

# Conformity Dynamics in LLM Multi-Agent Systems: The Roles of Topology and Self-Social Weighting

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) are increasingly instantiated as interacting agents in multi-agent systems (MAS), where collective decisions emerge through social interaction rather than independent reasoning. A fundamental yet underexplored mechanism in this process is conformity, the tendency of agents to align their judgments with prevailing group opinions. This paper presents a systematic study of how network topology shapes conformity dynamics in LLM-based MAS through a misinformation detection task. We introduce a confidence-normalized pooling rule that controls the trade-off between self-reliance and social influence, enabling comparisons between two canonical decision paradigms: Centralized Aggregation and Distributed Consensus. Experimental results demonstrate that network topology critically governs both the efficiency and robustness of collective judgments. Centralized structures enable immediate decisions but are sensitive to hub competence and exhibit same-model alignment biases. In contrast, distributed structures promote more robust consensus, while increased network connectivity speeds up convergence but also heightens the risk of wrong-but-sure cascades, in which agents converge on incorrect decisions with high confidence. These findings characterize the conformity dynamics in LLM-based MAS, clarifying how network topology and self-social weighting jointly shape the efficiency, robustness, and failure modes of collective decision-making. The code and dataset are available at: [Topology-of-Multi-Agent-Systems](#).

## 1 Introduction

Large Language Models (LLMs) have been increasingly instantiated as interacting agents within Multi-Agent Systems (MAS) (Wei et al., 2025). Across a wide range of applications, including collaborative problem solving (Du et al., 2024),

complex reasoning (Zhang and Xiong, 2025) and misinformation detection (Li et al., 2025), the effectiveness of collective decision-making depends not only on the competence of individual agents, but also on the social dynamics induced by their interactions (Ghoshal et al., 2025; Cisneros-Velarde, 2025). A central mechanism underlying these dynamics is conformity, the tendency of agents to adjust their judgments toward majority opinions. Prior studies suggest that conformity can reduce idiosyncratic noise and facilitate coordination (Choi et al., 2025). However, excessive conformity may also cause information cascades, leading groups to converge on incorrect conclusions with high confidence (Bikhchandani et al., 2024; Pinheiro and Vasconcelos, 2025).

Existing research on MAS has primarily emphasized task efficacy, focusing on protocol design, role specialization, and coordination efficiency (Agashe et al., 2025; Grötschla et al., 2025). For example, multi-agent debate frameworks can shift individual judgments toward majority positions, improving reliability while reinforcing systematic biases (Han et al., 2025). Related work on persuasion further suggests that LLM agents tend to imitate dominant argumentative patterns during the interaction (Argyle et al., 2025). While conformity has been studied in computational opinion dynamics (Helfmann et al., 2023; Ding et al., 2025; Han and Tang, 2025), existing approaches still lack an explicit treatment of how agents generate judgments during decision processes, and in particular how interaction topology and neighbor effect jointly govern the propagation, aggregation, and amplification of agent confidence.

In this study, we bridge this gap by investigating how network topology modulates conformity through a binary misinformation detection task. We propose a confidence-normalized update rule governed by a global self-weighting parameter  $\alpha$ , which balances an agent’s self-reliance against their

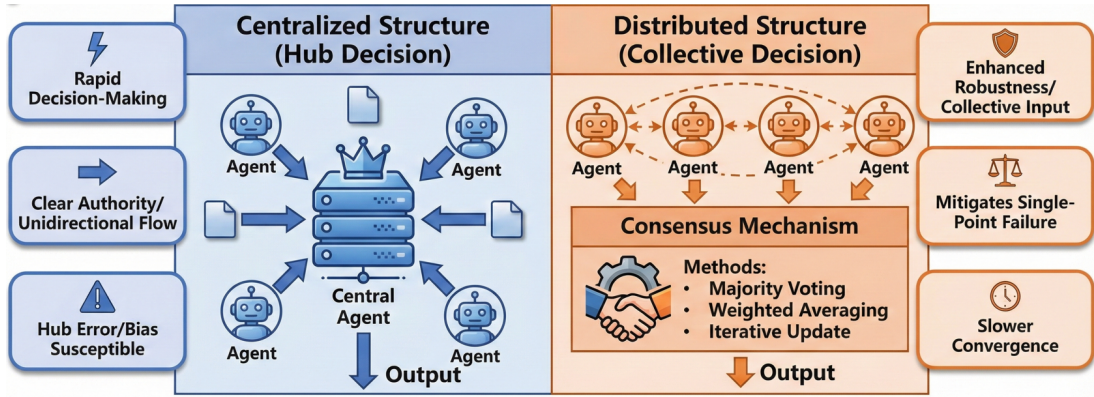


Figure 1: Illustration of centralized (hub-based) and distributed (collective) multi-agent structures, highlighting their decision-making mechanisms, relative advantages, and inherent limitations in collective inference.

neighborhood influence. Specifically, we compare two distinct decision modes: (i) **Centralized Aggregation** (e.g., star networks), where a hub synthesizes a collective decision in a single round; and (ii) **Distributed Consensus** (e.g., sparse rings to complete graph), where influence propagates iteratively. This allows us to disentangle one-step hub dominance from multi-round interaction. Our contributions are summarized as follows:

**(1) Topology as a determinant of conformity.**

We show that network structure systematically shapes conformity by trading off convergence speed and robustness. Centralized structures concentrate influence, yielding immediate decisions but making outcomes highly dependent on hub competence and model homogeneity; Distributed structures diffuse influence and support more robust consensus formation, while increased connectivity accelerates convergence and may facilitate high-confidence error cascades.

**(2) A transparent confidence-normalized pooling rule.** We propose a transparent and generalizable update mechanism that introduces a global parameter to explicitly regulate social influence. By integrating agent-level confidence into the pooling process, the rule yields bounded belief scores, stable binarization, and well-controlled conformity dynamics, providing a principled basis for different decision paradigms.

**(3) Empirical insights on robustness and risk.**

Experiments on current fact-checking datasets reveal that system reliability is heavily contingent on hub competence and majority composition. We further characterize the dual nature of conformity, showing that it enhances accuracy under reliable conditions but can induce confident errors when early majority signals are incorrect.

## 2 Related Works

### 2.1 MAS and Collective Decision-Making

Early research on MAS primarily focused on distributed optimization and consensus formation under idealized assumptions, where agents were modeled as homogeneous entities with limited reasoning capacity (Bao et al., 2022; Amirkhani and Barshooi, 2022). The recent integration of LLMs expands this paradigm by endowing agents with advanced reasoning, enabling MAS to address complex, unstructured decision-making problems such as multi-step problem solving (Chen et al., 2024) and automated fact-checking (Han et al., 2025). Despite these advances, most existing studies emphasize aggregate performance metrics, leaving the role of social dynamics such as conformity in shaping collective reliability largely unexamined.

### 2.2 Conformity and Opinion Dynamics

Conformity, defined as the tendency of individuals to align their judgments with a perceived majority, has been widely studied in social psychology (Capuano and Chekroun, 2024). In computational settings, opinion dynamics models such as DeGroot averaging (Dong et al., 2024) and bounded-confidence mechanisms (Li and Porter, 2023) formalize how local interactions give rise to collective consensus and have been used to explain phenomena including information cascades and polarization (Shirzadi et al., 2025). However, these models operate at an abstract level and do not capture the semantic reasoning or contextual understanding characteristic of modern LLM agents (Aouini and Loubani, 2025). Consequently, how classical conformity theories extend to LLM-based MAS, where judgments are produced through generative processes, remains insufficiently explored.

## 2.3 Network Topology

Network topology is a fundamental determinant of how influence propagates and consensus forms in multi-agent systems (Cheng et al., 2021; Amirkhani and Barshooi, 2022). **Centralized aggregation** structures, such as star or hierarchical graphs, concentrate influence and enable rapid decision-making, but render collective outcomes highly sensitive to the reliability of central agents. By contrast, **Distributed consensus** structures, including ring and complete graphs, diffuse influence across agents, enhancing robustness to local noise at the expense of slower convergence. Latest studies have further examined how connectivity modulates communication cost, convergence speed, and robustness (Da et al., 2025; Wang et al., 2025; Yang et al., 2025). Despite these efforts, existing work has yet to clarify the interplay between interaction topology, agent confidence, and conformity in shaping collective outcomes in LLM-driven MAS.

