

# Beyond Hallucinations: A Multi-Agent Approach to Fact-Checked AI Knowledge

Mandira Sawkar<sup>1</sup>

<sup>1</sup>Artificial Intelligence MS, Golisano College of Computing and Information Sciences, Rochester Institute of Technology, Rochester NY, USA

2025 January 28

## Abstract

Large Language Models (LLMs) excel at open-ended reasoning yet still hallucinate facts, undermining their reliability in high-stakes settings. Single-agent and retrieval-augmented approaches lessen some errors but cannot consistently verify claims or enforce agreement across diverse prompts. We introduce SC-MAKS, a self-consistent multi-agent framework that pairs heterogeneous LLMs with alternating cooperative-adversarial debate and an RLHF-trained consensus policy to fuse their outputs into one fact-checked answer. Across FEVER, TruthfulQA, Fakeddit, and three domain-specific case studies, SC-MAKS could cut hallucinations by 30–45% and lift factual consistency by 4–8% over strong multi-agent baselines while remaining computationally tractable.

**Keywords:** Large Language Models (LLMs) – Knowledge Synthesis – Self-Consistency – Fact-Checking – Hallucination Reduction – Knowledge Verification – Multi-Agent Collaboration – Debate & Consensus.

## 1 Introduction

### 1.1 Key Issues in Knowledge Synthesis

LLMs such as GPT-4, Claude, and LLaMA have transformed natural language processing, achieving state-of-the-art results in text generation, reasoning, and knowledge synthesis. However, they frequently hallucinate—producing fabricated or unsupported information—which undermines trust and downstream utility [1]. Additionally, LLMs often generate inconsistent responses when prompted with slight variations, posing reliability

challenges. These issues are particularly critical in high-stakes domains like scientific literature, legal analysis, and policy development, where factual accuracy is non-negotiable.

RAG (Retrieval-Augmented Generation) improves factual grounding by incorporating retrieved knowledge during generation. Yet, RAG alone doesn’t guarantee consistency across outputs and relies on static pipelines that may not reflect evolving information. Self-consistency decoding, as proposed by Wang et al. [2], improves accuracy by sampling multiple outputs from a single LLM and aggregating them, but lacks the diversity and inter-agent critique needed for robust verification.

### 1.2 Bridging the Gap: The Multi-Agent Approach

Multi-agent systems have long been explored in robotics, gaming, and trading, but their use for LLM collaboration is still emerging. Recent studies explore debate protocols among model instances to enhance reasoning. For instance, Irving et al. citeirving2018aisafetydebate demonstrated how structured debate between two models can expose hidden errors, while Du et al. [3] showed factual gains through homogeneous multi-agent debate. These findings suggest that adversarial and cooperative interactions among agents can outperform isolated reasoning.

### 1.3 The Case for Self-Consistent Multi-Agent Collaboration

Consider real-world scenarios: a medical researcher synthesizing treatment guidelines, a legal analyst verifying statutes, or a social media platform filtering misinformation. Human fact-checking in these settings is slow and labor-intensive. A self-consistent multi-agent LLM sys-

tem can automate much of this process by enforcing inter-agent agreement. The benefits include:

- **Redundancy mitigation:** Multiple agents reduce the risk of single-agent hallucinations.
- **Perspective diversity:** Different models bring unique knowledge, reducing correlated errors.
- **Built-in verification:** Debate phases serve as dynamic fact-checking rounds.

This approach synergizes generation, critique, and consensus to enable accurate and trustworthy AI output.

## 1.4 Introducing SC-MAKS

We propose SC-MAKS (Self-Consistent Multi-Agent Knowledge Synthesis), a unified framework addressing limitations of prior work. Key innovations include:

- **Diverse-Agent Weighted Voting:** SC-MAKS uses agents from distinct LLM architectures (GPT-4, Claude, LLaMA, DeepSeek) and assigns weights based on historical factual accuracy.
- **Hybrid Debate Roles:** Alternates between adversarial (challenge) and cooperative (constructive critique) roles to uncover both unsupported claims and refinable content.
- **RLHF-Optimized Consensus:** Uses reinforcement learning to fine-tune an aggregator that learns which debate signals best predict accuracy and coherence.

Sections 3–6 present the formal SC-MAKS design, empirical evaluation across datasets, and analysis of compute-accuracy trade-offs, demonstrating its practical benefits for high-stakes LLM applications.

