

# MAD-SHERLOCK: MULTI-AGENT DEBATES FOR OUT-OF-CONTEXT MISINFORMATION DETECTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

One of the most challenging forms of misinformation involves the out-of-context (OOC) use of images paired with misleading text, creating false narratives. Existing AI-driven detection systems lack explainability and require expensive finetuning. We address these issues with MAD-Sherlock: a **Multi-Agent Debate** system for OOC Misinformation Detection. MAD-Sherlock introduces a novel multi-agent debate framework where multimodal agents collaborate to assess contextual consistency and request external information to enhance cross-context reasoning and decision-making. Our framework enables explainable detection with state-of-the-art accuracy even without domain-specific fine-tuning. Extensive ablation studies confirm that external retrieval significantly improves detection accuracy, and user studies demonstrate that MAD-Sherlock boosts performance for both experts and non-experts. These results position MAD-Sherlock as a powerful tool for autonomous and citizen intelligence applications.

## 1 INTRODUCTION

Our growing dependence on online channels for news and social networking has been complemented by a surge in exploits of digital misinformation (Aslett et al., 2024; Hasher et al., 1977; Brashier & Marsh, 2020). While many manipulation techniques pose serious threats, one of the most prevalent methods for creating fake online content is the out-of-context (OOC) use of images (pbs). This involves using unaltered images in a misleading, false context to convey deceptive information, a strategy that requires minimal technical expertise. Indeed, the problem of OOC misinformation detection requires a complex understanding of the relationship between the text and image and the ability to identify when they do not go together. Identifying these minute inconsistencies is a time-consuming and high-effort task for humans. A study by Sultan et al. (2022) shows that time pressure reduces the ability of human beings to detect misinformation effectively, further adding to the scalability issues in human expert detection.

Therefore, attention has turned to AI-driven tools that can help human experts recognise instances of OOC image-based misinformation at scale. Unfortunately, conventional deep learning forensic techniques (Castillo Camacho & Wang, 2021; Heidari et al., 2024; Zhu et al., 2018; Amerini et al., 2021; Hina et al., 2021), which target detecting manipulations such as PhotoShop editing (Tolosana et al., 2020; Masood et al., 2023; Farid, 2016; Wang et al., 2019) and AI-generated (or manipulated) fake images called Deepfakes (mit), rely on spotting artifacts from image or text tampering. In contrast, OOC detection demands cross-contextual reasoning, as the deception arises from the misalignment between the legitimate image and its falsely associated textual content.

Pretrained Large Multimodal Models (Liu et al., 2024b; OpenAI & et al., 2024; Li et al., 2019; Radford et al., 2021, LMMs) provide a promising direction for detecting OOC use of images for their ability to process both text and image content in tandem. However, using LMMs directly for OOC detection presents several challenges, particularly in the news domain. For instance, news articles often include images that are not directly related to the article’s content. An article about the 2024 U.S. presidential candidates, for example, might feature a close-up of Donald Trump from an unrelated online database. Although the image was taken outside the election period, it is not considered OOC since it doesn’t misrepresent the article’s context. Such cases complicate LMMs’ ability to accurately identify OOC usage based solely on their pre-trained knowledge, as this knowledge may be outdated or insufficiently detailed.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

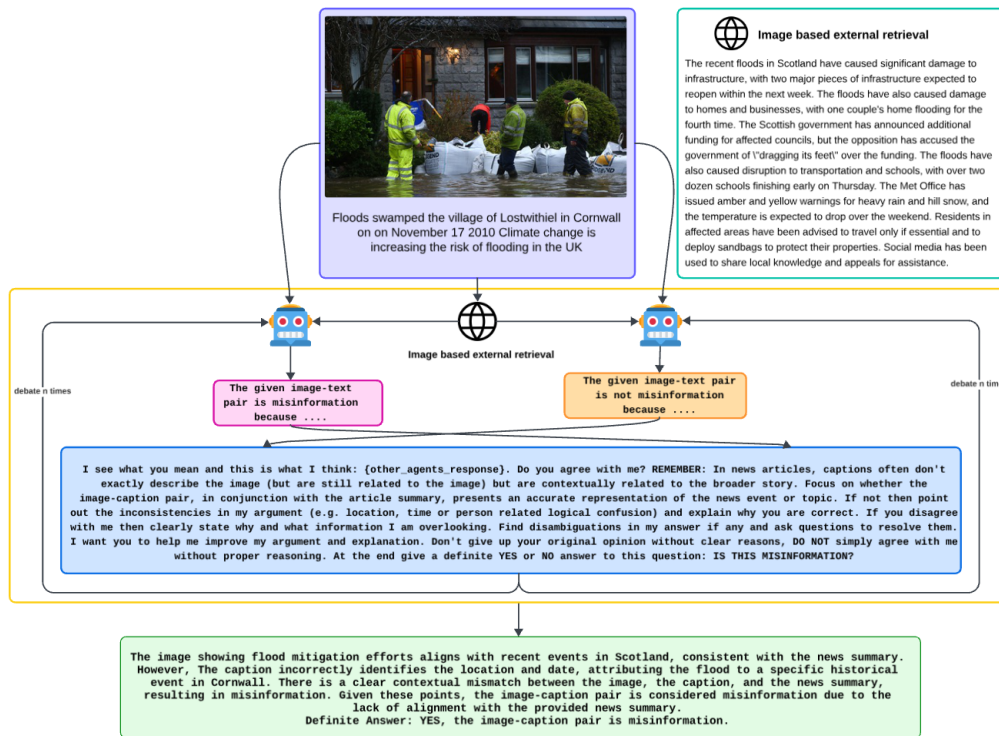


Figure 1: **Overview of MAD-Sherlock:** Two or more independent agents see the same image-text input and are tasked with detecting whether the input is misinformation or not. After the agents form their independent opinions, they participate in a debate until they converge on the same response or when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

Moreover, even with recent advancements, LMMs are capable of hallucinating and, hence, generating false information (Bai et al., 2024; Liu et al., 2024a). While rapidly improving, they sometimes fail to understand user instructions and intent correctly. We show that off-the-shelf LMMs indeed suffer from these issues, reducing their ability to detect OOC misinformation in practice. While prior work (Qi et al., 2024) has shown that off-the-shelf models can be improved using task-specific fine-tuning, this approach is resource-intensive and requires continual updating to keep up with recent events. Moreover, detecting OOC images only solves part of the problem. The real value lies in being able to *explain* the OOC use of pictures in human-readable form. It can be instrumental for human validators to observe the model’s line of logic and gain better insight into, and trust in, the classification process.

