

# Demystifying Multi-Agent Debate: The Role of Confidence and Diversity

Xiaochen Zhu<sup>\*</sup> Caiqi Zhang<sup>\*</sup> Yizhou Chi<sup>\*</sup> Tom Stafford<sup>†</sup> Nigel Collier<sup>†</sup> Andreas Vlachos<sup>†</sup>  
<sup>†</sup>University of Cambridge <sup>†</sup>University of Sheffield  
{xz479, cz391, yc697, nhc30, av308}@cam.ac.uk  
t.stafford@sheffield.ac.uk

## Abstract

Multi-agent debate (MAD) is widely used to improve large language model (LLM) performance through test-time scaling, yet recent work shows that vanilla MAD often underperforms simple majority vote despite higher computational cost. Studies show that, under homogeneous agents and uniform belief updates, debate preserves expected correctness and therefore cannot reliably improve outcomes. Drawing on findings from human deliberation and collective decision-making, we identify two key mechanisms missing from vanilla MAD: (i) diversity of initial viewpoints and (ii) explicit, calibrated confidence communication. We propose two lightweight interventions. First, a diversity-aware initialisation that selects a more diverse pool of candidate answers, increasing the likelihood that a correct hypothesis is present at the start of debate. Second, a confidence-modulated debate protocol in which agents express calibrated confidence and condition their updates on others' confidence. We show theoretically that diversity-aware initialisation improves the prior probability of MAD success without changing the underlying update dynamics, while confidence-modulated updates enable debate to systematically drift to the correct hypothesis. Empirically, across six reasoning-oriented QA benchmarks, our methods consistently outperform vanilla MAD and majority vote. Our results connect human deliberation with LLM-based debate and demonstrate that simple, principled modifications can substantially enhance debate effectiveness.

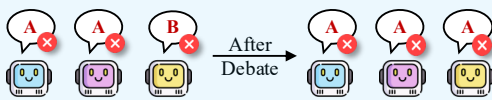
## 1 Introduction

Multi-agent debate (MAD) has rapidly become a popular technique for improving large language models' (LLMs) reasoning and test-time scaling (TTS) (Liang et al., 2024; Du et al., 2024; Yang

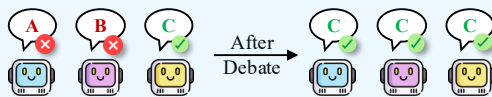
<sup>\*</sup>Equal contribution. Our code is available at: <https://github.com/SpaceHunterInf/DMAD>

### (a) The Role of Initial Diversity

Scenario 1: **Low Diversity** is susceptible to limitations such as conformity and local optima!



Scenario 2: **High Diversity** increases the probability that the correct answer is in the pool!



### (b) The Role of Confidence Communication

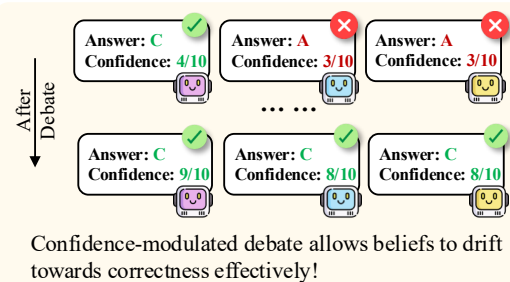


Figure 1: Illustration of two human-inspired ingredients for effective MAD: (a) diverse initial answers increase the chance the correct hypothesis is present; (b) explicit confidence sharing enables confidence-weighted updates that steer beliefs toward the correct answer.

et al., 2025c; Guo et al., 2024). Inspired by studies of group deliberation, MAD operationalises the idea that independent agents exchanging reasons should, in principle, outperform individual sampling or simple voting (Mercier and Sperber, 2011; Moshman and Geil, 1998).

Despite growing enthusiasm and substantial token budgets devoted to multi-agent debate (MAD), its benefits are far from guaranteed (Choi et al., 2025; Wu et al., 2025). In practice, vanilla MAD often underperforms a simple majority vote despite incurring substantially higher computational cost. Choi et al. (2025) formally characterise this limitation by showing that, under homogeneous agents

and unweighted belief updates, debate preserves expected correctness over time. Mathematically, such behaviour is described as a martingale over belief trajectories. In other words, while agents may exchange arguments, the probability of converging to the correct answer does not systematically increase.

Previous work in cognitive science has identified several factors that affect the success of human deliberation:

1. **Diversity of initial viewpoints.** Human group deliberation benefits from heterogeneous priors and perspectives (Moshman and Geil, 1998). Exposure to differing ideas increases the search space over possible solutions and thus benefits the deliberation (Smaldino et al., 2024). Empirical work shows that the diversity of positions raised during discussion is positively correlated with performance gains (Karadzhov et al., 2024).
2. **Confidence communication and alignment.** In human discussion, participants do not only share answers; they also communicate uncertainty, negotiate and align their confidence expressions, and use these signals to weight others’ contributions when forming collective decisions (Bahrami et al., 2010; Dezeccache et al., 2022; Mercier and Sperber, 2017).

These two mechanisms are largely absent in standard MAD implementations. Firstly, agents are sampled from the same or similar models without regard for diversity, a problem exacerbated by post-training and alignment that can induce diversity collapse (Padmakumar and He, 2024; Kirk et al., 2024). Secondly, unlike humans who convey confidence through tone, prosody, and facial expression (Guyer et al., 2021), LLMs can only communicate confidence in an explicitly verbalised manner. Such signals are rarely requested or used in current MAD protocols, and LLMs are known to exhibit overconfident expressions even when they are wrong (Sun et al., 2025).

Recent work has started to analyse these failure modes and the role of diversity in MAD. Estornell and Liu (2024) show that when agents lack diversity, debate dynamics can quickly become static, causing the procedure to collapse back to the initial majority. Wynn et al. (2025) and Wu et al. (2025) empirically validate that, in such settings, MAD often provides little or no gain over majority vote.

However, these works mostly diagnose the problem; they do not offer mechanisms that can reliably increase diversity and thereby improve debate outcomes. In parallel, several studies attempt to incorporate numerical confidence expressions into debates (Yoffe et al., 2024; Lin and Hooi, 2025) and report promising improvements. Yet their implementations typically rely on post-hoc calibration and focus only on having agents state their confidence, while overlooking how agents interpret and use others’ confidence, which is a central aspect of human group decision-making (Zarnoth and Sniezek, 1997; Fusaroli et al., 2012).

In this paper, we examine MAD through the lens of these two missing ingredients: **diversity** and **confidence**. By analysing debate dialogues, we show that initial answer diversity correlates with performance gains with statistical significance, and that MAD is more effective on harder datasets where such diversity arises naturally. Following this, we introduce two simple interventions. First, a **diversity-aware initialization** that selects a more diverse subset of candidate answers as the initial debate pool, improving performance without additional training. Second, a **confidence-modulated debate protocol** in which agents express calibrated confidence and condition their updates on others’ confidence. Crucially, we view them as complementary mechanisms acting on different stages of deliberation: diversity initialisation affects the initial belief distribution (*i.e.*, what hypotheses are present), while confidence affects the aggregation dynamics (*i.e.*, how each contribution influences others). We provide theoretical guarantees showing that diversity increases the prior probability of the debate success, while confidence-weighted updates break the martingale limitation and allow beliefs to drift toward correctness. Empirically, our methods consistently outperform vanilla MAD and majority vote across six reasoning-oriented question-answering (QA) benchmarks. Conceptually, this work integrates human deliberation theory with LLM-based MAD, offering a principled explanation and practical guidance for designing more effective multi-agent systems.

## 2 Related Work

**Multi-Agent Debate.** Human group deliberation often improves collective performance because discussion allows individuals to pool information, detect mistakes of each other, and benefit from

diverse cognitive perspectives (Kerr and Tindale, 2004). MAD aims to replicate these advantages with LLMs, revising low-quality responses, reducing hallucinations, and eliciting reasoning that may be overlooked by a single model (Chan et al., 2024; Liang et al., 2024; Estornell and Liu, 2024). The standard MAD protocol asks each agent to answer the same question and then revise their response given previous agents' arguments (Yang et al., 2025c). While MAD sometimes yields modest gains over single-model sampling, a growing body of work shows that it often fails to inherit the benefits of human debate. Several studies report that MAD can perform worse than simple majority vote, especially when debates extend over multiple rounds (Wynn et al., 2025; Wu et al., 2025). Wu et al. (2025) argue that, unlike human groups, LLM agents struggle to recognise and disregard low-quality contributions, causing noise to propagate rather than dissipate. Most notably, Choi et al. (2025) show that homogeneous LLM agents produce debate dynamics that behave empirically like a martingale over agents' expressed beliefs; they further prove that MAD cannot surpass the accuracy of majority vote in expectation. These findings reveal a fundamental discrepancy between human deliberation and LLM debate and motivate reconsideration of the assumptions and protocols underlying MAD.

**Diversity in Collective Problem Solving.** Diversity is a central factor in the success of human deliberation. Psychological work shows that heterogeneous viewpoints help mitigate conformity pressure and “groupthink,” where groups converge to majority opinions regardless of correctness (Janis, 1972; Asch, 1955). Maintaining a diversity of candidate solutions and exploring disagreement improves group reasoning outcomes (Karadzhev et al., 2024; Smaldino et al., 2024). MAD exhibits similar vulnerabilities. Without sufficient diversity, debates collapse to the initial majority and fail to recover from incorrect early answers (Zhu et al., 2025; Estornell and Liu, 2024). The problem is exacerbated by properties of contemporary LLMs: post-training alignment reduces sampling diversity (Kirk et al., 2024; Padmakumar and He, 2024; Hu et al., 2025), and models often exhibit sycophantic tendencies toward user prompts (Sharma et al., 2024). Recent analyses of MAD suggest that greater diversity in initial reasoning is associated with improved performance (Wynn et al., 2025; Liang et al., 2024; Estor-

nell and Liu, 2024). However, this line of work primarily attributes diversity to heterogeneity across model families, parameterisations, or prompting styles. As a result, it remains unclear what fundamentally governs answer diversity or how it can be reliably increased in a simple, training-free manner. For example, Liang et al. (2024) and Choi et al. (2025) encourage diversity by prompting agents with different personas, while Estornell and Liu (2024) propose a KL-divergence-based pruning strategy to remove redundant answers during debate. These approaches alter prompting or debate structure but do not directly increase the diversity of the initial answer pool, nor do they explicitly target the prior probability of including a correct hypothesis.

