for**

$\text{sub\_answer}_i \leftarrow \text{FinalizeSubAnswer}(\text{candidate\_sub\_answers})$  % Aggregate or vote on candidate sub-answers

$\text{action}_i \leftarrow \text{concat}(\text{sub\_question}_i, \text{sub\_answer}_i)$  % Form the final action by concatenation

$\text{new\_actions.add}(\text{action}_i)$  % Add the action to the set of new actions

**end for**

**return**  $\text{new\_actions}$

the final sub-answer for each sub-question is determined by a heuristic function, e.g., majority voting.

**Reward Function** Following Hao et al. (2023); Qi et al. (2024), we consider a simple yet effective reward function: actions that frequently lead to correct final answers are assigned higher rewards. Specifically,  $Q(s, a)$ , the reward value for node  $s$  created by action  $a$ , receives a positive reward if a trajectory containing node  $s$  reaches a correct final answer, and no reward otherwise. Since the gold answer is not available during testing, the confidence given by *majority voting* is regarded as an approximation of the reward value.

**MCTS Iterations** Typically, each MCTS iteration involves four steps: *Selection*, *Expansion*, *Simulation*, and *Back-propagation*. To balance exploration and exploitation, we adopt the widely-used *Upper Confidence Bounds for Trees* (UCT) algorithm (Kocsis & Szepesvári, 2006) for *Selection*. Formally, a node  $s$  is selected to maximize:

$$\text{UCT}(s, a) = \frac{Q(s, a)}{N(s, a)} + c \sqrt{\frac{\ln N_{\text{parent}}(s)}{N(s, a)}} \quad (1)$$

where  $N_{\text{parent}}(s)$  is the number of times the parent node of  $s$  has been visited,  $N(s, a)$  is the number of times node  $s$  has been visited, and  $c$  is a constant. Once the node  $s$  is selected, an *Expansion* step is performed to add child nodes to  $s$ . After that, starting from a random child node, a *Simulation* is performed using the default rollout policy until a terminal node is obtained or a predefined maximum depth is reached. The outcome of the simulation determines the reward, which is then propagated back up the tree during

the *Back-propagation* step. Upon multiple iterations, we consider each leaf node as a solution. In this work, we focus on *Expansion*, which aims to effectively expand the search space.

**Sampling Diversity** Applying stochastic sampling techniques in LLM generation is essential for introducing diversity to MCTS. As presented in the lower part of Figure 2, given the selected state  $s_0$ , the sub-questions and the sub-answer candidates are all stochastically sampled using temperature scaling, top- $k$  sampling and nucleus sampling (Holtzman et al., 2020). In § 4.1, we empirically alter search diversity by manipulating generation temperature for single-LLM search.

## 2.2. Mixture-of-Search-Agents

Conventional Monte Carlo Tree Search (MCTS) methods utilizing a single model face two significant limitations: (1) Encouraging search diversity while maintaining generation quality is challenging (Zhang et al., 2020), necessitating meticulous tuning of sampling parameters to balance the trade-off between these aspects; (2) using heuristic metrics like majority voting to determine the final sub-answer can be less accurate when the model favors incorrect search directions. To this end, we explore a simple yet effective alternative, Mixture-of-Search-Agents (MoSA), which employs multiple agents to perform search algorithms like MCTS and utilizes a neural function to refine the candidate step-wise outputs. Firstly, leveraging the distinct distributions from different models intrinsically yields better generation

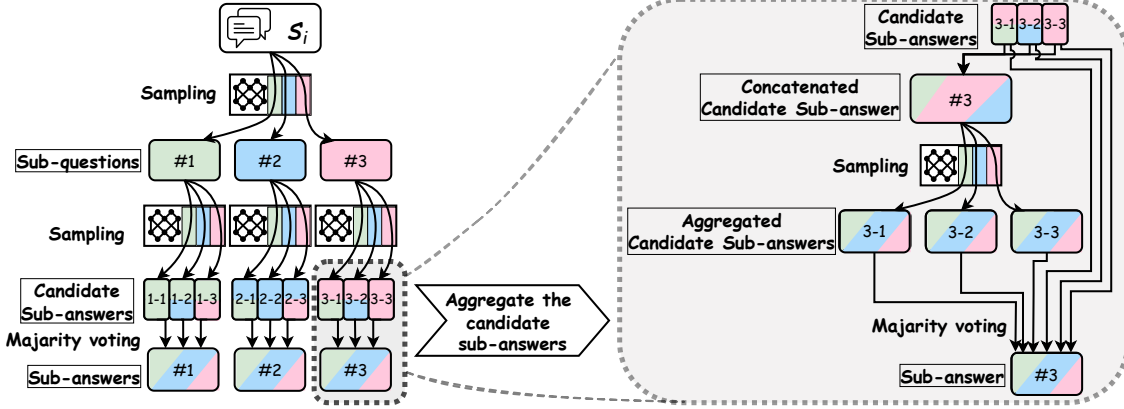


Figure 3: Generate three new actions using MOSA. *Left:* Use MOSA to propose sub-questions and sub-answers. *Right:* Use MOSA to aggregate candidate sub-answers.

diversity, alleviating the necessity for sampling parameters optimization. Additionally, incorporating a neural function enhances the robustness of answer aggregation.

Figure 3 illustrates how our method generates three new actions starting from the current node  $s_i$ . Unlike vanilla single-model search, MOSA employs multiple agents (denoted by distinct colors) to explore diverse actions, such as sub-questions and sub-answers. In the remainder of this section, we illustrate two roles performed by MOSA when generating new actions in MCTS. Specifically, we will start with the straightforward improvement, *MOSA as Proposers*, where multiple agents are involved for sampling actions; then we will introduce the more intricate *MOSA as Aggregators*, which extends the heuristic majority voting method to an aggregating phase where multiple LLMs read and refine the answers given by all.

**MOSA as Proposers to Diversify Actions** The left side of Figure 3 shows *MOSA as Proposers*, where the single search agent adopted by the baseline MCTS method (as in Figure 2) is replaced by MOSA. MOSA leverages multiple LLMs to enhance action diversity by fulfilling two sub-roles: multi-agent proposers that generate *sub-questions* and *sub-answers*.

Generating a new search action begins with sampling a sub-question from the current state  $s_i$ . The sub-question proposing phase essentially controls the directions of the current search step because whatever follows within this step is constrained by the scope of that sub-question. Because of this, we consider maintaining the independence among sub-questions, ensuring that the initial search direction indicated by each sub-question is independent of others. As shown in the upper-left part of Figure 3, this effectively diversifies the sampled sub-questions as the same  $s_i$  is colored with distinct characteristics after going through different LLMs.