In this work, we propose a novel LMM-based post-training approach for scalable OOC misinformation detection that simultaneously improves contextual reasoning, provides in-built explainability, and achieves state-of-the-art detection accuracy even without task-specific fine-tuning (see Section 3). Specifically, our framework *MAD-Sherlock: a Multi-Agent Debate system for OOC Misinformation Detection* frames the detection problem as a dialectic debate between multiple LMM agents, where, in contrast to prior work (Minsky, 1988; Li et al., 2023a; Du et al., 2023a; Khan et al., 2024)), agents have access to external information retrieval.

Compared to single-agent chain-of-thought approaches (Wei et al., 2024), the use of multiple agents allows for a clean separation of agent contexts, decentralisation of action spaces, and opportunities for parallel computation (Schroeder de Witt et al., 2020; Du et al., 2023b). In addition, due to its compositional nature, both additional human and autonomous agents can be dynamically added to the multi-agent reasoning process, allowing the use of MAD-Sherlock as an interactive tool for human experts. To the best of our knowledge, no prior work has used debating LMMs for detecting and *explaining* OOC image use.

108 We perform a comprehensive empirical evaluation of our method (see Section 4), including the  
109 study of multiple debate configurations. To optimise OpenAI API use, we utilize an experimental  
110 pipeline where preliminary experiments are performed using the open-source LLaVA model (Liu  
111 et al., 2024b), which is only later replaced with GPT-4o (OpenAI) to achieve state-of-the-art per-  
112 formance. We find that *MAD-Sherlock* outperforms both prior work and novel baselines that we  
113 introduce, is more robust to various failure modes, and produces coherent explanations that help  
114 both human experts and non-experts significantly improve their detection accuracy in user studies.  
115 We identify both access to external information retrieval and complete freedom of opinion as key  
116 ingredients to *MAD-Sherlock*’s performance. Finally, we discuss current limitations of our method  
117 and propose future work toward overcoming the scalability challenges in large-scale online OOC  
118 misinformation detection.

## 119 2 RELATED WORK

122 Recent work has focused on using joint image-text representations to classify an instance as OOC.  
123 Aneja et al. (2022) follow a self-supervised approach to assess whether two captions accompanying  
124 an image are contextually similar. They enforce image-text matching during training by formulating  
125 a scoring function to align objects in the image with the caption. During inference, they use the se-  
126 mantic similarity between the two captions to classify them as OOC or not. The increased reliance  
127 on textual content limits the capabilities of this approach. This work also does not provide expla-  
128 nations for model predictions and is, therefore not interpretable. Moreover, this method works for  
129 image caption pairs where captions have information about objects in the image. This is not always  
130 the case with news articles (our domain of application), where captions can often just be related to  
131 the main content of an article rather than precisely describing the objects in the image. Appendix  
132 A.2 shows an example of the same.

133 Abdelnabi et al. (2022) present the Consistency Checking Network (CCN) in which they emulate  
134 different aspects of human reasoning across modalities for misinformation detection. This method  
135 uses evidence related to the image-text pair aggregated from the Internet. The CCN consists of  
136 memory networks to assess the consistency of the image-caption pair against the retrieved evidence  
137 and a CLIP (Radford et al. (2021)) component to evaluate the consistency between the image and  
138 caption pair. The use of external evidence to better inform model decisions is an important idea  
139 and also explains the superior classification performance of CCN when compared to other methods.  
140 This method also lacks the explainability component.

141 Zhang et al. (2024) extend the neural symbolic method (Yi et al. (2019); Zhu et al. (2022)) to propose  
142 an interpretable cross-modal misinformation detection model to provide supporting evidence for  
143 the output prediction. They use symbolic graphs based on the Abstract Meaning Representation  
144 (Banarescu et al. (2013)) of textual and visual information to detect OOC image use. Zhou et al.  
145 (2020) introduce Similarity Aware Fake news detection (SAFE), where neural networks are used  
146 to learn features of text and visual news representations. Their representations and relationships  
147 are jointly learned and used to predict fake news. Wang et al. (2018) introduce EANN: Event  
148 Adversarial Neural Networks to derive event invariant features which can be used to detect fake  
149 news that has recently been generated. EANN uses adversarial training to learn multi-modal features  
150 independent of news events. These methods require pretraining from scratch and, therefore, don’t  
benefit from the advanced reasoning capabilities and world knowledge of large pretrained models.

151 Shalabi et al. (2023) use synthetic multi-modal data to establish the authenticity of image-text pairs.  
152 They use BLIP-2 (Li et al. (2023b)) to generate a caption for the original image and Stable Diffusion  
153 (Rombach et al. (2022)) to generate an image for the given original caption. This synthetic data is  
154 then used to reason that if the original image and caption are OOC, then the original and generated  
155 images should also be OOC as well as the original and generated text. This method relies on syn-  
156 thetic multi-modal data generation, which not only adds an additional computational overhead but  
157 also increases dependence on often unreliable synthetically generated data. Therefore, this method  
158 can suffer from issues related to generation models, including potential biases that these models may  
159 possess. This method also lacks interpretability.

160 Sniffer (Qi et al. (2024)) is the closest to our work. It uses the InstructBLIP (Dai et al. (2023)) model  
161 to detect OOC image use and provide an explanation for its prediction. It makes use of internal and  
external knowledge using entity extraction APIs and image-based web searches. Information from

all the sources is given to an LLM to predict and explain if an image has been used OOC. Sniffer only uses basic textual information, such as news article titles from websites, to form its external knowledge base. It also requires extensive training to adapt the model to the news domain which adds additional computational overhead and also restricts the generalization abilities of the model to other domains.

### 3 METHODOLOGY

We present an explainable misinformation detection system, MAD-Sherlock, which jointly predicts and explains instances of misinformation. Figure 1 illustrates our approach. To the best of our knowledge, all prior work except Qi et al. (2024) provide predictions without explanations, and no prior work uses multiple models to approach this problem. We present a novel methodology that involves multiple multi-modal models debating against each other in order to decide if an image-text pair is misinformation or not. In this work, we aim to answer the question:

*Can debating multi-modal models, when equipped with external context, be used to solve the problem of explainable misinformation detection by picking up on minute contextual inconsistencies?*

We use detailed external information to inform the model’s predictions through the external retrieval module, which utilises reverse image-based search in order to provide the agents with external real-world context related to the image-text pair. We carry out our experiments with the GPT-4o (OpenAI) model to achieve state-of-the-art performance on the misinformation detection task while also providing detailed and coherent explanations for the predictions. We achieve this without any domain-specific fine-tuning, thus ensuring easier and faster generalization to other domains in addition to low 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. Instead of simple back-and-forth conversations, we opt for a debating set-up in which agents are asked to frame their own opinions and then defend them to other agent(s). We observe this facilitates more involved and detailed discussions among the models.

