

# 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

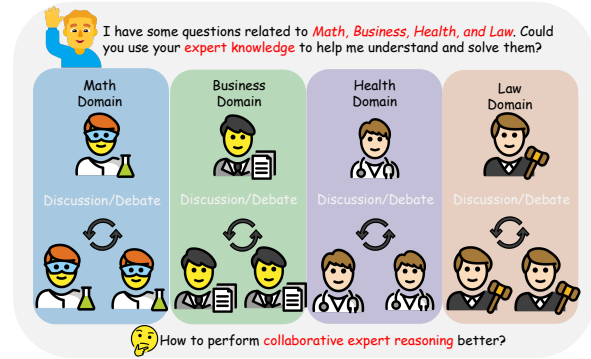


Figure 1: Workflow diagram for a multi-agent reasoning system with specialized agents.

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.

069	We first focus on expertise-domain align-	<b>cover non-linear dynamics; specifically, adding</b>	121
070	ment—where existing practice reveals significant	<b>more experts tends to improve the collective rea-</b>	122
071	gaps. Although expertise specialization is widely	<b>soning ability of the system. This positive trend</b>	123
072	adopted in multi-agent systems (Wang et al., 2024a;	<b>holds regardless of whether the larger system</b>	124
073	Li et al., 2024a), the impact of collaborative ex-	<b>scale contains greater viewpoint diversity or a</b>	125
074	pertise configuration on downstream scenarios re-	<b>more comprehensive workflow structure,</b> indi-	126
075	remains underexplored. This ambiguity creates prac-	icating a general benefit to increasing the number	127
076	tical challenges in configuring expert roles for task	of expert agents and encouraging such designs for	128
077	domains. To address this gap, we empirically evalu-	enhanced system performance. However, our anal-	129
078	ate the influence of different collaborative expertise	ysis of the computational trade-offs associated with	130
079	configurations on task performance across four rep-	system scaling reveals that, while the system would	131
080	resentative domains from MMLU-pro (Wang et al.,	benefit from the expansion, there remains a critical	132
081	2024d). <b>Our findings in Section 4 demonstrate a</b>	need for more efficient communication protocols	133
082	<b>positive correlation between task performance</b>	between agents for more scalable and cost-effective	134
083	<b>and the alignment of group expertise with the</b>	multi-agent reasoning process.	135
084	<b>task domain, underscoring the necessity of ac-</b>		
085	<b>curately matching the multi-agent system exper-</b>	<b>2 Related Works</b>	136
086	<b>tise with downstream tasks.</b>		
087	Having established the critical role of expertise-	<b>2.1 Multi-Agent Collaboration</b>	137
088	domain alignment, we then examine how collabor-	Multi-Agent Collaboration adopts multiple LLMs	138
089	ation paradigms structure interactions between	to solve the problem collaboratively. Abundant	139
090	specialized agents. Currently, the collaboration	researches have investigated the multi-agent col-	140
091	paradigm predominantly used in recent studies	laboration framework to improve decision-making	141
092	could be categorized into two kinds: (1) Diversity-	capability of the system (Wang et al., 2024b; Liang	142
093	Driven Perspective Integration, where agents, often	et al., 2024; Du et al., 2024). In addition to collabor-	143
094	embodying different viewpoints or roles, are en-	oration among LLMs, several researchers instruct	144
095	couraged to generate diverse responses to enrich	the agents to cooperate in a workflow to study the	145
096	the solution space (Wang et al., 2024b; Chen et al.,	multi-agent systems’ ability of solving real world	146
097	2024b; Hu et al., 2025). (2) Structured Workflow	challenges (Li et al., 2024b; Xu et al., 2024b; Chen	147
098	Cooperation, where different agents are assigned	et al., 2024a). While Qian et al. (2024), Yang et al.	148
099	distinct sub-tasks within a predefined pipeline to	(2024) and Wang et al. (2024c) has investigate the	149
100	collaboratively construct a solution (Chen et al.,	effect of varying the scale of multi-agent system	150
101	2024c; Hong et al., 2024; Zhang et al., 2025). We	on reasoning and simulation, prior researches have	151
102	design comparative experiments to unveil the per-	not systematically examined the interplay between	152
103	formance differences between paradigms. <b>Our</b>	collective expertise specialization, collaboration	153
104	<b>observations in Section 5 reveal a consistent ad-</b>	mechanisms, and the impact of system scale simul-	154
105	<b>vantage for diversity-driven collaboration over</b>	taneously. In this work, we conduct extensive ex-	155
106	<b>structured workflow collaboration, suggesting</b>	periments to formally analyze the influence of these	156
107	<b>the superiority of the diversity-driven paradigm.</b>	three critical dimensions on multi-agent collabor-	157
108	Finally, constructing large-scale multi-agent sys-	ative reasoning. Our findings provide actionable	158
109	tem has become a critical, yet often enigmatic as-	insights toward more effective system design.	159
110	pect of multi-agent system design (Chen et al.,		
111	2024c; Piao et al., 2025). While intuition and	<b>2.2 LLMs as Domain Experts</b>	160
112	some preliminary studies (Qian et al., 2024; Li	The rapid evolution of LLMs has endowed them	161
113	et al., 2023) suggest that larger groups would lead	with vast repositories of domain-specific knowl-	162
114	to better reasoning performance, the actual effec-	edge, enabling their application across a wide	163
115	tiveness of scaling within the context of collabor-	range of expert tasks. Recent researches have ex-	164
116	orative expertise specialization and the potential	explored the potential of LLMs to emulate specific	165
117	computation-performance trade-off, are not well	personas by conditioning them on detailed char-	166
118	understood. Our systematic experiments involve in-	acter profiles (Chan et al., 2024; Samuel et al.,	167
119	crementally increasing the system scale to examine	2024; Xu et al., 2023). These studies demon-	168
120	potential scaling laws. <b>The results in Section 6 un-</b>	strate that by providing LLMs with demographic or	169

