# LLM-Consensus: Multi-Agent Debate for Visual Misinformation Detection

**Kumud Lakara** [1]  **Georgia Channing** [1]  **Juil Sock** [2]  **Christian Rupprecht** [1]  **Philip Torr** [1]  **John Collomosse** [3]
**Christian Schroeder de Witt** [1]

## 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 LLM-Consensus, a multi-agent debate system for OOC misinformation detection. LLM-Consensus 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 LLM-Consensus boosts performance for both experts and non-experts. These results position LLM-Consensus 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.

OOC misinformation detection requires a nuanced understanding of the relationship between text and images and the ability to identify when they are misaligned. Detecting these subtle inconsistencies is time-consuming for humans, and (Sultan et al., 2022) shows that time pressure further reduces detection accuracy, exacerbating scalability issues.

As a result, AI-driven tools have gained attention for scaling detection efforts. However, 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), designed to detect manipulations like Photo-Shop editing (Tolosana et al., 2020; Masood et al., 2023; Farid, 2016; Wang et al., 2019) and AI-generated Deepfakes (mit), focus on spotting artifacts from tampering. In contrast, OOC detection requires cross-contextual reasoning, as the deception stems from the misalignment of legitimate images with misleading text.

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.

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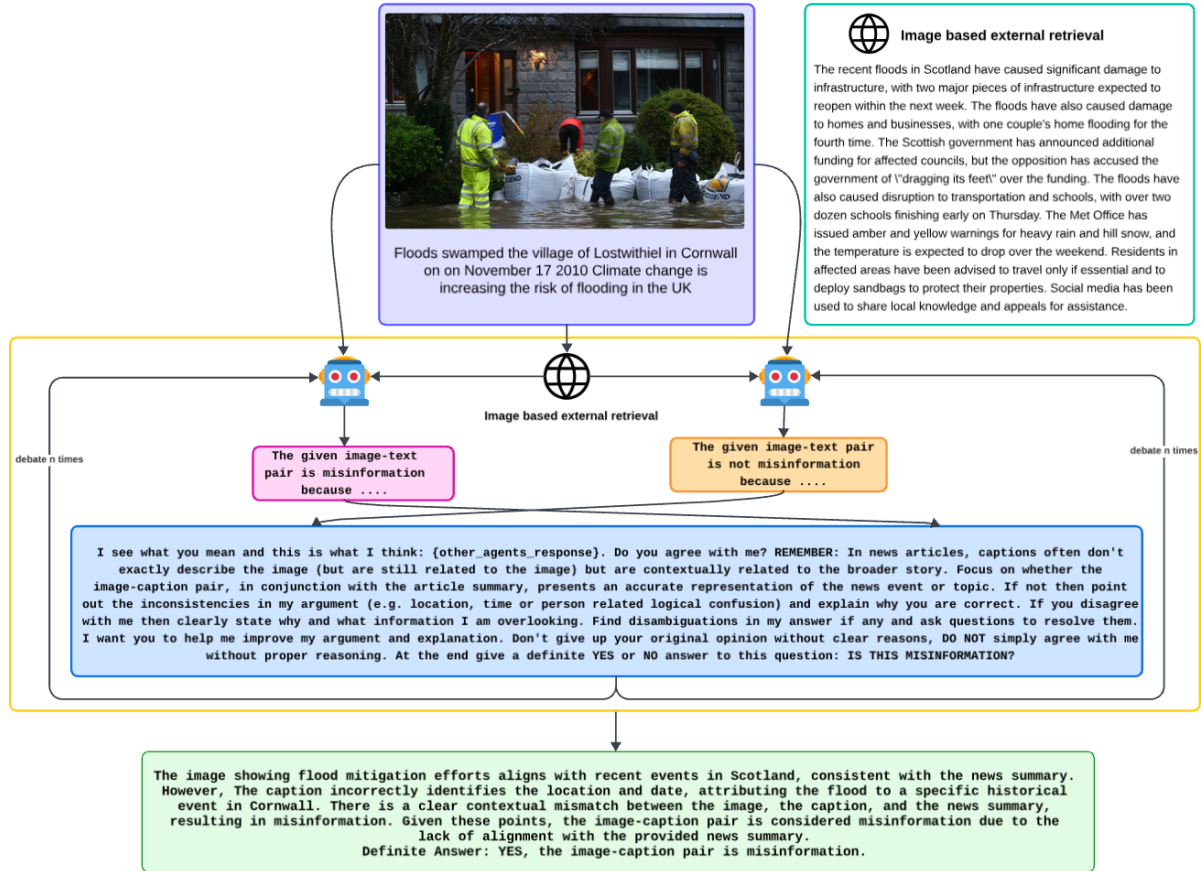
*Figure 1.* **Overview of LLM-Consensus:** Two or more independent agents see the same image-text input and are tasked with detecting whether the input is misinformation or not. After the agents form their independent opinions, they participate in a debate until they converge on the same response or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

