

Debate Only When Necessary: Adaptive Multiagent Collaboration for Efficient LLM Reasoning

Anonymous ACL submission

Abstract

Multiagent collaboration has emerged as a promising framework for enhancing the reasoning capabilities of large language models (LLMs). Despite improvements in reasoning, the approach introduces substantial computational overhead resulting from iterative agent interactions. Furthermore, engaging in unnecessary debates increases the risk of generating erroneous responses. To address these challenges, we propose Debate Only When Necessary (DOWN), an adaptive multiagent collaboration framework that integrates a deterministic gating module conditioned on the initial response. Debate is activated exclusively for queries that necessitate further deliberation, wherein agents refine their outputs by leveraging peer responses and their associated confidence scores. Evaluations on benchmarks show that DOWN improves efficiency by up to *six times* while preserving or even outperforming the performance of existing methods. Further analysis indicates that DOWN effectively mitigates the risk of error propagation stemming from the redundant debate process. These findings demonstrate the effectiveness of our approach in delivering high-performance LLM solutions at a lower computational cost.

1 Introduction

Building on the remarkable advancements in large language models (LLMs), recent research has increasingly focused on extending their capabilities to address complex real-world problems (Yao et al., 2023; Fan et al., 2024; Chen et al., 2024b). Among various research directions, multiagent collaboration has emerged as a promising approach, inspired by human decision-making processes in complex problem-solving (Minsky, 1988; Li et al., 2023; Chen et al., 2024a; Wang et al., 2025; Wu et al., 2023; Du et al., 2024). By engaging in structured debate, LLM agents systematically exchange perspectives and iteratively cross-examine each other’s

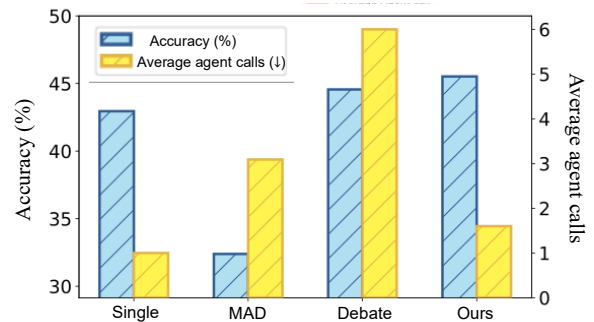


Figure 1: Comparison of accuracy and average agent calls across various multiagent debate methods

reasoning to refine their responses. This collaborative process facilitates divergent thinking and enhances the reasoning capabilities of LLMs (Liang et al., 2024; Chen et al., 2024b; Chan et al., 2024).

Despite these advantages, multiagent collaboration systems exhibit several key limitations. From an efficiency perspective, iterative interactions among agents inherently require multiple agent calls, leading to increased latency and higher inference costs (Snell et al., 2024; Kapoor et al., 2024). Meanwhile, when agents engage in redundant debate, there is a higher likelihood of generating errors that may subsequently be propagated by other agents in the system (Wang et al., 2024). Figure 1 illustrates both of these issues by plotting the accuracy and average agent calls of different multiagent debate methods. In the Debate system, additional debate rounds improve performance while incurring a sixfold increase in computational overhead. In the MAD framework, despite using more agent calls than the single agent baseline, its accuracy paradoxically declines. Regarding the practical application of multiagent collaboration systems, these challenges emphasize the need for an optimized collaboration approach (Kapoor et al., 2024; Tran et al., 2025).

To address these limitations, we propose Debate

070 Only When Necessary (DOWN), an adaptive multi-
071 agent collaboration framework designed to activate
072 debate selectively. At the core of DOWN is a de-
073 terministic gating module that operates solely on
074 the initial agent response and determines, on the
075 basis of its associated confidence score, whether
076 further debate is necessary. The gating module is
077 designed with the objective of minimizing compu-
078 tational cost while preserving performance and is
079 optimized through a utility score that balances effi-
080 ciency and predictive accuracy. This conditional ac-
081 tivation allows the system to invoke debate only for
082 queries that benefit from deeper collaborative rea-
083 soning, reducing computation without sacrificing
084 quality. Upon debate activation, agents engage in a
085 confidence-guided collaboration, drawing on their
086 peers' contributions and integrating the most cred-
087 ible and compelling elements to produce a more
088 accurate and coherent final outcome.

089 We employ models of varying scales, includ-
090 ing approximately 8B and 70B parameter mod-
091 els, evaluated on the MUSR, StrategyQA, Com-
092 monsenseQA, and MMLU benchmarks. The find-
093 ings demonstrate that adaptive debate engagement
094 achieves up to a sixfold reduction in computational
095 overhead while maintaining or even surpassing the
096 performance of full-debate baselines. Notably, we
097 reveal that this conditional debate serves as a safe-
098 guard against cascading errors, effectively enhanc-
099 ing the advantages of agent collaboration. This
100 work offers the following principal contributions:

- 101 • We propose Debate Only When Necessary
102 (DOWN), an adaptive multiagent framework
103 that selectively initiates debate based on the
104 initial output, achieving up to a sixfold cost
105 savings on representative reasoning bench-
106 marks while preserving competitive perfor-
107 mance.
- 108 • Extensive analysis reveals the strong robust-
109 ness and generalization capability of our ap-
110 proach, substantiating its broad applicability.
111 Moreover, the adaptive gating mechanism is
112 found to effectively mitigate error propaga-
113 tion.
- 114 • To the best of our knowledge, this is the first
115 study to investigate the gating mechanism in
116 multiagent systems grounded in the initial re-
117 sponse, maximizing efficiency while preserv-
118 ing the benefits of multiagent collaboration.

2 Related Work 119

2.1 LLM-based Multiagent Collaboration 120

121 Drawing inspiration from human collaborative
122 problem-solving behavior, multiagent collabora-
123 tion systems leverage collective intelligence to im-
124 prove decision-making. Studies have demonstrated
125 that LLM-powered multiagent systems promote
126 divergent thinking (Xiong et al., 2023; Liu et al.,
127 2024; Liang et al., 2024) and improve reasoning
128 capabilities (Li et al., 2023; Yin et al., 2023; Zhuge
129 et al., 2023). With these advantages, multiagent
130 collaboration is leveraged for diverse NLP applica-
131 tions: mitigating hallucinations (Fang et al., 2025),
132 aggregating knowledge across multiple specialized
133 LLMs (Wang et al., 2025), generating novel scien-
134 tific ideas and insights (Su et al., 2024), evaluat-
135 ing LLM-generated responses (Chan et al., 2024),
136 and refining datasets for instruction fine-tuning (Li
137 et al., 2024a). These advancements highlight the
138 growing impact of multiagent collaboration.

2.2 Debate Structures in Multiagent Collaboration Systems 139

140 Recent studies have developed debate structures
141 to optimize the benefits of multiagent collabora-
142 tion. For instance, Du et al. (2024) introduces a
143 framework in which agents iteratively refine their
144 responses based on peer-generated outputs. Liang
145 et al. (2024) propose a structured debate format
146 that assigns distinct roles to encourage divergent
147 thinking. Wang et al. (2024) develop a discussion
148 system in which agents are organized into multiple
149 groups to engage in discussions. 150

151 However, iterative debate frameworks face a fun-
152 damental trade-off between efficiency and accu-
153 racy (Kapoor et al., 2024; Snell et al., 2024). To
154 this end, Li et al. (2024b) restrict discussions to lo-
155 cal neighbors, while Zhou et al. (2025) and Liu et al.
156 (2024) introduce a shortcut mechanism that shares
157 similarities with our approach. However, both meth-
158 ods still incur non-trivial agent calls. The former ne-
159 cessitates non-trivial interactions among neighbor-
160 ing agents, and the latter relies on a self-consistency
161 criterion that requires at least five agent calls per
162 sample. Additionally, iterative debates may propa-
163 gate errors or introduce a trade-off between diver-
164 sity and quality (Wang et al., 2024; Kapoor et al.,
165 2024; Li et al., 2025). Our approach focuses on
166 these limitations by proposing a resource-efficient
167 and performance-effective debate framework.