**Confidence, Calibration, and Uncertainty Communication.** In human deliberation, confidence sharing is crucial for effective collaboration: individuals communicate uncertainty, align and adjust confidence expressions through interaction, and use others' confidence cues to weight arguments when forming joint decisions (Bahrami et al., 2010; Fusaroli et al., 2012; Mercier, 2020). Such communication is facilitated by paralinguistic signals, tone, prosody, facial expression (Guyer et al., 2021), which are unavailable in text-only LLM settings. Although recent work has sought to teach LLMs to proactively express uncertainty during generation (Zhang et al., 2025b; Yang et al., 2025b,a), these methods have not yet been systematically studied in debate or multi-turn settings (Zhang et al., 2026). Several recent studies instead introduce numerical confidence into MAD. Lin and Hooi (2025) and Yoffe et al. (2024) use post-hoc calibration methods to assign confidence scores to agents' answers and use these scores as additional signals during debate. While these approaches yield limited improvements, they diverge from psychological theories that emphasise aligned confidence communication and shared metacognition. Moreover, they overlook the gap between confidence expression and perception in LLMs: models are known to be systematically overconfident in their verbalised confidence judgments (Sun et al., 2025), and there is no guarantee that LLM agents interpret or appropriately weight others' confidence. As a result, current confidence-based approaches do not fully capture the mechanisms that make human deliberation successful.

### 3 Formulation

Following prior theoretical analyses of multi-agent debate (MAD), we adopt the framework of [Choi et al. \(2025\)](#) and consider a controlled setting with homogeneous LLM agents performing question answering (QA). This abstraction isolates the structural properties of debate dynamics without conflating them with model heterogeneity or explicit role specialisation.

Let  $\{a_1, \dots, a_N\}$  denote  $N$  agents. For an input question  $x \in \mathcal{X}$ , each agent samples an answer

$$y_i \sim a_i(x), \quad a_i(x) := f(x; \theta_i), \quad (1)$$

where  $f$  is the conditional distribution defined by an LLM with parameters  $\theta_i$ . In the homogeneous setting, all agents share identical parameters ( $\theta_i = \theta$ ), and stochasticity arises solely from sampling.

**Majority Vote.** Given an initial pool of sampled responses  $\{y_i\}_{i=1}^N$ , the majority vote baseline selects

$$V(\{y_i\}_{i=1}^N) = \arg \max_y \text{freq}(y). \quad (2)$$

This corresponds to test-time ensembling without interaction and provides an important empirical baseline for vanilla MAD.

**Multi-Agent Debate.** A debate consists of  $T$  rounds of iterative revision. At round  $t = 1$ , each agent independently samples an initial answer  $y_{i,1} \sim a_i(x)$ . For each subsequent round  $t > 1$ , every agent observes all the answers produced by other agents and itself in the previous round:

$$\mathcal{R}_t = \{y_{j,t-1} \mid j = 1, \dots, N\}. \quad (3)$$

Each agent then revises its answer by applying an answer-update operator

$$y_{i,t} = \mathcal{D}(x, \mathcal{R}_t), \quad (4)$$

where  $\mathcal{D}$  is instantiated as a prompting-based revision operator executed by the LLM.

After  $T$  rounds, the ensemble-level MAD output is obtained by majority vote over the terminal responses:

$$y_{\text{MAD}} = V(\{y_{i,T}\}_{i=1}^N). \quad (5)$$

### 4 Methodology

Incentivized by the role of diversity and confidence in human deliberation, we incorporate diversity-aware initialization and confidence modulated debate as additional features on top of vanilla MAD.

#### 4.1 Diversity-aware Initialization

Studies of human group deliberation suggest that considering a more diverse set of candidate solutions is correlated with higher collective answer correctness ([Karadzhov et al., 2024](#)). One proposed explanation is that broader exploration of the solution space increases the likelihood that at least one high-quality hypothesis is discovered, rather than converging prematurely to suboptimal consensus ([Smaldino et al., 2024](#)). We adopt this intuition as a design principle for initializing multi-agent debate.

Formally, given a set of answers  $S$ , we define its diversity as the number of distinct answers:

$$\text{div}(S) = |\text{unique}(S)|. \quad (6)$$

Vanilla MAD initializes debate by independently sampling  $N$  answers  $\{y_{i,1}\}_{i=1}^N$ . Instead, we first sample a larger pool of  $N_{\text{cand}} \geq N$  candidate answers  $\{y^{(1)}, \dots, y^{(N_{\text{cand}})}\}$ , and select a subset of size  $N$  to initialise the debate by solving

$$S_{\text{div}} = \arg \max_{S \subseteq \{1, \dots, N_{\text{cand}}\}, |S|=N} \text{div}(S). \quad (7)$$

We apply a simple, training-free greedy approximation that iteratively selects the candidate with the largest marginal contribution to diversity until  $N$  answers are chosen. Compared to vanilla MAD, this procedure incurs additional test-time cost by sampling  $N_{\text{cand}}$  rather than  $N$  initial answers, but requires no model training or architectural modification. Unlike approaches that prune similar responses within a single debate round ([Estornell and Liu, 2024](#)), our method preserves the standard MAD protocol and increases the probability that the initial debate state contains at least one correct hypothesis.

#### 4.2 Confidence-Modulated Debate

[Bahrami et al. \(2010\)](#) show that when debate participants can directly communicate their uncertainty, joint decisions improve. In human deliberation, such uncertainty cues are often conveyed through paralinguistic signals (e.g., facial expressions and tone) ([Guyer et al., 2021](#)). However, these cues are unavailable in text-based LLM debates, where confidence must be expressed verbally. Moreover, estimating confidence for black-box models such as LLMs is itself non-trivial. While confidence can be inferred indirectly through sampling-based uncertainty estimates ([Zhang et al., 2024, 2025a](#)) or linguistic markers ([Yang et al., 2025a,b](#)), we focus on verbalised numerical confidence, as it is the

most direct signal agents can exchange during debate and can be explicitly calibrated (Zhang et al., 2025b).

To incorporate confidence into debate, we extend each agent’s output at round  $t$  to

$$y_{i,t}^{\text{conf}} = (y_{i,t}, w_{i,t}), \quad (8)$$

where  $w_{i,t} \in \{0, \dots, 10\}$  is a discrete confidence score indicating the agent’s self-assessed certainty, where 0 indicates the agent is unsure about its answer and 10 indicates total confidence. The debate-update operator becomes

$$y_{i,t}^w = \mathcal{D}^w(x, \mathcal{R}_t^w), \quad (9)$$

where  $\mathcal{R}_t^w$  contains each agent’s answer and confidence.

We use the term *modulated* to emphasise that confidence does not change the content of an agent’s answer, but instead scales its influence during aggregation, analogous to how modulation alters the strength of a carrier signal without changing the underlying message (Crecraft and Gergely, 2002). However, leveraging such confidence signals introduces two **key challenges**. First, verbally reported scores are often miscalibrated: LLMs may express high confidence even when incorrect (Sun et al., 2025). Second, even with calibrated scores, agents may fail to *use* confidence effectively during interaction—*i.e.*, they may ignore confidence cues or update their beliefs inappropriately. We therefore (i) train agents to express calibrated numerical confidence, and (ii) teach agents to perceive and exploit confidence signals during debate via reinforcement learning.

**Calibrated Confidence Expression.** Rather than relying on post-hoc calibration (Lin and Hooi, 2025; Yoffe et al., 2024), we directly train the model to generate calibrated confidence scores during debate using reinforcement learning (RL).

Similar to Zhang et al. (2025b) and Stangel et al. (2025), we use a binary cross-entropy style log-based reward encouraging alignment between correctness and expressed confidence. Let  $z_{i,t} \in \{0, 1\}$  denote the correctness indicator of  $y_{i,t}$  and map  $w_{i,t}$  to  $(0, 1)$ , we adopt the following reward:

$$r^{\text{conf}}(y_{i,t}, w_{i,t}) = \frac{\lambda}{R_{\max}} [z_{i,t} \log(w_{i,t}) + (1 - z_{i,t}) \log(1 - w_{i,t})]. \quad (10)$$

$\lambda, R_{\max}$  are scaling coefficients. This encourages the model to assign higher confidence to correct predictions and lower confidence to incorrect ones.

**Confidence Perception and Usage.** Expression alone does not guarantee that agents use confidence appropriately during debate. As emphasised in Fusaroli et al. (2012), effective deliberation requires confidence alignment, where individuals adjust their reasoning strategies based on others’ confidence signals. To train LLM agents to perceive and use confidence, we augment the reward with correctness signal and apply RL to the debate-update operator:

$$r^{\text{total}} = \lambda_1 z_{i,t} + \lambda_2 r^{\text{conf}}(y_{i,t}, w_{i,t}), \quad (11)$$

where  $\lambda_1, \lambda_2$  control the trade-off between accuracy and calibration. In this case, the learned policy is encouraged to explicitly condition on others’ confidence. This captures both calibrated expression and active alignment and usage of confidence, two components necessary for mirroring human deliberation dynamics.

