

Evaluate, Scale, and Credit: A Comprehensive Study on Multi-Agent Collaboration of Large Language Models

Anonymous ACL submission

Abstract

Large Language Models based Multi-Agent Systems (LLM-MAS) perform well in many domains, but we still lack a clear understanding of the collaboration mechanism among multiple LLM-based agents. This study aims to explore three key issues: (1) Can multi-agent outperform single-agent systems? (2) Is scaling better for multi-agent systems? (3) How to credit agents and find potential effective structures? Specifically, we design five collaboration architectures and evaluate their effectiveness across different LLMs and tasks. Our findings offer significant insights for understanding the collaboration within MAS, building collaboration architectures among agents, and reducing system costs. Furthermore, our conclusion will inspire and provide new perspectives for future studies on LLM-MAS.

1 Introduction

Large Language Model-based Multi-agent Systems (LLM-MAS) specialize multiple LLMs into different agents and effectively simulate complex real-world environments through the interaction among these diverse agents (Guo et al., 2024). With proven outstanding abilities in contextual understanding, reasoning, and generation, LLMs empower agents to collaboratively plan, discuss, and make decisions, imitating human team cooperation to solve real world problems (Li et al., 2023; Hong et al., 2023; Wu et al., 2023).

Recent research efforts have focused on exploring and optimizing the collaboration mechanisms of MAS driven by LLMs (Liang et al., 2023; Du et al., 2023; Chan et al., 2023), revealing two critical challenges: *architecture scaling* and *contribution crediting*. The challenge of *architecture scaling* encompasses expanding the number of agents and increasing their interaction frequency to solve more complex tasks (Zhang et al., 2023b; Chan et al., 2023; Li et al., 2024). However, while enhancing system capabilities, scaling also leads to

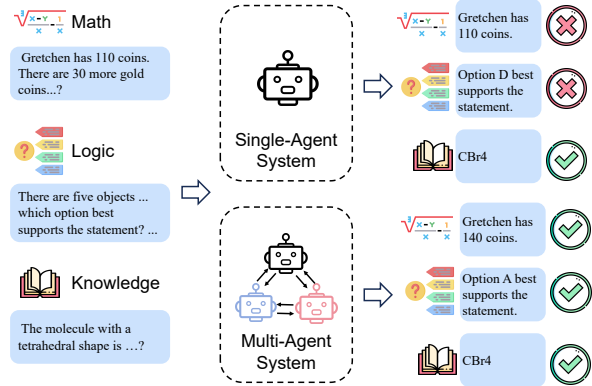


Figure 1: Illustrations of single-agent system and multi-agent system.

a substantial rise in communication overhead, presenting a notable challenge in maintaining system efficiency (Zhang et al., 2023b; Yin et al., 2023). At the same time, the challenge of *contribution crediting* involves the accurate allocation of contributions among agents, which is crucial for promoting collaboration and ensuring interpretability and robustness within LLM-MAS systems (Liu et al., 2023b). Evaluating LLM-MAS systems from the perspectives of scaling and crediting not only diagnoses their current shortcomings and limitations but also directs future developments toward more efficient, effective, and scalable multi-agent collaborations.

In this paper, to comprehensively evaluate the multi-agent collaboration of large language models, we designed a unified evaluation procedure and conducted systematic evaluations on 9 datasets across 3 tasks. Specifically, we design five collaboration architectures that reflect different communication patterns and the diversity of agent collaboration. This paper primarily investigates three research questions (RQ):

RQ1: Can multi-agent outperform single-agent systems? Different from single-agent systems, LLM-MAS involves multiple agents that influence each other with frequent and complex agent

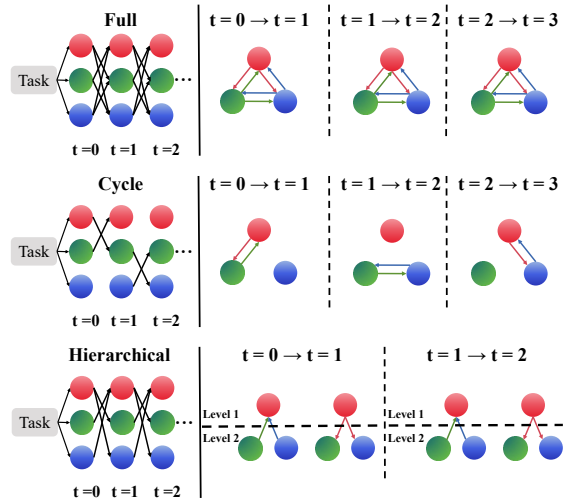
068	interactions. The interaction or communication pat-	systems, and single-agent performance does not	119
069	terns between agents, which we refer to as the col-	determine multi-agent benefit. 2) More agents	120
070	laboration architectures, can significantly affect the	will bring more benefits, and achieving agreement	121
071	system performance. Some researchers have ex-	among agents is crucial for better performance. 3)	122
072	explored several optimal collaboration architectures	Pruning communication paths based on crediting is	123
073	(Chan et al., 2023; Chen et al., 2023b) and designed	possible and reveals some aggregating information	124
074	various LLM-based multi-agent systems. Yin et al.	and self-reflection structures among agents in the	125
075	(2023) explored integrating different collaboration	collaboration strategy	126
076	architectures to enhance system performance. How-	Generally, our contributions are as follows:	127
077	ever, past studies primarily focused on exploring		
078	specific systems, lacking a comprehensive study on	• We introduce five collaborative multi-agent ar-	128
079	the general properties of LLM-MAS. Inspired by	chitectures and conduct extensive experiments	129
080	traditional multi-agent theory, we construct several	in various scenarios to explore three crucial	130
081	collaboration architectures and use these architec-	questions regarding the multi-agent collabora-	131
082	tures to build multiple multi-agent systems and	tion of large language models.	132
083	conduct systematic studies in different scenarios.		
084	RQ2: Is scaling better for multi-agent sys-	• We investigate the connection between the	133
085	tems? Cost is a crucial but often overlooked limit-	scale and performance of LLM-MAS and	134
086	ing factor in LLM-based multi-agent research. In	provide an in-depth study of the agreement	135
087	this study, we analyze the scale of MAS, including	changes of the system and the early stopping	136
088	time step, agent number, and the threshold of early	mechanism.	137
089	stopping, etc. Yin et al. (2023) discussed the the-		
090	oretical costs of some collaboration architectures.	• By quantifying the credits of communication	138
091	Li et al. (2024) systematically studied the effect of	paths, we propose a Shapley value-based prun-	139
092	the agent number in a sampling-and-voting method.	ing approach for LLM-MAS. This pruned	140
093	However, they did not consider the communication	structure significantly reduces communication	141
094	between agents. A question remains: How do we	costs across various datasets while achieving	142
095	decide the scale or cost-related factors? Therefore,	superior performance and implies some poten-	143
096	we systematically analyze the relationship between	tially effective structures.	144
097	scale and performance in multi-agent.		
098	RQ3: How to credit agents and find potential	2 Collaboration Architectures	145
099	effective structures We also address the credit-	Traditional multi-agent research (Esmaeili et al.,	146
100	ing of communication paths in LLMs-based MAS,	2016; Damba and Watanabe, 2007; Dorri et al.,	147
101	which has received less attention than agent role	2018; Horling and Lesser, 2004) has identified and	148
102	assignments. Current strategies primarily utilize	delineated various effective multi-agent architec-	149
103	LLMs for evaluating agent outputs through rank-	tures, including Flat, Hierarchical, Holonic, and	150
104	ing or rating (Liu et al., 2023b; Jiang et al., 2023b;	Team. Each architecture possesses distinct advan-	151
105	Qin et al., 2023). This type of method, despite its	tages and is suitable for specific scenarios.	152
106	prevalence, faces challenges in accuracy. In con-	Inspired by the multi-agent theory and recent	153
107	trast, traditional multi-agent reinforcement learning	multi-agent research, we designed five unique col-	154
108	(MARL) offers insights into collaboration through	laboration architectures that reflect different com-	155
109	credit assignment, focusing on the distribution of	munication patterns and the diversity of agent col-	156
110	rewards among agents based on their contributions.	laboration. Figure 2 contains five types of collab-	157
111	Inspired by MARL principles (Minsky, 1961; Sune-	oration architectures. There are three static col-	158
112	hag et al., 2018), we explore an LLM-independent	laboration architectures: FULL, CYCLE, and HI-	159
113	method using <i>Shapley value</i> to quantify the contri-	ERARCHICAL in Figure 2(a), and two dynamic	160
114	butions of each collaboration in the system proce-	architectures, TEAM and RANK, in Figure 2(b).	161
115	dures and pruning multi-agent systems.		
116	Our experiments provide insights into the multi-	• FULL Inspired by the Flat structure (Dorri	162
117	agent collaboration of large language models: 1)	et al., 2018) of traditional multi-agent theory,	163
118	Multi-agent systems often outperform single-agent	information can be freely passed from one	164
		agent to another. In particular, when there are	165