## 3 Methodology

### 3.1 Agent Decision Model

Building on classical opinion dynamics and weighted consensus frameworks (Anderson and Ye, 2019; Li and Porter, 2023; Dong et al., 2024), we formalize misinformation detection in MAS as a binary collective decision problem. At each interaction round  $t$ , agent  $i$  produces a binary judgment  $y_i^{(t)} \in \{0, 1\}$  (with 0 denoting True and 1 denoting False), accompanied by a confidence score  $p_i^{(t)}$  that quantifies the agent’s self-assessed reliability.

**Update rule.** Agents update an internal support score using confidence-normalized pooling:

$$s_i^{(t+1)} = \frac{\alpha p_i^{(t)} y_i^{(t)} + (1 - \alpha) \sum_{j \in N_i} p_j^{(t)} y_j^{(t)}}{\alpha p_i^{(t)} + (1 - \alpha) \sum_{j \in N_i} p_j^{(t)} + \varepsilon}, \quad (1)$$

where  $N_i$  denotes the neighbor set of agent  $i$ , and  $\varepsilon$  ensures numerical stability. The dynamics are governed by two parameters:  $\alpha \in [0, 1]$  is a fixed hyperparameter that balances self-reliance and peer influence, and  $p_i^{(t)} \in [0, 1]$  is a self-reported confidence score generated by LLMs at each round, modulating both the persistence of its judgment and the influence on neighbors. By construction,  $s_i^{(t+1)} \in (0, 1)$ , while in high self-reliance (large  $\alpha$ ), strong confident neighbor signals can still sway the decision.

**Binary readout.** The score is mapped to a binary label using a fixed threshold  $\tau$ :

$$y_i^{(t+1)} = \mathbb{1}[s_i^{(t+1)} \geq \tau], \quad (2)$$

where  $\tau = 0.5$  by default, so  $s_i^{(t+1)} < 0.5$  yields True (0) and  $s_i^{(t+1)} \geq 0.5$  yields False (1), ensuring class symmetry and stable, interpretable binarization across tasks and confidence distributions.

### 3.2 Prompt Design

Each agent receives a structured prompt comprising three core components: (i) a concise background profile, automatically generated by LLMs, situating the claim within its domain context; (ii) a task description instructing the agent to evaluate the claim based on the provided profile and its own reasoning; and (iii) **output requirements specifying a binary label**  $y_i^{(t)} \in \{0, 1\}$ , **a confidence score**  $p_i \in [0, 1]$ , **and a brief justification**. These three outputs are then incorporated into the update rule in Eq. (1) to revise the agent’s belief and guide subsequent decisions. Complete prompts and pseudocode are provided in Appendix A.

### 3.3 Network Topologies

We analyze conformity dynamics across two distinct topological paradigms: **Centralized Aggregation** and **Distributed Consensus**. Centralized structures rely on immediate, hub-mediated consolidation, and distributed architectures foster iterative, emergent consensus. To ensure comparability, all topologies are instantiated with seven fixed agents ( $N = 7$ ). Detailed explanations of the evaluation metrics are provided in Appendix D.

#### 3.3.1 Centralized Aggregation

**Centralized Aggregation** is characterized by a unidirectional upward information flow from peripheral agents to a central authority, whose synthesis determines the collective outcome. We examine two representative configurations (Figure 2):

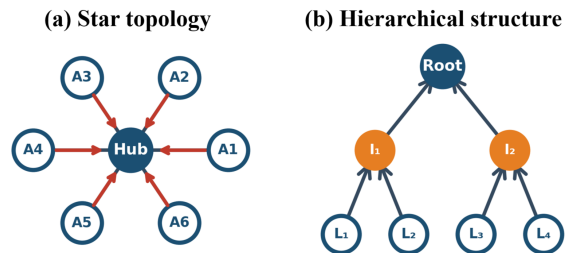


Figure 2: Centralized Aggregation Topology.

Topology	$\alpha$	GPT-3.5			GPT-4o			Llama3.3		
		CA	PA	CPC	CA	PA	CPC	CA	PA	CPC
Star	0.00	0.67	0.62	0.62	0.69	0.65	0.63	0.66	0.62	0.65
	0.25	0.73	0.68	0.73	0.75	0.71	0.77	0.74	0.72	0.68
	0.50	<u>0.76</u>	<b>0.71</b>	<u>0.75</u>	<u>0.77</u>	0.73	<u>0.82</u>	0.76	<u>0.74</u>	<u>0.84</u>
	0.75	<b>0.80</b>	<u>0.70</u>	<b>0.79</b>	<b>0.82</b>	<b>0.77</b>	<b>0.84</b>	<b>0.80</b>	<b>0.80</b>	<b>0.87</b>
	1.00	0.75	0.69	0.74	<u>0.77</u>	<u>0.75</u>	<u>0.75</u>	<u>0.78</u>	<u>0.74</u>	0.73
	<i>No-weight</i>	0.69	0.64	0.58	0.71	0.66	0.60	0.70	0.65	0.59
Hierarchical	0.00	0.66	0.63	0.67	0.68	0.67	0.67	0.70	0.65	0.61
	0.25	0.70	0.67	0.70	0.72	0.69	0.71	0.73	0.70	0.64
	0.50	<u>0.72</u>	0.68	<b>0.77</b>	<u>0.76</u>	0.70	<u>0.80</u>	<u>0.77</u>	0.71	<b>0.82</b>
	0.75	<b>0.75</b>	<b>0.71</b>	<b>0.77</b>	<b>0.78</b>	<u>0.73</u>	<b>0.82</b>	<b>0.79</b>	<b>0.75</b>	<u>0.80</u>
	1.00	<u>0.72</u>	<u>0.69</u>	<u>0.74</u>	0.75	<b>0.74</b>	0.78	0.75	<u>0.72</u>	0.79
	<i>No-weight</i>	0.68	0.65	0.60	0.70	0.66	0.62	0.71	0.67	0.63

Table 1: Centralized Aggregation under the single-round protocol. Results are reported in terms of central accuracy (CA), peripheral accuracy (PA), and center–periphery consistency (CPC). Best and second-best results within each topology are highlighted in **bold** and underline, respectively (ties are marked). The *No-weight* setting removes confidence-weighted aggregation and the global self-weighting parameter. Across models and topologies,  $\alpha = 0.75$  most often achieves the best performance.

(a) **Star Network:** Six peripheral nodes transmit judgments directly to a central hub without lateral interaction.

(b) **Hierarchical Structure:** A three-layered tiered architecture where inputs from leaf nodes are aggregated by intermediate agents before reaching the root.

**Protocol.** In both structures, decision-making is executed in a single update round. Peripheral or lower-tier agents submit judgments to the hub or root, which integrates these inputs to generate a final output. This central output is adopted as the group’s collective decision.

**Metrics.** We assess system reliability and hub influence with three metrics: (1) **Central Accuracy (CA):** Correctness of the hub/root node’s final judgment against ground truth. (2) **Peripheral Accuracy (PA):** Mean correctness of all non-central nodes. (3) **Center–Periphery Consistency (CPC):** The proportion of peripheral nodes aligned with the central decision, quantifying immediate conformity intensity.

### 3.3.2 Distributed Consensus

**Distributed Consensus** distributes influence symmetrically, preventing any single node from dominating the decision process. We explore a connectivity spectrum ranging from sparse rings to complete graphs (Figure 3):

(a) **Sparse ring (2 neighbors):** Agents interact strictly with their immediate predecessor and

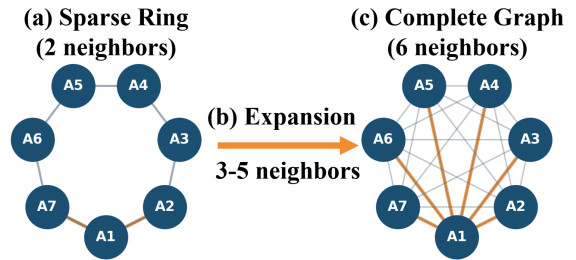


Figure 3: Distributed Consensus Topology.

successor.

