

Divide-and-Conquer Is What LLM-Based Multi-Agent System Need

Anonymous ACL submission

Abstract

Large language model (LLM) based multi-agent systems offer promising capabilities in social simulation and complex task solving, yet face key challenges in system design, generalizability, and scalability. We introduce AGENTGROUPCHAT-V2, a novel framework featuring: (1) a fully parallel divide-and-conquer architecture for efficient task decomposition and distributed processing; (2) an adaptive collaboration engine that dynamically selects heterogeneous LLMs and interaction strategies; (3) agent organization optimization for effective problem breakdown. Experiments show that AGENTGROUPCHAT-V2 achieves state-of-the-art results across several benchmarks, with substantial improvements on tasks such as GSM8K, AIME, and HumanEval, especially as task complexity increases. Our results demonstrate that AGENTGROUPCHAT-V2 enables the construction of robust and general-purpose LLM multi-agent systems, excelling in complex reasoning scenarios.

1 Introduction

Interest in multi-agent systems based on large language models (LLMs) has grown rapidly (Guo et al., 2024; Li et al., 2023; Xi et al., 2023; Liang et al., 2024), driven by their promise in domains such as social simulation (Gao et al., 2024; Park et al., 2023; Gu et al., 2024b) and complex task resolution (Hong et al., 2023; Wu et al., 2023). In social simulations, LLM agents provide new ways to model human-like interactions and study collective behavior (Gu et al., 2024b; Lee et al., 2023); for problem-solving, collaborative multi-agent networks enhance reasoning and planning capabilities through distributed intelligence (Huang et al., 2024; Wei et al., 2022; Du et al., 2023). Compared to single-agent systems, multi-agent architectures foster emergent strategies and superior collective intelligence (Ferber and Weiss, 1999; Dafoe et al., 2021; Gu et al., 2024b).

Despite their progress, current LLM multi-agent frameworks face major barriers regarding generalizability, scalability, and real-world utility (Ouyang et al., 2022; Cemri et al., 2025). Many popular systems target either social simulation (Gao et al., 2023; Gu et al., 2024b) or specific tasks such as software development (Hong et al., 2023; Wu et al., 2023), lacking methods to support diverse, complex problem types (Shoham and Leyton-Brown, 2008). Furthermore, they frequently adopt sequential execution patterns, or to say workflow (Sapkota et al., 2025; Zhang et al., 2025), resulting in inefficient computation and high resource overhead (Kumar, 2025). Notably, increased computational resources do not always yield better performance, and multi-agent approaches may sometimes underperform compared to single-agent baselines (Zhang et al., 2023; Cemri et al., 2025; Sapkota et al., 2025).

To overcome these limitations, it is important to design general and efficient collaboration mechanisms among agents. A promising direction is to rethink the organization of both tasks and agent interactions to fully exploit the benefits of multi-agent collaboration. In this context, we identify that the divide-and-conquer paradigm—which systematically decomposes complex tasks and distributes subtasks among specialized agents—is fundamental for improving both efficiency and adaptability in multi-agent systems (Shoham and Leyton-Brown, 2008; Ferber and Weiss, 1999). Through structured decomposition, responsibilities are clearly assigned, concurrency are maximized, and emergent behaviors are managed more effectively, enabling scalable and robust multi-agent intelligence.

To this end, we propose AGENTGROUPCHAT-V2, centered on three innovations: (1) a fully parallel architecture with hierarchical manager modules supporting scalable, distributed operation; (2) dynamic task-level divide-and-conquer via hierarchical task decomposition and dependency management; and (3) execution-level specialization

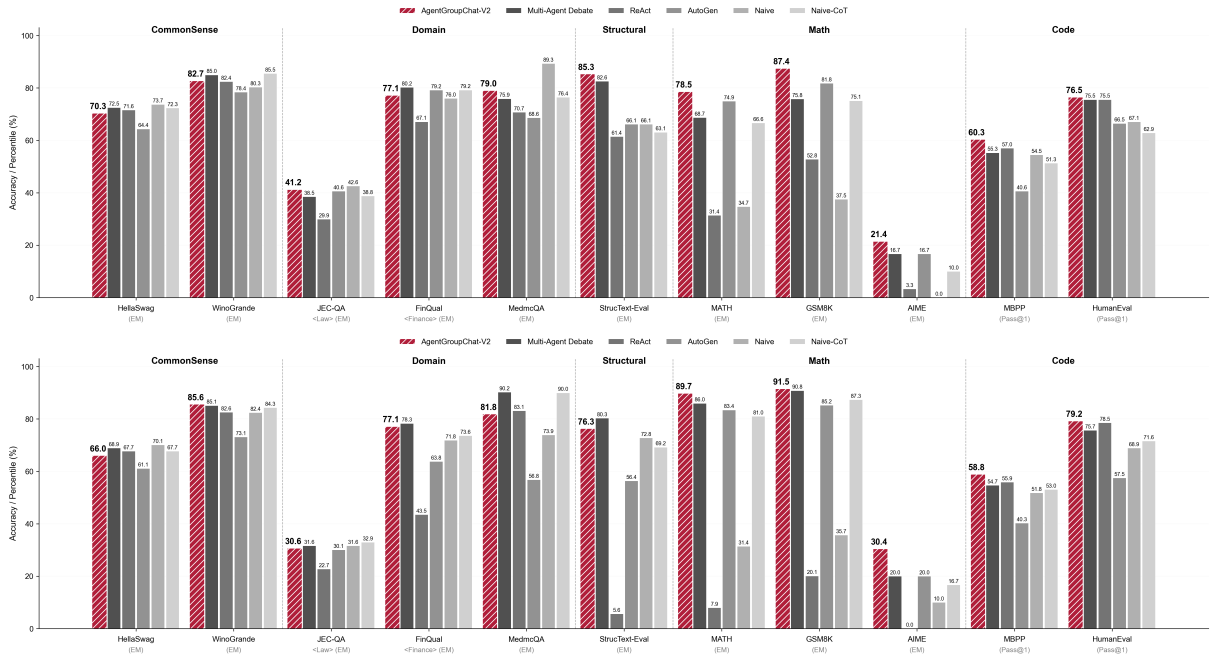


Figure 1: **Upper:** The performance of AGENTGROUPCHAT-V2 powered by Qwen2.5-72B. **Downer:** The performance of AGENTGROUPCHAT-V2 powered by Llama3.1-70B. Both models are evaluated across five diverse reasoning domains: commonsense reasoning, domain-specific knowledge, structural text understanding, mathematics, and code generation. AGENTGROUPCHAT-V2 consistently outperforms existing multi-agent approaches and baseline methods across both models and all benchmark categories, demonstrating the effectiveness, robustness and promising of multi-agent approach.

through adaptive role assignment to heterogeneous LLM agents. Our experiments confirm that AGENTGROUPCHAT-V2 significantly boosts performance and generalizability, especially on high-difficulty reasoning and generation tasks (see Figure 1).

