

# MMFAKEBENCH: A MIXED-SOURCE MULTIMODAL MISINFORMATION DETECTION BENCHMARK FOR LVLMS

Anonymous authors

Paper under double-blind review

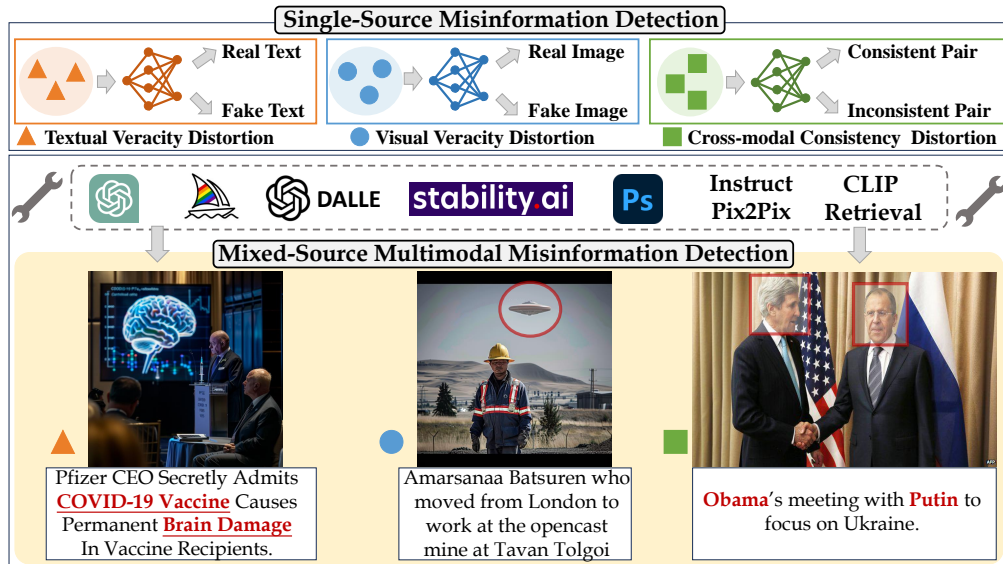


Figure 1: **Top:** Previous methods often assume a single misinformation source and conduct single-source detection. **Bottom:** We collaborate generative models and AI tools to build a mixed-source multimodal misinformation benchmark and achieve mixed-source detection.

## ABSTRACT

Current multimodal misinformation detection (MMD) methods often assume a single source and type of forgery for each sample, which is insufficient for real-world scenarios where multiple forgery sources coexist. The lack of a benchmark for mixed-source misinformation has hindered progress in this field. To address this, we introduce MMFAKEBENCH, the first comprehensive benchmark for mixed-source MMD. MMFAKEBENCH includes 3 critical sources: textual veracity distortion, visual veracity distortion, and cross-modal consistency distortion, along with 12 sub-categories of misinformation forgery types. We further conduct an extensive evaluation of 6 prevalent detection methods and 15 Large Vision-Language Models (LVLMS) on MMFAKEBENCH under a zero-shot setting. The results indicate that current methods struggle under this challenging and realistic mixed-source MMD setting. Additionally, we propose a new approach MMD-Agent, which integrates the reasoning, action, and tool-use capabilities of LVLMS agents, significantly enhancing accuracy and generalization. We believe this study will catalyze future research into more realistic mixed-source multimodal misinformation and provide a fair evaluation of misinformation detection methods. Code and a portion of the data are accessible in supplementary materials.

## 1 INTRODUCTION

Recent advances in generative models for texts (Brown et al., 2020; Touvron et al., 2023) and images (Dhariwal & Nichol, 2021; Rombach et al., 2022) have significantly lowered the barrier to producing diverse multimodal misinformation, posing threats to politics, finance, and public health. For instance, the misinformation “COVID-19 vaccine causes brain damage”, shown in Fig. 1, accompanied by a highly convincing image, can lead to public distrust in medical treatments and vaccine refusal. Therefore, identifying multimodal misinformation on social media is urgent.

Most current multimodal misinformation detection (MMD) methods (Abdelnabi et al., 2022; Qi et al., 2024; Ying et al., 2023; Huang et al., 2023; Zhang & Gao, 2023; Lee et al., 2021) typically assume that each sample has a single, known forgery source. As depicted in Fig. 1 **Top**, these forgery sources involve either textual veracity with fake news text, visual veracity with fake images, or inconsistency between the text and image. However, the single-source assumption is overly simplistic and fails to capture the complexity of real-world scenarios, where misinformation often stems from multiple, random sources. To address this mixed-source MMD problem, two key challenges need to be solved. First, existing datasets primarily consist of single-source misinformation, lacking misinformation from multiple sources. This limitation prevents comprehensive evaluation of MMD methods. Second, there is a lack of general detectors capable of handling mixed-source misinformation. Hence, we present MMFakeBench, encompassing the mixed-source MMD benchmark, evaluations, and framework.