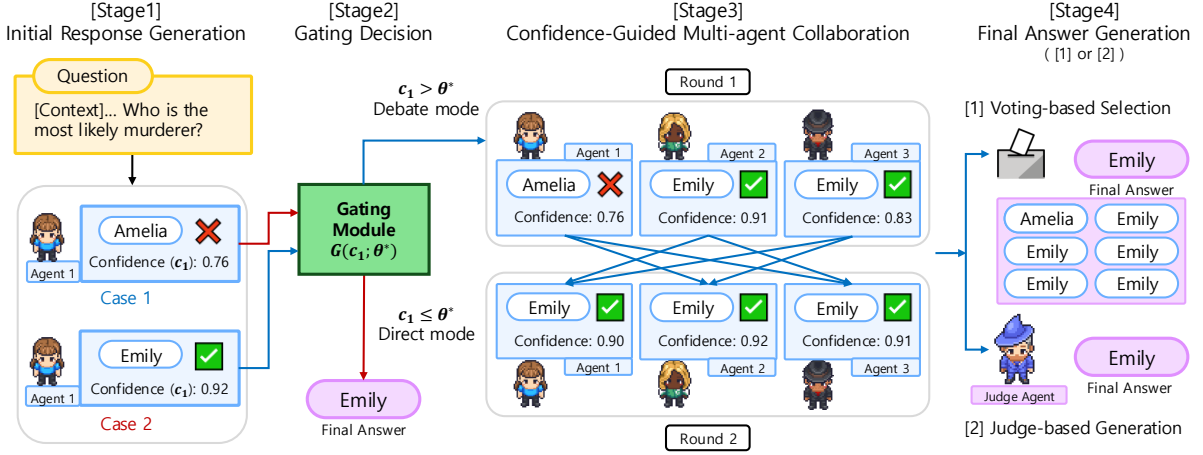


Figure 2: The DOW framework begins with an initial agent generating a response, after which the gating module determines whether further collaboration is required based on the confidence score of the initial output c_1 . Queries deemed acceptable ($G(c_1; \theta^*) = 0$) proceed via the direct mode without invoking additional agents, whereas those considered uncertain ($G(c_1; \theta^*) = 1$) trigger a collaborative process in which agents iteratively refine their answers by referencing the outputs and confidence scores of their peers. The final answer is subsequently determined either through majority voting or by a designated judge agent.

3 Debate Only When Necessary

We introduce Debate Only When Necessary (DOW), an adaptive multiagent collaboration framework that engages in debate based on the gating module. This section provides a detailed description of the debate procedure and the optimization of the gating module. Figure 2 illustrates the overall architecture.

3.1 System Pipeline

Stage 1: Initial Response Generation Given an input query q , an initial agent \mathcal{A}_1 from a pool of agents $\{\mathcal{A}\}_{i=1}^M$ generates an initial response r_1 , accompanied by a reason supporting its answer. We explore two configurations to ensure broad applicability: (1) A homogeneous-agent configuration where all collaborating agents use the same model (e.g., Llama) and (2) a heterogeneous-agent configuration where the initial agent is randomly selected for each query (e.g., Llama, Qwen, GPT-4o-mini). Along with the response, a confidence score is extracted to quantify the certainty of the prediction, a commonly used measure in prior work (Razghandi et al., 2025; Taubenfeld et al., 2025; Chen et al., 2024a). We first obtain the token logits $L(t_i)$, where $t_i \in r_1$, from the hidden representation of the generated response. These logits are then passed through a softmax function to yield the probability distribution $P(t_i)$ over tokens. To obtain a robust estimate of confidence, we apply normalization by averaging the token probabilities

over the generated response. The confidence score c_1 , is defined as:

$$c_1 = \frac{1}{|r_1|} \sum_{i=1}^{|r_1|} P(t_i), P(t_i) = \frac{\exp(L(t_i))}{\sum_j \exp(L(t_j))}. \quad (1)$$

For models without access to internal logits, we employ verbalized confidence as a proxy.

Stage 2: Gating Decision The gating module determines whether to engage in the collaborative debate guided by the initial response’s confidence score c_1 , which is formalized as:

$$G(c_1; \theta) = \begin{cases} 0, & \text{if } c_1 > \theta \quad (\text{Direct mode}), \\ 1, & \text{if } c_1 \leq \theta \quad (\text{Debate mode}), \end{cases} \quad (2)$$

where θ is the gating parameter, and its optimization procedure is described in the subsequent subsection. When the gating function assigns a query to $G(c_1; \theta) = 0$, the system directly accepts the initial response as the final output without invoking further agents, reducing computational overhead. Conversely, when $G(c_1; \theta) = 1$, the query is routed to the debate process, enabling agents to collaboratively refine and enhance the response.

Stage 3: Confidence-Guided Multiagent Collaboration The collaboration process involves response generation and refinement. In our experiments, we set up an environment with two rounds, each consisting of three agents. To conduct multiagent collaboration, in round 1, we obtain responses

r_2 and r_3 from additional agents \mathcal{A}_2 and \mathcal{A}_3 , along with their respective confidence scores c_2 and c_3 . These confidence scores are explicitly concatenated with each response to convey the agent’s certainty level in subsequent rounds.

After all agents have completed their first-round responses, the collaboration proceeds to the next round. Each agent \mathcal{A}_i gains access to the responses produced by the other agents while excluding its own, denoted as $\mathcal{R}_{-i} = \{r_j | j \neq i\}$. Leveraging this additional context \mathcal{R}_{-i} , each agent refines its reasoning and generates an updated second-round response $\{r_i^{(2)}\}_{i=1}^3$, capitalizing on more confident and persuasive arguments.

Stage 4: Final Answer Generation The responses generated in Stage 3 serve as ingredients for deriving the final answer. We explore two distinct strategies for final answer output: voting-based selection and judge-based generation. The voting-based approach determines the final answer by selecting the most frequent response among all agent-generated outputs. We design this majority voting to enhance robustness by leveraging consensus across multiple perspectives, effectively mitigating individual agent biases. The judge-based approach introduces an additional judging agent, which generates the final response based on all agent outputs. We induce a judging mechanism to update the final response, prioritizing well-supported and coherent arguments.

3.2 Gating Module Optimization

The performance of DOWN hinges on properly designing the gating module. We optimize the gating parameter θ by optimizing a utility-driven objective that controls the trade-off between predictive fidelity and computational efficiency.

Balancing Quality and Efficiency Let $A_\theta = \frac{k}{N}$ and $S_\theta = \frac{m}{N}$ denote the accuracy and debate skip rate (i.e. efficiency) for a given threshold value θ in the discrete candidate set Θ . To prevent sacrificing quality for efficiency, we identify a high-performance region using the one-sided 95% Wilson lower bound of the maximum observed accuracy. Given the highest observed accuracy A_θ^{max} obtained from k^{max} correct predictions, the Wilson lower bound \tilde{A}_θ^{max} is computed as:

$$\tilde{A}_\theta^{max} = \frac{2k^{max} + z^2 - z\sqrt{z^2 + 4k^{max}\left(1 - \frac{k^{max}}{N}\right)}}{2(N + z^2)}, \quad (3)$$

where $z = \Phi^{-1}(0.95) \approx 1.645$. Threshold candidates whose accuracy lies within this region are considered statistically indistinguishable from the best-performing value, whereas those outside the region $A_\theta < \tilde{A}_\theta^{max}$ are subject to soft penalization.

In parallel, the design prevents convergence to degenerate behaviors, such as engaging in debate either almost constantly or almost never. A penalty applies whenever the skip rate lies outside the interval $[s_{min}, 1 - s_{min}]$, where s_{min} denotes the proportion symmetrically excluded from the lower and upper ends of the unit interval $[0, 1]$. The total penalty term is defined as: $P_\theta = \max(0, \tilde{A}_\theta^{max} - A_\theta) + \max(0, s_{min} - S_\theta) + \max(0, S_\theta - (1 - s_{min}))$.

Utility-based Optimization Min-max scaling is employed for normalization, balancing the contributions of the two components to the final score:

$$\hat{A}_\theta = \frac{A_\theta - A_\theta^{min}}{A_\theta^{max} - A_\theta^{min} + \epsilon}, \hat{S}_\theta = \frac{S_\theta - S_\theta^{min}}{S_\theta^{max} - S_\theta^{min} + \epsilon}, \quad (4)$$

where $\epsilon \ll 1$ ensures numerical stability. The gating parameter is then optimized by maximizing the utility score: $\theta^* = \arg \max_{\theta \in \Theta} (\hat{A}_\theta + \hat{S}_\theta - \lambda \cdot P_\theta)$, where λ controls the relative strength of the penalty.

At inference time, the optimized gate module $G(c_1; \theta^*)$ governs the processing path. Queries with acceptable outputs proceed through the direct mode, whereas uncertain cases are directed to the debate mode. Though optimized, the gating module further permits flexible adjustment to user preferences, with higher thresholds favoring accuracy and lower thresholds emphasizing efficiency.