## 5 Theoretical Analysis

We now analyse how does our interventions improve over the vanilla MAD dynamics of Choi et al. (2025). Unlike prior work that studies diversity or confidence in isolation, we analyse them as complementary mechanisms acting on different stages of the same deliberative process. We work within their Dirichlet–categorical model (DCM) for MAD, and state only the ingredients needed for our results; a full description is given in the Appendix D.1.

Let  $A = \{1, \dots, K\}$  be the finite answer set and assume, without loss of generality, that option 1 is the unique correct answer. At debate round  $t$ , agent  $i$  is parameterised by a Dirichlet vector

$$\alpha_{i,t} = (\alpha_{i,t}^{(1)}, \dots, \alpha_{i,t}^{(K)}) \in \mathbb{R}_{>0}^K, \quad (12)$$

and generates its answer  $y_{i,t} \in A$  by

$$\theta_{i,t} \sim \text{Dir}(\alpha_{i,t}), \quad y_{i,t} \mid \theta_{i,t} \sim \text{Cat}(\theta_{i,t}).$$

We write

$$p_{i,t} := \Pr(y_{i,t} = 1 \mid \alpha_{i,t}) = \frac{\alpha_{i,t}^{(1)}}{\sum_{k=1}^K \alpha_{i,t}^{(k)}} \quad (13)$$

for agent  $i$ ’s belief on the correct option at round  $t$ .

In the given DCM model, agents update their Dirichlet parameters by adding one count for each agent’s answer in the previous round. Let  $\mathcal{F}_t$  denote the filtration generated by all belief parameters up to round  $t$ . In the homogeneous, fully connected

setting where all agents share the same prior and observe the same multiset of answers, Choi et al. (2025) show that

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] = p_{i,t-1}, \quad (14)$$

so  $\{p_{i,t}\}_{t>0}$  is a martingale. This means during the debate, the expected belief at the next round always equals to the current belief. In expectation, vanilla MAD neither helps nor hurts the performance.

### 5.1 Diversity Improves What Is Debated

Diversity-aware initialisation affects the support of the debate by increasing the likelihood that the initial answer pool contains at least one correct hypothesis. Importantly, this intervention operates entirely at initialisation and does not modify the subsequent debate dynamics.

**Proposition 1** (Diversity-aware initialization improves prior success). *Let  $A_T$  be the event that debate outputs the correct answer at the final round  $T$  under the unweighted DCM dynamics. Let  $S$  denote the number of distinct informative hypotheses in the initial answer pool (e.g., the number of distinct options that are sampled at least once). Suppose that the conditional success probability  $\mathbb{P}(A_T \mid S = s)$  is nondecreasing in  $s$ , and that the distribution of  $S$  under our diversity-aware initialiser first-order stochastically dominates that under random i.i.d. sampling. Then*

$$\mathbb{P}(A_T \mid \text{diverse init}) \geq \mathbb{P}(A_T \mid \text{random init}).$$

Proposition 1 formalises that diversity-aware initialization improves MAD by shifting the distribution of the initial debate state toward pools with broader hypothesis coverage. Conditional on this state, the unweighted debate dynamics remain a martingale, but the prior probability that debate begins with a useful hypothesis is strictly increased. A full proof is provided in Appendix D.2.

### 5.2 Confidence Improves How Debate Aggregates

Confidence-modulated debate changes how information is aggregated. Answers that are both high-confidence and positively correlated with correctness receive more weight in the Dirichlet updates, causing the expected belief in the correct answer to drift upward over rounds instead of remaining

flat as in the unweighted martingale case. Consider the DCM debate model described above and extend it so that, at each round, agents exchange answers together with scalar confidence scores and update their Dirichlet parameters using confidence-weighted counts.

**Theorem 1** (Confidence-weighted debate yields a submartingale). *Assume that: (i) agents are homogeneous and fully connected; and (ii) confidence is positively correlated with correctness, in the sense that higher-confidence answers are, on average, more likely to be correct than lower-confidence ones. Let  $p_{i,t}$  denote agent  $i$ 's belief on the correct answer at round  $t$ , and let  $\mathcal{F}_t$  be the filtration generated by the debate history. Then the belief process  $\{p_{i,t}\}_{t \geq 0}$  becomes a strict submartingale:*

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] \geq p_{i,t-1},$$

with strict inequality on a set of positive probability.

Conceptually speaking, confidence-modulated debate breaks the martingale symmetry of vanilla MAD and, in expectation, strictly improves correctness over rounds. Whereas a martingale preserves the expected belief in the correct answer, a submartingale has increasing expected belief, implying systematic progress toward the correct hypothesis. Appendix D.3 defines the confidence-weighted Dirichlet-categorical update, and Appendix D.4 proves that this yields a (strict) submartingale.

## 6 Experiment Setup

We first conduct a series of experiments to study the dynamics of vanilla MAD and provide additional analysis under the setting of Choi et al. (2025). Specifically, we initialise 5 homogeneous agents using the same language model, and run debates for 5 turns. To obtain a comprehensive understanding, we evaluate MAD across a wide range of tasks using multiple language models.

**Datasets.** We cover a series of QA datasets with varying levels of difficulty. Our in-domain datasets include GSM8K (Cobbe et al., 2021), CommonsenseQA (Talmor et al., 2019), HellaSwag (Zellers et al., 2019), and MMLU (Hendrycks et al., 2021), which contains 57 sub-tasks across diverse knowledge domains. We perform confidence-modulated debate training on these datasets. To assess the

Model Strategy	GSM8K	CSQA	HellaSwag	MMLU-FL	MMLU-PM	MMLU-All	GPQA	ARC-C
<i>Qwen-2.5-7B-Instruct</i>								
Single Model	0.908	0.784	0.752	0.547	0.765	0.735	0.298	0.805
Debate 5 × 5	0.847	0.727	0.698	0.571	0.734	0.727	0.339	0.882
Majority Vote	0.908	0.808	0.762	0.548	0.806	0.764	0.335	0.880
High Diversity	0.920	<b>0.836</b>	0.779	0.587	0.822	0.743	0.334	0.884
HD + Learn2Agg	0.893	0.776	0.777	0.603	0.805	0.752	0.353	0.889
Confidence	0.910	0.810	<b>0.782</b>	0.587	0.818	<b>0.769</b>	0.350	<b>0.911</b>
HD + Confidence	<b>0.932</b>	0.828	0.781	<b>0.619</b>	<b>0.825</b>	<b>0.769</b>	<b>0.361</b>	0.890
<i>Llama-3.1-8B-Instruct</i>								
Single Model	0.797	0.703	0.676	0.503	0.798	0.701	0.266	0.780
Debate 5 × 5	0.850	0.688	0.667	0.544	0.772	0.680	0.316	0.802
Majority Vote	0.870	0.718	0.686	0.532	0.809	0.711	0.290	0.832
High Diversity	0.863	0.721	<b>0.701</b>	0.548	0.805	0.687	0.321	0.877
HD + Learn2Agg	<b>0.891</b>	0.747	0.677	0.551	0.797	0.700	0.350	0.832
Confidence	0.882	0.734	0.685	0.548	0.818	<b>0.720</b>	<b>0.352</b>	<b>0.882</b>
HD + Confidence	0.890	<b>0.752</b>	0.688	<b>0.555</b>	<b>0.842</b>	0.713	<b>0.352</b>	0.854

Table 1: Benchmark performance for Qwen-2.5-7B-Instruct and Llama-3.1-8B-Instruct across major datasets. Cell background indicates improvement (green) relative to the Majority Vote baseline within each model group; color intensity scales with the magnitude of the change. Best results per model group are **bolded**. We specifically look into MMLU Formal Logic (FL) and Professional Medicine (PM) to keep consistent with Choi et al. (2025).

Model	r	p
Qwen-2.5-7B-Instruct	0.070	<0.001
Llama-3.1-8B-Instruct	0.019	0.002

Table 2: Pearson correlation (r) and significance levels (p) for initial answer diversity vs. debate performance gain (n=26,624).

generalization capabilities of our methods, we additionally evaluate on GPQA-Main (Rein et al., 2024) and ARC-Challenge (Clark et al., 2018), which we treat as out-of-domain (OOD) benchmarks.

**Language Models and Prompts.** We evaluate Llama-3.1-8B-Instruct (Dubey et al., 2024) and Qwen-2.5-7B-Instruct (Team, 2024) for both vanilla MAD and our trained variants. In the MAD setting, we elicit chain-of-thought (CoT) reasoning (Wei et al., 2022) by prompting the model to output its reasoning steps followed by its final answer. To incorporate confidence, we additionally prompt the model to produce a numerical confidence value. Detailed prompts are provided in Appendix A.

**Sampling and Training Schema.** For agent initialisation and all inference steps, we apply temperature 1.0 and nucleus sampling with  $p = 0.9$  (Holtzman et al., 2020). For diversity injection, we set  $N_{\text{cand}} = 10$ . For confidence expression calibration and perception training, we use the GRPO algorithm combined with LoRA adaptation (Shao et al., 2024; Hu et al., 2022). Further details on the sampling strategy and training configuration

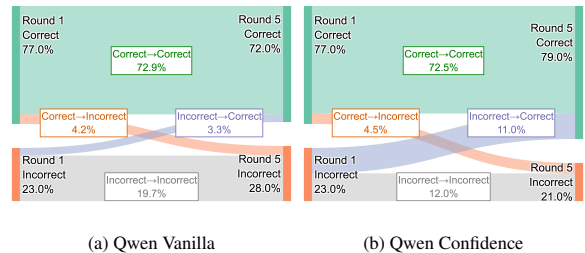


Figure 2: Sankey plots of correctness transitions from round 1 to round 5 (Qwen): The confidence-modulated debate (b) shows improvement on answer revision (blue), indicating better belief drift to correctness.

are provided in the Appendix B. We evaluate our methods against single model, majority vote, and majority vote after debate as baselines.