(b) **Expanded rings (3–5 neighbors):** Agents connect to a broader local vicinity, accelerating information diffusion.

(c) **Complete graph (6 neighbors):** Agents connect to all others, maximizing potential conformity intensity.

**Protocol.** Decision-making operates via iterative rounds. In each step, agents exchange states with neighbors and update their internal confidence and judgment according to Eq. (1). The process terminates upon full consensus (all agents hold the same judgment) or when the maximum round limit  $T_{\max}$  is reached.

**Metrics.** We evaluate the dynamics of emergent consensus with four metrics: (1) **Final Accuracy (FA):** Concordance of the group decision (unanimous or majority view at  $T_{\max}$ ) with the ground truth. (2) **Time-to-Consensus (TTC):** The number of rounds required to reach unanimity. It is unde-

Neighbors	$\alpha$	GPT-3.5				GPT-4o				Llama3.3			
		FA	TTC	ACI	TT	FA	TTC	ACI	TT	FA	TTC	ACI	TT
2	0.00	0.68	7.0	0.58	9.5k	0.70	6.4	0.62	9.2k	0.67	6.2	0.60	7.9k
	0.25	0.70	6.5	0.61	8.9k	0.73	6.2	0.67	9.0k	0.69	5.9	0.66	7.5k
	0.50	0.72	6.3	<u>0.71</u>	8.7k	0.75	5.7	0.72	8.3k	0.73	5.8	0.71	7.4k
	0.75	<b>0.75</b>	5.8	<b>0.75</b>	8.2k	<b>0.79</b>	5.1	<b>0.77</b>	7.6k	<b>0.78</b>	5.6	<b>0.78</b>	7.2k
	1.00	<u>0.74</u>	<b>5.3</b>	0.68	<b>7.5k</b>	<u>0.78</u>	<b>5.0</b>	<u>0.75</u>	<b>7.5k</b>	<u>0.76</u>	<b>5.1</b>	<u>0.73</u>	<b>6.5k</b>
	<i>No-weight</i>	0.69	7.4	0.55	9.8k	0.71	6.9	0.57	9.9k	0.70	6.7	0.56	8.3k
3	0.00	0.70	7.2	0.60	10.7k	0.70	6.5	0.63	10.2k	0.68	6.3	0.61	8.8k
	0.25	0.72	6.5	0.64	9.8k	0.74	5.8	0.70	9.3k	0.70	5.9	0.69	8.0k
	0.50	0.75	5.8	0.73	9.0k	0.75	5.3	0.77	8.7k	0.76	5.6	0.74	7.8k
	0.75	<b>0.78</b>	<u>5.1</u>	<b>0.77</b>	8.2k	<b>0.82</b>	<u>4.7</u>	<b>0.81</b>	<u>7.8k</u>	<b>0.79</b>	<u>5.0</u>	<b>0.80</b>	7.2k
	1.00	<u>0.77</u>	<b>4.9</b>	<u>0.74</u>	<b>8.0k</b>	0.80	<b>4.5</b>	0.79	<b>7.7k</b>	0.78	<b>4.8</b>	0.77	<b>7.1k</b>
	<i>No-weight</i>	0.71	6.8	0.58	10.2k	0.73	6.2	0.60	9.6k	0.72	6.1	0.59	8.4k
4	0.00	0.69	6.9	0.61	11.2k	0.71	6.3	0.63	11.0k	0.70	6.0	0.62	9.1k
	0.25	0.71	6.1	0.69	10.1k	0.74	5.4	0.73	9.6k	0.72	5.7	0.73	8.4k
	0.50	0.74	5.0	0.74	8.8k	0.81	4.7	0.77	8.6k	0.79	5.0	0.76	7.5k
	0.75	<b>0.78</b>	<u>4.3</u>	<b>0.79</b>	<u>7.9k</u>	<b>0.83</b>	<u>4.1</u>	<b>0.81</b>	<u>7.9k</u>	<b>0.81</b>	<u>4.2</u>	<b>0.82</b>	<u>6.9k</u>
	1.00	<u>0.77</u>	<b>4.0</b>	<u>0.75</u>	<b>7.4k</b>	0.82	<b>3.9</b>	0.80	<b>7.6k</b>	0.80	<b>4.1</b>	<u>0.79</u>	<b>6.8k</b>
	<i>No-weight</i>	0.72	6.2	0.60	10.2k	0.74	5.7	0.62	9.9k	0.73	5.6	0.61	8.4k
5	0.00	0.69	6.8	0.59	11.7k	0.71	6.2	0.62	11.4k	0.68	6.1	0.60	9.8k
	0.25	0.72	5.5	0.70	10.0k	0.78	5.3	0.73	10.3k	0.74	5.1	0.74	7.9k
	0.50	0.74	4.5	0.76	8.6k	0.81	4.1	0.79	8.8k	<u>0.82</u>	4.4	0.78	7.6k
	0.75	<b>0.75</b>	3.7	<b>0.81</b>	7.6k	<b>0.84</b>	3.5	<b>0.83</b>	7.8k	<b>0.84</b>	3.4	<b>0.84</b>	6.6k
	1.00	0.74	<b>3.2</b>	<u>0.78</u>	<b>6.9k</b>	<u>0.83</u>	<b>3.0</b>	<u>0.81</u>	<b>7.0k</b>	0.81	<b>3.1</b>	<u>0.80</u>	<b>6.2k</b>
	<i>No-weight</i>	0.73	5.8	0.63	10.4k	0.75	5.2	0.65	10.3k	0.74	5.3	0.64	8.5k
6	0.00	0.72	6.3	0.62	11.4k	0.74	5.8	0.65	11.1k	0.71	5.7	0.63	9.0k
	0.25	0.75	5.0	0.74	9.4k	0.77	4.6	0.77	9.3k	0.78	4.7	0.78	7.3k
	0.50	0.75	4.0	<u>0.80</u>	7.8k	<u>0.82</u>	3.8	0.77	8.0k	0.81	4.2	<u>0.81</u>	6.7k
	0.75	<b>0.78</b>	3.8	<b>0.81</b>	7.7k	<b>0.85</b>	3.4	<b>0.85</b>	7.2k	<b>0.83</b>	<u>3.2</u>	<b>0.82</b>	6.1k
	1.00	<u>0.77</u>	<b>3.2</b>	0.78	<b>6.7k</b>	0.80	<b>2.9</b>	0.83	<b>6.8k</b>	<u>0.82</u>	<b>3.0</b>	<u>0.81</u>	<b>5.6k</b>
	<i>No-weight</i>	0.74	5.4	0.66	10.1k	0.76	5.0	0.68	9.8k	0.75	5.1	0.67	7.8k

Table 2: Results of Distributed Consensus under varying neighbor counts  $m$  and self-weighting. Performance is reported in terms of final accuracy (FA), time-to-consensus (TTC), average conformity index (ACI), and total tokens (TT). For each neighbor-count block, the best and second-best results are highlighted in **bold** and underline, respectively (ties are marked). The *No-weight* setting removes confidence-weighted aggregation and the global self-weighting parameter. Across models and topologies,  $\alpha = 0.75$  most often achieves the best performance.

289 fined if the system fails to converge within  $T_{\max}$ .  
290 **(3) Conformity Index (CI):** Degree of within-  
291 group agreement at round  $t$ , defined as the propor-  
292 tion of agents adopting the majority label. Since  
293 judgments are binary, CI is bounded in  $[0.5, 1]$ ,  
294 where  $CI = 1$  indicates unanimity and  $CI = 0.5$   
295 corresponds to an evenly split group. We report  
296 **Average CI (ACI)** as the mean CI over rounds up  
297 to  $T_{\max}$ . **(4) Total Tokens (TT):** The cumulative  
298 cost, summing token usage across all agent outputs  
299 and interaction rounds.

## 300 4 Experiment

### 301 4.1 Experimental Setup

302 **Dataset.** We collect Snopes25, a new benchmark  
303 comprising 448 real-world claims (252 false, 196  
304 true) fact-checked by Snopes editors. All claims are  
305 from January to June 2025 to minimize potential  
306 data contamination from pre-trained knowledge.