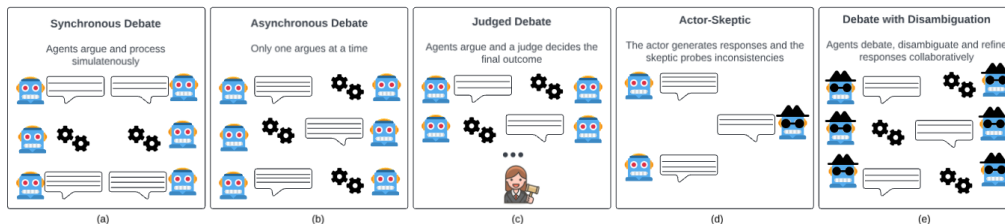


Figure 2: **Debating Strategies:** We experiment with multiple debating strategies. The asynchronous debate setup where agents argue one after the other and take turns presenting their arguments is the best configuration.

**Asynchronous Debate (not) against Human:** This is one of the core setups we test in our experiments. We define an asynchronous debating strategy in which models wait for the other participants’ responses before generating their own. Figure 2 (a) and (b) show synchronous and asynchronous debating structures, respectively. While synchronous debates, where all participants speak at once, can be faster and computationally more efficient, we opt for an asynchronous setting where each model response is based on previous responses of other models. We observe that this method of structuring the debate works better since models are able to pick up on contextual ambiguities in their responses in a more organised and structured way which is crucial to the process of misinformation detection.

216 An important point to note for this setup is the way we structure model prompts. The debating mod-  
217 els are not aware that they are debating other AI agents. The prompts are structured such that each  
218 debating model believes it is talking to a human.  
219

220 **Judged Debate:** We also experiment with an asynchronous debate setup with a judge. Figure 2  
221 (c) shows this setup. In this setup, models participate in an asynchronous debate as usual however,  
222 the final decision is made by a judge at the end of the debate. Models are incentivised to structure  
223 their arguments in a way that makes them most convincing to the judge. We structure this debate  
224 configuration similar to Khan et al. (2024), where the judge does not have access to all the external  
225 information and has to rely only on the debate transcript to decide the final answer.  
226

227 **Actor-Skeptic:** In this setup, only one agent, the *actor*, is tasked with deciding whether a given  
228 image-text pair is misinformation. The agent generates a response which a *skeptic* then evaluates.  
229 The skeptic is tasked with finding logical errors in the actor’s argument and asking follow-up ques-  
230 tions to disambiguate the actor’s response. It is important to note that neither the skeptic nor the  
231 actor has access to the ground truth. This setup does not benefit from an ensemble since both mod-  
232 els assume different roles and only one agent is tasked with generating the final answer.  
233

234 **Debate with Disambiguation:** Improving on the actor-skeptic method, in this setup, we allow  
235 all agents to act as actors *and* skeptics. Models are tasked with not only generating their own  
236 responses but also disambiguation queries to refine further or refute the other agents’ responses.  
237 These disambiguation queries are then used to search the Internet to obtain information to refine  
238 model outputs further. We are the first to propose this debate setup and while it does not achieve the  
239 best results in this work, we believe future research can greatly benefit by refining this setup further.

240 Through empirical testing of all the described debate set ups, we identify asynchronous debate—  
241 where the model believes it is debating against a human rather than an AI agent—as the most effec-  
242 tive configuration.  
243

### 244 3.2 PROMPT ENGINEERING

245  
246 The debate structure is substantiated through prompt engineering. Figure 1 shows that the first stage  
247 of our method requires for each AI agent to generate an independent response to whether the given  
248 image-text pair is misinformation. Each agent must take into account the external context related to  
249 the image obtained through the external information retrieval module. Specifications of the various  
250 prompts used in this work can be found in Appendix A.3. An initial prompt provides the agent with  
251 a summary of the news articles related to the image and, based on it, asks the agent to classify the  
252 image-text pair as misinformation or not. The prompt asks the agent to focus on certain details in the  
253 image, such as watermarks, flags, etc. We observe that images used in news articles often contain  
254 minute yet crucial details which can be used to inform the final decision about whether the image  
255 actually belongs to the news articles. Therefore, prompting the agent to pay special attention to these  
256 details further helps detect inconsistencies.

256  
257 Once the agents have formed independent opinions about the image-text pair, they must then partic-  
258 ipate in a debate. While the same prompt can be used for each debate round, we provide a different  
259 prompt for the first round to allow each agent to understand the changing nature of the conversation.  
260 Responses from other agents are provided as a part of the prompt, and the agent is asked to agree or  
261 disagree. The agent must also clearly state the reasoning process behind its argument. The prompt  
262 further requires the agent to identify ambiguities in the other agents’ reasoning. This allows the  
263 agent to closely analyse the responses from other agents and use them to inform its own decision.

263  
264 A separate prompt is then used to facilitate all rounds of the debate following round one. It takes into  
265 account the other agent’s suggestions from the previous round and, after refining its own response,  
266 asks the agent to point out any inconsistencies in this new response. We note that in such a scenario,  
267 agents are prone to simply agreeing with each other and repeating the other agent’s response. This is  
268 especially the case when agents believe they are conversing with a human. However, this tendency  
269 of the agents to easily give up their own opinions and simply agree with the other participants does  
not facilitate an advantageous debate, and the agents do not discover any new information related to  
the image-text pair. Therefore, the debate prompt also contains clear and explicit instructions asking

the agent not to simply agree with the presented response unless it has an acceptable reason to do so. We find this helps the agents develop stronger stances and reluctance to blind agreement.

### 3.3 EXTERNAL INFORMATION RETRIEVAL

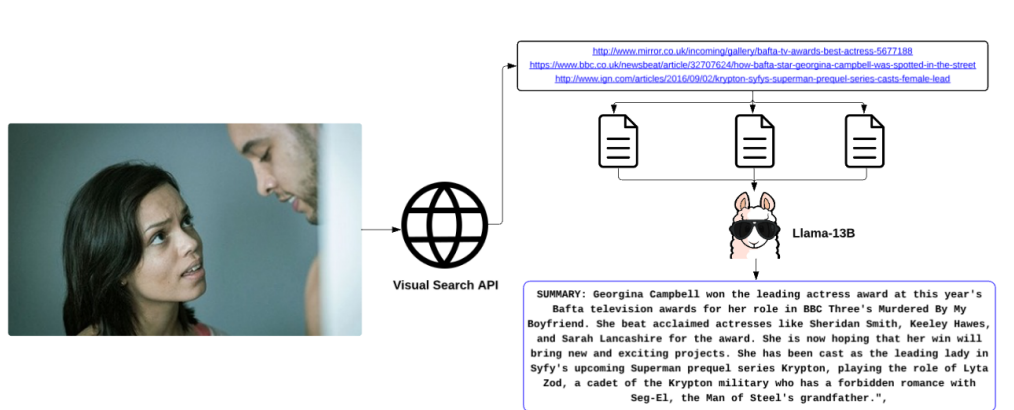


Figure 3: **Structure of the external information retrieval module:** We use the Bing Visual Search API (vis) to obtain web pages related to a given image, which are then summarised using Llama-13B (Touvron et al. (2023)). This summary is then passed to the debating agents as a part of the initial prompt.