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 LLM-Consensus 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 LLM-Consensus 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.

We perform a comprehensive empirical evaluation of our method (see Section 4), including the study of multiple debate configurations. To optimise OpenAI API use, we utilize an experimental pipeline where preliminary experiments are performed using the open-source LLaVA model (Liu et al., 2024b), which is only later replaced with GPT-4o (OpenAI) to achieve state-of-the-art performance. We find that LLM-Consensus outperforms both prior work and novel baselines that we introduce, is more robust to various failure modes, and produces coherent explanations that help both human experts and non-experts significantly improve their detection accuracy in user studies. We identify both access to external information retrieval and complete freedom of opinion as key ingredients to the performance of LLM-Consensus. Finally, we discuss current limitations

of our method and propose future work toward overcoming the scalability challenges in large-scale online OOC misinformation detection.

## 2. Related Work

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

(Abdelnabi et al., 2022) introduce the Consistency Checking Network (CCN), which models aspects of human reasoning across modalities for misinformation detection using external evidence aggregated from the Internet. The CCN employs memory networks to evaluate the consistency of image-caption pairs against retrieved evidence and a CLIP component (Radford et al., 2021) to assess the alignment between the image and caption. The use of external evidence enhances CCN's classification performance compared to other methods. However, it lacks an explainability component, providing only binary decisions without further explanation.

(Zhang et al., 2024) extend the neural symbolic method (Yi et al., 2019; Zhu et al., 2022) to develop an interpretable cross-modal misinformation detection model that provides supporting evidence for its predictions. They use symbolic graphs based on Abstract Meaning Representation (Banarescu et al., 2013) to detect out-of-context image use. Similarly, (Zhou et al., 2020) introduce Similarity Aware Fake news detection (SAFE), where neural networks jointly learn and analyze features and relationships between text and visual news representations to predict fake news. (Wang et al., 2018) introduce EANN: Event Adversarial Neural Networks to derive event invariant features which can be used to detect fake news that has recently been generated. EANN uses adversarial training to learn multi-modal features 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.

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

Sniffer (Qi et al., 2024) is the closest to our work. It uses the InstructBLIP (Dai et al., 2023) model 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 LLM-Consensus, 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. 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.

**Asynchronous Debate (not) against Human:**   One of the core setups we test in our experiments is an asynchronous debating strategy, where models wait for the responses of other participants before generating their own. Figure 2 (a) and (b) illustrate the differences between synchronous and asynchronous debate structures. While synchronous debates, where all participants respond simultaneously, can be faster and more computationally efficient, we find the asynchronous setup more effective. This structure allows models to better identify contextual ambiguities in a more organized and systematic manner, which is crucial for misinformation detection. A key aspect of this setup is the design of model prompts: the debating models are not aware that they are interacting with other AI agents but are instead led to believe they are debating with a human.

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

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

**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.

## 3.2. Prompt Engineering

The debate structure is enabled through prompt engineering. As shown in Figure 1, the first stage involves each AI agent independently determining whether the image-text pair is misinformation, using context from the external retrieval module (details in Appendix A.3). The initial prompt summarizes related news articles and instructs the agent to focus on key image details, like watermarks and flags, to identify inconsistencies.