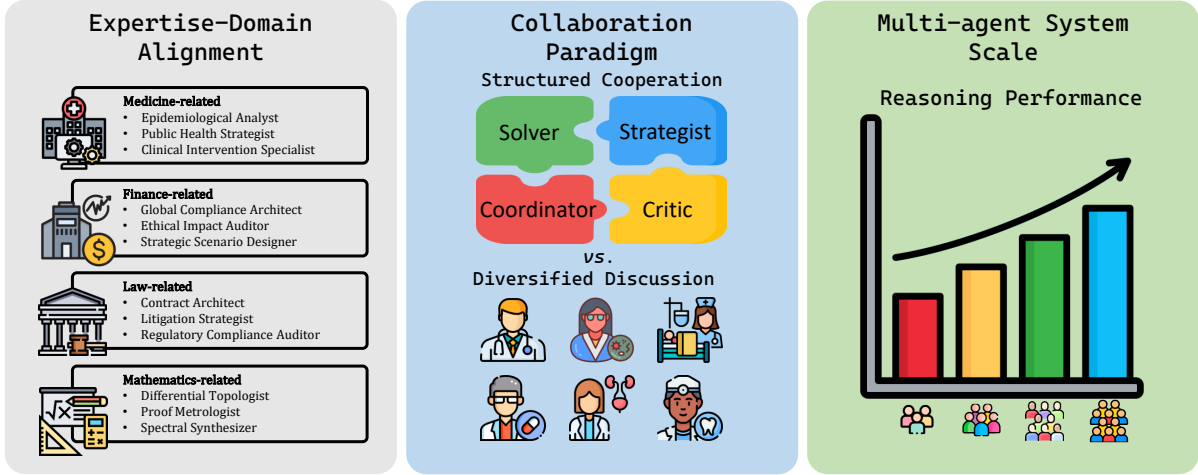


Figure 2: Demonstration of three key factors characterizing research on multi-agent collaborative reasoning systems. (1) expertise-domain alignment, (2) collaboration paradigm, and (3) scale of the multi-agent system.

