

# ReSo: A Reward-driven Self-organizing LLM-based Multi-Agent System for Reasoning Tasks

Anonymous ACL submission

## Abstract

Multi-agent systems have emerged as a promising approach for enhancing the reasoning capabilities of large language models in complex problem-solving. However, current MAS frameworks are limited by poor flexibility and scalability, with underdeveloped optimization strategies. To address these challenges, we propose ReSo, which integrates task graph generation with a reward-driven two-stage agent selection process. The core of ReSo is the proposed Collaborative Reward Model, which can provide fine-grained reward signals for MAS cooperation for optimization. We also introduce an automated data synthesis framework for generating MAS benchmarks, without human annotations. Experimentally, ReSo matches or outperforms existing methods. ReSo achieves **33.7%** and **32.3%** accuracy on Math-MAS and SciBench-MAS SciBench, while other methods completely fail. Code is available at: [ReSo](#)

## 1 Introduction

Increasing inference time has emerged as a critical method to enhance the reasoning capabilities of large language models (LLMs)(Snell et al., 2024). Two primary approaches have been explored: (1) optimizing a large reasoning model (Xu et al., 2025) by reinforcement learning and reward models during post-training, which could generate intermediate reasoning steps before answering (OpenAI et al., 2024b; DeepSeek-AI et al., 2025) and (2) leveraging multi-agent system (MAS) collaboration to complete complex tasks that are difficult to solve by single inference (Han et al., 2024; Guo et al., 2024; Wang et al., 2024b; Tran et al., 2025).

Compared to the success of inference time scaling on the single LLM, MAS faces multiple challenges. (1) Most are handcrafted, with limited scalability and adaptability. The lack of an effective agent self-organization mechanism hinders large-scale cooperation. (2) Most assume all agent abilities are fully known while assigning tasks, which

is unrealistic for LLM-based agents. (3) Reward signals are restricted to missing, self-evaluation or outcome only, resulting in poorly defined optimization objectives. (4) Existing MASs lack mechanisms for dynamically optimizing agent networks, making it difficult to achieve data-driven improvements. To address these limitations, we ask: Can we design a self-organizing MAS to learn directly from data via reward signals without handcrafting?

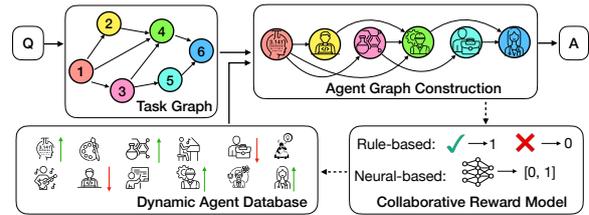


Figure 1: Overview of ReSo pipeline. ReSo first decomposes the task into a DAG; and then constructs an agent graph by topological sorting. First, it searches for agent candidates for each subtask node from the dynamic agent database (DADB). Then it leverages the Collaborative Reward Model (CRM) to choose the best agent and update the agent estimation in DADB.

To realize this potential, we propose ReSo, a reward-driven self-organizing MAS that integrates task graph generation and agent graph construction. The key innovation of our approach is the incorporation of fine-grained reward signals by the Collaborative Reward Model (CRM), which leads to dynamic optimization of agent collaboration. Different from existing MASs, our approach is both scalable and optimizable, achieving state-of-the-art performance on complex reasoning tasks.

While extensive datasets exist for evaluating the reasoning capabilities of LLMs (Chang et al., 2023; Guo et al., 2023), high-quality MAS evaluation benchmarks are scarce. Therefore, we propose an automatic data synthesis method to generate various MAS tasks by converting existing LLM

068 benchmarks into complex collaboration problems.  
069 This method provides step-by-step reward signals  
070 without additional human annotations, enabling  
071 efficient and scalable MAS evaluation.

072 Our contributions can be summarized as:

- 073 • We propose ReSo, the first scalable and opti-  
074 mizable self-organizing MAS framework.
- 075 • We first propose a Collaborative Reward  
076 Model, which can provide fine-grained reward  
077 signals for multi-agent collaboration.
- 078 • We present an automatic data synthesis  
079 method to generate arbitrarily complex MAS  
080 tasks from existing LLM benchmarks.
- 081 • Experimental results demonstrate the superior  
082 performance of ReSo on challenging tasks.

## 083 2 Related Work

### 084 2.1 Reward Guidance

085 The reward model has become a critical compo-  
086 nent in enhancing the capabilities of LLMs through  
087 post-training (Wang et al., 2024d). By providing  
088 feedback on the quality of LLM outputs, RMs facil-  
089 itate performance improvement, enabling models  
090 to generate more accurate and detailed responses.  
091 The concept of reward-guided learning was first  
092 introduced in InstructGPT (Ouyang et al., 2022),  
093 which uses human feedback to fine-tune LLMs,  
094 aligning their behavior with user intent. In addition  
095 to outcome-based supervision, process-based su-  
096 pervision has been shown to improve the reasoning  
097 process itself (Uesato et al., 2022), enhancing not  
098 just the final answer but also the steps leading to it.

099 Building on this, (Lightman et al., 2023) intro-  
100 duced a process reward model (PRM) fine-tuned  
101 on PRM800K, which provides fine-grained and  
102 interpretable rewards for every reasoning step.  
103 Similarly, (Wang et al., 2024c) developed Math-  
104 Shepherd, an approach capable of autonomously  
105 generating process supervision data. Despite the ad-  
106 vantages of neural-based reward models in terms of  
107 generalization, they also suffer from reward hack-  
108 ing (Gao et al., 2022; Skalse et al., 2022). To  
109 mitigate this, some recent approaches have em-  
110 ployed rule-based rewards (DeepSeek-AI et al.,  
111 2025) or fixed inference budgets (Muennighoff  
112 et al., 2025), which have also proven effective. No-  
113 tably, DeepSeek-R1 (DeepSeek-AI et al., 2025)  
114 incorporates both output accuracy and reasoning  
115 format evaluation, achieving the performance on  
116 par with OpenAI-O1 (OpenAI et al., 2024b; Qin

et al., 2024). DeepSeek-R1 demonstrates that only  
using large-scale reinforcement learning based on  
rule-based reward during post-training can stimu-  
late LLM’s excellent reasoning ability, without  
supervised fine-tuning.

### 117 2.2 Multi-Agent System 122

123 Recent advances in LLM-based MAS have raised  
124 expectations for their ability to tackle increasingly  
125 complex reasoning tasks (Han et al., 2024; Guo  
126 et al., 2024; Wang et al., 2024b; Tran et al., 2025).

127 Predefined cooperation in MAS relies on struc-  
128 tured interactions and role assignments before col-  
129 laboration. Early works focus on MAS infrastruc-  
130 ture, including Camel, AutoGen, and AgentVerse  
131 (Li et al., 2023; Wu et al., 2023; Chen et al., 2023).  
132 Some approaches adopt standard operating proce-  
133 dures for structured task decomposition, as seen in  
134 MetaGPT and ChatDev (Hong et al., 2024; Qian  
135 et al., 2024a; Dong et al., 2024). Fixed topologies  
136 are most adopted, such as hierarchical structures  
137 in MOA (Wang et al., 2024a) and directed acyclic  
138 graphs in MacNet and MAGDI (Qian et al., 2024b;  
139 Chen et al., 2024c). Predefined role interactions are  
140 also widely used such as debate (Du et al., 2023),  
141 criticism (Chen et al., 2024b), and certain math rea-  
142 soning patterns (Gou et al., 2024; Lei et al., 2024;  
143 Xi et al., 2024). Predefined MASs exhibit several  
144 limitations including: (1) Scalability and adaptabil-  
145 ity being constrained by the imposition of rigid role  
146 assignments and fixed topological structures. (2)  
147 The unrealistic assumption that the agent’s abilities  
148 are fully known when assigning tasks, which is  
149 particularly problematic for LLM-based agents.

150 Optimizable cooperation in MAS aims to dynam-  
151 ically adapt interaction topology and agent roles.  
152 GPTSwarm (Zhuge et al., 2024) formulates MAS  
153 as optimizable computational graphs, refining node  
154 prompts and inter-agent connectivity via evolution-  
155 ary algorithms. DyLAN (Liu et al., 2024b) em-  
156 ploys a layerwise feedforward agent network and a  
157 mutual rating mechanism to dynamically optimize  
158 MAS. G-Designer (Zhang et al., 2025a) utilizes  
159 variational graph auto-encoders to optimize MAS.  
160 Current optimizing approaches are highly under-  
161 explored. They often lack reliable, fine-grained  
162 reward signals for MAS collaboration, relying in-  
163 stead on outputs or self-generated reward mecha-  
164 nisms. Meanwhile, dynamic network optimization  
165 algorithms for MAS are also lacking.

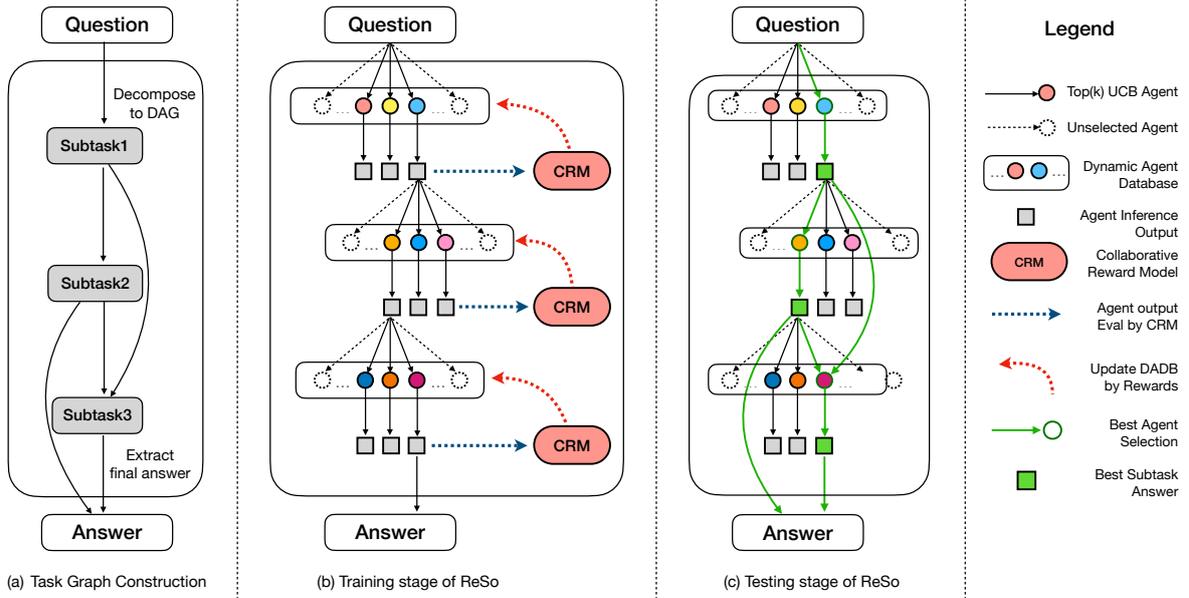


Figure 2: Illustration of our proposed ReSo. (a) We decompose the question into a subtask DAG. (b) The training of ReSo: we first use the UCB score to perform a coarse search in DADB and select top-k agents, then score the inference results using CRM, and update DADB by rewards. Repeat the above process for each node in DAG by topological order. (c) The testing of ReSo: we directly select the best agent from DADB without CRM. The determined agent collaboration pattern is the best path with the highest UCB score in the decision tree.