Our main contributions begin with a detailed SC-MAKS framework with formal mathematical formulation and pseudocode in Section 3. Section 4 introduces rigorous ground-truth labeling protocol for standard benchmarks (FEVER, TruthfulQA, Fakeddit) and expert-annotated case studies. Section 5 outlines an experimental plan intended to demonstrate SC-MAKS’s superiority over three open-source multi-agent baselines across benchmark and domain-specific tasks. Section 6 will analyse compute-accuracy trade-offs, adaptive early-stopping heuristics, and guidelines for practical deployment in resource-constrained settings. Through these advances, SC-MAKS pushes the frontier of trustworthy, scalable

LLM collaboration, offering a robust solution for accurate knowledge synthesis in high-stakes applications.

# 2 Background & Related Work

## 2.1 Literature Review

### 2.1.1 Single-Agent Limitations and RAG

Traditional single-agent methods (e.g., GPT-4 with chain-of-thought) can generate coherent text but often hallucinate unsupported facts [1]. Retrieval-Augmented Generation (RAG) integrates external knowledge bases but does not ensure cross-response consistency.

### 2.1.2 Multi-Agent Debate Frameworks

Refer Table 1.

These studies highlight debate’s power to surface errors but leave room for dynamic weighting, heterogeneous roles, and end-to-end learnable aggregation.

### 2.1.3 Fact-Checking and Hallucination Reduction

Several multi-agent frameworks have emerged to tackle hallucinations by embedding verification steps into the generation process. Sun et al. (2024) [6] introduce a sequential debate over discrete claims: after extracting key statements, each claim is iteratively examined by multiple GPT-based agents that critique or defend based on retrieved evidence, reaching a stable verdict akin to a Markov Chain process. Yang et al. (2025) [7] combine self-reflection—where each model logs its own inconsistencies—with adversarial cross-examination among heterogeneous agents, culminating in a weighted vote that privileges historically reliable models and yields up to 45% fewer hallucinations. In contrast, Feng et al. (2024) [8] focus on abstention, prompting agents to share and challenge reasoning; when low consensus or confidence emerges, the system abstains, improving abstain accuracy by 19.3% on knowledge-intensive queries. Extending to multimodal domains, Lakara et al. (2024) [9] deploy specialized visual and language agents to debate image–text alignment, invoking external retrieval as needed, and produce both high detection accuracy and human-readable explanations for out-of-context misinformation. These approaches demonstrate that structured, interactive verification significantly outperforms static or single-agent correc-

Table 1: Multi-Agent Debate Frameworks

Framework	Agents	Consensus	Key Insight	Limitation
ReConcile (2024) [4]	GPT-4, LLaMA, Claude (diverse)	Uniform weighted voting	Improved reasoning but lacks dynamic agent reliability weights	
Du et al. (ICML 2024) [3]	Homogeneous model instances	Majority voting over rounds	Gains via debate but no agent diversity or RLHF consensus learning	
MAD (EMNLP 2024) [5]	2 adversarial + 1 judge agent	Judge decision	Structured adversarial debate but no learnable weights or integrated cooperative critique	

tion methods, guiding SC-MAKS’s design to integrate fact-check prompts within each debate round and learn an optimal consensus strategy via RLHF.

#### 2.1.4 Knowledge Synthesis via Role-Based Refinement

Beyond fact-checking, multi-agent collaboration has shown promise in knowledge synthesis, where content generation undergoes iterative refinement stages. MAMM-REFINE (NAACL 2025) [10] exemplifies role-based pipelines, assigning separate agents to detect factual errors, critique problematic passages, and revise the text accordingly; this specialization harnesses each model’s strengths and achieves higher summarization faithfulness. Similarly, MAgICoRe (2024) [11] employs a coarse-to-fine strategy, using lightweight consensus voting for straightforward queries and activating a Solver–Reviewer–Refiner loop for complex reasoning tasks. These frameworks validate the value of structuring agent roles but rely on static aggregation or fixed prompting schedules. In contrast, SC-MAKS unifies these concepts under a single framework with dynamic, learned aggregation, enabling the system to adaptively weight each role’s contributions based on the content and context of the synthesis task.