Since a model’s world knowledge is limited to its training data (and hence a particular time frame), incorporating external retrieval allows the model to access information beyond this training data (and time frame). Previous work makes use of pre-existing external retrieval-based datasets (Abdelnabi et al. (2022)) to supplement external information related to an image-caption pair. However, we find this information lacking in detail since it is limited to the title of a news article. Agents can greatly benefit from the knowledge of the entire news article and its content rather than just the title when making a decision about whether a given image-caption pair, when considered in the context of the news article, is misinformation. To this end, we propose our own external information retrieval module. We observe a significant improvement in accuracy after incorporation of the external information retrieval module into the pipeline. The module is implemented in two stages:

#### 3.3.1 API-BASED INFORMATION RETRIEVAL

The Bing Visual Search API (vis) is used for the task of obtaining web pages related to a given image. A given image from the dataset is used to obtain a list of web pages completely and partially related to the image. We take the top three matching web pages in which the image appears. We believe these web pages contain sufficient information to allow the agent to develop a general understanding of the context in which the image is originally used. Since the community-accepted dataset for this task; NewsCLIPpings (Luo et al. (2021)), contains images from articles published more than ten years ago, some images do not result in any web pages or viable search results. In such a scenario, we simply do not pass any external context to the agent and only rely on the agent’s existing knowledge base. Since this is not the case for a significant fraction of examples in the dataset, it does not adversely affect system performance.

#### 3.3.2 SUMMARIZATION USING LLM

Once the top three web pages have been identified, we scrape the text from the web pages to obtain the textual information related to the context in which the image appears on the Internet. The compiled textual information is often too long to be passed directly to the agent, and hence, we use the Llama-13B (Touvron et al. (2023)) language model to summarize this information. The summaries obtained from the LLM only focus on the most important parts of the text and hence also allow agents to develop a more focused understanding of the external context. While this method of summarization works for most samples in our dataset, there are some examples where the obtained web

pages are not in the English language, and the LLM struggles with summarization. In this regard, we add an additional check that ignores text from web pages in languages other than English. While this restricts our system to the English language, it does not adversely affect system performance due to the distribution of the dataset, which consists of images mostly taken from English-language news articles. Multi-lingual support can be achieved by first translating text to English and then summarizing it.

### 3.4 COHERENT REASONING

All the different components of MAD-Sherlock are brought together in this stage of the pipeline. Each multi-modal agent is employed to participate in the best-debating set-up with the relevant prompts and is asked to detect a given image-text pair as misinformation and provide an explanation for the same. The agents also have access to external information related to the image through the external retrieval module. The final decision of the system is obtained once the debate terminates, which is after a certain number of debate rounds or after all agents converge to a common response, whichever is earlier.

## 4 EXPERIMENTS AND RESULTS

### 4.1 DATASET

We perform a series of experiments and report results on the NewsCLIPpings dataset (Luo et al. (2021)). The dataset is built based on the VisualNews (Liu et al. (2020)) dataset, which consists of image-caption pairs from four news agencies: BBC, USA Today, The Guardian and The Washington Post. The NewsCLIPpings dataset is created by generating OOC samples by replacing an image in one image-caption pair with a semantically related image from a different image-caption pair. CLIP (Radford et al. (2021)) is used to retrieve semantically similar images for a given caption. We report results on the Merged-Balanced version of the dataset, which has balanced proportions of all the retrieval strategies and positive/negative samples. The training, validation and test sets have 71,072, 7,024 and 7,264 samples, respectively.

### 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 is \$0.24 and it takes between 5 to 15 seconds to do so.

**Debate Setup:** We conduct experiments to select the best debating configuration using the LLaVA model (Liu et al. (2024b)). The experiments are carried out on a smaller subset containing 1000 test samples of the main NewsCLIPpings test dataset. 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 (vis) 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 (new). 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:** We 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)) and LLaVA (Liu et al. (2024b)). We also compare performance against GPT-4o (OpenAI & et al. (2024); OpenAI). 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 also show results for two baseline methods trained from scratch, namely EANN (Wang et al. (2018)) and SAFE (Zhou et al. (2020)). We further compare MAD-Sherlock to DT-Transformer (Papadopoulos et al. (2023)), CCN (Abdelnabi et al. (2022)), Sniffer (Qi et al. (2024)), VINVL (Huang et al. (2024)), SSDL (Mu et al. (2023)) and Neuro-Sym (Zhu et al. (2022)).



### 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 model, to select the best one for comparison with other works and further experimentation.

Debate Setup	Accuracy	Precision	Recall
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
<b>Async_Debate<sub>human</sub> (w external info)</b>	<b>86.2</b>	<b>82.6</b>	<b>90.6</b>
Actor-Skeptic	69.5	66.1	69.4
Judged Debate	66.7	66.7	61.5
Debate with Disambiguation	77.8	74.7	82.6

Table 1: **Performance comparison between different debate setups:** The Async\_Debate<sub>human</sub> where the model has external context and believes it is debating a human being is the best setup.

Table 1 shows that the Async\_Debate<sub>human</sub> setup where the agent has access to external information performs the best of all the debating configurations. We also report results for the Asynchronous Debate setup without access to external information to emphasise the importance of external information for the problem of misinformation detection in the news domain. The external retrieval of information significantly boosts performance. All following debate set-ups, therefore, use external information related to the image-caption pair as a part of their initial prompt. We also observe a significant performance increase when the agent believes it is conversing with a human instead of another AI agent. Qualitatively, the agent considers the other agent’s responses more critically and with more seriousness when it believes that the agent is a human. Further, the Asynchronous Debate setup benefits from the ensemble of agents which is not present in the actor-skeptic setup, where only one agent is responsible for generating the responses. The generation of disambiguation queries within the same response, confuses the agents and even deviates them from their own chain of thought. We believe this accounts for the counter-intuitively performance of this method where agents perform worse with more information. The judged debate setup focuses on enforcing agents to structure their responses in a way that will convince a judge. The agents also debate with opposite stances and do not have the option of changing their stance mid-debate. This can further confuse the judge and lead to incorrect decisions. This is resolved in the Async\_Debate<sub>human</sub> set up where agents are given complete freedom over their initial opinions, as well as their opinions during the debate. If they believe they are convinced by the other agents’ arguments, they can choose to change their response and the debate ends. Based on the results from Table 1, we choose the best-performing debate set-up, i.e. Async\_Debate<sub>human</sub> (with external information) as the debate configuration for further experimentation and comparison.