After the initial search directions are created, the target is to comprehensively explore each search direction. To achieve this, each sub-question is answered by various LLMs, generating a diverse set of candidate sub-answers. These candidates are then aggregated to reach a finalized sub-answer. A simple yet effective aggregating method is majority voting, leveraging the principle of self-consistency (Wang et al., 2023).

### MOSA as Aggregators for Collaborative Refinement

We introduce a neural function, termed “aggregator”, to mitigate the limitation of majority voting for selecting the final answer. An aggregator leverages the innate capability of the LLM to critique, compare and aggregate multiple answers into a final answer. Specifically, we prompt each LLM to consolidate all responses into an aggregated answer (see Appendix B for detailed prompts), resulting in a new set of aggregated answers as illustrated in the right section of Figure 3. The underlying intuition is that this aggregation process enhances the likelihood of producing correct answers by facilitating comparisons among different responses, thereby increasing the overall success rates for correct answers under majority voting. We present an example below to illustrate this intuition.

In the previous section, we consider majority-voting after obtaining candidate sub-answers from diverse proposers. Let us consider a sub-question that requires 3 sub-answers to be generated and an MOSA component consisting of 3 distinct LLMs. We simply assume that each LLM proposes one sub-answer. If there are  $k$  LLMs that are proficient at this sub-question and the other  $3 - k$  are not, then it is likely that we would have  $k$  *good* sub-answer candidates and  $3 - k$  *bad* candidates<sup>1</sup>. With majority voting, a *bad*

<sup>1</sup>For clarity in illustrating our motivation, we simplify the correctness of candidate sub-answers into two groups: *good* and *bad*. This abstraction helps explain the role of aggregators in improving



finalized answer is likely if  $k \geq 2$ .

Now we turn to use MOSA to aggregate the candidate sub-answers and then include the aggregated sub-answers into majority-voting. The inputs for all three aggregator LLMs are the same, which concatenates the sub-question and all the three candidate sub-answers. We hypothesize that a *bad* aggregator that receives at least a *good* sub-answer could yield a sub-answer that is at least better than its original *bad* sub-answer. Such a hypothesis has been empirically verified in the case of instruction following by Wang et al. (2024b), who showed that many LLMs can generate higher-quality responses by building upon outputs from other LLMs. Thus, if the two *bad* aggregator LLMs can learn from the *good* sub-answer and generate *good* aggregated sub-answers, then we will have 4 *good* sub-answers and 2 *bad* ones, which lead to a *good* finalized sub-answer.

### 3. Experiments

#### 3.1. Baselines

**Few-shot Chain-of-Thought (CoT)** (Wei et al., 2023) feeds the LLM with a few demonstrations followed by the input question. Since we are using instruction-tuned LLMs, we format the demonstrations as multi-turn dialogues. In each turn, the human asks a question and then the assistant answers it.

**Self-Consistency@ $n$**  (Wang et al., 2023) also adopts the few-shot CoT prompting scheme, but it samples  $n$  independent answers per instance. The final answer is then given by majority voting over the  $n$  candidate answers. Except for the conventional single-LLM self-consistency experiments, we also evaluate self-consistency with multiple different LLMs. Such a multi-LLM self-consistency setting can be regarded as a simplified version of Wang et al. (2024b), which collects direct answers from various agents and aggregates them with majority voting.

**Reasoning-via-Planning (RAP)** (Hao et al., 2023) is a representative LLM-based reasoning method using MCTS. We use it as the foundation to apply MOSA. In each search step, RAP generates one or more sub-questions along with their sub-answers. The original RAP paper adopted different reward functions for different types of tasks. In this work, we use the simple self-consistency score as the reward value, which has been shown to be competitive with those manually designed ones in Appendix A.1 of Qi et al. (2024). Note that we ensure the total number of LLM forward calls of a single-LLM method are approximately the same as its multi-LLM counterpart, e.g., RAP *versus* RAP + MOSA as Proposers in Table 1.

answer correctness, though actual correctness exists on a spectrum depending on task complexity and evaluation criteria.

**rStar** (Qi et al., 2024) is one of the recent SoTA MCTS-based LLM reasoning methods. The authors proposed a comprehensive set of search actions, which we have introduced in § 2.1. We adopt their innovative set of actions to evaluate the effects brought by the scope of action set on MOSA in § 4.3.

#### 3.2. Experimental Settings

**Benchmarks** We perform evaluation on four reasoning benchmarks covering different scopes, including three mathematical reasoning datasets (GSM8K (Cobbe et al., 2021), SVAMP (Patel et al., 2021), MATH-500 (Hendrycks et al., 2021; Lightman et al., 2023)) and one commonsense reasoning dataset (StrategyQA (Geva et al., 2021)).

**Models** We adopt four open-sourced instruction-following LLMs to formulate the LLM pool of MOSA: Llama-3.1-8B-Instruct (Grattafiori et al., 2024), Qwen-2-7B-Instruct (Yang et al., 2024), Ministral-8B-Instruct-2410 (Mistral, 2024), and GLM-4-9B-Chat (GLM et al., 2024). The number of LLMs could also be made larger or smaller, depending on customized choices. Our later experiments will show that benchmark performances are positively correlated with the number of distinct LLMs.

**Implementation Details** For few-shot CoT baselines, we report the results of all four LLMs. For other single-LLM baselines, like Self-Consistency@ $n$  and RAP, we adopt **Llama-3.1-8B-Instruct** due to its competitiveness and robustness across various benchmarks. For all experiments regarding sampling from multiple LLMs, we try to maintain a pseudo uniform distribution for the `SelectLLM` function in Algorithm 1. That is, if 7 completions need to be sampled and there are 4 distinct LLMs, we manually assign each LLM to sample one completion and then uniformly sample 3 LLMs out of 4 without replacement to finish the remaining 3 completions. Hyper-parameter settings are listed in Appendix A.

#### 3.3. Main Results

We report the main results on the four benchmarks in Table 1. Below we highlight our key findings.

**MOSA Leads in Reasoning Tasks** RAP + MOSA as Proposers & Aggregators consistently yields superior performances across all datasets (GSM8K, SVAMP, MATH-500, StrategyQA), reaching an average performance (Avg.) of 79.97%. Specifically, it obtains exceptional improvements (+1.8%) over the best baseline on the challenging MATH-500 benchmark, suggesting it is effective at handling complex reasoning problems.