## 7 Results & Analysis

### Diversity alone improves MAD performance.

We observe that vanilla MAD performs similarly or worse than simple majority vote, echoing findings reported in prior work (Choi et al., 2025). However, further analysis reveals several meaningful patterns. As shown in Table 2, even without any intervention, initial answer diversity in vanilla MAD is significantly correlated with final answer accuracy. While the correlation is weak, considering its consistency and direction in addition to the diversity collapse of the instruct model, they together suggest that diversity exerts a systematic effect in vanilla MAD, motivating mechanisms that explicitly leverage this

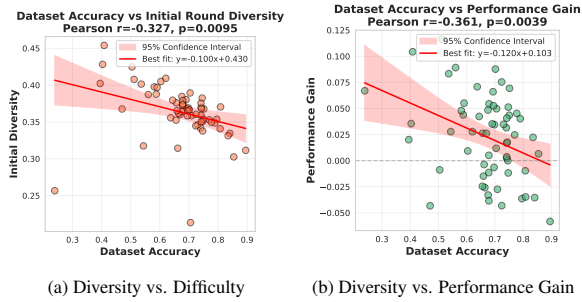


Figure 3: **(a)** Correlation between initial answer diversity and dataset difficulty (measured by model accuracy) over 62 datasets (including 57 MMLU sub-tasks). **(b)** Relationship between dataset difficulty and debate performance gain. MAD tends to yield larger gains on more difficult datasets.

signal. This also mirrors findings in human group deliberation, where diversity promotes collective accuracy (Karadzhov et al., 2024).

We additionally study how diversity interacts with task difficulty. We define the difficulty as a model’s accuracy on a given dataset. Figure 3a shows that MAD naturally produces a more diverse initial answer pool on harder datasets (*i.e.*, low accuracy). Consequently, as illustrated in Figure 3b, MAD tends to be more beneficial on more challenging tasks. These results provide an empirical explanation for when MAD helps and when it fails.

As shown in Table 1, across both models, our diversity-aware initialisation consistently improves accuracy on nearly all datasets and surpasses simple majority vote. To quantify this effect, we directly measure (i) the number of unique answers and (ii) Pass@5 in the initial round (*i.e.*, whether any agent initially produces a correct answer). As shown in Table 3, with diversity-aware initialisation, both metrics increase across the board, demonstrating that our intervention raises the prior probability of MAD success by ensuring more debates begin with at least one correct hypothesis. This directly supports the theoretical conclusion of Proposition 1.

**Confidence-modulated debate improves aggregation.** Our confidence-modulated debate includes an additional training phase in which agents learn to express calibrated confidence and to use others’ confidence when revising their answers. To isolate the effect of confidence, we compare against an ablation that removes the confidence reward and simply trains a model to select the best answer via aggregation (Learn2Agg), similar to Zhao et al. (2025). Even without diversity-aware initialisation,

<b>(a) Overall Pass@5</b>		
Model	Random	Diverse
Qwen-2.5-7B-Instruct	0.7921	0.9097
Llama-3.1-8B-Instruct	0.7423	0.9026
<b>(b) Number of Unique Answers@5</b>		
Model	Random	Diverse
Qwen-2.5-7B-Instruct	1.45	1.61
Llama-3.1-8B-Instruct	1.92	2.32

Table 3: Comparison of Pass@5 and answer diversity across models and initialisation regimes.

the confidence-modulated debate consistently outperforms simple aggregation across all evaluated datasets, including the two OOD benchmarks. This confirms that incorporating confidence signals provides a reliable mechanism for improving debate quality, consistent with the theoretical prediction of Theorem 1. We further analyse how the dynamics change with in the debate. We conducted a controlled experiment using the initial turn output from vanilla MAD to initialise the confidence modulated debate. As shown in Figure 2, confidence modulation increases correction of initially incorrect answers and stabilises correct ones. This indicates that confidence-modulated debate improved belief propagation with better drift to correctness, aligning with our theoretical analysis of Theorem 1.

## 8 Conclusion

We revisit multi-agent debate (MAD) through the lens of human group deliberation and recent theory showing that vanilla MAD induces a martingale over agents’ beliefs. Drawing on psycholinguistic and socio-cognitive insights, we identify two missing ingredients in standard LLM debates: diversity of initial viewpoints and explicit, calibrated confidence communication. We introduce a diversity-aware initialization that increases the probability that debate starts with a correct hypothesis, and a confidence-modulated debate protocol in which agents express and use calibrated confidence. Theoretically, we show that diversity improves prior probability of MAD success without affecting martingale dynamics, while confidence-weighted updates break the martingale symmetry and yield a submartingale over correctness. Empirically, our methods consistently outperform vanilla MAD and majority vote across six benchmarks. Overall, our results demonstrate that simple, principled changes inspired by human deliberation can substantially improve LLM multi-agent debate.

## Limitations

Our work has several limitations. First, our theoretical analysis builds on the Dirichlet–categorical debate model and homogeneity assumptions of [Choi et al. \(2025\)](#). While this abstraction enables clear insights into debate dynamics, it does not fully capture the complexity of real LLM behavior. In particular, the DCM model reduces belief updates to simple counting, ignoring that agents may receive different inputs or contextual positions, which can induce effects such as position bias ([Wang et al., 2024](#)). Moreover, the martingale analysis abstracts away settings in which LLMs may refine their own beliefs through self-reflection or internal reasoning ([Huang et al., 2023](#)).

Second, although our diversity-aware initialization is training-free and effective, it relies on a heuristic selection of diverse outputs from a fixed candidate pool. We do not claim this strategy to be optimal, and more principled or adaptive objectives for promoting diversity may further improve debate outcomes ([He et al., 2025](#); [Anschel et al., 2025](#)).

Third, our confidence-modulated debate assumes that expressed confidence is positively correlated with correctness after calibration. While our training procedure improves this alignment empirically, confidence miscalibration remains a known challenge for LLMs, and failures in confidence estimation could reduce the effectiveness of confidence-weighted aggregation.

Fourth, our experiments focus on homogeneous agents and fully connected debate graphs. Extending the framework to heterogeneous models, asymmetric communication topologies, or larger agent populations may introduce additional dynamics not captured in this work.

Finally, we evaluate MAD primarily on English language question-answering benchmarks. While these tasks are standard for studying reasoning and debate, future work is needed to assess whether our findings generalize to more open-ended, interactive, or real-world multi-agent settings.

## Ethics Statement

Our research adheres to strict ethical guidelines. We verified the licenses of all the software and datasets used in this study to ensure full compliance with their terms. No privacy concerns have been identified. We have conducted a thorough assessment of the project and do not anticipate any

further risks. We only used AI assistance for grammar checking when writing the paper.

## Acknowledgement

We thank Zhaobo Han, PhD candidate in pure mathematics at UCLA, for reviewing our proofs. Andreas Vlachos is supported by the DARPA program SciFy.

## References

- Oron Anschel, Alon Shoshan, Adam Botach, Shunit Haviv Hakimi, Asaf Gendler, Emanuel Ben Baruch, Nadav Bhonker, Igor Kviatkovsky, Manoj Aggarwal, and Gerard Medioni. 2025. Group-aware reinforcement learning for output diversity in large language models. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 32382–32403.
- Solomon E Asch. 1955. Opinions and social pressure. *Scientific american*, 193(5):31–35.
- Bahador Bahrami, Karsten Olsen, Peter E Latham, Andreas Roepstorff, Geraint Rees, and Chris D Frith. 2010. Optimally interacting minds. *Science*, 329(5995):1081–1085.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2024. [Chateval: Towards better llm-based evaluators through multi-agent debate](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Hyeong Kyu Choi, Xiaojin Zhu, and Yixuan Li. 2025. [Debate or vote: Which yields better decisions in multi-agent large language models?](#) *ArXiv preprint*, abs/2508.17536.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#). *ArXiv preprint*, abs/1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *ArXiv preprint*, abs/2110.14168.
- David Crecraft and Stephen Gergely. 2002. *Analog Electronics: circuits, systems and signal processing*. Elsevier.
- Guillaume Dezeache, Martin Dockendorff, Dardo N Ferreiro, Ophelia Deroy, and Bahador Bahrami. 2022. Democratic forecast: Small groups predict the future better than individuals and crowds. *Journal of Experimental Psychology: Applied*, 28(3):525.

- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2024. [Improving factuality and reasoning in language models through multiagent debate](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Andrew Estornell and Yang Liu. 2024. [Multi-llm debate: Framework, principals, and interventions](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Riccardo Fusaroli, Bahador Bahrami, Karsten Olsen, Andreas Roepstorff, Geraint Rees, Chris Frith, and Kristian Tylén. 2012. Coming to terms: Quantifying the benefits of linguistic coordination. *Psychological science*, 23(8):931–939.
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. [Large language model based multi-agents: A survey of progress and challenges](#). In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI 2024, Jeju, South Korea, August 3-9, 2024*, pages 8048–8057. ijcai.org.
- Joshua J Guyer, Pablo Briñol, Thomas I Vaughan-Johnston, Leandre R Fabrigar, Lorena Moreno, and Richard E Petty. 2021. Paralinguistic features communicated through voice can affect appraisals of confidence and evaluative judgments. *Journal of nonverbal behavior*, 45(4):479–504.
- Andre Wang He, Daniel Fried, and Sean Welleck. 2025. Rewarding the unlikely: Lifting grpo beyond distribution sharpening. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 25559–25571.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text degeneration](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Tiancheng Hu, Benjamin Minixhofer, and Nigel Collier. 2025. Navigating the alignment-calibration trade-off: A pareto-superior frontier via model merging. *arXiv preprint arXiv:2510.17426*.
- Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023. Large language models can self-improve. In *Proceedings of the 2023 conference on empirical methods in natural language processing*, pages 1051–1068.
- Irving L Janis. 1972. Victims of groupthink: A psychological study of foreign-policy decisions and fiascoes.
- Georgi Karadzhov, Andreas Vlachos, and Tom Stafford. 2024. [The effect of diversity on group decision-making](#). *ArXiv preprint*, abs/2402.01427.
- Norbert L Kerr and R Scott Tindale. 2004. Group performance and decision making. *Annu. Rev. Psychol.*, 55(1):623–655.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. 2024. [Understanding the effects of RLHF on LLM generalisation and diversity](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. Encouraging divergent thinking in large language models through multi-agent debate. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pages 17889–17904.
- Zijie Lin and Bryan Hooi. 2025. [Enhancing multi-agent debate system performance via confidence expression](#). *ArXiv preprint*, abs/2509.14034.
- Hugo Mercier. 2020. Not born yesterday: The science of who we trust and what we believe. In *Not born yesterday*. Princeton University Press.
- Hugo Mercier and Dan Sperber. 2011. Why do humans reason? arguments for an argumentative theory. *Behavioral and brain sciences*, 34(2):57–74.
- Hugo Mercier and Dan Sperber. 2017. *The enigma of reason*. Harvard University Press.
- David Moshman and Molly Geil. 1998. Collaborative reasoning: Evidence for collective rationality. *Thinking & Reasoning*, 4(3):231–248.
- Vishakh Padmakumar and He He. 2024. [Does writing with language models reduce content diversity?](#) In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.

- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *ArXiv preprint*, abs/2402.03300.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, 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. 2024. [Towards understanding sycophancy in language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Paul E Smaldino, Cody Moser, Alejandro Perez Velilla, and Mikkel Werling. 2024. Maintaining transient diversity is a general principle for improving collective problem solving. *Perspectives on Psychological Science*, 19(2):454–464.
- Paul Stangel, David Bani-Harouni, Chantal Pellegrini, Ege Özsoy, Kamilia Zaripova, Matthias Keicher, and Nassir Navab. 2025. [Rewarding doubt: A reinforcement learning approach to calibrated confidence expression of large language models](#). *ArXiv preprint*, abs/2503.02623.
- Fengfei Sun, Ningke Li, Kailong Wang, and Lorenz Goette. 2025. [Large language models are overconfident and amplify human bias](#). *ArXiv preprint*, abs/2505.02151.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Qwen Team. 2024. [Qwen2.5: A party of foundation models](#).
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and 1 others. 2024. Large language models are not fair evaluators. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9440–9450.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Haolun Wu, Zhenkun Li, and Lingyao Li. 2025. [Can llm agents really debate? a controlled study of multi-agent debate in logical reasoning](#). *ArXiv preprint*, abs/2511.07784.
- Andrea Wynn, Harsh Satija, and Gillian Hadfield. 2025. [Talk isn't always cheap: Understanding failure modes in multi-agent debate](#). *ArXiv preprint*, abs/2509.05396.
- Ruihan Yang, Caiqi Zhang, Zhisong Zhang, Xinting Huang, Sen Yang, Nigel Collier, Dong Yu, and Deqing Yang. 2025a. [LoGU: Long-form generation with uncertainty expressions](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 18947–18968, Vienna, Austria. Association for Computational Linguistics.
- Ruihan Yang, Caiqi Zhang, Zhisong Zhang, Xinting Huang, Dong Yu, Nigel Collier, and Deqing Yang. 2025b. [UNCLE: Benchmarking uncertainty expressions in long-form generation](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 30328–30344, Suzhou, China. Association for Computational Linguistics.
- Yongjin Yang, Euiin Yi, Jongwoo Ko, Kimin Lee, Zhijing Jin, and Se-Young Yun. 2025c. [Revisiting multi-agent debate as test-time scaling: A systematic study of conditional effectiveness](#). *ArXiv preprint*, abs/2505.22960.
- Luke Yoffe, Alfonso Amayuelas, and William Yang Wang. 2024. [Debunc: Improving large language model agent communication with uncertainty metrics](#). *ArXiv preprint*, abs/2407.06426.
- Paul Zarnoth and Janet A Sniezek. 1997. The social influence of confidence in group decision making. *Journal of Experimental Social Psychology*, 33(4):345–366.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. 2024. [LUQ: Long-text uncertainty quantification for LLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5244–5262, Miami, Florida, USA. Association for Computational Linguistics.

- Caiqi Zhang, Ruihan Yang, Zhisong Zhang, Xinting Huang, Sen Yang, Dong Yu, and Nigel Collier. 2025a. [Atomic calibration of LLMs in long-form generations](#). In *Proceedings of the 14th International Joint Conference on Natural Language Processing and the 4th Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics*, pages 148–169, Mumbai, India. The Asian Federation of Natural Language Processing and The Association for Computational Linguistics.
- Caiqi Zhang, Ruihan Yang, Xiaochen Zhu, Chengzu Li, Tiancheng Hu, Yijiang River Dong, Deqing Yang, and Nigel Collier. 2026. [Confidence estimation for llms in multi-turn interactions](#). *Preprint*, arXiv:2601.02179.
- Caiqi Zhang, Xiaochen Zhu, Chengzu Li, Nigel Collier, and Andreas Vlachos. 2025b. [Reinforcement learning for better verbalized confidence in long-form generation](#). *ArXiv preprint*, abs/2505.23912.
- Wenting Zhao, Pranjal Aggarwal, Swarnadeep Saha, Asli Celikyilmaz, Jason Weston, and Ilia Kulikov. 2025. [The majority is not always right: RL training for solution aggregation](#). *ArXiv preprint*, abs/2509.06870.
- Xiaochen Zhu, Caiqi Zhang, Tom Stafford, Nigel Collier, and Andreas Vlachos. 2025. Conformity in large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3854–3872.

## A Prompts

Here we provide the prompt we used for MAD. We use the prompt in Table 4 to elicit the initial pool of answers for the debate. For the subsequent rounds, we use the prompt in Table 5 to encourage debate. For a single agent, the other agents’ responses in the previous round will be concatenated together and presented together with its own previous response. For settings not incorporating confidence, the confidence related descriptions are simply removed. All the interactions are done using chat template.

Initial Round Prompt with Confidence
<p>Answer the question. Think step by step first, at the end of your reasoning, provide your final answer and confidence level, ranging from 0 to 10, where 0 means no confidence at all and 10 means complete confidence.</p> <p>Output format:            &lt;reasoning&gt; YOUR DETAILED REASONING HERE &lt;/reasoning&gt;</p> <p>&lt;answer&gt; YOUR FINAL ANSWER &lt;/answer&gt; {extra_info}</p> <p>&lt;confidence&gt; INTEGER &lt;/confidence&gt;</p> <p>{question}</p>

Table 4: Initial Round Prompt with Confidence.

## B Sampling and Training Schema

All experiments are conducted on NVIDIA A100 GPUs (80GB), with a total consumption of approximately 1,500 GPU hours. We use vLLM (v0.11.0) for all inference. Both training and inference are performed in float16 precision. For all fine-tuning stages, we employ LoRA with rank 64.

### B.1 Confidence Expression Training

We first initialize the model with a small number of supervised fine-tuning (SFT) steps to improve format adherence for confidence expression. To construct gold labels, we subsample 5k data points from the combined training datasets and use self-consistency (Wang et al., 2023) to derive numerical confidence targets. We then proceed with reinforcement learning (RL) training on a subsampled set of 10k difficult examples from the training data. Table 6 demonstrate that trained models exhibits well-calibrated verbalised confidence. Model also main-

MAD Prompt with Confidence
<p>You are revising your answer after reviewing other agents’ reasoning.            Question: {question}            Other agents’ responses: {other_agents}            Your previous reasoning and answer: {reasoning}            Instructions:</p> <ul style="list-style-type: none"> <li>- Reflect on how others reasoned.</li> <li>- You may revise your answer if someone’s reasoning provides stronger evidence.</li> <li>- However, if you believe all of them missed something important, propose a better or alternative answer — clearly explain why.</li> <li>- Be concise and clear.</li> <li>- Update your confidence level to reflect how certain you are now. Confidence level is ranging from 0 to 10, where 0 means no confidence at all and 10 means complete confidence.</li> </ul> <p>Think step by step first, at the end of your reasoning, provide your final answer and confidence level in the following format:            &lt;reasoning&gt; YOUR DETAILED REASONING HERE &lt;/reasoning&gt;            &lt;answer&gt; YOUR FINAL ANSWER &lt;/answer&gt; {extra_info}            &lt;confidence&gt; INTEGER &lt;/confidence&gt;</p>

Table 5: MAD Prompt with Confidence.

tains a similar accuracy demonstrating the training does not introduce new knowledge that causing data leakage.

Llama-3.1-8B-Instruct				
Method	Brier	ECE	AUROC	Acc
Vanilla	0.217	0.257	0.644	0.491
SC	0.051	0.157	0.852	N/A
VC	0.069	0.167	0.763	0.502
Qwen-2.5-7B-Instruct				
Method	Brier	ECE	AUROC	Acc
Vanilla	0.212	0.217	0.576	0.567
SC	0.078	0.156	0.855	N/A
VC	0.046	0.153	0.754	0.578

Table 6: Calibration and performance metrics for Llama-3.1-8B-Instruct and Qwen-2.5-7B-Instruct on hard sub-sampled datasets. Lower is better for Brier score and ECE; higher is better for AUROC and accuracy. All values rounded to three significant digits.

