

Towards Multi-Agent Reasoning Systems for Collaborative Expertise Delegation: An Exploratory Design Study

Anonymous ACL submission

Abstract

Designing effective collaboration structure for multi-agent systems to stimulate collective reasoning capability is crucial yet remains under-explored. In this paper, we systematically investigate how collaborative reasoning performance is affected by three key design factors: (1) expertise-domain alignment, (2) collaboration paradigm, and (3) system scale. Our findings reveal that expertise alignment benefits are highly domain-contingent, proving most effective for contextual reasoning tasks. Furthermore, collaboration focused on integrating diverse responses consistently outperforms sequential functional cooperation. Finally, we empirically explore the impact of scaling the multi-agent system with expertise specialization and analyze the resulting performance-computational cost trade-off, highlighting the need for more efficient communication protocol design. Our work provides concrete guidelines for configuring multi-agent reasoning system with expertise role delegation.

1 Introduction

Collective intelligence, the emergent problem-solving capability arising from structured group interactions, has long been recognized as a cornerstone of complex human decision-making (Surowiecki, 2004). Through mechanisms like deliberative debate and systematic knowledge integration, human collectives consistently outperform individual experts in tasks requiring multi-perspective analysis and contextual synthesis.

The recent evolution of large language and reasoning models (LLMs/LRMs) has spurred parallel investigations into machine collective intelligence through multi-agent systems—artificial analogs of human collaboration patterns (Yang et al., 2025; Jaech et al., 2024; Team et al., 2025). A common technique predominantly being deployed in this area is called **expertise role delegation**, where LLMs are instructed to simulate specific expert

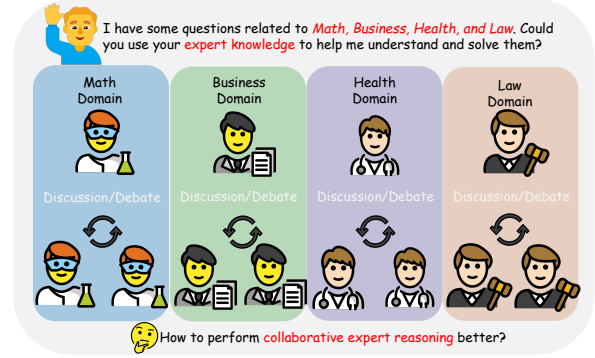


Figure 1: Workflow diagram for a multi-agent reasoning system with specialized agents.

personas (Li et al., 2024a; Xu et al., 2024a). Despite broad adoption in various multi-agent frameworks, the usage of expertise role delegation remains heuristic without rigorous analysis on the assignment of expertise and agent interplay dynamics. This methodological gap mirrors challenges in organizational science, where empirical social theories have systematically decoded the principles governing effective group collaboration. Durkheim’s *division of labor theory* emphasizes that aligning expertise with downstream tasks is critical for group efficacy (Durkheim, 1893). Complementarily, research on *process loss*—information degradation during collective action—demonstrates that inefficiency in system arises from structural configurations and group magnitude (Steiner, 1972).

Anchored in these theoretical foundations for effective collaboration, we decompose multi-agent system design with expertise role delegation into three critical dimensions: (1) expertise-domain alignment, (2) collaboration paradigm, and (3) system scale. Rather than proposing another task-oriented framework, this study serves as the first exploratory analysis investigating how these critical dimensions impact multi-agent system efficacy, providing empirical foundations for future expertise-driven multi-agent architectures.

We first focus on expertise-domain alignment—where existing practice reveals significant gaps. Although expertise specialization is widely adopted in multi-agent systems (Wang et al., 2024a; Li et al., 2024a), the impact of collaborative expertise configuration on downstream scenarios remains underexplored. This ambiguity creates practical challenges in configuring expert roles for task domains. To address this gap, we empirically evaluate the influence of different collaborative expertise configurations on task performance across four representative domains from MMLU-pro (Wang et al., 2024d). **Our findings in Section 4 demonstrate a positive correlation between task performance and the alignment of group expertise with the task domain, underscoring the necessity of accurately matching the multi-agent system expertise with downstream tasks.**

Having established the critical role of expertise-domain alignment, we then examine how collaboration paradigms structure interactions between specialized agents. Currently, the collaboration paradigm predominantly used in recent studies could be categorized into two kinds: (1) Diversity-Driven Perspective Integration, where agents, often embodying different viewpoints or roles, are encouraged to generate diverse responses to enrich the solution space (Wang et al., 2024b; Chen et al., 2024b; Hu et al., 2025). (2) Structured Workflow Cooperation, where different agents are assigned distinct sub-tasks within a predefined pipeline to collaboratively construct a solution (Chen et al., 2024c; Hong et al., 2024; Zhang et al., 2025). We design comparative experiments to unveil the performance differences between paradigms. **Our observations in Section 5 reveal a consistent advantage for diversity-driven collaboration over structured workflow collaboration, suggesting the superiority of the diversity-driven paradigm.**

Finally, constructing large-scale multi-agent system has become a critical, yet often enigmatic aspect of multi-agent system design (Chen et al., 2024c; Piao et al., 2025). While intuition and some preliminary studies (Qian et al., 2024; Li et al., 2023) suggest that larger groups would lead to better reasoning performance, the actual effectiveness of scaling within the context of collaborative expertise specialization and the potential computation-performance trade-off, are not well understood. Our systematic experiments involve incrementally increasing the system scale to examine potential scaling laws. **The results in Section 6 un-**

cover non-linear dynamics; specifically, adding more experts tends to improve the collective reasoning ability of the system. This positive trend holds regardless of whether the larger system scale contains greater viewpoint diversity or a more comprehensive workflow structure, indicating a general benefit to increasing the number of expert agents and encouraging such designs for enhanced system performance. However, our analysis of the computational trade-offs associated with system scaling reveals that, while the system would benefit from the expansion, there remains a critical need for more efficient communication protocols between agents for more scalable and cost-effective multi-agent reasoning process.

2 Related Works

2.1 Multi-Agent Collaboration

Multi-Agent Collaboration adopts multiple LLMs to solve the problem collaboratively. Abundant researches have investigated the multi-agent collaboration framework to improve decision-making capability of the system (Wang et al., 2024b; Liang et al., 2024; Du et al., 2024). In addition to collaboration among LLMs, several researchers instruct the agents to cooperate in a workflow to study the multi-agent systems’ ability of solving real world challenges (Li et al., 2024b; Xu et al., 2024b; Chen et al., 2024a). While Qian et al. (2024), Yang et al. (2024) and Wang et al. (2024c) has investigate the effect of varying the scale of multi-agent system on reasoning and simulation, prior researches have not systematically examined the interplay between collective expertise specialization, collaboration mechanisms, and the impact of system scale simultaneously. In this work, we conduct extensive experiments to formally analyze the influence of these three critical dimensions on multi-agent collaborative reasoning. Our findings provide actionable insights toward more effective system design.