**Synergistic Effect between Multi-Agent Collaboration and Search-based Reasoning** MOSA integrates two research paradigms: multi-agent collaboration and search-

Table 1: Main results. Those rows marked by ‡ were reported by the rStar paper (Qi et al., 2024) using Llama-3-8B-Instruct. All other results are reported by our experiments. MULTI refers to multi-LLM while STG represents StrategyQA. The highest number on each dataset is marked in bold while the secondary high is underlined.

METHOD	MULTI?	SEARCH?	GSM8K	SVAMP	MATH	STG	AVG.
FEW-SHOT CoT:							
◦ LLAMA-3.1-8B-INSTRUCT	×	×	84.00	86.80	41.60	67.39	69.95
◦ MINISTRAL-8B-INSTRUCT-2410	×	×	82.41	89.20	40.00	70.60	70.55
◦ QWEN-2-7B-INSTRUCT	×	×	84.00	88.60	24.20	66.67	65.87
◦ GLM-4-9B-CHAT	×	×	83.85	89.70	40.00	71.32	71.22
w/ LLAMA-3.1-8B-INSTRUCT:							
◦ SELF-CONSISTENCY@4	×	×	88.02	89.70	43.80	69.43	72.74
◦ SELF-CONSISTENCY@32	×	×	90.37	92.40	44.80	70.89	74.62
◦ SELF-CONSISTENCY@128	×	×	90.98	93.30	52.20	71.32	76.95
◦ SELF-CONSISTENCY@256	×	×	90.90	92.90	53.20	71.03	77.01
w/ ALL FOUR LLMs:							
◦ SELF-CONSISTENCY@4	✓	×	90.45	92.20	46.20	70.45	74.82
◦ SELF-CONSISTENCY@32	✓	×	90.75	93.20	52.60	71.76	77.08
◦ SELF-CONSISTENCY@128	✓	×	91.21	93.70	53.80	72.78	77.87
◦ SELF-CONSISTENCY@256	✓	×	90.98	93.50	54.20	71.47	77.54
RAP ‡	×	✓	80.59	85.70	18.80	68.71	63.45
RAP	×	✓	90.52	91.60	53.00	75.40	77.63
+ SINGLE-LLM AS AGGREGATOR	×	✓	90.05	92.50	54.80	75.69	78.26
+ MOSA AS PROPOSERS	✓	✓	91.13	94.50	54.60	75.69	78.98
+ MOSA AS PROPOSERS & AGGREGATORS	✓	✓	<b>91.96</b>	<b>94.90</b>	<b>56.60</b>	<b>76.42</b>	<b>79.97</b>

based reasoning. When applied independently, each achieves moderate improvements, but their combination yields significantly enhanced results due to synergy effects. (1) Transitioning from single-agent to multi-agent: Across all four benchmarks, transitioning from a single LLM to multiple LLMs with the best non-search baseline (Self-consistency) results in an average absolute improvement of +0.53%. By contrast, transitioning from single-agent search (RAP) to multi-agent search (MOSA as Proposers) yields a larger average absolute improvement of +1.35%. Augmenting with aggregators further increases the improvement from single-agent search (RAP + Single-LLM as Aggregator) to multi-agent search (MOSA as Proposers and Aggregators), achieving +1.71%. (2) Transitioning from non-search to search-based reasoning: Using a single LLM, the performance gap between non-search (Self-consistency@256) and search (RAP) is +0.62%. This gap widens to +1.44% when employing multiple LLMs, showcasing the synergy between multi-agent collaboration and search. These results highlight that combining multi-agent collaboration with search-based reasoning yields significantly greater performance gains than applying either approach in isolation.

**Boosting Search-based Reasoning with MOSA as Aggregators** While vanilla RAP performs well, the inclusion of aggregators, particularly with MOSA as Aggregators, significantly enhances performance. For instance, augmenting RAP with a single-LLM aggregator yields an average improvement of +0.63%. This improvement increases to +0.99% when MOSA as Proposers is further enhanced with

MOSA as Aggregators.

**Search-Based Methods Excel in Complex Reasoning Tasks** The best accuracy numbers on GSM8K and SVAMP, both exceeding 90%, suggest these datasets are relatively easier. In contrast, MATH-500 and StrategyQA, with best scores around 55% and 80%, respectively, are more challenging. Notably, search-based methods demonstrate a clear advantage on these more complex datasets, underscoring their effectiveness in tackling intricate reasoning tasks. Take StrategyQA as an example, the best accuracy number with non-search methods (Self-consistency) is 72.78%, which is significantly lower than the best search counterpart (RAP) accuracy (75.69%).

## 4. Analysis

We perform a comprehensive analysis on MOSA in this section. Specifically, we scale the diversity of the single-LLM search baseline in § 4.1 and compare it with MOSA. In § 4.2, we vary the number of distinct LLMs in MOSA. In § 4.3, we combine MOSA with the rich set of actions proposed by Qi et al. (2024). Finally, we evaluate variations of MOSA by ablating the numbers of proposers and aggregators in § 4.4.

### 4.1. Diversity versus Performance

For single-LLM search, a common technique to increase generation diversity is to manipulate with decoding hyper-

Table 2: The results of MoSA combined with rStar (Qi et al., 2024), a recent SoTA MCTS-based reasoning method that extends the set of actions. Those results marked with ‡ were reported by the rStar paper using Llama-3-8B-Instruct. All other results are reported by our experiments. The definitions of  $A\{1, 2, 3, 4, 5\}$  are in § 2.1. **The highest number on each dataset is marked in bold while the secondary high is underlined.**

METHOD	ACTION SET	GSM8K	SVAMP	MATH	STG	Avg.
RAP ‡	$A\{2, 3\}$	80.59	85.70	18.80	68.71	63.45
RAP	$A\{2, 3\}$	90.52	91.60	53.00	75.40	77.63
+ MoSA AS PROPOSERS	$A\{2, 3\}$	91.13	94.50	54.60	<u>75.69</u>	78.98
+ MoSA AS PROPOSERS & AGGREGATORS	$A\{2, 3\}$	<u>91.96</u>	94.90	56.60	<b>76.42</b>	79.97
RSTAR ‡	$A\{1, 2, 3, 4, 5\}$	88.70	91.89	38.30	71.47	72.59
RSTAR	$A\{1, 2, 3, 4, 5\}$	91.36	93.30	59.00	74.96	79.66
+ MoSA AS PROPOSERS	$A\{1, 2, 3, 4, 5\}$	<u>91.96</u>	<b>95.60</b>	<u>63.20</u>	75.11	<u>81.47</u>
+ MoSA AS PROPOSERS & AGGREGATORS	$A\{1, 2, 3, 4, 5\}$	<b>92.04</b>	<u>95.10</u>	<b>63.60</b>	75.40	<b>81.54</b>

parameters, e.g., the sampling temperature. We modify the temperature of the RAP + Single-LLM as Aggregator baseline on two datasets, with  $T = \{0.25, 0.5, 0.75, 1.0, 1.25\}$ . Diversity is assessed by calculating the  $\{1, 2, 3, 4\}$ -gram Vendi Score (Friedman & Dieng, 2023) across search trajectories. Figure 4 illustrates that while the reasoning accuracy of RAP initially improves with increasing diversity, it subsequently declines. More importantly, the two benchmarks favor different temperature values. This suggests that attaining an optimal balance between diversity and reasoning performance requires careful tuning, as balancing diversity and quality can be challenging (Zhang et al., 2020). In contrast, RAP + MoSA with the default sampling parameters consistently holds an advantageous position.

