

Generalization of Diverse Tool-Augmented Multi-Agent Debate Frameworks to Out-of-Domain Fact Verification

Assignee Research

June 11, 2026

Abstract

Large Language Models (LLMs) suffer from hallucinations and factual inaccuracies, especially in complex reasoning and fact verification tasks. Multi-Agent Debate (MAD) systems aim to improve answer accuracy by enabling multiple LLM agents to engage in dialogue, promoting diverse reasoning and mutual verification. However, existing MAD frameworks primarily rely on internal knowledge or static documents, making them vulnerable to hallucinations. While MADKE introduces external evidence to mitigate this, its one-time retrieval mechanism limits adaptability to new arguments or emerging information

1 Introduction

This paper examines: Tool-MAD: A Multi-Agent Debate Framework for Fact Verification with Diverse Tool Augmentation and Adaptive Retrieval. Research question: How does the performance of diverse tool-augmented multi-agent debate frameworks generalize to out-of-domain fact verification tasks compared to single-agent baselines on standard NLP benchmarks?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.1/10.

3 Results

15 papers retrieved. 19 claims extracted; 14 independently verified. Quality review score: 7.1/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
FactScore attempts to quantify factuality at the claim level by decomposing model outputs.	✓	0.18
Consistency-based factuality approaches, including self-consistency and re-asking, evaluate factual correctness by measu	✓	0.27
RAGAS introduces two complementary metrics: faithfulness, measuring whether claims match retrieved evidence, and answer	✓	0.22
Prior metrics like faithfulness and answer relevance have predominantly been applied as post-hoc evaluators for fully ge	✓	0.22
No prior work integrates metric-guided evaluation into a multi-agent debate structure.	✓	0.15
No prior work uses factuality signals to modulate argument selection, judge decisions, or evidence refinement across rou	✓	0.22
Tool-MAD incorporates faithfulness and answer relevance as round-level stability indicators internally.	✓	0.15
Toolformer enables models to learn API-calling behaviors.	×	0.12
HuggingGPT and GEAR treat the LLM as a coordinator for heterogeneous models or systems.	×	0.14
Domain-specific tools in agents like ChemCrow and biomedical retrieval agents can dramatically improve reasoning fidelit	✓	0.19
Existing tool-augmented systems overwhelmingly adopt a single-agent perspective and use tools in a one-shot or sequentia	✓	0.27
Existing retrieval systems typically perform a single retrieval step at the beginning of a task without adapting to evol	✓	0.25
Tool-MAD uses faithfulness and answer relevance collectively as a stability score to serve as an auxiliary signal for th	✓	0.17
Tool-MAD was evaluated across four fact verification benchmark datasets.	×	0.12
Tool-MAD outperforms the MAD framework with performance improvements of up to 35.5%.	×	0.13
Tool-MAD outperforms the MADKE framework with performance improvements of up to 5.5%.	×	0.11
Tool-MAD maintains robust performance under different retrieval tools and corpus configurations in medical QA settings.	✓	0.22
Tool-MAD leverages a diverse set of external tools, including real-time search APIs and Retrieval-Augmented Generation (✓	0.22
Tool-MAD introduces an adaptive query formulation mechanism enabling agents to iteratively refine evidence retrieval bas	✓	0.23

References

- <http://arxiv.org/abs/2408.04114v1>
- <http://arxiv.org/abs/2601.19151v1>
- <http://arxiv.org/abs/2601.04742v1>