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

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

### 240 3.3 Collaborative Expertise Specialization

241 In this paper, we primarily studied the effect of col-  
242 laborative expertise specialization on better multi-  
243 agent system design from the perspective of expert-  
244 domain alignment, collaboration paradigms and  
245 system scale. To formalize the role and responsi-  
246 bility of the agents in the multi-agent system, we  
247 define each expert to be of the following format:

$$248 \mathcal{A}_i \leftarrow (EG, FR, R, ID)$$

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

### 253 3.4 General Experiment Setup

254 As detailed in Section 3.2, we select 4 represen-  
255 tative domains from MMLU-pro to investigate  
256 the effects of collaborative expertise specializa-  
257 tion. To be consistent with all experiments, we  
258 utilize DeepSeek-R1-Distill-Qwen-7B (DeepSeek-  
259 AI et al., 2025) as the foundational model for all  
260 agents. Each agent is initialized with its specific  
261 expert description and responsibilities via its sys-  
262 tem prompt, while the task instance is provided  
263 through the user prompt. The detailed prompts  
264 could be found in Appendix B. All experiments  
265 adopt accuracy as the evaluation metric.

## 266 4 Leveraging the “Right” Agent

267 Expertise specialization is a widely adopted tech-  
268 nique in agent research, demonstrably enhancing  
269 the reasoning capabilities of LLMs within specific  
270 domains (Li et al., 2024b). While the benefits  
271 of specialization for individual agents are well-  
272 established, the effect of collaborative expertise  
273 specialization on the collective reasoning perfor-  
274 mance of multi-agent systems remains underex-  
275 plored. This section presents our experimental in-  
276 vestigation into this critical area, designed to unveil  
277 how different collaborative expertise specialization  
278 configurations influence the reasoning capabilities  
279 of multi-agent systems.

### 280 4.1 Setup

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

283 agent viewpoints and integrate them in the final  
284 answer (Liang et al., 2024), in our experiments,  
285 we adopt diversity-driven collaboration paradigm  
286 where we distribute each agent with a specific do-  
287 main expert configuration and instruct them to gen-  
288 erate responses based on their expertise. At this  
289 stage, we fix the size of the multi-agent reason-  
290 ing system to be 3 for controllable computational  
291 cost. We employ GPT-4o (OpenAI et al., 2024)  
292 for expert configuration generation. The detailed  
293 prompts utilized for this automated role generation  
294 process are provided in Appendix C.

### 295 4.2 “Right” Expertise Helps Reasoning

296 **Our experiments demonstrate a clear perfor-**  
297 **mance advantage when the collaborative ex-**  
298 **pertise specialization of the multi-agent system**  
299 **aligns with the domains of the downstream task.**

300 Misaligned expertise configurations often under-  
301 perform compared to aligned ones. This primary  
302 finding is quantitatively supported by the results  
303 presented in Table 1. Specifically, in 75% of  
304 the aligned cases (diagonal entries), the system  
305 achieves the highest accuracy compared to configu-  
306 rations where the agent group simulates expertise  
307 from other domains for the same task.

308 To gain a more nuanced understanding of when  
309 expertise alignment is most beneficial, we analyze  
310 the system performance according to the primary  
311 reasoning type required by each domain, as cate-  
312 gorized in Section 3.2. Our analysis reveals that  
313 the benefits of expertise alignment are most pro-  
314 nounced for tasks demanding contextual reason-  
315 ing—Health and Law. Systems operating on these  
316 two domains exhibit an average relative perfor-  
317 mance improvement of 6.75% when expertise is  
318 correctly aligned, compared to the misaligned con-  
319 figurations which perform the second best for those  
320 tasks. Conversely, for domains requiring mathemat-  
321 ical reasoning—Math and Business, the specialized  
322 experts yield only marginal gains or even degrada-  
323 tion relative to misaligned configurations. We hy-  
324 pothesize this divergence stems from the inherent  
325 strengths of LLMs on math. These models often  
326 possess robust mathematical reasoning capabilities  
327 due to extensive pre-training, potentially reducing  
328 the added value of specialized agents. Contextual  
329 reasoning tasks, however, appear to benefit more  
330 from the structured integration of specialized per-  
331 spectives provided by the multi-agent reasoning  
332 system since applying domain knowledge in these  
333 contexts often requires nuanced interpretation, syn-