2 Related Work

LLM-based multi-agent systems have demonstrated transformative potential in two key domains (Li et al., 2023; Wooldridge, 2009):

2.1 Social Dynamics Simulation

Recent advances integrate cognitive architectures with behavioral economics principles (Bates et al., 1994). Xie et al. (Xie et al., 2024a) and Han et al. (Han et al., 2023) develop game-theoretic frameworks capturing trust dynamics and market competition through iterative belief-updating mechanisms (Dafoe et al., 2021). Sociological simulations leverage event-driven architectures, with Park et al. (Park et al., 2023) modeling opinion evolution through social interaction cascades, while Zhang et al. (Zhang et al., 2024) employ hierarchical Bayesian networks for electoral behavior prediction (Lee et al., 2023). Emerging platforms like Gu’s group chat simulator (Gu et al., 2024b) and Liu’s rumor-control Twitter emulator (Liu et al.,

2024) demonstrate practical applications in digital social dynamics (Gao et al., 2023).

2.2 Collaborative Problem-Solving

Cutting-edge systems employ structured debate protocols and knowledge fusion mechanisms (Du et al., 2023). Xiong’s FORD framework (Xiong et al., 2023) enhances reasoning through tri-phase argumentation processes, complemented by Du’s knowledge graph-based consensus formation (Du et al., 2023). Software engineering innovations like Qian’s dialogue-driven development (Qian et al., 2024) and Hong’s documentation-centric workflow (Hong et al., 2023) establish new paradigms for AI-assisted programming (Wang et al., 2023). Domain-specific implementations showcase methodological cross-pollination, from Sun’s legal argumentation system with adversarial validation (Sun et al., 2024) to Yu’s cognitive conflict-driven MOOC platform (Yu et al., 2024), demonstrating versatile problem-solving architectures (Ferber and Weiss, 1999).

Although impressive progress has been made in both social simulation and collaborative problem-solving, previous LLM-based multi-agent systems still lack solutions that are truly effective and broadly applicable for general-purpose tasks. In

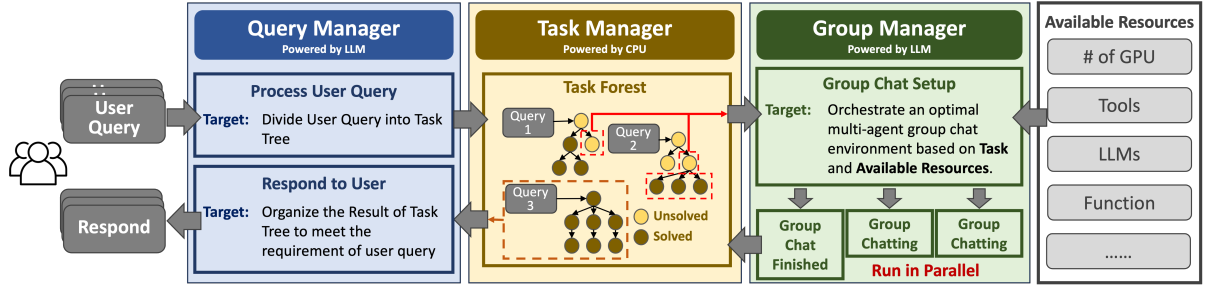


Figure 2: Illustration of AgentGroupChat-V2 framework, which composes of three main components: Query Manager, Task Manager and Group Manager. The framework illustrates the complete workflow from user query processing through task decomposition and management to multi-agent group chat execution, with arrows indicating data flow between components. Task Forest visualization demonstrates how queries are transformed into hierarchical task structures with solved (brown) and unsolved (yellow) nodes, while parallel group chats are carrying out in Group Manager.

contrast, the core insight of AGENTGROUPCHAT-V2 is to allow agents to autonomously adopt a divide-and-conquer paradigm—dynamically decomposing tasks and organizing cooperation according to the problem structure. This enables a flexible and scalable multi-agent system that can generalize across domains, addressing a key gap left by prior approaches.

3 Framework of AGENTGROUPCHAT-V2

This work introduces AGENTGROUPCHAT-V2, featuring a modular architecture composed of three main components: *Query Manager*, *Task Manager*, and *Group Manager* (see Figure 2). This design enables flexible data exchange and function invocation through standardized interfaces, allowing each module to scale independently with computational demand and supporting horizontal system expansion.

3.1 Query Manager

The Query Manager serves as the system’s front-end, handling user interaction and integrating LLMs as inference engines. It receives user queries, performs semantic analysis to decompose them into task tree structures, and forwards these to the Task Manager. After task processing is complete, the Query Manager collects and standardizes results, ensuring that responses accurately and appropriately address user requirements.

3.2 Task Manager

The Task Manager acts as the central coordination hub, managing overall task execution. Typically implemented as a single instance, it maintains a task forest representing multiple task trees derived from user input. The Task Manager tracks bidi-

rectional relationships between tasks, forming a comprehensive dependency graph. It allocates and schedules tasks based on their structure and system resource availability, enabling parallel execution where possible. For hierarchical dependencies, it ensures information flow between child and parent tasks, and upon completion of all sub-tasks, returns consolidated results to the Query Manager.

3.3 Group Manager

The Group Manager is responsible for orchestrating multi-agent collaboration to execute assigned tasks. It manages the preparation and deployment of agent groups, including the selection of suitable LLMs, resource allocation, and scheduling. Multiple Group Manager instances can operate in parallel, each supervising independent agent groups, which allows the system to efficiently scale its computational resources according to demand.

4 Group Chat Design

This section details the organizational and implementation mechanisms of group chat in AGENTGROUPCHAT-V2, the core of multi-agent collaboration and problem-solving. As illustrated in Figure 3, the system adopts a structured approach to both task management and agent group orchestration.

4.1 Task

In AGENTGROUPCHAT-V2, each *Task* is distinctly defined by an identifier, description, parent and child relationships, and processing result. Root tasks are parentless, and leaf tasks have empty child sets. Tasks progress through several states: initialization (structure created), waiting (awaiting prerequisites), execution (delegated to Group Man-

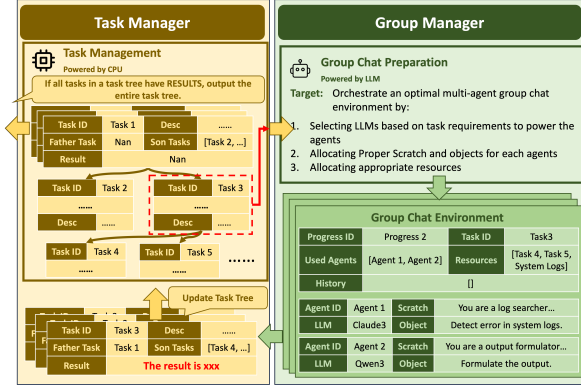


Figure 3: Detailed implementation of Task Manager and Group Manager modules showing: (1) Task Management structure with hierarchical organization of tasks, including Task ID, descriptions, parent-child relationships, and result tracking; (2) Group Chat Preparation process outlining agent selection criteria and resource allocation strategies; and (3) Group Chat Environment configuration displaying the agent infrastructure with Claude3 and Qwen3 LLMs assigned specific roles, workspace allocations, and specialized functions for systematic problem-solving through collaborative interactions.

ager), completion (result produced) or failure (error encountered). This structure forms hierarchical task trees that support flexible decomposition and dependency tracking.