2.2 LLMs as Domain Experts

The rapid evolution of LLMs has endowed them with vast repositories of domain-specific knowledge, enabling their application across a wide range of expert tasks. Recent researches have explored the potential of LLMs to emulate specific personas by conditioning them on detailed character profiles (Chan et al., 2024; Samuel et al., 2024; Xu et al., 2023). These studies demonstrate that by providing LLMs with demographic or

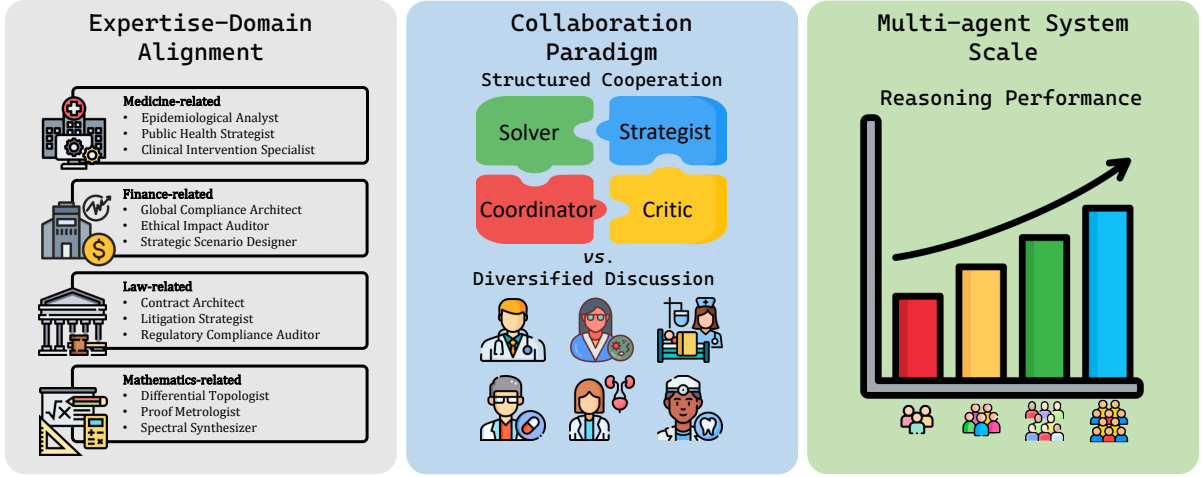


Figure 2: Demonstration of three key factors characterizing research on multi-agent collaborative reasoning systems. (1) expertise-domain alignment, (2) collaboration paradigm, and (3) scale of the multi-agent system.

role-specific prompts, they can effectively exhibit human-like personality traits and behaviors. Furthermore, Kong et al. (2024) and Xu et al. (2023) have shown that instructing LLMs to simulate domain experts can enhance their reasoning capabilities in specialized contexts, underscoring the necessity of introducing expert knowledge into reasoning process. Despite these advancements and the growing prominence of multi-agent systems in research, the specific impact of collaborative expertise specialization on reasoning performance remains underexplored. In this paper, through meticulously designed experiments, we systematically investigate the impact of expertise specialization within multi-agent reasoning systems. Our findings reveal that simulating specialized roles significantly enhances performance on tasks requiring contextual reasoning, while showing limited influence on those primarily dependent on factual recall or mathematical deduction.

3 Preliminary

3.1 Problem Setup

Formally, given a multi-agent system $\mathcal{M}_n = \{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n\}$ where n indicates the number of agents inside the system and \mathcal{A}_i represents the i -th agent of the system, a query \mathcal{Q} , and a set of candidate options \mathcal{S} . A multi-agent system reasoning process is expressed as:

$$\mathcal{Y} = \mathcal{F}(\mathcal{A}_1(\mathcal{Q}, \mathcal{S}), \mathcal{A}_2(\mathcal{Q}, \mathcal{S}), \dots, \mathcal{A}_n(\mathcal{Q}, \mathcal{S}))$$

where \mathcal{Y} stands for the final answer generated by the system. $\mathcal{A}_i(\mathcal{Q}, \mathcal{S})$ represents the answer of agent i , \mathcal{F} stands for the communication proto-

col manually customized by the design of the system which aggregate the answer of each agents into the final answer. Typically, it could be majority vote, debate, etc (Kaesberg et al., 2025; Liu et al., 2024a). In our specific setup, we adopt a sequential processing communication mechanism inspired by Qian et al. (2024) to prevent context explosion (Liu et al., 2024b; Xu et al., 2024c). In this mechanism, for $i = 2, \dots, n$, agent \mathcal{A}_i receives the complete output generated by the immediately preceding agent \mathcal{A}_{i-1} . In contrast, from the preceding agents $\{\mathcal{A}_1, \dots, \mathcal{A}_{i-2}\}$, \mathcal{A}_i receives only the final answers. The detailed communication algorithm could be found in Appendix A Algorithm 1.

3.2 Dataset

For our experiments, we select four distinct domains from MMLU-pro (Wang et al., 2024d): Math, Health, Business, and Law. These four domains are selected for being representative and frequently studied in contemporary multi-agent reasoning research (Cui et al., 2023; Lei et al., 2024; Ghezloo et al., 2025). We further classify these four domains into three categories based on the primary reasoning type required: (1) **Mathematical Reasoning**: Domains requiring formal mathematical deduction to derive the answer. (2) **Factual Recall Reasoning**: Domains primarily requiring the recall of domain-specific factual knowledge, seldom needing extensive reasoning steps other than simple mathematical calculations. (3) **Contextual Reasoning**: Domains requiring not only the retrieval of relevant expert knowledge but also its application within the reasoning process of specific scenarios or contexts. This choice of evaluation domains and

fine-grained classification of their reasoning types allow us to investigate the effects of collaborative expertise specialization on multi-agent system from a more systematic manner.

3.3 Collaborative Expertise Specialization

In this paper, we primarily studied the effect of collaborative expertise specialization on better multi-agent system design from the perspective of expert-domain alignment, collaboration paradigms and system scale. To formalize the role and responsibility of the agents in the multi-agent system, we define each expert to be of the following format:

$$\mathcal{A}_i \leftarrow (EG, FR, R, ID)$$

where \mathcal{A}_i stands for agent i , EG , FR and R represent Expert Group, Formal Role and Responsibility respectively. ID represents an agent’s index within the group of all agents who share the same role.

3.4 General Experiment Setup

As detailed in Section 3.2, we select 4 representative domains from MMLU-pro to investigate the effects of collaborative expertise specialization. To be consistent with all experiments, we utilize DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI et al., 2025) as the foundational model for all agents. Each agent is initialized with its specific expert description and responsibilities via its system prompt, while the task instance is provided through the user prompt. The detailed prompts could be found in Appendix B. All experiments adopt accuracy as the evaluation metric.

4 Leveraging the “Right” Agent

Expertise specialization is a widely adopted technique in agent research, demonstrably enhancing the reasoning capabilities of LLMs within specific domains (Li et al., 2024b). While the benefits of specialization for individual agents are well-established, the effect of collaborative expertise specialization on the collective reasoning performance of multi-agent systems remains underexplored. This section presents our experimental investigation into this critical area, designed to unveil how different collaborative expertise specialization configurations influence the reasoning capabilities of multi-agent systems.

4.1 Setup

Considering the primary principle of multi-agent reasoning system is to incorporate more diverse

agent viewpoints and integrate them in the final answer (Liang et al., 2024), in our experiments, we adopt diversity-driven collaboration paradigm where we distribute each agent with a specific domain expert configuration and instruct them to generate responses based on their expertise. At this stage, we fix the size of the multi-agent reasoning system to be 3 for controllable computational cost. We employ GPT-4o (OpenAI et al., 2024) for expert configuration generation. The detailed prompts utilized for this automated role generation process are provided in Appendix C.

4.2 “Right” Expertise Helps Reasoning

Our experiments demonstrate a clear performance advantage when the collaborative expertise specialization of the multi-agent system aligns with the domains of the downstream task.

