

# Talk Isn't Always Cheap: Understanding Failure Modes in Multi-Agent Debate

**Andrea Wynn\***  
Department of Computer Science  
Johns Hopkins University  
awynn13@jhu.edu

**Harsh Satija\***  
Vector Institute

**Gillian K. Hadfield**  
Department of Computer Science  
Johns Hopkins University

## Abstract

While multi-agent debate has been proposed as a promising strategy for improving AI reasoning ability, we find that debate can sometimes be harmful rather than helpful. Prior work has primarily focused on debates within homogeneous groups of agents, whereas we explore how diversity in model capabilities influences the dynamics and outcomes of multi-agent interactions. Through a series of experiments, we demonstrate that debate can lead to a decrease in accuracy over time — even in settings where stronger (i.e., more capable) models outnumber their weaker counterparts. Our analysis reveals that models frequently shift from correct to incorrect answers in response to peer reasoning, favoring agreement over challenging flawed reasoning. We perform additional experiments investigating various potential contributing factors to these harmful shifts — including sycophancy, social conformity, and model and task type. These results highlight important failure modes in the exchange of reasons during multi-agent debate, suggesting that naive applications of debate may cause performance degradation when agents are neither incentivised nor adequately equipped to resist persuasive but incorrect reasoning.

## 1 Introduction

Large Language Model (LLM) agents have demonstrated remarkable problem-solving abilities across a wide array of complex reasoning tasks [Brown et al., 2020]. Recently, a new line of research on interactive reasoning among multiple LLMs through debate has promoted the multi-agent debate framework as a promising approach to enhancing the reasoning and decision-making capabilities of LLM agents [Du et al., 2023, Chan et al., 2023, Liang et al., 2023, Khan et al., 2024]. Various forms of multi-agent debate have been shown to improve performance on multiple arithmetic and strategic reasoning benchmarks [Du et al., 2023, Subramaniam et al., 2025], produce more truthful answers and evaluations [Chan et al., 2023, Khan et al., 2024], and enhance tasks such as machine translation [Liang et al., 2023] and negotiation [Fu et al., 2023]. The core concept of these studies is that by engaging LLM agents through structured argumentation or discourse, we can facilitate the exchange of reasoning among different agents and guide them toward more accurate answers. Intuitively, greater exchanges of reasoning should lead to better decisions—allowing multiple agents to challenge flawed reasoning, highlight overlooked details, and reduce individual biases. But is this always the case?

In this work, we show that the benefits of multi-agent debate are not as universal as commonly assumed. Through a series of empirical studies, we show that multi-agent debate can sometimes degrade performance, leading to worse final answers than those generated by a single agent acting alone. These failures are not rare edge cases, but arise systematically in settings where agents amplify

---

\*equal contribution

each other’s errors – agreeing reflexively rather than challenging flawed reasoning. These findings hold even when there is variation in the abilities of the participating LLM agents. For instance, we discover that introducing a weak or less capable (lower-performing) LLM agent into a debate with a strong or more capable (higher-performing) agent can detrimentally affect the debate outcome, producing results worse than if the agents had not engaged in discussion. The presence of a weaker agent disrupts the performance of the stronger agent. Moreover, in certain cases, the longer a debate continues, the more performance can degrade. In other words: talk isn’t always cheap – and in some cases, it’s actively harmful.

We present a systematic evaluation of multi-agent debate across multiple tasks, showing that debate can sometimes *harm* group performance, particularly with heterogeneous LLM agents engaged in debate. Our findings challenge the prevailing narrative that more discussion between agents is inherently beneficial. Instead, we uncover several key factors that mediate the success or failure of debate, including task type and complexity, agent diversity and capability, and social influence. In doing so, we offer a nuanced view of when and why debate helps – and when it hurts. Together, these results suggest that while debate remains a promising tool for improving model reasoning, it should be applied carefully to ensure safety on the task and setting of interest.

**Statement of Contributions.** We summarize our contributions as follows:

1. We conduct a comprehensive evaluation of the effectiveness of the multi-agent debate framework across three different datasets. In Section 5.1, we demonstrate that debate may degrade performance compared to majority voting. Additionally, in Section 5.2 we show that the performance may progressively deteriorate as the debate progresses.
2. We extend the debate framework beyond the homogeneous setting to investigate the effect of diverse agent populations. Contrary to the prevalent belief that LLM agents are inherently collaborative – suggesting that a mixture of diverse models improves response quality when they can access outputs from other models, even if those outputs are of lower quality [Wang et al., 2024] – we find this assumption does not hold in multi-agent debates. In Section 5, we observe that the presence of weaker agents can negatively affect performance.
3. Our analysis in Section 6 reveals that a significant portion of correct answers become corrupted during debate. We investigate this phenomenon from various perspectives, such as the effect of sequential revision (Section 6.1), social influence (Section 6.2), and sycophancy (Section 6.3). These insights opens avenues for future research focused on enhancing reasoning exchange in multi-agent systems.

## 2 Related Work

**Debate and multi-agent reasoning.** Multi-agent debate was initially proposed as method for the scalable oversight problem where a judge or verifier can interject and elicit hidden contradictions, using structured back-and-forth conversations [Irving et al., 2018, Khan et al., 2024, Michael et al., 2023, Kenton et al., 2024]. More recently, another form of multi-agent debate, sometimes referred to as multi-agent deliberation, investigates leveraging different LLM agents to surface better answers by having them exchange reasoning via iterative discussion [Chan et al., 2023, Liang et al., 2023, Subramaniam et al., 2025]. Most of these studies focus on a homogeneous setting where all LLM agents utilize the same underlying language model [Du et al., 2023] or a model of similar ability [Yao et al., 2025], finding that this approach enhances accuracy across various Question-Answering (QA) tasks.

