

---

# MAD-Sherlock: Multi-Agent Debate for Visual Misinformation Detection

---

Kumud Lakara<sup>\*1</sup> Georgia Channing<sup>\*1</sup> Christian Rupprecht<sup>1</sup> Juil Sock<sup>2</sup> Philip Torr<sup>1</sup> John Collomosse<sup>3</sup>  
Christian Schroeder de Witt<sup>1</sup>

## Abstract

One of the most challenging forms of misinformation involves pairing images with misleading text to create false narratives. Existing AI-driven detection systems often require domain-specific finetuning, limiting generalizability, and offer little insight into their decisions, hindering trust and adoption. We introduce MAD-Sherlock, a multi-agent debate system for out-of-context misinformation detection. MAD-Sherlock frames detection as a multi-agent debate, reflecting the diverse and conflicting discourse found online. Multimodal agents collaborate to assess contextual consistency and retrieve external information to support cross-context reasoning. Our framework is domain- and time-agnostic—requiring no finetuning—yet achieves state-of-the-art accuracy with in-depth explanations. Evaluated on NewsCLIPPings, VERITE, and MMFakeBench, it outperforms prior methods by 2%, 3%, and 5%, respectively. Ablation and user studies show that the debate and resultant explanations significantly improve detection performance and improves trust for both experts and non-experts, positioning MAD-Sherlock as a robust tool for autonomous citizen intelligence.

## 1. Introduction

The rise of online news and social media has been paralleled by a surge in digital misinformation (Aslett et al., 2024; Hasher et al., 1977; Brashier and Marsh, 2020). Among the most widespread tactics is out-of-context (OOC) image use (Fazio, 2020), where unaltered images are paired with misleading text to deceive, requiring little technical skill.

---

<sup>\*</sup>Equal contribution <sup>1</sup>University of Oxford, Oxford, UK  
<sup>2</sup>BBC AI Research, London, UK <sup>3</sup>University of Surrey, Guildford, UK. Correspondence to: Georgia Channing <cgeorgia@robots.ox.ac.uk>, Christian Schroeder de Witt <cs@robots.ox.ac.uk>.

Accepted at the ICML 2025 Workshop on Collaborative and Federated Agentic Workflows (CFAgentic@ICML'25), Vancouver, Canada. July 19, 2025. Copyright 2025 by the author(s).

OOC detection demands nuanced reasoning to identify misalignment between images and accompanying text. This task is time-consuming for humans, and detection accuracy drops under time pressure (Sultan et al., 2022), limiting scalability.

AI tools offer a path forward, but traditional forensic methods (Castillo Camacho and Wang, 2021; Heidari et al., 2024; Zhu et al., 2018; Amerini et al., 2021; Hina et al., 2021) focus on tampering artifacts (e.g., Photoshop edits (Farid, 2016; Wang et al., 2019) or Deepfakes (Tolosana et al., 2020; Masood et al., 2023; Somers, 2020)) and are ill-suited for OOC detection, which hinges on cross-contextual reasoning rather than visual manipulation.

Pretrained large multimodal models (LMMs) offer a promising foundation for detecting OOC image use by jointly interpreting text and images (Liu et al., 2024b; OpenAI and et al., 2024; Li et al., 2019; Radford et al., 2021). However, applying them to news content is challenging. News articles often feature loosely related but contextually appropriate images—e.g., a pre-2024 photo of Donald Trump used in election coverage—which makes OOC detection difficult using pretraining alone. LMMs also suffer from hallucinations, misinterpret intent (Bai et al., 2024; Liu et al., 2024a), and lack up-to-date knowledge. Off-the-shelf models exhibit these limitations, reducing their reliability. Fine-tuning helps (Qi et al., 2024), but is costly and requires frequent updates. Crucially, beyond detection, models must provide clear, human-readable explanations to support understanding and trust in their decisions.

In this work, we propose a novel post-training framework for scalable OOC misinformation detection that improves contextual reasoning, offers built-in explainability, and achieves state-of-the-art accuracy without task-specific fine-tuning (see Section 3). Our method frames the detection problem as a dialectic debate between LMM agents, augmented with external information retrieval.

Unlike single-agent chain-of-thought methods (Wei et al., 2024), our multi-agent setup enables context separation, decentralized reasoning, and parallelization (Schroeder de Witt et al., 2020; Du et al., 2023). Prior work (Miresghalah et al., 2024) shows agents struggle to maintain diverse perspectives within a single context window; our approach

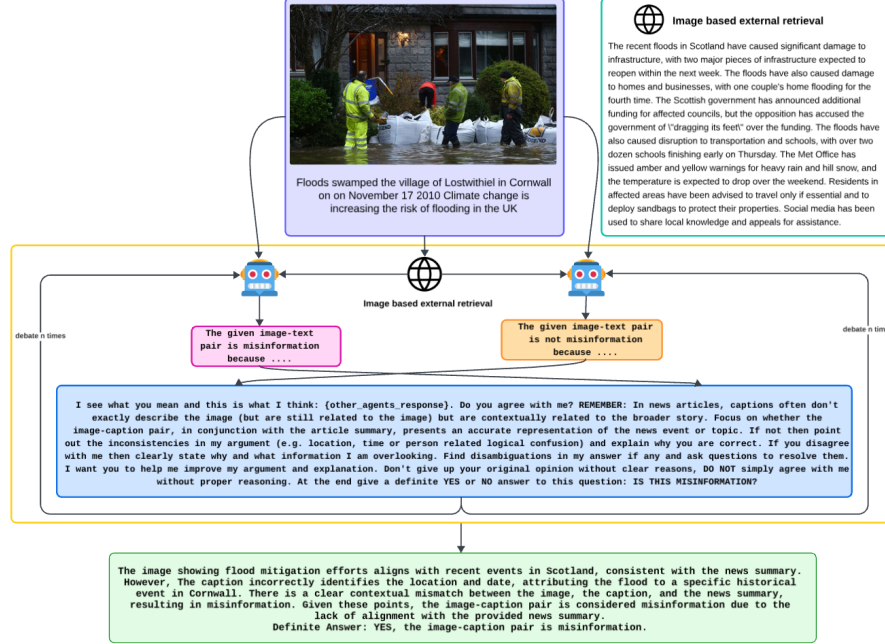


Figure 1. **Overview of MAD-Sherlock:** Two or more independent agents see the same image-text input and are tasked with detecting whether the input is misinformation or not. After the agents form their independent opinions, they participate in a debate until they converge on the same response or when  $n$  debate rounds are completed (whichever is earlier).

addresses this via modular, compositional interaction. It also allows seamless integration of human or autonomous agents, making MAD-Sherlock a flexible tool for expert oversight. To our knowledge, this is the first use of multi-agent LMM debate for both detecting and *explaining* OOC image use.