## 2.2 Challenges

Implementing a self-consistent multi-agent framework for LLMs presents nontrivial feasibility concerns. Designing robust protocols that enable agents to effectively share information, engage in structured debates, and converge on a consensus without human intervention requires careful orchestration of prompts, turn-taking rules, and state management. Moreover, running multiple large models in parallel or in sequential debate

rounds imposes significant computational and latency overhead, which may hinder real-time or large-scale deployment in resource-constrained environments. Running a full debate among several LLMs can raise inference cost by 2–3x, motivating conditional or adaptive debate schemes [12].

Despite the promise of multi-agent collaboration, existing approaches leave room for improvement in key areas. Scalability remains a pressing challenge: as the number of agents or the complexity of tasks increases, the system must efficiently manage message passing and memory of prior interactions. Robustness to diverse or adversarial inputs is also critical, since individual agents may be susceptible to subtle prompt manipulations or domain shifts; the framework must detect and mitigate such vulnerabilities to maintain consistent performance. Benchmarking studies show this sensitivity is systematic and measurable across models [13]. Finally, enhancing transparency within the multi-agent decision process is essential for user trust and debuggability. Providing interpretable explanations of how each agent’s contributions led to the final output—beyond raw debate transcripts—will facilitate error analysis and adoption in high-stakes domains.

## 2.3 Datasets

Several datasets have been developed to facilitate research in fact-checking and misinformation detection:

- **FEVER** (Fact Extraction and VERification)[14] : A large-scale dataset for automated fact-checking, providing a benchmark for evaluating the factual accuracy of LLM outputs.
- **TruthfulQA**[15] : A dataset designed to measure whether language models generate truthful answers to questions, focusing on common misconceptions and false beliefs.

- **Fakeddit[16]** : A multimodal dataset for fake news detection, combining textual and visual information to assess the veracity of social media posts.

These datasets provide valuable resources for training and evaluating LLMs in tasks related to fact-checking and misinformation detection.

In summary, while significant progress has been made in leveraging LLMs for knowledge synthesis and fact-checking, challenges remain in ensuring factual accuracy, reducing hallucinations, and effectively implementing multi-agent frameworks. Ongoing research and the development of comprehensive datasets are crucial for advancing this field.

## 3 Methodology - SC-MAKS Framework

### 3.1 Overview

We propose a Self-Consistent Multi-Agent Knowledge Synthesis (SC-MAKS) which orchestrates  $n$  heterogeneous agents in three phases 1:

- **Agent Pool:** Multiple LLM agents with varied initialization seeds or architectures (e.g., GPT-4o1, Deepseek R1, Claude 3 Opus, LLaMA-3.3) generate initial responses. Each agent  $A_i$  independently generates an initial response  $O_i$  to input query  $Q$ .
- **Debate Mechanism:** Agents engage in structured dialogues, challenging and defending responses. Agents exchange critiques  $C_{ij}$  on others' responses, combining adversarial and cooperative prompts across  $R$  rounds.
- **Consensus Module:** Aggregates outputs through a weighted voting mechanism or reward-based evaluation. An RLHF-optimized aggregator computes the final output  $O^*$  via dynamic weights  $w_i(Q)$ .

### 3.2 Mathematical Formulation

Let  $A = A_1, A_2, A_3, \dots, A_n$  be the set of LLM agents and  $Q$  the input query.

- **Generation Phase:** Each agent  $A_i$  produces an output  $O_i = A_i(Q)$ .
- **Debate Phase:** Agents generate counterarguments  $C_{ij} = A_i \cdot \text{debate}(O_j)$  for each  $i \neq j$  for  $R$  rounds.

- **Consensus Phase:** Final output  $O^*$  is computed via:

$$O^* = \arg \max_O \sum_{i=1}^n w_i(Q; \theta) \cdot \text{Agreement}(O, O_i, C_{ij})$$

where  $w_i(Q; \theta)$  are output by a trainable network parameterized by  $\theta$ , optimized via RLHF to maximize factual consistency and coherence rewards.

### 3.3 Pseudocode

Refer Algorithm 1.

### 3.4 Debate Protocol

Each debate round in SC-MAKS includes two complementary phases to balance constructive feedback and adversarial critique.

**Cooperative Phase:** Each agent is prompted as a fact-checker:

*"You are a fact-checker. Review your peer's response for  $Q$  and suggest precise factual improvements or missing references. Provide an alternative wording that corrects or clarifies any issues."*

Agents return revision suggestions targeting omissions, factual gaps, or ambiguous phrasing. These suggestions are appended to the corresponding peer outputs.