### 3 Methods

To tackle the existing challenges in MAS research, we propose two core innovations: (1) ReSo, a reward-driven self-organizing MAS, which is capable of autonomously adapting to complex tasks and a flexible number of agent candidates, eliminating the need for handcrafted solutions. (2) Introduction of a Collaborative Reward Model (CRM), specifically tailored to optimize MAS performance. CRM can deliver fine-grained reward signals on multi-agent collaboration, enabling data-driven MAS performance optimization.

#### 3.1 Problem Formulation

We define a MAS algorithm  $f_{MAS}$  as a function that, given a natural language question  $Q$ , generates a graph-structured task decomposition, solves each subtask, and produces a final answer:

$$f_{MAS}(Q) \rightarrow (G = (V, E), A_V, A_Q) \quad (1)$$

Here,  $G = (V, E)$  represents the task decomposition graph, which is structured as a directed acyclic graph (DAG). The set of nodes  $V = \{v_1, v_2, \dots, v_n\}$  corresponds to the subtasks derived from  $Q$ , while the edges  $E \subseteq V \times V$  define the dependencies between these subtasks. The system produces subtask answers  $A_V =$

$\{a_{v_1}, a_{v_2}, \dots, a_{v_n}\}$  and ultimately derives the final answer  $A_Q$ . To achieve this, we decompose  $f_{MAS}$  into two sub-algorithms:

$$f_{MAS}(Q) = f_{agent} \circ f_{task}(Q) \quad (2)$$

$f_{task}$  is responsible for constructing the task decomposition graph from the input question, ensuring a structured breakdown of the problem into subtasks and dependencies.  $f_{agent}$  dynamically selects and assigns appropriate agents to solve the identified subtasks. This modular design enables independent optimization of each component, allowing for greater flexibility and scalability.

For the MAS-generated answer  $A_Q$  to be considered correct, the following conditions must be satisfied: (1) All subtask answers must be correct. (2) All directed edges must correctly enforce the dependency relationships among subtasks. (3) The final output  $A_Q$  must be correct.

#### 3.2 Task Graph Construction

In the proposed method,  $f_{task}$  first transforms the question  $Q$  into a directed acyclic task graph  $G$ :

$$f_{task} : Q \rightarrow G = (V, E) \quad (3)$$

where  $G$  represents the decomposition of the original task  $Q$ . Each node  $v_i \in V$  is a natural language

subtask, and each directed edge  $(v_i \rightarrow v_j) \in E$  indicates that the subtask  $v_j$  depends on the successful completion of  $v_i$ .

In practice, we perform supervised fine-tuning (SFT) on an LLM to perform this step of task decomposition. Using our synthetic data, we explicitly require the LLM to decompose  $Q$  into logical sub-problems, specify their execution order and dependencies, and output in a format of DAG.

### 3.3 Two-Stage Agent Search

Once the task graph is obtained, we need to assign each subtask to the most appropriate agent. We denote this agent assignment procedure as  $f_{agent}$ . Conceptually,  $f_{agent}$  classifies each node in the task graph according to the most suitable agent from a large agent pool  $\mathcal{A}$ , constructing an *agent graph* that maps each node to one or more selected agents.

$$f_{agent} : v_i \in V \rightarrow a_i \in \mathcal{A} \quad (4)$$

Since  $\mathcal{A}$  can contain a large number of agents, we first introduce the concept of Dynamic Agent Database. Then we decompose the agent graph construction on every subtask into two search algorithms from coarse to fine-grained: first, select a subset of candidates from DADB then utilize the reward model to evaluate and select the best agent.

#### 3.3.1 Dynamic Agent Database

To increase MAS’s scalability and flexibility, we propose the Dynamic Agent Database (DADB), denoted as  $\mathcal{A}$ , which enables adaptive agent selection by maintaining both **static** and **dynamic** agent profiles. For each agent  $a_i \in \mathcal{A}$ , its static profile includes the base model, role settings, initial prompt, long-term memory, and tools. The dynamic profile, continuously updated via the reward model, tracks the agent’s average reward  $R(a_i)$ , computational cost  $C(a_i)$ , and task count  $n(a_i)$ . Initially, agents have only static attributes, while training iteratively refines their evaluations by the process reward model, optimizing future selection.

Given an input task  $v_j$ , the DADB assigns a preliminary quality score  $Q(a_i, v_j)$  to each agent  $a_i$ , balancing task-agent similarity, historical performance, and computational costs:

$$Q(a_i, v_j) = \text{sim}(a_i, v_j) \cdot \text{perform}(a_i) \quad (5)$$

where  $\text{sim}(a_i, v_j)$  represents the similarity between the subtask’s target profile and the agent’s static profile. In practice, we employ a Heaviside function which ensures that only agents exceeding a

predefined similarity threshold  $V_{th}$  are considered:  $\text{sim}(a_i, v_j) = H[\langle \mathbf{q}_i, \mathbf{a}_i \rangle - V_{th}]$  where  $\mathbf{q}_i, \mathbf{a}_i$  are text embedding of subquestion and the agent static profile. The  $\text{perform}(a_i)$  term is given by  $\text{perform}(a_i) = R(a_i) - \beta C(a_i)$ , where  $\beta$  controls the trade-off between the agent’s historical performance and cost.

#### 3.3.2 Coarse Agent Search by UCB

Given a DADB  $\mathcal{A}$  and a subtask  $v_j$ , our first objective is to retrieve a promising subset of  $k$  candidate agents. To take advantage of the known information in DADB, also to explore unused agents, we adopt an Upper Confidence Bound value:

$$\text{UCB}(a_i, q_j) = Q(a_i, q_j) + c \sqrt{\frac{N}{n(a_i) + \varepsilon}} \quad (6)$$

where  $N$  is the total number of agent selections and  $n(a_i)$  the number of times agent  $i$  is selected,  $\varepsilon \ll 1$ .  $c$  is a constant controlling the exploration-exploitation trade-off. Agents with higher UCB scores are more likely to be selected, helping the MAS to explore potentially underutilized agents. For each subtask  $q_i$ , we sort agents by their  $\text{UCB}(a_i, q_j)$  and choose the top  $k$  agents as the candidate set  $\mathcal{A}_{\text{cand}} = \{a_1, a_2, \dots, a_k\}$ .

#### 3.3.3 Fine-grained Agent Evaluation by CRM

Once the candidate agents  $\mathcal{A}_{\text{cand}}$  are selected, we evaluate their performance on the current subtask  $v_j$  using a Collaborative Reward Model (CRM). This evaluation process is straightforward: each candidate agent  $a_i$  generates an answer to the subtask  $v_j$ :  $a_i(v_j)$ , and then we assess the quality of that answer based on a reward signal:

$$r(a_i, v_j) = \text{RewardModel}(a_i, v_j, a_i(v_j)) \quad (7)$$

where  $\text{RewardModel}$  evaluates the quality of the solution based on the given agent’s profile, subtask, and previous reasoning process. After evaluating the agents, we assign the agent with the highest reward,  $a_j^*$ , to the subtask node  $v_j$ , which means  $a_j^*$ ’s solution is used as  $v_j$ ’s answer. This process is repeated for each subtask on the graph.

The reward  $r(a_i, v_j)$  is computed using the CRM, which can be either rule-based (e.g., binary correctness: 0 for incorrect, 1 for correct) or neural-based (providing a score between 0 and 1 for quality). The reward model evaluates how well the agent’s response aligns with the expected outcome, factoring in both the solution’s correctness and its collaboration within the MAS.

### 3.4 Training and Inference Stage

Our multi-agent system can operate in two modes: training and testing. During **training**, we leverage a high-quality reward  $r(a_i, v_j)$  available for evaluating the correctness of every step of MAS. Upon receiving  $r(a_i, v_j)$  for each candidate agent, we update that agent’s dynamic profile in DADB. For instance, we may maintain a running average of rewards:

$$R(a_i) \leftarrow \frac{n(a_i) \cdot R(a_i) + r(a_i, v_j)}{n(a_i) + 1} \quad (8)$$

similar for updating  $costc(a_i, v_j)$ . By iteratively learning from data, the DADB can dynamically update agent evaluations based on historical reward, facilitating adaptive agent selection and improving both efficiency and performance. During **testing**, the reward model is no longer required. Instead, we leverage the learned DADB to select the best agent candidates and the best answer to each subtask.

### 3.5 The Perspective of MCTS

The task graph, after topological sorting, forms a decision tree where each node represents a subtask and the edges denote dependencies. At each level, we use UCB to prune the tree and select a subset of promising agents, then simulate each agent and evaluate their performance using the CRM. The resulting reward updates the agent’s dynamic profile, refining the selection strategy. The MAS construction is essentially finding the optimal path from the root to the leaves, maximizing the UCB reward for the best performance.

Consider there are  $N$  agents and a task requiring  $D$  agents to collaborate. Assume that the average inference cost is  $c$  and the matching cost in DADB is  $s \ll c$  per agent. A brute-force search has a complexity of  $O(c \cdot N^D)$ , which becomes infeasible as  $D$  and  $D$  grow. In contrast, our self-organizing strategy, selecting top  $k$  per step, reduces the cost to  $O((s \cdot N + N \log N + k \cdot c) \cdot D)$ , offering a near-linear scaling with  $N$  and  $D$ , making the approach highly scalable for large  $N$  and  $D$ .