4 Experiments

Experimental Setup The evaluation considers two agent configurations: homogeneous and heterogeneous. The homogeneous setting employs a single model across all agents, using Llama-3.1 8B, Ministral 8B, Qwen-2.5 72B, Llama-3.3 70B, and GPT-4o-mini. The heterogeneous-model configuration combines Qwen-2.5 72B, Llama-3.3 70B, and GPT-4o-mini, with the model selection order randomized for each query. Each debate round involves three agents, and the debate process consists of two rounds, where initial responses are generated in the first round and refined in the second. Our experiments are conducted on three benchmarks of MUSR (Sprague et al., 2024), StrategyQA (Geva et al., 2021), and CommonsenseQA (CSQA) (Talmor et al., 2019), with MMLU (Hendrycks et al.,

Benchmark	Method	Ministral 7B		Llama3.1-7B		GPT-4o-mini		Llama-3.3-70B		Qwen-2.5-70B		Average	
		Acc	AC (↓)	Acc	AC(↓)	Acc	AC(↓)	Acc	AC(↓)	Acc	AC(↓)	Acc	AC(↓)
MUSR	Single CoT	51.06	1.00	42.95	1.00	55.75	1.00	56.33	1.00	57.80	1.00	52.78	1.00
	Self-refine	36.90	9.00	39.46	9.00	54.29	9.00	53.67	9.00	58.47	9.00	48.56	9.00
	Self-Consistency	48.24	6.00	44.70	6.00	55.88	6.00	58.18	6.00	58.29	6.00	53.06	6.00
	MAD	28.67	3.01	32.39	3.09	43.23	3.02	51.22	3.00	49.13	3.04	40.93	3.03
	Debate	48.54	6.00	44.56	6.00	57.32	6.00	57.28	6.00	58.69	6.00	53.28	6.00
	D <small>OWN</small> (Vote)	53.71	1.48	<u>45.51</u>	1.50	57.09	1.80	<u>57.80</u>	1.02	<u>59.39</u>	1.28	<u>54.70</u>	1.42
	D <small>OWN</small> (Judge)	53.71	1.57	45.52	1.60	57.35	1.96	57.80	1.03	59.52	1.34	54.78	1.50
StrategyQA	Single CoT	67.69	1.00	70.74	1.00	78.17	1.00	80.35	1.00	78.60	1.00	75.11	1.00
	Self-refine	67.69	9.00	69.54	9.00	76.42	9.00	77.73	9.00	78.17	9.00	73.91	9.00
	Self-Consistency	<u>68.12</u>	6.00	68.56	6.00	79.48	6.00	<u>80.79</u>	6.00	77.29	6.00	74.85	6.00
	MAD	57.64	3.73	44.54	4.66	<u>70.51</u>	3.38	79.04	3.07	<u>73.80</u>	3.24	65.07	3.62
	Debate	70.74	6.00	70.08	6.00	79.04	6.00	80.35	6.00	79.91	6.00	<u>76.02</u>	6.00
	D <small>OWN</small> (Vote)	<u>68.12</u>	3.16	71.18	2.53	80.79	1.92	82.53	1.07	<u>77.73</u>	2.64	76.07	2.26
	D <small>OWN</small> (Judge)	<u>68.12</u>	<u>3.59</u>	69.87	<u>2.83</u>	<u>79.91</u>	<u>2.10</u>	82.53	<u>1.08</u>	77.73	<u>2.97</u>	75.63	<u>2.51</u>
CSQA	Single CoT	70.60	1.00	70.02	1.00	82.15	1.00	83.70	1.00	84.60	1.00	78.21	1.00
	Self-refine	70.84	9.00	70.02	9.00	82.72	9.00	83.78	9.00	82.39	9.00	77.95	9.00
	Self-Consistency	69.78	6.00	73.55	6.00	82.39	6.00	<u>85.18</u>	6.00	85.42	6.00	79.26	6.00
	MAD	42.42	<u>4.30</u>	19.66	4.68	<u>78.54</u>	3.00	<u>78.30</u>	3.00	<u>74.20</u>	3.00	58.62	3.60
	Debate	67.73	6.00	70.52	6.00	83.05	6.00	85.42	6.00	84.60	6.00	78.26	6.00
	D <small>OWN</small> (Vote)	71.33	4.11	<u>70.76</u>	1.59	<u>82.96</u>	<u>3.70</u>	<u>84.68</u>	<u>3.80</u>	<u>84.93</u>	2.27	78.93	3.09
	D <small>OWN</small> (Judge)	<u>71.17</u>	4.73	<u>71.33</u>	<u>1.71</u>	82.80	4.24	84.68	4.36	<u>85.01</u>	<u>2.52</u>	<u>79.00</u>	<u>3.51</u>

Table 1: Comparison of accuracy (Acc) and average agent calls (AC) across single model strategies, multiagent debate systems, and our proposed approach on three benchmarks

2021) employed for further analysis of generalization performance. Accuracy is used as an evaluation metric, while efficiency is measured in terms of the average number of agent calls (AC) and debate skip rate. Further details of the models, benchmarks, and prompts are provided in Appendix A and Appendix B.

Baselines We conduct comparisons against multiple established reasoning systems, including a single Chain-of-Thought (CoT) (Kojima et al., 2022) baseline. (1) Self-Refine (Madaan et al., 2023): Each round consists of generating an initial response, receiving feedback, and producing a refined response. (2) Self-Consistency (Wang et al., 2023): Sampling responses multiple times and deriving the final answer by selecting the most consistent one through majority voting. (3) MAD (Liang et al., 2024): Conducting a debate between two agents with opposing perspectives while a moderator selects the most plausible solution or continues the debate if needed. (4) Debate (Du et al., 2024): Agents engage in a structured debate, iteratively updating their own responses by incorporating insights derived from other agents’ responses in the previous round.

We additionally introduce baselines that incorporate various efficient collaboration strategies. For a fair comparison, each proposed efficiency mechanism is integrated into the standard Debate or

MAD frameworks. The Early Stopping (Liu et al., 2024) terminates the debate when two out of three agent outputs are consistent in the initial round. Similarly, the Efficiency Optimization (Zhou et al., 2025) applies a consistency-based criterion that halts the debate when all of the outputs are consistent. Both methods utilize output consistency as the decision signal and are directly applied to the Debate method. In contrast, SparseMAD (Li et al., 2024b) improves efficiency by structuring debates through local neighbor interactions instead of fully connected discussions, which we incorporate into the MAD framework.

5 Results and Analysis

5.1 Main Results

Homogeneous-Model Configuration Table 1 presents the results on three benchmark datasets. Remarkably, our framework attains performance on par with or even outperforms single and multiagent baselines across the MUSR, StrategyQA, and CSQA benchmarks, while requiring only 1.42 to 3.09 AC on average, resulting in up to a sixfold improvement in efficiency. The framework demonstrates notable robustness across the Mistral, Llama, GPT, and Qwen families at parameter scales of 7B and 70B, attesting to the wide-reaching applicability of DOWN across diverse architectures. These findings highlight that conditional debate, rather



Figure 3: Performance comparison under a heterogeneous-model configuration. For single model approaches, we report GPT-4o-mini results.

374 than applying it uniformly or relying solely on a
 375 single agent response, offers a more principled balance
 376 between accuracy and efficiency. Interestingly,
 377 MAD exhibits lower average accuracy on all bench-
 378 marks. Consistent with findings from Wang et al.
 379 (2024), we attribute this to its inherent tendency
 380 toward contradictory reasoning. While constructive
 381 disagreement encourages divergent insights, it also
 382 intensifies erroneous reasoning, deteriorating the
 383 quality of final responses.

384 **Heterogeneous-Model Configuration** Under
 385 the heterogeneous-model configuration, external
 386 feedback originates from models trained on distinct
 387 datasets and scales, which facilitates the emergence
 388 of diverse and complementary reasoning behaviors.
 389 As shown in Figure 3, DOWN consistently outper-
 390 forms strong baselines, achieving higher accuracy
 391 while demanding far fewer agent calls across both
 392 datasets. DOWN surpasses the Debate in perfor-
 393 mance, with a markedly reduced computational
 394 burden. These results demonstrate that DOWN re-
 395 mains effective beyond homogeneous setups, adapt-
 396 ing seamlessly to heterogeneous-model collabora-
 397 tion. This highlights the robustness and scalability
 398 of our system in diverse reasoning environments.

399 **Final Answer Generation Strategies** Voting-
 400 based selection and judge-based evaluation are two
 401 strategies for consolidating debate responses into

Method	Collaboration	Acc.	AC
<i>Case1: Debate</i>			
Early Stopping	[R1]	82.53	<u>5.04</u>
	[R2]		
Efficiency Optimization	[R1]	<u>82.97</u>	10.00
	[R2]		
DOWN	[R1] [R2]	83.41	2.57
<i>Case2: MAD</i>			
SparseMAD (D= 5/5)	+	79.04	21.00
SparseMAD (D= 4/5)	+	<u>79.91</u>	18.00
SparseMAD (D= 3/5)	+	80.35	15.00
SparseMAD (D= 2/5)	+	79.04	<u>12.00</u>
DOWN	[R1] [R2]	<u>79.91</u>	3.01

Table 2: Results comparing efficient debate strategies on the StrategyQA dataset using the Llama-70B model. Red and blue circles indicate dynamic debate control: red when debate is skipped because consistent responses reach the specified count, blue when skipped by the gating module. Yellow circles denote standard agent outputs. Each round is denoted by R, and the communication graph density by D.

402 a final decision. Experimental results indicate that
 403 both approaches yield similar response patterns,
 404 with marginal differences depending on the spe-
 405 cific setting. However, the judge-based method re-
 406 quires an additional agent call, making it slightly
 407 less efficient. Given this trade-off, the voting-based
 408 approach is preferable when prioritizing computa-
 409 tional efficiency, as it achieves comparable accu-
 410 racy with fewer computational resources.