Misaligned expertise configurations often underperform compared to aligned ones. This primary finding is quantitatively supported by the results presented in Table 1. Specifically, in 75% of the aligned cases (diagonal entries), the system achieves the highest accuracy compared to configurations where the agent group simulates expertise from other domains for the same task.

To gain a more nuanced understanding of when expertise alignment is most beneficial, we analyze the system performance according to the primary reasoning type required by each domain, as categorized in Section 3.2. Our analysis reveals that the benefits of expertise alignment are most pronounced for tasks demanding contextual reasoning—Health and Law. Systems operating on these two domains exhibit an average relative performance improvement of 6.75% when expertise is correctly aligned, compared to the misaligned configurations which perform the second best for those tasks. Conversely, for domains requiring mathematical reasoning—Math and Business, the specialized experts yield only marginal gains or even degradation relative to misaligned configurations. We hypothesize this divergence stems from the inherent strengths of LLMs on math. These models often possess robust mathematical reasoning capabilities due to extensive pre-training, potentially reducing the added value of specialized agents. Contextual reasoning tasks, however, appear to benefit more from the structured integration of specialized perspectives provided by the multi-agent reasoning system since applying domain knowledge in these contexts often requires nuanced interpretation, syn-

Dom.\Exp.	Math	Fina	Med	Law	Δ_h	Δ_{abs}
Math	78.0	76.3	76.3	<u>76.4</u>	2.1%	1.6 \uparrow
Business	65.4	<u>64.3</u>	62.4	62.4	-1.7%	1.1 \downarrow
Health	<u>28.9</u>	26.8	30.4	26.1	5.2%	1.5 \uparrow
Law	18.3	<u>19.2</u>	18.5	20.8	8.3%	1.6 \uparrow

Table 1: This table shows the impact of collaborative expertise specialization for different expert groups across various domains. "Dom." and "Exp." abbreviate Domain and Expert Group, respectively. $\Delta_{rel}/\Delta_{abs}$ indicate the relative/absolute performance improvement of the domain-aligned expert group compared to the best-performing alternative group respectively.

thesis of information, and reasoning beyond direct mathematical deduction.

4.3 Analysis on Expert-Domain Alignment

Furthermore, our experimental results reveal a positive correlation between how well the simulated group expertise aligns with the downstream task domain and the observed performance gain. This relationship is visualized in the expertise-domain correlation heatmap presented in Figure 3. Specifically, configurations where the simulated expertise is more relevant to the target task domain tend to yield greater performance improvements compared to less relevant configurations.

To quantify this expertise-domain relevance, we first establish a relevance matrix. We randomly sample 100 instances from each of the four primary task domains. For each instance, we prompt Deepseek-V3 (DeepSeek-AI et al., 2025) to identify a list of 2-3 key expertise domains pertinent to solving the task. We then aggregate these identified candidate domains across all instances within each primary task domain. The relevance scores are calculated by counting the occurrences where a specific knowledge domain (e.g., Business) is deemed relevant for tasks in a primary domain (e.g., Math). These frequencies form a relevance matrix, visualized as a heatmap in Figure 3, where deeper color indicate higher relevance scores.

Comparing this relevance heatmap with the results in Table 1, we observe a consistent pattern supporting our initial finding—Higher expertise-domain relevance, indicated by deeper colors in the heatmap entries, generally corresponds to better reasoning performance. Many cells with high relevance scores in Figure 3 correspond to performance that are bolded or underlined in Table 1, signify-

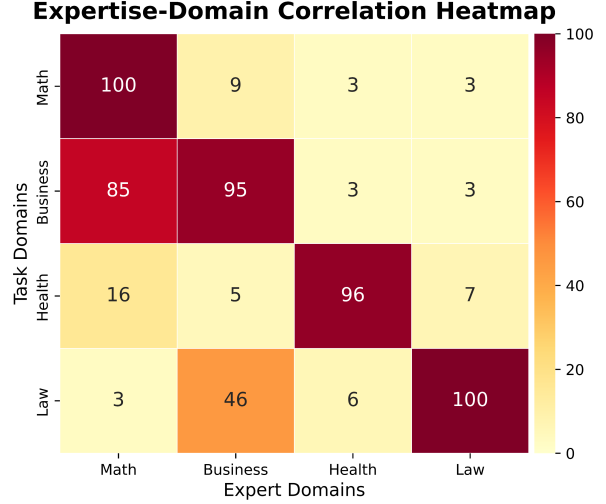


Figure 3: Heatmap illustrating the correlation between specialized group expertise and task domains. Deeper colors indicate stronger correlations.

ing the best or second-best performance among group expertise specialization performance for that task domain. Conversely, low relevance scores typically correspond to misaligned configurations which barely demonstrate distinct advantages conferred by their specific (misaligned) expertise.

Our findings further support the established use of collective expertise specialization in multi-agent reasoning systems, while simultaneously highlighting the critical importance of aligning expertise design with the specific requirements of the target downstream domains, paving a fundamental guidance for future specialization technique application in multi-agent reasoning system design.

5 Collaborate in Efficient Way

Process loss theory elevates collaboration paradigm selection as a critical determinant in multi-agent system efficacy (Steiner, 1972). Crucially, even with optimal domain expertise among agents, the interaction mechanism governing their collaboration fundamentally mediates collective performance through information degradation pathways. In this section, we present comparative experiments designed to analyze these distinct collaboration paradigms. Our objective is to investigate their potential advantages, thereby providing empirically grounded insights for effective collaboration paradigm choice in multi-agent system design.

5.1 Setup

Our analysis leverages the results presented in Figure 4, where we demonstrate both domain-wise

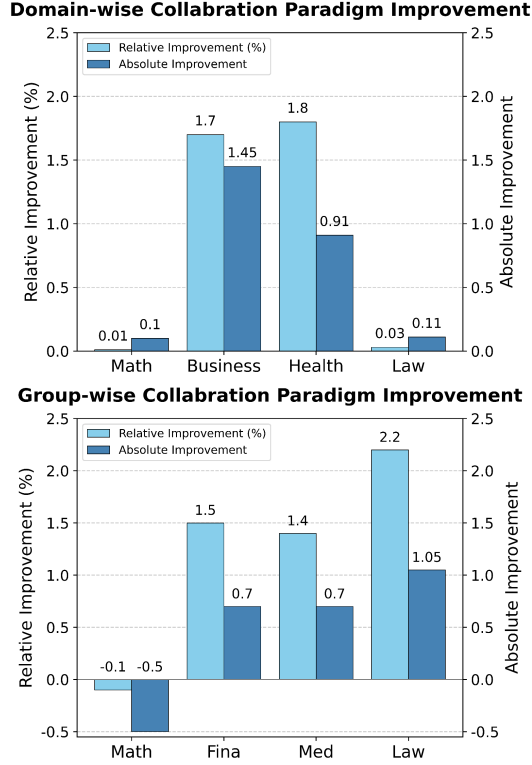


Figure 4: Comparative analysis of diversity-driven versus structured workflow collaboration paradigms. Positive values signify Diversity-Driven’s advantage over Structured Workflow.

and group-wise comparisons for a comprehensive overview. The detailed distinction between paradigms are illustrated as follows:

Diversity-Driven Collaboration: This paradigm emphasizes assigning agents highly specialized, fine-grained expertise within a broader domain (e.g., specific sub-fields of Laws). The objective is to foster collaboration through the integration of diverse, complementary knowledge perspectives during the reasoning process. Each agent contributes deep expertise from a narrow viewpoint.