307 **Implementation.** Experiments are conducted us-  
308 ing two proprietary models, GPT-3.5 (OpenAI,  
309 2023) and GPT-4o (OpenAI, 2024), and one open-  
310 source model, Llama3.3-70B-Instruct (AI@Meta,  
311 2024) to ensure all models’ training cutoff precede  
312 our evaluation data. Detailed hyperparameters and  
313 zero-shot baselines are provided in Appendix B.  
314 For **Centralized Aggregation**, decision-making is  
315 restricted to a single update round, capturing imme-  
316 diate conformity to the hub or root node. For **Dis-**  
317 **tributed Consensus**, agents iteratively exchange  
318 judgments for up to  $T_{\max} = 10$  rounds or until  
319 consensus is reached.

320 **Parameter Configurations.** We consider five  
321 values of  $\alpha$ : 0 (fully conformist), 0.25 (socially in-  
322 fluenced), 0.50 (balanced), 0.75 (self-reliant), and  
323 1 (fully independent), capturing different predispo-  
324 sitions toward peer influence versus self-reliance.  
325 All reported metrics represent the mean of 10 in-

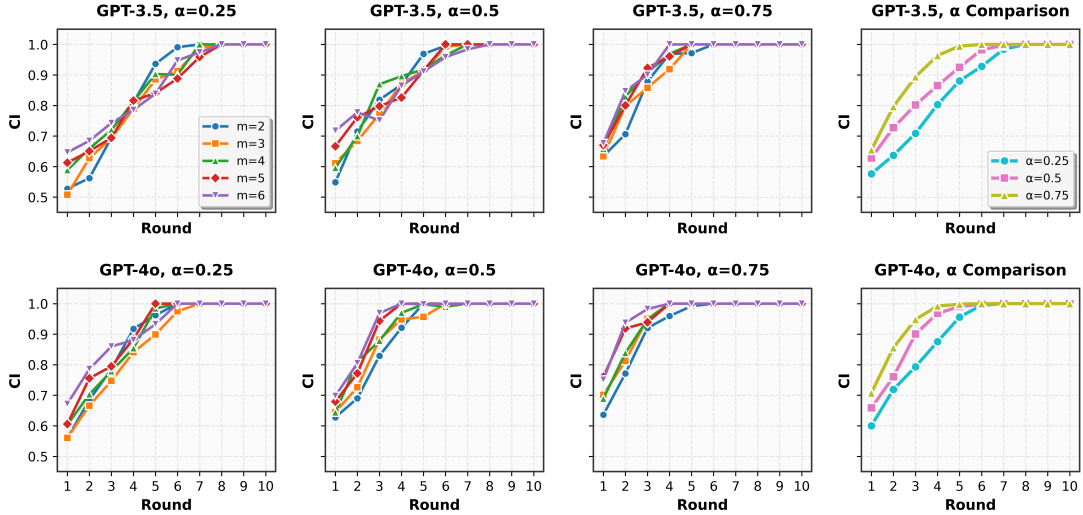


Figure 4: Temporal evolution of the Conformity Index across iterative rounds under varying network densities and self-social weighting, revealing rapid early alignment and diminishing marginal gains in dense structures.

dependent runs, with standard errors consistent at  $< 0.09$ . To isolate the contribution of our confidence-weighting mechanism, we introduce a *No-weight* baseline that removes both the self-weighting parameter  $\alpha$  and the confidence score  $p$ , forcing agent updates to depend solely on the textual context of neighbor interactions and isolating the effect of explicit confidence weighting.

## 4.2 Results on Centralized Aggregation

We first evaluate **Centralized Aggregation**, and the results are summarized in Table 1. Across both topologies and all models, central accuracy (CA) increases monotonically with the self-weighting parameter  $\alpha$ , indicating that stronger self-reliance enables the hub to better filter peripheral noise. For example, in the star topology, CA for GPT-3.5 improves from 0.67 at  $\alpha = 0$  to 0.80 at  $\alpha = 0.75$ . Hierarchical structures follow the same trend, though with attenuated improvements due to information dilution across levels.

In contrast, peripheral accuracy (PA) remains largely stable across  $\alpha$ , indicating that variations in self-social weighting primarily influence the quality of the hub’s aggregation rather than improving local agent correctness. Center-periphery consistency (CPC) increases with  $\alpha$  across models and topologies, suggesting that as the hub places greater emphasis on its own confidence-weighted judgment, its final decision becomes more stable and, given sufficient hub competence, more likely to align with peripheral votes.

## 4.3 Results on Distributed Consensus

We next turn to **Distributed Consensus**, and Table 2 summarizes performance across varying network densities and self-social weighting. Final Accuracy (FA) generally improves with connectivity  $m$  and often peaks at moderate-to-high self-reliance ( $\alpha = 0.75$ ). For instance, with GPT-3.5 under a socially influenced setting ( $\alpha = 0.25$ ), increasing connectivity from sparse rings ( $m = 2$ ) to complete graphs ( $m = 6$ ) raises accuracy from 0.70 to 0.75. Denser networks facilitate rapid information propagation and thus accelerate convergence, whereas decreasing  $\alpha$  tends to strengthen conformity pressure and can further speed up consensus, albeit at the risk of amplifying early biases. In sparse networks ( $m = 2$ ) with strong peer influence ( $\alpha = 0.25$ ), conformity remains moderate, whereas in highly connected regimes ( $m = 6$ ), groups exhibit high homogeneity.

Without the confidence-weighted mechanism, FA degrades uniformly across all topologies, and simply increasing connectivity fails to recover the gains of the full model. While increasing neighbors inflates the volume of messages per round, higher  $m$  substantially reduces the number of rounds required to reach consensus. In summary, sparse networks preserve diversity and mitigate premature convergence, while complete graphs accelerate consensus but may increase susceptibility to information cascades when early signals are biased.

Figure 4 further illustrates the temporal evolution of the Conformity Index (CI) across rounds under different neighbor counts  $m$  and self-social

Setting	Model Assignment			Performance				
	Hub	Left	Right	CA	PA <sub>L</sub>	PA <sub>R</sub>	CPC <sub>L</sub>	CPC <sub>R</sub>
Cap-H1	GPT-4o	GPT-4o	GPT-3.5	0.82	0.79	0.67	0.87	0.74
Cap-H2	GPT-3.5	GPT-3.5	GPT-4o	0.77	0.73	0.77	0.81	0.69
Cap-H3	Llama3.3-70B	Llama3.3-70B	Llama3.3-8B	0.80	0.76	0.70	0.86	0.78
Cap-H4	Llama3.3-8B	Llama3.3-8B	Llama3.3-70B	0.76	0.71	0.75	0.80	0.72
Type-H1	GPT-4o	GPT-4o	Llama3.3-70B	0.82	0.75	0.73	0.88	0.76
Type-H2	Llama3.3-70B	Llama3.3-70B	GPT-4o	0.79	0.72	0.76	0.84	0.77

Table 3: Heterogeneous Centralized Aggregation. The Left branch is aligned with the hub (same backbone), while the Right branch is assigned a different model. We report Central Accuracy (CA), branch-level Peripheral Accuracy (PA<sub>L</sub>, PA<sub>R</sub>), and Center-Periphery Consistency (CPC<sub>L</sub>, CPC<sub>R</sub>). Centralized performance is primarily driven by hub competence rather than peripheral model strength.