only two agents, the collaboration architecture degenerates into a typical debate architecture. This kind of architecture simulates information propagation in unrestricted discussions, facilitating the fast spread of information. However, such networks may lead to high costs.

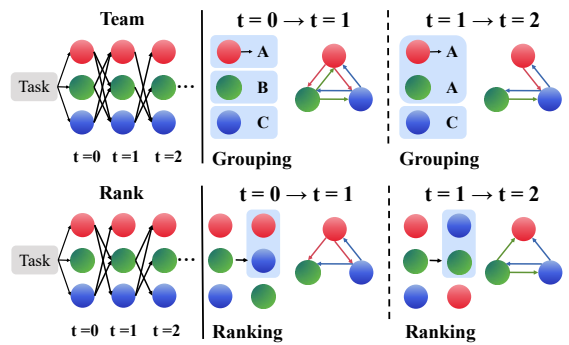
- **CYCLE** Inspired by Multi-Agent Debate (Liang et al., 2023), information is propagated among pairs of agents to reach a final agreement. This architecture simulates private conversation. It emphasizes how information gradually evolves and spreads over a limited number of interactions. This type of architecture has less costs, but the time required for the system to reach an agreement may be longer.
- **HIERARCHICAL** Inspired by the Hierarchical structure (Damba and Watanabe, 2007) of the traditional multi-agent theory, information is propagated between nodes at different levels. This architecture simulates the Delphi method¹ in expert groups. This kind of architecture emphasizes aggregation and processing of the information.

- **TEAM** Inspired by the Team structure (Parker, 1993) of the traditional multi-agent, information flows between agents with different viewpoints (answers). This architecture simulates the propagation of information during a team discussion. This kind of architecture has no interaction between agents with the same viewpoint.
- **RANK** Inspired by the idea of agent optimization in DyLAN (Liu et al., 2023b), information and messages are sorted before it is delivered, and only top-k information can be passed to the next time step. This architecture simulates a review or screening process, such as editorial review or administrator approval, emphasizing the concern for information quality.

¹Delphi method: soliciting experts' opinions on a problem, organizing and summarizing them, then anonymously feeding them back to the experts, and soliciting opinions again until they reach an agreement



(a) The static Collaboration architectures



(b) The dynamic collaboration architectures

Figure 2: Collaboration Architectures

3 Experiments

This section introduces the dataset and LLM we used, providing a data foundation for subsequent problem analysis.

Tasks and Datasets. In our experiments, we used a general evaluation procedure to assess the performance of five architectures across three tasks, including: 1) Math: GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021b), and SVAMP (Patel et al., 2021) datasets; 2) Knowledge: MMLU (Hendrycks et al., 2021a), CommonsenseQA (Talmor et al., 2019), and CommonsenseQA 2.0 (Talmor et al., 2022); 3) Logic: LogiQA (Liu et al., 2020), LogiQA2.0 (Liu et al., 2023a), and ReClor (Yu et al., 2020).

Model Details. We tested the proposed collaboration mechanism based on different models. Considering cost and effectiveness, we selected open-source models, e.g., Llama2-7b-Chat (Touvron et al., 2023), Mistral-7b-Instruction (Jiang et al., 2023a), and Starling-LM-7B-alpha (Zhu

et al., 2023), for our experiments. Specifically, we downloaded the corresponding open-source models on hugging face and deployed the APIs using Fastchat and vLLM (Kwon et al., 2023). These three LLMs will be combined with the five architectures to form 15 multi-agent systems. The maximum time step is six if not explicitly stated. Additionally, we conducted experiments on the quantized model of the larger model Mixtral-8x7B model², which can be found in Appendix A.4.

System Details. To reflect the difference between the Agents, we set the temperature of each Agent to a different value between 0 and 1 during generation. By default, we used 3 Agents with temperatures of 1, 0.6, and 0.4. Inspired by social comparison theory and review collaboration (Xu et al., 2023c), we considered generating solutions, final answers, and reviewing other agents’ answers during generation. Complete prompt examples can be found in the appendix. Motivated by Liu et al. (2023b) and Practical Byzantine Fault Tolerance, when 2/3 of the agents in the system reached a consensus (i.e., the answer is the same), we made the system early stop, and the process stopped.

4 Can multi-agent systems outperform single-agent systems?

This section evaluates the multi-agent benefit. We conducted experiments with multi-agent systems composed of three LLMs and five collaboration architectures across nine datasets and analyzed the MAS performance according to the relative improvement of multi-agent systems. Moreover, we investigated the impact of the possible factors of multi-agent synergy, i.e., collaboration architecture, LLM, and task.