Structured Workflow Collaboration: Conversely, this paradigm assigns roles based on distinct functional responsibilities within a predefined problem-solving process, in our case, solver, critic and coordinator. Collaboration centers on agents executing specific steps and refining intermediate outputs based on their functional role, rather than primarily contributing unique domain knowledge specializations. The differentiation between agents stems from their function within the workflow.

To ensure a plausible, accurate generation of expert role descriptions, we continue to employ GPT-4o with collaboration paradigm as extra input.

5.2 Diversity Matters in Collaboration

Our primary finding is that the diversity-driven paradigm generally yields superior performance compared to the structured workflow paradigm. This advantage holds true both when considering performance from both domain-wise and group-wise perspectives.

A domain-wise analysis, depicted in Figure 4, confirms this trend. Irrespective of the domain’s primary reasoning type categorized in Section 3.2, the diversity-driven approach consistently results in performance gains over structured workflow. Notably, the most substantial improvements are observed in business and health domains, which demonstrate an average relative performance increase of 1.75% under diversity-driven paradigm. This indicates the potential of expertise with finer-granularity perform well across different domains.

Examining the results from group-wise perspective further supports this conclusion. With the exception of math expert group, all other specialized groups achieve higher average performance across all task domains when employing diversity-driven paradigm. When including the math group, the overall average relative performance improvement facilitated by the diversity-driven approach across all groups is 1.25%, indicating consistent benefits regardless of the task domain encountered.

Synthesizing these observations, the diversity-driven collaboration paradigm demonstrates a consistent performance advantage over structured workflow collaboration paradigm across both different tested domains and distinct expertise configurations. This suggests that multi-agent systems could benefit significantly from collaboration structures that emphasize fine-grained expertise allocation which stimulates viewpoint diversity, providing a solid empirical basis for future research directions in designing multi-agent reasoning system’s collaboration pattern.

5.3 Analysis on Response Diversity

To quantitatively characterize how the collaboration paradigm influences the diversity of agent contributions, we further design a response diversity analysis. We leverage semantic embeddings derived from Sentence-BERT (Reimers and Gurevych, 2019). For each task instance solved by the multi-agent system, we generate embeddings for the output of each agent and measure the internal diversity of the system’s responses by calcu-

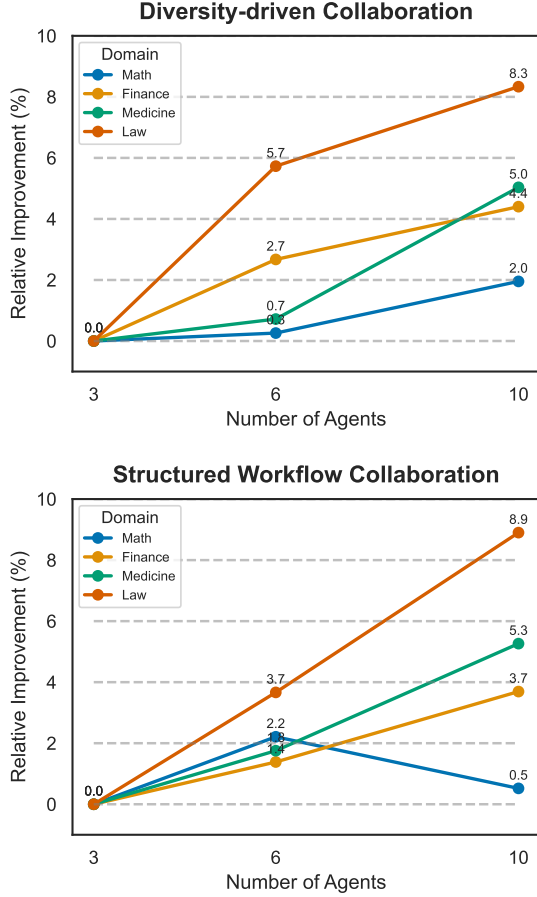


Figure 5: Domain-wise relative performance improvement by scaling up the multi-agent system (3, 6, and 10 agents), shown for different collaboration mechanisms.

lating the pairwise cosine similarity between the embeddings of outputs from different agents. This serves as a measure of how semantically distinct the contributions are at different stages.

The results clearly indicate that, the pairwise cosine similarity values are consistently lower for the diversity-driven collaboration paradigm compared to the structured workflow paradigm. This finding demonstrates that the diversity-driven approach, which emphasizes fine-grained expertise, fosters greater semantic diversity among agent responses throughout the collaborative reasoning, confirming the hypothesis that response diversity matters in multi-agent system. The distribution of the similarity scores could be found in Appendix D

6 Scaling Up Reasoning Experts

Finally, another dimension mentioned by process loss theory is the system scale. While the deployment of large-scale multi-agent systems for simulating social behaviors has received considerable attention, the implications of scaling under collab-

orative expertise specialization setup remain unexplored.

This section details our investigation into the effects of varying system scale on both the reasoning performance of multi-agent systems and the associated computational trade-offs. We aim to elucidate how increasing the number of agents influences collective reasoning efficacy and to call for a better communication protocol design through our performance/token overhead trade-off analysis.

6.1 Setup

We expand our experimental setup from 3 agents to systems comprising 6 and 10 agents. For these larger systems, we systematically replicate the experiments previously introduced, allowing for a direct comparison across different scales.

Generating coherent and appropriately specialized expert role configurations for these larger systems requires extending the initial configurations of the 3 agent system and we continue to leverage GPT-4o for this purpose. The detailed prompts employed for this role augmentation process are provided in Appendix C

6.2 More Experts, More Intelligent System

We evaluate the effect of system scale on reasoning performance by comparing the results from larger agent systems against the baseline 3 agent system. Specifically, we calculate the domain-wise relative performance difference for the system size of 6 and 10 with respect to system of size 3. These relative performance differences are illustrated in Figure 5.

Our findings reveal a consistent trend: increasing the number of agents generally enhances the multi-agent system’s reasoning performance across the evaluated domains, regardless of whether diversity-driven or structured workflow paradigm is employed. However, the magnitude of this improvement varies significantly by domain. Corroborating our earlier observations regarding domain-specific analysis in Section 4, the performance gains within math domain are marginal, even when scaling up to 10 agents. Conversely, domains that necessitate substantial contextual reasoning and knowledge application demonstrate significantly larger performance improvements with increased system scale. This disparity suggests that the benefits derived from incorporating additional agents are most pronounced for tasks requiring the integration of diverse knowledge perspectives or complex, case-specific analysis inherent in non-mathematical rea-

soning. For domains characterized by intense mathematical reasoning, simply increasing the number of agents could barely yield diminishing returns. We believe our finding offers valuable insight for constructing large-scale multi-agent systems intended for diverse domains.

6.3 Token-Performance Trade-off

We further explore the token-performance trade-off inherent in scaling multi-agent reasoning systems by calculating the ratio of performance improvement over token overhead (PoT) with quantitative results presented in Figure 6. We use the sum of reasoning token and answer token for the calculation of token overhead. All the performance improvement and token consumption overhead are counted relatively against system of size 3.

Our analysis reveals distinct trends both across and within domains. Cross-domain comparisons demonstrate that tasks requiring substantial contextual reasoning, such as those in health and law, yield higher PoT ratios. This suggests that increasing agent collaboration is particularly beneficial in these areas, as greater token consumption during the reasoning process leads to higher performance improvements. Conversely, mathematical reasoning tasks exhibit only marginal performance gains with additional agents, which implies smaller ensembles can achieve comparable performance with lower computational overhead, making large-scale multi-agent systems unnecessary for these tasks.