Estornell and Liu [2024] examine the theoretical properties and effects of opinion diversity within the debate framework, reporting a “tyranny of the majority” effect. They found that if the majority of agents provide the same answer – regardless of its correctness – minority agents tend to conform, creating an echo chamber effect. Additionally, they propose a theoretical result indicating that diversity of model abilities should *improve* overall debate performance; we show empirical evidence that this is often not true in practice.

Some works explore debates between agents of diverse nature or abilities. For instance, Estornell et al. [2025], Subramaniam et al. [2025] propose training LLM agents to debate collaboratively with distinct roles (actors/generators and critics), demonstrating that this approach can surpass previous unsupervised debate setups in reasoning benchmarks. Finally, studies such as Amayuelas et al. [2024]

investigate whether the collaborative nature holds if an explicit adversary is introduced into the debate process—where the adversary actively seeks to reduce performance. These works suggest that, when agents constructively challenge each other, answers can improve. However, other studies caution that debate can fail when agents emphasize persuasion over truth. Agarwal and Khanna [2025] introduce a single-round debate on factual questions (using TruthfulQA) where one agent states a true answer calmly and another delivers a confident, emotional false answer. They show the LLM judge often chooses the persuasive falsehood with high confidence, suggesting that a vivid but incorrect argument can override a correct one. These results highlight a risk: unless the judge is well-calibrated, debate may amplify bluster, mirroring human misinformation scenarios. Yao et al. [2025] focus on investigating sycophancy in the debate setting, showing that agent disagreement rate decreases as debate progresses, and that this observation is correlated with performance degradation.

**Collaborative multi-agent frameworks.** Beyond explicit debates, many recent systems assume collaboration among LLMs improves reasoning [Li et al., 2023, Wang et al., 2024, Tran et al., 2025]. For instance, Wu et al. [2023] provides a general multi-agent conversation framework: developers can define many agents (assistant, user-proxy, tools, etc.) that autonomously chat to solve complex tasks. These role-based and decentralized systems often yield richer interactions (e.g. multi-turn planning) than a single LLM alone.

**Sequential revision.** Frameworks that utilize interactive reasoning at inference time in a sequential manner are employed for either self-refinement [Madaan et al., 2023] or self-consistency [Wang et al., 2022]. Self-refinement involves iteratively revising or adapting a model’s responses based on previous outputs, prompting the model to intentionally reflect on its existing responses and correct any mistakes [Kamoi et al., 2024]. Works based on self-consistency often explicitly run multiple reasoning paths in parallel. For example, Wang et al. [2023] sample numerous independent chain-of-thought answers and select the most common answer. This aggregation of diverse reasoning paths significantly enhances the accuracy of arithmetic and commonsense QA by minimizing uncertainty. He et al. [2025] extend this idea: they spawn multiple “reactive” and “reflection” agent pairs, each exploring a different reasoning path, and then use a separate summarizer to aggregate them. Likewise, Yang et al. [2025] build a decentralized multi-agent planner where each LLM agent maintains its own memory (a hierarchical knowledge graph) and communicates via structured prompts. They find that these collaborative agents reach goals with 60% fewer steps than a lone agent, underscoring the value of structured cooperation.

Collectively, prior work shows that multi-agent debate and collaboration can amplify reasoning abilities, but also highlights key failure modes: judges may be fooled by rhetoric, and human-like dynamics can bias group outputs. Our work builds on these insights by expanding the analysis to heterogeneous debate settings, then performing a deep dive into the potential factors affecting when debate helps – or hurts – LLM performance on a variety of tasks.

### 3 Setting: Multi-agent debate

Let  $\mathcal{Q}$  denote a dataset of questions related to a task, where each question  $q \in \mathcal{Q}$  in the task is natural language text. The objective is to generate an answer  $a \in \mathcal{A}$  for any given input question  $q \in \mathcal{Q}$ , where  $\mathcal{A}$  is the set of possible answers for that task. We assume that there exists a ground truth answer  $a^* \in \mathcal{A}$  for each question  $q \in \mathcal{Q}$  which is denoted by  $f^{gt} : \mathcal{Q} \rightarrow \mathcal{A}$ , i.e.,  $f^{gt}(q) = a^*$ .

**Single-Agent Setting:** In the single agent setting, each agent uses an underlying LLM  $l : \mathcal{Q} \rightarrow \mathcal{A}$  to generate answer  $a \sim l(q)$  where  $a \in \mathcal{A}$  denotes the answer generated by the LLM in response to the question  $q$ . The main metric of interest is accuracy which is calculated as  $\frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \mathbb{1}[l(q) = f^{gt}(q)]$ .

Usually, various task specific prompts are also given as input to the LLM in addition to the input question for generating an answer. Let  $\mathcal{P}_T$  denote the task-specific prompt, then the resulting LLM can be described as  $l_{\mathcal{P}_T} : \mathcal{P}_T \times \mathcal{Q} \rightarrow \mathcal{A}$  or  $a \sim l(\mathcal{P}_T(q))$ .

**Multi-Agent Debate:** We follow the multi-agent debate framework from Du et al. [2023], Subramaniam et al. [2025] which involves initially posing a question to a group of LLM agents. After the initial responses, which include answers plus reasoning, are generated, the debate process then iteratively revisits the question for each agent, contextualizing input through the collective responses from previous rounds. Specifically, each agent generates a new response to the question based on its

prior response and the summarized responses from the other agents. The final answer is based on the majority vote among all agents.