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

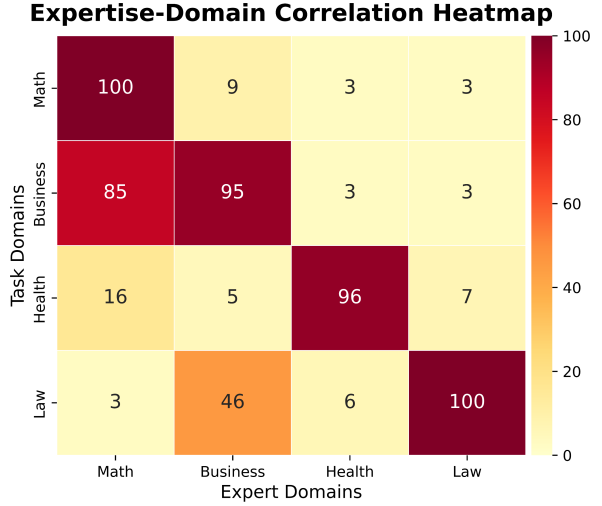


Figure 3: Heatmap illustrating the correlation between specialized group expertise and task domains. Deeper colors indicate stronger correlations.

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

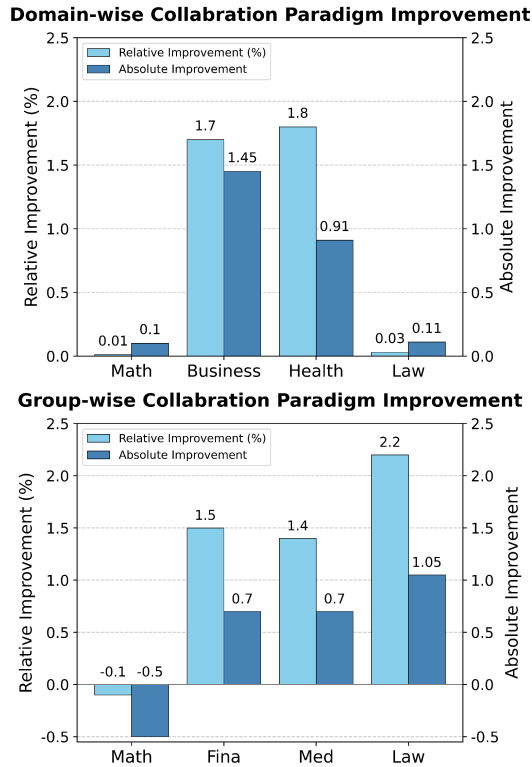


Figure 4: Comparative analysis of diversity-driven versus structured workflow collaboration paradigms. Positive values signify Diversity-Driven’s advantage over Structured Workflow.

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-

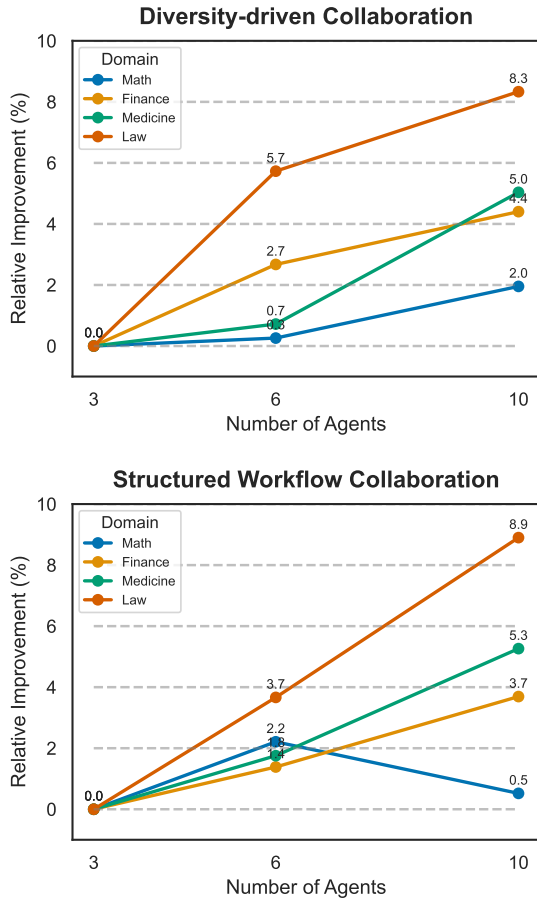


Figure 5: Domain-wise relative performance improvement by scaling up the multi-agent system (3, 6, and 10 agents), shown for different collaboration mechanisms.