For intra-domain analysis, while structured workflows improved PoT in 75% of domains and diversity-driven approaches in 50% respectively, the critical finding is that neither collaboration paradigm guarantees an enhanced PoT across all domains tested. This widespread inconsistency in scaling behavior, regardless of the collaboration paradigm, highlights the pressing need for advancements in multi-agent communication protocols to achieve more stable and predictable performance enhancements as system complexity increases.

7 Implications on System Design

In this section, we provide implications concluded from our observations for future expertise-driven multi-agent framework designs based on different downstream task category introduced in section 3.2. **For mathematical reasoning and factual recall tasks, minimalist configurations prove optimal: 3 agents with broad role delegation and**

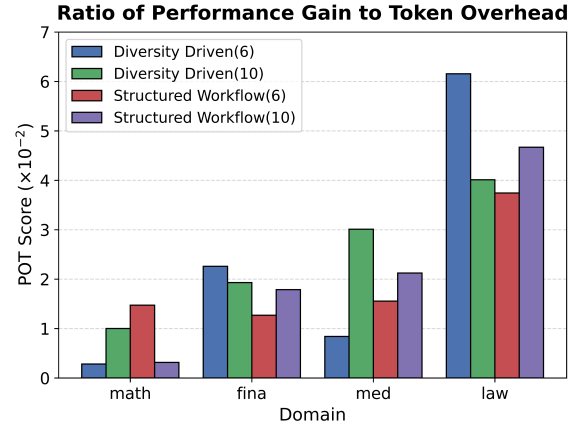


Figure 6: Performance improvement versus token overhead ratio across different domains. Both performance and token overhead are measured as relative increases compared to the system of size 3.

lightweight coordination maximize efficiency. Since domain-specialized LLMs inherently excel at these tasks, complex collaboration introduces unnecessary process loss without performance gains. Simple knowledge verification through brief discussion suffices. **Conversely, contextual and abstract reasoning tasks (e.g., law/medical domains) necessitate structured frameworks. We recommend larger multi-agent system with carefully curated expertise-task alignment and collaboration paradigm promoting viewpoint diversity.** This compensates for LLMs’ limitations in contextual reasoning by enabling complementary knowledge integration and reducing solution-space collapse through deliberate debate protocols. These empirically-derived principles establish a concrete foundation for future multi-agent system frameworks, providing designers with task-aware design heuristics for expertise-driven architectures.

8 Conclusions

In conclusion, this paper systematically investigates the three factors of multi-agent system expertise specialization on collective reasoning intelligence: expertise-domain alignment, collaboration paradigm, and system scale. Our experiments verify the advantage brought by expertise specialization in multi-agent reasoning system, demonstrate the superiority of diversity-driven collaboration and indicate the existence of scaling law in multi-agent reasoning system with experts. These findings provide actionable insights for designing specialized multi-agent reasoning systems in future researches and underscore the need for developing more efficient coordination protocol as systems scale.

Limitations

Our adoption of MMLU-pro for evaluating specialized multi-agent reasoning system across diverse domains, while leveraging its strength in assessing varied domain-specific knowledge, inherently limits our assessment scope. Specifically, its focus on these reasoning paradigms means other crucial multi-agent capabilities, such as coding, might be overlooked. Apart from that, to enhance alignment with real-world scenarios, our evaluation concentrates on four key domains: Math, Business, Health, and Law, selected for their prominence in mainstream research. A direct limitation of this focused approach is that other potentially relevant domains would remain underexplored in the present study. Moreover, To simplify the research setup and promote more stable conclusions, we exclusively utilize one message propagation mechanism. This methodological choice, however, means that the potential influence of diverse communication strategies on system performance remains an unexplored aspect in our current study. Finally, We select DeepSeek-R1-Distilled-Qwen-7B as the base model for all experiments to ensure controllable computational overhead. This decision, while practical, limits our current investigation, deferring the study of multi-agent system architectures with larger-scale models to future research.

Ethics Statement

Our study involves publicly available datasets and use Large Language Models through APIs. Consequently, the ethical considerations of this paper could be listed as follow:

Datasets: We use publicly available datasets only for academic research purpose. We guarantee no personal data has been involved.

LLMs API: Our application of LLMs conform API provider’s policy strictly, maintaining fair use and respecting intellectual property.

Transparency: We provide detailed descriptions of our method and the prompts used in our experiments, in line with standard practices in the research community. We will also make our code publicly available upon acceptance.

References

Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. [Scaling synthetic data creation with 1,000,000,000 personas](#). *CoRR*, abs/2406.20094.

Guhong Chen, Liyang Fan, Zihan Gong, Nan Xie, Zixuan Li, Ziqiang Liu, Chengming Li, Qiang Qu, Shiwen Ni, and Min Yang. 2024a. [Agentcourt: Simulating court with adversarial evolvable lawyer agents](#). *CoRR*, abs/2408.08089.

Justin Chih-Yao Chen, Swarnadeep Saha, and Mohit Bansal. 2024b. [Reconcile: Round-table conference improves reasoning via consensus among diverse llms](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 7066–7085. Association for Computational Linguistics.

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. 2024c. [Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. 2023. [Chatlaw: Open-source legal large language model with integrated external knowledge bases](#). *CoRR*, abs/2306.16092.

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, and S. S. Li. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *CoRR*, abs/2501.12948.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2024. [Improving factuality and reasoning in language models through multiagent debate](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

737	Émile Durkheim. 1893. <i>The Division of Labor in Society</i> . Free Press, New York, NY. Reprinted 1984.	797
738		798
739	Fatemeh Ghezloo, Mehmet Saygin Seyfioglu, Rustin Soraki, Wisdom Oluchi Ikezogwo, Beibin Li, Tejoram Vivekanandan, Joann G. Elmore, Ranjay Krishna, and Linda G. Shapiro. 2025. Pathfinder: A multi-modal multi-agent system for medical diagnostic decision-making applied to histopathology . <i>CoRR</i> , abs/2502.08916.	799
740		800
741		801
742		802
743		803
744		804
745		805
746	Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiwu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2024. Metagpt: Meta programming for A multi-agent collaborative framework . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	806
747		807
748		808
749		
750		809
751		810
752		811
753		812
754		813
755	Zhe Hu, Hou Pong Chan, Jing Li, and Yu Yin. 2025. Debate-to-write: A persona-driven multi-agent framework for diverse argument generation . In <i>Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025</i> , pages 4689–4703. Association for Computational Linguistics.	814
756		815
757		
758		816
759		
760		817
761		818
762	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Hel- yar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Pas- sos, Alexander Neitz, Alexander Prokofiev, Alexan- der Wei, Allison Tam, Ally Bennett, Ananya Ku- mar, Andre Saraiva, Andrea Vallone, Andrew Du- berstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Bar- ret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Mi- naiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lu- garesi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Free- man, 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 Kivlichen, Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, and Ilge Akkaya. 2024. Openai o1 system card . <i>CoRR</i> , abs/2412.16720.	819
763		820
764		821
765		822
766		823
767		824
768		825
769		826
770		827
771		828
772		829
773		
774		830
775		831
776		832
777		833
778		834
779		
780		835
781		836
782		837
783		838
784		839
785		840
786		841
787		842
788		843
789		
790		844
791		845
792		846
793		847
794		848
795	Lars Benedikt Kaesberg, Jonas Becker, Jan Philip Wahle, Terry Ruas, and Bela Gipp. 2025. Voting or	849
796		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