Formally, we have a group of  $N$  agents, each with their own LLM  $l_i$ , that are presented with a question  $q \in \mathcal{Q}$  and tasked with generating an answer  $a \in \mathcal{A}$ . We use  $d(l_1, \dots, l_n) : \mathcal{Q} \rightarrow \mathcal{A}$  to denote the debate procedure that takes as input a question and a set of LLM agents and generates an answer. The debate procedure runs over multiple rounds, where in any given round  $t$ , each agent  $i$  uses their underlying LLM  $l_i$  to iteratively generate an answer in the following manner:

1. **Starting Round** ( $t = 0$ ): Each agent  $i$  generates a response  $g_i \in \mathcal{G}$  via  $l_i$  using a starting prompt  $\mathcal{P}_{\text{starting}}$ , i.e.  $g_i^0 \sim l_i(\mathcal{P}_{\text{starting}}(q))$ . Here  $\mathcal{G}$  denotes the set of possible generations from LLM.
2. **Debate Rounds** ( $t = 1, \dots, T$ ): For the subsequent rounds, each agent  $i$  is given the question  $q$  and responses from the other agents from the previous round. Let  $o_i^t = \{g_j^{t-1}\}_{j \neq i}$  denote the outputs from the other agents from previous round. Then the goal of agent  $i$  for this round is to generate an updated response  $g_i^t \in \mathcal{G}$  that takes into account the responses from the other agents and its own response from the previous round  $g_i^{t-1}$ , i.e.,

$$g_i^t \sim l_i(\mathcal{P}_{\text{debate}}(q, o_i^t, g_i^{t-1})),$$

where  $\mathcal{P}_{\text{debate}}$  is the debate prompt.<sup>2</sup>

This procedure runs for  $T$  rounds, where the final output of the debate is the set of responses  $\{g_i^T\}_{i=1}^n$ . Then the majority response across all agents is selected as the final answer. Let  $\text{majority} : \mathcal{G}^N \rightarrow \mathcal{A}$  denote the majority voting and filtering for responses across all agents, then the final answer is given by  $a = \text{majority}(\{g_i^T\}_{i=1}^n)$ . Accuracy is evaluated with respect to the final for the group of debating agents.

## 4 Experimental Setup

In this section, we describe the set of tasks, models, and prompts used in our experiments.

### 4.1 Datasets

**CommonSenseQA (CSQA):** The CommonSenseQA dataset [Talmor et al., 2019] consists of multiple-choice questions with complex semantics that often require prior knowledge to answer correctly. The dataset is intended to test for prior common-sense knowledge encoded within LLMs and checks for common misconceptions. Additionally, this task (common-sense reasoning) was not evaluated in prior work on multi-agent debate, making it a good testbed for evaluating generalization of these debate approaches to different tasks.

**MMLU:** Massive Multitask Language Understanding, or MMLU [Hendrycks et al., 2021], is a widely used multiple-choice dataset covering 57 domains including elementary mathematics, US history, computer science, law, and more. To perform well on MMLU, models need robust world knowledge and problem solving ability.

**GSM8K:** GSM8K [Cobbe et al., 2021] is a dataset of linguistically diverse grade school math word problems which require multi-step mathematical reasoning to solve. This dataset is not multiple-choice and instead requires open-ended generation of the answer to the math questions, potentially with intermediate reasoning steps.

### 4.2 Models

We used three models from distinct model families in our experiments: GPT-4o-mini [OpenAI, 2024], LLaMA-3.1-8B-Instruct [et al., 2024] and Mistral-7B-Instruct-v0.2 [Jiang et al., 2023]. To align with prior work on multi-agent debate [Du et al., 2023, Subramaniam et al., 2025], we ran all experiments and models with the default temperature parameter, the  $\text{top\_p} = 0.9$ , maximum generation length

<sup>2</sup>If the responses are too large to fit in the context window, then they are summarized via making another LLM call, i.e.,  $o_i^t = \{l_i(\mathcal{P}_{\text{summarize}}(g_j^{t-1}))\}_{j \neq i}$ .

of 2048 tokens, and  $T = 2$  rounds of debate. We take 100 random samples for each task from the dataset and report the result over 5 random seeds.

### 4.3 Prompts

We provide prompts on an example question for each task below:

- **CommonSenseQA:** Can you answer the following question as accurately as possible? If a product doesn't last, what does it have a reputation of doing?: A) disintegrating, B) wearing out, C) dissolving, D) falling apart, E) dissipating Explain your answer by providing a bullet point summary of your reasoning, putting the answer in the form (X) at the end of your response.
- **MMLU:** Can you answer the following question as accurately as possible? What is the value of p in  $24 = 2p$ ?: A) p = 4, B) p = 8, C) p = 12, D) p = 24 Explain your answer by providing a bullet point summary of your reasoning, putting the answer in the form (X) at the end of your response.
- **GSM8K:** Can you solve the following math problem? Mark is trying to choose between two venues for a surprise party for his wife. The first venue charges a flat fee of \$200, regardless of how many guests attend; the second charges \$25 per person who attends. However, the first venue does not include food, which Mark estimates will cost \$5 for each person who attends. At the second venue, food for each guest is already included in the price. How many guests are necessary for the two venues to be equal in cost? Provide a bullet point summary of your reasoning. Your final answer should be a single numerical number, in the form `answer`, at the end of your response.

For all models and tasks, we provide the following system prompt,  $\mathcal{P}_{\text{system}}$ :

System prompt,  $\mathcal{P}_{\text{system}}$

You are a helpful assistant that can answer questions and provide helpful information.