5.2 Comparison to Efficient Debate Methods

411 Table 2 presents experimental results where exist-
 412 ing efficiency strategies and our DOWN are inte-
 413 grated into the Debate (Case 1) and MAD (Case 2)
 414 frameworks. As a result, in Case 1, DOWN outper-
 415 forms existing efficiency methods in both accuracy
 416 and efficiency. DOWN achieves the highest accu-
 417 racy of 83.41 despite requiring only 2.57 agent
 418 calls. We attribute this to the fact that the Early
 419 Stopping and Efficiency Optimization require full
 420 agent response generation in the first round to eval-
 421 uate consistency, limiting their efficiency gains.
 422 In Case 2, SparseMAD reduces communication
 423 density through neighbor-based interactions but
 424 continues to incur a substantial number of agent
 425 calls. DOWN, by contrast, attains comparable accu-
 426

LLM Agent	Response Shift	MAD	Debate	Ours
GPT-4o-mini	✓→✗	70.59	50.00	33.59
	✗→✓	29.41	50.00	66.41
Llama-3.3 70B	✓→✗	48.91	60.09	12.57
	✗→✓	51.09	39.91	87.43
Qwen-2.5 72B	✓→✗	63.87	50.00	39.91
	✗→✓	36.13	50.00	60.09
Mix	✓→✗	70.85	30.81	47.35
	✗→✓	29.15	69.19	52.65

Table 3: Proportions of correct and incorrect response shifts induced by debate across multiagent collaboration systems. We denote a correct response as ✓ and an incorrect response as ✗.

racy with only 3.01 agent calls, thereby matching the performance of the 4-out-of-5 (4/5) neighbor-connected configuration. These findings highlight the persistent reliance of existing efficient debate methods on a substantial number of agent calls. DOWN, by discerning the need for debate solely from the initial response, delivers markedly higher efficiency while sustaining robust performance.

5.3 Response Shifts in Accuracy

Table 3 reports the proportions of correct and incorrect changes in answers, computed over StrategyQA samples where the final answer differs from the initial prediction. To deepen our investigation, we set the threshold parameters to 0.9, increasing the number of cases where debate is triggered. Across all model configurations, DOWN consistently achieves a higher successful correction rate than baseline debate systems. Notably, for the Llama 70B model, 87.43% of changed answers reflect successful corrections of initially incorrect predictions. MAD shows a high rate of incorrect changes, erroneously revising correct initial answers into incorrect ones. The Debate baseline shows a similar rate of correct and incorrect changes, indicating a limited capacity to prioritize reliable peer inputs. The high success rate of corrections achieved by DOWN stems from two core design principles. It initiates debate when the gating module directs the query to debate mode, avoiding unnecessary revisions of correct answers. Once debate is engaged, it draws on peer responses with their associated confidence scores, enabling the model to integrate reliable information.

5.4 Ablation Study

To analyze the contribution of each component in the DOWN framework, we conduct an ablation

Method	Acc.	AC
DOWN (θ^*)	71.18	2.53
- Confidence-guided debate	68.12	2.46
- Gating module (single response only)	70.74	1.00
- Gating module (debate only)	70.74	6.00
+ Lowered gating params ($\theta^* - 0.1$)	69.87	1.20
+ Raised gating params ($\theta^* + 0.1$)	71.62	3.53

Table 4: Ablation study using the StrategyQA dataset

study on the StrategyQA. The results are shown in Table 4. We begin by examining the role of confidence scores shared among agents during debate. Eliminating these signals results in a 3.06% point decline in accuracy, underscoring that confidence serves as a crucial cue for selectively incorporating peer responses throughout the debate process.

The subsequent ablation analysis scrutinizes the efficacy of the gating module. When the gating module is removed and the debate process is either entirely suppressed or universally enforced for all inputs, accuracy drops to 70.74 in both scenarios. This outcome highlights that the gating module mediates a principled balance between performance gains and computational cost by reducing the average number of agent calls. The findings further reveal that indiscriminate debate engagement is redundant, conferring no benefit over the single-agent baseline and indicating that existing multi-agent collaboration strategies remain suboptimal with respect to both accuracy and efficiency.

To further validate the efficacy of the gating parameter, we perturb θ^* by ± 0.1 and examine the resulting impact on performance. Deviations from the optimized value in either direction expose a trade-off between accuracy and the average number of agent calls. By contrast, the optimized threshold θ^* maintains accuracy within a statistically reliable range, as quantified by the Wilson lower bound, while markedly reducing computational cost.

5.5 Generalization Analysis

While DOWN has demonstrated strong effectiveness on reasoning tasks, we evaluate its generalization capability across diverse knowledge domains using the MMLU benchmark to substantiate its broader applicability. The performance of multiagent collaboration systems across six domains in MMLU is shown in Figure 5. Experimental results show that DOWN achieves performance comparable to that of the Debate method across tasks. Given that the Debate approach con-

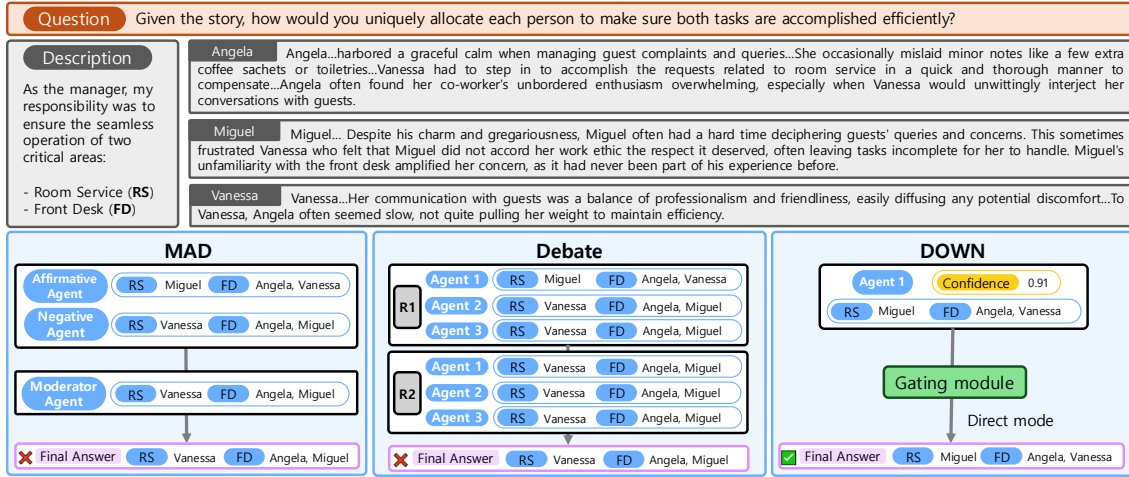


Figure 4: Qualitative comparison of multiagent collaboration methods evaluated on the MUSR dataset

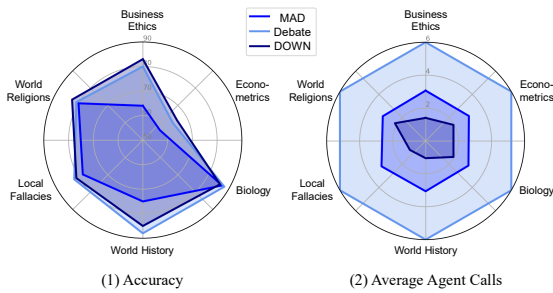


Figure 5: Accuracy and average agent calls (AC) of multiagent debate methods across six MMLU domains

504 sistically requires an average of six agent calls, the
 505 DOWN framework shows significantly higher effi-
 506 ciency. The results suggest that DOWN generalizes
 507 well, achieving robust performance on benchmarks
 508 evaluating both reasoning capabilities and factual
 509 knowledge across diverse domains.

510 From another perspective, *DOWN* supports both
 511 dataset-specific construction and broadly applica-
 512 ble construction of the gating module, depending
 513 on the user preference. As shown in Table 1 and
 514 Figure 5, gating modules produced under both set-
 515 tings consistently achieve strong and stable im-
 516 provements in efficiency and accuracy. These re-
 517 sults demonstrate that the gating module effectively
 518 accommodates a wide range of data conditions
 519 rather than overfitting to any particular dataset. This
 520 highlights DOWN’s potential to serve as a scalable
 521 and effective alternative to full multiagent debate
 522 systems.

5.6 Qualitative Analysis

524 Figure 4 presents a qualitative comparison of re-
 525 sponses generated by different debate systems on

the MUSR dataset. The results show that existing
 526 debate mechanisms introduce unnecessary modi-
 527 fications. Although the initial responses produced
 528 by the Debate and MAD methods are correct, sub-
 529 sequent iterative revisions lead to an incorrect fi-
 530 nal prediction. This implies the potential risk of
 531 error propagation when redundant debate occurs.
 532 In contrast, with an initial confidence score of
 533 0.91, DOWN skips the debate and directly adopts
 534 the initial response as the final answer. This sug-
 535 gests that the adaptive gating module under high-
 536 confidence conditions prevents unnecessary modi-
 537 fications while maintaining efficiency. To further
 538 substantiate these findings, we present additional
 539 qualitative analysis in the Appendix C.

6 Conclusion

542 This work proposed the DOWN framework to ad-
 543 dress the computational inefficiencies and error-
 544 propagation challenges inherent in multiagent
 545 collaboration. By conditionally activating debate
 546 based on the gating module, we significantly re-
 547 duced computational overhead while maintaining
 548 or even improving performance. The results indi-
 549 cated that conditional debate improved efficiency
 550 and mitigated cascading errors, resulting in more
 551 stable and reliable reasoning dynamics. Further-
 552 more, the confidence-guided multiagent debate
 553 amplified the influence of reliable responses on
 554 final decisions. Taken together, these findings po-
 555 sition DOWN as a robust and principled optimiza-
 556 tion approach, delivering high-performance and
 557 computationally efficient solutions for multi-agent
 558 collaborative systems.