After forming initial opinions, the agents engage in a debate. The first round uses a unique prompt to help agents adapt to the discussion, while subsequent rounds use a standard prompt incorporating suggestions from previous rounds. Agents must agree or disagree with others' responses while identifying ambiguities and refining their reasoning.

To prevent blind agreement, the prompt explicitly instructs agents to provide valid reasoning before accepting others' responses. This promotes stronger, independent stances and enhances the discovery of new information.

## 3.3. External Information Retrieval

A model's world knowledge is limited by its training data and time frame, but external retrieval allows access to up-to-date information beyond these constraints. Previous work relies on external datasets (Abdelnabi et al., 2022), but these often provide only news article titles, which lack sufficient detail. 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. Full access to the article content can significantly enhance an agent's ability to assess whether an image-caption pair constitutes misinformation. To address this, we propose a dedicated external information retrieval module, which improves accuracy when integrated into the pipeline. The module operates in two stages:

### 3.3.1. API-BASED INFORMATION RETRIEVAL

The Bing Visual Search API is used for the task of obtaining web pages related to a given image (vis) . 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
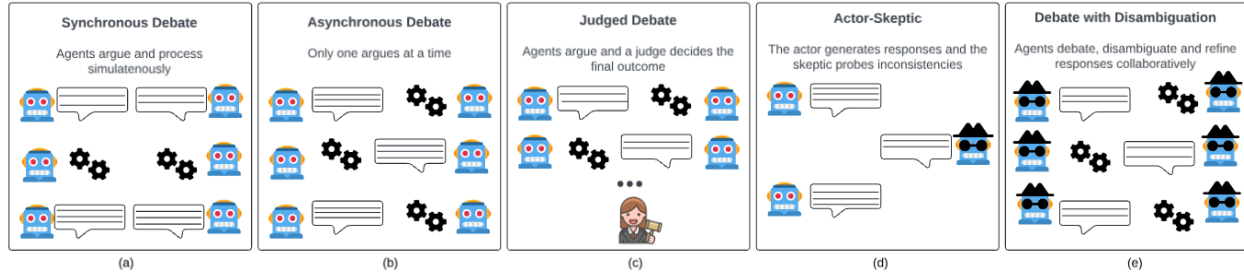
*Figure 2.* **Debating Strategies:** We experiment with multiple debating strategies. The asynchronous debate setup where agents argue one after the other and take turns presenting their arguments is the best configuration.
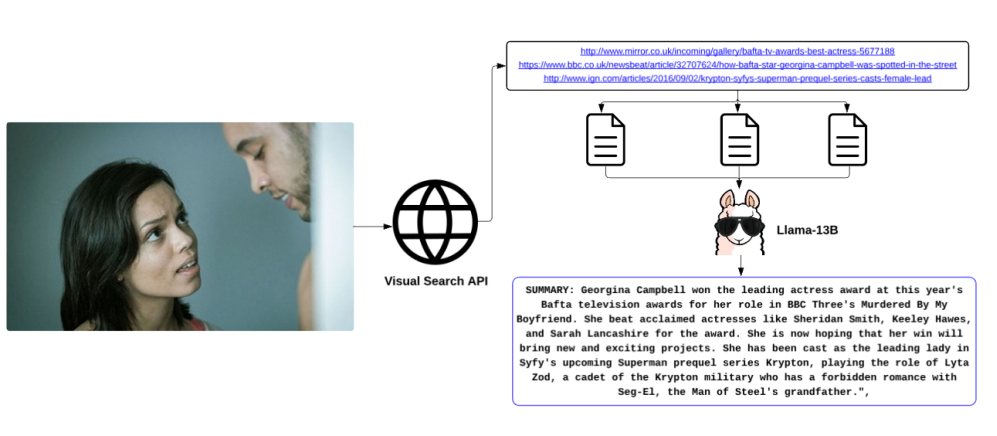


*Figure 3.* **Structure of the external information retrieval module**: We use the Bing Visual Search API (vis) to obtain web pages related to a given image, which are then summarised using Llama-13B (Touvron et al., 2023). This summary is then passed to the debating agents as a part of the initial prompt.

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 relevant for only a very small number of examples in the dataset, it does not adversely affect system performance.