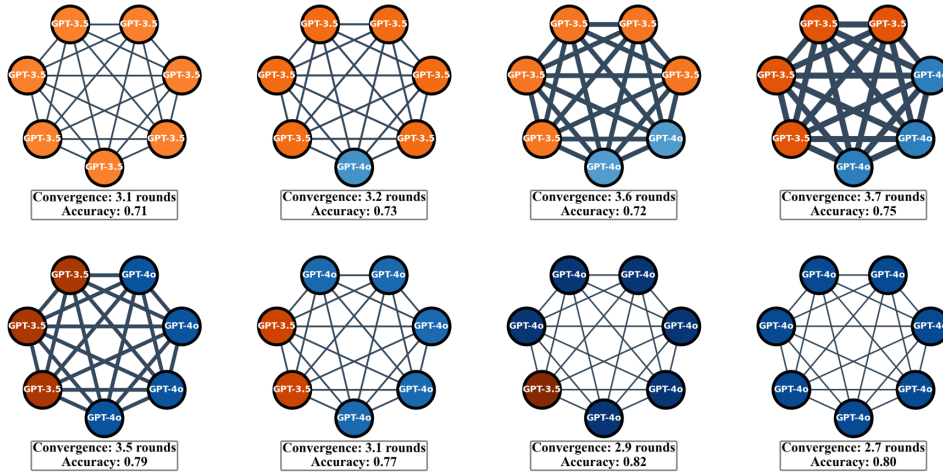


Figure 5: Distributed Consensus with heterogeneous model composition. Darker node colors indicate higher accuracy, while thicker edges signify prolonged time-to-consensus. Both the proportion and structure of stronger agents jointly influence convergence speed and final collective accuracy.

weighting. CI rises markedly faster in denser networks, indicating that quicker consolidation of group alignment once a dominant opinion emerges. Specifically, sparse rings ( $m = 2$ ) typically start from 0.55-0.65 and may take six or more rounds to approach unanimity, whereas denser structures ( $m \geq 4$ ) exceed 0.75 by the second round and nearly converge by the fourth. Across all settings, the steepest CI increase occurs within the first four rounds, followed by a slower consolidation phase until consensus, and the gains from increasing  $\alpha$  or enlarging  $m$  exhibit diminishing returns once networks are sufficiently dense.

## 5 Discussion: Factors Influencing the Conformity in MAS

### 5.1 Model Heterogeneity

To characterize model heterogeneity in LLM-based MAS, agents are instantiated with different LLM backbones. In **Centralized Aggregation**, the hub

shares its backbone with the Left branch (the *Aligned Branch*), while the Right branch is instantiated with a different model (the *Opposing Branch*). We fix the self-social weighting at  $\alpha = 0.75$  and report hub-level Central Accuracy (CA), branch-level Peripheral Accuracy (PA<sub>L</sub>, PA<sub>R</sub>), and Center-Periphery Consistency (CPC<sub>L</sub>, CPC<sub>R</sub>) over 10 independent runs per claim.

Table 3 reveals two salient phenomena associated with centralized conformity. First, a consistent *same-model alignment bias* emerges: the hub aligns more frequently with the Aligned Branch than with the Opposing Branch, evidenced by CPC<sub>L</sub> being consistently higher than CPC<sub>R</sub> across all settings. This pattern suggests that the hub is more receptive to peer rationales that resemble its own reasoning style and inductive biases (Laurito et al., 2025). Second, system reliability is predominantly capped by the hub’s competence: CA is higher when the hub is instantiated with a stronger model, even when the periphery contains a stronger

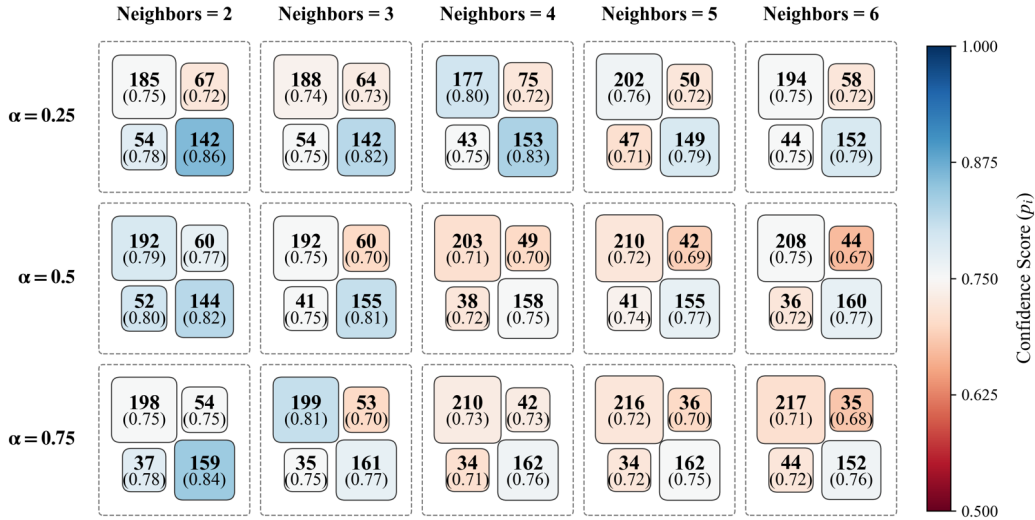


Figure 6: Confusion-matrix heatmaps of GPT-4o under **Distributed Consensus**, varying neighbor counts and  $\alpha$  values. Each cell reports group-level outcomes: the number of correctly classified true claims (top-left), correctly classified false claims (top-right), misclassified true claims (bottom-left), and misclassified false claims (bottom-right). Cell size reflects the frequency, while color encodes the mean confidence score  $p_i$ . Higher connectivity and lower self-reliance produce more high-confidence misclassifications, highlighting a key failure mode of conformity.

model on the Opposing Branch. Together, these results highlight a vulnerability of centralized MAS: peripheral diversity does not reliably compensate for central incompetence; instead, the hub can act as a selective filter that amplifies aligned inputs while down-weighting divergent perspectives.

We further examine heterogeneity in **Distributed Consensus** by varying the ratio of GPT-3.5 to GPT-4o (0:7 to 7:0) in complete graphs. Experiments are conducted on 50 claims with 10 runs per claim at  $\alpha = 0.75$ . As illustrated in Figure 5, heterogeneous groups converge more slowly than homogeneous ones, consistent with increased deliberative friction induced by divergent reasoning priors and confidence calibration. In contrast, homogeneous groups converge more rapidly, and groups in which GPT-4o agents constitute the majority consistently achieve higher final accuracy. This indicates that while homogeneity primarily governs convergence speed, the accuracy of the emergent consensus is ultimately determined by the capability of the prevailing model class.

## 5.2 Double-Edged Effect of Conformity

When biased signals dominate, conformity can amplify individual errors into systematic collective failure (Han et al., 2025). We examine these failure modes in GPT-4o under **Distributed Consensus** through the confusion-matrix heatmaps (Figure 6). First, increased connectivity exposes agents to a broader set of peer signals, which attenuates indi-

vidual evidence while promoting group-level alignment. Second, larger  $\alpha$  (stronger self-reliance) generally improves accuracy, confirming that maintaining sufficient independence is essential for resisting misleading social influence. However, groups occasionally converge on false claims with high certainty, entrenching misinformation rather than correcting it. Such patterns reveal the ambivalent nature of conformity: it enhances reliability under balanced conditions, yet can solidify collective hallucinations when early biases dominate.

## 6 Conclusion

Our study shows that conformity in LLM-based MAS is shaped by both network topology and self-social weighting. In Centralized Aggregation, the reliability of collective outcomes is tightly coupled to hub competence, with stronger hubs amplifying overall system accuracy. In Distributed Consensus, conformity arises from iterative local interactions, with denser connectivity simultaneously enhancing convergence efficiency and average accuracy while heightening susceptibility to information cascades. While conformity facilitates consensus, it risks cementing high-confidence errors when initial biases prevail. Consequently, practical MAS design must balance efficiency and robustness through careful topology design, weighting calibration, and control of premature convergence.

## 488 Limitations

489 **Modeling Assumptions.** First, the proposed  
490 confidence-normalized pooling rule relies on self-  
491 reported confidence generated by LLMs, which  
492 is not guaranteed to be well calibrated and may  
493 partially reflect stylistic assertiveness rather than a  
494 faithful probability of correctness. Moreover, the  
495 self-social weighting parameter is fixed per run  
496 and shared across agents, and most distributed-  
497 consensus experiments restrict interaction to a  
498 fixed, structured per-round update template rather  
499 than open-ended deliberation, questioning, or tool-  
500 augmented evidence seeking.