#### 4.3.2 PERFORMANCE COMPARISON

We present our results on the NewsCLippings dataset against existing out-of-context detection methods discussed in section 4.2.

Table 2 shows the comparison between our system and existing methods. We report state-of-the-art performance when using our proposed debate configuration with the GPT-4o (OpenAI & et al. (2024); OpenAI) model. Sniffer (Qi et al., 2024), being the only work comparable in performance to ours, is finetuned extensively to adapt it to the NewsCLippings dataset. While we do not provide a quantitative assessment of explanations by MAD-Sherlock, we do believe our system produces more coherent, detailed and comprehensive explanations when compared to other baselines. This is attributed to the fact that in a multi-agent setup, we have multiple context windows which leads to more coherent and relevant final explanations. We leave the detailed analysis of these explanations and the development of the associated metrics as future work. We also note that the debate paradigm



	Model	Accuracy $\uparrow$
432		
433		
434	SAFE	50.7
435	EANN	58.1
436	VisualBERT	54.8
437	CLIP	62.6
438	InstructBLIP	48.6
439	LLaVA	57.1
440	GPT-4o	70.7
441	DT-Transformer	77.1
442	CCN	84.7
443	SSDL	65.6
444	VINVL	65.4
444	Neuro-Sym	68.2
445	GPT-4o <sup>#</sup> (w internet access)	86.00
446	Sniffer (w finetuning)	88.4
447	Sniffer (w/o finetuning)	84.5
448	<b>MAD-Sherlock (ours)</b>	<b>90.17</b>

449 Table 2: **Performance comparison between our model and baselines:** MAD-Sherlock (with GPT-  
450 4o) out performs all related work. Note: the GPT-4o<sup>#</sup> setup is identical to MAD-Sherlock with the  
451 absence of a multi-agent debate, here only a single agent which has access to external information  
452 is considered and the results are reported on a smaller heldout test set of 1000 samples.

453  
454  
455 in itself is essential to the system performance. We observe a drop in performance and quality of  
456 explanations when using an identical system configuration but with a single model.

457 We also note that single multi-modal models, including VisualBERT, CLIP, InstructBLIP, LLaVA  
458 and GPT-4o do not perform at par with other related work. This can be attributed to the neces-  
459 sity for external context for misinformation detection in the news domain and the lack of diverse  
460 perspectives that arise naturally in a multi-agent framework. Therefore, these standalone models,  
461 while promising, currently are unable to detect misinformation effectively. These models require ad-  
462 ditional integration into more comprehensive pipelines, as done in this work. In line with previous  
463 work, we also note that baselines trained from scratch, such as SAFE (Zhou et al. (2020)) and EANN  
464 (Wang et al. (2018)) perform worse than pretrained multi-modal models. This further concretizes  
465 the fact that image-based OOC detection in the news domain requires strong world knowledge as  
466 well as advanced multi-modal reasoning capabilities.

## 467 5 USER STUDY

468  
469  
470 We conducted a user study to evaluate the effectiveness of our system in detecting and explaining  
471 misinformation. While it is easy to quantify model performance in terms of misinformation detec-  
472 tion, there are no effective metrics to assess the quality of the explanations generated by the model.  
473 Therefore, in order to perform a thorough analysis of the system performance, a user study is essen-  
474 tial. For a deeper analysis we further grouped the participants based on their profession into three  
475 groups, namely: Journalists, AI Academics (studying AI) and Others. Further details regarding the  
476 study setup and participant groups can be found in Appendix A.4.

477 In the study, participants were shown ten image-text pairs and were asked to decide if the image and  
478 caption, when considered together, were misinformation or not. They were also asked to provide  
479 a confidence rating for their answer on a scale of 0-10, with 10 being the highest confidence level.  
480 For each image-text pair, after the participants provided their initial answers, they were shown AI  
481 insights about the same image-text pair. These AI insights were the final output explanations from  
482 MAD-Sherlock. Participants were then asked to reconsider their answers in light of the new infor-  
483 mation from the AI agent. Table 3 shows that average system performance is better than the average  
484 human performance for both cases where the participants have access to AI insights and where they  
485 do not. Therefore, MAD-Sherlock can be used as a reliable assistive tool for OSINT research for  
detecting and explaining misinformation with little or no human intervention.

Study Setup	Average Accuracy $\uparrow$
Humans	60.3 $\pm$ 13.5
Humans+MAD-Sherlock	76.7 $\pm$ 12.2
<b>MAD-Sherlock</b>	<b>80.0 <math>\pm</math> 0.0</b>

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

We further observe that the average human accuracy for the misinformation detection task increases by more than 27% concretizing the fact that AI insights from our model do actually improve human efficiency in detecting misinformation. We also observe interesting patterns for group-wise analysis which we believe would be valuable for future work. Table 4 shows that the performance of all groups improves significantly and is not far off from that of professional journalists. The average confidence level (out of 10) is comparable across all the groups before and after considering MAD-Sherlock insights and generally increases. Therefore, we conclude that MAD-Sherlock can significantly uplift non-expert performance and hence can be useful in citizen intelligence applications.

Group	Avg acc $\uparrow$ (only human)	Avg conf $\uparrow$ (only human)	Avg acc $\uparrow$ (with MAD-Sherlock)	Avg conf $\uparrow$ (with MAD-Sherlock)
Journalists	70.0 $\pm$ 1.4	4.3 $\pm$ 2.1	82.2 $\pm$ 0.9	5.3 $\pm$ 1.3
AI Academics	60.7 $\pm$ 1.4	3.2 $\pm$ 0.8	79.3 $\pm$ 1.3	5.8 $\pm$ 1.4
Others	56.7 $\pm$ 1.5	3.9 $\pm$ 1.2	71.7 $\pm$ 1.1	5.8 $\pm$ 1.4

Table 4: **Performance comparison between different participant groups:** All groups show performance improvement with MAD-Sherlock. AI Academics are able to perform nearly at par with professional journalists after considering insights from MAD-Sherlock.

## 6 CONCLUSION AND FUTURE WORK

Misinformation detection has become a pressing issue in recent times. With the ever-advancing capabilities of vision and language models, the detection of OOC image use has become a very difficult task. In this work, we explore the question of whether it is possible for multiple AI agents to pool their contextual knowledge and converge to a common prediction in order to identify instances of misinformation. We identify `Asynchronous_Debatehuman` as the most optimal communication setup for AI models. We observe significant performance improvement when the models believe they are debating against a human instead of another AI agent. We observe that in this setup, models tend to be more involved and open to changing their opinions. Our method also allows for agents to have freedom of opinion which they may change mid-debate. Agents in such a setting show enhanced abilities to critically evaluate an argument and pick up on minute inconsistencies.