4.2 Group

A *Group* represents a collaborative work unit managed by the Group Manager. Each group contains a progress ID, a linked task, assigned resources, and a set of agents. Agents within the group are specified by identity, LLM engine (e.g., Claude3, Qwen3), workspace (Scratch), assigned object (Object), and History. Scratch is for agent-internal computation, Object documents the role for interpretability, and History preserves conversational context for persistent and multi-turn collaboration. Groups are assembled with diverse agent roles, and their lifecycle includes preparation (configuration), activity (discussion), result integration, and termination.

4.3 Group Environment Configuration

Group environment configuration determines the initial state and capabilities of each group. This includes progress and task IDs, the list and order of agents, and shared resources. Group Manager orchestrates configuration, assigning agents and resources by analyzing the task needs and available LLMs. Each agent receives an appropriate iden-

Algorithm 1 StartGroupChat (GroupManager)

```

1: Input: max_action_turn, agent_ids, initial_env
2: Output: final_env, task_result
3: env ← initial_env
4: for turn = 1 to max_action_turn do
5:   for current_agent in agent_ids do
6:     env ← current_agent.perceive(env)
7:     (message, target_agent) ← cur-
      rent_agent.decide_action(env)
8:     dialogue_history ← EXECUTEAC-
      TION(current_agent, env, message, target_agent)
9:     env ← UpdateEnvironment(dialogue_history, env)
10:  end for
11:  discussion_summary ← SummarizeGroupMes-
      sages(env)
12:  env ← UpdateEnvironment(discussion_summary,
      env)
13:  is_complete ← TaskMan-
      ager.CheckTaskCompletion(env)
14:  if is_complete then
15:    return env, ExtractTaskResult(env)
16:  end if
17: end for
18: return env, ExtractTaskResult(env)

```

tity, Scratch space, task Object, and an initialized History.

4.4 Group Chat Orchestration

Group chat proceeds according to Algorithm 1, coordinated by Group Manager over a series of action turns. At each turn, agents sequentially perceive the environment, generate messages, interact with others, and update the environment and dialogue history. After each round, the system summarizes the discussion and checks if the task has been completed; results can be output early upon successful completion.

4.5 Agent Interaction

Agent interactions (Algorithm 2) support both broadcast and pairwise communication. Upon receiving an action, an agent may broadcast to all members or engage in an alternating dialogue loop with a specific target agent until a stopping criterion is met.

4.6 Chat Results Processing

To ensure effective utilization of group discussions, the Group Manager summarizes dialogues, extracts key conclusions, and stores them in agents' memories to inform future turns. It also validates the quality and completeness of outcomes, and formats results for downstream consumption by the Task Manager.

Baseline	LLM	GSM8K	MATH					AIME (2024)
			Level1	Level2	Level3	Level4	Level5	
Naive	Qwen2.5-72B	37.52	60.17	49.41	37.16	26.49	21.78	0.0
	Llama-3.1-70B	35.70	50.44	44.11	30.97	28.20	19.84	10.0
Naive-CoT	Qwen2.5-72B	75.13	85.84	77.65	75.22	63.68	48.25	10.0
	Llama-3.1-70B	87.33	91.15	88.82	87.17	80.77	67.32	16.7
ReAct	Qwen2.5-72B	52.76	53.10	35.88	36.73	28.21	19.46	3.3
	Llama-3.1-70B	20.09	28.32	11.76	7.52	4.27	2.33	0.0
AutoGen	Qwen2.5-72B	<u>81.80</u>	<u>89.38</u>	<u>85.29</u>	88.05	71.79	<u>54.86</u>	<u>16.7</u>
	Llama-3.1-70B	85.21	96.46	<u>92.94</u>	<u>91.15</u>	<u>82.48</u>	66.93	<u>20.0</u>
Multi-Agent Debate	Qwen2.5-72B	75.81	85.84	79.41	75.66	<u>66.24</u>	52.14	<u>16.7</u>
	Llama-3.1-70B	<u>90.82</u>	<u>97.35</u>	<u>92.94</u>	92.48	85.90	<u>71.98</u>	<u>20.0</u>
AGENTGROUPCHAT-V2	Qwen2.5-72B	87.41	92.92	90.00	<u>84.07</u>	71.79	59.10	21.4
	Llama-3.1-70B	91.50	98.23	94.12	88.94	81.20	83.54	30.4

Table 1: Experiment result on Math Problem.

Baseline	LLM	MBPP			HumanEval		
		pass@1	pass@3	pass@5	pass@1	pass@3	pass@5
Naive	Qwen2.5-72B	54.46	60.09	62.02	67.07	78.34	81.60
	Llama-3.1-70B	51.80	60.82	63.65	68.90	80.92	83.52
Naive-CoT	Qwen2.5-72B	51.30	61.52	65.08	62.86	78.60	82.43
	Llama-3.1-70B	53.04	60.51	62.75	71.58	81.23	83.56
ReAct	Qwen2.5-72B	<u>57.02</u>	65.34	68.50	<u>75.54</u>	84.39	86.84
	Llama-3.1-70B	<u>55.88</u>	<u>62.39</u>	<u>64.84</u>	<u>78.53</u>	85.71	87.66
AutoGen	Qwen2.5-72B	40.60	54.33	59.27	66.46	80.83	84.22
	Llama-3.1-70B	40.26	55.05	60.07	57.50	75.20	80.26
Multi-Agent Debate	Qwen2.5-72B	55.28	<u>63.76</u>	<u>66.67</u>	<u>75.54</u>	<u>83.19</u>	<u>85.41</u>
	Llama-3.1-70B	54.68	63.02	65.75	75.67	<u>84.81</u>	<u>87.65</u>
AGENTGROUPCHAT-V2	Qwen2.5-72B	60.34	<u>63.76</u>	64.45	76.46	82.31	84.15
	Llama-3.1-70B	58.84	60.35	60.84	79.20	80.38	80.91

Table 2: Experiment result on Code Generation.

5 Experiment Setup

5.1 Task & Benchmark

Mathematical Reasoning: We evaluate on **GSM8K** (Cobbe et al., 2021), **MATH** (Hendrycks et al., 2021), and **AIME** (Art of Problem Solving, 2024). Accuracy on all math benchmarks is verified through symbolic answer equivalence (see Appendix A for detailed evaluation criteria).

Code Generation: Benchmarks include **MBPP** (Austin et al., 2021) and **HumanEval** (Chen et al., 2021). Pass rate calculation and details on pass@k methodology are presented in Appendix A.

Domain-Specific Tasks: We test financial reasoning (**FinQual** (Xie et al., 2024b)), legal QA (**JEC-QA**), and medical QA (**MedmcQA** (Pal

et al., 2022)). Full descriptions of dataset construction are provided in Appendix B.

Structural Text Understanding: **StrucTextEval** (Gu et al., 2024a) assesses structured text processing (details in Appendix B).

Commonsense Reasoning: **HellaSwag** (Zellers et al., 2019) and **WinoGrande** (Sakaguchi et al., 2021); see Appendix B for dataset details.