### B.2 Confidence Perception and Usage Training

Since the input–output format between debate rounds is identical, the debate process can be viewed as Markovian, with each turn conditioned only on agents’ responses from the previous round.

GRPO Confidence Calibration Setup	
LoRA	( $r=64$ , $\alpha=32$ ; q/k/v/o proj)
Seq. Length	2048
Epochs	1
Batch Size	4 ( $\times 2$ grad acc.)
Learning Rate	$5 \times 10^{-6}$
Generations / Prompt	8
Sampling	$T=1.0$ , top- $p=0.9$
s KL Coefficient	$\beta=0.01$

Table 7: Training configuration for confidence calibration.

GRPO Multi-Agent Debate Setup	
LoRA	( $r=64$ , $\alpha=32$ ; q/k/v/o proj)
Epochs	1
Batch Size	4 ( $\times 4$ grad acc.)
Learning Rate	$5 \times 10^{-5}$
Generations / Prompt	8
Sampling	$T=1.0$ , top- $p=1.0$
KL Coefficient	$\beta=0.01$
Max Completion Length	1024
Confidence Used	Yes
Reward Scales	Correct: 10, Conf: 3, Engage: 5
Wrong-Format Penalty	-30

Table 8: GRPO configuration for multi-agent debate training.

We therefore apply GRPO to single-turn inputs and outputs using the combined correctness and confidence rewards defined in Equation (11).

During training, we observed reward hacking behavior, where agents stopped engaging with others’ reasoning or confidence signals and instead output only a final answer. To mitigate this, we introduce a simple engagement constraint that encourages constructive interaction by rewarding the presence of discourse cues (e.g., “agent”, “I agree”, “but I think”, “convinced”, “skeptical”). This heuristic helps ensure that agents actively reference and respond to others during debate.

## C What We Tried but Did Not Work

**Increasing temperature does not increase answer diversity.** In early experiments, we attempted to promote diversity by increasing the sampling temperature. However, as shown in Table 9, raising the temperature from 1.0 to 1.2 did not yield more diverse answers or improved MAD performance. We attribute this to the fact that higher temperature primarily induces diversity in *reasoning paths* rather than in final answers. For instance,

Qwen often produced the same answer expressed through different languages (e.g., Chinese or Korean) or stylistic variations. A higher temperature also harms instruction following and lead to worse output quality. This observation suggests that naive temperature scaling is insufficient for inducing answer-level diversity and motivates the need for more targeted diversity-aware initialization strategies.

Qwen-2.5-7B-Instruct		
Temperature	Unique Answers@5	Accuracy
$t = 0.8$	1.18	0.722
$t = 1.0$	1.45	0.753
$t = 1.2$	1.42	0.733

Table 9: Effect of temperature on answer diversity and average accuracy for Qwen-2.5-7B-Instruct. Values rounded to three significant digits.

**Training on harder data improves confidence calibration.** For confidence expression training, we initially subsampled 12k data points uniformly from the training datasets. In practice, the model struggled to converge under this setting. We attribute this to the model’s already high accuracy on these datasets: during GRPO training, many trajectories resulted in correct answers with maximal confidence (10), leading to a collapse in reward diversity and poor sampling efficiency. To address this issue, we manually curated a harder subset of data on which model accuracy was closer to 50%. Training on this subset substantially improved convergence and led to more stable confidence calibration.

## D Proofs

In this appendix we show that both our interventions of diversity injection and confidence-modulated debate could benefit the result of the vanilla MAD which is essentially a martingale process as established in (Choi et al., 2025). For diversity injection, we are not changing the dynamics of the martingale process, but increases the prior probability of success. For confidence-modulated debate, the incorporation of a confidence score that is positively correlated with correctness breaks martingale process into submartingale.

### D.1 Baseline DCM Debate Model

We briefly recall the Bayesian model of multi-agent debate used in Choi et al. (2025). Let  $A = \{1, \dots, K\}$  denote the finite answer set and let answer 1 be the unique correct option. There are  $N$  agents indexed by  $i \in \{1, \dots, N\}$ . At debate round  $t$ , agent  $i$  is parameterized by a Dirichlet belief vector

$$\boldsymbol{\alpha}_{i,t} = (\alpha_{i,t}^{(1)}, \dots, \alpha_{i,t}^{(K)}) \in \mathbb{R}_{>0}^K, \quad (15)$$

and generates its answer  $y_{i,t} \in A$  as

$$\begin{aligned} \boldsymbol{\theta}_{i,t} &\sim \text{Dir}(\boldsymbol{\alpha}_{i,t}), \\ y_{i,t} \mid \boldsymbol{\theta}_{i,t} &\sim \text{Cat}(\boldsymbol{\theta}_{i,t}). \end{aligned}$$

The induced marginal probability of choosing option  $k$  is

$$\Pr(y_{i,t} = k \mid \boldsymbol{\alpha}_{i,t}) = \frac{\alpha_{i,t}^{(k)}}{\sum_{j=1}^K \alpha_{i,t}^{(j)}}. \quad (16)$$

We write

$$p_{i,t} := \Pr(y_{i,t} = 1 \mid \boldsymbol{\alpha}_{i,t}) = \frac{\alpha_{i,t}^{(1)}}{\sum_{j=1}^K \alpha_{i,t}^{(j)}} \quad (17)$$

for agent  $i$ 's belief (probability) on the correct option at round  $t$ .

Let  $G$  be the (undirected) debate graph and  $\mathcal{N}(i)$  denote the neighbors of  $i$  (possibly including  $i$  itself). In the *unweighted* model, agent  $i$  aggregates its neighbors' answers  $\{y_{j,t-1} : j \in \mathcal{N}(i)\}$  into a count vector  $\mathbf{c}_{i,t} = (c_{i,t}^{(1)}, \dots, c_{i,t}^{(K)})$ , where

$$c_{i,t}^{(k)} = \sum_{j \in \mathcal{N}(i)} \mathbf{1}\{y_{j,t-1} = k\}. \quad (18)$$

The Dirichlet parameters are then updated by Bayesian conjugacy:

$$\boldsymbol{\alpha}_{i,t} = \boldsymbol{\alpha}_{i,t-1} + \mathbf{c}_{i,t}. \quad (19)$$

Let  $\mathcal{F}_t$  denote the filtration generated by all belief parameters up to round  $t$ ,  $\mathcal{F}_t := \sigma(\{\boldsymbol{\alpha}_{i,s} : i \leq N, s \leq t\})$ . In the homogeneous, fully-connected setting where all agents share the same prior and observe the same multiset of answers, Choi et al. (2025) show that

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] = p_{i,t-1}, \quad (20)$$

so  $\{p_{i,t}\}_{t>0}$  is a martingale. Intuitively, each neighbour contributes one *unweighted* ‘‘ball’’ to a Pólya-urn-like update, and symmetry ensures that, in expectation, debate neither helps nor hurts.

## D.2 Effect of Diversity Injection

Recall that the martingale result of Choi et al. (2025) states that, conditional on the initial debate state  $\mathcal{F}_1$  (the first-round pool and its arguments), the expected belief does not drift:

$$\mathbb{E}[p_t \mid \mathcal{F}_1] = p_1 \quad \forall t, \quad (21)$$

so the probability that debate outputs the correct answer is determined entirely by the quality of  $\mathcal{F}_1$ .

We now formalize how diversity-aware initialization improves this prior success probability without altering the martingale dynamics.

**Theorem 2** (Effect of diversity-aware initialization). *Let  $A_T$  be the event that debate outputs the correct answer at the final round  $T$ , and let  $S \in \mathbb{N}$  denote the number of distinct informative hypotheses (e.g., distinct correct or partially-correct reasoning chains) present in the initial pool  $\mathcal{F}_1$ . Assume:*

(i) (Monotonicity in coverage) *The success probability is weakly increasing in  $S$ , i.e.*

$$\Pr(A_T \mid S = s) \quad \text{is nondecreasing in } s. \quad (22)$$

(ii) (Diversity increases coverage) *Let  $S_{\text{rand}}$  and  $S_{\text{div}}$  denote the values of  $S$  under random sampling and diversity-aware initialization respectively. Then  $S_{\text{div}}$  first-order stochastically dominates  $S_{\text{rand}}$ , i.e.*

$$\Pr(S_{\text{div}} \geq s) \geq \Pr(S_{\text{rand}} \geq s) \quad \forall s, \quad (23)$$

*with strict inequality for some  $s$ .*

*Then the overall probability of eventual correctness is higher under diversity-aware initialization:*

$$\Pr(A_T \mid \text{diverse}) \geq \Pr(A_T \mid \text{random}), \quad (24)$$

*with strict inequality whenever the dominance in (23) is strict.*

*Proof.* Define the function

$$f(s) := \Pr(A_T \mid S = s), \quad s \in \mathbb{N}. \quad (25)$$

Assumption (22) states that  $f$  is nondecreasing.

Under a given initialization scheme (random or diverse), we can write the overall success probability using the law of total probability as

$$\Pr(A_T \mid \text{scheme}) = \mathbb{E}[f(S_{\text{scheme}})], \quad (26)$$

where  $S_{\text{scheme}}$  is the corresponding random variable ( $S_{\text{rand}}$  or  $S_{\text{div}}$ ).

By first-order stochastic dominance (23), we know that for any nondecreasing function  $f$ ,

$$\mathbb{E}[f(S_{\text{div}})] \geq \mathbb{E}[f(S_{\text{rand}})], \quad (27)$$

with strict inequality if dominance is strict on a set where  $f$  is strictly increasing.