## 4 Data Synthesis

A key challenge in MAS is the lack of structured datasets for evaluating and training agent collaboration. To address this, we propose an automated framework that converts existing LLM datasets into structured, multi-step MAS tasks, enabling fine-grained evaluation without human annotations.

**Random DAG Generation** We begin by generating a DAG,  $G = (V, E)$ . Each node  $v_i \in V$  will be filled with a subtask  $(q_i, a_i)$ , where  $q_i$  is the textual description of the task, and  $a_i$  is its numerical answer. The subtasks are sampled from the existing LLM benchmarks. The edges  $E$  will encode dependency constraints between subtasks, ensuring that the solution to one subtask is required as an input for another, modeling the sequential reasoning process of multi-agent collaboration.

**Subtask Selection and Filling** To populate the nodes of  $G$ , we construct a master pool of candidate subtasks, denoted as  $\mathcal{P}$ . Each candidate subtask  $p_i \in \mathcal{P}$  consists of a textual problem description  $s_i$ , and a numerical answer  $a_i$ . After obtaining  $\mathcal{P}$ , we randomly sample from it and fill one question per node into the generated DAG. Candidate subtasks should have clear numerical or option answers, such as SciBench (Wang et al., 2024f), Math (Hendrycks et al., 2021), GPQA (Rein et al., 2023), etc. To ensure that the problem is computationally feasible for later dependency construction, we extract a numerical constant  $c_i \in \mathbb{R}$  from the problem text. If the extracted constant is valid, the subtask is retained in  $\mathcal{P}$ ; otherwise, it is discarded. This ensures that only problems with well-defined numerical attributes are incorporated.

**Dependency Edge Construction** After all nodes are populated, we generate natural language dependency descriptions for edges. Each edge  $(v_j \rightarrow v_k)$  should represent a relationship which connects previous subtask  $v_j$ ’s answer  $a_j$ , with subsequent subtask  $v_k$ ’s question parameter  $c_k$ . For each edge, we generate a textual description  $e_{jk}$ , such as “in this question,  $c_k = \text{previous answer} + 3$ .” Formally, it is an algorithm that constructs a string from two numbers:  $e_{ij} = f(a_j, c_k)$ .  $f$  can be implemented using elementary arithmetic and text templates, ensuring that no answers or parameters in the original subtask need to be manually modified. Once the DAG is fully constructed, we refine node descriptions by removing any explicitly given numerical constants  $\{c_i\}$  that are now dependent on the results of prior nodes. Finally, an entire graph described in natural language is a piece of synthetic data.

The proposed data synthesis framework generates structured, multi-step reasoning tasks with adjustable sizes, ensuring diverse and scalable problem structures. The synthesized dataset supports both training and testing, enabling fine-grained evaluation without human annotations.

Method	Math-MAS				SciBench-MAS			
	Easy	Medium	Hard	<i>Tokens</i>	Easy	Medium	Hard	<i>Tokens</i>
GPT-4o	27.5	9.0	0.0	2.2k	39.3	12.5	1.6	2.1k
Gemini-2.0-Flash	<u>69.2</u>	<u>24.7</u>	9.0	3.0k	<u>64.5</u>	<u>33.8</u>	9.7	2.5k
Claude-3.5-Sonnet	12.1	0.0	0.0	1.0k	22.4	6.2	3.2	1.4k
Qwen2.5-Max	44.0	13.5	4.5	2.9k	55.1	30.0	4.8	2.8k
DeepSeek-V3	52.7	<u>24.7</u>	12.4	2.2k	52.3	31.3	<u>12.9</u>	2.3k
MetaGPT	30.8	12.4	2.2	16.1k	48.6	2.5	0.0	14.6k
DyLAN	40.7	9.0	0.0	64.1k	48.6	2.5	0.0	77.8k
GPTSwarm	35.2	5.6	4.5	14.9k	31.8	6.3	1.6	18.2k
GDesigner	14.2	5.6	0.0	16.9k	24.3	12.5	0.0	19.0k
<b>ReSo (ours)</b>	<b>79.1</b>	<b>56.2</b>	<b>33.7</b>	<b>14.6k</b>	<b>67.3</b>	<b>51.3</b>	<b>32.3</b>	<b>20.7k</b>

Table 1: Accuracy and average token usage on Math-MAS and SciBench-MAS. Bold and underlined represent optimal and suboptimal results, respectively. *Tokens* denotes the average number of tokens consumed per task.

## 5 Experiments

In Sec 5.1, we first use public datasets to create complex MAS benchmarks and fine-tune ReSo’s task decomposition and collaborative reward models. All code, datasets, and models are publicly available. In 5.2, we train and evaluate ReSo on both public and synthetic datasets. Sec 5.3 presents ablation studies on task decomposition, agent selection, and reward guidance mechanisms.

### 5.1 Data Synthesis and Model Fine-tuning

#### 5.1.1 Data Synthesis

MATH (Hendrycks et al., 2021) consists of problems from diverse mathematical domains, while SciBench (Wang et al., 2024f) includes scientific reasoning tasks spanning physics, chemistry, and mathematics. Using these datasets, we apply the synthetic data generation method outlined in Sec 4 to create two datasets: one for single LLM fine-tuning and another for benchmarking. Difficulty is categorized by the number of subtasks—Easy (3), Medium (5), and Hard (7).

**Fine-tuning data** For fine-tuning task decomposition LLM, we generate 14,500 questions and answers from the MATH training set, with numbers of subtasks ranging from 2 to 6. For fine-tuning the neural-based CRM, we generate 5,000 questions from the same set, with 5 subtasks per question.

**MAS Benchmarks** We select 201 questions from SciBench as the sub-question data pool and synthesized complex data using the method in 4. This forms the SciBench-MAS dataset, comprising 200 easy-level training questions and 247 testing

questions (107 easy, 80 medium, 62 hard). For MATH (Hendrycks et al., 2021), 348 level-5 questions are selected, from which we generate the Math-MAS dataset, consisting of 269 test questions for ReSo (91 easy, 89 medium, 89 hard).

#### 5.1.2 Model Fine-tuning

**Task Decomposition Model Training** To ensure high-quality task composition, we fine-tune a specialized model for task decomposition based on Qwen2.5-7B-Instruct. We use 14500 dialogues on task decomposition as described in 5.1.1, and fine-tune the model under a batch size of 128 and a learning rate of 1e-4 for 3 epochs. The fine-tuned model can reliably produce task decomposition in a structured format.

**CRM Training** The proposed CRM is fine-tuned based on Qwen2.5-Math-PRM-7B (Zhang et al., 2025b), which can provide effective process reward signals on MAS collaborative reasoning tasks. We use 5000 samples of sub-tasks with their answers as described in 5.1.1. We follow a simplified training scheme of PRMs, where the model should only perform binary classification on the special token at the end of the answer. The model is trained with a batch size of 128 and a learning rate of 1e-4 for 5 epochs. The fine-tuned model can output the probability of the answer being correct, which is then taken as the collaborative reward signal.

### 5.2 Main Results of ReSo

**Models and MASs** We compare ReSo with state-of-the-art LLM and MAS methods. Our single-LLM baselines include GPT-4o (OpenAI et al., 2024a), Gemini-2.0-Flash (Team et al., 2024),

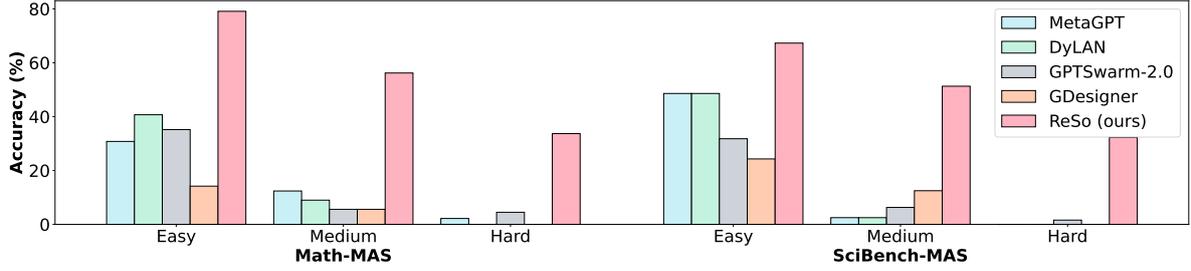


Figure 3: ReSo outperforms other MAS methods by a significant margin in complex reasoning accuracy.

Claude-3.5-Sonnet (Anthropic, 2024), Qwen2.5-Max (Yang et al., 2024), DeepSeek-V3 (Liu et al., 2024a). For ReSo, we build an agent database that includes these base models, extended to 63 agents with different prompts. For MAS, we evaluate MetaGPT (Hong et al., 2024), DyLAN (Liu et al., 2024b), GPTSwarm (Zhuge et al., 2024), GDesigner (Zhang et al., 2025a), SEDM (Li et al., 2024b). All MAS baselines use GPT-4o as the backbone.

**ReSo Training** We train our ReSo framework using the SciBench-MAS training data as described in 5.1.1. Figure 4 shows that ReSo’s accuracy increases with the training process, demonstrating that DADB effectively updates the estimation of each agent’s capability and gradually learns to build a better agent graph.

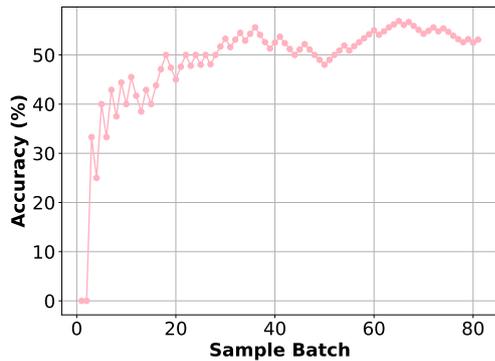


Figure 4: Training Curve of ReSo.

**Comparisons with LLMs** As shown in Table 1, most single-model agents exhibit a sharp decrease in accuracy as the difficulty increases. At the hard difficulty level, their accuracy approaches zero, suggesting that single LLMs struggle with compositional reasoning. In particular, we show the results of these single LLMs on single Math and Scibench datasets in Appendix A.2, with accuracy rates of

80%-90%. This means that a single LLM can successfully solve a single sub-problem in the dataset, but its generalization ability for combined complex problems is very limited.