Our final system, MAD-Sherlock, achieves state-of-the-art performance on the misinformation detection task. Further, owing to our advanced external retrieval module, MAD-Sherlock provides clear, coherent and detailed explanations. As a result, MAD-Sherlock significantly improves the OOC misinformation detection performance of both human experts, and non-experts.

We identify several promising avenues for future research in this field. The research community would benefit from a continuously updated benchmark dataset, incorporating more recent news articles and subtler inconsistencies. A direct extension of this work involves applying our methods to video-text pairs and supporting multi-lingual content. Future extensions of this work could further validate our findings by leveraging more advanced and refined models in the summarization pipeline which we believe would further improve system performance. It is also worth comparing MAD-Sherlock to systems using multi-agent collaboration with external information retrieval.

Finally, while we conducted extensive user studies with MAD-Sherlock, deploying it on a larger scale in professional environments and within the citizen intelligence community will provide valuable insights into its real-world performance, uncovering new opportunities for improvement. For an analysis of limitations, please refer to Appendix A.1.

## REFERENCES

- 540  
541  
542 Deepfakes, explained — MIT Sloan — mitsloan.mit.edu. [https://mitsloan.mit.edu/](https://mitsloan.mit.edu/ideas-made-to-matter/deepfakes-explained)  
543 [ideas-made-to-matter/deepfakes-explained](https://mitsloan.mit.edu/ideas-made-to-matter/deepfakes-explained). [Accessed 28-09-2024].
- 544  
545 Newspaper3k: Article scraping & curation — newspaper 0.0.2 documentation. [https://](https://newspaper.readthedocs.io/en/latest/)  
546 [newspaper.readthedocs.io/en/latest/](https://newspaper.readthedocs.io/en/latest/).
- 547 Out-of-context photos are a powerful low-tech form of misinforma-  
548 tion — pbs.org. [https://www.pbs.org/newshour/science/](https://www.pbs.org/newshour/science/out-of-context-photos-are-a-powerful-low-tech-form-of-misinformation)  
549 [out-of-context-photos-are-a-powerful-low-tech-form-of-misinformation](https://www.pbs.org/newshour/science/out-of-context-photos-are-a-powerful-low-tech-form-of-misinformation).  
550 [Accessed 28-09-2024].
- 551  
552 Visual Search API — Microsoft Bing. [https://www.microsoft.com/en-us/bing/](https://www.microsoft.com/en-us/bing/apis/bing-visual-search-api)  
553 [apis/bing-visual-search-api](https://www.microsoft.com/en-us/bing/apis/bing-visual-search-api).
- 554 Sahar Abdelnabi, Rakibul Hasan, and Mario Fritz. Open-domain, content-based, multi-modal fact-  
555 checking of out-of-context images via online resources, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2112.00061)  
556 [abs/2112.00061](https://arxiv.org/abs/2112.00061).
- 557  
558 Irene Amerini, Aris Anagnostopoulos, Luca Maiano, Lorenzo Ricciardi Celsi, et al. Deep learning  
559 for multimedia forensics. *Foundations and Trends® in Computer Graphics and Vision*, 12(4):  
560 309–457, 2021.
- 561 Shivangi Aneja, Cise Midoglu, Duc-Tien Dang-Nguyen, Sohail Ahmed Khan, Michael Riegler, Pål  
562 Halvorsen, Chris Bregler, and Balu Adsumilli. Acm multimedia grand challenge on detecting  
563 cheapfakes, 2022. URL <https://arxiv.org/abs/2207.14534>.
- 564  
565 Kevin Aslett, Zeve Sanderson, William Godel, Nathaniel Persily, Jonathan Nagler, and Joshua A  
566 Tucker. Online searches to evaluate misinformation can increase its perceived veracity. *Nature*,  
567 625(7995):548–556, 2024.
- 568 Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng  
569 Shou. Hallucination of multimodal large language models: A survey, 2024. URL [https://](https://arxiv.org/abs/2404.18930)  
570 [arxiv.org/abs/2404.18930](https://arxiv.org/abs/2404.18930).
- 571  
572 Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin  
573 Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. Abstract Meaning Representa-  
574 tion for sembanking. In Antonio Pareja-Lora, Maria Liakata, and Stefanie Dipper (eds.),  
575 *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*,  
576 pp. 178–186, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL  
577 <https://aclanthology.org/W13-2322>.
- 578  
579 Nadia M Brashier and Elizabeth J Marsh. Judging truth. *Annual review of psychology*, 71(1):499–  
580 515, 2020.
- 581  
582 Ivan Castillo Camacho and Kai Wang. A comprehensive review of deep-learning-based methods for  
583 image forensics. *Journal of imaging*, 7(4):69, 2021.
- 584  
585 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,  
586 Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language  
587 models with instruction tuning, 2023. URL <https://arxiv.org/abs/2305.06500>.
- 588  
589 Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving  
590 factuality and reasoning in language models through multiagent debate, 2023a. URL [https://](https://arxiv.org/abs/2305.14325)  
591 [arxiv.org/abs/2305.14325](https://arxiv.org/abs/2305.14325).
- 592  
593 Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving  
594 Factuality and Reasoning in Language Models through Multiagent Debate. October 2023b. URL  
595 <https://openreview.net/forum?id=QAwaaLJNck>.
- 596  
597 Hany Farid. *Photo Forensics*. The MIT Press, 2016. ISBN 0262035340.