The benefit of MAS

Final Success Rate (e.g., accuracy) is the most commonly used metric for evaluating multi-agent systems (Du et al., 2023; Chan et al., 2023; Liu et al., 2023b; Chen et al., 2023a), which offers the advantages of simplicity and intuitiveness. However, the final success rate is highly correlated with the LLM and Task and does not reflect multi-agent synergy. To examine the benefits of multi-agent synergy, a natural idea is to consider the relative improvement in accuracy, which we refer to as accuracy improve-

²<https://huggingface.co/TheBloke/Mixtral-8x7B-Instruct-v0.1-GPTQ>

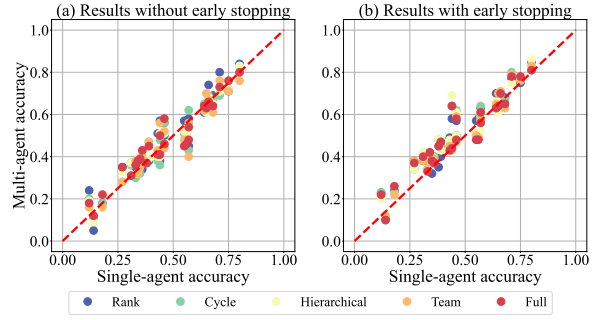


Figure 3: The scatter plot comparing multi-agent to single-agent performance.

ment (Δ_{acc}).

$$\Delta_{acc} = \frac{\text{Perf}_m - \text{Perf}_s}{\text{Perf}_s} \quad (1)$$

where Perf_s and Perf_m represent the performance (accuracy) of the single-agent³ and system, respectively.

Finding 1: *Multi-agent collaboration often help, and early stopping is necessary.*

We plotted single-agent and multi-agent accuracy for all possible <Architecture, LLM, Task> triplets, totaling 135 points, as shown in Figure 3. A point above the red line indicates that the multi-agent system outperforms the single-agent. As the chart shows, 55.6% MAS showed improvement compared to single-agent. With early stopping activated, this number increased to 80%. This finding suggests that multi-agent approaches generally offer improvements, and early stopping mechanisms are crucial for maximizing system performance. Detailed data are given in the Appendix A.2.

Finding 2: *Every factor related to multi-agent synergy influences the system significantly, and single-agent performance does not determine multi-agent benefit.*

In this part, we investigated the effect of the three factors: architecture, LLM, and task. To study the effect of architecture, we formed a vector of performances for all five architectures in every possible <LLM, Task>. We averaged these vectors to indicate the relative performance of architectures. To minimize the influence of LLM and task, we performed z-score normalization or Min-max normalization on all vectors before averaging. Let $\text{Perf}(a, m, t)$ be

³We use the results generated by greedy decoding to represent single-agent accuracy.

Metric		Single-Agent Accuracy		Multi-Agent Accuracy		Δ_{acc}	
		Min-Max	Z-score	Min-Max	Z-score	Min-Max	Z-score
Normalization							
Architecture	Full			0.41	-0.21	0.41	-0.21
	Cycle			0.51	0.10	0.51	0.10
	Hierarchical	\	\	0.52	0.04	0.52	0.04
	Rank			0.55	0.18	0.55	0.18
	Team			0.46	-0.11	0.46	-0.11
LLM	Llama2	0.12	-0.79	0.10	-0.97	0.60	0.30
	Mistral	0.31	-0.32	0.51	0.03	0.66	0.45
	Starling	0.93	1.11	0.90	0.94	0.15	-0.75
Task	Math	0.56	0.08	0.46	-0.10	0.63	0.31
	Knowledge	0.77	0.57	0.84	0.76	0.46	-0.10
	Logic	0.23	-0.65	0.22	-0.66	0.41	-0.20

Table 1: The analysis for the possible factors of multi-agent synergy. It is important to note that these values are **not the actual accuracy** of systems. They are the average values after Normalization across the different architectures (or LLMs, tasks).

the performance of a multi-agent system composed of architecture and LLM on task, $\text{Perf} = (\text{Perf}(\text{FULL}), \text{Perf}(\text{TEAM}), \dots, \text{Perf}(\text{CYCLE}))$

$$\bar{\text{Perf}} = \frac{\sum_{(m,t) \in \mathcal{M} \times \mathcal{T}} \text{Norm}(\tilde{\text{Perf}}(m,t))}{|\mathcal{M}| |\mathcal{T}|} \quad (2)$$

The experiment results in Table 1 show that (1) different architectures led to different improvements, and the Rank architecture achieved relatively the best results; (2) different LLMs led to different improvements, and Mistral achieved relatively the best results; (3) the effectiveness of the multi-agent approach also depended on the task. Math got the highest multi-agent benefits, which aligned with our expectations.

Notably, it is challenging to predict multi-agent benefits based on single-agent performance. For example, although Starling performed best with the single agent, its multi-agent benefits were less than Mistral. Knowledge tasks generally had the highest accuracy, but the multi-agent method improvement was less than Math. Besides, we plotted a scatter plot of single-agent performance and system improvement in the Appendix A.1, as shown in Figure 7, revealing no apparent correlation between Perf_s and Δ_{acc} .

5 Is scaling better for multi-agent systems?

This section examines and analyzes the relationship between scale and performance in MAS. In particular, we considered the agent number and maximum

communication rounds (time step) in MAS. Departing from Li et al. (2024), we focused on the scale of MAS with dynamic interactions among agents rather than the simple ensemble of answers.

Finding 3: *Many hands may make light work. More agents will bring more benefits.*

In this part, we explored the impact of different agent numbers in MAS. Due to max context length and expensive cost, we did not compare systems with more than five agents.

Table 2 shows the accuracy of different systems on different datasets. We observed an overall improvement in LLM-based agents, consistent with the findings of Li et al. (2024), which suggest that adding more agents can lead to better system performance. Although the performance did not continue to increase with five agents on the LOGIQA2 dataset, we believe that adding more and varied agents will improve its performance. It is worth pointing out that the only difference among the agents here is temperatures. Theoretically, adding agents with different roles or different LLMs will better improve performance (Chan et al., 2023; Liu et al., 2023b).

Finding 4: *Agreement is strength. Achieving agreement among agents is crucial for better performance.*

We calculated the system agreement at each time step and the proportion of the correct answer in each time step and shown the result in Figure 4. Generally, the higher system agreement could lead to better system performance. This observation

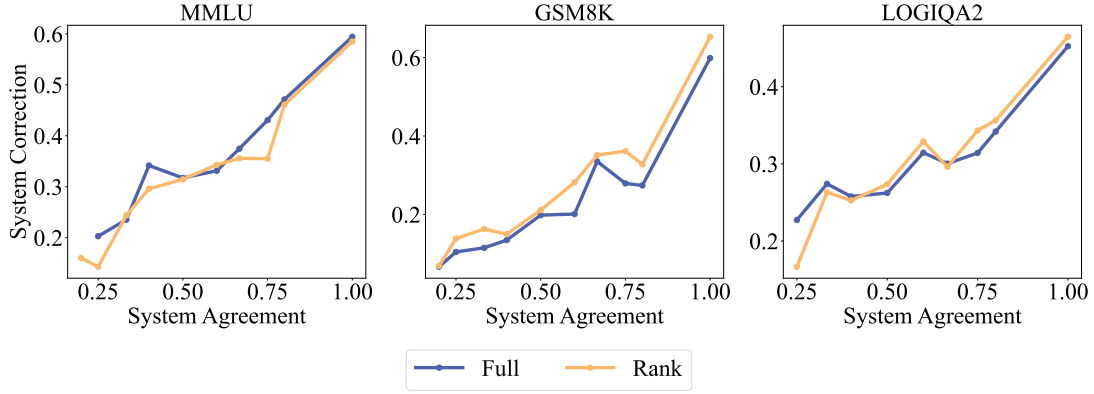


Figure 4: The agreement and system performance. This graph includes the results of 20 rounds of interaction on a dataset of 100 data points for systems ranging from 1 to 5 agents. The x-axis represents the system’s agreement, specifically, the proportion of agents that reached a consensus (calculated by the proportion of the most voted answer to the total number of answers).