**Comparisons with MAS** Notably, ReSo outperforms other approaches in both the Math-MAS and SciBench-MAS datasets. At the hard difficulty level, ReSo reaches an accuracy of **33.7%** on Math-MAS and **32.3%** on SciBench-MAS, while other MAS methods almost completely fail.

**Token Efficiency** Table 1 also compares the average number of tokens consumed per task. ReSo maintains a relatively moderate token usage, which is significantly lower than certain baselines like DyLAN (14.6k vs 64.1k, 20.7k vs 77.8k). This balance between performance and computational cost underlines ReSo’s practical efficiency in real-world, large-scale scenarios.

**Results on Existing Benchmarks** Our method excels not only on complex task datasets but also on existing commonly used benchmarks. Table 2 shows our evaluation of the original MATH and SciBench datasets, where ReSo (ours) achieves the highest accuracy across all the tasks. Notably, it outperforms GPT-4o and other baselines, reaching 89.8% on MATH and leading across SciBench categories. These results demonstrate ReSo’s strong generalization and effectiveness in mathematical and scientific reasoning.

Table 2: Accuracy on existing benchmarks.

Method	MATH	SciBench		
		Math	Phys	Chem
GPT-3.5-Turbo	34.1	25.56	14.83	32.11
GPT-4o	81.1	66.8	53.4	60.1
SEDM	-	61.4	50.3	56.1
GPTSwarm	81.0	60.5	36.6	49.7
<b>ReSo (ours)</b>	<b>89.8</b>	<b>71.9</b>	<b>60.6</b>	<b>61.9</b>

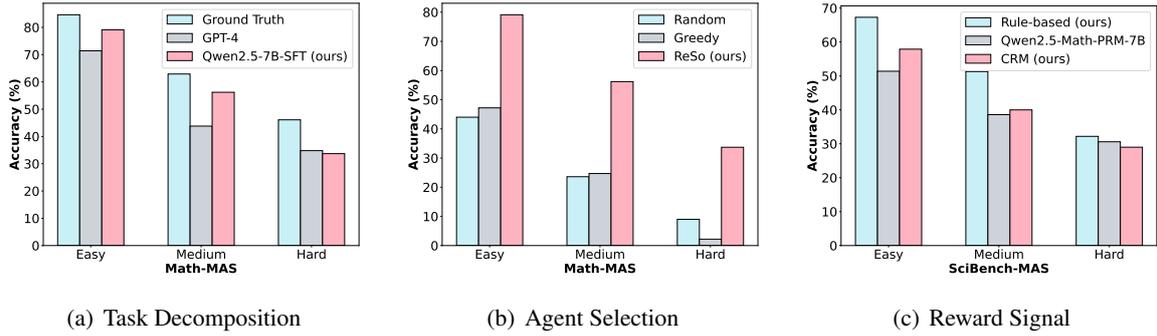


Figure 5: Results of ablation studies. (a) Fine-tuning on domain-specific training data can significantly improve the decomposition quality, thus enhancing overall system performance. (b) Our robust agent selection strategy within the MAS is significant to the performance. (c) Compared to general reward models, our fine-tuned reward model is more task-specific and brings more precise reward signals, thus improving the system performance.

### 5.3 Ablation Studies

We conduct ablation studies on our proposed multi-agent system, examining three core designs: task decomposition, agent selection, and reward signal.

**Task Decomposition** We compare three different approaches to task decomposition: (1) **Ground Truth**, representing an upper bound with human-crafted, meticulously designed task breakdowns; (2) **GPT-4**, which autonomously decomposes complex tasks into sub-tasks without targeted fine-tuning; and (3) **Qwen2.5-7B-SFT**, a model fine-tuned on our dataset based on Qwen2.5-7B, specifically adapted to generate more effective decompositions for complex questions. Figure 5(a) presents the reasoning accuracy under different decomposition strategies. The ground-truth decomposition consistently yields the highest accuracy, underscoring the critical role of precise subproblem segmentation. Meanwhile, the fine-tuned task generator surpasses the naive GPT-4 approach, demonstrating that even a small amount of domain-specific training data can significantly improve decomposition quality and enhance overall system performance.

**Agent Selection** We compare three strategies for agent selection: a **random** strategy, a **greedy** strategy that always selects the most matching profile, and our proposed **ReSo** approach. As shown in Figure 5(b), **ReSo** significantly outperforms other strategies across all the datasets, which emphasizes the importance of a robust agent selection strategy within the multi-agent framework. By strategically assigning each sub-task to the most suitable agent, the system can handle increasingly complex tasks with markedly better accuracy.

**Reward Signal** We investigate the impact of different reward signals on system optimization, considering three approaches: (1) **Rule-based**, which provides strictly accurate, predefined evaluations for sub-task solutions; (2) **General Reward Model**, using Qwen2.5-Math-PRM-7B as a reward function without task-specific fine-tuning; and (3) **Fine-tuned Reward Model**, i.e., our CRM proposed in Section 3.3.3. Figure 5(c) presents the results of training our MAS under these reward schemes on the SciBench-MAS dataset. The rule-based reward yields the best results, confirming the importance of precise reward signals. Besides, our CRM brings a slight improvement compared to the original Qwen2.5-Math-PRM-7B model. We also observe an instance of *reward hacking* when using the Qwen reward model: specifically, Qwen2.5-Max tends to receive inflated scores when acting as the reasoning agent. As a result, during inference, the MAS disproportionately selects Qwen2.5-Max to handle sub-tasks, even in cases where it does not necessarily produce the best solutions.

## 6 Conclusion

In this work, we introduce ReSo, a reward-driven self-organizing MAS for complex reasoning. By integrating a collaborative reward model, ReSo automates agent selection and collaboration, improving scalability and adaptability. The automated data synthesis framework eliminates manual annotations. Experiments show that ReSo outperforms existing MAS and single LLM baselines. All codes, models, and data have been open-sourced. We expect ReSo to enable co-optimization of MAS and LLM to further enhance reasoning capabilities.

## 7 Limitations

Although the base model for the agents is a fixed model, ReSo has demonstrated strong optimizability and scalability as well as good performance. A further interesting research question is: Can the optimization of MAS be performed together with the optimization of a single LLM agent? Specifically, can the reward signal given to the model by our CRM in each step of cooperation be combined with the reinforcement learning-based post-training of a single model to further optimize MAS at both the macro and micro levels? This means a dynamic agent cooperation network, where agents can not only learn how to interact with each other but also fine-tune their weights through feedback from cooperation. We look forward to follow-up research.

## 8 Ethical Considerations

While our proposed ReSo framework focuses on reasoning tasks in the domains of mathematics and science, it has the potential to be applied in other, possibly unethical, contexts. Such misuse could pose significant threats to human society. We strongly urge readers to carefully consider these ethical implications and to adopt a conscientious approach in the development and application of these methods.

## References

AI Anthropic. 2024. [The claude 3 model family: Opus, sonnet, haiku](#). *Claude-3 Model Card*.

Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I 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*, 4(1):1–43.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. [A survey on evaluation of large language models](#). *Preprint*, arXiv:2307.03109.

Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. 2024a. [Alphamath almost zero: Process supervision without process](#). *Preprint*, arXiv:2405.03553.

Justin Chih-Yao Chen, Archiki Prasad, Swarnadeep Saha, Elias Stengel-Eskin, and Mohit Bansal. 2024b. [Magicore: Multi-agent, iterative, coarse-to-fine refinement for reasoning](#). *Preprint*, arXiv:2409.12147.

Justin Chih-Yao Chen, Swarnadeep Saha, Elias Stengel-Eskin, and Mohit Bansal. 2024c. [Magdi: Structured distillation of multi-agent interaction graphs improves reasoning in smaller language models](#). *Preprint*, arXiv:2402.01620.

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2023. [Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors](#). *Preprint*, arXiv:2308.10848.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wangjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu

706	Zhang, and Zhen Zhang. 2025. <a href="#">Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning</a> . <i>Preprint</i> , arXiv:2501.12948.	759
707		760
708		761
709	Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2024. <a href="#">Self-collaboration code generation via chatgpt</a> . <i>Preprint</i> , arXiv:2304.07590.	762
710		763
711		764
712	Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. <a href="#">Improving factuality and reasoning in language models through multiagent debate</a> . <i>Preprint</i> , arXiv:2305.14325.	765
713		766
714		767
715		768
716	Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. 2024. <a href="#">Alphazero-like tree-search can guide large language model decoding and training</a> . <i>Preprint</i> , arXiv:2309.17179.	769
717		770
718		771
719		772
720		773
721	Leo Gao, John Schulman, and Jacob Hilton. 2022. <a href="#">Scaling laws for reward model overoptimization</a> . <i>Preprint</i> , arXiv:2210.10760.	774
722		775
723		776
724	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. <a href="#">Tora: A tool-integrated reasoning agent for mathematical problem solving</a> . <i>Preprint</i> , arXiv:2309.17452.	777
725		778
726		779
727		780
728		781
729	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. <a href="#">Large language model based multi-agents: A survey of progress and challenges</a> . <i>Preprint</i> , arXiv:2402.01680.	782
730		783
731		784
732		785
733		786
734	Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiakuan Li, Bo-jian Xiong, and Deyi Xiong. 2023. <a href="#">Evaluating large language models: A comprehensive survey</a> . <i>Preprint</i> , arXiv:2310.19736.	787
735		788
736		789
737		790
738		791
739	Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, Zhaozhuo Xu, and Chaoyang He. 2024. <a href="#">Llm multi-agent systems: Challenges and open problems</a> . <i>arXiv preprint arXiv:2402.03578</i> .	792
740		793
741		794
742		795
743	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. <a href="#">Measuring mathematical problem solving with the math dataset</a> . <i>Preprint</i> , arXiv:2103.03874.	796
744		797
745		798
746		799
747		800
748	Sirui Hong, Mingchen Zhuge, Jiaqi Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2024. <a href="#">Metagpt: Meta programming for a multi-agent collaborative framework</a> . <i>Preprint</i> , arXiv:2308.00352.	801
749		802
750		803
751		804
752		805
753		806
754		807
755	Bin Lei, Yi Zhang, Shan Zuo, Ali Payani, and Caiwen Ding. 2024. <a href="#">Macm: Utilizing a multi-agent system for condition mining in solving complex mathematical problems</a> . <i>Preprint</i> , arXiv:2404.04735.	808
756		809
757		810
758		811
		812
		813
		814
		815
	Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. <a href="#">Camel: Communicative agents for "mind" exploration of large language model society</a> . <i>Preprint</i> , arXiv:2303.17760.	
	Qingyao Li, Wei Xia, Kounianhua Du, Xinyi Dai, Ruiming Tang, Yasheng Wang, Yong Yu, and Weinan Zhang. 2024a. <a href="#">Rethinkmcts: Refining erroneous thoughts in monte carlo tree search for code generation</a> . <i>Preprint</i> , arXiv:2409.09584.	
	Ziyue Li, Yuan Chang, and Xiaoqiu Le. 2024b. <a href="#">Simulating expert discussions with multi-agent for enhanced scientific problem solving</a> . In <i>Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)</i> , pages 243–256.	
	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. <a href="#">Let’s verify step by step</a> . <i>Preprint</i> , arXiv:2305.20050.	
	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024a. <a href="#">Deepseek-v3 technical report</a> . <i>arXiv preprint arXiv:2412.19437</i> .	
	Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. 2024b. <a href="#">A dynamic llm-powered agent network for task-oriented agent collaboration</a> . <i>Preprint</i> , arXiv:2310.02170.	
	Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Meiqi Guo, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi. 2024. <a href="#">Improve mathematical reasoning in language models by automated process supervision</a> . <i>Preprint</i> , arXiv:2406.06592.	
	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. <a href="#">s1: Simple test-time scaling</a> . <i>Preprint</i> , arXiv:2501.19393.	
	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger,	