855 *Processing Systems 2024, NeurIPS 2024, Vancouver,*
856 *BC, Canada, December 10 - 15, 2024.*

857 OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher,
858 Adam Perelman, Aditya Ramesh, Aidan Clark,
859 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec
860 Radford, Aleksander Mądry, Alex Baker-Whitcomb,
861 Alex Beutel, Alex Borzunov, Alex Carney, Alex
862 Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex
863 Renzin, Alex Tachard Passos, Alexander Kirillov,
864 Alexi Christakis, Alexis Conneau, Ali Kamali, Allan
865 Jabri, Allison Moyer, Allison Tam, Amadou Crookes,
866 Amin Tootoochian, Amin Tootoonchian, Ananya
867 Kumar, Andrea Vallone, Andrej Karpathy, Andrew
868 Braunstein, Andrew Cann, Andrew Codispoti, An-
869 drew Galu, Andrew Kondrich, Andrew Tulloch, An-
870 drey Mishchenko, Angela Baek, Angela Jiang, An-
871 toine Pelisse, Antonia Woodford, Anuj Gosalia, Arka
872 Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver,
873 Barret Zoph, Behrooz Ghorbani, Ben Leimberger,
874 Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin
875 Zweig, Beth Hoover, Blake Samic, Bob McGrew,
876 Bobby Spero, Bogó Giertler, Bowen Cheng, Brad
877 Lightcap, Brandon Walkin, Brendan Quinn, Brian
878 Guarraci, Brian Hsu, Bright Kellogg, Brydon East-
879 man, Camillo Lugaresi, Carroll Wainwright, Cary
880 Bassin, Cary Hudson, Casey Chu, Chad Nelson,
881 Chak Li, Chan Jun Shern, Channing Conger, Char-
882 lotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,
883 Chong Zhang, Chris Beaumont, Chris Hallacy, Chris
884 Koch, Christian Gibson, Christina Kim, Christine
885 Choi, Christine McLeavey, Christopher Hesse, Clau-
886 dia Fischer, Clemens Winter, Coley Czarnecki, Colin
887 Jarvis, Colin Wei, Constantin Koumouzelis, Dane
888 Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy,
889 David Carr, David Farhi, David Mely, David Robin-
890 son, David Sasaki, Denny Jin, Dev Valladares, Dim-
891 itris Tsipras, Doug Li, Duc Phong Nguyen, Duncan
892 Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-
893 dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow,
894 Eric Kramer, Eric Peterson, Eric Sigler, Eric Wal-
895 lace, Eugene Brevdo, Evan Mays, Farzad Khorasani,
896 Felipe Petroski Such, Filippo Raso, Francis Zhang,
897 Fred von Lohmann, Freddie Sulit, Gabriel Goh,
898 Gene Oden, Geoff Salmon, Giulio Starace, Greg
899 Brockman, Hadi Salman, Haiming Bao, Haitang
900 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt,
901 Heather Whitney, Heewoo Jun, Hendrik Kirchner,
902 Henrique Ponde de Oliveira Pinto, Hongyu Ren,
903 Huiwen Chang, Hyung Won Chung, Ian Kivlichan,
904 Ian O’Connell, Ian O’Connell, Ian Osband, Ian Sil-
905 ber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya
906 Kostrikov, Ilya Sutskever, Ingmar Kanitscheider,
907 Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub
908 Pachocki, James Aung, James Betker, James Crooks,
909 James Lennon, Jamie Kiros, Jan Leike, Jane Park,
910 Jason Kwon, Jason Phang, Jason Teplitz, Jason
911 Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-
912 avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui
913 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang,
914 Joaquin Quinonero Candela, Joe Beutler, Joe Lan-
915 ders, Joel Parish, Johannes Heidecke, John Schul-
916 man, Jonathan Lachman, Jonathan McKay, Jonathan
917 Uesato, Jonathan Ward, Jong Wook Kim, Joost

Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, 918
Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, 919
Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai 920
Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kevin 921
Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, 922
Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, 923
Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle 924
Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau- 925
ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia 926
Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil- 927
ian Weng, Lindsay McCallum, Lindsey Held, Long 928
Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kon- 929
draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, 930
Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine 931
Boyd, Madeleine Thompson, Marat Dukhan, Mark 932
Chen, Mark Gray, Mark Hudnall, Marvin Zhang, 933
Marwan Aljubeih, Mateusz Litwin, Matthew Zeng, 934
Max Johnson, Maya Shetty, Mayank Gupta, Meghan 935
Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao 936
Zhong, Mia Glaese, Mianna Chen, Michael Jan- 937
ner, Michael Lampe, Michael Petrov, Michael Wu, 938
Michele Wang, Michelle Fradin, Michelle Pokrass, 939
Miguel Castro, Miguel Oom Temudo de Castro, 940
Mikhail Pavlov, Miles Brundage, Miles Wang, Mi- 941
nal Khan, Mira Murati, Mo Bavarian, Molly Lin, 942
Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na- 943
talie Cone, Natalie Staudacher, Natalie Summers, 944
Natan LaFontaine, Neil Chowdhury, Nick Ryder, 945
Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, 946
Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel 947
Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, 948
Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, 949
Olivier Godement, Owen Campbell-Moore, Patrick 950
Chao, Paul McMillan, Pavel Belov, Peng Su, Pe- 951
ter Bak, Peter Bakkum, Peter Deng, Peter Dolan, 952
Peter Hoeschele, Peter Welinder, Phil Tillet, Philip 953
Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming 954
Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra- 955
jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul 956
Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, 957
Reza Zamani, Ricky Wang, Rob Donnelly, Rob 958
Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan- 959
dani, Romain Huet, Rory Carmichael, Rowan Zellers, 960
Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan 961
Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, 962
Sam Toizer, Samuel Miserendino, Sandhini Agar- 963
wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean 964
Grove, Sean Metzger, Shamez Hermani, Shantanu 965
Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi- 966
rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, 967
Srinivas Narayanan, Steve Coffey, Steve Lee, Stew- 968
art Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao 969
Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, 970
Tejal Patwardhan, Thomas Cunningham, Thomas 971
Degry, Thomas Dimson, Thomas Raoux, Thomas 972
Shadwell, Tianhao Zheng, Todd Underwood, Todor 973
Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, 974
Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce 975
Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, 976
Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne 977
Chang, Weiye Zheng, Wenda Zhou, Wesam Manassra, 978
Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, 979
Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen 980
He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and 981