5.2 Baseline Methods

We consider five representative baselines: **Naive**, **Naive-CoT** (Wei et al., 2022), **ReAct** (Yao et al., 2022), **AutoGen** (Wu et al., 2023), and **Multi-Agent Debate** (Liang et al., 2023). Implementation and protocol details can be found in Appendix C.

Method	Model	Width=1			Width=2			Width=3		
		Depth=1	Depth=2	Depth=3	Depth=1	Depth=2	Depth=3	Depth=1	Depth=2	Depth=3
Naive	Llama3.1-70B	78.5	77.2	72.8	69.3	63.5	<u>54.2</u>	66.8	53.3	42.7
	Qwen2.5-72B	<u>81.9</u>	<u>80.7</u>	<u>77.9</u>	<u>72.5</u>	<u>68.0</u>	<u>55.4</u>	<u>70.3</u>	<u>56.8</u>	<u>45.0</u>
Naive-CoT	Llama3.1-70B	84.1	83.8	<u>76.2</u>	75.5	70.3	53.7	73.2	<u>57.1</u>	<u>43.5</u>
	Qwen2.5-72B	86.7	85.3	78.3	77.8	71.4	54.8	75.5	55.2	42.3
ReAct	Llama3.1-70B	5.2	3.8	2.1	4.1	2.9	1.3	3.5	1.8	0.5
	Qwen2.5-72B	8.1	6.7	4.9	6.3	4.2	2.8	5.1	3.1	1.2
AutoGen	Llama3.1-70B	28.4	24.7	19.8	22.1	17.3	12.6	18.9	13.2	8.7
	Qwen2.5-72B	34.2	31.1	26.5	27.8	22.9	17.4	24.3	18.1	12.8
Multi-Agent Debate	Llama3.1-70B	<u>83.3</u>	<u>82.2</u>	74.8	<u>74.1</u>	<u>68.9</u>	52.2	<u>72.4</u>	54.7	40.3
	Qwen2.5-72B	82.1	81.5	75.2	76.3	67.6	53.8	71.2	55.4	41.7
AGENTGROUPCHAT-V2	Llama3.1-70B	81.7	80.8	79.1	73.6	69.2	58.4	72.8	62.9	48.7
	Qwen2.5-72B	83.5	82.9	77.3	76.7	70.8	59.2	73.6	64.7	52.1

Table 3: Exact Match (EM) Accuracy for StrucText-Eval.

Role Type	Role Setting
General Role	Agent-001 is a math expert .
Specialized Role	Agent-001 is a error detection specialist focused on identifying calculation mistakes, logical inconsistencies, and reasoning flaws in mathematical solutions.
	Agent-002 is a logical reasoning specialist specialized in mathematical proof construction, step-by-step deduction, and logical argument validation.
	Agent-003 is a context comprehension specialist responsible for understanding problem statements, extracting key information, and summarizing lengthy mathematical contexts.
	Agent-004 is a computational specialist dedicated to performing accurate calculations, numerical analysis, and algebraic manipulations.
	Agent-005 is a solution verification specialist focused on checking final answers, validating solution paths, and ensuring mathematical correctness.

Table 4: Role configuration for multi-agent mathematical problem-solving collaboration.

5.3 Large Language Models

All experiments use **Qwen2.5-72B-Instruct** and **Llama-3.1-70B-Instruct-Turbo**, with full model usage and inference details available in Appendix D.

6 Experiment

6.1 Overall Performance

6.1.1 Mathematical Reasoning Analysis

Finding 1

AGENTGROUPCHAT-V2 demonstrates strong mathematical reasoning by leveraging effective task decomposition and agent specialization, particularly excelling on high-complexity problems.

Table 1 reports results across GSM8K, MATH, and AIME, where AGENTGROUPCHAT-V2 not

only achieves the highest overall accuracy, but its advantages notably expand with problem complexity. For example, on the demanding AIME(2024), AGENTGROUPCHAT-V2 attains 30.4% accuracy (Llama-3.1-70B), nearly double that of the best baseline. The gap is less pronounced on simpler datasets (e.g., GSM8K), but becomes striking on the more layered MATH and in the progression from lower to higher levels within each benchmark. The marked improvement from Naive to Naive-CoT highlights the benefit of explicit step-wise reasoning, but only AGENTGROUPCHAT-V2 further leverages dedicated agent roles for effective divide-and-conquer. Importantly, as the diversity and interdependence of solution steps grow, single-agent and even loosely-structured frameworks stall, while specialized agents in AGENTGROUPCHAT-V2 coordinate to solve highly non-trivial subtasks in parallel.

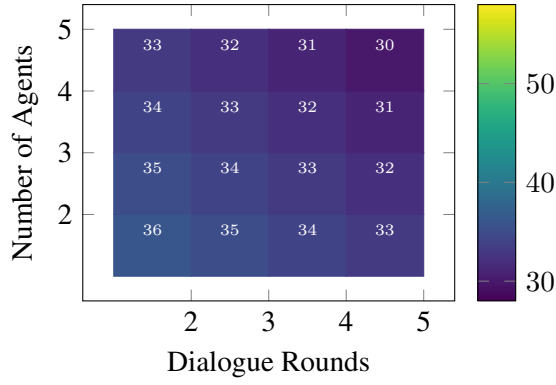


Figure 4: AGENTGROUPCHAT-V2 w/ General Role performance on MATH-100

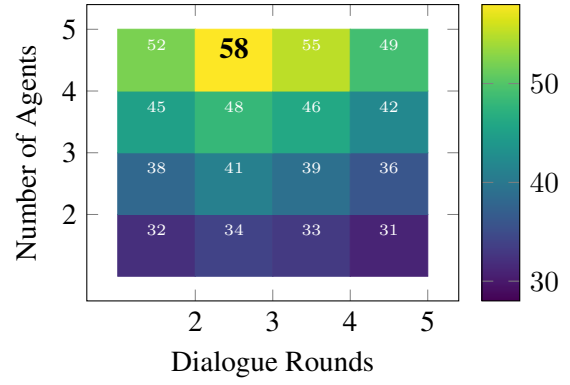


Figure 5: AGENTGROUPCHAT-V2 w/ Specified Role performance on MATH-100

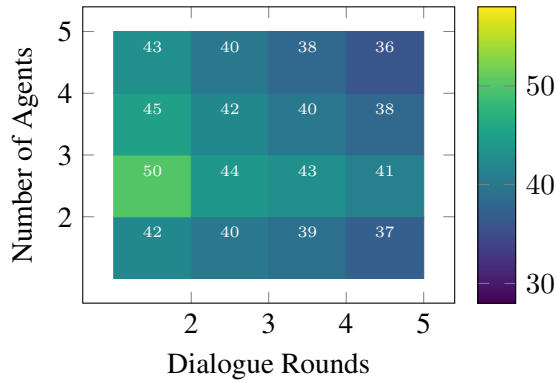


Figure 6: AutoGen performance on MATH-100

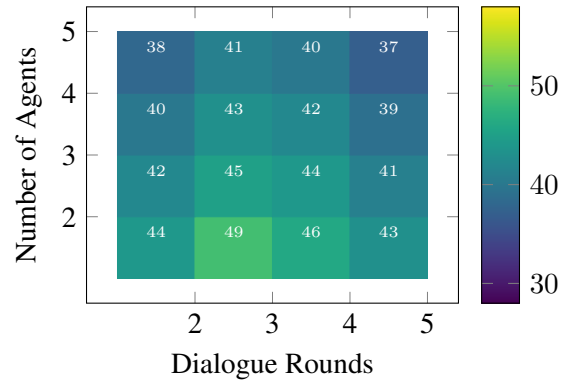


Figure 7: Multi-Agent Debate performance on MATH-100

6.1.2 Code Generation Analysis

Finding 2

Collaborative multi-agent analysis promotes rapid convergence to high-quality initial solutions in code generation, but performance improvements diminish for broader sampling, highlighting the need for task-specific exploration strategies.