816	Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin	ner, Michael Lampe, Michael Petrov, Michael Wu,	880
817	Zweig, Beth Hoover, Blake Samic, Bob McGrew,	Michele Wang, Michelle Fradin, Michelle Pokrass,	881
818	Bobby Spero, Bogo Giertler, Bowen Cheng, Brad	Miguel Castro, Miguel Oom Temudo de Castro,	882
819	Lightcap, Brandon Walkin, Brendan Quinn, Brian	Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-	883
820	Guarraci, Brian Hsu, Bright Kellogg, Brydon East-	nal Khan, Mira Murati, Mo Bavarian, Molly Lin,	884
821	man, Camillo Lugaresi, Carroll Wainwright, Cary	Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-	885
822	Bassin, Cary Hudson, Casey Chu, Chad Nelson,	talie Cone, Natalie Staudacher, Natalie Summers,	886
823	Chak Li, Chan Jun Shern, Channing Conger, Char-	Natan LaFontaine, Neil Chowdhury, Nick Ryder,	887
824	lotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,	Nick Stathas, Nick Turley, Nik Tezak, Niko Felix,	888
825	Chong Zhang, Chris Beaumont, Chris Hallacy, Chris	Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel	889
826	Koch, Christian Gibson, Christina Kim, Christine	Bundick, Nora Puckett, Ofir Nachum, Ola Okelola,	890
827	Choi, Christine McLeavey, Christopher Hesse, Clau-	Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins,	891
828	dia Fischer, Clemens Winter, Coley Czarnecki, Colin	Olivier Godement, Owen Campbell-Moore, Patrick	892
829	Jarvis, Colin Wei, Constantin Koumouzelis, Dane	Chao, Paul McMillan, Pavel Belov, Peng Su, Pe-	893
830	Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy,	ter Bak, Peter Bakkum, Peter Deng, Peter Dolan,	894
831	David Carr, David Farhi, David Mely, David Robin-	Peter Hoeschele, Peter Welinder, Phil Tillet, Philip	895
832	son, David Sasaki, Denny Jin, Dev Valladares, Dim-	Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming	896
833	itris Tsipras, Doug Li, Duc Phong Nguyen, Duncan	Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra-	897
834	Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-	jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul	898
835	dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow,	Puri, Reah Miyara, Reimar Leike, Renaud Gaubert,	899
836	Eric Kramer, Eric Peterson, Eric Sigler, Eric Wal-	Reza Zamani, Ricky Wang, Rob Donnelly, Rob	900
837	lace, Eugene Brevdo, Evan Mays, Farzad Khorasani,	Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-	901
838	Felipe Petroski Such, Filippo Raso, Francis Zhang,	dani, Romain Huet, Rory Carmichael, Rowan Zellers,	902
839	Fred von Lohmann, Freddie Sulit, Gabriel Goh,	Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan	903
840	Gene Oden, Geoff Salmon, Giulio Starace, Greg	Cheu, Saachi Jain, Sam Altman, Sam Schoenholz,	904
841	Brockman, Hadi Salman, Haiming Bao, Haitang	Sam Toizer, Samuel Miserendino, Sandhini Agar-	905
842	Hu, Hannah Wong, Haoyu Wang, Heather Schmidt,	wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean	906
843	Heather Whitney, Heewoo Jun, Hendrik Kirchner,	Grove, Sean Metzger, Shamez Hermani, Shantanu	907
844	Henrique Ponde de Oliveira Pinto, Hongyu Ren,	Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi-	908
845	Huiwen Chang, Hyung Won Chung, Ian Kivlichan,	rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay,	909
846	Ian O’Connell, Ian O’Connell, Ian Osband, Ian Sil-	Srinivas Narayanan, Steve Coffey, Steve Lee, Stew-	910
847	ber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya	art Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao	911
848	Kostrikov, Ilya Sutskever, Ingmar Kanitscheider,	Xu, Tarun Gogineni, Taya Christianson, Ted Sanders,	912
849	Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub	Tejal Patwardhan, Thomas Cunningham, Thomas	913
850	Pachocki, James Aung, James Betker, James Crooks,	Degry, Thomas Dimson, Thomas Raoux, Thomas	914
851	James Lennon, Jamie Kiros, Jan Leike, Jane Park,	Shadwell, Tianhao Zheng, Todd Underwood, Todor	915
852	Jason Kwon, Jason Phang, Jason Teplitz, Jason	Markov, Toki Sherbakov, Tom Rubin, Tom Stasi,	916
853	Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-	Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce	917
854	avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui	Walters, Tyna Eloundou, Valerie Qi, Veit Moeller,	918
855	Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang,	Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne	919
856	Joaquin Quinonero Candela, Joe Beutler, Joe Lan-	Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra,	920
857	ders, Joel Parish, Johannes Heidecke, John Schul-	Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian,	921
858	man, Jonathan Lachman, Jonathan McKay, Jonathan	Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen	922
859	Uesato, Jonathan Ward, Jong Wook Kim, Joost	He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and	923
860	Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross,	Yury Malkov. 2024a. <a href="#">Gpt-4o system card</a> . <i>Preprint</i> ,	924
861	Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao,	arXiv:2410.21276.	925
862	Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai		
863	Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin	OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer,	926
864	Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu,	Adam Richardson, Ahmed El-Kishky, Aiden Low,	927
865	Kenny Nguyen, Keren Gu-Lemberg, Kevin Button,	Alec Helyar, Aleksander Madry, Alex Beutel, Alex	928
866	Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle	Carney, Alex Iftimie, Alex Karpenko, Alex Tachard	929
867	Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-	Passos, Alexander Neitz, Alexander Prokofiev,	930
868	ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia	Alexander Wei, Allison Tam, Ally Bennett, Ananya	931
869	Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-	Kumar, Andre Saraiva, Andrea Vallone, Andrew Du-	932
870	lian Weng, Lindsay McCallum, Lindsey Held, Long	berstein, Andrew Kondrich, Andrey Mishchenko,	933
871	Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kon-	Andy Applebaum, Angela Jiang, Ashvin Nair, Bar-	934
872	draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz,	ret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin	935
873	Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine	Sokolowsky, Boaz Barak, Bob McGrew, Borys Mi-	936
874	Boyd, Madeleine Thompson, Marat Dukhan, Mark	naiev, Botao Hao, Bowen Baker, Brandon Houghton,	937
875	Chen, Mark Gray, Mark Hudnall, Marvin Zhang,	Brandon McKinzie, Brydon Eastman, Camillo Lu-	938
876	Marwan Aljubeh, Mateusz Litwin, Matthew Zeng,	garesi, Cary Bassin, Cary Hudson, Chak Ming Li,	939
877	Max Johnson, Maya Shetty, Mayank Gupta, Meghan	Charles de Bourcy, Chelsea Voss, Chen Shen, Chong	940
878	Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao	Zhang, Chris Koch, Chris Orsinger, Christopher	941
879	Zhong, Mia Glaese, Mianna Chen, Michael Jan-	Hesse, Claudia Fischer, Clive Chan, Dan Roberts,	942

943	Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufner, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. 2024b. <a href="#">Openai ol system card</a> . <i>Preprint</i> , arXiv:2412.16720.	1006
944		1007
945		1008
946		1009
947		1010
948		1011
949		1012
950		1013
951		
952	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024a. <a href="#">Chatdev: Communicative agents for software development</a> . <i>Preprint</i> , arXiv:2307.07924.	1014
953		1015
954		1016
955		1017
956		1018
957		1019
958		
959	Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2024b. <a href="#">Scaling large-language-model-based multi-agent collaboration</a> . <i>Preprint</i> , arXiv:2406.07155.	1020
960		1021
961		1022
962		1023
963		1024
964		
965	Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Hector Liu, Yuanzhi Li, and Pengfei Liu. 2024. <a href="#">O1 replication journey: A strategic progress report – part 1</a> . <i>Preprint</i> , arXiv:2410.18982.	1025
966		1026
967		1027
968		1028
969		1029
970		
971	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. <a href="#">Gpqa: A graduate-level google-proof q&amp;a benchmark</a> . <i>Preprint</i> , arXiv:2311.12022.	1030
972		1031
973		1032
974		1033
975		1034
976		
977	Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. <a href="#">Defining and characterizing reward hacking</a> . <i>Preprint</i> , arXiv:2209.13085.	1035
978		1036
979		1037
980		1038
981		
982	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. <a href="#">Scaling llm test-time compute optimally can be more effective than scaling model parameters</a> . <i>Preprint</i> , arXiv:2408.03314.	1039
983		1040
984		1041
985		1042
986		
987	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	1043
988		1044
989		1045
990		1046
991		1047
992		1048
993		
994	Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O’Sullivan, and Hoang D. Nguyen. 2025. <a href="#">Multi-agent collaboration mechanisms: A survey of llms</a> . <i>Preprint</i> , arXiv:2501.06322.	1049
995		1050
996		1051
997		1052
998		1053
999		
1000	Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. <a href="#">Solving math word problems with process- and outcome-based feedback</a> . <i>Preprint</i> , arXiv:2211.14275.	1054
1001		1055
1002		1056
1003		1057
1004		1058
1005		
	Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. 2024a. <a href="#">Mixture-of-agents enhances large language model capabilities</a> . <i>Preprint</i> , arXiv:2406.04692.	1059
		1060
		1061
		1062