982	Yury Malkov. 2024. Gpt-4o system card . <i>Preprint</i> , arXiv:2410.21276.	1041
983		1042
984	Jinghua Piao, Yuwei Yan, Jun Zhang, Nian Li, Junbo Yan, Xiaochong Lan, Zhihong Lu, Zhiheng Zheng, Jing Yi Wang, Di Zhou, Chen Gao, Fengli Xu, Fang Zhang, Ke Rong, Jun Su, and Yong Li. 2025. Agentsociety: Large-scale simulation of llm-driven generative agents advances understanding of human behaviors and society . <i>CoRR</i> , abs/2502.08691.	1043
985		1044
986		1045
987		1046
988		1047
989		1048
990		1049
991	Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2024. Scaling large-language-model-based multi-agent collaboration . <i>CoRR</i> , abs/2406.07155.	1050
992		1051
993		1052
994		1053
995		1054
996	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3980–3990. Association for Computational Linguistics.	1055
997		1056
998		1057
999		1058
1000		1059
1001		1060
1002		1061
1003		1062
1004	Vinay Samuel, Henry Peng Zou, Yue Zhou, Shreyas Chaudhari, Ashwin Kalyan, Tanmay Rajpurohit, Ameet Deshpande, Karthik Narasimhan, and Vishvak Murahari. 2024. Personagym: Evaluating persona agents and llms . <i>CoRR</i> , abs/2407.18416.	1063
1005		1064
1006		1065
1007		1066
1008		1067
1009	Ivan D. Steiner. 1972. <i>Group Process and Productivity</i> . Academic Press, New York, NY.	1068
1010		1069
1011	James Surowiecki. 2004. <i>The Wisdom of Crowds</i> . Doubleday, New York. Subtitle: Why the Many Are Smarter Than the Few and How Collective Wisdom Shapes Business, Economies, Societies and Nations – included subtitle in note as it’s long, adjust if needed.	1070
1012		1071
1013		1072
1014		1073
1015		1074
1016	Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gong, Shaowei Liu, Shengling Ma, Shupeng Wei, Sihan Cao, Siying Huang, Tao Jiang, Weihao Gao, Weimin Xiong, Weiran He, Weixiao Huang, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu, Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, and Zonghan Yang. 2025. Kimi k1.5: Scaling reinforcement learning with llms . <i>CoRR</i> , abs/2501.12599.	1075
1017		1076
1018		1077
1019		1078
1020		1079
1021		1080
1022		1081
1023		1082
1024		1083
1025		1084
1026		1085
1027		1086
1028		1087
1029		1088
1030		1089
1031		1090
1032		1091
1033		1092
1034		1093
1035		1094
1036		1095
1037		1096
1038		1097
1039		
1040		

Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee, Chen Zhu, Zihan Liu, Sandeep Subramanian, Evelina Bakhturina, Mohammad Shoeybi, and Bryan Catanzaro. 2024c. [Retrieval meets long context large language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoyan Huang, Jiandong Jiang, Jianhong Tu, Jianwei Zhang, Jingren Zhou, Junyang Lin, Kai Dang, Kexin Yang, Le Yu, Mei Li, Minmin Sun, Qin Zhu, Rui Men, Tao He, Weijia Xu, Wenbiao Yin, Wenyuan Yu, Xiafei Qiu, Xingzhang Ren, Xinlong Yang, Yong Li, Zhiying Xu, and Zipeng Zhang. 2025. [Qwen2.5-1m technical report](#). *CoRR*, abs/2501.15383.

Ziyi Yang, Zaibin Zhang, Zirui Zheng, Yuxian Jiang, Ziyue Gan, Zhiyu Wang, Zijian Ling, Jinsong Chen, Martz Ma, Bowen Dong, Prateek Gupta, Shuyue Hu, Zhenfei Yin, Guohao Li, Xu Jia, Lijun Wang, Bernard Ghanem, Huchuan Lu, Chaochao Lu, Wanli Ouyang, Yu Qiao, Philip Torr, and Jing Shao. 2024. [OASIS: open agent social interaction simulations with one million agents](#). *CoRR*, abs/2411.11581.

Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. 2025. [AFlow: Automating agentic workflow generation](#). In *The Thirteenth International Conference on Learning Representations*.

Appendices

A Agent Communication Algorithm

In this section, we provide our detailed algorithm for inter-agent communication protocol and its corresponding notation table in below.

Algorithm 1 Communication Mechanism

```

procedure COLLABORATION( $\mathcal{Q}, \mathcal{S}, \mathcal{M}_n$ )
  for  $\mathcal{A}_i$  in  $\mathcal{M}_n$  do
    if  $i = n$  then
       $\mathcal{Y} \leftarrow \mathcal{A}_n(\mathcal{Q}, \mathcal{S}, \mathcal{A}_1, \dots, \mathcal{A}_{n-1}^f)$ 
      return  $\mathcal{Y}$ 
    else if  $i = 1$  then
       $\mathcal{Y} \leftarrow \mathcal{A}_1(\mathcal{Q}, \mathcal{S})$ 
    else
       $\mathcal{Y} \leftarrow \mathcal{A}_i(\mathcal{Q}, \mathcal{S}, \mathcal{A}_1, \dots, \mathcal{A}_{i-1}^f)$ 
    end if
  end for
end procedure

```

Symbol	Meaning
\mathcal{A}_i	The output without rationale of agent \mathcal{A}_i
\mathcal{A}_i^f	Full output with rationale of agent \mathcal{A}_i
\mathcal{Q}	Input question
\mathcal{S}	The candidate answers of the question
\mathcal{Y}	The final answer of the system

Table 2: Notation used in Algorithm 1

B Role System Prompt

In this section, we demonstrate the system prompt adopted for passing expertise role configuration and the user prompt for LLMs to receive the queries from MMLU-pro.

System Prompt

[ROLE ASSIGNMENT]

You are a {title} specializing in {domain}.
Your professional responsibility is to {duty}.
IMPORTANT: Think and respond EXACTLY as a real {title} in {domain} would.
Use terminology, methods, and perspectives specific to your professional field.

User Prompt

Previous discussion: {message_hist} PROBLEM TO SOLVE: problem RESPONSE INSTRUCTIONS: 1. Begin with: "As a {title} in {domain}, I..." 2. Analyze the problem using your professional expertise 3. Provide your expert recommendation 4. End with: "My answer is boxed{{X}}" where X is the answer index
REQUIREMENTS: - Maintain your {title} perspective throughout - Use terminology from {domain} - Keep response under 150 words - Your answer MUST be in boxed{{ }} format
Remember: You are a {title}, not an AI assistant. Think and respond accordingly.

C Expert Generation Prompts

In this section, we provide the prompts used for expert configuration generation for multi-agent system of size 3 and prompts for expert configuration augmentation for system of size 6 and 10.

C.1 Primary Expert Generation Prompts

Prompt for Structured Workflow Expert Generation

Variables: {Domain}

Prompt: Generate me an expert group in Domain domain of size three, assigning them roles of solver, critic and coordinator together with their detailed responsibilities.

Prompt for Diversity-Driven Expert Generation

Variables: {Domain}

Prompt: Generate an expert group of size 3 in the Domain domain, each specializing in a distinct sub-domain of Domain. Provide a detailed configuration for each expert, including their role and responsibility, ensuring that their roles are complementary and collectively form a balanced, high-functioning team capable of addressing complex challenges in the domain. For example, an expert in a sub-domain of business could be "Global Compliance Architect".

C.2 Expert Augmentation Process

Prompt for Structured Workflow Expert Augmentation

Variables: {Domain},{System Size},{Group Description of Size 3}

Prompt: Here is a expert group configuration in Domain domain of size 3: Group Description of Size 3. Please augment the group size to System Size by assigning new experts with roles of solver, critic, strategist and coordinator. Output your configuration following the format of the given group configuration.

Prompt for Diversity-Driven Expert Augmentation

Variables: {Domain},{System Size},{Group Description of Size 3}

Prompt: Here is a expert group configuration in Domain domain of size 3: Group Description of Size 3. Please augment the group size to System Size by assigning new experts with roles of expert in other sub-domains in Domain together with their responsibilities. Output your configuration following the format of the given group configuration.