For multi-agent debate, we additionally use the following prompt as  $\mathcal{P}_{\text{debate}}$  for each round of debate, adjusted to the answer format of each task (following Du et al. [2023]):

Debate prompt,  $\mathcal{P}_{\text{debate}}$

These are the solutions to the problem from other agents: {AGENT\_RESPONSES} Using the reasoning from other agents as additional advice, can you give an updated answer? Explain your reasoning. Examine your solution and that of other agents. Put your answer in the form (X) at the end of your response.

### 4.4 Code

We provide the source code for all our experiments at <https://github.com/TheNormativityLab/talk-aint-cheap/>.

## 5 Results

### 5.1 Effectiveness of Debate

We present results showing that debate can sometimes be *harmful* rather than helpful – in particular, that sometimes agents perform better *without any debate* than after exchanging reasons with other agents. We present the results in Table 1. We find that in the case of CommonSenseQA, which was not studied in prior work on multi-agent debate [Du et al., 2023], debate always harms performance. A key insight is that even when groups include more “strong” models (e.g., GPT) than “weak” ones (e.g., Mistral), the process of debate does not always yield performance gains. Contrary to the prevailing narrative that debate improves collective reasoning, our experiments demonstrate that performance






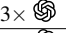






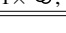
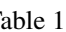

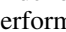
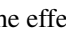
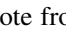

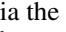
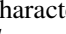
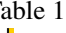
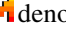
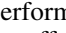
	CommonSense QA		MMLU		GSM8K	
1× 	74.8±1.9		82.6±2.1		93.2±1.8	
1× 	57.0±1.5		55.6±0.5		76.4±2.1	
1× 	41.6±2.3		34.0±1.9		34.2±1.8	
	w/o Debate	After Debate	w/o Debate	After Debate	w/o Debate	After Debate
3× 	44.4±2.7	39.4±3.9 ↓ 5.0	33.6±1.8	24.4±2.9 ↓ 9.2	43.6±1.5	46.4±1.4 ↑ 2.8
3× 	63.0±3.9	58.6±2.3 ↓ 4.4	61.6±2.4	57.8±1.8 ↓ 3.8	87.6±1.5	84.2±2.0 ↓ 3.4
3× 	75.6±2.2	74.8±2.1 ↓ 0.8	81.4±3.3	82.2±2.7 ↑ 0.8	94.0±0.9	94.4±1.5 ↑ 0.4
1×  , 2× 	66.2±2.2	64.4±2.1 ↓ 1.8	65.0±2.3	68.0±2.2 ↑ 3.0	88.4±1.3	92.8±1.7 ↑ 4.4
2×  , 1× 	74.8±1.3	74.0±0.7 ↓ 0.8	82.6±3.2	81.0±3.0 ↓ 1.6	93.6±0.8	94.6±1.3 ↑ 1.0
2×  , 1× 	58.2±3.8	50.2±3.9 ↓ 8.0	51.8±2.2	43.6±1.9 ↓ 8.2	82.6±1.9	75.8±2.1 ↓ 6.8
1×  , 2× 	53.4±2.7	46.8±2.5 ↓ 6.6	40.0±2.1	28.0±1.2 ↓ 12.0	61.0±2.4	64.8±1.5 ↑ 3.8
1×  , 2× 	62.4±1.1	59.4±1.9 ↓ 3.0	65.8±2.9	58.8±1.2 ↓ 7.0	90.2±0.8	87.8±1.9 ↓ 2.4
2×  , 1× 	74.6±1.6	72.4±2.7 ↓ 2.2	82.8±2.7	80.8±2.8 ↓ 2.0	93.4±1.3	93.0±1.3 ↓ 0.4
1×  , 1×  , 1× 	66.6±1.9	65.4±2.1 ↓ 1.2	57.8±3.1	63.4±2.1 ↑ 5.6	86.8±0.9	90.2±1.2 ↑ 3.4

Table 1: The first column shows the configuration of LLM agents, where  denotes GPT-4o-mini,  denotes Mistral-7B and  denotes the Llama-3.1-8B models. The top 3 rows benchmark the performance of how well a single agent performs on these tasks. The subsequent rows benchmark the effectiveness of debate procedure. The **w/o Debate** column represents the case where majority vote from multiple agents based on their initial responses is chosen, i.e., there is no exchange of reasoning. The **After Debate** denotes the result of majority vote after the exchange of reasons via the debate procedure. The arrows denote the difference in performance post debate procedure characterizing the benefit of exchange of reasons on the performance. A red arrow ↓ indicates a *decrease* in performance after debate, while a green arrow ↑ indicates an *increase* in performance after debate. All the experiments were done on 100 random samples and across 5 different seeds, reported are mean and standard error.

can actually decrease after agents engage in debate, even when stronger models outnumber weaker ones.

## 5.2 Performance Degradation during Debate

Figure 1 presents performance across three tasks – MMLU, CommonSenseQA, and GSM8K – as a function of debate rounds among groups of language models with varying individual performance on the underlying task. In fact, in many group configurations, we observe that performance *degrades* as the debate progresses. This trend appears across datasets, but is especially pronounced in MMLU and CommonSenseQA, where groups with mixed-capability models often suffer from group performance degradation during debate despite having a majority of stronger agents. These findings challenge the belief that deliberation or iterative reasoning among AI agents will always lead to better outcomes.