559 Limitations

560 While our proposed DOWN framework demon-
561 strates strong efficiency and robustness, several lim-
562 itations remain. While the confidence score widely
563 serves as a fundamental indicator of a model’s inter-
564 nal certainties (Razghandi et al., 2025; Taubenfeld
565 et al., 2025; Chen et al., 2024a) and is adopted
566 as the gating criterion in DOWN, it remains sus-
567 ceptible to over- or under-confidence. The use of
568 model-specific optimized thresholds ensures con-
569 sistent gating across models despite their varying
570 degrees of over- or under-confidence, though fur-
571 ther refinement through more sophisticated uncer-
572 tainty estimation remains a promising direction for
573 future work. This study primarily focuses on En-
574 glish debates, which may limit its applicability to
575 multilingual multiagent collaboration settings. Ex-
576 tending our method to multilingual LLMs would
577 provide deeper insights into its robustness across
578 diverse linguistic contexts. Due to computational
579 constraints, we employ the GPT-4o-mini instead of
580 the GPT-4o model. Applying the GPT-4o model to
581 our system could further deepen the understanding
582 of our DOWN framework.

583 References

584 Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu,
585 Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu.
586 2024. [Chateval: Towards better LLM-based eval-
587 uators through multi-agent debate](#). In *The Twelfth
588 International Conference on Learning Representa-
589 tions*.

590 Justin Chen, Swarnadeep Saha, and Mohit Bansal.
591 2024a. [ReConcile: Round-table conference improves
592 reasoning via consensus among diverse LLMs](#). In
593 *Proceedings of the 62nd Annual Meeting of the As-
594 sociation for Computational Linguistics (Volume 1:
595 Long Papers)*, pages 7066–7085, Bangkok, Thailand.
596 Association for Computational Linguistics.

597 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang,
598 Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu,
599 Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong,
600 Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie
601 Zhou. 2024b. [Agentverse: Facilitating multi-agent
602 collaboration and exploring emergent behaviors](#). In
603 *The Twelfth International Conference on Learning
604 Representations*.

605 Yilun Du, Shuang Li, Antonio Torralba, Joshua B.
606 Tenenbaum, and Igor Mordatch. 2024. [Improving
607 factuality and reasoning in language models through
608 multiagent debate](#). In *Proceedings of the 41st Inter-
609 national Conference on Machine Learning*, volume
610 235 of *Proceedings of Machine Learning Research*,
611 pages 11733–11763. PMLR.

Lizhou Fan, Lingyao Li, Zihui Ma, Sanggyu Lee, Huizi
612 Yu, and Libby Hemphill. 2024. [A bibliometric re-
613 view of large language models research from 2017 to
614 2023](#). *ACM Transactions on Intelligent Systems and
615 Technology*, 15(5):1–25. 616

Yi Fang, Moxin Li, Wenjie Wang, Lin Hui, and Fuli
617 Feng. 2025. [Counterfactual debating with preset
618 stances for hallucination elimination of LLMs](#). In
619 *Proceedings of the 31st International Conference on
620 Computational Linguistics*, pages 10554–10568, Abu
621 Dhabi, UAE. Association for Computational Linguis-
622 tics. 623

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot,
624 Dan Roth, and Jonathan Berant. 2021. [Did aristotle
625 use a laptop? a question answering benchmark with
626 implicit reasoning strategies](#). *Transactions of the
627 Association for Computational Linguistics*, 9:346–
628 361. 629

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,
630 Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
631 2021. [Measuring massive multitask language under-
632 standing](#). In *International Conference on Learning
633 Representations*. 634

Sayash Kapoor, Benedikt Stroebel, Zachary S Siegel,
635 Nitya Nadgir, and Arvind Narayanan. 2024. [Ai
636 agents that matter](#). *arXiv preprint arXiv:2407.01502*. 637

Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu-
638 taka Matsuo, and Yusuke Iwasawa. 2022. [Large lan-
639 guage models are zero-shot reasoners](#). In *Advances in
640 Neural Information Processing Systems*, volume 35,
641 pages 22199–22213. Curran Associates, Inc. 642

Huaoli, Yu Chong, Simon Stepputtis, Joseph Camp-
643 bell, Dana Hughes, Charles Lewis, and Katia Sycara.
644 2023. [Theory of mind for multi-agent collaboration
645 via large language models](#). In *Proceedings of the
646 2023 Conference on Empirical Methods in Natural
647 Language Processing*, pages 180–192, Singapore. As-
648 sociation for Computational Linguistics. 649

Renhao Li, Minghuan Tan, Derek F. Wong, and Min
650 Yang. 2024a. [CoEvol: Constructing better responses
651 for instruction finetuning through multi-agent coop-
652 eration](#). In *Proceedings of the 2024 Conference on
653 Empirical Methods in Natural Language Processing*,
654 pages 4703–4721, Miami, Florida, USA. Association
655 for Computational Linguistics. 656

Wenzhe Li, Yong Lin, Mengzhou Xia, and Chi Jin. 2025.
657 [Rethinking mixture-of-agents: Is mixing different
658 large language models beneficial?](#) *arXiv preprint
659 arXiv:2502.00674*. 660

Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Pe-
661 ter Grabowski, Yeqing Li, and Eugene Ie. 2024b.
662 [Improving multi-agent debate with sparse commu-
663 nication topology](#). In *Findings of the Association
664 for Computational Linguistics: EMNLP 2024*, pages
665 7281–7294, Miami, Florida, USA. Association for
666 Computational Linguistics. 667

668	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen,	723
669	Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and	Quoc-Viet Pham, Barry O’Sullivan, and Hoang D	724
670	Zhaopeng Tu. 2024. Encouraging divergent thinking	Nguyen. 2025. Multi-agent collaboration mech-	725
671	in large language models through multi-agent debate.	anisms: A survey of llms. <i>arXiv preprint</i>	726
672	In <i>Proceedings of the 2024 Conference on Empiri-</i>	<i>arXiv:2501.06322.</i>	727
673	<i>cal Methods in Natural Language Processing</i> , pages		
674	17889–17904, Miami, Florida, USA. Association for		
675	Computational Linguistics.		
676	Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi	Junlin Wang, Jue WANG, Ben Athiwaratkun, Ce Zhang,	728
677	Yang. 2024. A dynamic LLM-powered agent net-	and James Zou. 2025. Mixture-of-agents enhances	729
678	work for task-oriented agent collaboration. In <i>First</i>	large language model capabilities. In <i>The Thirteenth</i>	730
679	<i>Conference on Language Modeling.</i>	<i>International Conference on Learning Representa-</i>	731
680	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong,	733
681	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	and Yangqiu Song. 2024. Rethinking the bounds of	734
682	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang,	LLM reasoning: Are multi-agent discussions the key?	735
683	Shashank Gupta, Bodhisattwa Prasad Majumder,	In <i>Proceedings of the 62nd Annual Meeting of the</i>	736
684	Katherine Hermann, Sean Welleck, Amir Yazdan-	<i>Association for Computational Linguistics (Volume 1:</i>	737
685	bakhsh, and Peter Clark. 2023. Self-refine: Itera-	<i>Long Papers)</i> , pages 6106–6131, Bangkok, Thailand.	738
686	tive refinement with self-feedback. In <i>Advances in</i>	Association for Computational Linguistics.	739
687	<i>Neural Information Processing Systems</i> , volume 36,		
688	pages 46534–46594. Curran Associates, Inc.	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	740
689	Marvin Minsky. 1988. <i>Society of mind.</i> Simon and	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	741
690	Schuster.	and Denny Zhou. 2023. Self-consistency improves	742
691	Ali Razghandi, Seyed Mohammad Hadi Hosseini, and	chain of thought reasoning in language models. In	743
692	Mahdieh Soleymani Baghshah. 2025. Cer: Confi-	<i>The Eleventh International Conference on Learning</i>	744
693	dence enhanced reasoning in llms. <i>arXiv preprint</i>	<i>Representations.</i>	745
694	<i>arXiv:2502.14634.</i>		
695	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Ku-	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	746
696	mar. 2024. Scaling llm test-time compute optimally	Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang,	747
697	can be more effective than scaling model parameters.	Xiaoyun Zhang, and Chi Wang. 2023. Auto-	748
698	<i>arXiv preprint arXiv:2408.03314.</i>	gen: Enabling next-gen llm applications via multi-	749
699	Zayne Rea Sprague, Xi Ye, Kaj Bostrom, Swarat Chaud-	agent conversation framework. <i>arXiv preprint</i>	750
700	huri, and Greg Durrett. 2024. MuSR: Testing the lim-	<i>arXiv:2308.08155.</i>	751
701	its of chain-of-thought with multistep soft reasoning.	Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing	752
702	In <i>The Twelfth International Conference on Learning</i>	Qin. 2023. Examining inter-consistency of large lan-	753
703	<i>Representations.</i>	guage models collaboration: An in-depth analysis via	754
704	Haoyang Su, Renqi Chen, Shixiang Tang, Xinzhe	debate. In <i>Findings of the Association for Computa-</i>	755
705	Zheng, Jingzhe Li, Zhenfei Yin, Wanli Ouyang, and	<i>tional Linguistics: EMNLP 2023</i> , pages 7572–7590,	756
706	Nanqing Dong. 2024. Two heads are better than	Singapore. Association for Computational Linguis-	757
707	one: A multi-agent system has the potential to im-	tics.	758
708	prove scientific idea generation. <i>arXiv preprint</i>	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	759
709	<i>arXiv:2410.09403.</i>	Shafraan, Karthik R Narasimhan, and Yuan Cao. 2023.	760
710	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	React: Synergizing reasoning and acting in language	761
711	Jonathan Berant. 2019. CommonsenseQA: A ques-	models. In <i>The Eleventh International Conference</i>	762
712	tion answering challenge targeting commonsense	<i>on Learning Representations.</i>	763
713	knowledge. In <i>Proceedings of the 2019 Conference</i>	Zhangyue Yin, Qiushi Sun, Cheng Chang, Qipeng	764
714	<i>of the North American Chapter of the Association for</i>	Guo, Junqi Dai, Xuanjing Huang, and Xipeng Qiu.	765
715	<i>Computational Linguistics: Human Language Tech-</i>	2023. Exchange-of-thought: Enhancing large lan-	766
716	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	guage model capabilities through cross-model com-	767
717	4149–4158, Minneapolis, Minnesota. Association for	munication. In <i>Proceedings of the 2023 Conference</i>	768
718	Computational Linguistics.	<i>on Empirical Methods in Natural Language Process-</i>	769
719	Amir Taubenfeld, Tom Sheffer, Eran Ofek, Amir Feder,	<i>ing</i> , pages 15135–15153, Singapore. Association for	770
720	Ariel Goldstein, Zorik Gekhman, and Gal Yona. 2025.	Computational Linguistics.	771
721	Confidence improves self-consistency in llms. <i>arXiv</i>	Wei Zhou, Mohsen Mesgar, Annemarie Friedrich, and	772
722	preprint arXiv:2502.06233.	Heike Adel. 2025. Efficient multi-agent collabora-	773
		tion with tool use for online planning in complex	774
		table question answering. In <i>Findings of the Associ-</i>	775
		<i>ation for Computational Linguistics: NAACL 2025</i> ,	776
		pages 945–968, Albuquerque, New Mexico. Associa-	777
		tion for Computational Linguistics.	778

Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dylan R. Ashley, Róbert Csordás, Anand Gopalakrishnan, Abdullah Hamdi, Hasan Abed Al Kader Hamoud, Vincent Herrmann, Kazuki Irie, Louis Kirsch, Bing Li, Guohao Li, Shuming Liu, Jinjie Mai, Piotr Piękos, Aditya Ramesh, Imanol Schlag, Weimin Shi, Aleksandar Stanić, Wenyi Wang, Yuhui Wang, Mengmeng Xu, Deng-Ping Fan, Bernard Ghanem, and Jürgen Schmidhuber. 2023. *Mindstorms in natural language-based societies of mind*. In *RO-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*.

A Implementation Details

We utilize the pre-trained model parameters provided by HuggingFace, with the exact models listed as follows: Llama-3.1 8B (meta-llama/Llama-3.1-8B-Instruct), Mistral 8B (mistralai/Mistral-8B-Instruct-2410), Qwen-2.5 72B (Qwen/Qwen2.5-72B-Instruct), Llama-3.3 70B (meta-llama/Llama-3.3-70B-Instruct), and GPT-4o-mini (gpt-4o-mini). For the optimization of the gating module, we define the candidate set Θ to commence at 0.7 with an interval of 0.1, corresponding to a threshold that delivers a debate skip rate S_θ close to 1. The symmetric margin s_{\min} and the penalty scaling factor λ are set to 0.1 and 15. For reproducibility and deterministic response generation, the temperature is set to 0.0, and the random seed is fixed at 777. The maximum generated sequence length is limited to 512 tokens. All experiments are conducted under consistent configurations across models using four A6000 GPUs with 48GB of memory.

For the evaluation, we utilize benchmarks designed to assess reasoning capabilities. MUSR (Sprague et al., 2024) evaluates multi-step soft reasoning over free-text narratives, offering a more complex yet realistic reasoning challenge compared to synthetic benchmarks. Assessments are performed on 756 samples drawn from the murder mysteries, object placements, and team allocation subsets. StrategyQA (Geva et al., 2021) requires implicit reasoning, where inference steps must be strategically derived rather than explicitly provided. Evaluations are conducted on the development set comprising 229 samples. CommonsenseQA (CSQA) (Talmor et al., 2019) is a multiple-choice dataset requiring commonsense reasoning beyond textual information. Evaluations are conducted on the development set of 1,221 samples. To further assess the generalization capability of the DOWN framework, we conduct

additional experiments on the MMLU (Hendrycks et al., 2021) knowledge benchmark. We consider six domains of business ethics, econometrics, biology, world history, logical fallacies, and world religions that cover a broad spectrum of knowledge areas.

B Prompts leveraged in DOWN framework

The prompts employed in our experiments are presented in the following textboxes. These prompts are applied throughout the DOWN framework to ensure structured and consistent response generation.

Initial Answer Generation

[debate topic]

Please output your answer in json format, with the format as follows:
`{"base_reason": "", "base_answer": ""}`.
 Please strictly output in JSON format, do not output irrelevant content.

Initial Answer Generation (w Confidence)

[debate topic]

Please output your answer in json format, with the format as follows:
`{"base_reason": "", "base_answer": "", "confidence_score": range of 0-1}`.
 Please strictly output in JSON format, do not output irrelevant content.

Answer Update

These are the solutions to the problem from other agents:
 One agent solution: [agent1 response]
 One agent solution: [agent2 response]
 Using the solutions from other agents as additional information, can you provide your answer to the problem?

Judge Prompt

Based on the following responses, generate an updated response that most accurately addresses the given query.

Query: **[debate topic]**

Responses: **[agent responses]**.

Summarize your reasons for selecting this side and provide the final answer.

C Qualitative Analysis of Debate Systems

Additional qualitative analysis for each multiagent collaboration method is presented in Table 5, Table 6, and Table 7.

Table 5 presents the responses for different debate methods on the StrategyQA dataset. While the initial response in the MAD system is correct, the system revises its answer after being influenced by the opposing stance of the negative debater. While this process encourages divergent thinking, it ultimately leads to incorrect conclusions due to error propagation. In contrast, the Debate methodology consistently generates correct responses across all cases. Although this approach required six agent calls, it enhances response reliability. The DOWN methodology, on the other hand, produces an initial response with a confidence score of 0.95, leading the system to skip the debate process. This result highlights the efficiency gains achieved through conditional debate.

Table 6 shows the results on the MUSR dataset. Experimental results indicate that both the MAD and Debate methodologies fail to conduct correct reasoning. The initial response in DOWN shows a confidence score of 0.89, while subsequent responses report 0.84 and 0.92, respectively. Notably, in the second round, the initial agent revised its response based on higher-confidence outputs from other agents, ultimately arriving at the correct answer. Model confidence scores also generally increase in the second round. This finding underscores the importance of confidence scores as a guiding metric, allowing the system to integrate reliable responses.

Table 7 follows a trend similar to Table 5. The DOWN system strategically skips the debate process, enhancing efficiency while maintaining answer accuracy.

Method	Agent	Category	Content
Query	Could a hundred thousand lolcats fit on a first generation iPhone?		
Answer	Yes		
MAD	Affirmative Debater	Reason	The storage capacity of a first generation iPhone is 4GB, 8GB, or 16GB, and the average size of an image file is approximately 100KB. ... However, ... it's theoretically possible for a hundred thousand lolcats to fit on a first generation iPhone with sufficient storage.
		Answer	Yes
	Negative Debater	Reason	No. The first generation iPhone had a maximum storage capacity of 16GB for the 16GB model, but also consider other factors such as the operating system, apps, and other data that would occupy some of that space. ... Considering these constraints, it's unlikely that a hundred thousand lolcats could fit on a first generation iPhone.
		Answer	No
	Moderator	Reason	While the affirmative side presents a theoretical calculation that a hundred thousand lolcats could fit within the storage capacity of a higher-end first generation iPhone, the negative side raises ... by the operating system, other data, and the file system. These factors significantly reduce the available space for storing images, making it unlikely that a hundred thousand lolcats could be accommodated.
		Answer	No
	Final Answer		No
Debate	Round 1 Agent 1	Reason	The first generation iPhone had a storage capacity of 4GB, 8GB, ... but exceeds the lower-end models. ... Assuming a very small size for each lolcat to fit on the screen, it's theoretically possible to display them, but it would be extremely impractical and unlikely.
		Answer	Yes
	Round 1 Agent 2	Reason	The first generation iPhone has a storage capacity of 4GB, 8GB, or 16GB, and assuming an average size of a lolcat image (around 100KB), we can estimate the maximum number of images that can fit. ... Still, theoretically, a hundred thousand lolcats could fit on a 16GB first-generation iPhone if we only consider storage space for images and nothing else.
		Answer	Yes
	Round 1 Agent 3	Reason	The first generation iPhone had a storage capacity of 4GB, 8GB, or 16GB, and assuming an average image size of 100KB, we can estimate the maximum number of images that can be stored. ... Nevertheless, from a purely storage-based perspective, it is theoretically possible for a hundred thousand lolcats to fit on a 16GB first generation iPhone.
		Answer	Yes