Combining these facts gives

$$\Pr(A_T \mid \text{diverse}) = \mathbb{E}[f(S_{\text{div}})] \geq \mathbb{E}[f(S_{\text{rand}})] = \Pr(A_T \mid \text{random}), \quad (28)$$

with strict inequality under strict dominance. This is exactly the claimed result.  $\square$

**Simple illustration (signal-coverage model).** As a concrete example, suppose each distinct informative hypothesis in the initial pool independently has correctness probability  $p > 0$ . Then, conditional on observing  $S = s$  such hypotheses,

$$\Pr(A_T \mid S = s) = 1 - (1 - p)^s, \quad (29)$$

which is strictly increasing in  $s$ , verifying the monotonicity assumption (22). If diversity-aware initialization increases the effective spread of hypotheses so that  $S_{\text{div}}$  stochastically dominates  $S_{\text{rand}}$ , then

$$\Pr(A_T \mid \text{diverse}) = \mathbb{E}[1 - (1 - p)^{S_{\text{div}}}] \geq \mathbb{E}[1 - (1 - p)^{S_{\text{rand}}}] = \Pr(A_T \mid \text{random}), \quad (30)$$

in agreement with Theorem 2.

**Relation to the martingale result.** Note that all improvement arises from a better initial distribution over  $\mathcal{F}_1$ . The martingale property of Choi et al. (2025) states that, once  $\mathcal{F}_1$  is fixed, debate cannot raise the conditional expectation of correctness in subsequent rounds. Diversity-aware initialization operates *before* this point: it shifts the distribution of  $\mathcal{F}_1$  toward states that already contain useful hypotheses, thereby increasing the overall probability of correctness without altering the debate dynamics themselves.

### D.3 Confidence-Weighted Debate

We now extend this model by allowing agents to emit not only an answer  $y_{i,t}$  but also a scalar confidence  $w_{i,t} \in (0, 1]$ . At each round  $t$ , agent  $i$  therefore outputs a pair  $(y_{i,t}, w_{i,t})$ .

**Confidence-weighted counts.** Given neighbors' outputs  $\{(y_{j,t-1}, w_{j,t-1}) : j \in \mathcal{N}(i)\}$ , agent  $i$  constructs a *confidence-weighted* count vector  $\tilde{\mathbf{c}}_{i,t} \in \mathbb{R}_{\geq 0}^K$  with components

$$\tilde{c}_{i,t}^{(k)} := \sum_{j \in \mathcal{N}(i)} w_{j,t-1} \mathbf{1}\{y_{j,t-1} = k\}, \quad k = 1, \dots, K. \quad (31)$$

We then update the Dirichlet parameters by adding the weighted counts:

$$\boldsymbol{\alpha}_{i,t} = \boldsymbol{\alpha}_{i,t-1} + \tilde{\mathbf{c}}_{i,t}. \quad (32)$$

This preserves conjugacy:

$$\boldsymbol{\theta}_{i,t} \mid \mathcal{F}_t \sim \text{Dir}(\boldsymbol{\alpha}_{i,t}), \quad (33)$$

but the increments are now real-valued rather than integer counts.

For notational convenience, write the total Dirichlet mass as

$$S_{i,t} := \sum_{k=1}^K \alpha_{i,t}^{(k)}, \quad (34)$$

and define, for the update from  $t - 1$  to  $t$ ,

$$\Delta_{i,t} := \sum_{k=1}^K \tilde{c}_{i,t}^{(k)} = \sum_{j \in \mathcal{N}(i)} w_{j,t-1}, \quad (35)$$

and the confidence-weighted ‘‘mass on the correct answer’’

$$\Delta_{i,t}^{(1)} := \tilde{c}_{i,t}^{(1)} = \sum_{j \in \mathcal{N}(i)} w_{j,t-1} \mathbf{1}\{y_{j,t-1} = 1\}. \quad (36)$$

Under the confidence-weighted update (32), the belief of agent  $i$  in the correct answer at round  $t$  can be written as

$$p_{i,t} = \frac{\alpha_{i,t-1}^{(1)} + \Delta_{i,t}^{(1)}}{S_{i,t-1} + \Delta_{i,t}}. \quad (37)$$

### D.4 From Martingale to Submartingale

We now study the conditional expectation of  $p_{i,t}$  given  $\mathcal{F}_{t-1}$ . The key observation is that, under (37), the next-round belief can be written as a convex combination of the current belief  $p_{i,t-1}$  and a confidence term.

**Lemma 1** (Convex combination form). Assume  $\Delta_{i,t} > 0$ . Define the realised fraction of confidence mass that lands on the correct answer as

$$\hat{q}_{i,t-1} := \frac{\Delta_{i,t}^{(1)}}{\Delta_{i,t}}. \quad (38)$$

Then, for each agent  $i$  and round  $t$ ,

$$p_{i,t} = \bar{\lambda}_{i,t} p_{i,t-1} + (1 - \bar{\lambda}_{i,t}) \hat{q}_{i,t-1}, \quad (39)$$

where

$$\bar{\lambda}_{i,t} := \frac{S_{i,t-1}}{S_{i,t-1} + \Delta_{i,t}} \in (0, 1). \quad (40)$$

In particular,  $p_{i,t}$  is (pathwise) a convex combination of  $p_{i,t-1}$  and  $\hat{q}_{i,t-1}$ .

*Proof.* Starting from (37) and using  $\alpha_{i,t-1}^{(1)} = p_{i,t-1} S_{i,t-1}$  and  $\Delta_{i,t}^{(1)} = \hat{q}_{i,t-1} \Delta_{i,t}$ , we obtain

$$p_{i,t} = \frac{p_{i,t-1} S_{i,t-1} + \hat{q}_{i,t-1} \Delta_{i,t}}{S_{i,t-1} + \Delta_{i,t}}. \quad (41)$$

Define  $x := S_{i,t-1} / \Delta_{i,t}$ . Then

$$p_{i,t} = \frac{x p_{i,t-1} + \hat{q}_{i,t-1}}{x + 1} = \frac{x}{x + 1} p_{i,t-1} + \frac{1}{x + 1} \hat{q}_{i,t-1}, \quad (42)$$

which is exactly (39) with  $\bar{\lambda}_{i,t} = x / (x + 1) \in (0, 1)$ .  $\square$

Lemma 1 shows that, on each realisation, the new belief is a convex combination of the previous belief and the realised confidence-weighted fraction of correct answers among the neighbors. To understand the effect on the *expected* belief, we now take conditional expectations and formalise a positive-correlation assumption.

**Assumption 1** (Confidence is positively correlated with correctness). For each agent  $j$  and round  $t - 1$ , define

$$p_{j,t-1} := \Pr(y_{j,t-1} = 1 \mid \mathcal{F}_{t-1}), \quad (43)$$

and let

$$\rho_{j,t-1} := \frac{\mathbb{E}[w_{j,t-1} \mathbf{1}\{y_{j,t-1} = 1\} \mid \mathcal{F}_{t-1}]}{\mathbb{E}[w_{j,t-1} \mid \mathcal{F}_{t-1}]}. \quad (44)$$

We say confidence is positively correlated with correctness if

$$\rho_{j,t-1} \geq p_{j,t-1} \quad \text{for all } j, t, \quad (45)$$

with strict inequality on a set of positive probability (for at least one agent and round). In words, among the mass of answers that receive high confidence, the fraction that are correct is at least as large as the agent's marginal accuracy, and strictly larger sometimes.

We also define the conditional expectation of the confidence-weighted correctness fraction:

**Definition 1** (Confidence-weighted expected correctness). For each agent  $i$  and round  $t$ , define

$$q_{i,t-1} := \frac{\mathbb{E}[\Delta_{i,t}^{(1)} \mid \mathcal{F}_{t-1}]}{\mathbb{E}[\Delta_{i,t} \mid \mathcal{F}_{t-1}]}. \quad (46)$$

Under mild regularity conditions (e.g., finite second moments),  $q_{i,t-1}$  captures the expected fraction of total confidence mass that lands on the correct answer.

We can now relate  $q_{i,t-1}$  to the neighbour-level  $\rho_{j,t-1}$ 's:

**Lemma 2** (Decomposition of confidence-weighted correctness). *For each agent  $i$  and round  $t$ ,*

$$q_{i,t-1} = \sum_{j \in \mathcal{N}(i)} \omega_{j,t-1} \rho_{j,t-1}, \quad (47)$$

where

$$\omega_{j,t-1} := \frac{\mathbb{E}[w_{j,t-1} \mid \mathcal{F}_{t-1}]}{\sum_{\ell \in \mathcal{N}(i)} \mathbb{E}[w_{\ell,t-1} \mid \mathcal{F}_{t-1}]} \in (0, 1), \quad \sum_{j \in \mathcal{N}(i)} \omega_{j,t-1} = 1. \quad (48)$$

*Proof.* By definition,

$$\begin{aligned} \mathbb{E}[\Delta_{i,t}^{(1)} \mid \mathcal{F}_{t-1}] &= \mathbb{E} \left[ \sum_{j \in \mathcal{N}(i)} w_{j,t-1} \mathbf{1}\{y_{j,t-1} = 1\} \mid \mathcal{F}_{t-1} \right] \\ &= \sum_{j \in \mathcal{N}(i)} \mathbb{E}[w_{j,t-1} \mathbf{1}\{y_{j,t-1} = 1\} \mid \mathcal{F}_{t-1}]. \end{aligned}$$

Similarly,

$$\mathbb{E}[\Delta_{i,t} \mid \mathcal{F}_{t-1}] = \sum_{j \in \mathcal{N}(i)} \mathbb{E}[w_{j,t-1} \mid \mathcal{F}_{t-1}]. \quad (49)$$