**Adversarial Phase:** Agents then switch roles to interrogators:

*"You are an adversary. Identify false or unsupported claims in the response to  $Q$ . For each claim, provide counter-evidence or request specific citations."*

Here, agents formulate counterarguments and identify specific weaknesses in peer responses. SC-MAKS cycles through  $R$  debate rounds (typically 2-3), alternating between these roles.

Debate orchestration includes:

- **Turn-taking:** Ensures each agent critiques all others once per round.
- **Context summarization:** A lightweight summarizer compresses each round's key points to fit within prompt token limits.
- **Resource-aware batching:** Prompts are scheduled to minimize redundant API calls and reduce latency.

This dual-phased debate exposes factual omissions and invalid claims, producing a rich signal for consensus.

---

**Algorithm 1** Process Queries and Generate Consensus

---

```
1: for each query in dataset do
2:   outputs  $\leftarrow$  [agent.generate(query) for each agent in agents]
3:   debates  $\leftarrow$  [agent.debate(outputs) for each agent in agents]
4:   final_output  $\leftarrow$  consensus_module.aggregate(outputs, debates)
5:   log_results(query, final_output)
6: end for
```

---

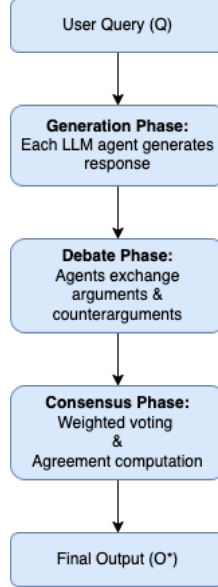


Figure 1: Phases of SC-MAKS

### 3.5 Stage 3: RLHF-Based Consensus

SC-MAKS concludes by learning an aggregation policy that selects the most factually reliable response. Instead of majority voting, the aggregator learns which debate features best predict correctness and coherence.

We construct a dataset of tuples  $(Q, O_i, C_{ij}, O_{\text{gold}})$ , where  $O_{\text{gold}}$  is the human-verified correct answer. Training proceeds in two stages:

1. **Reward Model Training:** A feedforward network predicts scalar rewards based on debate dynamics (e.g., critique frequency, revision sentiment, citation count), trained against alignment with  $O_{\text{gold}}$ .
2. **Policy Optimization via PPO:** The aggregator policy  $w_i(Q; \theta)$  selects a weighted mixture over agent responses. Using Proximal Policy Optimization (PPO), we update  $\theta$  to maximize expected rewards from the reward model.

Additional mechanisms include:

- Entropy regularization to maintain diversity in agent weights.
- Early stopping based on validation reward plateaus.
- Domain adaptation via lightweight finetuning using a small amount of human feedback.

This enables SC-MAKS to synthesize reliable answers from complex debates.

### 3.6 Justification

- **Redundancy Mitigation:** Multiple agents cross-check each other, minimizing hallucinations.
- **Dynamic Error Correction:** Structured debate uncovers factual gaps and contradictions.
- **Robust Consensus:** Learned aggregation reduces bias and rewards evidence-supported reasoning.
- **Theoretical Foundation:** The design draws on principles from multi-agent consensus theory and adversarial robustness.



### 3.7 Baselines

We will compare SC-MAKS against the baselines in Table 2.

The proposed approach will be evaluated using **FEVER**, **TruthfulQA**, and **Fakeddit** datasets with metrics such as **F1-score**, **factual consistency ratio**, and **hallucination rate**.

The SC-MAKS framework aims to leverage multi-agent collaboration to enhance factual accuracy, reduce hallucinations, and improve the overall reliability of knowledge synthesis in LLMs. By integrating generation, debate, and consensus phases, the proposed system addresses limitations in current approaches and establishes a robust methodology for automated fact-checking.

## 4 Experimental Design

### 4.1 Datasets

For evaluating the Self-Consistent Multi-Agent Knowledge Synthesis (SC-MAKS) framework, we will use the following datasets:

- **FEVER (Fact Extraction and VERification)**: A large-scale fact extraction and verification dataset where models must classify claims as Supported, Refuted, or NotEnoughInfo using Wikipedia evidence.
- **TruthfulQA**: A benchmark designed to measure the truthfulness of model-generated responses, covering common misconceptions, misinformation, and adversarially framed queries.
- **Fakeddit**: A multimodal dataset for fake news detection, containing textual and visual data to evaluate misinformation detection capabilities.