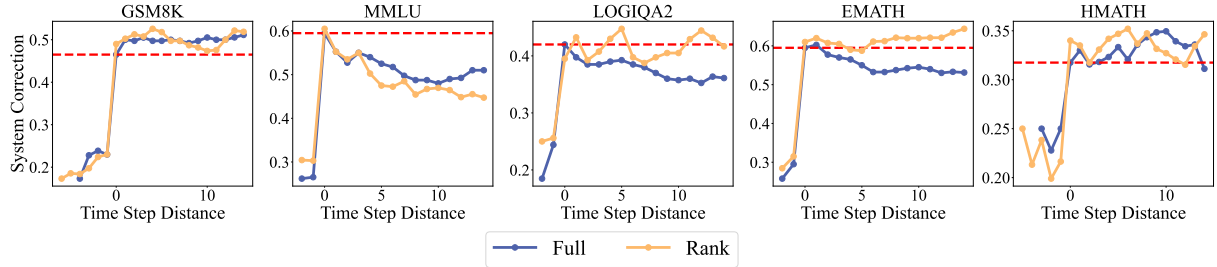


Figure 5: The agreement and system performance. The x-axis represents the distance from the current time step to the early stopping time step. For example, 0 represents the early stopping time step, 1 represents the next time step after early stopping, and -1 represents the time step before early stopping. These data come from a MAS composed of four agents based on Mistral.

Dataset	MMLU		GSM8K		LOGIQA2	
Architecture	Full	Rank	Full	Rank	Full	Rank
1 agent	44.0%	44.0%	46.0%	46.0%	39.0%	39.0%
2 agents	55.0%	55.0%	45.0%	41.0%	39.0%	44.0%
3 agents	64.0%	57.0%	47.0%	50.0%	47.0%	40.0%
4 agents	67.0%	67.0%	51.0%	57.0%	44.0%	42.0%
5 agents	67.0%	68.0%	52.0%	51.0%	45.0%	46.0%

Table 2: The performance of systems with different agents. Every system here is conducted with Mistral and applied early stopping.

may indicate that the benefit of MAS comes from the procedure in which agents collaborate and ultimately reach a consensus. Additionally, we found that different datasets had different performance-increasing speeds. Therefore, we wondered if the agreement threshold for early stopping is unique for different datasets.

Considering that 95% of the data reached early stopping within ten time steps, we examined the ten time steps before and after reaching early stopping. As shown in Figure 5, we found that both

MMLU and LogiQA2 reached their best performance at the early stopping time step. At the same time, GSM8K could further improve performance after early stopping, suggesting that using 2/3 as the early stopping threshold for GSM8K may not be reasonable. To determine the source of this observation, we additionally tested 100 sampled data of High school Mathematics Problems and Elementary Mathematics Problems in MMLU (named EMATH and HMATH), and the results revealed that EMATH showed a relatively small decrease with FULL and fluctuating correction with RANK, while HMATH showed a fluctuating increase in both architectures. We speculated the threshold might related to the task and its complexity. Math problems had a higher threshold, and the more challenging the tasks were, the higher the threshold was.

6 Crediting and Pruning

In this section, we credit the relative importance of communication paths to prune less significant

380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400

Architecture	Math			Knowledge			Logic			Avg \uparrow	$C_{rel}\downarrow$
	MATH	SVAMP	GSM8K	CSQA	CSQA2	MMLU	LogiQA	LogiQA2	ReClor		
OPTIMIZED (<i>Ours</i>)	0.24	0.72	0.56	0.60	0.67	0.66	0.37	0.47	0.41	0.52	0.46
FULL	0.22	0.70	0.47	0.58	0.64	0.64	0.38	0.47	0.43	0.50	1.00
RANK	0.19	0.71	0.49	0.57	0.66	0.65	0.42	0.45	0.50	0.52	0.66

Table 3: The performance of systems conducted with FULL, RANK, and OPTIMIZED architecture on different datasets. These systems were based on Mistral and built with 3 agents. C_{rel} indicates the relative number of communication paths, assume the path number of Full architecture to be 1.

ones and investigate the potential effective structures within the MAS collaboration. Specifically, we sampled an additional 200 data points from GSM8K to prune the FULL architecture, which is chosen because other architectures are essentially subsets of it, and pruning it is most likely to reveal potentially effective structures.

Credit assignment in MAS

Recent LLM-MAS use LLMs to rank or rate the information output of agents, calculating contributions based on these rankings or scores. While this type of approach has achieved certain results in many related studies (Chan et al., 2023; Zhang et al., 2023b; Jiang et al., 2023b), ranking or rating text by LLMs remains an unsolved problem (Wang et al., 2023a; Shen et al., 2023). Inspired by Credit Assignment in MARL, we broke down the pruning of the collaboration architecture into identifying the relative importance and reward of each communication path between agents at any single time step.

We use the Shapley value (Shapley and Corporation, 1951) to indicate the relative importance. The Shapley value is a concept from cooperative game theory that offers a fair distribution of the total gains to the players (agents) based on their contributions to the alliance (MAS).

Suppose the set of communication paths to Agent n at time step t is S . We defined the value function $v(S)$ as the accuracy difference of Agent n ⁴ between time step $t-1$ and t . Given N paths, the formula for the Shapley value of path i is:

$$\sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (v(S \cup \{i\}) - v(S)) \quad (3)$$

$v(S \cup \{i\})$ is the value of the alliance contains path i and $v(S)$ is the value of the alliance without path i . A higher Shapley value suggests a more significant importance or contribution of this path.

⁴we calculated the accuracy in the picked 200 training data

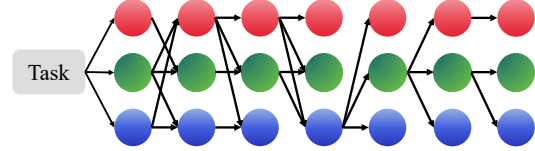


Figure 6: The pruned architecture. Contains only 46% communication paths in the Full architecture.

After calculating the Shapley values of the communication paths, we removed those paths where Shapley values were lower than a certain threshold (we took the threshold as 0.002 to eliminate those paths with a small effect). This ensures an overall improvement of each time step. After pruning a time step, we used the pruned structure to optimize the next step, continuing this process until there was no positive reward path or the maximum time step was reached.