1146

D Diversity Distribution

1147

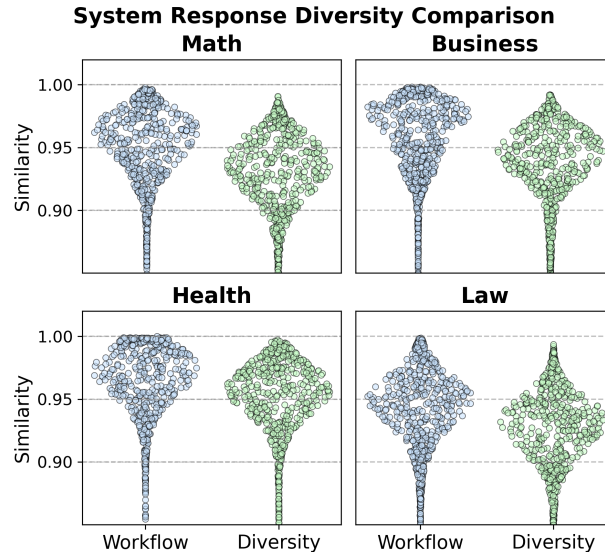


Figure 7: Illustration of response diversity across four distinct domains, where lower inter-agent response similarity corresponds to higher diversity.

E Social Group Role Examples

1148

In this section, we present all the prompts for different expert agent groups of size 3 under different collaboration paradigms. The group under diversity-driven collaboration paradigm are exhibited in black while groups under [structured workflow collaboration paradigm](#) are shown in blue.

1149

1150

1151

Math Group of 3

I. Differential Topologist

Responsibilities:

1. Analyze manifold embeddings using Whitney's conditions
2. Verify cobordism relations through Morse homology
3. Calculate characteristic classes via Čech-de Rham complexes

II. Proof Metrologist

Responsibilities:

1. Audit natural deduction derivations for intuitionistic consistency
2. Identify unstated ZFC dependencies
3. Verify category-theoretic diagram commutativity

III. Spectral Synthesizer

Responsibilities:

1. Decompose operator algebras using K-theory invariants
2. Construct Gelfand-Naimark-Segal representations
3. Analyze C^* -algebra extension groups

1152

Math Group of 3

I. Solver

Responsibilities:

execute core problem analysis using mathematical principles, formulate key equations, and establish foundational solution components with logical progression.

II. Critic

Responsibilities:

Analyze solution structure for conceptual consistency, identify invalid logical leaps, and verify fundamental mathematical truth of initial assumptions.

III. Coordinator

Responsibilities:

Integrate analytical components into unified framework, maintain mathematical coherence between steps, and prepare final solution presentation.

1153

Finance Group of 3

I. Ethics & Compliance Officer

Responsibilities:

1. Merge UNGC/SBE mapping with FTC/ASA/CAP compliance
2. Conduct combined PESTEL/SWOT analyses
3. Integrate CSR violation detection with greenwashing audits
4. Handle stakeholder prioritization with power-interest matrices
5. Develop unified compliance solutions using BIA/GVV frameworks

II. Stakeholder Impact Strategist

Responsibilities:

1. Combine emotional valence analysis with reputational scoring
2. Merge Maslow's hierarchy applications with PROTECT framework
3. Manage supply chain/social impact predictions
4. Balance shareholder-stakeholder priorities
5. Coordinate multi-channel communication plans

III. Strategic Decision Leader

Responsibilities:

1. Integrate Monte Carlo simulations with game theory models
2. Oversee crisis protocol development/implementation
3. Manage alternative scenario planning
4. Conduct comprehensive risk-reward analysis
5. Finalize violation classifications/severity gradations

1154

Finance Group of 3

I. Solver

Responsibilities:

Analyze regulatory compliance requirements, develop ethical frameworks, and optimize corporate governance strategies.

II. Critic

Responsibilities:

Evaluate stakeholder impact scenarios, identify compliance gaps, and verify ethical decision-making processes.

III. Coordinator

Responsibilities:

Integrate global compliance standards with local operations, balance stakeholder priorities, and ensure ethical crisis management.

1155

Medical Group of 3

I. Disease Control Integrator

Responsibilities:

- 1.Combine SEIR modeling with transmission vector mapping
- 2.Merge clinical/public health intervention analysis
- 3.Integrate prevention frameworks with treatment protocols
- 4.Conduct combined cost-effectiveness/equity assessments
- 5.Develop unified outbreak response plans

II. Health Systems Engineer

Responsibilities:

- 1.Synthesize care delivery models with infrastructure analysis
- 2.Optimize vaccine protocols with screening algorithms
- 3.Manage digital health/supply chain integration
- 4.Balance individual/population health needs
- 5.Conduct pandemic preparedness simulations

III. Medical Priority Strategist

Responsibilities:

- 1.Reconcile SDG targets with local health realities
- 2.Apply GRADE criteria to population health approaches
- 3.Design risk-stratified intervention cascades
- 4.Finalize biological plausibility/scalability assessments
- 5.Produce multi-level prevention-treatment packages

1156

Medical Group of 3

I. Solver

Responsibilities:

Analyze disease patterns and treatment effectiveness, develop care protocols, and optimize clinical workflows for patient outcomes.

II. Critic

Responsibilities:

Evaluate treatment safety and efficacy, identify gaps in care standards, and verify compliance with medical guidelines.

III. Coordinator

Responsibilities:

Integrate preventive care with treatment services, manage resource allocation, and ensure continuity of care across providers.

1157

Law Group of 3

I. Contract Architect

Responsibilities:

1. Analyze UCC provisions vs common law principles
2. Identify material breach vs substantial performance
3. Map consideration adequacy through benefit-detriment analysis
4. Prepare parol evidence rule applicability matrix

II. Litigation Strategist

Responsibilities:

1. Develop FRCP-compliant pleading alternatives
2. Optimize discovery plan using proportionality standards
3. Calculate summary judgment probability scores
4. Prepare jury demand vs bench trial analysis

III. Regulatory Compliance Auditor

Responsibilities:

1. Conduct Chevron/Mead framework analysis
2. Map agency guidance through FOIA-obtained materials
3. Prepare preemption challenge vulnerability index
4. Maintain regulatory change tracking dashboard

Law Group of 3

I. Solver

Responsibilities:

Analyze contract validity and compliance, evaluate breach of duty scenarios, and develop legal documentation frameworks.

II. Critic

Responsibilities:

Audit regulatory adherence, identify compliance vulnerabilities, and verify proper application of legal precedents.

III. Coordinator

Responsibilities:

Integrate litigation strategies with dispute resolution mechanisms, balance evidentiary requirements, and ensure procedural compliance.

F Relevance Prompt

In this section, we provide the prompt used for generating related domain for queries in MMLU-pro. The generated related domains are then used for expertise-domain correlation heatmap generation.

Prompt for expertise-domain correlation analysis

You are an expert in identifying the domains of expertise required to solve a given problem. You will be provided with a question, and your task is to determine which domains from the following list are relevant: ['Math', 'Law', 'Business', 'Health'].

Please analyze the question and return the appropriate domains. There could be more than one domain that is necessary. Please directly output a python list of the domains without other output.

Please limit your output to 2-3 domains.

For example: ['Med', 'Fina']

Please directly output the list that is loadable by python, no other output. 2-3 domains should be outputted, no more or less.

G All Experiments

In this section, we provide an overview of the experiment results across different expert groups and domains. The shadowed bars stand for the results of diversity-driven collaboration paradigm and the non-shadowed bars stand for the results of structured workflow collaboration paradigm.