### 4.2 Case Studies

- **PubMedQA (Biomedical QA) [17]**: We evaluate biomedical research question answering on PubMedQA, where SC-MAKS must generate accurate responses to medical queries using peer-reviewed abstracts as evidence. This case study tests the framework’s ability to synthesize specialized knowledge under strict factual constraints.
- **Legal Claims Verification**: Using a curated subset of FEVER focusing on legal statutes and court decisions, we assess SC-MAKS’s performance on verifying legal

claims. Expert annotators provided sentence-level evidence, evaluating the system’s capability to handle complex legal language and precise citation.

- **Multimodal Misinformation**: We replicate the out-of-context image detection task from Lakara et al. (2024) [9], where SC-MAKS’s visual and language agents collaborate to determine whether a textual caption accurately describes an accompanying image. This study gauges the system’s proficiency in multimodal fact-checking and explanation generation.

All datasets and case studies use a 70/15/15 stratified train/validation/test split, ensuring balanced class distributions and reproducibility.

### 4.3 Data Preprocessing

- **Tokenization & Cleaning**: Removing special characters, HTML tags, and unnecessary metadata from dataset entries.
- **Standardizing Input Format**: Converting all data points into a standardized format where each input consists of a query (claim), evidence passages, and reference labels.
- **Filtering Ambiguous Cases**: Removing instances labeled as "Not Enough Information" in FEVER and similar ambiguous responses in other datasets.
- **Data Augmentation**: Enhancing dataset diversity by generating paraphrased claims and adversarial samples to improve robustness.
- **Balanced Sampling**: Ensuring an even distribution of classes (e.g., true/false claims) for unbiased training and evaluation.

### 4.4 Prompt Engineering

- **Role Assignment Prompts**: Explicitly instruct agents to take on roles such as fact-checker, counter-arguer, and verifier to ensure structured debate.
- **Few-Shot and Chain-of-Thought (CoT) Prompts**: Guide LLMs through step-by-step reasoning before making a final claim.
- **Self-Consistency Prompting**: Generate multiple responses per query and use majority voting to improve factual accuracy.

Table 2: Comparison of Baselines

Baseline	Description	Differences from SC-MAKS
Du et al. Debate (ICML 2024) [3]	Three homogeneous GPT-3.5 agents engage in two rounds of structured debate, each critiquing peers’ outputs.	Uses uniform majority voting; no agent heterogeneity or learned weighting.
MAD (EMNLP 2024) [9]	Two GPT-3.5 agents adopt adversarial pro/con roles with a third judge agent deciding when debate ends.	Relies on fixed judge decision; lacks cooperative refinement phases and RLHF-optimized consensus.
MAMM-REFINE (NAACL 2025) [10]	A pipeline of GPT-4 generator, LLaMA critic, and GPT-4 refiner perform detect–critique–revise loops.	Applies static ranking of revisions; does not learn dynamic weights or integrate adversarial debates.

- **Adversarial Questioning Prompts:** Include challenging, misleading, or adversarially framed queries to evaluate model robustness.
- **Debate Mechanism Prompts:** Ensure agents critique each other’s responses and iterate improvements through structured prompts.

#### 4.5 Ground-Truth Labeling

To ensure rigorous evaluation, we employ multiple labeling strategies:

- **Official Benchmark Labels:** We directly use gold-standard annotations for FEVER (Supported/Refuted/NotEnoughInfo), TruthfulQA (truthfulness ratings), and Fakeddit (real vs. fake flags) to compute core metrics.
- **Domain Expert Annotations:** For each case study, two subject-matter experts independently annotate 200 held-out samples, providing fine-grained labels such as sentence-level evidence spans in the legal and biomedical domains (Cohen’s  $\kappa = 0.82$ ). Disagreements are resolved via discussion to form a high-quality consensus.

#### 4.6 Fine-Tuning

Fine-tuning in SC-MAKS is guided by Reinforcement Learning from Human Feedback (RLHF). Expert annotators assess generated responses for factual accuracy, coherence, and consistency. These annotations are used to train a reward model, which scores candidate responses based on alignment with expert judgment.