### 3.3.2. SUMMARIZATION USING LLM

Once the top three web pages are identified, we scrape their text to gather contextual information about the image. This textual data is often too lengthy for direct input to the agent, so we use the Llama-13B (Touvron et al., 2023) model to generate concise summaries, focusing on key details to help the agent better understand the external context. For non-English web pages, the LLM struggles with summarization,

so we add a filter to exclude them. Although this restricts the system to English, it does not impact performance due to the dataset's predominantly English-language content. Multi-lingual support could be added by translating text to English before summarization.

### 3.4. Coherent Reasoning

All the different components of LLM-Consensus 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. We use the NewsCLIPpings dataset for evaluating our model's efficacy in identifying out-of-context images and information due to its well-established use in prior work, its relative scale compared to other available datasets, and the lack of more recent large-scale benchmarks in this domain. Though no single dataset can capture the full diversity of real-world misinformation, NewsCLIPpings provides a standardized evaluation setting that facilitates comparability with existing methods.

## 4.2. Experimental Setup

All experiments were run on 8 A40 (46GB) Nvidia GPU server. The estimated cost of processing one data sample using LLM-Consensus 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). 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 LLM-Consensus 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 LLM-Consensus 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 LLM-Consensus 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$_{AI}$ (believes debating AI) | 75.2 | 54.5 | 86.4 |
| Async_Debate$_{human}$ (w/o external info) | 77.1 | 68.4 | 89.3 |
| Async_Debate$_{human}$ (w external info) | **86.2** | **82.6** | **90.6** |
| 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$_{human}$ 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$_{human}$ setup where the agent has access to external information performs the best of all the debating configurations. We also report results for the Async_Debate$_{human}$ 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. 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_human` 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_human` (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.

| Model | Accuracy ↑ |
|---|---|
| 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[#] (w retrieval) | 86.0 |
| Sniffer (w finetuning) | 88.4 |
| Sniffer (w/o finetuning) | 84.5 |
| **LLM-Consensus (ours)** | **90.8** |

*Table 2.* **Performance comparison between our model and baselines:** LLM-Consensus (with GPT-4o) out performs all related work. Note: the GPT-4o[#] setup is identical to LLM-Consensus with the absence of a multi-agent debate, here only a single agent with access to external information is considered.

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 NewsCLIP-

| 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 | **90.3** | 90.1 |
| ✗ | ✓ | ✓ | ✓ | **90.8** | 85.5 | **99.0** |

*Table 3.* **Ablation:** Quantitative evaluation of each component of LLM-Consensus on classification performance.

pings dataset. While we do not provide a quantitative assessment of explanations by LLM-Consensus, 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 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 also note that single multi-modal models, including VisualBERT, CLIP, InstructBLIP, LLaVA and GPT-4o do not perform at par with other related work. This can be attributed to the necessity for external context for misinformation detection in the news domain and the lack of diverse perspectives that arise naturally in a multi-agent framework. Therefore, these standalone models, while promising, currently are unable to detect misinformation effectively. These models require additional integration into more comprehensive pipelines, as done in this work. In line with previous work, we also note that baselines trained from scratch, such as SAFE (Zhou et al., 2020) and EANN (Wang et al., 2018) perform worse than pretrained multi-modal models. This further concretizes the fact that image-based OOC detection in the news domain requires strong world knowledge as well as advanced multi-modal reasoning capabilities.

### 4.3.3. ABLATIONS

To analyze the importance of each component of the LLM-Consensus 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 2) to 90.2, underscoring its role in enabling contextual reasoning and refining

| Study Setup | Average Accuracy↑ |
|---|---|
| Humans | $60.3 \pm 13.5$ |
| Humans+LLM-Consensus | $76.7 \pm 12.2$ |
| **LLM-Consensus** | $\mathbf{80.0 \pm 0.0}$ |

*Table 4.* **Performance comparison between different study setups:** LLM-Consensus outperforms humans with and without AI assistance.

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.