Method	Agent	Category	Content	
	Round 2 Agent 1	Reason	Considering the storage capacity of the first generation iPhone, with models ranging from 4GB to 16GB, ... would require approximately 10GB of storage space. ... Given this interpretation and focusing solely on storage capacity, it is theoretically possible for a hundred thousand lolcats to fit on a 16GB first-generation iPhone, especially when considering only the storage aspect and disregarding the practicality of displaying them.	
		Answer	Yes	
	Round 2 Agent 2	Reason	Considering the storage capacity of the first generation iPhone, with models ranging from 4GB to 16GB, and estimating an average size of 100KB per lolcat image, a hundred thousand lolcats would require approximately 10GB of storage space. ... Therefore, focusing solely on storage, a hundred thousand lolcats could theoretically fit on a 16GB first-generation iPhone, albeit with the understanding that displaying them all at once would be highly impractical.	
		Answer	Yes	
	Round 2 Agent 3	Reason	Considering the storage capacity of the first generation iPhone, with models ranging from 4GB to 16GB, and assuming an average size of 100KB per lolcat image, a hundred thousand lolcats would require approximately 10GB of storage space. ... Taking into account the solutions from other agents, which also consider the theoretical possibility of storing a hundred thousand lolcats on a 16GB model, it appears that from a storage perspective, it is indeed possible.	
		Answer	Yes	
	Final Answer		Yes	
	DOWN	Round 1 Agent 1	Reason	The first generation iPhone had a storage capacity of 4GB, 8GB, or 16GB, and each lolcat image is approximately 100KB in size. A hundred thousand lolcats would require 10GB of storage space, which exceeds the capacity of the 4GB and 8GB models but fits within the 16GB model.
			Answer	Yes
Confidence Score			0.95	
Final Answer		Yes		

Table 5: Qualitative analysis of the StrategyQA dataset

Method	Agent	Category	Content
Query			<p>Given the story, how would you uniquely allocate each person to make sure both tasks are accomplished efficiently?</p> <p>Choices: ‘Seeking Sponsors and Donations: Amelia, Organizing Event Details: George and Lily’, ‘Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily’, ‘Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George’</p> <p>As the clock ticked closer to our most significant fundraising event, three key players - George, Amelia, and Lily - paced around the office, eager to know their assignments. As their manager, I held the responsibility of delegating two critical tasks among them: the challenging pursuit of Seeking Sponsors and Donations, and the intricate duty of Organizing Event Details.</p> <p>Amelia, our vibrant member, often found herself with Lily, brainstorming creative ideas in the cafeteria or around the picnic table outside. ... She was infamous for her love for detailed schedules and checklists, which often found her volunteering to plan her family reunions that ran smoothly under her supervision. ...</p> <p>However, the presence of George often marred her high spirits. George, indeed, was a hard one to deal with. Even though he had this endless passion ... , he often disagreed with Lily’s ideas during meetings, publicly criticized Amelia’s work enough to dampen her spirit, and worst of all, forgot crucial tasks that he was responsible for. ... But he often confessed feeling overwhelmed managing multiple tasks at once.</p> <p>Then there was Lily, a networking maven. ... Yet, her true forte was her meticulousness that reflected in her attention to detail. ... While she tended to dismiss George’s inputs during meetings, she showed a great deal of respect to Amelia and considered her suggestions before making a decision.</p> <p>Now, the ball was in my court. With the fundraising event nearing, it was my responsibility to ensure that these three unique individuals were correctly placed to utilize their strengths best. ... As the manager, I knew that this decision would shape the fate of our event and define our team’s dynamic in the foreseeable future.</p>
Answer			Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily
MAD	Affirmative Debater	Reason	Based on the characteristics and strengths of each individual, Amelia’s experience in attending charity galas and her ability to make connections, as well as her attention to detail, make her a strong candidate for Seeking Sponsors and Donations. ... George’s tendency to disagree with others and forget tasks suggests that he may not be the best fit for a role that requires independent work or a high level of organization, but his passion for fundraising can still be utilized in a team setting.
		Answer	Seeking Sponsors and Donations: Amelia, Organizing Event Details: George and Lily

Method	Agent	Category	Content	
	Negative Debater	Reason	I disagree with your answer. Lily's experience in sales and her ability to ask people for money, combined with her connections with wealthy friends and acquaintances, make her a strong candidate for Seeking Sponsors and Donations. ... It allows Lily to leverage her sales experience and connections to secure sponsors and donations, while Amelia and George work together to organize the event details, with Amelia's organizational skills helping to keep George on track.	
		Answer	Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George	
	Moderator	Reason	The negative side's allocation takes into account the strengths and weaknesses of each individual more effectively. ... George's passion for fundraising can still be utilized in a team setting, and his tendency to disagree with others and forget tasks is mitigated by working with Amelia and Lily.	
		Answer	Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George	
	Final Answer			Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George
	Debate	Round 1 Agent 1	Reason	The task of Seeking Sponsors and Donations requires someone with strong connections and the ability to ask for money, which Lily possesses due to her past experience in sales and her network of wealthy friends and acquaintances. ... However, his experience in a marketing firm could still be valuable in a supporting role or in a different capacity, but for the two main tasks at hand, Amelia and Lily are the better fit.
Answer			Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George	
Round 1 Agent 2		Reason	Based on the characteristics and interaction dynamics of the three individuals, Amelia's strengths in seeking sponsors and donations due to her connections and experience, as well as her attention to detail, make her a strong candidate for Seeking Sponsors and Donations. ... Instead, considering George's potential to disrupt the team's dynamic, it would be more efficient to have him focus on Seeking Sponsors and Donations, where his interactions are more limited, and utilize Amelia and Lily's collaborative and organized approach for Organizing Event Details.	
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily	
Round 1 Agent 3		Reason	To efficiently allocate tasks, we need to consider the strengths and weaknesses of each person. ... George, despite his flaws, has a passion for fundraising and experience in managing sponsorships, making him a better fit for Seeking Sponsors and Donations, where his skills can be utilized without disrupting the team's dynamics.	
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily	

Method	Agent	Category	Content	
	Round 2 Agent 1	Reason	Considering the characteristics and interaction dynamics of the three individuals, as well as the solutions provided by other agents, it becomes clear that Amelia and Lily should be paired for Organizing Event Details due to their complementary strengths, meticulousness, and positive working relationship. ... This allocation allows for the efficient utilization of each person's strengths while minimizing potential conflicts and disruptions to the team's dynamics.	
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily	
	Round 2 Agent 2	Reason	After considering the characteristics, strengths, and weaknesses of each individual, as well as the solutions provided by other agents, it becomes clear that the most efficient allocation of tasks would be to utilize Lily's strong network and sales experience for Seeking Sponsors and Donations, and pair Amelia and George for Organizing Event Details. ... Therefore, the most suitable allocation would be to have Lily focus on Seeking Sponsors and Donations, where her skills can be maximized without causing disruptions, and have Amelia and George work together on Organizing Event Details, with close supervision to mitigate any potential issues caused by George's behavior.	
		Answer	Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George	
	Round 2 Agent 3	Reason	After considering the characteristics, strengths, and weaknesses of each individual, as well as the solutions provided by other agents, it becomes clear that the most efficient allocation of tasks would be to utilize Amelia and Lily's collaborative and organized approach for Organizing Event Details. ... This allocation allows each individual to contribute their strengths to the tasks, minimizing the impact of their weaknesses and ensuring the success of the event.	
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily	
	Final Answer		Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily	
	DOWN	Round 1 Agent 1	Reason	The decision is based on the unique skills and strengths of each team member. Amelia's experience in attending charity galas and her ability to make connections, as well as her attention to detail, make her a strong candidate for Seeking Sponsors and Donations. ... Given their positive working relationship and respect for each other's opinions, pairing Amelia and Lily for Organizing Event Details would be an effective way to utilize their skills and minimize potential conflicts.
			Answer	Seeking Sponsors and Donations: Lily, Organizing Event Details: Amelia and George
		Confidence Score	0.89	

Method	Agent	Category	Content
	Round 1 Agent 2	Reason	To efficiently accomplish both tasks, it's crucial to allocate individuals based on their strengths and working dynamics. Amelia's experience with charity events, connections, and organizational skills make her ideal for Seeking Sponsors and Donations. ... By placing George with Lily, his enthusiasm can be channeled constructively under her detail-oriented supervision, minimizing the impact of his weaknesses.
		Answer	Seeking Sponsors and Donations: Amelia, Organizing Event Details: George and Lily
		Confidence Score	0.84
	Round 1 Agent 3	Reason	The decision is based on the unique strengths and weaknesses of each team member. Amelia's experience in attending charity galas and her ability to create detailed schedules make her a strong candidate for Organizing Event Details. ... George's passion for fundraising and experience in managing sponsorships make him a good candidate for Seeking Sponsors and Donations, but his tendency to disagree with others and forget tasks make it important to pair him with a task that has fewer collaborative elements.
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily
		Confidence Score	0.92
	Round 2 Agent 1	Reason	After considering the strengths, weaknesses, and dynamics of each team member, as well as the solutions provided by other agents, it becomes clear that the most efficient allocation of tasks would leverage the unique skills of each individual while minimizing potential conflicts. ... Despite his challenges with task management and interpersonal conflicts, his strengths can be effectively utilized in a role that plays to his advantages.
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily
		Confidence Score	0.91
	Round 2 Agent 2	Reason	After considering the solutions from other agents and re-evaluating the strengths and weaknesses of each team member, it's clear that Amelia's skills in attending charity galas and creating detailed schedules make her a strong fit for Organizing Event Details. ...Considering these points and the fact that Lily has a strong ability to seek sponsors and donations due to her sales experience and network, and Amelia and Lily have a positive working relationship, the most efficient allocation would be to have Amelia and Lily handle Organizing Event Details, utilizing their organizational skills and positive dynamic, and George, with his passion for fundraising, focus on Seeking Sponsors and Donations, where his strengths can be maximized with minimal disruption.