MAD-Sherlock is backbone-agnostic, compatible with open- and closed-source models. We prototype and tune with LLaVA (Liu et al., 2024b) to reduce API costs, then deploy GPT-4o (OpenAI) for final results. By avoiding task-specific fine-tuning, our method remains broadly applicable across domains and time periods.

We evaluate on three benchmarks—NewsCLIPPings, VERITE, and MMFakeBench—and achieve new state-of-the-art performance across all. MAD-Sherlock outperforms prior methods and baselines, demonstrates robustness, and produces detailed, interpretable explanations. In user studies, these explanations significantly improve detection accuracy for both experts and non-experts.

Key contributors to performance include external retrieval and agent independence. We discuss limitations and outline future directions toward scalable, general-purpose AI for public good.

## 2. Related Work

Recent work has explored joint image-text representations for OOC classification. Aneja et al. (2022) use a self-supervised approach to enforce image-text alignment via an object-caption matching score, classifying OOC instances by caption similarity. However, the method is heavily text-dependent and lacks interpretability.

Abdelnabi et al. (2022) introduce the Consistency Checking Network (CCN), which leverages memory networks and CLIP for image-caption consistency against external evidence, improving accuracy but offering only binary outputs without explanations.

Zhang et al. (2024) employ AMR-based symbolic graphs for interpretable OOC detection with evidence. Similarly, Zhou et al. (2020) propose SAFE for joint text-visual reasoning in fake news detection. Wang et al. (2018) present EANN, which uses adversarial training to extract event-invariant multimodal features. These methods require training from scratch and lack the advanced reasoning and knowledge of large pretrained models.

Sniffer (Qi et al., 2024) is most similar to our work. It uses InstructBLIP (Dai et al., 2023) to detect OOC image use and generate explanations by aggregating internal and external knowledge from entity extraction APIs and image-based

web searches. However, its external knowledge is limited to basic textual content (e.g., article titles), and adapting it to the news domain requires extensive training, reducing generalizability and increasing computational overhead.

The CRAVE framework (Dey et al., 2025) clusters retrieved multimodal evidence into coherent narratives, then uses an LLM to produce fact-checking judgments with natural language explanations. While effective, treating all clusters equally may amplify fringe perspectives and distort the evidence base.

Liu et al. (2025) introduce MMD-Agent, a single-agent framework that hierarchically decomposes misinformation detection into textual, visual, and cross-modal subtasks before reasoning over their outputs to make a final prediction.

### 3. Methodology

We present MAD-Sherlock, an explainable misinformation detection system that jointly predicts and explains instances of misinformation (Figure 1). Unlike prior work, which largely provides predictions without explanations, our approach uses multiple multi-modal models debating to determine whether an image-text pair constitutes misinformation. We address the question:

*Can debating multi-modal models, equipped with external context, detect misinformation by identifying subtle contextual inconsistencies?*

Our external retrieval module uses reverse image search to provide agents with real-world context, enhancing their predictions. Using the GPT-4o (OpenAI) model, we achieve state-of-the-art performance with detailed, coherent explanations, without requiring domain-specific fine-tuning. This ensures faster generalization to new domains with minimal computational overhead.

#### 3.1. Debate Modelling

Analogous to real-world conversations, communication between two AI agents can also be structured in a myriad of ways. We explore multiple debating strategies to structure the conversation between agents, all of which are tested and evaluated in our experiments and informed by the work of Brown-Cohen et al. (2023).

**Asynchronous Debate (not) Against Human:** We experiment with an asynchronous debating strategy, where models wait for others’ responses before generating their own (Figure 2a vs. b). While synchronous debates are faster and more efficient, the asynchronous setup proves more effective for identifying contextual ambiguities—critical in misinformation detection. Notably, models are prompted to believe they are debating a human rather than another AI.

**Judged Debate:** We also explore an asynchronous debate setup with a judge (Figure 2c), where models debate as usual, but a final decision is made by a judge based solely on the debate transcript. Following Khan et al. (2024), the judge lacks access to external information, incentivizing models to structure arguments that are maximally persuasive.

**Actor-Skeptic:** In this setup, a single *actor* determines whether an image-text pair constitutes misinformation. A *skeptic* then critiques the actor’s reasoning, probing for logical flaws and ambiguities. Neither agent has access to the ground truth, and since their roles are distinct, this configuration does not benefit from ensembling.

#### 3.2. Prompt Engineering

The debate is structured via prompt engineering (Figure 1). In the first stage, each agent independently assesses whether the image-text pair is misinformation, using context from an external retrieval module (see Appendix A.3). Prompts summarize relevant articles and emphasize visual cues (e.g., watermarks, flags) to guide initial judgments. Agents then debate: the first round uses a tailored prompt to initiate discussion; later rounds use a standard prompt incorporating prior responses. Agents must agree or disagree with peers while refining reasoning, and prompts explicitly discourage blind agreement by requiring justification for any alignment.

#### 3.3. External Information Retrieval

A model’s world knowledge is bounded by its training data and time frame, but external retrieval enables access to up-to-date context. Prior work leverages retrieval-based datasets (Abdelnabi et al., 2022), though these typically include only news headlines, which are often too sparse for reliable inference. Full article access can significantly improve misinformation detection. To this end, we introduce a two-stage external retrieval module that enhances accuracy when integrated into the pipeline:

**API-Based Information Retrieval** We use the Bing Visual Search API to retrieve web pages related to each image. For each image, we collect the top three matching pages in which it appears, assuming these provide sufficient context for understanding the image’s original use. Since NewsCLippings (Luo et al., 2021) includes articles over a decade old, some images yield no search results. In these rare cases, we omit external context and rely solely on the agent’s prior knowledge, which has minimal impact on overall performance.