## 5. User Study

We conducted a user study to evaluate the effectiveness of our system in detecting and explaining misinformation. While model performance in misinformation detection is easy to quantify, there are no established metrics for assessing the quality of its explanations, making the user study essential. Participants were grouped into three categories based on their professions: Journalists, AI Academics, and Others (see Appendix A.4 for details).

During the study, participants reviewed ten image-text pairs and decided whether each pair constituted misinformation. They also rated their confidence on a scale of 0 to 10. After submitting their initial answers, participants were shown LLM-Consensus's AI-generated explanations and asked to reconsider their responses. As shown in Table 4, the system outperformed the average human performance in both cases, with and without AI insights. These results demonstrate LLM-Consensus's potential as a reliable assistive tool for OSINT research, capable of detecting and explaining misinformation with minimal human intervention.

Group-wise analysis, shown in Table 5, 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 LLM-Consensus insights. Thus, LLM-Consensus can substantially boost non-expert performance, making it valuable for citizen intelligence applications.

## 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

| Metric | Journalists | Academics | Others |
|---|---|---|---|
| Accuracy (only human) | $70.0 \pm 1.4$ | $60.7 \pm 1.4$ | $56.7 \pm 1.5$ |
| Confidence (only human) | $4.3 \pm 2.1$ | $3.2 \pm 0.8$ | $3.9 \pm 1.2$ |
| Accuracy (with LLM-Consensus) | $82.2 \pm 0.9$ | $79.3 \pm 1.3$ | $71.7 \pm 1.1$ |
| Confidence (with LLM-Consensus) | $5.3 \pm 1.3$ | $5.8 \pm 1.4$ | $5.8 \pm 1.4$ |

*Table 5.* **Performance comparison:** LLM-Consensus improves performance across all participant groups.

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_Debate_human` 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, LLM-Consensus, achieves state-of-the-art performance in misinformation detection while offering clear, detailed explanations, thanks to our advanced external retrieval module. As a result, we observe significant improvements in out-of-context (OOC) detection for both experts and non-experts.

We identify several promising directions for future research. The community would benefit from a continuously updated benchmark dataset incorporating recent news articles and subtler inconsistencies. We plan to extend our methods to video-text pairs and multi-lingual content. We also believe that integrating more advanced models in the summarization pipeline will further enhance performance and are interested in comparing LLM-Consensus with multi-agent systems using external retrieval.

Finally, we conducted extensive user studies but recognize the importance of large-scale deployment in professional settings and the citizen intelligence community. This will provide valuable real-world feedback and uncover further opportunities for improvement. For a discussion of limitations, see Appendix A.1.

# References

Deepfakes, explained — MIT Sloan — mit-sloan.mit.edu. https://mitsloan.mit.edu/ideas-made-to-matter/deepfakes-explained. [Accessed 28-09-2024].

Newspaper3k: Article scraping & curation — newspaper 0.0.2 documentation. https://newspaper.readthedocs.io/en/latest/.

Out-of-context photos are a powerful low-tech form of misinformation — pbs.org. URL https://www.pbs.org/newshour/science/. [Accessed 28-09-2024].

Visual Search API — Microsoft Bing. https://www.microsoft.com/en-us/bing/apis/bing-visual-search-api.

Abdelnabi, S., Hasan, R., and Fritz, M. Open-domain, content-based, multi-modal fact-checking of out-of-context images via online resources, 2022. URL https://arxiv.org/abs/2112.00061.

Amerini, I., Anagnostopoulos, A., Maiano, L., Celsi, L. R., et al. Deep learning for multimedia forensics. *Foundations and Trends® in Computer Graphics and Vision*, 12 (4):309–457, 2021.

Aneja, S., Midoglu, C., Dang-Nguyen, D.-T., Khan, S. A., Riegler, M., Halvorsen, P., Bregler, C., and Adsumilli, B. Acm multimedia grand challenge on detecting cheapfakes, 2022. URL https://arxiv.org/abs/2207.14534.

Aslett, K., Sanderson, Z., Godel, W., Persily, N., Nagler, J., and Tucker, J. A. Online searches to evaluate misinformation can increase its perceived veracity. *Nature*, 625 (7995):548–556, 2024.