Finding 6: Information aggregation and Self-reflection.

An interesting phenomenon occurs in pruned architecture: Information Aggregation and Self-Reflection. Information aggregates to specific agents and then spreads back to all agents, consistent with hierarchical architecture. Furthermore, we found that in the pruned architecture, agents tend to communicate with others at early time steps and tend to make a self-reflection, which aligns with the method mentioned in Wang et al. (2023b), at the later time steps. This may reduce the propagation of misinformation after multiple rounds of interaction.

Finding 5: The pruning architecture reduces cost and outperforms other architectures in many instances.

We extracted 200 data points from the GSM8K training dataset and pruned the FULL architecture with 3 agents for 8 time steps. We applied these pruned architectures on all datasets, with the results shown in Table 3. To align with other architectures,

we used only the first 6 time steps for evaluation. The pruned architectures performed well on the GSM8K and exhibited a certain degree of transfer ability on other datasets. Specifically, it outperforms FULL and RANK on 7 datasets. It is worth noting that we deleted those paths with smaller benefits during pruning, which further reduces the cost. The pruned architecture only contains 46.2% communication paths in the FULL architecture.

7 Related Work

LLM-based multi-agent. In the last few years, researchers have conducted numerous studies on LLM-MAS. Some studies focus on approaching collaborative mechanisms to enhance systems. These studies, e.g., Debate (Du et al., 2023), MAD (Liang et al., 2023), Deepwide (Zhang et al., 2023b), and ChatEval (Chan et al., 2023), concentrated on continuous debates among agents. Other studies focus on the decomposition of complex tasks, such as Camel (Li et al., 2023), ChatDev (Qian et al., 2023), AutoGen (Wu et al., 2023), and MetaGPT (Hong et al., 2023), exploring MAS for task division where different agents responsible for different sub-tasks. Additionally, a series of studies have explored how to use LLMs to simulate human behavior. This includes strategic and sandbox games like Werewolf (Xu et al., 2023a,b), Avalon (Lan et al., 2023), Minecraft (Chen et al., 2023b; Gong et al., 2023), game theory simulation (Fu et al., 2023; Mao et al., 2023; Guo et al., 2023), and sociological simulation (Park et al., 2023; Zhang et al., 2023a). However, the scale, agent credit, and factors related to multi-agent synergy have also not been comprehensively studied.

Collaboration Architecture of multi-agent. Traditional multi-agent research has proposed a variety of possible structures (Horling and Lesser, 2004) such as Flat, Hierarchical, Holonic (Esmaeili et al., 2016), Team, and Congregation (Brooks and Durfee, 2003). In the past few years, some studies have leveraged the capabilities of LLMs to construct more complex MAS. Shi et al. (2023); Du et al. (2023); Liang et al. (2023) organized multiple LLM-based agents for fixed rounds of debates. Chen et al. (2023a) organized agents in the form of a Round-Table Conference. ChatLLM (Hao et al., 2023) and WideDeep (Zhang et al., 2023b) organized agents into linear layers to enhance system capabilities. Zhang et al. (2023c) adopted a dynamic acyclic graph structure during the reasoning

process. Liu et al. (2023b) proposed a dynamic architecture that can adjust according to different queries. Yin et al. (2023) proposed four architectures based on network topology.

Contribution of Agents. Evaluating the contribution of LLM agents is crucial for optimizing MAS. Credit assignment (Agogino and Tumer, 2004), introduced from traditional multi-agent, studies how to measure the impact of actions on global rewards. Extensive research has been delving into this problem, including implicit methods like policy gradients and Q-learning algorithms and explicit methods such as the Shapley value and actor-critic architecture. LLM-MAS studies primarily use extra LLMs for evaluation. Jiang et al. (2023b); Qin et al. (2023); Liu et al. (2023b) ranking outputs of agents to determine contributions. Others calculate contributions based on LLM’s intermediate outcomes, such as the confidence evaluation proposed by (Yin et al., 2023), which calculates the model’s confidence based on the variation in responses.

8 Conclusion and Future Direction

This paper focuses on three main questions: exploring the performance of multi-agent systems under various scenarios, investigating the influence of scale-related factors, and finding potential effective structures by crediting communication paths. Our empirical study offers significant insights for collaboration within MAS, finding that single-agent performance does not decide the performance of multi-agent synergy. Furthermore, our study at scale suggests that adding more agents can lead to better system performance, aligning with the conclusions from (Li et al., 2024). We observed that the system agreement gradually increases as the time step increases. We also pruned the FULL architecture based on the Shapley value, which achieved the best results and demonstrated certain transferability and implied potential effective architecture. Our empirical study on scaling and crediting can be helpful in future studies of LLM-based multi-agent systems.

Limitations

Our study also has some limitations. First, we did not experiment with a MAS consisting of more than five agents due to the limited context length of the open-source model. We plan to use models that support longer contexts for systems with

569	more agents in the future. Besides, an interesting problem arises in Q2: Why does MAS show a performance decline after reaching early stopping on some datasets? According to our case study, this problem came from the accidentally generated error messages and the fast spreading of misinformation. We plan to analyze this phenomenon systematically in the future. Lastly, considering the extra computational costs of Shaley value, using Information Gain and a simplified method from MARL may be better.	
570		
571		
572		
573		
574		
575		
576		
577		
578		
579		
580	References	
581	Adrian K. Agogino and Kagan Tumer. 2004. Unifying temporal and structural credit assignment problems . In <i>3rd International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS 2004)</i> , 19-23 August 2004, New York, NY, USA, pages 980–987. IEEE Computer Society.	
582		
583		
584		
585		
586		
587	Christopher H. Brooks and Edmund H. Durfee. 2003. Congregation formation in multiagent systems . <i>Auton. Agents Multi Agent Syst.</i> , 7(1-2):145–170.	
588		
589		
590	Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate .	
591		
592		
593		
594	Justin Chih-Yao Chen, Swarnadeep Saha, and Mohit Bansal. 2023a. Reconcile: Round-table conference improves reasoning via consensus among diverse llms .	
595		
596		
597		
598	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. 2023b. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors .	
599		
600		
601		
602		
603		
604	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems . <i>ArXiv preprint</i> , abs/2110.14168.	
605		
606		
607		
608		
609		
610	Ariuna Damba and Shigeyoshi Watanabe. 2007. Hierarchical control in a multiagent system . In <i>Second International Conference on Innovative Computing, Informatio and Control (ICICIC 2007)</i> , pages 111–111. IEEE.	
611		
612		
613		
614		
615	Ali Dorri, Salil S. Kanhere, and Raja Jurdak. 2018. Multi-agent systems: A survey . <i>IEEE Access</i> , 6:28573–28593.	
616		
617		
618	Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate . <i>ArXiv preprint</i> , abs/2305.14325.	
619		
620		
621		
	Ahmad Esmaeili, Nasser Mozayani, Mohammad Reza Jahed-Motlagh, and Eric T. Matson. 2016. The impact of diversity on performance of holonic multi-agent systems . <i>Eng. Appl. Artif. Intell.</i> , 55:186–201.	622
		623
		624
		625
	Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023. Improving language model negotiation with self-play and in-context learning from AI feedback . <i>ArXiv preprint</i> , abs/2305.10142.	626
		627
		628
		629
	Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, and Jianfeng Gao. 2023. Mindagent: Emergent gaming interaction . <i>ArXiv preprint</i> , abs/2309.09971.	630
		631
		632
		633
		634
	Jiaxian Guo, Bo Yang, Paul Yoo, Bill Yuchen Lin, Yusuke Iwasawa, and Yutaka Matsuo. 2023. Suspicion-agent: Playing imperfect information games with theory of mind aware GPT-4 . <i>ArXiv preprint</i> , abs/2309.17277.	635
		636
		637
		638
		639
	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges . <i>ArXiv preprint</i> , abs/2402.01680.	640
		641
		642
		643
		644
	Rui Hao, Linmei Hu, Weijian Qi, Qingliu Wu, Yirui Zhang, and Liqiang Nie. 2023. Chatllm network: More brains, more intelligence . <i>ArXiv preprint</i> , abs/2304.12998.	645
		646
		647
		648
	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	649
		650
		651
		652
		653
		654
	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the MATH dataset . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	655
		656
		657
		658
		659
		660
		661
		662
	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, and Chenglin Wu. 2023. Metagpt: Meta programming for multi-agent collaborative framework . <i>ArXiv preprint</i> , abs/2308.00352.	663
		664
		665
		666
		667
		668
	Bryan Horling and Victor R. Lesser. 2004. A survey of multi-agent organizational paradigms . <i>Knowl. Eng. Rev.</i> , 19(4):281–316.	669
		670
		671
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L�elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth�ee Lacroix, and William El Sayed. 2023a. Mistral 7b .	672
		673
		674
		675
		676
		677
		678