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**

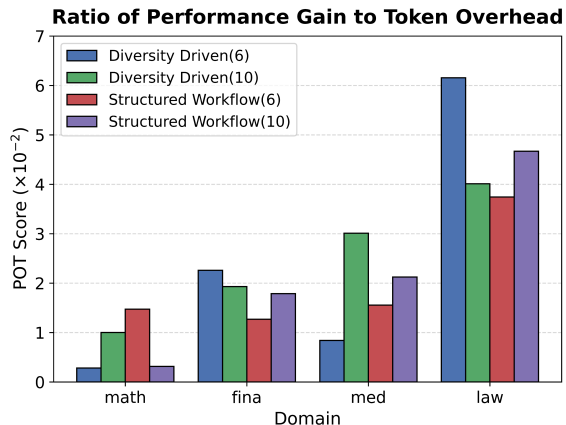


Figure 6: Performance improvement versus token overhead ratio across different domains. Both performance and token overhead are measured as relative increases compared to the system of size 3.

**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.

## 629 Limitations

630 Our adoption of MMLU-pro for evaluating special-  
631 ized multi-agent reasoning system across diverse  
632 domains, while leveraging its strength in assess-  
633 ing varied domain-specific knowledge, inherently  
634 limits our assessment scope. Specifically, its focus  
635 on these reasoning paradigms means other crucial  
636 multi-agent capabilities, such as coding, might be  
637 overlooked. Apart from that, to enhance align-  
638 ment with real-world scenarios, our evaluation con-  
639 centrates on four key domains: Math, Business,  
640 Health, and Law, selected for their prominence  
641 in mainstream research. A direct limitation of  
642 this focused approach is that other potentially rele-  
643 vant domains would remain underexplored in the  
644 present study. Moreover, To simplify the research  
645 setup and promote more stable conclusions, we ex-  
646 clusively utilize one message propagation mecha-  
647 nism. This methodological choice, however, means  
648 that the potential influence of diverse communi-  
649 cation strategies on system performance remains  
650 an unexplored aspect in our current study. Finally,  
651 We select DeepSeek-R1-Distilled-Qwen-7B as the  
652 base model for all experiments to ensure control-  
653 lable computational overhead. This decision, while  
654 practical, limits our current investigation, deferring  
655 the study of multi-agent system architectures with  
656 larger-scale models to future research.

## 657 Ethics Statement

658 Our study involves publicly available datasets and  
659 use Large Language Models through APIs. Con-  
660 sequently, the ethical considerations of this paper  
661 could be listed as follow:

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

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

668 **Transparency:** We provide detailed descriptions  
669 of our method and the prompts used in our ex-  
670 periments, in line with standard practices in the  
671 research community. We will also make our code  
672 publicly available upon acceptance.

## 673 References

674 Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and  
675 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. 677 678 679 680 681

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. 682 683 684 685 686 687 688 689

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. 690 691 692 693 694 695 696 697 698

Jiayi 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. 699 700 701 702

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. 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730