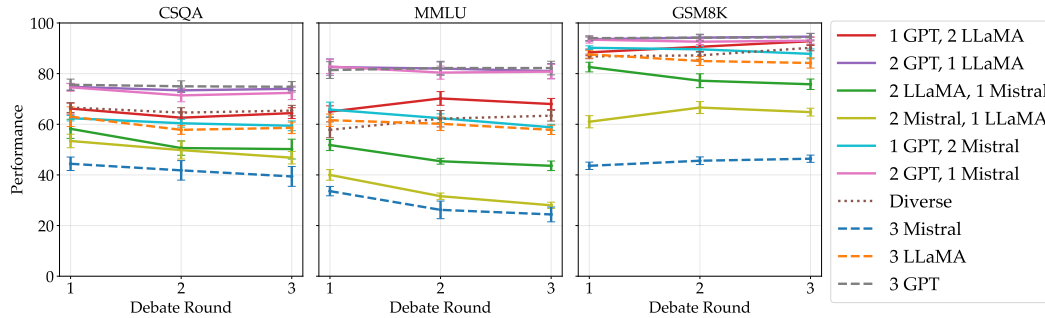





Figure 1: In many cases, we find that group accuracy frequently *degrades* over the course of debate, rather than improving performance. Diverse refers to the case (1× , 1× , 1× ).

## 6 Failure Modes of Debate

Multi-agent debate is intended to improve reasoning by leveraging disagreement—encouraging models to refine their arguments through interaction and converge on correct conclusions. Yet, debates between LLM agents often fail to reach majority agreement on the true answer, even when at least one agent is initially correct. Examining why these failures occur is crucial for both understanding and designing these systems. We show that there are many other factors influencing failure modes of debate, which we explore in depth in this section. This suggests a richer and more complex interplay of factors influencing the effectiveness of debate. In this section, we dissect these dynamics empirically, examining when and how agents revise their answers, what social conditions are correlated with undesirable answer revisions, and whether an explicit intervention targeting sycophancy can meaningfully improve debate performance.

### 6.1 Does exchange of reasoning help in sequential revision?

If an LLM agent can reflect and correct mistakes based on the reasoning of other agents, we would expect the model to improve its answer and the collective performance of the group. We know the self-correction capability of single-agent LLMs does not easily work out of the box [Huang et al., 2023], and we want to evaluate if the debate procedure helps with this or not.

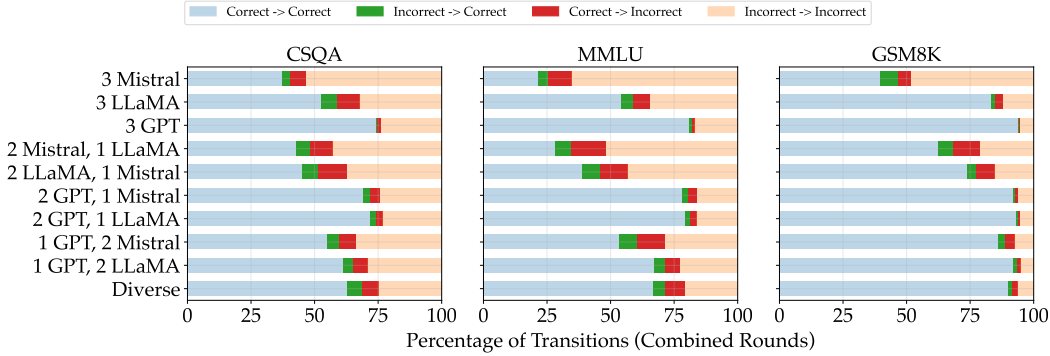


Figure 2: Breakdown of how agent change answers for different agent settings; results are aggregated over all debate rounds. We observe that most agents with incorrect initial answers do not improve their overall performance (peach bars), and, of those that do change their answers, more change from a correct answer to an incorrect one (red region) than from an incorrect to a correct one (green region).

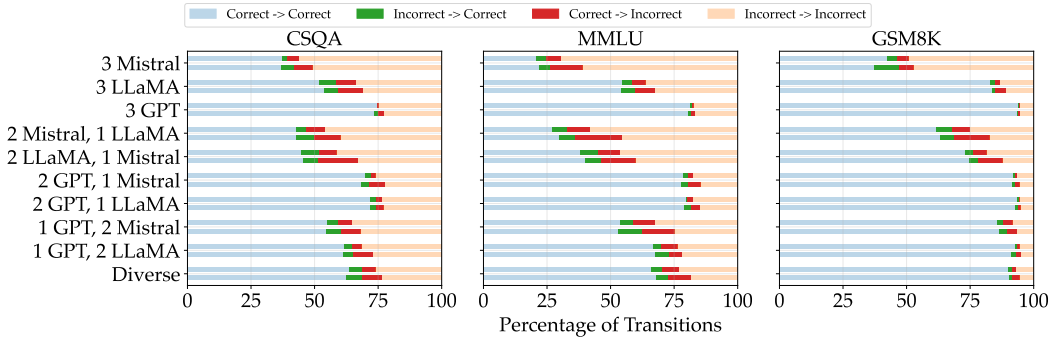


Figure 3: Breakdown of how agents change answers between debate rounds for different agent settings. The top row denotes the first round, and the row below denotes the second round. We find that the social effect dominates: agents that can resist flipping originally correct answers in round 1 have lower resistance to the social pressure from disagreement after round 2.