As shown in Table 2, AGENTGROUPCHAT-V2 consistently delivers the best performance in pass@1 accuracy, achieving 79.20% on HumanEval and 76.46% on MBPP. These results suggest that multi-agent collaborative analysis quickly converges to high-quality first solutions, outperforming both naive and chain-of-thought approaches. For larger values of k (e.g., pass@5), ReAct and similar iterative frameworks gain ground, as their exploration strategy generates more diverse candidate programs. This confirms that while divide-and-conquer collaboration is ideal for tasks emphasizing correctness and logic in minimal samples, a more exploratory approach may be benefi-

cial when breadth and coverage are prioritized.

6.1.3 Structural Text Understanding Analysis

Finding 3

AGENTGROUPCHAT-V2 robustly handles increasing structural complexity, maintaining higher accuracy as task depth and breadth grow.

Performance on StrucText-Eval (Table 3) further demonstrates AGENTGROUPCHAT-V2's resilience to context and structure complexity. When both the width and depth of structured text expand, most methods—including sophisticated single-agent and multi-agent baselines—exhibit sharp performance degradation (sometimes $\geq 30\%$ absolute drop). In contrast, AGENTGROUPCHAT-V2's divide-and-conquer strategy achieves the highest accuracy in the most challenging (width=3, depth=3) scenarios, showing only moderate loss with increased context. This underscores the benefits of explicit recursive decomposition and hierarchical collaboration for tasks involving layered,

Algorithm 2 EXECUTEACTION{current_agent, env, message, target_agent}

```

1: Input: current_agent, env, message, target_agent, max_chat_turn
2: Output: dialogue_history
3: turn_count  $\leftarrow$  1
4: dialogue_history  $\leftarrow$  []
5: if target_agent is AllGroupMembers then
6:   message  $\leftarrow$  GenerateMessage(current_agent, env, message, target_agent)
7:   Append (current_agent, target_agent, message) to dialogue_history
8:   return dialogue_history
9: end if
10: sender  $\leftarrow$  current_agent
11: receiver  $\leftarrow$  target_agent
12: while turn_count  $\leq$  2  $\cdot$  max_chat_turn do
13:   response  $\leftarrow$  GenerateResponse(receiver, dialogue_history, env)
14:   if not response then
15:     return dialogue_history
16:   end if
17:   Append (receiver, sender, response) to dialogue_history
18:   temp  $\leftarrow$  sender
19:   sender  $\leftarrow$  receiver
20:   receiver  $\leftarrow$  temp
21:   turn_count  $\leftarrow$  turn_count + 1
22: end while
23: return dialogue_history

```

interdependent information.

6.2 Ablation Study

To elucidate the impact of team design and interaction strategy, we conduct systematic ablation experiments on MATH-100, varying agent count and dialogue rounds, and distinguishing between *general* and *specialized* role assignments, see Table 4 for definitions.

Finding 4

Specialized agent roles combined with moderate-scale collaboration significantly outperform homogeneous teams, scaling positively with agent diversity and demonstrating that targeted division of labor avoids redundancy and amplifies collective reasoning.

In AGENTGROUPCHAT-V2 with specialized roles (Fig. 5), accuracy grows robustly with team size: with 5 specialized agents and 3 dialogue rounds, accuracy peaks at 58%, a striking improvement of 64.6% over minimal settings. Notably, each additional agent consistently brings non-trivial gains, with collaboration between distinct expert perspectives enabling broader coverage and

more rigorous validation of solution steps.

Finding 5

Homogeneous (general) agent groups show stagnant or declining performance as either agent count or rounds increase, due to cumulative information redundancy and lack of complementary reasoning.

General-role configurations (Fig. 4) see accuracy plateau or degrade as more agents are added, with negligible improvement from increasing dialogue rounds. This effect highlights that, absent specialized functions, collaboration fails to create new insights and mostly produces redundant or even conflicting reasoning.

Finding 6

Traditional multi-agent frameworks (AutoGen, Debate) suffer diminishing returns or negative scaling at larger team sizes, illustrating the necessity of proactive role differentiation and dialogue orchestration to unlock benefits in large-scale collaborative problem solving.

AutoGen and Multi-Agent Debate (Figs. 6, 7) serve as external controls—both show early improvements with more agents or dialogue depth, but soon flatten or reverse, likely due to coordination overhead without sufficient specialization.

Further results. Experiments on **domain-specific QA** (Finance, Law, Medical), as well as common-sense reasoning (HellaSwag, WinoGrande), are provided in Appendix E. These results reveal that for tasks with low logical complexity or heavy reliance on external knowledge, the divide-and-conquer scheme may introduce unnecessary deliberation and overhead, and simpler strategies can be preferable. Rich qualitative **examples and discussion** of task decomposition can be found in Appendix F.

7 Conclusion

We introduce AGENTGROUPCHAT-V2 in this paper, a framework for LLM-based multi-agent systems featuring hierarchical task division and adaptive collaboration. Experiments show strong gains in reasoning and code generation, setting a new standard for complex multi-agent applications.

Limitations

AGENTGROUPCHAT-V2, as a multi-agent system framework based on large language models (LLMs), has achieved notable progress in broad compatibility and efficient task collaboration. Nevertheless, our evaluations reveal two main limitations. First, experimental results indicate that for tasks heavily reliant on extensive external knowledge or with relatively low logical complexity (such as commonsense reasoning or domain-specific question answering), introducing multi-agent division of labor and parallel discussion may lead to unnecessary inference redundancy and increased computational overhead. In such cases, a single model invocation or a simple reasoning chain may outperform the multi-agent approach. Second, the system’s performance is highly contingent upon the quality of task decomposition and the appropriateness of agent role assignment. For problems that are difficult to structurally decompose or involve highly coupled steps, the divide-and-conquer strategy does not necessarily yield solutions superior to those provided by a single agent.

Ethical Concerns

This study centers on algorithmic and system architectural design, with all evaluation tasks selected from publicly available and widely-used standard datasets. No sensitive personal information, simulated societal opinion, or real-world decision-making applications are involved. The multi-agent interactions are solely employed to enhance reasoning and problem-solving abilities, explicitly excluding the generation of misleading conclusions, group manipulation, or the deployment of autonomous systems with real-world operational authority. During the research process, neither the codebase nor experimental logs store any user privacy information; only aggregated evaluation data are released externally. Potential ethical risks, such as information bias or uncontrollable behavior emerging from automated collaboration, are primarily determined by the underlying LLM and task setup, rather than AGENTGROUPCHAT-V2 itself. All experiments are conducted within closed, secure environments. Should the framework be applied in scenarios involving human users or open environments in the future, further assessments and governance will strictly adhere to prevailing ethical guidelines in academia and industry.