679	Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023b.	Bernstein. 2023. Generative agents: Interactive simu-	733
680	LLM-blender: Ensembling large language models	lacra of human behavior . In <i>Proceedings of the 36th</i>	734
681	with pairwise ranking and generative fusion . In <i>Pro-</i>	<i>ceedings of the 36th Annual ACM Symposium on User Interface Software</i>	735
682	<i>ceedings of the 61st Annual Meeting of the Associa-</i>	<i>and Technology, UIST 2023, San Francisco, CA, USA,</i>	736
683	<i>tion for Computational Linguistics (Volume 1: Long</i>	<i>29 October 2023- 1 November 2023, pages 2:1–2:22.</i>	737
684	<i>Papers)</i> , pages 14165–14178, Toronto, Canada. As-	ACM.	738
685	sociation for Computational Linguistics.		
686	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	Lynne E. Parker. 1993. Designing control laws for co-	739
687	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.	operative agent teams . In <i>Proceedings of the 1993</i>	740
688	Gonzalez, Hao Zhang, and Ion Stoica. 2023. Effi-	<i>IEEE International Conference on Robotics and Au-</i>	741
689	cient memory management for large language model	<i>tomation, Atlanta, Georgia, USA, May 1993, pages</i>	742
690	serving with pagedattention. In <i>Proceedings of the</i>	<i>582–587</i> . IEEE Computer Society Press.	743
691	<i>ACM SIGOPS 29th Symposium on Operating Systems</i>		
692	<i>Principles</i> .	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	744
693	Yihuai Lan, Zhiqiang Hu, Lei Wang, Yang Wang, De-	2021. Are NLP models really able to solve simple	745
694	heng Ye, Peilin Zhao, Ee-Peng Lim, Hui Xiong, and	math word problems? In <i>Proceedings of the 2021</i>	746
695	Hao Wang. 2023. Llm-based agent society investi-	<i>Conference of the North American Chapter of the</i>	747
696	gation: Collaboration and confrontation in avalon	<i>Association for Computational Linguistics: Human</i>	748
697	gameplay . <i>ArXiv preprint</i> , abs/2310.14985.	<i>Language Technologies</i> , pages 2080–2094, Online.	749
698		Association for Computational Linguistics.	750
699	Guohao Li, Hasan Abed Al Kader Hammoud, Hani	Chen Qian, Xin Cong, Cheng Yang, Weize Chen,	751
700	Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023.	Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong	752
701	CAMEL: communicative agents for "mind" explo-	Sun. 2023. Communicative agents for software de-	753
702	ration of large scale language model society . <i>ArXiv</i>	velopment . <i>ArXiv preprint</i> , abs/2307.07924.	754
703	<i>preprint</i> , abs/2303.17760.	Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang,	755
704	Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and	Junru Wu, Jiaming Shen, Tianqi Liu, Jialu Liu, Don-	756
705	Deheng Ye. 2024. More agents is all you need .	ald Metzler, Xuanhui Wang, and Michael Bender-	757
706	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	sky. 2023. Large language models are effective text	758
707	Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and	rankers with pairwise ranking prompting . <i>ArXiv</i>	759
708	Shuming Shi. 2023. Encouraging divergent thinking	<i>preprint</i> , abs/2306.17563.	760
709	in large language models through multi-agent debate .	L.S. Shapley and Rand Corporation. 1951. Notes on the	761
710	<i>CoRR</i> , abs/2305.19118.	N-person Game . Notes on the N-person Game. Rand	762
711	Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan	Corporation.	763
712	Duan, Ming Zhou, and Yue Zhang. 2023a. Logiqa	Chenhui Shen, Liying Cheng, Xuan-Phi Nguyen, Yang	764
713	2.0 - an improved dataset for logical reasoning in	You, and Lidong Bing. 2023. Large language mod-	765
714	natural language understanding . <i>IEEE ACM Trans.</i>	els are not yet human-level evaluators for abstrac-	766
715	<i>Audio Speech Lang. Process.</i> , 31:2947–2962.	tive summarization . In <i>Findings of the Association</i>	767
716	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang,	<i>for Computational Linguistics: EMNLP 2023, pages</i>	768
717	Yile Wang, and Yue Zhang. 2020. Logiqa: A chal-	<i>4215–4233</i> , Singapore. Association for Computa-	769
718	lenge dataset for machine reading comprehension	<i>tional Linguistics</i> .	770
719	with logical reasoning . In <i>Proceedings of the Twenty-</i>	Zijing Shi, Meng Fang, Shunfeng Zheng, Shilong Deng,	771
720	<i>Ninth International Joint Conference on Artificial</i>	Ling Chen, and Yali Du. 2023. Cooperation on the	772
721	<i>Intelligence, IJCAI 2020</i> , pages 3622–3628. ijcai.org.	fly: Exploring language agents for ad hoc teamwork	773
722	Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi	in the avalon game . <i>ArXiv preprint</i> , abs/2312.17515.	774
723	Yang. 2023b. Dynamic llm-agent network: An llm-	Peter Sunehag, Guy Lever, Audrunas Gruslys, Woj-	775
724	agent collaboration framework with agent team opti-	ciech Marian Czarnecki, Vinicius Zambaldi, Max	776
725	mization .	Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z	777
726	Shaoguang Mao, Yuzhe Cai, Yan Xia, Wenshan Wu,	Leibo, Karl Tuyls, et al. 2018. Value-decomposition	778
727	Xun Wang, Fengyi Wang, Tao Ge, and Furu Wei.	networks for cooperative multi-agent learning based	779
728	2023. ALYMPICS: language agents meet game the-	on team reward . In <i>Proceedings of the 17th Interna-</i>	780
729	ory . <i>ArXiv preprint</i> , abs/2311.03220.	<i>tional Conference on Autonomous Agents and Multi-</i>	781
730	Marvin Minsky. 1961. Steps toward artificial intelli-	<i>Agent Systems</i> , pages 2085–2087.	782
731	gence . <i>Proceedings of the IRE</i> , 49(1):8–30.	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	783
732	Joon Sung Park, Joseph C. O’Brien, Carrie Jun Cai,	Jonathan Berant. 2019. CommonsenseQA: A ques-	784
	Meredith Ringel Morris, Percy Liang, and Michael S.	tion answering challenge targeting commonsense	785
		knowledge . In <i>Proceedings of the 2019 Conference</i>	786
		<i>of the North American Chapter of the Association for</i>	787
		<i>Computational Linguistics: Human Language Tech-</i>	788
		<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	789