**Summarization using LLM** After retrieving the top three web pages, we scrape their text and use LLaMA-13B (Touvron et al., 2023) to generate concise summaries focused on

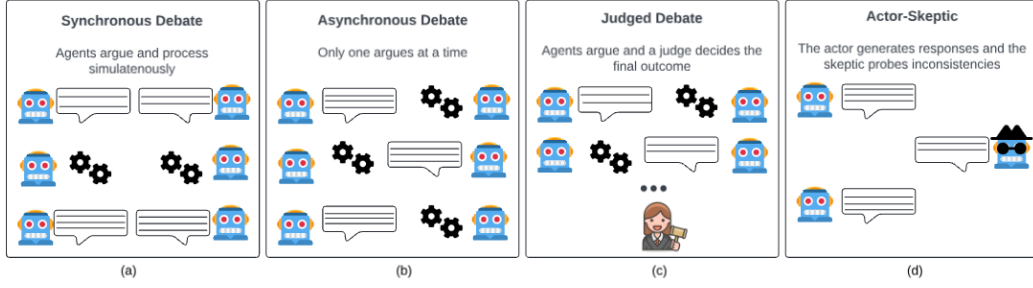


Figure 2. **Debating Strategies:** We experiment with multiple debating strategies to evaluate which performs best on our task.

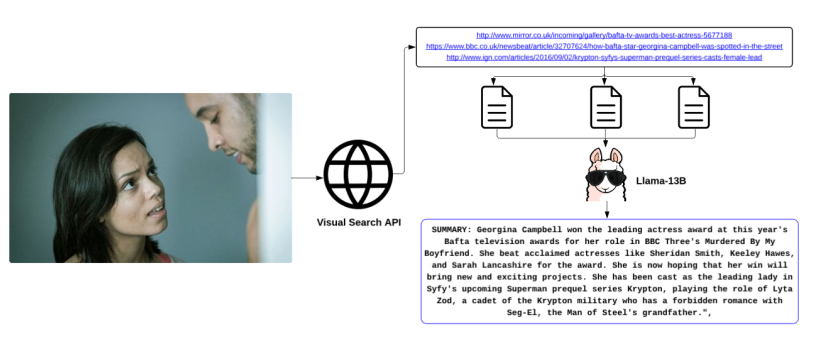


Figure 3. **Structure of the external information retrieval module:** We use the Bing Visual Search API to obtain web pages related to a given image, which are then summarised using Llama-13B (Touvron et al., 2023). This summary is then passed to the debating agents as a part of the initial prompt.

key contextual details. Since the model struggles with non-English content, we filter out non-English pages—limiting the system to English but without affecting performance, as the dataset is predominantly English. Multilingual support could be added via translation prior to summarization.

**Coherent Reasoning** This stage of the pipeline integrates all components of MAD-Sherlock. Each multimodal agent engages in the optimal debate setup using relevant prompts and is tasked with determining and explaining whether a given image-text pair constitutes misinformation. Agents also access contextual information via the external retrieval module. The system’s final decision is produced when the debate concludes, which is either after a fixed number of rounds or once all agents converge on a shared response, whichever occurs first.

## 4. Experiments and Results

### 4.1. Dataset

We conduct experiments on three datasets: NewsCLippings, VERITE, and MMFakeBench (Luo et al., 2021; Papadopoulos et al., 2023b; Liu et al., 2025). This selection balances well-established benchmarks like NewsCLippings with newer, more diverse datasets such as VERITE and

MMFakeBench, providing both continuity with prior work and coverage of recent advances.

The NewsCLippings dataset builds on VisualNews (Liu et al., 2020), which contains image-caption pairs from BBC, USA Today, The Guardian, and The Washington Post. OOC samples are created by replacing an image with a semantically related one from a different pair. We use the Merged-Balanced version, with 71,072 training, 7,024 validation, and 7,264 test samples, offering balanced retrieval strategies and label distributions. Its scale and prior adoption make it well-suited for evaluating our model’s ability to detect out-of-context content.

We also evaluate on the VERITE dataset (Papadopoulos et al., 2023b), a recent benchmark for multimodal misinformation detection that controls for unimodal bias. Each image-caption pair is constructed so that neither modality alone determines veracity. VERITE includes three subsets—All, True vs. OOC, and True vs. MC—each testing different aspects of multimodal reasoning. Captions and images appear in both truthful and misleading contexts, promoting balanced, joint reasoning. The test split contains 338 true, 324 OOC, and 338 MC samples. We report results on the combined misinformation set (MC + OOC), with disaggregated results in Appendix A.4.



We also evaluate on MMFakeBench (Liu et al., 2025), a benchmark for multimodal misinformation detection in mixed-source scenarios. The validation set includes 1,000 image-caption pairs spanning three major categories—textual, visual, and cross-modal distortions—each with 12 subtypes, capturing the complexity of real-world misinformation.

## 4.2. Experimental Setup

All experiments were run on 8 A40 (46GB) Nvidia GPU server. The estimated cost of processing one data sample using MAD-Sherlock with a GPT-4o backbone is \$0.24. Inference times range from 5 to 15 seconds.

**Debate Setup** We conduct experiments to select the best debating configuration using the LLaVA model (Liu et al., 2024b). All experiments are run for  $k = 3$  rounds or until the agents converge (whichever is earlier).

**External Retrieval Module** We use the Bing Visual Search API to run an image-based reverse search. Using the API we select the top  $k = 3$  pages in which the image appears and scrape the text from them using the Newspaper3k library. Finally, we use Llama-13B (Touvron et al., 2023) to summarise the text obtained from the top  $k = 3$  web pages. This step is crucial since the web pages are usually news articles which contain large amounts of text which, when scraped and passed directly to the model, can exceed its maximum token length.

**Baselines and Prior Work** The models are presented with the image and caption pair and asked if the pair is misinformation. The models are further prompted to explain their reasoning. We present comparisons to other explainable methods, specifically Sniffer (Qi et al., 2024), CRAVE (Dey et al., 2025), and MMD-Agent (Liu et al., 2025) in the following section. We also compare MAD-Sherlock to existing pretrained multi-modal baselines including CLIP (Radford et al., 2021), VisualBERT (Li et al., 2019), InstructBLIP (Dai et al., 2023), LLaVA (Liu et al., 2024b), and GPT-4o (OpenAI and et al., 2024; OpenAI) on the NewsCLippings dataset in Appendix A.4. We show results for two baseline methods trained from scratch on the NewsCLippings, namely EANN (Wang et al., 2018) and SAFE (Zhou et al., 2020), in Appendix A.4. We further compare MAD-Sherlock to DT-Transformer (Papadopoulos et al., 2023a), CCN (Abdelnabi et al., 2022), VINVL (Huang et al., 2024), SSDL (Mu et al., 2023), and Neuro-Sym (Zhu et al., 2022) in Appendix A.4.

## 4.3. Results

We present results for the experiments conducted to select the best debate setup as well compare the performance of

MAD-Sherlock against existing methods. We use classification accuracy as the primary performance metric for comparison based on quantitative analysis.