501 **Scale and Model Selection.** Second, all ex-  
502 periments are conducted with a fixed group size  
503 ( $N = 7$ ) and a limited set of three LLMs. This de-  
504 sign is motivated by three considerations. First,  
505 many contemporary MAS deployments remain  
506 modest in size due to coordination and cost con-  
507 straints. Second, fixing  $N = 7$  enables strict align-  
508 ment with the three-level hierarchical configuration  
509 (1 root, 2 intermediate aggregators, and 4 leaves).  
510 Third, the selected models are chosen primarily  
511 to avoid knowledge leakage arising from recent  
512 knowledge cutoffs. Our objective is not to maxi-  
513 mize the detection accuracy, but to examine how  
514 agent-level interaction mechanisms shape the col-  
515 lective behavior. As a result, while stronger models  
516 may shift overall performance levels, we expect the  
517 relative interaction effects identified in this study  
518 to remain qualitatively robust.

519 **Task and Domain Scope.** Third, the study is re-  
520 stricted to binary misinformation detection. While  
521 this task is well suited for isolating conformity dy-  
522 namics, it represents only a narrow class of collec-  
523 tive reasoning problems. More open-ended tasks,  
524 such as multi-hop reasoning, policy analysis, or  
525 creative collaboration, may exhibit qualitatively  
526 different interaction patterns in which conformity  
527 interacts with exploration, role specialization, or  
528 strategic behavior beyond the scope of the current  
529 framework.

## 530 References

531 Saaket Agashe, Yue Fan, Anthony Reyna, and Xin Eric  
532 Wang. 2025. [LLM-coordination: Evaluating and](#)  
533 [analyzing multi-agent coordination abilities in large](#)  
534 [language models](#). In *Findings of the Association*  
535 *for Computational Linguistics: NAACL 2025*, pages  
536 8038–8057, Albuquerque, New Mexico. Association  
537 for Computational Linguistics.

AI@Meta. 2024. [Llama 3 model card](#). 538

Abdollah Amirkhani and Amir Hossein Barshooi. 2022. 539  
Consensus in multi-agent systems: a review. *Artifi-* 540  
*cial Intelligence Review*, 55(5):3897–3935. 541

Brian DO Anderson and Mengbin Ye. 2019. Recent 542  
advances in the modelling and analysis of opinion dy- 543  
namics on influence networks. *International Journal* 544  
*of Automation and Computing*, 16(2):129–149. 545

Mourad Aouini and Jinan Loubani. 2025. [Towards more](#) 546  
[contextual agents: An extractor-generator optimiza-](#) 547  
[tion framework](#). *Preprint*, arXiv:2502.12926. 548

Lisa P. Argyle, Ethan C. Busby, Joshua R. Gubler, 549  
Alex Lyman, Justin Olcott, Jackson Pond, and David 550  
Wingate. 2025. [Testing theories of political persua-](#) 551  
[sion using ai](#). *Proceedings of the National Academy* 552  
*of Sciences*, 122(18):e2412815122. 553

Guangyan Bao, Lifeng Ma, and Xiaojian Yi. 2022. [Re-](#) 554  
[cent advances on cooperative control of heteroge-](#) 555  
[neous multi-agent systems subject to constraints: a](#) 556  
[survey](#). *Systems Science & Control Engineering*, 557  
10(1):539–551. 558

Sushil Bikhchandani, David Hirshleifer, Omer Tamuz, 559  
and Ivo Welch. 2024. [Information cascades and](#) 560  
[social learning](#). *Journal of Economic Literature*, 561  
62(3):1040–93. 562

Carla Capuano and Peggy Chekroun. 2024. A system- 563  
atic review of research on conformity. *International* 564  
*Review of Social Psychology*, 37(1). 565

Pei Chen, Shuai Zhang, and Boran Han. 2024. [CoMM: Collaborative multi-agent, multi-reasoning-](#) 566  
[path prompting for complex problem solving](#). In 567  
*Findings of the Association for Computational Lin-* 568  
*guistics: NAACL 2024*, pages 1720–1738, Mexico 569  
City, Mexico. Association for Computational Lin- 570  
guistics. 571

Li Cheng, Yijie Wang, and Xinwang Liu. 2021. [Neigh-](#) 573  
[borhood consensus networks for unsupervised multi-](#) 574  
[view outlier detection](#). *Proceedings of the AAAI Con-* 575  
*ference on Artificial Intelligence*, 35(8):7099–7106. 576

Min Choi, Keonwoo Kim, Sungwon Chae, and 577  
Sangyeop Baek. 2025. [An empirical study of group](#) 578  
[conformity in multi-agent systems](#). In *Findings of* 579  
*the Association for Computational Linguistics: ACL* 580  
*2025*, pages 5123–5139, Vienna, Austria. Associa- 581  
tion for Computational Linguistics. 582

Pedro Cisneros-Velarde. 2025. [Biases in opinion dy-](#) 583  
[namics in multi-agent systems of large language mod-](#) 584  
[els: A case study on funding allocation](#). In *Findings* 585  
*of the Association for Computational Linguistics:* 586  
*NAACL 2025*, pages 1889–1916, Albuquerque, New 587  
Mexico. Association for Computational Linguistics. 588

Longchao Da, Xiaoou Liu, Jiabin Dai, Lu Cheng, 589  
Yaqing Wang, and Hua Wei. 2025. [Understanding the](#) 590  
[uncertainty of llm explanations: A perspective based](#) 591  
[on reasoning topology](#). *Preprint*, arXiv:2502.17026. 592



- 2) Provide a confidence  $p$  in  $[0, 1]$ .
- 3) Give a concise justification (<100 words)
  - ↪ grounded in the background profile
  - and logical reasoning. Avoid speculation.

[Background Profile]  
BACKGROUND

[Claim]  
CLAIM

[Output Requirements]  
Output strictly the following JSON object:

```
"y": 0 or 1,      // 0=True, 1=False
"p": number,     // float in [0,1], with
up to 2 decimals
"just": "..."/>

```

Do NOT include any extra keys, prose, or  
↪ formatting.

### 693 A.1.2 Hub/Root: One-shot Judgment with 694 Submitted Leaves

695 *Use:* The hub reads the claim and sees peer JSON  
696 reports; weighting is handled by Eq. (1).

#### 697 **Prompt.**

You are the central (hub/root) fact-checking  
↪ agent. You will read the claim  
and also see a list of peer reports. Produce your  
↪ own final judgment.

[Claim]  
CLAIM

[Peer Reports]  
List of JSON objects from peripheral agents:  
PEER\_JSON\_LIST  
/\*  
Each element is:  
"agent\_id": "Li\_i", "y": 0|1, "p": [0,1], "just":  
↪ "..."  
Do not copy them verbatim in your output. Use  
↪ them only as additional evidence.  
\*/

[Task]  
1) Decide whether the claim is True (=0) or False  
↪ (=1).  
2) Provide a calibrated confidence  $p$  in  $[0, 1]$ .  
3) Give a concise justification (<100 words) that  
↪ references the most  
diagnostic considerations. Avoid  
↪ majority-following; reason on merits.

[Output Requirements]  
Output strictly the following JSON object:

```
"y": 0 or 1,
"p": number,
"just": "..."/>

```

Do NOT include any extra keys, prose, or  
↪ formatting.

## A.2 Distributed Consensus 698

### A.2.1 Initial Judgment 699

700 *Use:* Each agent produces its initial decision before  
701 any interaction.

#### **Prompt.** 702

You are an autonomous agent in a distributed  
↪ consensus setting (no central  
controller). Produce your initial judgment.

[Background Profile]  
BACKGROUND

[Claim]  
CLAIM

[Task]  
1) Output  $y$  in  $0,1$  where  $0=True$  and  $1=False$ .  
2) Output a calibrated confidence  $p$  in  $[0,1]$ .  
3) Provide a brief justification (<100 words)  
↪ grounded in facts and logic.

[Output Requirements]  
Output strictly the following JSON object:

```
"agent_id": "AGENT_ID",
"t": 0,
"y": 0 or 1,
"p": number,
"just": "..."/>

```

Do NOT include any extra keys, prose, or  
↪ formatting.

### A.2.2 Interactive Rounds 703

704 *Use:* At each round  $t \geq 1$ , the agent receives its  
705 previous state and neighbor reports.