#### 4.2. Ablation of LLM Collaboration

To evaluate the impact of varying the number of different LLMs in MoSA, we conduct an analysis using 1 to 4 LLMs across four benchmarks, prioritizing them in the following order: Llama, GLM, Qwen, Ministral. All four variants utilize approximately the same number of LLM forward calls, ensuring that the only variable is the number of distinct LLMs involved. Figure 5 shows that increasing the number of different LLMs generally correlates with higher reasoning accuracy, except for a slight decrease in performance when the number of LLMs increases from 3 to 4 on MATH-500. This trend indicates that the diverse expertise contributed by different LLMs significantly enhances search-based reasoning performance.

#### 4.3. Support for Extended Action Set

rStar (Qi et al., 2024) proposes using a comprehensive set of actions in MCTS-based LLM reasoning. Since enriching the action set is orthogonal to our method, we hypothesize that MoSA is compatible with the enlarged action set. The results in Table 2 support our hypothesis. For example, rStar combined with MoSA boosts the reasoning accuracy

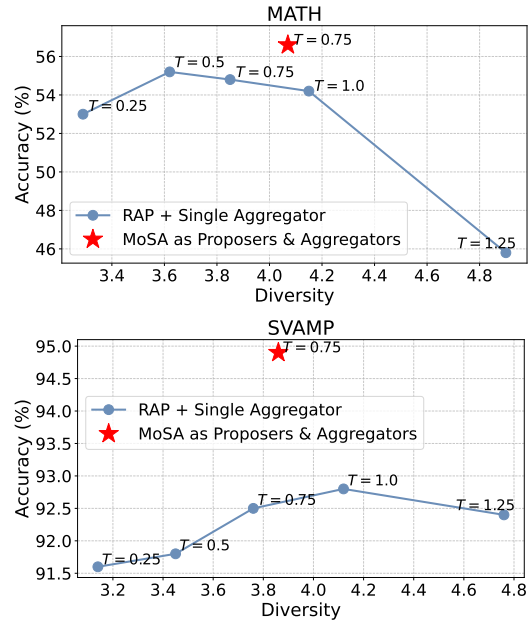


Figure 4: Diversity *versus* accuracy.  $T$  = Temperature.

on MATH-500 from 59.00% to 63.20% (+ MoSA as Proposers) and 63.60% (+ MoSA as Proposers & Aggregators). We also found that enriching the action set is not always beneficial. On StrategyQA, the expanded action set yielded inferior performance; however, we note that MoSA still demonstrates improvements.

#### 4.4. Ablation of Proposers & Aggregators

We consider to isolate the effects of MoSA as Proposers and MoSA as Aggregators by ablating the number of distinct LLMs for those two roles. As shown in Table 3, changing the number of distinct proposers to be single yields a larger decrease comparing with ablating the number of aggregators (-1.23% *versus* -0.47%), suggesting that MoSA brings more benefits as proposers.

Table 3: Ablation analysis to isolate the effects of MOSA as proposers and as aggregators, respectively, for search-based reasoning. By multi, we are referring to the default setting in our experiments, i.e., 4 distinct LLMs. **The highest number** on each dataset is marked in bold while the secondary high is underlined.

#PROPOSER(S)	#AGGREGATOR(S)	GSM8K	SVAMP	MATH	STRATEGYQA	AVG.
SINGLE	NONE	90.52	91.60	53.00	75.40	77.63
SINGLE	SINGLE	90.05	92.50	54.80	75.69	78.26
SINGLE	MULTI	91.05	91.90	55.60	<u>76.42</u>	78.74
MULTI	NONE	91.13	<u>94.50</u>	54.60	75.69	78.98
MULTI	SINGLE	<u>91.66</u>	94.20	<u>56.00</u>	76.13	<u>79.50</u>
MULTI	MULTI	<b>91.96</b>	<b>94.90</b>	<b>56.60</b>	<b>76.42</b>	<b>79.97</b>

## 5. Related Work

### 5.1. Reasoning with LLMs

The recent focus on large language models is partly due to their exceptional performance in solving complex reasoning tasks. A prominent example is Chain-of-Thought (CoT) reasoning (Wei et al., 2023). Recent advancements include self-consistency (Wang et al., 2023), problem decomposition (Zhou et al., 2023b), the use of tools (Gao et al., 2023; Chen et al., 2023), and search-based methods (Hao et al., 2023; Yao et al., 2024; Qi et al., 2024). Among these approaches, MOSA is most closely aligned with search-based reasoning methods.

**Search-based Reasoning** Search-based reasoning has demonstrated effectiveness, particularly for solving complex, multi-step problems (Hao et al., 2023; Yao et al., 2024; Chen et al., 2024c; Zhang et al., 2024a; Chen et al., 2024a; Qi et al., 2024; Zhang et al., 2024b; Zhou et al., 2023a; Koh et al., 2024). One of the recent state-of-the-art systems in this domain is rStar (Qi et al., 2024). rStar introduces two key innovations: (1) expanding the Monte Carlo Tree Search (MCTS) action space from one or two actions to five; and (2) employing a secondary LLM to verify the reasoning trajectories generated by the primary LLM through MCTS. In § 4.3, we empirically demonstrate that our method is complementary to the enriched action set of rStar.

### 5.2. LLM Ensemble

Ensembling, a widely used technique for leveraging the strengths of multiple models, remains highly effective in the era of LLMs. Jiang et al. (2023) proposed pairwise reranking of LLM outputs and fusing multiple responses using a trained generative model. Several studies have proposed training routing functions to match queries with appropriate LLMs (Lu et al., 2023; Shnitzer et al., 2023; Wang et al., 2024a). Others have proposed averaging the output distributions of multiple LLMs (Huang et al., 2024).