- 594 Lynn Hasher, David Goldstein, and Thomas Toppino. Frequency and the conference of referential  
595 validity. *Journal of verbal learning and verbal behavior*, 16(1):107–112, 1977.  
596
- 597 Arash Heidari, Nima Jafari Navimipour, Hasan Dag, and Mehmet Unal. Deepfake detection using  
598 deep learning methods: A systematic and comprehensive review. *Wiley Interdisciplinary Reviews:  
599 Data Mining and Knowledge Discovery*, 14(2):e1520, 2024.
- 600 Maryam Hina, Mohsin Ali, Abdul Rehman Javed, Fahad Ghabban, Liaqat Ali Khan, and Zunera  
601 Jalil. Sefaced: Semantic-based forensic analysis and classification of e-mail data using deep  
602 learning. *IEEE Access*, 9:98398–98411, 2021.
- 603 Mingzhen Huang, Shan Jia, Zhou Zhou, Yan Ju, Jialing Cai, and Siwei Lyu. Exposing text-image  
604 inconsistency using diffusion models. In *The Twelfth International Conference on Learning Rep-  
605 resentations*, 2024.
- 607 Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward  
608 Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with more per-  
609 suasive llms leads to more truthful answers, 2024. URL [https://arxiv.org/abs/2402.  
610 06782](https://arxiv.org/abs/2402.06782).
- 611 Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem.  
612 Camel: Communicative agents for "mind" exploration of large language model society, 2023a.  
613 URL <https://arxiv.org/abs/2303.17760>.
- 614 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image  
615 pre-training with frozen image encoders and large language models. In Andreas Krause, Emma  
616 Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.),  
617 *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Pro-  
618 ceedings of Machine Learning Research*, pp. 19730–19742. PMLR, 23–29 Jul 2023b. URL  
619 <https://proceedings.mlr.press/v202/li23q.html>.
- 621 Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple  
622 and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- 623 Hongzhan Lin, Ziyang Luo, Wei Gao, Jing Ma, Bo Wang, and Ruichao Yang. Towards explainable  
624 harmful meme detection through multimodal debate between large language models, 2024. URL  
625 <https://arxiv.org/abs/2401.13298>.
- 626 Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visualnews : Benchmark and  
627 challenges in entity-aware image captioning, 2020.
- 629 Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating  
630 hallucination in large multi-modal models via robust instruction tuning, 2024a. URL <https://arxiv.org/abs/2306.14565>.
- 632 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.  
633 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL [https://  
634 llava-vl.github.io/blog/2024-01-30-llava-next/](https://llava-vl.github.io/blog/2024-01-30-llava-next/).
- 636 Grace Luo, Trevor Darrell, and Anna Rohrbach. Newsclippings: Automatic generation of out-of-  
637 context multimodal media. *arXiv:2104.05893*, 2021.
- 638 Momina Masood, Mariam Nawaz, Khalid Mahmood Malik, Ali Javed, Aun Irtaza, and Hafiz Malik.  
639 Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way  
640 forward. *Applied intelligence*, 53(4):3974–4026, 2023.
- 641 Marvin Minsky. *Society of mind*. Simon and Schuster, 1988.
- 642 Michael Mu, Sreyasee Das Bhattacharjee, and Junsong Yuan. Self-supervised distilled learning for  
643 multi-modal misinformation identification. In *Proceedings of the IEEE/CVF Winter Conference  
644 on Applications of Computer Vision (WACV)*, pp. 2819–2828, January 2023.
- 646 OpenAI. GPT-4o. <https://openai.com/index/hello-gpt-4o/>. [Accessed 28-08-  
647 2024].

- 648 OpenAI and Josh Achiam et al. GPT-4 Technical Report, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2303.08774)  
649 [abs/2303.08774](https://arxiv.org/abs/2303.08774).  
650
- 651 Stefanos-Iordanis Papadopoulos, Christos Koutlis, Symeon Papadopoulos, and Panagiotis Petran-  
652 tonakis. Synthetic misinformers: Generating and combating multimodal misinformation. In  
653 *Proceedings of the 2nd ACM International Workshop on Multimedia AI against Disinforma-*  
654 *tion*, MAD '23, pp. 36–44, New York, NY, USA, 2023. Association for Computing Machin-  
655 ery. ISBN 9798400701870. doi: 10.1145/3592572.3592842. URL [https://doi.org/10.](https://doi.org/10.1145/3592572.3592842)  
656 [1145/3592572.3592842](https://doi.org/10.1145/3592572.3592842).
- 657 Peng Qi, Zehong Yan, Wynne Hsu, and Mong Li Lee. Sniffer: Multimodal large language model  
658 for explainable out-of-context misinformation detection, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2403.03170)  
659 [abs/2403.03170](https://arxiv.org/abs/2403.03170).
- 660 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
661 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
662 models from natural language supervision. In *International conference on machine learning*, pp.  
663 8748–8763. PMLR, 2021.
- 664 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
665 resolution image synthesis with latent diffusion models, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2112.10752)  
666 [abs/2112.10752](https://arxiv.org/abs/2112.10752).
- 667 Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip H. S.  
668 Torr, Mingfei Sun, and Shimon Whiteson. Is Independent Learning All You Need in the Star-  
669 Craft Multi-Agent Challenge?, November 2020. URL [https://arxiv.org/abs/2011.](https://arxiv.org/abs/2011.09533v1)  
670 [09533v1](https://arxiv.org/abs/2011.09533v1).
- 671 Fatma Shalabi, Huy H. Nguyen, Hichem Felouat, Ching-Chun Chang, and Isao Echizen. Image-text  
672 out-of-context detection using synthetic multimodal misinformation. In *2023 Asia Pacific Signal*  
673 *and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE,  
674 October 2023. doi: 10.1109/apsipaasc58517.2023.10317336. URL [http://dx.doi.org/](http://dx.doi.org/10.1109/APSIPAASC58517.2023.10317336)  
675 [10.1109/APSIPAASC58517.2023.10317336](http://dx.doi.org/10.1109/APSIPAASC58517.2023.10317336).  
676
- 677 Mubashir Sultan, Alan N Tump, Michael Geers, Philipp Lorenz-Spreen, Stefan M Herzog, and  
678 Ralf HJM Kurvers. Time pressure reduces misinformation discrimination ability but does not  
679 alter response bias. *Scientific Reports*, 12(1):22416, 2022.
- 680 Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-  
681 Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information*  
682 *Fusion*, 64:131–148, 2020.
- 683 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
684 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
685 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 686 Sheng-Yu Wang, Oliver Wang, Andrew Owens, Richard Zhang, and Alexei A. Efros. Detecting  
687 photoshopped faces by scripting photoshop. In *Proceedings of the IEEE/CVF International Con-*  
688 *ference on Computer Vision (ICCV)*, October 2019.
- 689 Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing  
690 Gao. Eann: Event adversarial neural networks for multi-modal fake news detection. In *Pro-*  
691 *ceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data*  
692 *Mining*, KDD '18, pp. 849–857, New York, NY, USA, 2018. Association for Computing Machin-  
693 ery. ISBN 9781450355520. doi: 10.1145/3219819.3219903. URL [https://doi.org/10.](https://doi.org/10.1145/3219819.3219903)  
694 [1145/3219819.3219903](https://doi.org/10.1145/3219819.3219903).  
695
- 696 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
697 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
698 models. In *Proceedings of the 36th International Conference on Neural Information Processing*  
699 *Systems*, NIPS '22, pp. 24824–24837, Red Hook, NY, USA, April 2024. Curran Associates Inc.  
700 ISBN 978-1-71387-108-8.  
701

702 Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Joshua B. Tenenbaum.  
703 Neural-symbolic vqa: Disentangling reasoning from vision and language understanding, 2019.  
704 URL <https://arxiv.org/abs/1810.02338>.

706  
707 Yizhou Zhang, Loc Trinh, Defu Cao, Zijun Cui, and Yan Liu. Interpretable detection of out-of-  
708 context misinformation with neural-symbolic-enhanced large multimodal model, 2024. URL  
709 <https://arxiv.org/abs/2304.07633>.