### 4.3.1. COMPARING DEBATE SETUPS

We compare multiple debating setups using the LLaVA-v1.6-34B model, to select the best one for comparison with other works and further experimentation.

Debate Setup	Accuracy	Precision	Recall
Actor-Skeptic	69.5	66.1	69.4
Judged Debate	66.7	66.7	61.5
Async_Debate <sub>AI</sub> (believes debating AI)	75.2	54.5	86.4
Async_Debate <sub>human</sub> (w/o external info)	77.1	68.4	89.3
Async_Debate <sub>human</sub> (w external info)	<b>86.2</b>	<b>82.6</b>	<b>90.6</b>

**Table 1. Performance comparison between different debate setups:** We observe best performance with the The Async\_Debate<sub>human</sub> set-up, in which the model has external context and believes it is debating a human.

Table 1 shows that Async\_Debate<sub>human</sub> with external information outperforms all other debate configurations. To highlight the role of external retrieval, we also report results without it and observe a significant performance drop.

We find that agents perform better when they believe they are debating a human, engaging more critically with peer responses. The asynchronous setup also benefits from ensemble effects, unlike the actor-skeptic setup where only one agent generates responses.

In judged debates, agents must adhere to fixed stances and persuade a judge, which can lead to confusion and errors. In contrast, Async\_Debate<sub>human</sub> allows agents to revise their views mid-debate, leading to clearer outcomes. Based on these findings, we adopt this setup—with external information—for all further experiments.

### 4.3.2. PERFORMANCE COMPARISON

We present our results on the NewsCLippings, VERITE, and MMFakeBench datasets against existing out-of-context detection methods discussed in Section 4.2. We include additional comparisons against legacy methods in Appendix A.4.

Table 2 shows the comparison between our method and other explainable methods on the NewsCLippings, VERITE, and MMFakeBench datasets. We report state-of-the-art performance when using our proposed debate configuration with the GPT-4o (OpenAI and et al., 2024; OpenAI) model across all datasets. As a reminder, Sniffer (Qi et al., 2024) is finetuned extensively to adapt it to the NewsCLippings dataset. MMD-Agent represents the single-agent framework proposed by the authors of the MMFakeBench dataset (Liu et al., 2025). The CRAVE framework introduced by Dey

Dataset	Model	Performance		
		Accuracy	Precision	Recall
NewsCLIPpings	SNIFFER (Qi et al., 2024)	88.4	91.8	86.9
	CRAVE (Dey et al., 2025)	85.0	85.0	85.0
	LLM-Consensus (ours)	<b>90.8</b>	85.5	99.0
VERITE	CRAVE (Dey et al., 2025)	82.0	82.0	82.0
	LLM-Consensus (ours)	<b>85.2</b>	83.6	96.0
MMFakeBench	MMD-Agent (Liu et al., 2025)	62.1	67.8	59.3
	CRAVE (Dey et al., 2025)	78.0	83.6	69.7
	LLM-Consensus (ours)	<b>83.8</b>	87.2	90.1

**Table 2. Performance comparison of LLM-Consensus against other explainable methods on various datasets:** Performance of LLM-Consensus and baselines across NewsCLIPpings, VERITE, and MMFakeBench. Our model shows consistent generalization and superior performance across datasets.

et al. (2025) represents the most recent work across all three datasets, but our framework shows markedly superior performance across all datasets.

We find that our system produces coherent, detailed and comprehensive explanations. We also note that the debate paradigm in itself is essential to the system performance. We observe a drop in performance and quality of explanations when using an identical system configuration but with a single model. We believe that this follows from work by Mireshghallah et al. (2024), which introduces and demonstrates the importance of contextual privacy. Introducing multiple context windows allows each agent to maintain its own role and informational perspective without leakage, which is not possible in single-agent systems.

We provide a qualitative sample of our framework’s explanations in Appendix A.5. We provide ablation studies on each component of the MAD-Sherlock pipeline in Appendix A.6.

We also note that single multi-modal models that do not do retrieval, including VisualBERT, CLIP, InstructBLIP, LLaVA and GPT-4o do not perform at par with other related work. We demonstrate these results in Appendix A.4. This can be attributed to the necessity for external context for misinformation detection and the lack of diverse perspectives that occur naturally in a multi-agent framework. These models require additional integration into more comprehensive pipelines, as done in this work.

## 5. User Study

We conducted a user study to assess our system’s ability to detect and explain misinformation—particularly important given the lack of standard metrics for evaluating explanation quality. Participants were grouped by profession: Journalists, Academics, and Others (see Appendix A.7).

Each participant reviewed ten image-text pairs, judged their veracity, and rated their confidence (0–10). After submitting initial responses, they viewed MAD-Sherlock’s explanations and revised their answers. As shown in Table 3, MAD-

Study Setup	Average Accuracy
Humans	60.3 ± 13.5
Humans+MAD-Sherlock	76.7 ± 12.2
<b>MAD-Sherlock</b>	<b>80.0 ± 0.0</b>

**Table 3. Performance comparison between different study setups:** MAD-Sherlock outperforms humans with and without AI assistance.

Metric	Journalists	Academics	Others
Accuracy (only human)	70.0 ± 1.4	60.7 ± 1.4	56.7 ± 1.5
Confidence (only human)	4.3 ± 2.1	3.2 ± 0.8	3.9 ± 1.2
Accuracy (with MAD-Sherlock)	82.2 ± 0.9	79.3 ± 1.3	71.7 ± 1.1
Confidence (with MAD-Sherlock)	5.3 ± 1.3	5.8 ± 1.4	5.8 ± 1.4

**Table 4. Performance comparison:** MAD-Sherlock improves performance across all participant groups.

Sherlock outperformed average human accuracy both with and without AI assistance, highlighting its potential for improving public safety and trust.

Group-wise analysis, shown in Table 4, reveals significant performance gains across all groups, with results approaching those of professional journalists. Confidence levels (out of 10) are comparable across groups and generally increase after using MAD-Sherlock insights. Thus, MAD-Sherlock can substantially boost non-expert performance, making it valuable for citizen intelligence applications.

## 6. Conclusion and Future Work

Out-of-context (OOC) image misuse is an increasing challenge for misinformation detection, especially as vision-language models grow more powerful yet less interpretable. We explore whether multiple AI agents can collaboratively reason about context to improve prediction accuracy. Our strongest results come from the `AsynchronousDebatehuman` setup, where agents believe they are debating a human. This setting promotes engagement, mid-debate revision, and better identification of subtle inconsistencies.

Our final system, MAD-Sherlock, achieves state-of-the-art performance while offering interpretable, evidence-based explanations—enabled by our external retrieval module. We observe substantial gains in OOC detection across both expert and non-expert users.