790	4149–4158, Minneapolis, Minnesota. Association for	Zhenran Xu, Senbao Shi, Baotian Hu, Jindi Yu, Dong-	846
791	Computational Linguistics.	fang Li, Min Zhang, and Yuxiang Wu. 2023c. Towards reasoning in large language models via multi-agent peer review collaboration.	847
792	Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bha-		848
793	gavatula, Yoav Goldberg, Yejin Choi, and Jonathan		849
794	Berant. 2022. Commonsenseqa 2.0: Exposing the limits of AI through gamification. <i>CoRR</i> ,	Zhangyue Yin, Qiushi Sun, Cheng Chang, Qipeng	850
795	abs/2201.05320 .	Guo, Junqi Dai, Xuanjing Huang, and Xipeng Qiu.	851
796		2023. Exchange-of-thought: Enhancing large language model capabilities through cross-model communication. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15135–15153, Singapore. Association for	852
797	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Computational Linguistics.	853
798	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		854
799	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti		855
800	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton		856
801	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,		857
802	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	Weihaio Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng.	858
803	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	2020. Reclor: A reading comprehension dataset requiring logical reasoning. In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenRe-	859
804	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	view.net.	860
805	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,		861
806	Isabel Kloumann, Artem Korenev, Punit Singh Koura,		862
807	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-		863
808	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Jintian Zhang, Xin Xu, and Shumin Deng. 2023a. Exploring collaboration mechanisms for llm agents: A social psychology view.	864
809	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-		865
810	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-		866
811	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	Xinghua Zhang, Bowen Yu, Haiyang Yu, Yangyu Lv,	867
812	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	Tingwen Liu, Fei Huang, Hongbo Xu, and Yongbin	868
813	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	Li. 2023b. Wider and deeper llm networks are fairer llm evaluators.	869
814	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,		870
815	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	Yifan Zhang, Jingqin Yang, Yang Yuan, and An-	871
816	Melanie Kambadur, Sharan Narang, Aurelien Ro-	drew Chi-Chih Yao. 2023c. Cumulative reasoning with large language models. <i>ArXiv preprint</i> ,	872
817	driguez, Robert Stojnic, Sergey Edunov, and Thomas	abs/2308.04371 .	873
818	Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.		874
819		Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu,	875
820	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu,	and Jiantao Jiao. 2023. Starling-7b: Improving llm helpfulness harmfulness with rlaiif.	876
821	Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and		877
822	Zhifang Sui. 2023a. Large language models are not fair evaluators.		
823			
824	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.		
825	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-		
826	hery, and Denny Zhou. 2023b. Self-consistency improves chain of thought reasoning in language models. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.		
827			
828			
829			
830			
831	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,		
832	Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,		
833	Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadal-		
834	lah, Ryen W White, Doug Burger, and Chi Wang.		
835	2023. Autogen: Enabling next-gen llm applications via multi-agent conversation.		
836			
837	Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xi-		
838	aolong Wang, Weidong Liu, and Yang Liu. 2023a.		
839	Exploring large language models for communication games: An empirical study on werewolf. <i>ArXiv preprint</i> , abs/2309.04658 .		
840			
841			
842	Zelai Xu, Chao Yu, Fei Fang, Yu Wang, and Yi Wu.		
843	2023b. Language agents with reinforcement learning for strategic play in the werewolf game. <i>ArXiv preprint</i> , abs/2310.18940 .		
844			
845			

878
879
880
881
882
883
884
885
886
887

A Appendix

A.1 The relation between single-agent performance and multi-agent benefit

In section 4, we propose the finding that single-agent performance does not determine multi-agent benefit. To further verify this finding, we made a scatter plot of single-agent system’s accuracy with multi-agent benefit, as shown in Fig. 7. It can be found that there is no obvious correlation between them, which supports the conclusion of section 4.

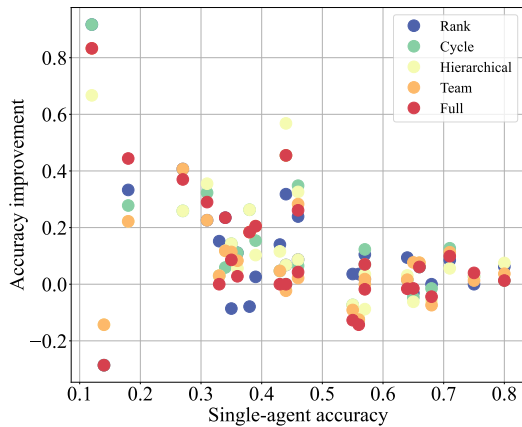


Figure 7: Single-agent accuracy and system improvement

A.2 System performance in every architecture, LLM, and dataset

In section 4, we calculated the average influence of different factors, i.e., architecture, LLM, and dataset, but the absolute performance of each factor was not shown. For this reason, we present all data in Table 4. Keep in mind that the table only contains results for the 3-agent system, considering the cost, we did not conduct such extensive experiments for systems consisting of more agents.

A.3 Shapley value of every path

In section 6, we optimized FULL architecture with Shapley value, but we didn’t present the middle value of the optimization. Here, we show the Shapley value of every path in each optimization time step in Table 5. Noticing that each column depend on the optimized architecture at that time step.

A.4 Results on bigger model

We have done some experiments on larger systems and more capable LLM Mixtral-8x7B-Instruct-v0.1-GPTQ. The results are as follows:

A.5 Agent prompt

We show the role prompt for each agent in Table 7 and Table 8.