Another line of research focuses on multi-agent collaboration, where multiple LLMs interact to discuss or debate

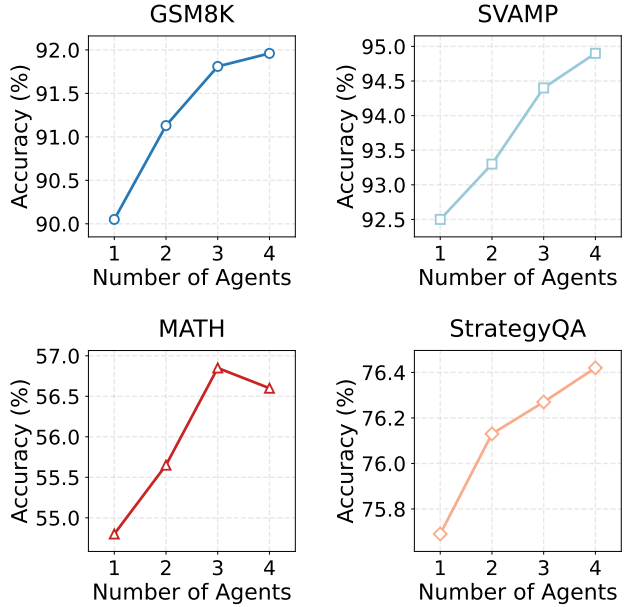


Figure 5: Reasoning accuracy with different number of distinct LLMs as search agents.

specific topics (Du et al., 2023; Liang et al., 2024; Chan et al., 2023; Xu et al., 2023; Liu et al., 2024; He et al., 2023; Chen et al., 2024b; Zhang et al., 2024c). Common design variations in this paradigm include role assignments for LLMs (e.g., debaters and judges) and discussion mechanisms (e.g., symmetric versus asymmetric interactions).

## 6. Conclusion

In this work, we investigated a novel paradigm called MOSA. MOSA combines independent exploration and iterative refinement among multiple LLMs to enhance reasoning diversity and accuracy. Experiments across benchmarks demonstrate its consistent advantages over single-LLM and multi-agent baselines, especially in complex tasks. This work underscores the potential of multi-agent collaboration in advancing search-based reasoning.



## Impact Statement

This work aims to contribute to the advancement of reasoning with LLMs. While our research could have various societal implications, none are deemed significant enough to warrant specific mention at this stage.

## References

- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
- Chan, C.-M., Chen, W., Su, Y., Yu, J., Xue, W., Zhang, S., Fu, J., and Liu, Z. Chateval: Towards better llm-based evaluators through multi-agent debate, 2023. URL <https://arxiv.org/abs/2308.07201>.
- Chen, G., Liao, M., Li, C., and Fan, K. Alphamath almost zero: Process supervision without process, 2024a. URL <https://arxiv.org/abs/2405.03553>.
- Chen, J., Saha, S., and Bansal, M. ReConcile: Round-table conference improves reasoning via consensus among diverse LLMs. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7066–7085, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.381. URL <https://aclanthology.org/2024.acl-long.381/>.
- Chen, Q., Wang, X., Mondorf, P., Hedderich, M. A., and Plank, B. Understanding when tree of thoughts succeeds: Larger models excel in generation, not discrimination, 2024c. URL <https://arxiv.org/abs/2410.17820>.
- Chen, W., Ma, X., Wang, X., and Cohen, W. W. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=YfZ4ZPt8zd>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- Coulom, R. Efficient selectivity and backup operators in monte-carlo tree search. In van den Herik, H. J., Ciancarini, P., and Donkers, H. H. L. M. J. (eds.), *Computers and Games*, pp. 72–83, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg. ISBN 978-3-540-75538-8.
- Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. Improving factuality and reasoning in language models through multiagent debate, 2023. URL <https://arxiv.org/abs/2305.14325>.
- Friedman, D. and Dieng, A. B. The vendi score: A diversity evaluation metric for machine learning. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- Gao, L., Madaan, A., Zhou, S., Alon, U., Liu, P., Yang, Y., Callan, J., and Neubig, G. Pal: program-aided language models. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- Geva, M., Khashabi, D., Segal, E., Khot, T., Roth, D., and Berant, J. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021. doi: 10.1162/tacl-a.00370. URL <https://aclanthology.org/2021.tacl-1.21/>.
- GLM, T., Zeng, A., Xu, B., Wang, B., Zhang, C., Yin, D., Rojas, D., Feng, G., Zhao, H., Lai, H., Yu, H., Wang, H., Sun, J., Zhang, J., Cheng, J., Gui, J., Tang, J., Zhang, J., Li, J., Zhao, L., Wu, L., Zhong, L., Liu, M., Huang, M., Zhang, P., Zheng, Q., Lu, R., Duan, S., Zhang, S., Cao, S., Yang, S., Tam, W. L., Zhao, W., Liu, X., Xia, X., Zhang, X., Gu, X., Lv, X., Liu, X., Liu, X., Yang, X., Song, X., Zhang, X., An, Y., Xu, Y., Niu, Y., Yang, Y., Li, Y., Bai, Y., Dong, Y., Qi, Z., Wang, Z., Yang, Z., Du, Z., Hou, Z., and Wang, Z. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A., Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang, B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra, C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong, C., Ferrer, C. C., Nikolaidis, C., Allonsius, D., Song, D., Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary, D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes, D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan, E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F., Synnaeve, G., Lee, G., Anderson, G. L., Thattai, G., Nail, G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Korevaar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I. A., Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J., Lee, J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J.,