We see several avenues for future work, including a continuously updated benchmark with recent news and more nuanced inconsistencies, extending to video-text and multi-lingual inputs, and large-scale deployments in professional and citizen intelligence settings. We hope to engage with the community to better understand how agentic workflows can enhance online safety. For limitations, see Appendix A.1.

## References

- Sahar Abdelnabi, Rakibul Hasan, and Mario Fritz. Open-domain, content-based, multi-modal fact-checking of out-of-context images via online resources, 2022. URL <https://arxiv.org/abs/2112.00061>.
- Irene Amerini, Aris Anagnostopoulos, Luca Maiano, Lorenzo Ricciardi Celsi, et al. Deep learning for multimedia forensics. *Foundations and Trends® in Computer Graphics and Vision*, 12(4):309–457, 2021.
- Shivangi Aneja, Cise Midoglu, Duc-Tien Dang-Nguyen, Sohail Ahmed Khan, Michael Riegler, Pål Halvorsen, Chris Bregler, and Balu Adsumilli. Acm multimedia grand challenge on detecting cheapfakes, 2022. URL <https://arxiv.org/abs/2207.14534>.
- Kevin Aslett, Zeve Sanderson, William Godel, Nathaniel Persily, Jonathan Nagler, and Joshua A Tucker. Online searches to evaluate misinformation can increase its perceived veracity. *Nature*, 625(7995):548–556, 2024.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. Hallucination of multimodal large language models: A survey, 2024. URL <https://arxiv.org/abs/2404.18930>.
- Nadia M Brashier and Elizabeth J Marsh. Judging truth. *Annual review of psychology*, 71(1):499–515, 2020.
- Jonah Brown-Cohen, Geoffrey Irving, and Georgios Pilioras. Scalable ai safety via doubly-efficient debate, 2023. URL <https://arxiv.org/abs/2311.14125>.
- Ivan Castillo Camacho and Kai Wang. A comprehensive review of deep-learning-based methods for image forensics. *Journal of imaging*, 7(4):69, 2021.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. URL <https://arxiv.org/abs/2305.06500>.
- Arka Ujjal Dey, Muhammad Junaid Awan, Georgia Channing, Christian Schroeder de Witt, and John Collomosse. Fact-checking with contextual narratives: Leveraging retrieval-augmented llms for social media analysis, 2025. URL <https://arxiv.org/abs/2504.10166>.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving Factuality and Reasoning in Language Models through Multiagent Debate. October 2023. URL <https://openreview.net/forum?id=QAwaaLJNCK>.
- Hany Farid. *Photo Forensics*. The MIT Press, 2016. ISBN 0262035340.
- Lisa Fazio. Out-of-context photos are a powerful low-tech form of misinformation — pbs.org, 2020. URL <https://www.pbs.org/newshour/science/>. [Accessed 28-09-2024].
- Lynn Hasher, David Goldstein, and Thomas Toppino. Frequency and the conference of referential validity. *Journal of verbal learning and verbal behavior*, 16(1):107–112, 1977.
- Arash Heidari, Nima Jafari Navimipour, Hasan Dag, and Mehmet Unal. Deepfake detection using deep learning methods: A systematic and comprehensive review. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 14(2):e1520, 2024.
- Maryam Hina, Mohsin Ali, Abdul Rehman Javed, Fahad Ghabban, Liaqat Ali Khan, and Zunera Jalil. Sefaced: Semantic-based forensic analysis and classification of e-mail data using deep learning. *IEEE Access*, 9:98398–98411, 2021.
- Mingzhen Huang, Shan Jia, Zhou Zhou, Yan Ju, Jialing Cai, and Siwei Lyu. Exposing text-image inconsistency using diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Akbar Khan, John Hughes, Dan Valentine, Laura Ruis, Kshittij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with more persuasive llms leads to more truthful answers, 2024. URL <https://arxiv.org/abs/2402.06782>.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- Hongzhan Lin, Ziyang Luo, Wei Gao, Jing Ma, Bo Wang, and Ruichao Yang. Towards explainable harmful meme detection through multimodal debate between large language models, 2024. URL <https://arxiv.org/abs/2401.13298>.
- Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visualnews : Benchmark and challenges in entity-aware image captioning, 2020.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating hallucination in large multi-modal models via robust instruction tuning, 2024a. URL <https://arxiv.org/abs/2306.14565>.

- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- Xuannan Liu, Zekun Li, Peipei Li, Huaibo Huang, Shuhan Xia, Xing Cui, Linzhi Huang, Weihong Deng, and Zhaofeng He. Mmfakebench: A mixed-source multimodal misinformation detection benchmark for llms, 2025. URL <https://arxiv.org/abs/2406.08772>.
- Grace Luo, Trevor Darrell, and Anna Rohrbach. Newsclippings: Automatic generation of out-of-context multimodal media. *arXiv:2104.05893*, 2021.
- Momina Masood, Mariam Nawaz, Khalid Mahmood Malik, Ali Javed, Aun Irtaza, and Hafiz Malik. Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward. *Applied intelligence*, 53(4):3974–4026, 2023.
- Niloofar Mireshghallah, Hyunwoo Kim, Xuhui Zhou, Yulia Tsvetkov, Maarten Sap, Reza Shokri, and Yejin Choi. Can llms keep a secret? testing privacy implications of language models via contextual integrity theory, 2024. URL <https://arxiv.org/abs/2310.17884>.
- Michael Mu, Sreyasee Das Bhattacharjee, and Junsong Yuan. Self-supervised distilled learning for multi-modal misinformation identification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2819–2828, January 2023.
- OpenAI. GPT-4o. <https://openai.com/index/hello-gpt-4o/>. [Accessed 28-08-2024].
- OpenAI and Josh Achiam et al. GPT-4 Technical Report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- Stefanos-Iordanis Papadopoulos, Christos Koutlis, Symeon Papadopoulos, and Panagiotis Petrantonakis. Synthetic misinformers: Generating and combating multimodal misinformation. In *Proceedings of the 2nd ACM International Workshop on Multimedia AI against Disinformation*, MAD ’23, page 36–44, New York, NY, USA, 2023a. Association for Computing Machinery. ISBN 9798400701870. doi: 10.1145/3592842.3592842. URL <https://doi.org/10.1145/3592842.3592842>.
- Stefanos-Iordanis Papadopoulos, Christos Koutlis, Symeon Papadopoulos, and Panagiotis C. Petrantonakis. Verite: A robust benchmark for multimodal misinformation detection accounting for unimodal bias, 2023b. URL <https://arxiv.org/abs/2304.14133>.
- Peng Qi, Zehong Yan, Wynne Hsu, and Mong Li Lee. Sniffer: Multimodal large language model for explainable out-of-context misinformation detection, 2024. URL <https://arxiv.org/abs/2403.03170>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- Christian Schroeder de Witt, Tarun Gupta, Denys Makovychuk, Viktor Makovychuk, Philip H. S. Torr, Mingfei Sun, and Shimon Whiteson. Is Independent Learning All You Need in the StarCraft Multi-Agent Challenge?, November 2020. URL <https://arxiv.org/abs/2011.09533v1>.
- Meredith Somers. Deepfakes, explained — MIT Sloan — mitsloan.mit.edu. <https://mitsloan.mit.edu/ideas-made-to-matter/deepfakes-explained>, 2020. [Accessed 28-09-2024].
- Mubashir Sultan, Alan N Tump, Michael Geers, Philipp Lorenz-Spreen, Stefan M Herzog, and Ralf HJM Kurvers. Time pressure reduces misinformation discrimination ability but does not alter response bias. *Scientific Reports*, 12(1):22416, 2022.
- Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information Fusion*, 64:131–148, 2020.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Sheng-Yu Wang, Oliver Wang, Andrew Owens, Richard Zhang, and Alexei A. Efros. Detecting photoshopped faces by scripting photoshop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. Eann: Event adversarial neural networks for multi-modal fake news detection. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD ’18, page 849–857, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450355520. doi: 10.1145/3219819.3219903. URL <https://doi.org/10.1145/3219819.3219903>.



Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, pages 24824–24837, Red Hook, NY, USA, April 2024. Curran Associates Inc. ISBN 978-1-71387-108-8.

Yizhou Zhang, Loc Trinh, Defu Cao, Zijun Cui, and Yan Liu. Interpretable detection of out-of-context misinformation with neural-symbolic-enhanced large multi-modal model, 2024. URL <https://arxiv.org/abs/2304.07633>.

Xinyi Zhou, Jindi Wu, and Reza Zafarani. Safe: Similarity-aware multi-modal fake news detection, 2020. URL <https://arxiv.org/abs/2003.04981>.

Wang Zhu, Jesse Thomason, and Robin Jia. Generalization differences between end-to-end and neuro-symbolic vision-language reasoning systems, 2022. URL <https://arxiv.org/abs/2210.15037>.

Xinshan Zhu, Yongjun Qian, Xianfeng Zhao, Biao Sun, and Ya Sun. A deep learning approach to patch-based image inpainting forensics. *Signal Processing: Image Communication*, 67:90–99, 2018.

## A. Appendix

### A.1. Limitations

Despite the strong performance of MAD-Sherlock, several limitations remain. First, while our model excels at detecting out-of-context image–text pairs, its reliance on external retrieval can lead to reduced accuracy when relevant context is unavailable or difficult to retrieve. Moreover, our framework cannot independently verify the factual correctness of externally retrieved news articles; the truthfulness of any source may be debated, potentially introducing misinformation into the reasoning process. Nevertheless, we observe that the model’s judgments align closely with the human-created labels used in these widely accepted datasets, underscoring its practical utility. Second, the quality of explanations is constrained to textual outputs, limiting multi-modal explanation capabilities such as image or video integration. Third, the system’s performance is sensitive to hyperparameter tuning, including the number of debate rounds and agents, which may require further optimization for broader use cases.

Additionally, while our user studies provided valuable insights, large-scale deployment in diverse, real-world settings, such as professional or citizen intelligence environments, is necessary to fully assess the method’s robustness and scalability. Finally, our dataset, though comprehensive, primarily focuses on English-language news, limiting the generalizability of the system across non-English contexts.

Another important limitation is the potential risk that open-sourcing MAD-Sherlock might allow adversaries to train models specifically designed to counter or evade detection by our system. As adversarial actors gain access to the source code, they could exploit its known strengths and weaknesses to develop countermeasures that diminish its effectiveness. However, despite these risks, we believe that open-sourcing remains the right path forward. Open-sourcing encourages transparency, collaboration, and rapid innovation, enabling the broader community to contribute improvements, detect vulnerabilities, and build on the system.

Moreover, by engaging the community, we can foster the development of more resilient and adaptive models that evolve in response to emerging adversarial techniques, thus maintaining MAD-Sherlock’s effectiveness in the long term. The collective strength of a diverse, open-source community can outweigh the potential threats posed by adversarial exploitation.

Future work will need to address these limitations to enhance the practical utility, robustness, and long-term resilience of MAD-Sherlock.

### A.2. Sample Image-Caption Pair in the News Domain



Figure 4. Russian President Vladimir Putin has called Ukraine’s move into Kursk a “major provocation”. Image and caption taken from the BBC article here (Accessed at 17:43 on Aug 11, 2024): <https://www.bbc.co.uk/news/articles/cze5pkg5jwlo>

### A.3. Prompts for MAD-Sherlock

You are a misinformation detection expert in the news domain. You will look at image-caption pairs and decide if the given image is rightly used in the given news context. To further assist you, a summary of news articles related to the image will be provided. Based on this, you need to decide if the caption belongs to the image or if it is being used to spread false information to mislead people. Note that the image is real. It has not been digitally altered. Carefully examine the image for any known entities, people, watermarks, dates, landmarks, flags, text, logos and other details which could give you important information to better explain your answer. Remember in news articles images and captions are often related contextually and the caption need not exactly describe the image. The goal is to consider the contextual relationship between the image and caption based on the news articles and correctly identify if the image caption pair is misinformation or not and to explain your answer in detail. Think step by step and plan a detailed explanation for your answer.

Figure 5. Model initialization prompt.

This is a summary of news articles related to the image: {}  
Based on this, you need to decide if the caption given below belongs to the image or if it is being used to spread false information to mislead people.  
CAPTION: {}  
Note that the image is real. It has not been digitally altered.  
Carefully examine the image for any known entities, people, watermarks, dates, landmarks, flags, text, logos and other details which could give you important information to better explain your answer.  
The goal is to correctly identify if this image caption pair is misinformation or not and to explain your answer in detail.  
At the end give a definite YES or NO answer to this question:  
IS THIS MISINFORMATION?

Figure 6. Initial prompt for independent opinion formation and response generation with GPT-4o.

This is what I think: {}.  
Do you agree with me?  
If you think I am wrong then convince me why you are correct.  
Clearly state your reasoning and tell me if I am missing out on some important information or am making some logical error.  
Do not describe the image.  
At the end give a definite YES or NO answer to this question:  
IS THIS MISINFORMATION?