To assess this, we analyze how agent responses change between debate rounds across all our tasks. Note that there are four possible types of transitions: correct  $\rightarrow$  incorrect, correct  $\rightarrow$  correct, incorrect

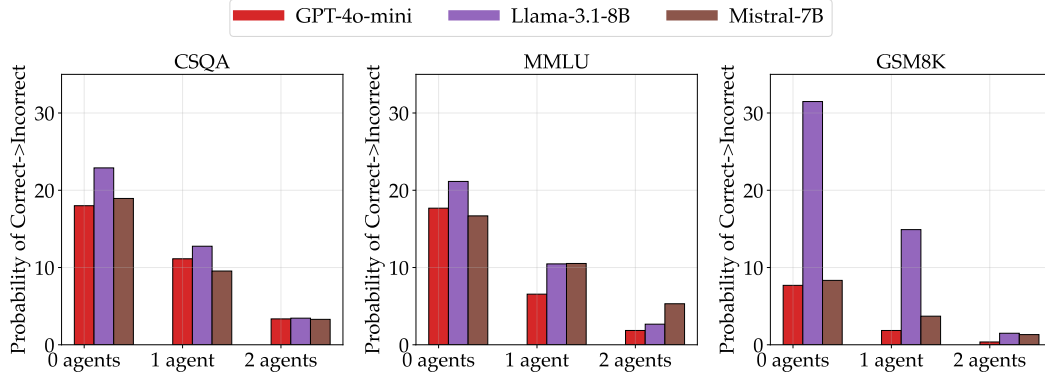


Figure 4: The likelihood of an answer being flipped from correct to incorrect (the *undesirable* flip direction), plotted against the number of agents who agree with the ego agent, averaged across all rounds of debate. We find that the number of other agents who agree with the agent appears to be correlated with the frequency with which the agents flip their answers, indicating that the models may be influenced by social effects. We further observe significant variance in this answer-flipping behavior conditioned on the specific dataset or model in question.

→ incorrect, and incorrect → correct. As shown in Figure 2, we observe that most agents with initially incorrect answers do not improve their performance (peach bars), and, of those that do change their answers, there exists a larger shift in agent responses from correct → incorrect answers (red) than incorrect → correct (green), indicating that debate can actively mislead agents who started with correct answers. Further, Figure 3 shows that the proportion of correct-to-incorrect transitions exceeds the proportion of incorrect-to-correct ones in subsequent rounds, corroborating our earlier results that debate performance degrades over rounds. We further find that there is a dominating effect of social pressure: agents that can originally resist flipping their correct answers to incorrect ones in round 1 have lower resistance to the social pressure from disagreement after round 2 (larger red bars in round 2). Together, these results highlight an important observation: to mitigate the harmful effects of debate, we need to find a way to reduce the number of *undesirable* answer flips – i.e., models changing their answers from correct → incorrect – and address the underlying social effects that cause these flips.

## 6.2 Are agents influenced by social factors?

We next proceed to investigate whether models, particularly stronger and more capable models, tend to be subject to social influence from disagreement with peers by examining the frequency of correct → incorrect answer flips. Figure 4 shows how frequently an agent changes its answer from correct → incorrect as a function of how many other agents initially agreed with it. Across all datasets, we observe that the likelihood of an undesirable answer-flip is highest when the ego agent is isolated – i.e. when no other agents agree with its answer – and decreases as more peers agree. Interestingly, this effect is highly variable across different datasets and models; for instance, on GSM8K, the strongest and weakest models (GPT and Mistral) are far less likely to make undesirable answer-flips than the third model (LLaMA). This provides strong preliminary evidence that LLM reasoning is influenced by social factors: models are sensitive to patterns of agreement and disagreement with peers, and their internal mechanisms for balancing correctness vs consensus may be quite fragile in the face of disagreement with peers.

## 6.3 Are agents sycophantic?

We next investigate another potential contributing factor to failure within debate: models may exhibit *sycophantic behavior* [Sharma et al., 2023] meaning the tendency, potentially due to RLHF [Kaufmann et al., 2024], of LLMs to prefer answers that match users’ beliefs over correct answers. Specifically, we explore whether the degradation in debate performance could be due to a sycophancy bias expressed towards other LLM agents as opposed to users (that is, generalizing from the effects