- Shah, J., van der Linde, J., Billock, J., Hong, J., Lee, J., Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton, J., Spisak, J., Park, J., Rocca, J., Johnstun, J., Saxe, J., Jia, J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li, K., Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik, K., Chiu, K., Bhalla, K., Lakhotia, K., Rantala-Yearly, L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L., Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat, L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh, M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham, M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M., Si, M., Singh, M. K., Hassan, M., Goyal, N., Torabi, N., Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N., Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P., Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan, P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan, R., Ganapathy, R., Calderer, R., Cabral, R. S., Stojnic, R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R., Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva, R., Hou, R., Wang, R., Hosseini, S., Chennabasappa, S., Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang, S., Raparthy, S., Shen, S., Wan, S., Bhosale, S., Zhang, S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S., Collot, S., Gururangan, S., Borodinsky, S., Herman, T., Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speckbacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V., Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do, V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong, W., Fu, W., Meers, W., Martinet, X., Wang, X., Wang, X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Goldschlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang, Y., Li, Y., Mao, Y., Coudert, Z. D., Yan, Z., Chen, Z., Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey, A., Shajnfeld, A., Gangidi, A., Victoria, A., Goldstand, A., Menon, A., Sharma, A., Boesenberg, A., Baevski, A., Feinstein, A., Kallet, A., Sangani, A., Teo, A., Yunus, A., Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poulton, A., Ryan, A., Ramchandani, A., Dong, A., Franco, A., Goyal, A., Saraf, A., Chowdhury, A., Gabriel, A., Bharambe, A., Eisenman, A., Yazdan, A., James, B., Maurer, B., Leonhardi, B., Huang, B., Loyd, B., Paola, B. D., Paranjape, B., Liu, B., Wu, B., Ni, B., Hancock, B., Wasti, B., Spence, B., Stojkovic, B., Gamido, B., Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C., Wang, C., Kim, C., Zhou, C., Hu, C., Chu, C.-H., Cai, C., Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty, D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine, D., David, D., Parikh, D., Liskovich, D., Foss, D., Wang, D., Le, D., Holland, D., Dowling, E., Jamil, E., Montgomery, E., Presani, E., Hahn, E., Wood, E., Le, E.-T., Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun, F., Kreuk, F., Tian, F., Kokkinos, F., Ozgenel, F., Cagioni, F., Kanayet, F., Seide, F., Florez, G. M., Schwarz, G., Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov, G., Guangyi, Zhang, Lakshminarayanan, G., Inan, H., Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H., Rudolph, H., Suk, H., Aspegren, H., Goldman, H., Zhan, H., Damla, I., Molybog, I., Tufanov, I., Leontiadis, I., Veliche, I.-E., Gat, I., Weissman, J., Geboski, J., Kohli, J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J., Chan, J., Zhen, J., Reizenstein, J., Teboul, J., Zhong, J., Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard, J., McPhie, J., Torres, J., Ginsburg, J., Wang, J., Wu, K., U, K. H., Saxena, K., Khandelwal, K., Zand, K., Matosich, K., Veeraraghavan, K., Michelena, K., Li, K., Jagadeesh, K., Huang, K., Chawla, K., Huang, K., Chen, L., Garg, L., A, L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L., Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M., Bhatt, M., Mankus, M., Hasson, M., Lennie, M., Reso, M., Groshev, M., Naumov, M., Lathi, M., Keneally, M., Liu, M., Seltzer, M. L., Valko, M., Restrepo, M., Patel, M., Vyatskov, M., Samvelyan, M., Clark, M., Macey, M., Wang, M., Hermoso, M. J., Metanat, M., Rastegari, M., Bansal, M., Santhanam, N., Parks, N., White, N., Bawa, N., Singhal, N., Egebo, N., Usunier, N., Mehta, N., Laptev, N. P., Dong, N., Cheng, N., Chernoguz, O., Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P., Saab, P., Balaji, P., Rittner, P., Bontrager, P., Roux, P., Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P., Liang, Q., Alao, R., Rodriguez, R., Ayub, R., Murthy, R., Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan, R., Battey, R., Wang, R., Howes, R., Rinott, R., Mehta, S., Siby, S., Bondu, S. J., Datta, S., Chugh, S., Hunt, S., Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma, S., Yamamoto, S., Ramaswamy, S., Lindsay, S., Lindsay, S., Feng, S., Lin, S., Zha, S. C., Patil, S., Shankar, S., Zhang, S., Zhang, S., Wang, S., Agarwal, S., Sajuyigbe, S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satterfield, S., Govindaprasad, S., Gupta, S., Deng, S., Cho, S., Virk, S., Subramanian, S., Choudhury, S., Goldman, S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson, T., Li, T., Zhang, T., Matthews, T., Chou, T., Shaked, T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V., Kumar, V. S., Mangla, V., Ionescu, V., Poenaru, V., Mihailescu, V. T., Ivanov, V., Li, W., Wang, W., Jiang, W., Bouaziz, W., Constable, W., Tang, X., Wu, X., Wang, X., Wu, X., Gao, X., Kleinman, Y., Chen, Y., Hu, Y., Jia, Y., Qi, Y., Li, Y., Zhang, Y., Zhang, Y., Adi, Y., Nam, Y., Yu, Wang, Zhao, Y., Hao, Y., Qian, Y., Li, Y., He, Y., Rait, Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao, Z., and Ma, Z. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Hao, S., Gu, Y., Ma, H., Hong, J., Wang, Z., Wang, D., and Hu, Z. Reasoning with language model is planning with world model. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 8154–8173,

- Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.507. URL <https://aclanthology.org/2023.emnlp-main.507>.
- He, Z., Cao, P., Chen, Y., Liu, K., Li, R., Sun, M., and Zhao, J. LEGO: A multi-agent collaborative framework with role-playing and iterative feedback for causality explanation generation. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 9142–9163, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.613. URL <https://aclanthology.org/2023.findings-emnlp.613/>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- Holtzman, A., Buys, J., Du, L., Forbes, M., and Choi, Y. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rygGQyrFvH>.
- Huang, Y., Feng, X., Li, B., Xiang, Y., Wang, H., Qin, B., and Liu, T. Ensemble learning for heterogeneous large language models with deep parallel collaboration, 2024. URL <https://arxiv.org/abs/2404.12715>.
- Jiang, D., Ren, X., and Lin, B. Y. LLM-blender: Ensembling large language models with pairwise ranking and generative fusion. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14165–14178, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.792. URL <https://aclanthology.org/2023.acl-long.792/>.
- Kocsis, L. and Szepesvári, C. Bandit based monte-carlo planning. In Fürnkranz, J., Scheffer, T., and Spiliopoulou, M. (eds.), *Machine Learning: ECML 2006*, pp. 282–293, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg. ISBN 978-3-540-46056-5.
- Koh, J. Y., McAleer, S., Fried, D., and Salakhutdinov, R. Tree search for language model agents, 2024. URL <https://arxiv.org/abs/2407.01476>.
- Liang, T., He, Z., Jiao, W., Wang, X., Wang, Y., Wang, R., Yang, Y., Shi, S., and Tu, Z. Encouraging divergent thinking in large language models through multi-agent debate. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 17889–17904, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.992. URL <https://aclanthology.org/2024.emnlp-main.992/>.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Liu, Z., Zhang, Y., Li, P., Liu, Y., and Yang, D. A dynamic LLM-powered agent network for task-oriented agent collaboration. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=XII0Wp1XA9>.
- Lu, K., Yuan, H., Lin, R., Lin, J., Yuan, Z., Zhou, C., and Zhou, J. Routing to the expert: Efficient reward-guided ensemble of large language models, 2023. URL <https://arxiv.org/abs/2311.08692>.
- Mistral, A. t. Un ministral, des ministraux — mistral ai — frontier ai in your hands. <https://mistral.ai/news/ministaux/>, 2024.
- Patel, A., Bhattamishra, S., and Goyal, N. Are NLP models really able to solve simple math word problems? In Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., and Zhou, Y. (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.168. URL <https://aclanthology.org/2021.naacl-main.168/>.
- Qi, Z., Ma, M., Xu, J., Zhang, L. L., Yang, F., and Yang, M. Mutual reasoning makes smaller llms stronger problem-solvers. *arXiv preprint arXiv:2408.06195*, 2024.
- Shnitzer, T., Ou, A., Silva, M., Soule, K., Sun, Y., Solomon, J., Thompson, N., and Yurochkin, M. Large language model routing with benchmark datasets, 2023. URL <https://arxiv.org/abs/2309.15789>.
- Wang, H., Polo, F. M., Sun, Y., Kundu, S., Xing, E., and Yurochkin, M. Fusing models with complementary expertise. In *The Twelfth International Conference on Learning Representations*, 2024a. URL <https://openreview.net/forum?id=PhMrGCMIRL>.