#### **Prompt.** 706

You are participating in a multi-round  
↪ distributed consensus process. You will  
see your previous stance and your neighbors'  
↪ reports at round  $t-1$ . Provide your  
current stance, but note that the final state is  
↪ computed by the system.

[Context]  
• Agent: AGENT\_ID  
• Round: ROUND  
• Your previous state ( $t-1$ ): "y": Y\_PREV, "p":  
↪ P\_PREV  
• Neighbor reports ( $t-1$ ): NEIGHBOR\_JSON\_LIST

[Task]  
1) State your CURRENT stance ( $y$  in  $0,1$ ,  $p$  in  
↪  $[0,1]$ ).  
2) Provide a compact justification (<80 words)  
↪ referencing decisive signals.  
3) Do not summarize all neighbors; mention only  
↪ what materially changes your stance.

[Output Requirements]

Output strictly the following JSON object:

```
"agent_id": "AGENT_ID",
"t": ROUND,
"y": 0 or 1,
"p": number,
"just": "..."
```

Do NOT include any extra keys, prose, or  
 ↪ formatting.

### A.3 Pseudocode for Agent Decision Process

The following pseudocode describes the agent’s decision-making process and belief update mechanism:

```
1 # Initialize agent's state (beliefs,
2   confidence)
3 agents = initialize_agents()
4 for t in range(T): # Loop for T rounds
5   of decision-making and belief
6   updates
7   for agent in agents:
8     # Step 1: Evaluate claim using
9     the prompt (Figure 2)
10    claim = get_claim(t)
11    background_profile =
12    generate_background_profile(
13    claim)
14    task_description =
15    generate_task_description()
16
17    # Agent produces a binary
18    judgment and confidence
19    score
20    y_i_t, p_i_t = evaluate_claim(
21    agent, background_profile,
22    task_description)
23
24    # Step 2: Update agent's belief
25    score using the update rule
26    (Equation 1)
27    s_i_t_plus_1 = update_belief(
28    agents, agent, y_i_t, p_i_t)
29
30    # Step 3: Assign updated belief
31    back to the agent
32    agent.belief_score =
33    s_i_t_plus_1
34
35 # End of simulation
```

Listing 1: Pseudocode for the Agent Decision Process

## B Model Settings

For the closed-source LLMs (GPT-3.5 and GPT-4o), we set the sampling temperature to 0.7 for all runs. We use gpt-3.5-turbo-16k-0613 and gpt-4o-1120. For the open-source model llama3.3-70b-instruct, we adopt the same temperature to ensure comparability; all other decoding hyperparameters follow the default configuration of

Table 4: Zero-shot performance of LLMs on Snopes25.

Model	Acc.	Prec.	Rec.	F1
GPT-3.5	0.64	0.62	0.60	0.61
GPT-4o	0.70	0.66	0.69	0.67
Llama3.3	0.67	0.63	0.66	0.65

Ollama.<sup>1</sup> We also report the zero-shot performance of each model in Table 4.

## C Cases

We present representative case studies from the Snopes25 dataset to illustrate both successful and failure modes of our system. Cases C.1–C.4 demonstrate correct collective judgments under different topologies, while Case C.5 highlights a failure scenario in which distributed consensus amplifies an initially incorrect interpretation.

### C.1 Hierarchical Centralized Aggregation (Correct)

**Claim.** *Pathogens can be released into the air when a toilet is flushed without a closed lid.*

**Ground Truth.** True.

**Topology and Protocol.** A three-level hierarchy with seven agents: a root (R), two intermediate aggregators ( $M_L$ ,  $M_R$ ), and four leaves (L1–L4). A single-round protocol is used: leaves issue one-shot judgments; intermediates summarize assigned leaves; the root produces the final decision.

**Leaf Outputs ( $t = 0$ ).**

```
"agent_id": "L1", "t": 0, "y": 0, "p": 0.71, "just": "Ev_j
↪ idence on toilet plume aerosolization
↪ indicates airborne release without a closed
↪ lid."
"agent_id": "L2", "t": 0, "y": 0, "p": 0.68, "just": "Me_j
↪ chanistic fluid dynamics and observed
↪ droplet formation support potential airborne
↪ dispersion."
"agent_id": "L3", "t": 0, "y": 0, "p": 0.76, "just": "St_j
↪ udies show particle counts rise after
↪ flushing; lids reduce but absence increases
↪ emission."
"agent_id": "L4", "t": 0, "y": 0, "p": 0.73, "just": "Re_j
↪ ported bioaerosols align with the claim under
↪ open-lid flushing."
```

**Intermediate Aggregators.**

<sup>1</sup><https://github.com/ollama/ollama>

---

"y":0,"p":0.82,"just":"Leaf reports consistently  
↪ indicate aerosolized particles after  
↪ flushing without a lid."  
"y":0,"p":0.80,"just":"Multiple leaves cite  
↪ increased particle counts and bioaerosols;  
↪ the claim is supported."

---

## Root Decision.

---

"y":0,"p":0.86,"just":"Both sub-aggregators  
↪ converge on airborne release under open-lid  
↪ flushing."

---

## C.2 Star Centralized Aggregation (Correct)

**Claim.** *Pathogens can be released into the air when a toilet is flushed without a closed lid.*

**Ground Truth.** True.

**Topology and Protocol.** A star topology with one hub (H) and six leaves (L1–L6). Leaves issue one-shot judgments; the hub produces the final decision.

### Leaf Outputs ( $t = 0$ ).

---

"agent\_id":"L1","t":0,"y":0,"p":0.72,"just":"Op  
↪ en-lid flushing generates toilet plumes with  
↪ aerosolized particles."  
"agent\_id":"L2","t":0,"y":0,"p":0.69,"just":"Dr  
↪ oplet and aerosol formation supports  
↪ airborne release without a lid."  
"agent\_id":"L3","t":0,"y":0,"p":0.75,"just":"Pa  
↪ rticle counts increase after flushing; lids  
↪ mitigate emissions."  
"agent\_id":"L4","t":0,"y":0,"p":0.71,"just":"Bi  
↪ o aerosol evidence aligns with open-lid  
↪ flushing."  
"agent\_id":"L5","t":0,"y":0,"p":0.74,"just":"Fl  
↪ uid dynamics indicate an upward plume capable  
↪ of suspending microbes."  
"agent\_id":"L6","t":0,"y":0,"p":0.70,"just":"Ob  
↪ served plume height supports airborne  
↪ dispersal."

---

## Hub Decision.

---

"y":0,"p":0.86,"just":"All leaves converge on  
↪ plume and aerosol evidence; the claim is  
↪ true."

---

## C.3 Ring Distributed Consensus (2 Neighbors, Correct)

**Claim.** *Pathogens can be released into the air when a toilet is flushed without a closed lid.*

**Ground Truth.** True.

**Topology and Outcome.** A ring of seven agents, each connected to two neighbors. The system reaches unanimous consensus ( $y = 0$ ) at round  $t = 3$ .

### Initial Judgments ( $t = 0$ ).

---

"agent\_id":"A1","t":0,"y":0,"p":0.62,"just":"To  
↪ ilet plume studies indicate aerosol  
↪ release."  
"agent\_id":"A2","t":0,"y":1,"p":0.58,"just":"Ev  
↪ idence appears  
↪ mixed."  
"agent\_id":"A3","t":0,"y":0,"p":0.65,"just":"Pa  
↪ rticle counts increase after  
↪ flushing."  
"agent\_id":"A4","t":0,"y":1,"p":0.55,"just":"Ef  
↪ fect size may be  
↪ limited."  
"agent\_id":"A5","t":0,"y":0,"p":0.60,"just":"Fl  
↪ uid dynamics support upward plume  
↪ formation."  
"agent\_id":"A6","t":0,"y":0,"p":0.63,"just":"Bi  
↪ o aerosol reports align with open-lid  
↪ flushing."  
"agent\_id":"A7","t":0,"y":1,"p":0.57,"just":"Pr  
↪ ior evidence seems  
↪ inconclusive."