Method	Agent	Category	Content
		Answer	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily
		Confidence Score	0.89
	Round 2 Agent 3	Reason	After re-evaluating the strengths, weaknesses, and dynamics of the team members, and considering the solutions from other agents, it becomes clear that Amelia’s skills in attending charity galas, making connections, and attention to detail are highly valuable for Seeking Sponsors and Donations. ...Therefore, the most efficient allocation is to have Amelia focus on Seeking Sponsors and Donations, where her strengths can be fully leveraged, and to pair George with Lily for Organizing Event Details, where Lily’s attention to detail and event management skills can complement George’s passion for fundraising and help manage his weaknesses.
		Answer	Seeking Sponsors and Donations: Amelia, Organizing Event Details: George and Lily
		Confidence Score	0.90
		Final Answer-Vote	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily
		Final Answer-Judge	Seeking Sponsors and Donations: George, Organizing Event Details: Amelia and Lily. Reason: After carefully evaluating the strengths, weaknesses, and dynamics of each team member, as well as considering the solutions provided by other agents, it is clear that the most efficient allocation of tasks would leverage the unique skills of each individual while minimizing potential conflicts. ... Considering these points, the allocation that best utilizes the strengths of each team member while minimizing disruptions is to have Amelia and Lily handle Organizing Event Details, and George focus on Seeking Sponsors and Donations, where his strengths can be maximized with minimal disruption.

Table 6: Qualitative analysis of the StrategyQA dataset

Method	Agent	Category	Content
Query			<p>Given the story, how would you uniquely allocate each person to make sure both tasks are accomplished efficiently?</p> <p>Choices: 'Room Service: Vanessa, Front Desk: Angela and Miguel', 'Room Service: Angela, Front Desk: Miguel and Vanessa', 'Room Service: Miguel, Front Desk: Angela and Vanessa'</p> <p>In the vibrant heart of the city, a boutique hotel thrived, its pulse maintained by the constant hum of hospitality. As the manager, my responsibility was to ensure the seamless operation of two critical areas - Room Service and Front Desk. The task of assigning these duties fell to me, with the diligent trio of Miguel, Angela, and Vanessa under my supervision.</p> <p>Angela, with her seven-year-long tenure in various customer service roles, harbored a graceful calm when managing guest complaints and queries. ... But Angela was not without her quirks. ... At times, Vanessa had to step in to accomplish the requests related to room service in a quick and thorough manner to compensate. Angela often found her co-worker's unbordered enthusiasm overwhelming, especially when Vanessa would unwittingly interject her conversations with guests.</p> <p>Miguel, on the other hand, was the silver lining of good-natured humor and friendliness that added a personal touch to the hotel's room service deliveries. His go-with-the-flow demeanor, however, sometimes came off as disorganized, especially to Angela who appreciated structure. ... Despite his charm and gregariousness, Miguel often had a hard time deciphering guests' queries and concerns. This sometimes frustrated Vanessa who felt that Miguel did not accord her work ethic the respect it deserved, often leaving tasks incomplete for her to handle. Miguel's unfamiliarity with the front desk amplified her concern, as it had never been part of his experience before.</p> <p>Vanessa had perfected her role at a high-end hotel reception over a span of five meticulous years. ... Vanessa meticulously paid attention to every minor detail while performing room service tasks, sometimes stepping on Angela's toes without realizing it. To Vanessa, Angela often seemed slow, not quite pulling her weight to maintain efficiency.</p> <p>With these intricacies at play, as the hotel manager, I had to make conscious decisions assigning my team between Room Service and Front Desk duties. Balancing Angela, Miguel and Vanessa's distinctive skills, patience, and perspectives, the heart of hospitality continued to beat, painting vivid memories in the hotel guest's minds.</p>
Answer			Room Service: Miguel, Front Desk: Angela and Vanessa

Method	Agent	Category	Content
MAD	Affirmative Debater	Reason	Considering the strengths and weaknesses of each person, Angela's experience in customer service and her calm demeanor make her a good fit for the Front Desk. ... Vanessa's attention to detail and professionalism make her a good fit for Room Service as well, where she can ensure that tasks are completed efficiently.
		Answer	Room Service: Miguel, Front Desk: Angela and Vanessa
	Negative Debater	Reason	I disagree with your answer. My allocation would be: Room Service: Vanessa, Front Desk: Angela and Miguel Here's my reasoning: * Vanessa's attention to detail and professionalism make her a perfect fit for Room Service, where she can ensure that tasks are completed efficiently and accurately. ... * Angela's experience in customer service and her calm demeanor make her a good fit for the Front Desk, where she can handle guest complaints and queries in a professional and respectful manner. ... * Miguel's charm and friendliness are also well-suited for the Front Desk, where he can interact with guests and provide a warm welcome. ... Overall, ... Vanessa's efficiency and attention to detail in Room Service can help to compensate for Angela's occasional misplacement of minor notes, while Angela and Miguel's combined strengths at the Front Desk can provide a warm and professional welcome for guests.
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel
		Moderator	Reason
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel
	Final Answer		
Debate	Round 1 Agent 1	Reason	To allocate each person efficiently, we need to consider their strengths and weaknesses. ... Miguel's charm and friendliness are best utilized in room service, where he can interact with guests in a more relaxed environment and work at his own pace without the strict structure of the front desk.
		Answer	Room Service: Miguel, Front Desk: Angela and Vanessa

Method	Agent	Category	Content	
	Round 1 Agent 2	Reason	To efficiently allocate tasks, we need to consider the strengths and weaknesses of each person. ... Therefore, the most efficient allocation would be to assign Vanessa to Room Service, where her attention to detail and professionalism can shine, and Angela and Miguel to the Front Desk, where Angela's customer service skills can be utilized and Miguel's charm can be an asset.	
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel	
	Round 1 Agent 3	Reason	To allocate each person efficiently, we need to consider their strengths and weaknesses. ... However, considering the need for balance and the fact that Miguel might struggle with the structure of the Front Desk, it's better to pair Angela with Miguel at the Front Desk, where Angela's organizational skills can help Miguel, and Vanessa can handle Room Service, utilizing her strengths without overwhelming others.	
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel	
	Round 2 Agent 1	Reason	Considering the strengths and weaknesses of each person, Angela's customer service skills are best utilized at the Front Desk, but her disorganization could be mitigated by pairing her with a detail-oriented person. ... Given these considerations and the analysis from other agents, it's clear that Angela and Miguel's skills complement each other at the Front Desk, with Angela handling customer interactions and Miguel providing charm, while Vanessa's strengths are best utilized in Room Service, ensuring efficiency and attention to detail.	
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel	
	Round 2 Agent 2	Reason	Considering the strengths and weaknesses of each person, Angela's customer service skills make her a strong fit for the Front Desk, while Miguel's charm can be leveraged in Room Service. ... The most efficient allocation, considering all perspectives, would be to assign Vanessa to Room Service, where her attention to detail can ensure tasks are completed efficiently, and Angela and Miguel to the Front Desk, where Angela's customer service skills and organizational abilities can help mitigate Miguel's weaknesses.	
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel	
	Round 2 Agent 3	Reason	Considering the strengths and weaknesses of each person, Angela's calm demeanor and customer service skills make her a strong fit for the Front Desk. ... Assigning Vanessa to Room Service leverages her strengths, and pairing Angela with Miguel at the Front Desk utilizes Angela's customer service skills and Miguel's charm, with Angela's organizational skills potentially helping to mitigate Miguel's disorganization.	
		Answer	Room Service: Vanessa, Front Desk: Angela and Miguel	
	Final Answer			Room Service: Vanessa, Front Desk: Angela and Miguel

Method	Agent	Category	Content
DOWN	Round 1 Agent 1	Reason	To efficiently allocate tasks, it's essential to consider the strengths and weaknesses of each team member. Angela excels in customer service, particularly at the Front Desk, due to her calm demeanor and respectful tone. ... Additionally, her efficiency and thoroughness in room service tasks, although sometimes stepping on Angela's toes, demonstrate her capability to handle detailed work.
		Answer	Room Service: Miguel, Front Desk: Angela and Vanessa
		Confidence Score	0.91
		Final Answer	Room Service: Miguel, Front Desk: Angela and Vanessa

Table 7: Qualitative analysis of the MUSR dataset