References

- Art of Problem Solving. 2024. Aime problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions. Accessed: 2024-06-09.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Joseph Bates and 1 others. 1994. The role of emotion in believable agents. *Communications of the ACM*, 37(7):122–125.
- Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, and 1 others. 2025. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Allan Dafoe, Yoram Bachrach, Gillian Hadfield, Eric Horvitz, Kate Larson, and Thore Graepel. 2021. Co-operative ai: machines must learn to find common ground.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*.
- Jacques Ferber and Gerhard Weiss. 1999. *Multi-agent systems: an introduction to distributed artificial intelligence*, volume 1. Addison-wesley Reading.
- Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jingtao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. 2024. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications*, 11(1):1–24.
- Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023. S3: Social-network simulation system with large language model-empowered agents. *arXiv preprint arXiv:2307.14984*.

502	Zhouhong Gu, Haoning Ye, Xingzhou Chen, Zeyang	Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu,	557
503	Zhou, Hongwei Feng, and Yanghua Xiao. 2024a.	Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji,	558
504	Structext-eval: Evaluating large language model’s	Shaoguang Mao, and 1 others. 2024. Taskmatrix. ai:	559
505	reasoning ability in structure-rich text. <i>arXiv preprint</i>	Completing tasks by connecting foundation models	560
506	<i>arXiv:2406.10621</i> .	with millions of apis. <i>Intelligent Computing</i> , 3:0063.	561
507	Zhouhong Gu, Xiaoxuan Zhu, Haoran Guo, Lin Zhang,	Yuhan Liu, Xiuying Chen, Xiaoqing Zhang, Xing Gao,	562
508	Yin Cai, Hao Shen, Jiangjie Chen, Zheyu Ye,	Ji Zhang, and Rui Yan. 2024. From skepticism to	563
509	Yifei Dai, Yan Gao, and 1 others. 2024b. Agent-	acceptance: Simulating the attitude dynamics toward	564
510	groupchat: An interactive group chat simulacra for	fake news. <i>arXiv preprint arXiv:2403.09498</i> .	565
511	better eliciting emergent behavior. <i>arXiv preprint</i>		
512	<i>arXiv:2403.13433</i> .	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	566
513	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	567
514	Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-	Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	568
515	angliang Zhang. 2024. Large language model based	others. 2022. Training language models to follow in-	569
516	multi-agents: A survey of progress and challenges.	structions with human feedback. <i>Advances in neural</i>	570
517	<i>arXiv preprint arXiv:2402.01680</i> .	<i>information processing systems</i> , 35:27730–27744.	571
518	Xu Han, Zengqing Wu, and Chuan Xiao. 2023. ”guinea	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan	572
519	pig trials” utilizing gpt: A novel smart agent-based	Sankarasubbu. 2022. Medmcqa: A large-scale multi-	573
520	modeling approach for studying firm competition and	subject multi-choice dataset for medical domain ques-	574
521	collusion. <i>arXiv preprint arXiv:2308.10974</i> .	tion answering . In <i>Proceedings of the Conference</i>	575
522	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	<i>on Health, Inference, and Learning</i> , volume 174 of	576
523	Arora, Steven Basart, Eric Tang, Dawn Song, and	<i>Proceedings of Machine Learning Research</i> , pages	577
524	Jacob Steinhardt. 2021. Measuring mathematical	248–260. PMLR.	578
525	problem solving with the math dataset. <i>NeurIPS</i> .	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Mered-	579
526	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng	ith Ringel Morris, Percy Liang, and Michael S Bern-	580
527	Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven	stein. 2023. Generative agents: Interactive simulacra	581
528	Ka Shing Yau, Zijuan Lin, Liyang Zhou, and 1	of human behavior. In <i>Proceedings of the 36th an-</i>	582
529	others. 2023. Metagpt: Meta programming for	<i>annual acm symposium on user interface software and</i>	583
530	multi-agent collaborative framework. <i>arXiv preprint</i>	<i>technology</i> , pages 1–22.	584
531	<i>arXiv:2308.00352</i> .	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan	585
532	Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei	Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng	586
533	Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruim-	Su, Xin Cong, and 1 others. 2024. Chatdev: Com-	587
534	ing Tang, and Enhong Chen. 2024. Understanding	municative agents for software development. In <i>Pro-</i>	588
535	the planning of llm agents: A survey . <i>Preprint</i> ,	<i>ceedings of the 62nd Annual Meeting of the Associa-</i>	589
536	<i>arXiv:2402.02716</i> .	<i>tion for Computational Linguistics (Volume 1: Long</i>	590
537	Abhishek Kumar. 2025. Large language model based	<i>Papers)</i> , pages 15174–15186.	591
538	multi-agent system augmented complex event pro-	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-	592
539	cessing pipeline for internet of multimedia things.	ula, and Yejin Choi. 2021. Winogrande: An adver-	593
540	<i>arXiv preprint arXiv:2501.00906</i> .	sarial winograd schema challenge at scale. <i>Commu-</i>	594
541	Sanguk Lee, Tai-Quan Peng, Matthew H Goldberg,	<i>nications of the ACM</i> , 64(9):99–106.	595
542	Seth A Rosenthal, John E Kotcher, Edward W	Ranjan Sapkota, Konstantinos I Roulmeliotis, and Manoj	596
543	Maibach, and Anthony Leiserowitz. 2023. Can large	Karkee. 2025. Ai agents vs. agentic ai: A concep-	597
544	language models capture public opinion about global	tual taxonomy, applications and challenges. <i>arXiv</i>	598
545	warming? an empirical assessment of algorithmic	<i>preprint arXiv:2505.10468</i> .	599
546	fidelity and bias. <i>arXiv preprint arXiv:2311.00217</i> .	Yoav Shoham and Kevin Leyton-Brown. 2008. <i>Mul-</i>	600
547	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii	<i>tiagent systems: Algorithmic, game-theoretic, and</i>	601
548	Khizbullin, and Bernard Ghanem. 2023. Camel:	<i>logical foundations</i> . Cambridge University Press.	602
549	Communicative agents for” mind” exploration of	Jingyun Sun, Chengxiao Dai, Zhongze Luo, Yangbo	603
550	large language model society. <i>Advances in Neural</i>	Chang, and Yang Li. 2024. Lawluo: A chinese	604
551	<i>Information Processing Systems</i> , 36:51991–52008.	law firm co-run by llm agents. <i>arXiv preprint</i>	605
552	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	<i>arXiv:2407.16252</i> .	606
553	Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Man-	607
554	Zhaopeng Tu. 2023. Encouraging divergent thinking	dlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and	608
555	in large language models through multi-agent debate.	Anima Anandkumar. 2023. Voyager: An open-ended	609
556	<i>arXiv preprint arXiv:2305.19118</i> .	embodied agent with large language models. <i>arXiv</i>	610
		<i>preprint arXiv:2305.16291</i> .	611

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Michael Wooldridge. 2009. *An introduction to multiagent systems*. John Wiley & sons.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Auto-gen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.

Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, and 1 others. 2023. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*.

Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye, Kai Shu, Adel Bibi, Ziniu Hu, Philip Torr, Bernard Ghanem, and Guohao Li. 2024a. Can large language model agents simulate human trust behaviors? *arXiv preprint arXiv:2402.04559*.

Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong Li, Yongfu Dai, Duanyu Feng, and 1 others. 2024b. Finben: A holistic financial benchmark for large language models. *Advances in Neural Information Processing Systems*, 37:95716–95743.

Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. 2023. Examining inter-consistency of large language models collaboration: An in-depth analysis via debate. *arXiv preprint arXiv:2305.11595*.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.

Jifan Yu, Zheyuan Zhang, Daniel Zhang-li, Shangqing Tu, Zhanxin Hao, Rui Miao Li, Haoxuan Li, Yuanchun Wang, Hanming Li, Linlu Gong, and 1 others. 2024. From mooc to maic: Reshaping online teaching and learning through llm-driven agents. *arXiv preprint arXiv:2409.03512*.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.

Guibin Zhang, Kaijie Chen, Guancheng Wan, Heng Chang, Hong Cheng, Kun Wang, Shuyue Hu, and Lei Bai. 2025. Evoflow: Evolving diverse agentic workflows on the fly. *arXiv preprint arXiv:2502.07373*.

Xinnong Zhang, Jiayu Lin, Libo Sun, Weihong Qi, Yihang Yang, Yue Chen, Hanjia Lyu, Xinyi Mou, Siming Chen, Jiebo Luo, and 1 others. 2024. Electionsim:

Massive population election simulation powered by large language model driven agents. *arXiv preprint arXiv:2410.20746*.

Yifan Zhang, Jingqin Wang, Jianye Yu, and Joey Tianyi Wen. 2023. Multi-agent reinforcement learning: A comprehensive survey. *arXiv preprint arXiv:2312.10256*.

A Details of Experiment Setup

A.1 Mathematical Reasoning Evaluation

For GSM8K, MATH, and AIME, we use symbolic equivalence to verify result accuracy. For AIME, exact match evaluation is employed on numerical answers.

A.2 Code Generation Evaluation Criteria

MBPP is evaluated by the unit-test pass rate: a solution is deemed correct if it passes all provided test cases.

HumanEval is assessed using the pass@k metric¹:

$$\text{Pass}@k = 1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \quad (1)$$

where n is the total number of generated solutions and c is the number of correct solutions.

B Datasets

B.1 Mathematical Reasoning

GSM8K (Cobbe et al., 2021): 1,300 grade school math problems with step-by-step reasoning. **MATH** (Hendrycks et al., 2021): 12,000 problems covering algebra, geometry, combinatorics, etc., with five difficulty levels. **AIME** (Art of Problem Solving, 2024): Problems from the American Invitational Mathematics Examination; assesses complex mathematical reasoning.

B.2 Code Generation

MBPP (Austin et al., 2021): 500 Python problems, primarily algorithmic, each with multiple test cases. **HumanEval** (Chen et al., 2021): 164 Python function synthesis tasks, evaluated using random solution sampling (see Appendix A).

B.3 Domain-Specific Tasks

FinQual (Xie et al., 2024b): 1,000 finance-related questions from CFA and FinQA, including both MCQs and numerical tasks. **JEC-QA**: 26,365 legal exam items for the Chinese Bar exam, focusing

¹As in (Chen et al., 2021), we use $n = 10$ samples per problem and report pass@1, pass@3, and pass@5.

on legal comprehension and scenario judgment. **MedmcQA** (Pal et al., 2022): 194,000 medical MCQs from Indian medical entrance exams.

B.4 Structural Text Understanding

StrucText-Eval (Gu et al., 2024a): Evaluates comprehension and manipulation of structured, hierarchical textual data.

B.5 Commonsense Reasoning

HellaSwag (Zellers et al., 2019): 70,000 multiple-choice questions; each instance offers four options for next-sentence prediction. **WinoGrande** (Sakaguchi et al., 2021): 44,000 pronoun disambiguation tasks to test contextual commonsense reasoning.

C Baselines

Naive: Each task is directly assigned to a single LLM, without explicit intermediate reasoning or decomposition.

Naive-CoT (Wei et al., 2022): Tasks are completed by a single LLM with chain-of-thought prompts, eliciting reasoning chains.

ReAct (Yao et al., 2022): Combines reasoning, action, and observation in a single-agent loop to iteratively decompose and solve tasks.

AutoGen (Wu et al., 2023): Multi-agent conversation framework combining AssistantAgent and UserProxyAgent, facilitating modular dialogue-based task solving.

Multi-Agent Debate (Liang et al., 2023): Multiple agents collaboratively analyze and debate a problem, iterating until consensus.

D Large Language Models

Our experiments employ **Qwen2.5-72B-Instruct** and **Llama-3.1-70B-Instruct-Turbo**, two leading open-source LLMs. Qwen is selected for its high performance on math and reasoning tasks, Llama-3.1 for its balance of capabilities and generation speed. All LLMs are run in inference mode with default decoding settings unless otherwise stated.

E Additional Experiments

E.1 Commonsense Reasoning

Table 5 presents experimental results on commonsense reasoning tasks (HellaSwag, WinoGrande), where direct model prompting or naive approaches often outperform collaborative frameworks due to the straightforward nature of the problems.

Baselines	LLM	HellaSwag	WinoGrande
Naive	Qwen2.5-72B	73.7	80.3
	Llama-3.1-70B	70.1	82.4
Naive-CoT	Qwen2.5-72B	72.3	85.5
	Llama-3.1-70B	67.7	84.3
ReAct	Qwen2.5-72B	71.6	82.4
	Llama-3.1-70B	67.7	82.6
AutoGen	Qwen2.5-72B	64.4	78.4
	Llama-3.1-70B	61.1	73.1
Multi-Agent Debate	Qwen2.5-72B	<u>72.5</u>	<u>85.0</u>
	Llama-3.1-70B	<u>68.9</u>	<u>85.1</u>
AGENTGROUPCHAT-V2	Qwen2.5-72B	70.3	82.7
	Llama-3.1-70B	66.0	85.6

Table 5: Experiment result on Commonsense Reasoning.

E.2 Domain-Specific Tasks

Table 6 reports results for financial, legal, and medical professional QA (FinQual, JEC-QA, MedmcQA). AGENTGROUPCHAT-V2 shows domain-adaptive effectiveness but is sometimes matched or surpassed by simpler competitive baselines on knowledge retrievalintensive tasks.

F Case Study and Qualitative Analysis

We provide illustrative cases on task tree decomposition and inter-agent group chat. For example, Figure 8 and 9 show typical hierarchical breakdowns in software engineering and document writing scenarios, with agent assignments detailed in each phase. Our group chat example demonstrates how specialized agents contribute distinct perspectives, requirement analysis, code, and review, progressively refining solutions within structured dialogue rounds. These findings highlight the effective orchestration and benefit of specialization in complex tasks.