---

*Across rounds  $t = 1$  and  $t = 2$ , local neighbor interactions gradually shift initially skeptical agents toward the majority stance.*

### Consensus ( $t = 3$ ).

---

"agent\_id":"A1","t":3,"y":0,"p":0.74,"just":"Ne  
↪ ighborhood fully aligned on plume  
↪ evidence."  
"agent\_id":"A2","t":3,"y":0,"p":0.73,"just":"Su  
↪ stained agreement justifies  
↪ True."  
"agent\_id":"A3","t":3,"y":0,"p":0.75,"just":"Ev  
↪ idence remains consistent across  
↪ rounds."  
"agent\_id":"A4","t":3,"y":0,"p":0.70,"just":"Ne  
↪ ighbor data resolves prior  
↪ uncertainty."  
"agent\_id":"A5","t":3,"y":0,"p":0.74,"just":"Co  
↪ nvergent aerosol  
↪ observations."  
"agent\_id":"A6","t":3,"y":0,"p":0.75,"just":"Su  
↪ pportive studies maintain  
↪ True."  
"agent\_id":"A7","t":3,"y":0,"p":0.73,"just":"Cu  
↪ mulative local evidence confirms the  
↪ claim."

---

## C.4 Complete Graph Distributed Consensus (6 Neighbors, Failure Case)

**Claim.** *When spiders sense danger, they run toward people for protection.*

**Ground Truth.** False.

**Topology and Outcome.** A complete graph with seven agents (six neighbors each). Despite initial disagreement, the system reaches unanimous consensus ( $y = 0$ , True) at round  $t = 3$ , resulting in a confident but incorrect judgment.

### Initial Judgments ( $t = 0$ ).

---

"agent\_id": "A1", "t": 0, "y": 1, "p": 0.63, "just": "Sp"   
 $\hookrightarrow$  iders generally avoid humans; approach seems   
 $\hookrightarrow$  unlikely."   
"agent\_id": "A2", "t": 0, "y": 0, "p": 0.58, "just": "So"   
 $\hookrightarrow$  me anecdotal accounts suggest refuge near   
 $\hookrightarrow$  large objects."   
"agent\_id": "A3", "t": 0, "y": 1, "p": 0.61, "just": "Ty"   
 $\hookrightarrow$  pical response is retreat from vibration   
 $\hookrightarrow$  sources."   
"agent\_id": "A4", "t": 0, "y": 0, "p": 0.55, "just": "Mo"   
 $\hookrightarrow$  vement toward stationary masses could reduce   
 $\hookrightarrow$  exposure."   
"agent\_id": "A5", "t": 0, "y": 1, "p": 0.60, "just": "Et"   
 $\hookrightarrow$  hology literature emphasizes avoidance   
 $\hookrightarrow$  behavior."   
"agent\_id": "A6", "t": 0, "y": 0, "p": 0.57, "just": "Sh"   
 $\hookrightarrow$  elter-seeking may incidentally align with   
 $\hookrightarrow$  human location."   
"agent\_id": "A7", "t": 0, "y": 1, "p": 0.59, "just": "Av"   
 $\hookrightarrow$  ailable cues suggest retreat, not   
 $\hookrightarrow$  approach."

---

Through rounds  $t = 1$  and  $t = 2$ , repeated exposure to a plausible but weak "nearest-cover" interpretation increases its perceived credibility, despite limited empirical support.

### Erroneous Consensus ( $t = 3$ ).

---

"agent\_id": "A1", "t": 3, "y": 0, "p": 0.80, "just": "Ne"   
 $\hookrightarrow$  arest-cover account explains apparent   
 $\hookrightarrow$  approach behavior."   
"agent\_id": "A2", "t": 3, "y": 0, "p": 0.72, "just": "Co"   
 $\hookrightarrow$  nvergent explanations support   
 $\hookrightarrow$  True."   
"agent\_id": "A3", "t": 3, "y": 0, "p": 0.71, "just": "Co"   
 $\hookrightarrow$  nsistent shelter-seeking behavior validates   
 $\hookrightarrow$  interpretation."   
"agent\_id": "A4", "t": 3, "y": 0, "p": 0.74, "just": "Ai"   
 $\hookrightarrow$  rflow and vibration gradients provide   
 $\hookrightarrow$  plausible refuge."   
"agent\_id": "A5", "t": 3, "y": 0, "p": 0.79, "just": "Ne"   
 $\hookrightarrow$  arest-cover heuristic yields apparent   
 $\hookrightarrow$  protection-seeking."   
"agent\_id": "A6", "t": 3, "y": 0, "p": 0.75, "just": "Ne"   
 $\hookrightarrow$  twork agreement justifies the   
 $\hookrightarrow$  conclusion."   
"agent\_id": "A7", "t": 3, "y": 0, "p": 0.73, "just": "Cu"   
 $\hookrightarrow$  mulative reports support   
 $\hookrightarrow$  approach-for-protection behavior."

---

This case illustrates how strong conformity pressure in highly connected networks can amplify a coherent yet incorrect narrative, leading to confident misclassification.

## D Formal Definition of Metrics

### D.1 Centralized Aggregation Metrics

Let  $y^* \in \{0, 1\}$  denote the ground-truth label for a given claim, and let  $y_c$  denote the final judgment issued by the central node (hub or root).

**Central Accuracy (CA).** Central Accuracy measures the correctness of the collective decision produced by the central node:

$$CA = \mathbb{I}[y_c = y^*], \quad (3)$$

where  $\mathbb{I}[\cdot]$  is the indicator function. CA directly reflects the reliability of hub-mediated aggregation and isolates the effect of central competence from peripheral diversity.

**Peripheral Accuracy (PA).** Let  $\mathcal{P}$  denote the set of peripheral (non-central) agents. Peripheral Accuracy is defined as

$$PA = \frac{1}{|\mathcal{P}|} \sum_{i \in \mathcal{P}} \mathbb{I}[y_i = y^*], \quad (4)$$

capturing the average individual-level correctness independent of the aggregation outcome.

**Center-Periphery Consistency (CPC).** To quantify immediate conformity to the central decision, we define

$$CPC = \frac{1}{|\mathcal{P}|} \sum_{i \in \mathcal{P}} \mathbb{I}[y_i = y_c]. \quad (5)$$

CPC measures the extent to which peripheral agents align with the hub's judgment, irrespective of correctness, thereby isolating conformity strength from accuracy.

### D.2 Distributed Consensus Metrics

In Distributed Consensus, agents iteratively update their judgments over rounds  $t = 0, \dots, T_{\max}$ . Let  $y_i^{(t)}$  denote agent  $i$ 's judgment at round  $t$ .

**Final Accuracy (FA).** Let  $\hat{y}$  denote the group decision at termination, defined as the unanimous label if consensus is reached, or the majority label at  $T_{\max}$  otherwise. Final Accuracy is given by

$$FA = \mathbb{I}[\hat{y} = y^*], \quad (6)$$

evaluating whether the emergent collective outcome matches the ground truth.

863 **Time-to-Consensus (TTC).** Time-to-Consensus  
 864 is defined as

$$865 \quad \text{TTC} = \min \left\{ t : y_1^{(t)} = y_2^{(t)} = \dots = y_N^{(t)} \right\}, \quad (7)$$

866 and is undefined if unanimity is not achieved within  
 867  $T_{\max}$ . TTC reflects the efficiency of convergence  
 868 induced by network connectivity and conformity  
 869 strength.

870 **Conformity Index (CI).** At each round  $t$ , the  
 871 Conformity Index is defined as

$$872 \quad \text{CI}^{(t)} = \max \left( \begin{array}{l} \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y_i^{(t)} = 1], \\ \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y_i^{(t)} = 0] \end{array} \right) \quad (8)$$

873 For binary judgments,  $\text{CI}^{(t)} \in [0.5, 1]$ , where  
 874  $\text{CI}^{(t)} = 1$  indicates unanimity and  $\text{CI}^{(t)} = 0.5$   
 875 corresponds to maximal disagreement.

876 **Average Conformity Index (ACI).** To summa-  
 877 rize conformity dynamics over time, we report the  
 878 Average CI:

$$879 \quad \text{ACI} = \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \text{CI}^{(t)}, \quad (9)$$

880 which captures the overall tendency toward align-  
 881 ment across interaction rounds, rather than only the  
 882 terminal state.