**Benchmark:** We introduce MMFakeBench, the first comprehensive benchmark for evaluating mixed-source MMD. As shown in Fig. 1 **Bottom**, leveraging advanced AI tools, such as diffusion generators and ChatGPT, MMFakeBench provides 12 forgery types with 11,000 data pairs from three primary sources: *textual veracity distortion*, *visual veracity distortion*, and *cross-modal consistency distortion*. Textual veracity distortion encompasses three types of rumors: natural, artificial, and GPT-generated rumors. Unlike (Thorne et al., 2018; Shu et al., 2020; Hanselowski et al., 2019; Chen & Shu, 2024) that focus solely on single-source, text-only rumors, MMFakeBench incorporates text-image rumors using highly relevant real or AI-generated images. Visual veracity distortion filters existing PS-edited images (Da et al., 2021; Nakamura et al., 2020) according to misinformation standards and incorporates high-quality AI-generated images. Cross-modal consistency distortion integrates inconsistencies from both edited and repurposed perspectives into five distinct sub-categories.

**Evaluations:** To access the current advancements in mixed-source MMD, we build the fine-grained multi-class evaluation metric and conduct a comprehensive evaluation of 6 state-of-the-art detection methods and 15 large vision-language models (LVLMs) on MMFakeBench. Specifically, we evaluate 6 detection methods in a single-source setting and assess their combined performance (text, image, and cross-modal inconsistency detectors) in the mixed-source setting. Additionally, we evaluate 15 large vision-language models (LVLMs), including proprietary models such as GPT-4V (OpenAI, 2023). The results indicate that existing detection methods exhibit poor generalization. Although LVLMs show robust generalization capabilities, their overall performance still requires improvement.

**Framework:** Based on our analysis, we propose a simple yet effective LVLm-based framework called **MMD-Agent**, which enhances detection performance and serves as a new baseline for future research. MMD-Agent decomposes mixed-source detection into three stages: textual veracity check, visual veracity check, and cross-modal consistency reason. This decomposition ensures methodical and thorough reasoning. At each stage, MMD-Agent instructs LVLMs to generate multi-perspective reasoning traces, integrating model actions for coherent decisions. Additionally, the models interact with external knowledge sources via tools (e.g., Wikipedia) to incorporate supplementary information into their reasoning.

In summary, the main contributions are: (1) We introduce mixed-source multimodal misinformation detection (MMD), a challenging setting for detecting misinformation from diverse and uncertain sources, breaking free from single-source constraints, and advancing practical misinformation detection tasks. (2) We develop MMFakeBench, the first benchmark dataset for evaluating mixed-source MMD. The dataset contains 3 critical categories (textual veracity distortion, visual veracity distortion, and consistency reasoning) and 12 sub-categories of forgery types. (3) Using the newly collected dataset, we benchmark mixed-source MMD by evaluating 6 prevalent detection methods and 15 LVLMs. (4) We propose MMD-Agent, a simple yet effective LVLm-based framework. It outperforms previous methods and LVLMs on the MMFakeBench benchmark, highlighting the potential of mixed-source MMD and providing a new baseline for future research.

Table 1: Comparison of misinformation datasets. (k) denotes the number of rumor types.  $\otimes$  denotes editing methods may unintentionally introduce the fact-conflicting content.