- Wang, J., Wang, J., Athiwaratkun, B., Zhang, C., and Zou, J. Mixture-of-agents enhances large language model capabilities, 2024b. URL <https://arxiv.org/abs/2406.04692>.
- Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=1PL1NIMMrw>.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
- Xu, Z., Shi, S., Hu, B., Yu, J., Li, D., Zhang, M., and Wu, Y. Towards reasoning in large language models via multi-agent peer review collaboration, 2023. URL <https://arxiv.org/abs/2311.08152>.
- Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., Dong, G., Wei, H., Lin, H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J., Yang, J., Xu, J., Zhou, J., Bai, J., He, J., Lin, J., Dang, K., Lu, K., Chen, K., Yang, K., Li, M., Xue, M., Ni, N., Zhang, P., Wang, P., Peng, R., Men, R., Gao, R., Lin, R., Wang, S., Bai, S., Tan, S., Zhu, T., Li, T., Liu, T., Ge, W., Deng, X., Zhou, X., Ren, X., Zhang, X., Wei, X., Ren, X., Liu, X., Fan, Y., Yao, Y., Zhang, Y., Wan, Y., Chu, Y., Liu, Y., Cui, Z., Zhang, Z., Guo, Z., and Fan, Z. Qwen2 technical report, 2024. URL <https://arxiv.org/abs/2407.10671>.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zhang, D., Huang, X., Zhou, D., Li, Y., and Ouyang, W. Accessing gpt-4 level mathematical olympiad solutions via monte carlo tree self-refine with llama-3 8b, 2024a. URL <https://arxiv.org/abs/2406.07394>.
- Zhang, D., Zhoubian, S., Hu, Z., Yue, Y., Dong, Y., and Tang, J. Rest-mcts\*: Llm self-training via process reward guided tree search, 2024b. URL <https://arxiv.org/abs/2406.03816>.
- Zhang, H., Duckworth, D., Ippolito, D., and Neelakantan, A. Trading off diversity and quality in natural language generation. *CoRR*, abs/2004.10450, 2020. URL <https://arxiv.org/abs/2004.10450>.
- Zhang, J., Xu, X., Zhang, N., Liu, R., Hooi, B., and Deng, S. Exploring collaboration mechanisms for LLM agents: A social psychology view. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14544–14607, Bangkok, Thailand, August 2024c. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.782. URL <https://aclanthology.org/2024.acl-long.782/>.
- Zhou, A., Yan, K., Shlapentokh-Rothman, M., Wang, H., and Wang, Y.-X. Language agent tree search unifies reasoning acting and planning in language models. *arXiv preprint arXiv:2310.04406*, 2023a.
- Zhou, D., Schärli, N., Hou, L., Wei, J., Scales, N., Wang, X., Schuurmans, D., Cui, C., Bousquet, O., Le, Q. V., and Chi, E. H. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*, 2023b. URL <https://openreview.net/forum?id=WZH7099tgfM>.



## A. Additional Experimental Settings

### A.1. Hyper-parameters

The default sampling parameters for LLM generation are  $\{\text{temperature}=0.75, \text{top\_k}=40, \text{top\_p}=0.95\}$ . Across all MCTS experiments, we set the number of rollouts to 8, the number of sub-questions per node to 4, the number of candidate sub-answers per sub-question to 4, the maximum depth allowed to 5.

### A.2. Dataset Statistics

Since we make use of the rStar code base <sup>2</sup> to implement MOSA, we directly adopt the data files released in their git repository. There are 1,319 instances in GSM8K, 1,000 instances in SVAMP, 500 instances in MATH-500, and 687 instances in StrategyQA.

## B. Additional Implementation Details for Aggregators

In this section, we will show the basic instruction and several in-context learning demonstrations for aggregators.

---

<sup>2</sup><https://github.com/zhentingqi/rStar/>

### Basic Instruction for Aggregators

**\*\*TASK\*\*:**

You are an intelligent and supportive AI assistant. You will receive a collection of responses from various AI assistants regarding a query. Your goal is to synthesize these responses into a single, high-quality response. You should first write down your thoughts as "### Thoughts: <THOUGHTS>.", and then write down the synthesized response in the format of "### Synthesized Response: <SYNTHESIZED RESPONSE>.". If "The answer is: <ANSWER>" appears in the set of responses, you should also write down "The answer is: <ANSWER>"; otherwise, do not write it.

1. Critically assess the information in the responses, as some may be biased, incorrect, or nonsensical.
2. Your synthesized response should not merely echo the provided answers but should present a refined, accurate, and thorough response to the prompt.
3. Ensure that your response is well-organized, coherent, and meets the highest standards of accuracy and reliability.

## An Example ICL Demonstration of MATH-500 for Aggregators: User Query

User:

[Here is the user query that was fed into those AI assistants]:

Given a question, please decompose it into sub-questions. For each sub-question, please answer it in a complete sentence, ending with "The answer is <a numeric answer>". When the original question is answerable, please start the subquestion with "Now we can answer the question: <original question>".

Question 5: Find the product  $CD$  of the integers  $C$  and  $D$  for which

$$\frac{C}{x-3} + \frac{D}{x+8} = \frac{4x-23}{x^2+5x-24}$$

for all real values of  $x$  except  $-8$  and  $3$ .

Question 5.1: What is the equation obtained by multiplying both sides of the given equation by  $(x-3)(x+8)$ ?