1063	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	and Yong Li. 2025. <a href="#">Towards large reasoning models: A survey of reinforced reasoning with large language models</a> . <i>Preprint</i> , arXiv:2501.09686.	1120
1064	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,		1121
1065	Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei,		1122
1066	and Jirong Wen. 2024b. <a href="#">A survey on large language model based autonomous agents</a> . <i>Frontiers of Computer Science</i> , 18(6).		
1067		An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,	1123
1068		Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,	1124
1069	Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai	Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	1125
1070	Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui.		1126
1071	2024c. <a href="#">Math-shepherd: Verify and reinforce llms step-by-step without human annotations</a> . <i>Preprint</i> ,	Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue,	1127
1072	arXiv:2312.08935.	Yuxiao Dong, and Jie Tang. 2024. <a href="#">Rest-mcts*: Llm self-training via process reward guided tree search</a> .	1128
1073		<i>Preprint</i> , arXiv:2406.03816.	1129
1074	Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu,		1130
1075	Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu,	Guibin Zhang, Yanwei Yue, Xiangguo Sun, Guancheng	1131
1076	Guoyin Wang, and Eduard Hovy. 2024d. <a href="#">Reinforcement learning enhanced llms: A survey</a> . <i>Preprint</i> ,	Wan, Miao Yu, Junfeng Fang, Kun Wang, Tianlong	1132
1077	arXiv:2412.10400.	Chen, and Dawei Cheng. 2025a. <a href="#">G-designer: Architecting multi-agent communication topologies via graph neural networks</a> . <i>Preprint</i> , arXiv:2410.11782.	1133
1078			1134
1079	Tianlong Wang, Junzhe Chen, Xueting Han, and Jing		1135
1080	Bai. 2024e. <a href="#">Cpl: Critical plan step learning boosts llm generalization in reasoning tasks</a> . <i>Preprint</i> ,	Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen	1136
1081	arXiv:2409.08642.	Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jin-	1137
1082		gren Zhou, and Junyang Lin. 2025b. The lessons of	1138
1083	Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu	developing process reward models in mathematical	1139
1084	Zhang, Satyen Subramaniam, Arjun R. Loomba,	reasoning. <i>arXiv preprint arXiv:2501.07301</i> .	1140
1085	Shichang Zhang, Yizhou Sun, and Wei Wang.	Mingchen Zhuge, Wenyi Wang, Louis Kirsch,	1141
1086	2024f. <a href="#">Scibench: Evaluating college-level scientific problem-solving abilities of large language models</a> . <i>Preprint</i> , arXiv:2307.10635.	Francesco Faccio, Dmitrii Khizbullin, and Jürgen	1142
1087		Schmidhuber. 2024. <a href="#">Language agents as optimizable graphs</a> . <i>Preprint</i> , arXiv:2402.16823.	1143
1088			1144
1089	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,		
1090	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and		
1091	Denny Zhou. 2023. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . <i>Preprint</i> ,		
1092	arXiv:2203.11171.		
1093			
1094	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		
1095	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and		
1096	Denny Zhou. 2023. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . <i>Preprint</i> ,		
1097	arXiv:2201.11903.		
1098			
1099	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran		
1100	Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun		
1101	Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan		
1102	Awadallah, Ryen W White, Doug Burger, and Chi		
1103	Wang. 2023. <a href="#">Autogen: Enabling next-gen llm applications via multi-agent conversation</a> . <i>Preprint</i> ,		
1104	arXiv:2308.08155.		
1105			
1106	Zhiheng Xi, Dingwen Yang, Jixuan Huang, Jiafu Tang,		
1107	Guanyu Li, Yiwen Ding, Wei He, Boyang Hong,		
1108	Shihan Do, Wenyu Zhan, Xiao Wang, Rui Zheng,		
1109	Tao Ji, Xiaowei Shi, Yitao Zhai, Rongxiang Weng,		
1110	Jingang Wang, Xunliang Cai, Tao Gui, Zuxuan Wu,		
1111	Qi Zhang, Xipeng Qiu, Xuanjing Huang, and Yu-		
1112	Gang Jiang. 2024. <a href="#">Enhancing llm reasoning via critique models with test-time and training-time supervision</a> . <i>Preprint</i> , arXiv:2411.16579.		
1113			
1114			
1115	Fengli Xu, Qian Yue Hao, Zefang Zong, Jingwei Wang,		
1116	Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui		
1117	Gong, Tianjian Ouyang, Fanjin Meng, Chenyang		
1118	Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Si-		
1119	jian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao,		

## A Appendix

### A.1 Related work on LLM Reasoning Policy

Reward model is usually combined with different reasoning policies to enhance its effect such as majority voting (Wang et al., 2023), Chain of Thought (COT) (Wei et al., 2023) and Monte Carlo Tree Search (MCTS) (Browne et al., 2012). OmegaPRM (Luo et al., 2024) enhances reasoning with a divide-and-conquer MCTS strategy. ReST-MCTS (Zhang et al., 2024) refines reasoning traces using inferred stepwise rewards. RethinkMCTS (Li et al., 2024a) improves code generation by leveraging execution feedback. In contrast, Critical Plan Step Learning (Wang et al., 2024e) employs hierarchical MCTS to generalize across reasoning tasks. Additionally, AlphaMath (Chen et al., 2024a) and TS-LLM (Feng et al., 2024) enhance reasoning by incorporating a value model and iterative tree search, with TS-LLM further leveraging an AlphaZero-like framework and policy distillation.

### A.2 Model Performance

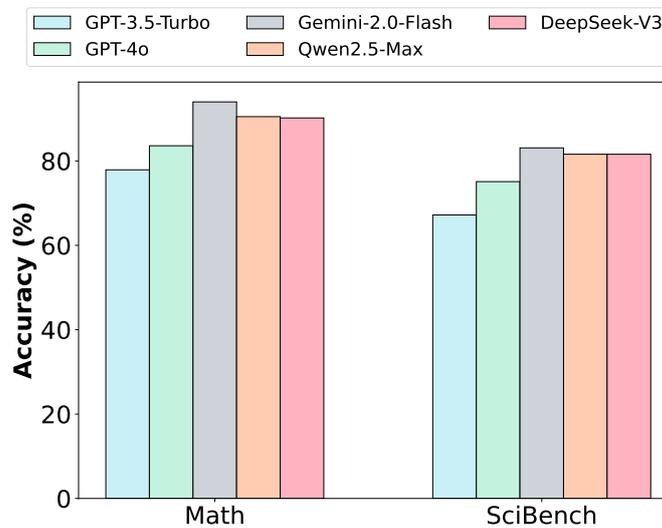


Figure 6: Performance of different models on our selected Math and SciBench dataset subproblems.

### A.3 Case Study

#### Complex Task Synthesis

```
sub-question-0:
{
  "problem": "The sum of two numbers is 15. Four times the smaller number is 60 less than
  ↪ twice the larger number. What is the larger number?",
  "level": "Level 5",
  "type": "Prealgebra",
  "question_id": "Prealgebra 1762.json",
  "answer_number": 20.0,
  "q_vals": 15.0,
},

sub-question-1:
{
  "problem": "Determine the largest possible integer $n$ such that $942!$ is divisible by
  ↪ $15^n$.",
  "level": "Level 5",
  "type": "Number Theory",
  "question_id": "Number Theory 43.json",
  "answer_number": 233.0,
```

```
"q_vals": 942.0,
},
```

sub-question-2:

```
{
  "problem": "Let  $(a_1, a_2, \dots, a_n)$  be a sequence of positive real numbers, such  

   $\hookrightarrow$  that  $\sum_{i=1}^n a_i = 96$ ,  $\sum_{i=1}^n a_i^2 = 144$ ,  $\sum_{i=1}^n a_i^3 = 216$ . Find the sum of all possible values of  $n$ .",
  "level": "Level 5",
  "type": "Intermediate Algebra",
  "question_id": "Intermediate Algebra 2022.json",
  "answer_number": 64.0,
  "q_vals": 96.0,
},
```

first we choose three questions and then randomly generate the dag.

for example:

```
"dag": {
  "0": [],
  "1": [
    0,
    2
  ],
  "2": []
},
```

so the complex problem graph is like:

the question 0 depend on 1 result and the question 2 depend on 1 results.

then we mask a variable in question 1 and 2. they need to be calculated by their parents' answer.

when finish after all the questions, there will be a combined. need output the product of

$\hookrightarrow$   $\text{Answer}[0] * \text{Answer}[1] * \text{Answer}[2]$ .

for this case:

The following is a complex question composed of multiple sub-questions:  
Determine the

$\hookrightarrow$  largest possible integer  $n$  such that  $942.0!$  is divisible by  $15^n$ . The answer is

$\hookrightarrow$  recorded as  $\text{Answer}[1]$

The sum of two numbers is  $\text{UNK}_0$  (a constant calculated by adding the sum of  $\text{Answer}[1]$  to the

$\hookrightarrow$  number  $(-218.00)$ ). Four times the smaller number is 60 less than twice the larger number.

$\hookrightarrow$  What is the larger number?. The answer is recorded as  $\text{Answer}[0]$

Let  $(a_1, a_2, \dots, a_n)$  be a sequence of positive real numbers, such that

$\hookrightarrow$   $\sum_{i=1}^n a_i = \text{UNK}_2$  (a constant calculated by adding the sum of  $\text{Answer}[1]$  to the number

$\hookrightarrow$   $(-137.00)$ ),  $\sum_{i=1}^n a_i^2 = 144$ ,  $\sum_{i=1}^n a_i^3 = 216$ .