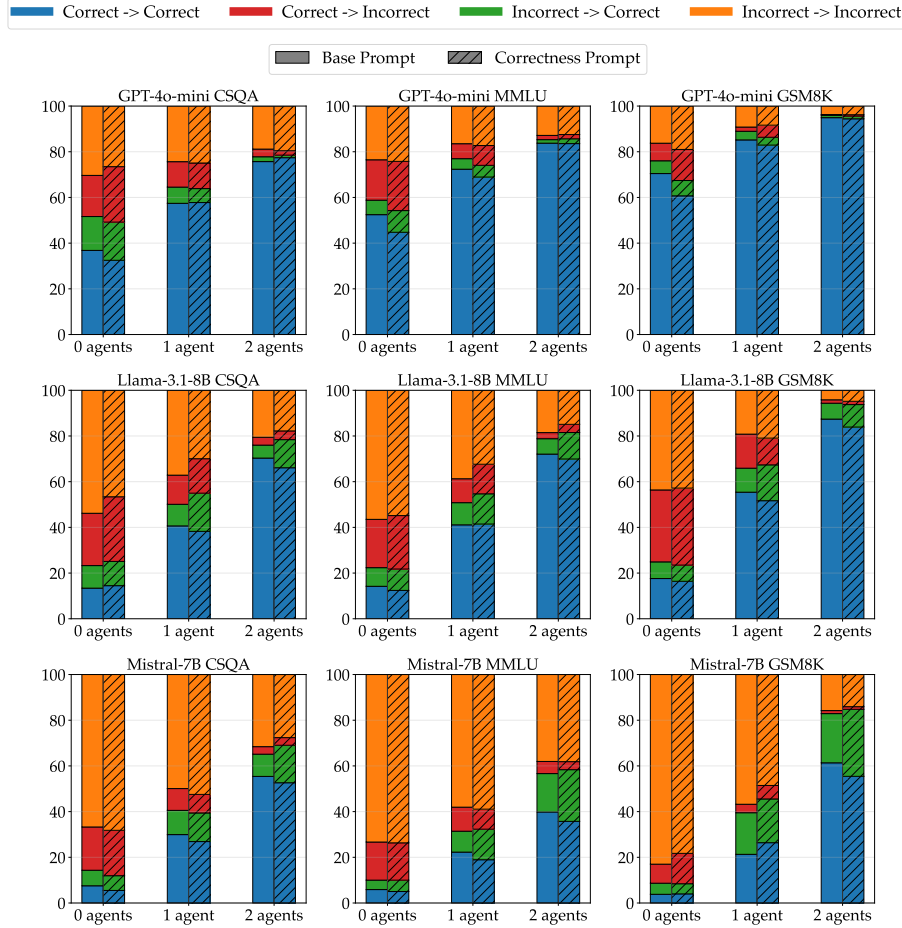


Figure 5: A comparison between the types of answer flips when agents are given the base prompt or the correctness payoff prompt. We find that introducing the correctness payoff intervention does not appear to decrease the number of undesirable **correct  $\rightarrow$  incorrect** flips – indicating that this intervention against sycophancy is insufficient on its own to resolve the issues with multi-agent debate.

of RLHF.) We consider how agents respond to seeing the reasoning and answers from other agents over the course of debate by studying when models change their answer from correct  $\rightarrow$  incorrect.<sup>3</sup>

**Correctness payoff prompt.** Political scientists have found that partisan bias effects on question-answering accuracy [Bartels, 2002] are significantly reduced when respondents are offered a monetary bonus for correct answers [Bullock et al., 2015, Prior et al., 2015]. Inspired by these results, we design a new experiment intended to encourage LLM agents to prefer correctness over agreement with other agents, countering the sycophancy bias. In particular, we define a payoff function scoring the models based on the correctness of their answers and explicitly instruct them to maximize this payoff. We modify our prompts as follows:

<sup>3</sup>Yao et al. [2025] also explore one particular form of sycophancy, showing that it does have some impact on the outcome of debate.

#### System prompt, $\mathcal{P}_{\text{system}}$

You are a helpful assistant that can answer questions and complete tasks. You will receive a fixed payoff  $p$  for each question given your response  $X$  and the true answer  $Y$ ,  $p = I[X == Y]$ , where  $I[\text{condition}]$  represents an indicator function which equals 1 if the condition is satisfied (i.e. your answer matches the true answer) and 0 otherwise. Your goal is to maximize your total payoff.

#### System prompt, $\mathcal{P}_{\text{debate}}$

These are the solutions to the problem from other agents: {AGENT\_RESPONSES} Analyze your solution and that of other agents, provide an updated answer to maximize your payoff  $p = I[X == Y]$ , and explain your reasoning. Put your answer in the form (X) at the end of your response.

We then investigate whether LLM agents are less likely to flip from correct to incorrect answers when they are prompted to maximize rewards for correct answers.

**Results.** We present results in Fig. 5. We find that adding a payoff for correctness in the model prompt does not significantly reduce the likelihood that LLM agents flip their answers from correct to incorrect, regardless of the number of other agents that agree with them. In fact, we find that in many cases, the number of correct  $\rightarrow$  incorrect transitions actually *increases* when using the correctness-payoff prompt – a counter-intuitive result that seems to suggest that asking models to prioritize correctness may not help resolve the issues observed in multi-agent debate.

## 6.4 Influence of Model and Task Type

Across all our analyses (Figures 2, 4, and 5), we observe substantial variation in behavior across individual models and tasks. Models of different individual abilities exhibit distinct patterns of answer change and sensitivity to peer disagreement, suggesting that an agent’s capability itself can influence how social dynamics unfold. Likewise, the magnitude and direction of these effects differ across tasks such as CommonSenseQA, MMLU, and GSM8K, implying that task structure and domain complexity also shape how debates evolve. Together, these findings indicate that no single mechanism, such as sycophancy alone, can fully explain debate failures: rather, the interplay between model capability, task characteristics, and social influence jointly determines how and when agents revise their beliefs.

## 7 Discussion

Our findings challenge the prevailing and natural view that deliberation among AI agents will always improve reasoning. In fact, our experiments reveal that multi-agent debate can sometimes degrade performance: group accuracy often declines over successive rounds of debate. This counterintuitive trend holds even when the majority of agents perform well individually on the task and even when weaker models have access to the reasoning of stronger models. In other words, additional exchange of reasons between agents does not always correct mistakes; instead, it may amplify them.

We study a number of correlated and contributing factors to these failure modes in debate, including sequential revision, social conditioning, and sycophancy. We additionally show across all experiments how models of different capabilities respond differently within multi-agent debate, and how the type of task or dataset can also have a significant influence on the results of debate. We find that beyond a simple single cause, multiple factors exist that likely contribute to failure modes in multi-agent debate, with the net effect that heterogeneous groups of models frequently converge on wrong answers together, potentially negating the benefits of multi-agent debate.

These results suggest that naive debate protocols can inadvertently amplify errors rather than correcting them, underscoring the need for more principled approaches to multi-agent debate. Future debate frameworks should incorporate mechanisms that promote critical evaluation over consensus – such as encouraging agents to assess the soundness of others’ reasoning, integrating confidence estimates

or credibility scores to weight contributions by expertise, and rewarding independent verification of claims. Training or incentive schemes might further discourage superficial agreement by penalizing unsupported conformity. Ultimately, improving the robustness of debate requires aligning models not merely to achieve consensus, but to engage in constructive epistemic disagreement – challenging peers when warranted and maintaining justified beliefs under social pressure.