Figure 7. Prompt for first round of debate.

I see what you mean and this is what I think: {}.

Do you agree with me?

If not then point out the inconsistencies in my argument (e.g. location, time or person related logical confusion) and explain why you are correct.

If you disagree with me then clearly state why and what information I am overlooking.

Find disambiguation in my answer if any and ask questions to resolve them.

I want you to help me improve my argument and explanation.

Don't give up your original opinion without clear reasons, DO NOT simply agree with me without proper reasoning.

At the end give a definite YES or NO answer to this question:

IS THIS MISINFORMATION?

Figure 8. Prompt for subsequent rounds of debate.

## A.4. Supplemental Results

### A.4.1. ADDITIONAL RESULTS ON VERITE

Model	External Retrieval	Data	Accuracy	Precision	Recall
4o	✓	VERITE (all)	85.2	83.6	96.0
4o	✓	VERITE (true vs. OOC)	79.5	69.7	96.4
4o	✓	VERITE (true vs. MC)	80.7	73.6	95.6
4o	✗	VERITE (all)	84.8	85.0	93.1
4o	✗	VERITE (true vs. OOC)	78.6	72.1	89.5
4o	✗	VERITE (true vs. MC)	82.6	75.6	96.4

Table 5. **External Retrieval Ablation on VERITE.** We compare performance of LLM-Consensus with and without external retrieval with a GPT-4o backbone. We also disaggregate our results by “miscaptioned” (MS) and “out-of-context” (OOC), two different classes of misinformation included in the VERITE dataset.



Model	Accuracy $\uparrow$
SAFE	50.7
EANN	58.1
VisualBERT	54.8
CLIP	62.6
InstructBLIP	48.6
LLaVA	57.1
GPT-4o	70.7
DT-Transformer	77.1
CCN	84.7
SSDL	65.6
VINVL	65.4
Neuro-Sym	68.2
GPT-4o <sup>#</sup> (w retrieval)	86.0
Sniffer (w finetuning)	88.4
Sniffer (w/o finetuning)	84.5
<b>MAD-Sherlock (ours)</b>	<b>90.8</b>

Table 6. Performance comparison on NewsCLIPPings against various legacy methods: MAD-Sherlock (with GPT-4o) outperforms all related work. Note: the GPT-4o<sup>#</sup> setup is an ablation in which there is a single agent using our external retrieval module (no debate).

## A.4.2. ADDITIONAL BASELINES FOR NEWSCLIPINGS

## A.5. Qualitative Analysis

Image Description: A banner that says ``STOP VICTIMISATION OF UNION MEMBERS END THE BLACKLIST`` in front of a city landscape.

Caption: Flashmob protest at Crossrail site in Central London by Unite union members in protest against blacklisting and sackings of union officials and members.

Label: This is not misinformation.

Framework prediction: This is not misinformation.

Explanation (Model Output): Here is a precise breakdown:

1. Image Analysis:

- The banner in the image clearly states ``STOP VICTIMISATION OF UNION MEMBERS END THE BLACKLIST``
- This message aligns with union-related grievances focused on stopping blacklisting and mistreatment of union members.

2. Caption Analysis:

- The caption explicitly mentions a flashmob protest by Unite union members at a Crossrail site in Central London.
- The topics of blacklisting and the sacking of union officials and members are mentioned in the caption, aligning with the grievances captured in the image.

3. News Summary Correlation:

- The true news articles discuss blacklisting practices within the construction industry, particularly concerning the Crossrail project and actions taken by the Consulting Association.
- Gail Cartmail from Unite the Union has provided evidence of blacklisting affecting union members working on the Crossrail project, which fits the context depicted in the banner.
- The articles also highlight broader issues of blacklisting, further supporting the theme reflected in the image and caption.

4. Contextual Match:

- The image and the caption are aligned in context. The message on the banner directly supports the protest described in the caption.
- The news article summary corroborates the issues of blacklisting faced by union members, as illustrated by the protest.

Your reasoning is logically sound, and no critical information is missing. Both the image and the caption are contextually connected to the theme of union protests against blacklisting, specifically within the Crossrail project setting.

Thus, I conclude:

Definite Answer: NO, the image-caption pair does not represent misinformation.

Figure 9. Qualitative sample of model explanation.

LLaVA	GPT-4o	Retrieval	Debate	Accuracy	Precision	Recall
✓	✗	✗	✓	77.1	68.4	89.3
✓	✗	✓	✓	86.2	82.6	90.6
✗	✓	✓	✗	86.0	80.2	95.6
✗	✓	✗	✓	90.2	<b>90.3</b>	90.1
✗	✓	✓	✓	<b>90.8</b>	85.5	<b>99.0</b>

Table 7. **Ablation:** Quantitative evaluation of each component of MAD-Sherlock on NewsCLiPPings classification performance.

## A.6. Ablations

### A.6.1. ABLATIONS

To analyze the importance of each component of the MAD-Sherlock framework, we conduct ablation experiments. Specifically, we evaluate the effect of using LLaVA against GPT-4o, the impact of the external retrieval module, and the power of the multi-agent debate framework.

We observe that the combination of GPT-4o, the external retrieval module, and the multi-agent debate framework yields the highest performance across all metrics, with 90.8% accuracy, 85.5% precision, and 99.0% recall, demonstrating the value of combining these components. The inclusion of debate alone significantly boosts accuracy from the GPT-4o baseline of 70.7% (as seen in Table 6) to 90.2%, underscoring its role in enabling contextual reasoning and refining predictions. Adding external retrieval to the GPT-4o with debate system primarily shifts the balance between precision and recall, where precision moves from 90.3% to 85.5% and recall from 90.1% to 99.0%. Meanwhile, retrieval contributes more substantially to LLaVA’s performance gains, likely due to GPT-4o’s broader world knowledge. Without external retrieval or the debate framework, the performance drops, emphasizing the critical role of these components in achieving state-of-the-art results. We show additional results on VERITE when ablating the retrieval module in Appendix A.4.

## A.7. User Study

We conduct a user study to assess the effectiveness of our model in detecting and explaining misinformation. Through this study, we aim to assess the persuasiveness of our system.

### A.7.1. STUDY DESIGN

The user study was designed to evaluate the effectiveness of our system in detecting and explaining misinformation. While it is easy to quantify model performance in terms of misinformation detection, there are no effective metrics to assess the quality of the explanations generated by the model. Therefore, in order to perform a thorough analysis of the system performance, a user study is essential.