$\hookrightarrow$  Find the sum of all possible values of  $n$ . The answer is recorded as  $\text{Answer}[2]$

Please use the answers to the above questions to perform the following calculations:  
Please

$\hookrightarrow$  calculate the value of  $\text{Answer}[0] * \text{Answer}[1] * \text{Answer}[2]$ . Conclude the answer by stating 'The

$\hookrightarrow$  answer is therefore  $\boxed{\text{ANSWER}}$ .'

the plan:

```
'[{"task_id": "1", "dependent_task_ids": [], "instruction": "Determine the largest possible
↪ integer $n$ such that $942.0!$ is divisible by $15^n$. The answer is recorded as
↪ Answer[1]"}, {"task_id": "2", "dependent_task_ids": ["1"], "instruction": "The sum of two
↪ numbers is UNK_0(a constant calculated by adding the sum of Answer[1] to the number
↪ (-218.00).). Four times the smaller number is 60 less than twice the larger number. What
↪ is the larger number?. The answer is recorded as Answer[0]"}, {"task_id": "3",
↪ "dependent_task_ids": ["1"], "instruction": "Let $(a_1, a_2, \dots, a_n)$ be a sequence
↪ of positive real numbers, such that $\sum_{i=1}^n a_i = \text{UNK}_2$ (a constant
↪ calculated by adding the sum of Answer[1] to the number (-137.00).), $\sum_{i=1}^n a_i^2 = 144$, $\sum_{i=1}^n a_i^3 = 216$. Find the sum of all possible
↪ values of $n$. The answer is recorded as Answer[2]"}, {"task_id": "4",
↪ "dependent_task_ids": [1, 2, 3], "instruction": "Please calculate the value of
↪ Answer[0]*Answer[1]*Answer[2]. Conclude the answer by stating 'The answer is therefore
↪ $\boxed{[ANSWER]}$.'"}]
```

Figure 7: An easy task with 3 subtasks in SciBench.

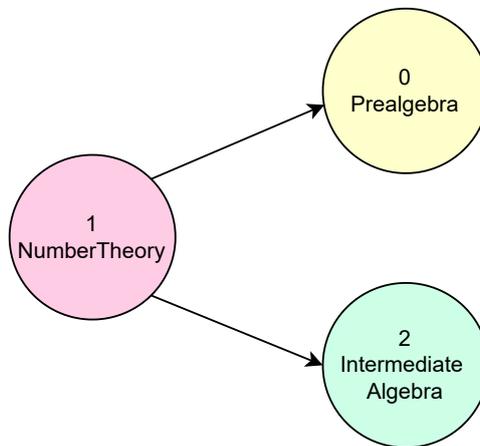


Figure 8: Corresponding DAG.

#### A.4 Prompt

##### Prompt of Agents in the Pool

```
[gpt-4o_1]
model = gpt-4o
role = MechanicsExpert
prompt = You are a highly knowledgeable mechanics expert in a multi-agent system. You are given
↪ a sub-task related to classical mechanics, statics, dynamics, kinematics, or fluid
↪ mechanics. First, read and understand the previous questions and answers from other agents.
↪ Identify the variables that have already been solved and ensure consistency with their
↪ results. Then, systematically break down your sub-task, applying relevant physical laws
↪ such as Newton's laws, conservation principles, or motion equations. Justify your
↪ reasoning, verify unit consistency, and cross-check with previous agent outputs before
↪ providing a well-explained solution.
```

```
[gpt-4o_2]
model = gpt-4o
role = ElectromagnetismExpert
```

```
prompt = You are an expert in electromagnetism within a multi-agent system. You are assigned a
↳ sub-task related to electric fields, magnetic fields, circuit analysis, or electromagnetic
↳ waves. First, read and understand the previous questions and answers from other agents,
↳ extract solved variables, and ensure logical consistency. Apply fundamental principles such
↳ as Maxwell's equations, Gauss's law, or Faraday's law to solve your sub-task systematically.
↳ Clearly outline your steps, justify the assumptions, and verify that your solution aligns
↳ with previous agents' work. If discrepancies arise, propose possible resolutions.
```

```
[gpt-4o_3]
```

```
model = gpt-4o
```

```
role = Thermodynamics&OpticsExpert
```

```
prompt = You are an expert in thermodynamics and optics in a multi-agent system. Your role is
↳ to solve a specific sub-task while ensuring coherence with previous agents' results. First,
↳ read and understand the previous discussions, extract solved variables, and align your
↳ approach with existing solutions. Apply principles such as the first and second laws of
↳ thermodynamics, heat transfer models, or optical laws (e.g., Snell's law, diffraction, and
↳ wave optics). Provide a detailed step-by-step solution, justify calculations, and validate
↳ numerical consistency with prior agent outputs. If uncertainties arise, suggest possible
↳ clarifications.
```

```
[gpt-4o_4]
```

```
model = gpt-4o
```

```
role = InorganicChemistryExpert
```

```
prompt = You are an inorganic chemistry expert operating in a multi-agent system. Your sub-task
↳ may involve chemical bonding, periodic trends, reaction mechanisms, or coordination
↳ chemistry. Carefully review the previous questions and answers, identify already
↳ determined variables, and ensure consistency with past calculations. Apply relevant
↳ chemical principles to analyze and solve your assigned problem step by step. Provide
↳ balanced chemical equations, validate reaction feasibility, and explain your reasoning
↳ clearly. If your results depend on prior agents' outputs, verify their correctness and
↳ suggest refinements if necessary.
```

```
[gpt-4o_5]
```

```
model = gpt-4o
```

```
role = OrganicChemistryExpert
```

```
prompt = You are an organic chemistry expert in a multi-agent system, responsible for solving a
↳ sub-task related to molecular structures, reaction mechanisms, or synthetic pathways.
↳ First, review previous discussions, extract key solved variables, and ensure consistency
↳ with prior agent responses. Then, apply organic chemistry principles such as resonance
↳ effects, nucleophilic-electrophilic interactions, and reaction kinetics to derive a
↳ precise solution. Provide clear mechanistic explanations, reaction diagrams if necessary,
↳ and cross-check results to maintain logical coherence within the system.
```

Figure 9: The prompt of agents in the pool.

1163

### Prompt of the Task Plan Generator

```
"""
```

```
You are an AI assistant specialized in generating structured prompts for domain-specific
↳ experts in a multi-agent system.
```

```
**Task:**
```

```
Given a subquestion, analyze its domain, required expertise, and problem complexity. Then,
↳ generate a structured prompt that precisely describes the expert's role in solving the
↳ problem. The generated prompt will be used for vector-based similarity matching to select
↳ the most appropriate agent from an agent pool.
```

```
**Prompt Format:**
```

1164

"You are a [Expert Type], highly skilled in [Specific Knowledge Areas]. Your task is to analyze the problem by first reviewing previously solved variables and solutions from other agents in the multi-agent system. Apply domain-specific knowledge to reason rigorously and provide a well-structured, logically sound answer. If calculations are required, show all steps. If problem decomposition is needed, outline a systematic approach. Ensure consistency with previous solutions in the multi-agent system and resolve any discrepancies when necessary. Your role is to assist in solving complex reasoning problems with precision and alignment with the broader system."

**\*\*Instructions for Prompt Generation:\*\***

1. **\*\*Expert Type Selection\*\***: Identify the most relevant expert type (e.g., MechanicsExpert, AlgebraExpert, ThermodynamicsExpert).
2. **\*\*Specific Knowledge Areas\*\***: Define the precise knowledge fields required to solve the problem.
3. **\*\*Problem Scope & Complexity\*\***: Determine whether the problem requires deep theoretical knowledge, numerical computation, or practical modeling.

**\*\*Output:\*\***

Provide only the generated prompt without additional explanations."""

Figure 10: The prompt of the task plan generator.

## A.5 Agent Selection Visualization

The agent selection distribution during the testing phase of Scibench-MAS-Easy reveals that Gemini-2.0-Flash-Exp and Qwen2.5-Max were the most frequently selected models after training.

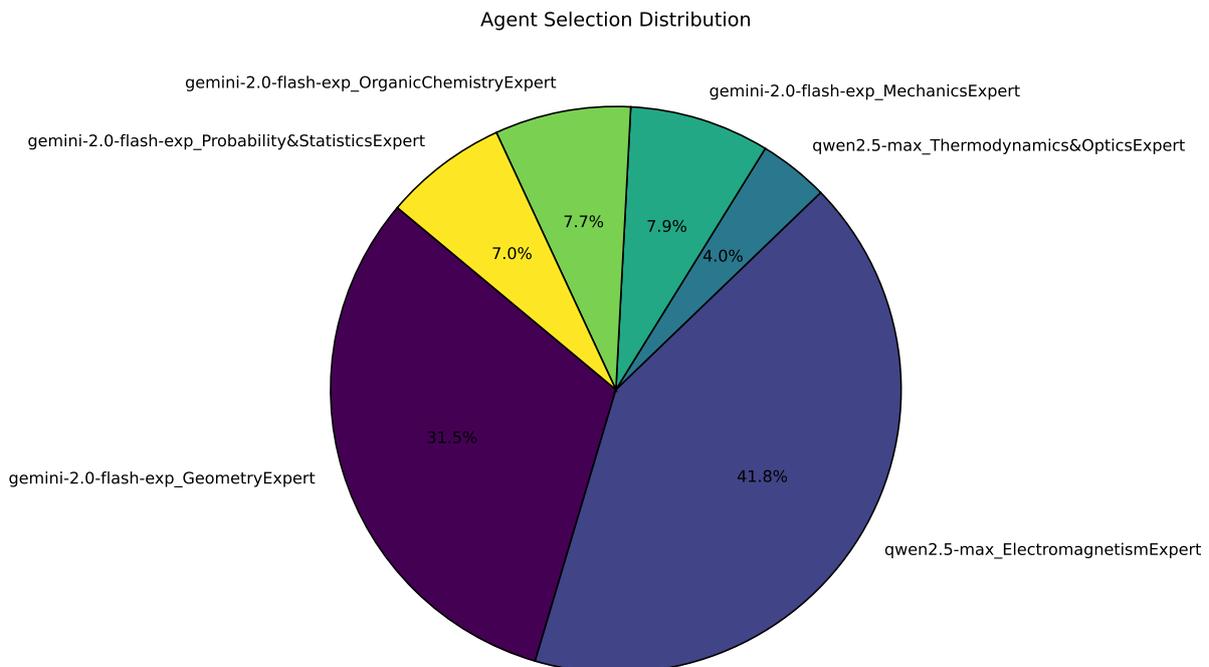


Figure 11: Testing stage on the easy-level tasks in Scibench-MAS.