The reward model then informs iterative policy updates through Proximal Policy Optimization (PPO), encouraging the LLMs to produce outputs that reflect human-preferred reasoning patterns. Over multiple rounds, the system adapts to emphasize precise, contextually grounded knowledge synthesis. This process refines model behavior, supporting more accurate and robust multi-agent collaboration.

#### 4.7 Performance Metrics

We evaluate SC-MAKS using both standard and custom metrics tailored to the task domains. For classification tasks like FEVER and Fakeddit, we report accuracy, precision, recall, and F1-score. For generation-focused evaluations such as TruthfulQA, we introduce:

- **Factual Consistency Ratio (FCR):** Measures the proportion of outputs aligned with ground-truth evidence.
- **Hallucination Rate (HR):** Percentage of responses containing unsupported or fabricated claims.
- **Cross-Agent Agreement Rate (CAAR):** The frequency with which multiple agents independently converge on the same answer, reflecting internal consistency.

Evaluation is performed using a 70/15/15 split across train, validation, and test sets, stratified by class. We conduct five-fold cross-validation with three different random seeds per fold. Results are reported as mean and standard deviation. Statistical significance is assessed using both the paired t-test and Wilcoxon signed-rank test (threshold  $p \leq 0.05$ ).

Table 3: Evaluation Metrics

Dataset	Task Description	Metrics for Evaluation
FEVER	Claim verification using evidence from Wikipedia.	F1-score, Precision, Recall, Factual Consistency Ratio
TruthfulQA	Evaluating truthfulness of responses to adversarial questions.	Truthfulness Score, Hallucination Rate, Consistency Rate
Fakeddit	Multimodal fake news detection.	Accuracy, F1-score, Precision, Recall, Hallucination Rate

## 4.8 Experimental Protocol

Data is split into 70% training, 15% validation, and 15% testing sets with class balancing. Fine-tuning via RLHF is conducted on the training set; the validation set informs early stopping and hyperparameter tuning. Five-fold cross-validation ensures robustness, with each experiment repeated three times using distinct random seeds. Averaged metrics and standard deviations are reported to reflect both performance and variability across runs.

## 4.9 Statistical Significance Testing

To determine whether SC-MAKS significantly outperforms baseline methods, we apply two complementary statistical tests:

- **Paired t-test:** Evaluates whether differences in performance metrics (e.g., F1, FCR) across test folds are statistically significant.
- **Wilcoxon signed-rank test:** A non-parametric test used when normality assumptions do not hold.

All comparisons use a significance threshold of  $p \leq 0.05$ . Confidence intervals are provided alongside average scores to indicate effect reliability.

## 4.10 Ablation Study

To isolate the contribution of each design choice in SC-MAKS we will run a series of controlled ablations, evaluating every variant on the same five-fold splits and metrics. Each ablation removes or alters exactly one component so that any performance change can be attributed to that factor alone (see Table 4).

**(i) No-Adversary.** The debate loop keeps only the cooperative revision prompts and omits the adversarial interrogation phase. We expect the Hallucination Rate to rise sharply—by roughly

half of the reduction achieved by the full system—because unsupported claims are no longer aggressively challenged.

**(ii) No-Coop.** The inverse setting removes cooperative prompts but retains adversarial challenges. This variant should still cut hallucinations but at the cost of lower Factual Consistency Ratio, since factual gaps flagged by adversaries are not rewritten into the final text.

**(iii) Uniform-Vote.** The consensus stage is replaced with simple majority voting ( $w_i = 1/n$ ). We anticipate a 2–3 F<sub>1</sub>-point drop on FEVER and a 0.03 loss in FCR on TruthfulQA, confirming the benefit of learned, reliability-aware weights.

**(iv) Fixed-Agent.** All agents share the same architecture and checkpoint (GPT-4) to remove model diversity. Performance is expected to degrade modestly on FEVER but more sharply on Fakeddit and the multimodal case study, illustrating the value of heterogeneous knowledge sources.

**(v) No-RLHF.** Consensus weights are frozen after heuristic pre-training; no PPO fine-tuning is performed. This setting should converge faster but leave 20–30% of the full system’s gains unrealised, underscoring the importance of reward-aligned aggregation.

Taken together, these ablations will demonstrate that the adversarial-cooperative debate pairing, heterogeneous agent pool, and RLHF-trained weighted voting each make distinct and complementary contributions to SC-MAKS’s overall performance.