Answer 5.1:  $(x-3)(x+8) \cdot \frac{C}{x-3} + (x-3)(x+8) \cdot \frac{D}{x+8} = (x-3)(x+8) \cdot \frac{4x-23}{x^2+5x-24}$

This simplifies to  $C(x+8) + D(x-3) = 4x-23$ .

The answer is  $C(x+8) + D(x-3) = 4x-23$ .

Question 5.2: What are the values of  $C$  and  $D$  that satisfy the equation  $C(x+8) + D(x-3) = 4x-23$  for all  $x$  except  $-8$  and  $3$ ?

Answer 5.2:

[Here are the responses from various AI assistants]:

Response 1:

To find the values of  $C$  and  $D$ , we can substitute specific values of  $x$  into the equation  $C(x+8) + D(x-3) = 4x-23$  to obtain a system of equations.

Question 5.2.1: What happens if we substitute  $x = 3$  into the equation  $C(x+8) + D(x-3) = 4x-23$ ?

Answer 5.2.1: Substituting  $x = 3$  gives  $C(3+8) + D(3-3) = 4 \cdot 3 - 23$ , which simplifies to  $11C = -11$ . The answer is  $11C = -11$ .

Question 5.2.2: What happens if we substitute  $x = -8$  into the equation  $C(x+8) + D(x-3) = 4x-23$ ?

Answer 5.2.2: Substituting  $x = -8$  gives  $C(-8+8) + D(-8-3) = 4 \cdot (-8) - 23$ , which simplifies to  $-11D = -59$ . The answer is  $-11D = -59$ .

Question 5.2.3: Now we can answer the question: What are the values of  $C$  and  $D$ ?

Answer 5.2.3: Solving the system of equations  $11C = -11$  and  $-11D = -59$  gives  $C = -1$  and  $D = 7/11 \cdot 11 = 7$ . The answer is  $-1$  and  $7$  respectively.

Response 2:

To find the values of  $C$  and  $D$ , we can compare the coefficients of the corresponding terms on both sides of the equation.

For the  $x^2$  terms, we have  $C + D = 0$ .

For the  $x$  terms, we have  $8C - 3D = 4$ .

Solving these equations simultaneously, we find that  $C = 3$  and  $D = -3$ .

Therefore, the product  $CD = 3 \cdot (-3) = \boxed{-9}$ .

Response 3:

To find  $C$  and  $D$ , we can substitute  $x = -8$  and  $x = 3$  into the equation  $C(x+8) + D(x-3) = 4x-23$ .

For  $x = -8$ :

$$C(-8+8) + D(-8-3) = 4(-8) - 23$$

$$-11D = -55$$

$$D = 5$$

For  $x = 3$ :

$$C(3+8) + D(3-3) = 4(3) - 23$$

$$11C = -13$$

$$C = -\frac{13}{11}$$

Now we can answer the question: What is the product  $CD$  of the integers  $C$  and  $D$ ?

$$CD = \left(-\frac{13}{11}\right)(5) = -\frac{65}{11}$$

The answer is  $-\frac{65}{11}$ .

Response 4:

To find  $C$  and  $D$ , we can set  $x = 3$  and  $x = -8$  to get two equations:

$$1. C(3+8) + D(3-3) = 4(3) - 23 \Rightarrow 11C = -11 \Rightarrow C = -1$$

$$2. C(-8+8) + D(-8-3) = 4(-8) - 23 \Rightarrow -11D = -55 \Rightarrow D = 5$$

The answer is  $C = -1$  and  $D = 5$ .

## An Example ICL Demonstration of MATH-500 for Aggregators: Assistant Response

**Assistant:****### Thoughts:**

1. The responses generally attempt to solve the equation  $C(x + 8) + D(x - 3) = 4x - 23$  by substituting specific values for  $x$ , namely  $x = 3$  and  $x = -8$ .
2. The correct method involves creating a system of linear equations by substituting these values and solving for  $C$  and  $D$ .
3. There is a discrepancy in the values of  $C$  and  $D$  obtained in the responses:
  - Some responses correctly find  $C = -1$  and  $D = 5$ .
  - Another response incorrectly finds  $C = -13/11$  and  $D = 5$ .
  - Another response finds  $C = 3$  and  $D = -3$ , which is not consistent with the given equation.
4. The correct values should satisfy both equations derived from the substitutions.

**### Synthesized Response:**

The equation  $C(x + 8) + D(x - 3) = 4x - 23$  can be solved by substituting specific values for  $x$ :

1. Substituting  $x = 3$ :

$$C(3 + 8) + D(3 - 3) = 4(3) - 23$$

$$11C = -11$$

$$C = -1$$

2. Substituting  $x = -8$ :

$$C(-8 + 8) + D(-8 - 3) = 4(-8) - 23$$

$$-11D = -55$$

$$D = 5$$

Thus, the values of  $C$  and  $D$  are  $C = -1$  and  $D = 5$ . The product  $CD$  is:

$$CD = (-1) \cdot 5 = -5$$

The answer is: -5



An Example ICL Demonstration of StrategyQA for Aggregators: User Query

**User:**

**[Here is the user query that was fed into those AI assistants]:**

Given a question, please decompose it into sub-questions. For each sub-question, please answer it in a complete sentence. When the original question is answerable, please start the sub-question with "Now we can answer the question: <original question>".

Question 5: Johnny Carson: American talk show host and comedian. Could Johnny Carson's children fill out a water polo team?

Question 5.1: How many children did Johnny Carson have?

Answer 5.1: Johnny Carson had three children: two daughters, Catherine and Christie, and a son, Jeff.

Question 5.2: How many players are needed to form a water polo team?

Answer 5.2:

**[Here are the responses from various AI assistants]:**

Response 1:

A standard water polo team consists of 13 players.

Response 2:

A water polo team consists of seven players, including one goalkeeper and six field players.

Response 3:

A water polo team consists of 13 players: 1 goalkeeper and 12 field players (7 field players and 5 substitutes).

Response 4:

A standard water polo team consists of 13 players, including 7 field players and 6 substitutes.

**An Example ICL Demonstration of StrategyQA for Aggregators: Assistant Response**

**Assistant:**

### Thoughts:

1. The responses contain valuable information about the composition of a water polo team.
2. There is a consensus that a water polo team consists of 13 players in total. However, there is some variation in how these 13 players are categorized.
3. The key point to clarify is how many players are in the water at one time and how many are substitutes.
4. It appears that in the water during play, there are 7 players (including the goalkeeper), and the remaining players are substitutes.

### Synthesized Response:

The synthesized response is:

To form a water polo team, a total of 13 players are required. This team includes 7 players who are in the water during play, comprising one goalkeeper and six field players. The remaining 6 players serve as substitutes.