Bai, Z., Wang, P., Xiao, T., He, T., Han, Z., Zhang, Z., and Shou, M. Z. Hallucination of multimodal large language models: A survey, 2024. URL https://arxiv.org/abs/2404.18930.

Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U., Knight, K., Koehn, P., Palmer, M., and Schneider, N. Abstract Meaning Representation for sembanking. In Pareja-Lora, A., Liakata, M., and Dipper, S. (eds.), *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pp. 178–186, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL https://aclanthology.org/W13-2322.

Brashier, N. M. and Marsh, E. J. Judging truth. *Annual review of psychology*, 71(1):499–515, 2020.

Castillo Camacho, I. and Wang, K. A comprehensive review of deep-learning-based methods for image forensics. *Journal of imaging*, 7(4):69, 2021.

Dai, W., Li, J., Li, D., Tiong, A. M. H., Zhao, J., Wang, W., Li, B., Fung, P., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. URL https://arxiv.org/abs/2305.06500.

Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. Improving factuality and reasoning in language models through multiagent debate, 2023a. URL https://arxiv.org/abs/2305.14325.

Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. Improving Factuality and Reasoning in Language Models through Multiagent Debate. October 2023b. URL https://openreview.net/forum?id=QAwaaLJNCk.

Farid, H. *Photo Forensics*. The MIT Press, 2016. ISBN 0262035340.

Hasher, L., Goldstein, D., and Toppino, T. Frequency and the conference of referential validity. *Journal of verbal learning and verbal behavior*, 16(1):107–112, 1977.

Heidari, A., Jafari Navimipour, N., Dag, H., and Unal, M. Deepfake detection using deep learning methods: A systematic and comprehensive review. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 14(2):e1520, 2024.

Hina, M., Ali, M., Javed, A. R., Ghabban, F., Khan, L. A., and Jalil, Z. Sefaced: Semantic-based forensic analysis and classification of e-mail data using deep learning. *IEEE Access*, 9:98398–98411, 2021.

Huang, M., Jia, S., Zhou, Z., Ju, Y., Cai, J., and Lyu, S. Exposing text-image inconsistency using diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024.

Khan, A., Hughes, J., Valentine, D., Ruis, L., Sachan, K., Radhakrishnan, A., Grefenstette, E., Bowman, S. R., Rocktäschel, T., and Perez, E. Debating with more persuasive llms leads to more truthful answers, 2024. URL https://arxiv.org/abs/2402.06782.

Li, G., Hammoud, H. A. A. K., Itani, H., Khizbullin, D., and Ghanem, B. Camel: Communicative agents for "mind" exploration of large language model society, 2023a. URL https://arxiv.org/abs/2303.17760.

Li, J., Li, D., Savarese, S., and Hoi, S. BLIP-2: Bootstrapping language-image pre-training with frozen image

encoders and large language models. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19730–19742. PMLR, 23–29 Jul 2023b. URL https://proceedings.mlr.press/v202/li23q.html.

Li, L. H., Yatskar, M., Yin, D., Hsieh, C.-J., and Chang, K.-W. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.

Lin, H., Luo, Z., Gao, W., Ma, J., Wang, B., and Yang, R. Towards explainable harmful meme detection through multimodal debate between large language models, 2024. URL https://arxiv.org/abs/2401.13298.

Liu, F., Wang, Y., Wang, T., and Ordonez, V. Visualnews : Benchmark and challenges in entity-aware image captioning, 2020.

Liu, F., Lin, K., Li, L., Wang, J., Yacoob, Y., and Wang, L. Mitigating hallucination in large multi-modal models via robust instruction tuning, 2024a. URL https://arxiv.org/abs/2306.14565.

Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., and Lee, Y. J. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.

Luo, G., Darrell, T., and Rohrbach, A. Newsclippings: Automatic generation of out-of-context multimodal media. *arXiv:2104.05893*, 2021.

Masood, M., Nawaz, M., Malik, K. M., Javed, A., Irtaza, A., and Malik, H. Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward. *Applied intelligence*, 53(4):3974–4026, 2023.

Minsky, M. *Society of mind*. Simon and Schuster, 1988.