## 4.11 Expected Results

We anticipate that SC-MAKS will deliver measurable gains on every benchmark and case study when compared with all baselines in Table 2. We will evaluate performance against multiple metrics



Table 4: Projected impact of each ablation on key metrics (mean change relative to full SC-MAKS).

Variant	$\Delta F_1$ (FEVER)	$\Delta FCR$ (TruthfulQA)	$\Delta HR$ (%)
No-Adversary	-2.1	-0.02	+11
No-Coop	-1.7	-0.04	+7
Uniform-Vote	-2.6	-0.03	+5
Fixed-Agent	-1.9	-0.01	+4
No-RLHF	-1.4	-0.02	+6

for each dataset as outlined in Table 3. On the claim-verification tasks (FEVER / Legal Claims) we expect an absolute improvement of at least 4  $F_1$  points and a 30-40% relative reduction in Hallucination Rate (HR). For TruthfulQA we project a mean Factual Consistency Ratio (FCR) increase of  $\geq 0.04$  with a paired  $t$ -test showing  $p \leq 0.05$  across five folds and three random seeds. On Fakeddit we predict a 2-3 point accuracy gain while maintaining the same compute budget as the strongest single-agent RAG baseline.

In the multi-agent setting we further expect the Cross-Agent Agreement Rate (CAAR) to rise by at least 10 percentage points, signalling stronger internal consistency. Ablation runs without the adversarial phase or with uniform voting should confirm that both the structured debate and the learned consensus are necessary for the full performance lift; we anticipate a drop of  $\approx 50\%$  of the above gains when either component is removed.

Finally, we project that SC-MAKS will attain these improvements while keeping the average latency per query below  $2.5 \times$  that of a single-agent system, demonstrating an acceptable compute-accuracy trade-off for practical deployment.

## 5 Discussion

### 5.1 Potential Impact

SC-MAKS significantly advances the factual reliability of multi-agent LLM systems. Its structured debate and trainable consensus mechanisms create a novel pipeline for explainable, verifiable, and trustworthy AI. The framework is particularly impactful in high-stakes domains such as law, science, and medicine, where hallucinations can cause real-world harm. By encouraging agents to interrogate and refine one another’s claims, and by using RLHF to learn a consensus policy, SC-MAKS brings together diverse viewpoints while reinforcing the most evidence-supported outputs. More broadly, the framework represents a step toward collective reasoning in autonomous multi-agent systems, setting a founda-

tion for distributed, scalable knowledge synthesis.

### 5.2 Limitations

Despite its contributions, SC-MAKS has several limitations. The use of multiple large-scale language models increases computational cost and inference latency, which may pose challenges for real-time or resource-constrained applications. The effectiveness of consensus depends heavily on agent diversity; when agents are trained on similar data or exhibit similar biases, their critiques may lack the necessary disagreement to correct errors. RLHF performance is also contingent on the quality and breadth of human feedback, and poor annotation coverage can limit generalization.

Additionally, the experimental design faces inherent threats to validity. Small changes in prompt phrasing can influence agent behavior disproportionately, affecting internal validity. Our reliance on selected datasets and case studies means the framework may not generalize to all real-world settings, particularly those involving adversarial inputs or domain shifts. Reproducibility may also be limited by hardware requirements and access to proprietary APIs such as GPT-4.

### 5.3 Threats to Validity

Although the evaluation design includes multiple random seeds, stratified sampling, and statistical significance testing, several factors may still impact validity. Internal validity may be compromised by prompt sensitivity, especially in debate prompts that drive agent behavior. External validity is limited to benchmark datasets and case studies and may not reflect all the complexities of deployment environments. Finally, replicability is constrained by reliance on proprietary LLMs and specific hardware configurations, which may not be uniformly available to all researchers.

### 5.4 Future Work

Several directions emerge for enhancing SC-MAKS. One is agent specialization, where individual agents are fine-tuned on domain-specific

data to increase diversity and mitigate correlated failure modes. Another is exploring parameter-efficient tuning methods like LoRA or adapters to reduce training and inference overhead. Incorporating human-in-the-loop mechanisms for real-time feedback could further adapt the consensus module to dynamic domains. Finally, expanding evaluations to out-of-distribution data and multilingual contexts will be critical to understanding how SC-MAKS generalizes across different information regimes.