Method	Model	JEC-QA (Law) EM Accuracy	FinQual (Finance) EM Accuracy	MedmcQA (Medical) EM Accuracy
Naive	Llama-3.1-70B	31.58	71.83	73.90
	Qwen2.5-72B	42.56	76.00	89.30
Naive-CoT	Llama-3.1-70B	32.93	73.60	89.97
	Qwen2.5-72B	38.80	79.20	76.40
ReAct	Llama-3.1-70B	22.71	43.50	83.10
	Qwen2.5-72B	29.92	67.10	70.70
AutoGen	Llama-3.1-70B	30.08	63.76	56.79
	Qwen2.5-72B	40.60	79.16	68.60
Multi-Agent Debate	Llama-3.1-70B	31.58	78.28	90.20
	Qwen2.5-72B	38.50	80.20	75.90
AGENTGROUPCHAT-V2	Llama-3.1-70B	30.62	77.09	81.82
	Qwen2.5-72B	41.20	77.11	79.00

Table 6: Experiment results on Law, Finance, and Medical multiple-choice QA tasks.

Table 7: Agent Assignment for Interactive Data Visualization Tool Development Tasks

Task Name	Participating Agents
Module Interface Design	Requirement Analyst Agent, Code Implementation Agent, Code Review Agent
File Parsing Implementation	Requirement Analyst Agent, Code Implementation Agent, Code Review Agent
Data Processing Implementation	Requirement Analyst Agent, Code Implementation Agent, Code Review Agent
Data Display Implementation	Requirement Analyst Agent, Code Implementation Agent, Code Review Agent
Interactive Function Implementation	Requirement Analyst Agent, Code Implementation Agent, Code Review Agent
Testing and Verification	Test Planning Agent, Test Execution Agent, Quality Assurance Agent

Table 8: Agent Assignment for Blockchain Technology Analysis Article Writing Tasks

Task Name	Participating Agents
Technology Survey	Research Planning Agent, Research Execution Agent, Content Review Agent
Case Collection	Research Planning Agent, Research Execution Agent, Content Review Agent
Market Analysis	Research Planning Agent, Research Execution Agent, Content Review Agent
Writing Technology Analysis Chapters	Writing Planning Agent, Writing Execution Agent, Content Review Agent
Writing Market Analysis Chapters	Writing Planning Agent, Writing Execution Agent, Content Review Agent
Overall Optimization	Integration Planning Agent, Integration Execution Agent, Quality Assurance Agent

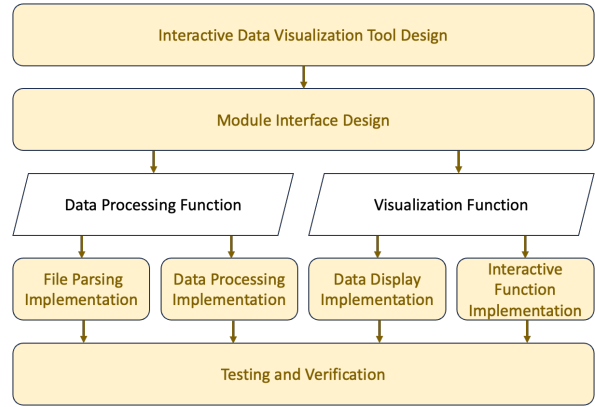


Figure 8: A schematic illustration of task decomposition for the “Interactive Data Visualization Tool Development” task by AGENTGROUPCHAT-V2. Rectangular boxes represent specific executable tasks, while parallelograms denote annotation labels for task groupings. Arrow connections indicate dependency relationships between tasks, and all tasks not directly connected can be executed in parallel after the completion of their prerequisite tasks.

Table 9: Two-Round Agent Dialogue for File Parsing Task

Round	Agent	Dialogue Content
Round 1	Requirement Analyst Agent	Let me analyze this file parsing task... We need to implement a parsing function supporting CSV and Excel files. Core requirements include: 1) File format validation; 2) Automatic character encoding detection; 3) Data parsing with JSON output; 4) Large file handling capability... I suggest implementing the basic framework first.
	Code Implementation Agent	<pre>def parse_file(file_path): try: if file_path.endswith('.csv'): df = pd.read_csv(file_path) elif file_path.endswith('.xlsx'): df = pd.read_excel(file_path) else: return {"error": "Unsupported format"} # ... data processing logic return df.to_dict('records') except Exception as e: return {"error": str(e)}</pre>
	Code Review Agent	Reviewing the implementation... The basic functionality exists, but several issues need attention: 1) No file size limitation, potential memory overflow; 2) Encoding detection relies entirely on pandas defaults, may cause garbled text; 3) Exception handling too simplistic, poor user experience; 4) Missing data type inference capability... Requires substantial improvements.
Round 2	Requirement Analyst Agent	Based on the code review feedback... Adding specific requirements: 1) 100MB file size limit; 2) Support UTF-8, GBK, ASCII encoding auto-detection; 3) Categorized error handling with user-friendly messages; 4) Intelligent data type recognition for dates, numbers, and text; 5) Chunked processing mechanism for large files... This should provide comprehensive coverage.
	Code Implementation Agent	<pre>import chardet import os from typing import Callable, Dict, Any class FileParser: def __init__(self, chunk_size=10000): self.chunk_size = chunk_size def parse_file(self, file_path: str) -> Dict[str, Any]: # File size validation if os.path.getsize(file_path) > 100*1024*1024: return {"error": "File exceeds 100MB limit"} # Enhanced encoding detection with open(file_path, 'rb') as f: raw_data = f.read(50000) detected = chardet.detect(raw_data) encoding = detected['encoding'] if detected['confidence'] > 0.7 else 'utf-8' try: if file_path.endswith('.csv'): # ... chunked CSV processing return self._parse_csv_chunked(file_path, encoding) elif file_path.endswith(('xlsx', 'xls')): # ... Excel processing with compatibility return self._parse_excel(file_path) except Exception as e: return {"error": f"Parse error: {str(e)}", "type": type(e).__name__} def _parse_csv_chunked(self, file_path, encoding): # ... streaming processing implementation chunks = [] for chunk in pd.read_csv(file_path, encoding=encoding, chunksize=self.chunk_size): chunks.append(chunk) df = pd.concat(chunks, ignore_index=True) return {"data": df.to_dict('records'), "dtypes": df.dtypes.to_dict()}</pre>
	Code Review Agent	Final review assessment...

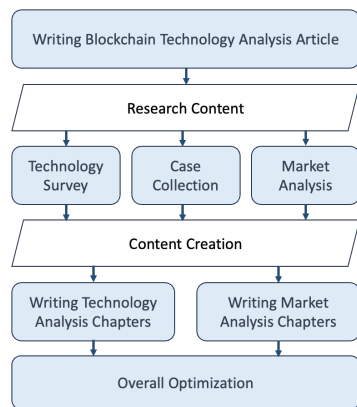


Figure 9: Another illustration of task decomposition for “Writing Blockchain Techonology Analysis Article” task by AGENTGROUPCHAT-V2