Mu, M., Das Bhattacharjee, S., and Yuan, J. Self-supervised distilled learning for multi-modal misinformation identification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 2819–2828, January 2023.

OpenAI. GPT-4o. https://openai.com/index/hello-gpt-4o/. [Accessed 28-08-2024].

OpenAI and et al., J. A. GPT-4 Technical Report, 2024. URL https://arxiv.org/abs/2303.08774.

Papadopoulos, S.-I., Koutlis, C., Papadopoulos, S., and Petrantonakis, P. Synthetic misinformers: Generating

and combating multimodal misinformation. In *Proceedings of the 2nd ACM International Workshop on Multimedia AI against Disinformation*, MAD '23, pp. 36–44, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701870. doi: 10.1145/3592572.3592842. URL https://doi.org/10.1145/3592572.3592842.

Qi, P., Yan, Z., Hsu, W., and Lee, M. L. Sniffer: Multimodal large language model for explainable out-of-context misinformation detection, 2024. URL https://arxiv.org/abs/2403.03170.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models, 2022. URL https://arxiv.org/abs/2112.10752.

Schroeder de Witt, C., Gupta, T., Makoviichuk, D., Makoviychuk, V., Torr, P. H. S., Sun, M., and Whiteson, S. Is Independent Learning All You Need in the StarCraft Multi-Agent Challenge?, November 2020. URL https://arxiv.org/abs/2011.09533v1.

Shalabi, F., Nguyen, H. H., Felouat, H., Chang, C.-C., and Echizen, I. Image-text out-of-context detection using synthetic multimodal misinformation. In *2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, October 2023. doi: 10.1109/apsipaasc58517.2023.10317336. URL http://dx.doi.org/10.1109/APSIPAASC58517.2023.10317336.

Sultan, M., Tump, A. N., Geers, M., Lorenz-Spreen, P., Herzog, S. M., and Kurvers, R. H. Time pressure reduces misinformation discrimination ability but does not alter response bias. *Scientific Reports*, 12(1):22416, 2022.

Tolosana, R., Vera-Rodriguez, R., Fierrez, J., Morales, A., and Ortega-Garcia, J. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information Fusion*, 64:131–148, 2020.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

Wang, S.-Y., Wang, O., Owens, A., Zhang, R., and Efros, A. A. Detecting photoshopped faces by scripting photoshop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

Wang, Y., Ma, F., Jin, Z., Yuan, Y., Xun, G., Jha, K., Su, L., and Gao, J. 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, pp. 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.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, pp. 24824–24837, Red Hook, NY, USA, April 2024. Curran Associates Inc. ISBN 978-1-71387-108-8.

Yi, K., Wu, J., Gan, C., Torralba, A., Kohli, P., and Tenenbaum, J. B. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding, 2019. URL https://arxiv.org/abs/1810.02338.

Zhang, Y., Trinh, L., Cao, D., Cui, Z., and Liu, Y. Interpretable detection of out-of-context misinformation with neural-symbolic-enhanced large multimodal model, 2024. URL https://arxiv.org/abs/2304.07633.

Zhou, X., Wu, J., and Zafarani, R. Safe: Similarity-aware multi-modal fake news detection, 2020. URL https://arxiv.org/abs/2003.04981.

Zhu, W., Thomason, J., and Jia, R. Generalization differences between end-to-end and neuro-symbolic vision-language reasoning systems, 2022. URL https://arxiv.org/abs/2210.15037.

Zhu, X., Qian, Y., Zhao, X., Sun, B., and Sun, Y. 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 LLM-Consensus, 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. 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 LLM-Consensus 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 LLM-Consensus'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 LLM-Consensus.

## 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 LLM-Consensus

```
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 5.* Initial prompt for independent opinion formation and response generation

```
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 6.* Prompt for Debate Round 1

```
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 7.* Prompt for Debate after Round 1

## A.4. 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.4.1. SETUP

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. The group of individuals included journalists from BBC as well as students and professors from the University of Oxford. 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 LLM-Consensus. 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 LLM-Consensus 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 LLM-Consensus. 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 LLM-Consensus 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.5. Multi-modal debates for hamful 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.