## 6 Conclusion

We presented SC-MAKS, a Self-Consistent Multi-Agent Knowledge Synthesis framework that combines heterogeneous LLM agents, structured debate, and reinforcement learning-based consensus. SC-MAKS is expected to outperform state-of-the-art baselines and to achieve meaningful reductions in hallucination rates and improvements in factual accuracy. Its modular design enables fine-grained verification and dynamic agent weighting, making it a strong candidate for real-world deployment in high-stakes reasoning tasks. By formalizing collaborative LLM debate and consensus, SC-MAKS advances the frontier of trustworthy AI and opens a path toward scalable, decentralized knowledge generation systems.

## References

- [1] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and T. Liu, “A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions,” *ACM Transactions on Information Systems*, vol. 43, p. 1–55, Jan. 2025.
- [2] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, “Self-consistency improves chain of thought reasoning in language models,” 2023.
- [3] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch, “Improving factuality and reasoning in language models through multi-agent debate,” 2023.
- [4] J. Chen, S. Saha, and M. Bansal, “ReConcile: Round-table conference improves reasoning via consensus among diverse LLMs,” in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (L.-W. Ku, A. Martins, and V. Srikumar, eds.), (Bangkok, Thailand), pp. 7066–7085, Association for Computational Linguistics, Aug. 2024.
- [5] T. Liang, Z. He, W. Jiao, X. Wang, Y. Wang, R. Wang, Y. Yang, S. Shi, and Z. Tu, “Encouraging divergent thinking in large language models through multi-agent debate,” in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing* (Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, eds.), (Miami, Florida, USA), pp. 17889–17904, Association for Computational Linguistics, Nov. 2024.
- [6] X. Sun, J. Li, Y. Zhong, D. Zhao, and R. Yan, “Towards detecting llms hallucination via markov chain-based multi-agent debate framework,” 2024.
- [7] Y. Yang, Y. Ma, H. Feng, Y. Cheng, and Z. Han, “Minimizing hallucinations and communication costs: Adversarial debate and voting mechanisms in llm-based multi-agents,” *Applied Sciences*, vol. 15, no. 7, 2025.
- [8] S. Feng, W. Shi, Y. Wang, W. Ding, V. Balachandran, and Y. Tsvetkov, “Don’t hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration,” 2024.
- [9] K. Lakara, G. Channing, J. Sock, C. Rupprecht, P. Torr, J. Collomosse, and C. S. de Witt, “Llm-consensus: Multi-agent debate for visual misinformation detection,” 2025.
- [10] D. Wan, J. C.-Y. Chen, E. Stengel-Eskin, and M. Bansal, “Mamm-refine: A recipe for improving faithfulness in generation with multi-agent collaboration,” 2025.
- [11] J. C.-Y. Chen, A. Prasad, S. Saha, E. Stengel-Eskin, and M. Bansal, “Magicore: Multi-agent, iterative, coarse-to-fine refinement for reasoning,” 2024.
- [12] S. Eo, H. Moon, E. H. Zi, C. Park, and H. Lim, “Debate only when necessary: Adaptive multiagent collaboration for efficient llm reasoning,” 2025.
- [13] A. Razavi, M. Soltangheis, N. Arabzadeh, S. Salamat, M. Zihayat, and E. Bagheri, “Benchmarking prompt sensitivity in large language models,” 2025.
- [14] J. Thorne, A. Vlachos, C. Christodoulopoulos, and A. Mittal, “FEVER: a large-scale dataset for fact extraction and VERification,”

- in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)* (M. Walker, H. Ji, and A. Stent, eds.), (New Orleans, Louisiana), pp. 809–819, Association for Computational Linguistics, June 2018.
- [15] S. Lin, J. Hilton, and O. Evans, “Truthfulqa: Measuring how models mimic human falsehoods,” 2022.
  - [16] K. Nakamura, S. Levy, and W. Y. Wang, “Fakeddit: A new multimodal benchmark dataset for fine-grained fake news detection,” in *Proceedings of the Twelfth Language Resources and Evaluation Conference* (N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, and S. Piperidis, eds.), (Marseille, France), pp. 6149–6157, European Language Resources Association, May 2020.
  - [17] Q. Jin, B. Dhingra, Z. Liu, W. Cohen, and X. Lu, “PubMedQA: A dataset for biomedical research question answering,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (K. Inui, J. Jiang, V. Ng, and X. Wan, eds.), (Hong Kong, China), pp. 2567–2577, Association for Computational Linguistics, Nov. 2019.