711  
712 Xinyi Zhou, Jindi Wu, and Reza Zafarani. Safe: Similarity-aware multi-modal fake news detection,  
713 2020. URL <https://arxiv.org/abs/2003.04981>.

715  
716 Wang Zhu, Jesse Thomason, and Robin Jia. Generalization differences between end-to-end and  
717 neuro-symbolic vision-language reasoning systems, 2022. URL <https://arxiv.org/abs/2210.15037>.

719  
720 Xinshan Zhu, Yongjun Qian, Xianfeng Zhao, Biao Sun, and Ya Sun. A deep learning approach to  
721 patch-based image inpainting forensics. *Signal Processing: Image Communication*, 67:90–99,  
722 2018.

## 723 A APPENDIX

### 724 A.1 LIMITATIONS

725  
726 Despite the strong performance of MAD-Sherlock, several limitations remain. First, while our  
727 model excels at detecting out-of-context image-text pairs, its reliance on external retrieval can lead  
728 to reduced accuracy when relevant context is unavailable or difficult to retrieve. Second, the quality  
729 of explanations is constrained to textual outputs, limiting multi-modal explanation capabilities such  
730 as image or video integration. Third, the system’s performance is sensitive to hyperparameter tun-  
731 ing, including the number of debate rounds and agents, which may require further optimization for  
732 broader use cases.

733  
734 Additionally, while our user studies provided valuable insights, large-scale deployment in diverse,  
735 real-world settings, such as professional or citizen intelligence environments, is necessary to fully  
736 assess the method’s robustness and scalability. Finally, our dataset, though comprehensive, primarily  
737 focuses on English-language news, limiting the generalizability of the system across non-English  
738 contexts.

739  
740 Another important limitation is the potential risk that open-sourcing MAD-Sherlock might allow  
741 adversaries to train models specifically designed to counter or evade detection by our system. As  
742 adversarial actors gain access to the source code, they could exploit its known strengths and weak-  
743 nesses to develop countermeasures that diminish its effectiveness. However, despite these risks, we  
744 believe that open-sourcing remains the right path forward. Open-sourcing encourages transparency,  
745 collaboration, and rapid innovation, enabling the broader community to contribute improvements,  
746 detect vulnerabilities, and build on the system.

747  
748 Moreover, by engaging the community, we can foster the development of more resilient and adap-  
749 tive models that evolve in response to emerging adversarial techniques, thus maintaining MAD-  
750 Sherlock’s effectiveness in the long term. The collective strength of a diverse, open-source commu-  
751 nity can outweigh the potential threats posed by adversarial exploitation.

752  
753 Future work will need to address these limitations to enhance the practical utility, robustness, and  
754 long-term resilience of MAD-Sherlock.

756 A.2 SAMPLE IMAGE-CAPTION PAIR IN THE NEWS DOMAIN  
757  
758  
759  
760



761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771 Figure 4: Russian President Vladimir Putin has called Ukraine’s move into Kursk a “major provo-  
772 cation”. Image and caption taken from the BBC article here (Accessed at 17:43 on Aug 11, 2024):  
773 <https://www.bbc.co.uk/news/articles/cze5pkg5jwlo>  
774  
775  
776

777 A.3 PROMPTS FOR MAD-SHERLOCK  
778  
779  
780  
781

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

795 Figure 5: Initial prompt for independent opinion formation and response generation  
796  
797  
798  
799  
800

801  
802 This is what I think: {}.  
803 Do you agree with me?  
804 If you think I am wrong then convince me why you are correct.  
805 Clearly state your reasoning and tell me if I am missing out on  
806 some important information or am making some logical error.  
807 Do not describe the image.  
808 At the end give a definite YES or NO answer to this question:  
IS THIS MISINFORMATION?  
809

Figure 6: Prompt for Debate Round 1



```

810 I see what you mean and this is what I think: {}.
811 Do you agree with me?
812 If not then point out the inconsistencies in my argument (e.g.
813 location, time or person related logical confusion) and explain
814 why you are correct.
815 If you disagree with me then clearly state why and what
816 information I am overlooking.
817 Find disambiguation in my answer if any and ask questions to
818 resolve them.
819 I want you to help me improve my argument and explanation.
820 Don't give up your original opinion without clear reasons, DO NOT
821 simply agree with me without proper reasoning.
822 At the end give a definite YES or NO answer to this question:
823 IS THIS MISINFORMATION?

```

Figure 7: Prompt for Debate after Round 1

#### 826 A.4 USER STUDY

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

##### 831 A.4.1 SETUP

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

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

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

860 We analyse two special cases, where MAD-Sherlock argues for the wrong answer. We include these  
861 results in order to observe how persuasive our system can be even when it is wrong. We note in the  
862 instance where the image-text pair was actually misinformation and the model argued that it was  
863 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

864 can appear even when it is wrong. While the case of false negatives is important, false positives  
865 are an even more concerning matter for our problem statement. In the case where MAD-Sherlock  
866 declared the given image-text pair to be misinformation when it was not, is important to analyse. In  
867 this setting 50% of the total participants changed their answer to the wrong one, therefore believing  
868 a piece of true information to be false. In some cases where participants chose the wrong response  
869 to begin with, their confidence in the response further increased after considering insights from the  
870 system. Finally, 4 participants did not change their answer to the wrong one after considering AI  
871 insights but their confidence in their response decreased.

872 The average time taken to complete the study was 12 minutes and 57 seconds. The average partici-  
873 pant was therefore able to go through 10 image-text pairs and decide if they were misinformation or  
874 not in under 13 minutes. The same task without AI insights would require extensive analysis and we  
875 project it would take between 30-45 minutes to decide if 10 image-text pairs were misinformation.  
876

#### 877 A.5 MULTI-MODAL DEBATES FOR HARMFUL MEME DETECTION

878 While this work relates to a different problem than OOC misinformation detection in the news do-  
879 main, we still find the approach taken by the authors a relevant related work and therefore include it  
880 here. Lin et al. (2024) use LMMs debating against each other to generate explanations for contra-  
881 dictory arguments regarding whether a given meme is harmful. These explanations are then used to  
882 train a small language model as a judge to determine whether the image and text that make up the  
883 meme are actually harmful. This work does not allow agents to have flexibility of opinion. There are  
884 always two agents, and each one is provided a stance to defend. Moreover, a judge decides the final  
885 outcome of the debate and needs to be trained on data from the debate. This method also does not  
886 benefit from external retrieval, and therefore, the debating agents are not aware of the crucial exter-  
887 nal context related to the input. Finally, this work is related to harmful *meme* detection and does not  
888 concern the problem of misinformation detection in the news domain, which likely requires more  
889 intricate contextual analysis, including of external context.  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917