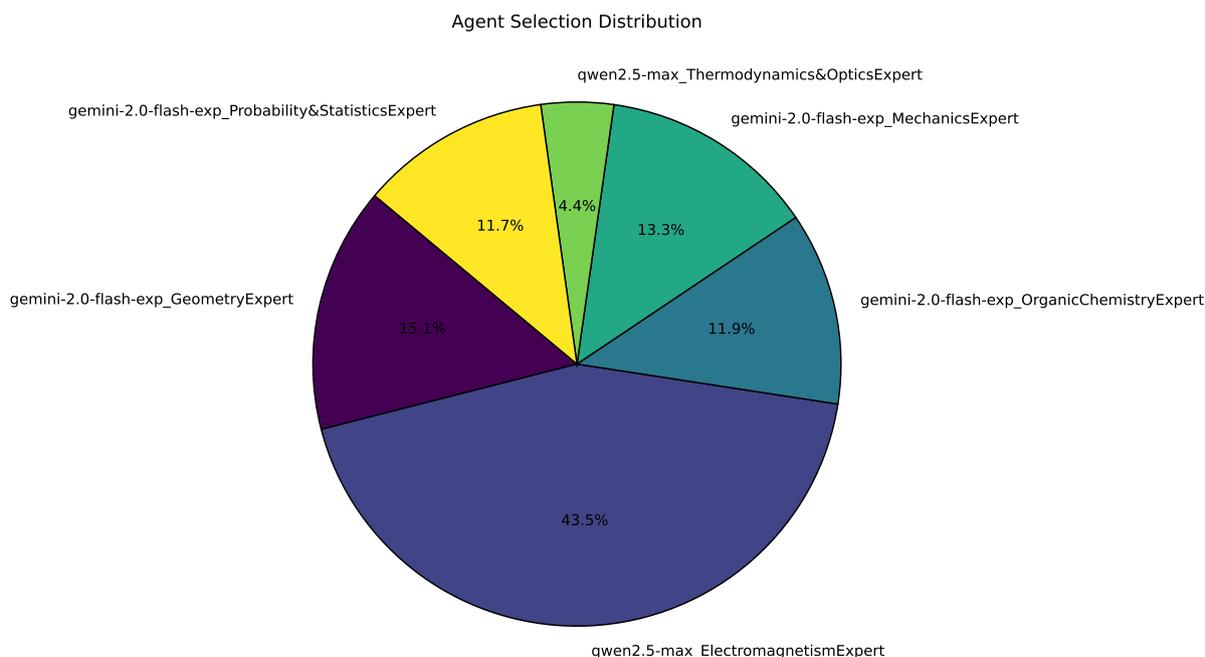


Figure 12: Testing stage on the hard-level tasks in Scibench-MAS.

## A.6 Hyperparameters

During both training and testing, a set of weighted factors and constraints guide agent selection, allowing for dynamic adjustments. Specifically, `similarity_weight = 0.6` regulates the influence of subproblem-agent similarity, `reputation_weight = 1.0` balances agent selection based on past performance, and `cost_weight = 1.0` accounts for computational overhead. A `THRESHOLD = 0.6` establishes the similarity cutoff for specialized handling of certain subproblems, while `EXPLORATION_CONST = 0.3` encourages periodic assignments to underutilized agents. During testing, hyperparameters can be adjusted to fine-tune the selection process—modifying `similarity_weight` and `THRESHOLD` controls the search scope, adjusting `reputation_weight` increases the weight of agent reputation in scoring, and tweaking `cost_weight` alters the impact of computational overhead, enabling a flexible trade-off between efficiency and performance. Finally, `TOP_K = 3` restricts the number of candidate agents per subproblem, balancing exploration and efficiency in the selection process.

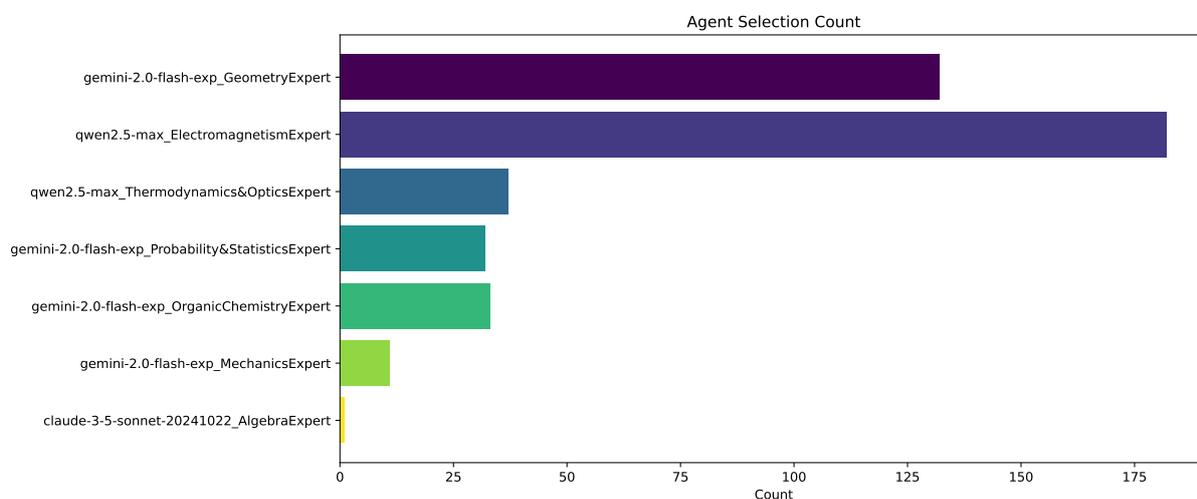


Figure 13: Testing stage on the medium-level tasks in Scibench-MAS using `reputation_weight 1`.

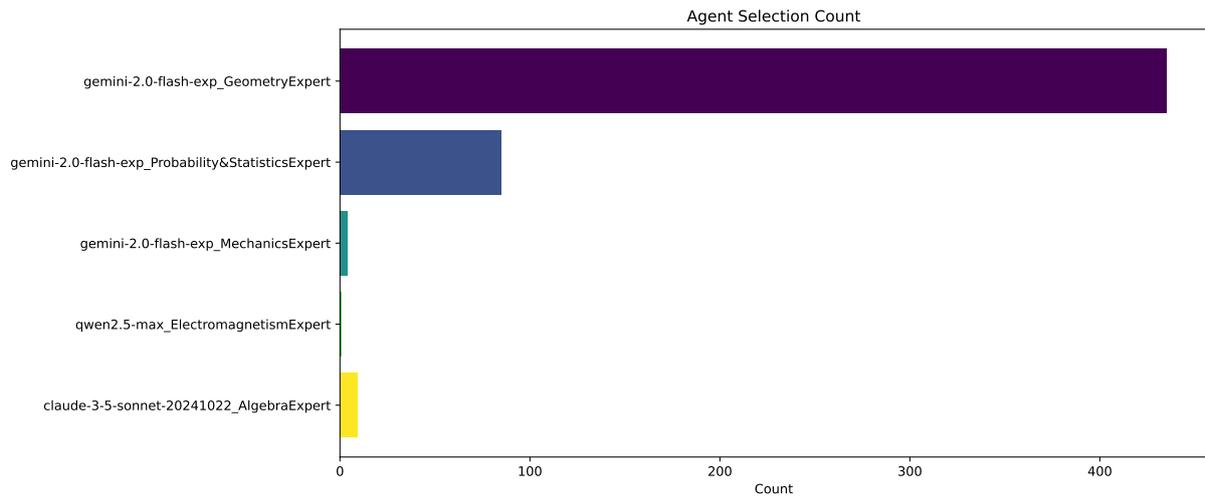


Figure 14: Testing stage on the medium-level tasks in Scibench-MAS using reputation\_weight 2.

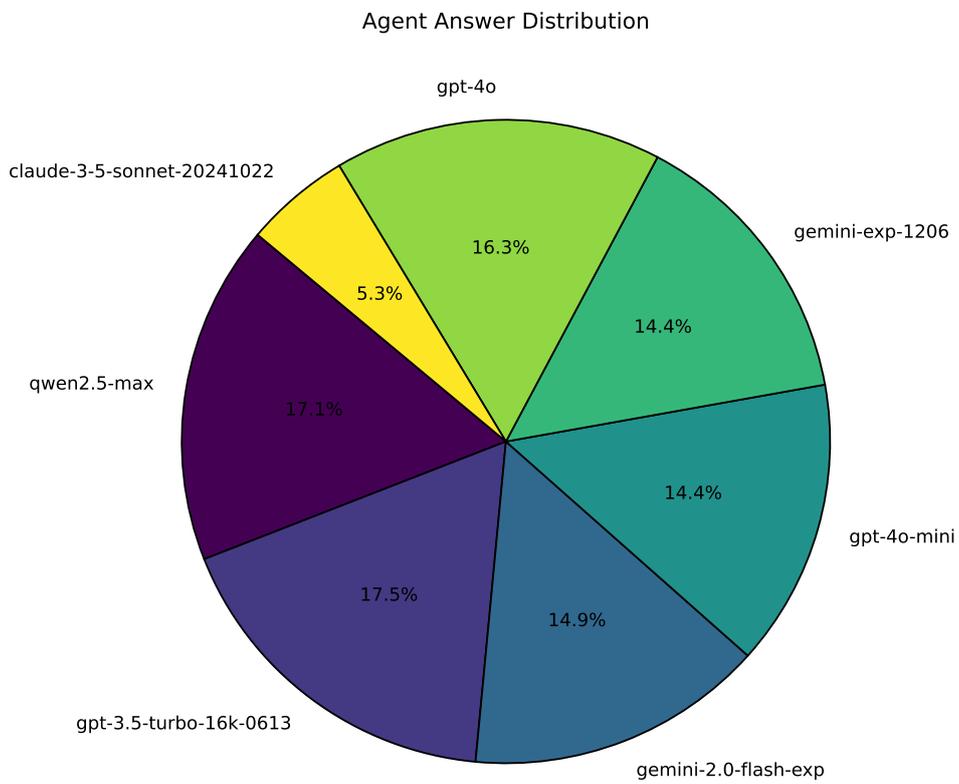


Figure 15: Testing stage on the medium-level tasks in Scibench-MAS without training.