## Acknowledgments and Disclosure of Funding

Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute [www.vectorinstitute.ai/#partners](http://www.vectorinstitute.ai/#partners).

## References

- Mahak Agarwal and Divyam Khanna. When persuasion overrides truth in multi-agent llm debates: Introducing a confidence-weighted persuasion override rate (cw-por), 2025. URL <https://arxiv.org/abs/2504.00374>.
- Alfonso Amayuelas, Xianjun Yang, Antonis Antoniadis, Wenyue Hua, Liangming Pan, and William Wang. Multiagent collaboration attack: Investigating adversarial attacks in large language model collaborations via debate. *arXiv preprint arXiv:2406.14711*, 2024.
- Larry M. Bartels. Beyond the running tally: Partisan bias in political perceptions. *Political Behavior*, 24(2):117–150, 2002. ISSN 01909320, 15736687. URL <http://www.jstor.org/stable/1558352>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- John G. Bullock, Alan S. Gerber, Seth J. Hill, and Gregory A. Huber. Partisan bias in factual beliefs about politics. *Quarterly Journal of Political Science*, 10(4):519–578, 2015. ISSN 1554-0626. doi: 10.1561/100.00014074. URL <http://dx.doi.org/10.1561/100.00014074>.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*, 2023.
- 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. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate, 2023. URL <https://arxiv.org/abs/2305.14325>.
- Andrew Estornell and Yang Liu. Multi-llm debate: Framework, principals, and interventions. *Advances in Neural Information Processing Systems*, 37:28938–28964, 2024.
- Andrew Estornell, Jean-Francois Ton, Yuanshun Yao, and Yang Liu. Acc-collab: An actor-critic approach to multi-agent llm collaboration, 2025. URL <https://arxiv.org/abs/2411.00053>.
- Aaron Grattafiori et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. Improving language model negotiation with self-play and in-context learning from ai feedback. *arXiv preprint arXiv:2305.10142*, 2023.
- Chengbo He, Bochao Zou, Xin Li, Jiansheng Chen, Junliang Xing, and Huimin Ma. Enhancing llm reasoning with multi-path collaborative reactive and reflection agents, 2025. URL <https://arxiv.org/abs/2501.00430>.

- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021. URL <https://arxiv.org/abs/2009.03300>.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. *arXiv preprint arXiv:2310.01798*, 2023.
- Geoffrey Irving, Paul Christiano, and Dario Amodei. Ai safety via debate, 2018. URL <https://arxiv.org/abs/1805.00899>.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han, and Rui Zhang. When can llms actually correct their own mistakes? a critical survey of self-correction of llms. *Transactions of the Association for Computational Linguistics*, 12:1417–1440, 2024.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke H  llermeier. A survey of reinforcement learning from human feedback, 2024. URL <https://arxiv.org/abs/2312.14925>.
- Zachary Kenton, Noah Y. Siegel, J  nos Kram  r, Jonah Brown-Cohen, Samuel Albanie, Jannis Bulian, Rishabh Agarwal, David Lindner, Yunhao Tang, Noah D. Goodman, and Rohin Shah. On scalable oversight with weak llms judging strong llms, 2024. URL <https://arxiv.org/abs/2407.04622>.
- Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R. Bowman, Tim Rockt  schel, and Ethan Perez. Debating with more persuasive llms leads to more truthful answers, 2024. URL <https://arxiv.org/abs/2402.06782>.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society, 2023. URL <https://arxiv.org/abs/2303.17760>.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*, 2023.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.
- Julian Michael, Salsabila Mahdi, David Rein, Jackson Petty, Julien Dirani, Vishakh Padmakumar, and Samuel R. Bowman. Debate helps supervise unreliable experts, 2023. URL <https://arxiv.org/abs/2311.08702>.
- OpenAI. Gpt-4o mini: Advancing cost-efficient intelligence. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>, 2024.
- Markus Prior, Gaurav Sood, and Kabir Khanna. You cannot be serious: The impact of accuracy incentives on partisan bias in reports of economic perceptions. *Quarterly Journal of Political Science*, 10(4):489–518, December 2015. doi: 10.1561/100.00014127. URL <https://ideas.repec.org/a/now/jlqjps/100.00014127.html>.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R. Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. Towards understanding sycophancy in language models, 2023. URL <https://arxiv.org/abs/2310.13548>.

- Vighnesh Subramaniam, Yilun Du, Joshua B Tenenbaum, Antonio Torralba, Shuang Li, and Igor Mordatch. Multiagent finetuning: Self improvement with diverse reasoning chains. *arXiv preprint arXiv:2501.05707*, 2025.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge, 2019. URL <https://arxiv.org/abs/1811.00937>.
- Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O’Sullivan, and Hoang D Nguyen. Multi-agent collaboration mechanisms: A survey of llms. *arXiv preprint arXiv:2501.06322*, 2025.
- Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities. *arXiv preprint arXiv:2406.04692*, 2024.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023. URL <https://arxiv.org/abs/2203.11171>.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation, 2023. URL <https://arxiv.org/abs/2308.08155>.
- Hanqing Yang, Jingdi Chen, Marie Siew, Tania Lorido-Botran, and Carlee Joe-Wong. Llm-powered decentralized generative agents with adaptive hierarchical knowledge graph for cooperative planning, 2025. URL <https://arxiv.org/abs/2502.05453>.
- Binwei Yao, Chao Shang, Wanyu Du, Jianfeng He, Ruixue Lian, Yi Zhang, Hang Su, Sandesh Swamy, and Yanjun Qi. Peacemaker or troublemaker: How sycophancy shapes multi-agent debate, 2025. URL <https://arxiv.org/abs/2509.23055>.