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. 731 732 733 734 735 736

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. <a href="#">Pathfinder: A multimodal multi-agent system for medical diagnostic decision-making applied to histopathology</a> . <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, Xiawu 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. <a href="#">Metagpt: Meta programming for A multi-agent collaborative framework</a> . 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			809
750			810
751			811
752			812
753			813
754			814
755	Zhe Hu, Hou Pong Chan, Jing Li, and Yu Yin. 2025. <a href="#">Debate-to-write: A persona-driven multi-agent framework for diverse argument generation</a> . 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.		815
756			816
757			817
758			818
759			819
760			820
761			821
762	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helvar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Pasos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, 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 Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, and Ilge Akkaya. 2024. <a href="#">Openai o1 system card</a> . <i>CoRR</i> , abs/2412.16720.		822
763			823
764			824
765			825
766			826
767			827
768			828
769			829
770			830
771			831
772			832
773			833
774			834
775			835
776			836
777			837
778			838
779			839
780			840
781			841
782			842
783			843
784			844
785			845
786			846
787			847
788			848
789			849
790			850
791			851
792			852
793			853
794			854
795	Lars Benedikt Kaesberg, Jonas Becker, Jan Philip Wahle, Terry Ruas, and Bela Gipp. 2025. <a href="#">Voting or</a>		855
796			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, Bogo 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 Quinero 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, Kavin 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 Feувrier, 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 Aljubeһ, 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, Weiyi 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. <a href="#">Gpt-4o system card</a> . <i>Preprint</i> , arXiv:2410.21276.	
983		
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. <a href="#">Agentsociety: Large-scale simulation of llm-driven generative agents advances understanding of human behaviors and society</a> . <i>CoRR</i> , abs/2502.08691.	
988		
989		
990		
991	Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2024. <a href="#">Scaling large-language-model-based multi-agent collaboration</a> . <i>CoRR</i> , abs/2406.07155.	
992		
993		
994		
995		
996	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert: Sentence embeddings using siamese bert-networks</a> . 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.	
997		
998		
999		
1000		
1001		
1002		
1003		
1004	Vinay Samuel, Henry Peng Zou, Yue Zhou, Shreyas Chaudhari, Ashwin Kalyan, Tanmay Rajpurohit, Ameet Deshpande, Karthik Narasimhan, and Vishvak Murahari. 2024. <a href="#">Personagym: Evaluating persona agents and llms</a> . <i>CoRR</i> , abs/2407.18416.	
1005		
1006		
1007		
1008		
1009	Ivan D. Steiner. 1972. <i>Group Process and Productivity</i> . Academic Press, New York, NY.	
1010		
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.	
1012		
1013		
1014		
1015		
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.	
1017		
1018		
1019		
1020		
1021		
1022		
1023		
1024		
1025		
1026		
1027		
1028		
1029		
1030		
1031		
1032		
1033		
1034		
1035		
1036		
1037		
1038		
1039		
1040		
		<a href="#">Kimi k1.5: Scaling reinforcement learning with llms</a> . <i>CoRR</i> , abs/2501.12599.
		1041
		1042
	Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. 2024a. <a href="#">Mixture-of-agents enhances large language model capabilities</a> . <i>CoRR</i> , abs/2406.04692.	1043
		1044
		1045
		1046
	Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. 2024b. <a href="#">Rethinking the bounds of LLM reasoning: Are multi-agent discussions the key?</a> In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 6106–6131. Association for Computational Linguistics.	1047
		1048
		1049
		1050
		1051
		1052
		1053
		1054
	Ruiyi Wang, Haofei Yu, Wenxin Sharon Zhang, Zhengyang Qi, Maarten Sap, Yonatan Bisk, Graham Neubig, and Hao Zhu. 2024c. <a href="#">Sotopia-<math>\pi</math>: Interactive learning of socially intelligent language agents</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 12912–12940. Association for Computational Linguistics.	1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024d. <a href="#">Mmlu-pro: A more robust and challenging multi-task language understanding benchmark</a> . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
	Baixuan Xu, Weiqi Wang, Haochen Shi, Wenxuan Ding, Huihao Jing, Tianqing Fang, Jiabin Bai, Xin Liu, Changlong Yu, Zheng Li, Chen Luo, Qingyu Yin, Bing Yin, Long Chen, and Yangqiu Song. 2024a. <a href="#">MIND: multimodal shopping intention distillation from large vision-language models for e-commerce purchase understanding</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 7800–7815. Association for Computational Linguistics.	1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
	Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhendong Mao. 2023. <a href="#">Expertprompting: Instructing large language models to be distinguished experts</a> . <i>CoRR</i> , abs/2305.14688.	1085
		1086
		1087
		1088
	Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z. Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su, Leander Maben, Raj Mehta, Wayne Chi, Lawrence Keunho Jang, Yiqing Xie, Shuyan Zhou, and Graham Neubig. 2024b. <a href="#">Theagentcompany: Benchmarking LLM agents on consequential real world tasks</a> . <i>CoRR</i> , abs/2412.14161.	1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097