A total of 30 participants volunteered to participate in this study. Participation was completely voluntary and no personal information was used for the purpose of analysis in this study. For a deeper analysis we further grouped the participants based on their profession into three groups, namely: Journalists, AI Academics and Others. The ‘others’ category included anyone who did not belong to the first two groups. The study was conducted through a Microsoft Form. Participants were shown 10 image-text pairs and were asked to decide if the image and caption when considered together was misinformation or not. They were also asked to provide a confidence rating for their answer on a scale of 0-10, with 10 being the highest confidence level. For each image-text pair, after the participants provided their initial answers, they were shown AI insights about the same image-text pair. These AI insights were the final outputs from MAD-Sherlock. Participants were then asked to reconsider their answer and again decide if the image-text pair was misinformation or not, in light of the new information from the AI agent. Participants were also required to re-evaluate their confidence score in this new answer. While it is not entirely avoidable, we did ask participants to keep aside their personal opinions of AI and consider all AI insights objectively. Participants were not allowed to access the Internet. This was done to ensure an unbiased estimate of average human performance.

The image-text pairs to include in the study were taken from the NewsCLIPpings (Luo et al., 2021) dataset. AI insights were taken from our best-performing setup involving the GPT-4o model. Of the 10 image-text pairs presented to the participants in the study, there were 5 instances of misinformation and 5 instances of true information. Further, all model insights were true except two of them. Therefore the model accuracy for the task was 80% and we use this as the baseline accuracy to compare human performance against.

We analyse two special cases, where MAD-Sherlock argues for the wrong answer. We include these results in order to observe how persuasive our system can be even when it is wrong. We note in the instance where the image-text pair was actually misinformation and the model argued that it was not, 6 participants changed their correct responses to those suggested by MAD-Sherlock. Although this is only 5% of the participants, it still gives a significant insight into how persuasive the model can appear even when it is wrong. While the case of false negatives is important, false positives are an even more concerning matter for our problem statement. In the case where MAD-Sherlock declared the given image-text pair to be misinformation when it was not, is important to analyse. In this setting 50% of the total participants changed their answer to the wrong one, therefore believing a piece of true information to be false. In some cases where participants chose the wrong response to begin with, their confidence in the response further increased after considering insights from the system. Finally, 4 participants did not change their answer to the wrong one after considering AI insights but their confidence in their response decreased.

The average time taken to complete the study was 12 minutes and 57 seconds. The average participant was therefore able to go through 10 image-text pairs and decide if they were misinformation or not in under 13 minutes. The same task without AI insights would require extensive analysis and we project it would take between 30-45 minutes to decide if 10 image-text pairs were misinformation.

#### A.8. Screenshots

We include representative screenshots of the user study.

#### A.9. Multi-modal debates for harmful meme detection

While this work relates to a different problem than OOC misinformation detection in the news domain, we still find the approach taken by the authors a relevant related work and therefore include it here. Lin et al. (2024) use LMMs debating against each other to generate explanations for contradictory arguments regarding whether a given meme is harmful. These explanations are then used to train a small language model as a judge to determine whether the image and text that make up the meme are actually harmful. This work does not allow agents to have flexibility of opinion. There are always two agents, and each one is provided a stance to defend. Moreover, a judge decides the final outcome of the debate and needs to be trained on data from the debate. This method also does not benefit from external retrieval, and therefore, the debating agents are not aware of the crucial external context related to the input. Finally, this work is related to harmful *meme* detection and does not concern the problem of misinformation detection in the news domain, which likely requires more intricate contextual analysis, including of external context.

#### A.10. Additional Experiments

**Debate with Disambiguation:** Building on the actor-skeptic method, we allow all agents to act as both actors and skeptics. Models generate their own responses and disambiguation queries to refine or challenge other agents’ outputs. These queries are used to retrieve additional information from the Internet, further improving model responses. The Debate with Disambiguation strategy achieves accuracy of 77.8, precision of 74.7, and recall of 82.6 when tested with a LLaVa backbone.



### Misinformation Survey

\* Required

#### Survey Instructions

In this survey, you will evaluate 10 image/caption pairs and report whether or not you believe they are misinformation. In some of the survey questions, we provide the output of an AI agent's evaluation of the image/caption pair for your reference. The AI agent has access to the internet.

P.S: This survey contains Karma to get free survey responses at [SurveySwap.io](https://SurveySwap.io)

What is your profession? \*

☐ Student

☐ Researcher


☐ Professional

☐ Other

#### Sample 1: No Internet

Please evaluate the image/caption pair to determine if they are misinformation.

Given this image/caption pair, do you believe that this is misinformation?  
Please **do not** consult the internet in formulating your response.



CAPTION: Flashmob protest at Crossrail site in Central London by Unite union members in protest against blacklisting and sackings of union officials and members. \*

☐ Yes, I believe it is misinformation.

☐ No, I do not believe it is misinformation.

How confident are you in your assessment? \*

0	1	2	3	4	5	6	7	8	9	10
---	---	---	---	---	---	---	---	---	---	----

Not at all confident Extremely confident


(a) Instructions screen

(b) No internet baseline

#### Sample 1: With Internet

Please evaluate the image/caption pair to determine if they are misinformation.

Given this image/caption pair, do you believe that this is misinformation?  
**Feel free** to use the internet in formulating your response.



CAPTION: Flashmob protest at Crossrail site in Central London by Unite union members in protest against blacklisting and sackings of union officials and members. \*

☐ Yes, I believe it is misinformation.

☐ No, I do not believe it is misinformation.

How confident are you in your assessment? \*

0	1	2	3	4	5	6	7	8	9	10
---	---	---	---	---	---	---	---	---	---	----

Not at all confident Extremely confident


(c) With internet baseline

#### Sample 1: With Retrieval

Please evaluate the image/caption pair to determine if they are misinformation.

Our AI agent has done research and summarized its findings on whether or not this is misinformation. See the summary here:

"Image Content Analysis: The image shows a group of people, likely workers, standing in front of a prominent government building (Palace of Westminster, indicative of London). The banner includes the words "Community," "the union steelworkers," and "For a better working world," which strongly suggests the involvement of steelworkers and their union. Caption Analysis: The caption states: "Flashmob protest at Crossrail site in Central London by Unite union members in protest against blacklisting and sackings of union officials and members." This suggests a specific protest related to the Crossrail project involving blacklisting and the sacking of union



(d) With retrieval-augmented summary

Figure 10. Screenshots of the user survey. The questions asked after AI summary is presented are the same as those following the other questions.