910
911
912

888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909

LLM	Architecture	MATH	GSM8K	SVAMP	CSQA	CSQA2	MMLU	LogicQA	LogiQA2	ReClor
llama2	Single Agent	0.14	0.18	0.57	0.31	0.57	0.44	0.38	0.33	0.27
	Full	0.10	0.26	0.61	0.40	0.56	0.44	0.45	0.33	0.37
	Cycle	0.10	0.23	0.64	0.41	0.61	0.44	0.48	0.34	0.34
	Hierarchical	0.12	0.22	0.59	0.42	0.52	0.47	0.48	0.34	0.34
	Team	0.12	0.22	0.58	0.38	0.57	0.43	0.45	0.34	0.38
	Rank	0.10	0.24	0.58	0.38	0.63	0.47	0.35	0.38	0.38
mistral	Single Agent	0.12	0.46	0.66	0.46	0.65	0.44	0.35	0.39	0.43
	Full	0.22	0.48	0.70	0.58	0.64	0.64	0.38	0.47	0.43
	Cycle	0.23	0.49	0.70	0.62	0.62	0.64	0.40	0.45	0.45
	Hierarchical	0.20	0.50	0.71	0.61	0.61	0.69	0.40	0.43	0.48
	Team	0.22	0.47	0.71	0.59	0.70	0.64	0.39	0.47	0.45
	Rank	0.23	0.50	0.70	0.57	0.63	0.58	0.32	0.40	0.49
starling	Single Agent	0.34	0.75	0.80	0.71	0.68	0.64	0.36	0.55	0.56
	Full	0.42	0.78	0.81	0.78	0.65	0.63	0.37	0.48	0.48
	Cycle	0.36	0.77	0.81	0.80	0.67	0.65	0.40	0.51	0.49
	Hierarchical	0.38	0.77	0.86	0.75	0.64	0.66	0.38	0.51	0.50
	Team	0.38	0.76	0.83	0.79	0.63	0.65	0.39	0.50	0.49
	Rank	0.42	0.75	0.85	0.77	0.68	0.70	0.40	0.57	0.58

Table 4: System accuracy on every system and dataset. Systems based on 3 agents. The max time step is 6.

Time Step	1	2	3	4	5	6	7	8
$Path(0,0)$	-0.043	0.009	0.006	-0.001	-0.009	0.008	0.010	-0.015
$Path(1,0)$	0.033	0.002	0.001	-0.008	0.018	-0.003	0.000	0.000
$Path(2,0)$	0.035	-0.031	-0.022	0.009	-0.004	-0.005	-0.005	0.000
$Path(0,1)$	0.013	0.016	0.007	-0.003	-0.032	-0.009	-0.006	-0.013
$Path(1,1)$	0.003	0.016	0.014	-0.005	0.021	0.016	-0.006	0.004
$Path(2,1)$	0.000	-0.017	-0.006	0.013	-0.004	-0.007	0.002	-0.011
$Path(0,2)$	-0.008	0.025	0.019	0.002	-0.008	-0.008	0.010	-0.001
$Path(1,2)$	0.015	0.010	0.014	-0.016	0.012	0.017	-0.020	-0.003
$Path(2,2)$	0.008	-0.020	-0.033	0.014	-0.008	-0.003	0.000	0.004

Table 5: The Shapley value for every path in every time step during optimizing, the $Path(i, j)$ denote the path from agent i to agent j

Agent Number	1	2	3	4	5	6	7	8	9	10
MMLU	0.67	0.80	0.73	0.73	0.75	0.76	-	-	-	-
GSM8k	0.68	0.74	0.77	0.77	0.76	0.79	0.81	0.78	0.79	0.8
LogiQA2	0.48	0.45	0.46	0.51	0.53	0.52	-	-	-	-

Table 6: The result of Mixtral-8x7B-Instrument-v0.1-GPTQ with Full architecture

[System Prompt]

You are an excellent and very capable domain question solver. You are now invited to an expert group of processing and solving domain application questions. Your codename in the expert group is Expert self.rriid. As a distinguished member of the expert group, you possess the capability to a broad range of domain disciplines, allowing you to adapt and apply the appropriate methodologies to the given questions.

[User Prompt]

Task Description

Your task is to systematically address the **domain** application question presented below, decipher complex question statements and elucidate your reasoning in a sequential, step-by-step fashion. Carefully utilize the provided information to work through the question. Your answer should be both concise and comprehensive, detailing the logical progression of your thought process. Besides, the expert group have provided some potential answers to this question, you should consider insights from these answers to enrich the quality and accuracy of your own answer.

Given Question

Question: **question**

Given Question Again

Read the given question again.

Question: **question**

Answers by Other Experts

There are some potential answers provided by different experts for the same question. Consider these responses to cross-verify your approach, broaden your understanding, and gain alternative perspectives with diverse approaches to the question-solving process. This may help you ensure consistency and accuracy in your methodology. However, we have not verified the correctness of these answers, so be careful of the quality and relevance of these answers.

messages

Output Format

start

Opinion: your opinion about other experts' answers

Solution: your detailed, step-by-step solution, final answer is formatted as "[final answer here]"

end

The output start with your opinion about other experts' answers, followed by your step-by-step solution in the next line.

Remember that your final answer in the solution is surrounded by '[' and ']', which is formatted as "[final answer here]".

Now take a deep breath and solve the question step by step.

Table 7: The prompt template for agent. We replace the colored slot with real text before querying the LLMs. Note that we use a similar template when conducting single-agent-based experiments and ignore the Answers by Other Experts.

[System Prompt]

You are an excellent and very capable domain question solver. You are now invited to an expert group of processing and solving domain application questions. Your codename in the expert group is Expert self.rrid. As a distinguished member of the expert group, you possess the capability to a broad range of domain disciplines, allowing you to adapt and apply the appropriate methodologies to the given questions.

[User Prompt]

Task Description

Your task is to systematically address the **domain** application question presented below, decipher complex question statements and elucidate your reasoning in a sequential, step-by-step fashion. Carefully utilize the provided information to work through the question. Your answer should be both concise and comprehensive, detailing the logical progression of your thought process. Besides, the expert group have provided some potential answers to this question, you should consider insights from these answers to enrich the quality and accuracy of your own answer.

Given Question

Question: **question**

Given Question Again

Read the given question again.

Question: **question**

Your Previous Answer

You have generated a previous answer ...
most recent answer, if it is transmitted.

Answers by Other Experts There are some potential answers provided by different experts for the same question. Consider these responses to cross-verify your approach, broaden your understanding, and gain alternative perspectives with diverse approaches to the question-solving process. This may help you ensure consistency and accuracy in your methodology. However, we have not verified the correctness of these answers, so be careful of the quality and relevance of these answers.

messages

Output Format

start

Opinion: your opinion about other experts' answers

Solution: your detailed, step-by-step solution, final answer is formatted as "[final answer here]"

end

The output start with your opinion about other experts' answers, followed by your step-by-step solution in the next line.

Remember that your final answer in the solution is surrounded by '[' and ']', which is formatted as "[final answer here]".

Now take a deep breath and solve the question step by step.

Table 8: The prompt template for interaction between agents.