1098 Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee,  
1099 Chen Zhu, Zihan Liu, Sandeep Subramanian, Evelina  
1100 Bakhturina, Mohammad Shoeybi, and Bryan Catan-  
1101 zaro. 2024c. [Retrieval meets long context large lan-  
1102 guage models](#). In *The Twelfth International Con-  
1103 ference on Learning Representations, ICLR 2024,  
1104 Vienna, Austria, May 7-11, 2024*. OpenReview.net.

1105 An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei  
1106 Huang, Haoyan Huang, Jiandong Jiang, Jianhong Tu,  
1107 Jianwei Zhang, Jingren Zhou, Junyang Lin, Kai Dang,  
1108 Kexin Yang, Le Yu, Mei Li, Minmin Sun, Qin Zhu,  
1109 Rui Men, Tao He, Weijia Xu, Wenbiao Yin, Wenyuan  
1110 Yu, Xiafei Qiu, Xingzhang Ren, Xinlong Yang, Yong  
1111 Li, Zhiying Xu, and Zipeng Zhang. 2025. [Qwen2.5-  
1112 1m technical report](#). *CoRR*, abs/2501.15383.

1113 Ziyi Yang, Zaibin Zhang, Zirui Zheng, Yuxian Jiang,  
1114 Ziyue Gan, Zhiyu Wang, Zijian Ling, Jinsong Chen,  
1115 Martz Ma, Bowen Dong, Prateek Gupta, Shuyue  
1116 Hu, Zhenfei Yin, Guohao Li, Xu Jia, Lijun Wang,  
1117 Bernard Ghanem, Huchuan Lu, Chaochao Lu, Wanli  
1118 Ouyang, Yu Qiao, Philip Torr, and Jing Shao. 2024.  
1119 [OASIS: open agent social interaction simulations  
1120 with one million agents](#). *CoRR*, abs/2411.11581.

1121 Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng,  
1122 Xiong-Hui Chen, Jiaqi Chen, Mingchen Zhuge, Xin  
1123 Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng,  
1124 Bang Liu, Yuyu Luo, and Chenglin Wu. 2025.  
1125 [AFlow: Automating agentic workflow generation](#). In  
1126 *The Thirteenth International Conference on Learning  
1127 Representations*.

## Appendices

1128

### A Agent Communication Algorithm

1129

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

1130

1131

1132

---

#### 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

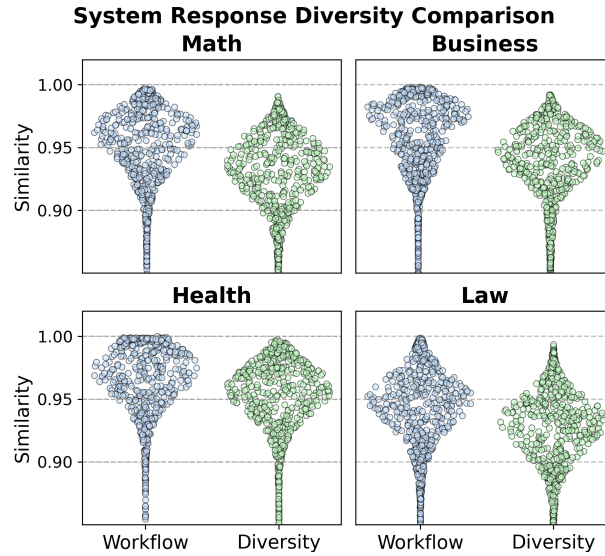


Figure 7: Illustration of response diversity across four distinct domains, where lower inter-agent response similarity corresponds to higher diversity.

## 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

1158

### 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.

1159

1160

## F Relevance Prompt

1161

1162

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.

1163