Plugging into (46) yields

$$q_{i,t-1} = \frac{\sum_{j \in \mathcal{N}(i)} \mathbb{E}[w_{j,t-1} \mathbf{1}\{y_{j,t-1} = 1\} \mid \mathcal{F}_{t-1}]}{\sum_{j \in \mathcal{N}(i)} \mathbb{E}[w_{j,t-1} \mid \mathcal{F}_{t-1}]} = \sum_{j \in \mathcal{N}(i)} \omega_{j,t-1} \rho_{j,t-1}, \quad (50)$$

with  $\omega_{j,t-1}$  and  $\rho_{j,t-1}$  as defined above.  $\square$

Combining Assumption 1 with Lemma 2 immediately yields:

**Corollary 1.** *Under Assumption 1, if the debate graph is fully connected and all agents are homogeneous so that  $p_{j,t-1} = p_{i,t-1}$  for all  $j \in \mathcal{N}(i)$ , then*

$$q_{i,t-1} \geq p_{i,t-1}, \quad (51)$$

*with strict inequality on a set of positive probability.*

*Proof.* By Lemma 2,

$$q_{i,t-1} = \sum_{j \in \mathcal{N}(i)} \omega_{j,t-1} \rho_{j,t-1}, \quad (52)$$

with  $\omega_{j,t-1} \geq 0$  and  $\sum_j \omega_{j,t-1} = 1$ . By Assumption 1,  $\rho_{j,t-1} \geq p_{j,t-1}$  for all  $j$ , with strict inequality for some  $j$  with positive probability. Under homogeneity  $p_{j,t-1} = p_{i,t-1}$  for all  $j$ , so

$$q_{i,t-1} \geq \sum_{j \in \mathcal{N}(i)} \omega_{j,t-1} p_{i,t-1} = p_{i,t-1}, \quad (53)$$

with strict inequality whenever at least one  $\rho_{j,t-1} > p_{i,t-1}$  with  $\omega_{j,t-1} > 0$ .  $\square$

We can now prove that confidence-weighted debate is a submartingale.

**Theorem 3** (Confidence-weighted debate is a submartingale). *Assume:*

1. the debate graph is fully connected and all agents are homogeneous, so  $p_{j,t-1} = p_{i,t-1}$  for all  $j$ ;
2. the Dirichlet update uses confidence-weighted counts (31)–(32); and
3. Assumption 1 holds.

Then for each agent  $i$  and each round  $t \geq 1$ ,

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] \geq p_{i,t-1}, \quad (54)$$

with strict inequality on a set of positive probability. That is,  $\{p_{i,t}\}_{t \geq 0}$  is a strict submartingale and therefore cannot be a martingale.

*Proof.* From Lemma 1 we have, pathwise,

$$p_{i,t} = \bar{\lambda}_{i,t} p_{i,t-1} + (1 - \bar{\lambda}_{i,t}) \hat{q}_{i,t-1}, \quad (55)$$

with  $\bar{\lambda}_{i,t} \in (0, 1)$ . Taking conditional expectations with respect to  $\mathcal{F}_{t-1}$  and using the fact that  $p_{i,t-1}$  is  $\mathcal{F}_{t-1}$ -measurable, we obtain

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] = \mathbb{E}[\bar{\lambda}_{i,t} \mid \mathcal{F}_{t-1}] p_{i,t-1} + \mathbb{E}[(1 - \bar{\lambda}_{i,t}) \hat{q}_{i,t-1} \mid \mathcal{F}_{t-1}]. \quad (56)$$

By definition of  $q_{i,t-1}$  and Jensen-type arguments, we can interpret  $q_{i,t-1}$  as the conditional expected fraction of confidence on the correct answer, so that

$$\mathbb{E}[\hat{q}_{i,t-1} \mid \mathcal{F}_{t-1}] = q_{i,t-1}. \quad (57)$$

Using this and linearity of expectation we get

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] = \lambda_{i,t} p_{i,t-1} + (1 - \lambda_{i,t}) q_{i,t-1}, \quad (58)$$

where

$$\lambda_{i,t} := \mathbb{E}[\bar{\lambda}_{i,t} \mid \mathcal{F}_{t-1}] \in (0, 1). \quad (59)$$

Finally, by Corollary 1,  $q_{i,t-1} \geq p_{i,t-1}$  with strict inequality on a set of positive probability. Hence

$$\mathbb{E}[p_{i,t} \mid \mathcal{F}_{t-1}] = \lambda_{i,t} p_{i,t-1} + (1 - \lambda_{i,t}) q_{i,t-1} \geq \lambda_{i,t} p_{i,t-1} + (1 - \lambda_{i,t}) p_{i,t-1} = p_{i,t-1}, \quad (60)$$

with strict inequality whenever  $q_{i,t-1} > p_{i,t-1}$  and  $1 - \lambda_{i,t} > 0$ . Thus  $\{p_{i,t}\}_{t \geq 0}$  is a strict submartingale and cannot be a martingale.  $\square$

**Discussion.** Theorem 3 makes precise the intuition that confidence-weighted debate breaks the martingale behavior of vanilla multi-agent debate. In the unweighted case, symmetry implies that the expected belief in the correct answer is preserved across rounds. Once we weight peer messages by confidence, and confidence is even mildly informative about correctness, the expected belief drifts upward. In other words, confidence-weighted debate can, in expectation, do strictly better than simple unweighted majority vote.

## E Example Output

Table 10: Example of Debate Trajectory

---

Model Input
< begin_of_text >< start_header_id >system< end_header_id >
Cutting Knowledge Date: December 2023 Today Date: 26 Jul 2024
< eot_id >< start_header_id >user< end_header_id >
You are revising your answer after reviewing other agents' reasoning.
<b>Question:</b> If the Moon is setting at noon the phase of the Moon must be (A) third quarter. (B) waning crescent. (C) waxing crescent. (D) full.
<hr/> <b>Other agents' responses:</b>
<b>Agent 1</b> provided the following reasoning and final answer: <reasoning> To determine the phase of the Moon when it is setting at noon, we need to understand the lunar cycle and how it relates to the position of the Earth, Sun, and Moon. 1. The lunar cycle is about 29.5 days, which is the time it takes for the Moon to orbit the Earth. (Steps 2 - 7 omitted here) 8. Since the Moon sets at noon, this means that the portion of the illuminated Moon that is visible from our location on Earth is decreasing. 9. Therefore, it must be in the waning phase of the lunar cycle. The phase of the Moon when it sets at noon is waning. So, the answer is (B) waning crescent.</reasoning> <answer>(B)</answer> <confidence>2</confidence>
<b>Agent 2</b> reasoning: <reasoning> The Moon is setting at noon, and it's in the waning phase. Since it's waning and visible as a crescent in the sky at noon, we are witnessing a waning crescent Moon.</reasoning> <answer>(B)</answer> <confidence>3</confidence>
<b>Agent 3</b> reasoning: <reasoning> Based on the relationship between the Earth, Sun, and Moon's positions, the Moon's phase when setting at noon can be determined. Since the Moon is on the opposite side of the Earth from the Sun at this time, it can either be waning crescent or third quarter. However, the only phase that meets the condition of the Moon being on the left side of the Sun and setting at noon is the third quarter.</reasoning> <answer>(A)</answer> <confidence>2</confidence>
<b>Agent 4</b> reasoning: <reasoning> When the Moon is setting at noon, it is on the opposite side of the Earth from the Sun, and it must be fully illuminated. Therefore, it cannot be in any waning phase, waxing crescent phase, or third quarter. This leaves us with the full moon phase as the only possibility.</reasoning> <answer>(D)</answer> <confidence>3</confidence>
<hr/> <b>Your previous reasoning and answer:</b> The Moon's position relative to the Sun at noon determines its phase. Since the Moon is in the western part of the sky at noon and the Sun is in the eastern part, the Moon must be on the opposite side of the Earth from the Sun. This can only be the case for a full moon. <answer>(D)</answer> <confidence>3</confidence>
<b>Instructions:</b> - Reflect on how others reasoned. - You may revise your answer if someone's reasoning provides stronger evidence. - However, if you believe all of them missed something important, propose a better or alternative answer — clearly explain why. - Be concise and clear. - Update your confidence level to reflect how certain you are now. Confidence level is ranging from 0 to 10, where 0 means no confidence at all and 10 means complete confidence
Think step by step first, at the end of your reasoning, provide your final answer and confidence level in the following format:

<reasoning>YOUR DETAILED REASONING HERE</reasoning>  
<answer>YOUR FINAL ANSWER</answer>, where YOUR FINAL ANSWER is one of (A), (B), (C) or (D).  
<confidence>INTEGER</confidence> <|eot\_id|>

---

### Model Output

<reasoning>Upon reviewing the reasoning provided by Agents 1, 2, 3, and 4, I notice that we all agree on the basic premise that the Moon's phase can be determined by its position relative to the Sun at noon. However, our conclusions differ slightly. Agent 1 and Agent 2 both consider the waning crescent Moon as a possibility, but Agent 3 and I initially agreed on the full moon phase. Agent 4, on the other hand, also concludes that the full moon phase is the only possibility.

Upon further reflection, I realize that Agent 3's reasoning is particularly strong, as they correctly identify that the Moon must be in the third quarter phase to meet the condition of being on the left side of the Sun and setting at noon. This is a more precise analysis than my initial reasoning, which focused solely on the full moon phase.

However, I also notice that Agent 4's reasoning is not entirely flawed. They correctly point out that the Moon must be fully illuminated when setting at noon, which does indeed suggest the full moon phase. The key difference lies in the interpretation of the Moon's position relative to the Sun. While Agent 4 assumes the Moon is on the right side of the Sun, which would indeed result in a full moon, Agent 3's analysis shows that the Moon can be on the left side of the Sun and still be in the third quarter phase.

Considering these points, I revise my initial answer to reflect the stronger evidence presented by Agent 3. While the full moon phase is still a possibility, the third quarter phase is a more accurate conclusion based on the Moon's position relative to the Sun at noon.</reasoning>

<answer>(A)</answer>

<confidence>8</confidence>

---