Dataset	Textual Veracity Distortion			Visual Veracity Distortion			Cross-modal Consistency Distortion	
	Text	Supporting Image		Text	Fact-conflicting Image		Image/Text	Image/Text
	(Rumor)	Repurposed	AI-generated	(Veracity)	PS-edited	AI-generated	Repurposing	Editing
FEVER (Thorne et al., 2018)	✓(1)	✗	✗	✗	✗	✗	✗	✗
Politifact (Shu et al., 2020)	✓(1)	✗	✗	✗	✗	✗	✗	✗
Gossipcop (Shu et al., 2020)	✓(1)	✗	✗	✗	✗	✗	✗	✗
Snopes (Hanselowski et al., 2019)	✓(1)	✗	✗	✗	✗	✗	✗	✗
MOCHEG (Yao et al., 2023)	✓(1)	✗	✗	✗	✗	✗	✗	✗
LLMFake (Chen & Shu, 2024)	✓(1)	✗	✗	✗	✗	✗	✗	✗
EMU (Da et al., 2021)	✗	✗	✗	✗	✓	✗	✗	✗
Fakeddit (Nakamura et al., 2020)	✗	✗	✗	✓	✓	✗	✗	✗
MAIM (Jaiswal et al., 2017)	✗	✗	✗	✗	✗	✗	✓	✗
MEIR (Sabir et al., 2018)	$\otimes$	✗	✗	✗	✗	✗	✗	✓
NewsCLIPpings (Luo et al., 2021)	✗	✗	✗	✗	✗	✗	✓	✗
COSMOS (Aneja et al., 2023)	✗	✗	✗	✗	✗	✗	✓	✗
DGM4 (Shao et al., 2023)	$\otimes$	✗	✗	✗	✗	✗	✗	✓
<b>MMFakeBench (Ours)</b>	✓(3)	✓	✓	✓	✓	✓	✓	✓

## 2 RELATED WORK

**Misinformation Benchmarks.** One group of misinformation datasets primarily focuses on distorting textual veracity. The FEVER (Thorne et al., 2018) dataset is constructed by manipulated Wikipedia sentences with manual annotation. Unlike these artificial rumors, other datasets, such as Snopes (Hanselowski et al., 2019), Politifact, Gossipcop (Shu et al., 2020), and MOCHEG (Yao et al., 2023), collect natural rumors from fact-checking websites. Recently, the LLMFake (Chen & Shu, 2024) instructs large language models (LLMs) to generate diverse misinformation. Apart from misleading text, the EMU (Da et al., 2021) and Fakeddit (Nakamura et al., 2020) collect Photoshop-edited images from the Reddit platform. Another group of misinformation datasets focuses on disrupting cross-modal consistency. The MAIM (Jaiswal et al., 2017) and MEIR (Sabir et al., 2018) datasets employ caption replacement and entity swapping, respectively. The NewsCLIPpings (Luo et al., 2021) and COSMOS (Aneja et al., 2023) datasets link out-of-context images to support certain narratives. The recent dataset DGM<sup>4</sup> (Shao et al., 2023) introduces global and local manipulation to alter semantics and sentiment. Different from these works containing only single-source misinformation, we propose the first benchmark dataset for evaluating mixed-source MMD, involving textual veracity distortion, visual veracity distortion, and cross-modal consistency distortion, shown in Table 1.

**Misinformation Detection.** Current misinformation detection approaches are mainly divided into two categories. The first is to check textual veracity by constructing features based on writing style (Przybyla, 2020), sentiment (Ghanem et al., 2021), user feedback (Min et al., 2022) and pre-trained language models (Huang et al., 2023; Zhang & Gao, 2023). The second is to fuse cross-modal features to detect semantic inconsistencies. Previous works focus on devising attention-based modules (Qian et al., 2021b; Ying et al., 2023; Wu et al., 2021) guided by diverse learning strategies (Chen et al., 2023). Recent works Shao et al. (2024); Qi et al. (2024); Liu et al. (2024c) capitalize the VLMS which benefit from large-scale pre-training for reasoning context cues. However, these works target a single-source problem, and evaluations are conducted in constrained scenarios. Our work is the first to introduce a comprehensive benchmark for mixed-source multimodal misinformation detection.

**Large Vision-Language Models.** Large language models (LLMs) such as GPT-3 (Brown et al., 2020) and Vicuna (Chiang et al., 2023) have demonstrated remarkable performance on various linguistic tasks. Inspired by LLMs, models like LLaVA (Liu et al., 2023a) and MiniGPT-4 (Zhu et al., 2023) facilitate image-text feature alignment by leveraging visual instruction tuning. More recently, the evolution of LVLMs has driven advancements in creating diverse and high-quality multimodal instruction datasets. Models such as InstructBLIP (Dai et al., 2023), mPLUG-Owl (Ye et al., 2023; 2024), LLaVA-1.5 (Liu et al., 2024a) exemplify these developments. In this paper, we explore the reasoning capabilities (Zheng et al., 2023; Zhang et al., 2024) of LVLMs to address the challenge of mixed-source multimodal misinformation by integrating reasoning, actions, and tool-use capabilities.

## 3 MMFAKEBENCH BENCHMARK

In MMFakeBench, we focus on multimodal misinformation involving both text and images, categorizing it into three distinct types based on the sources of falsified content:

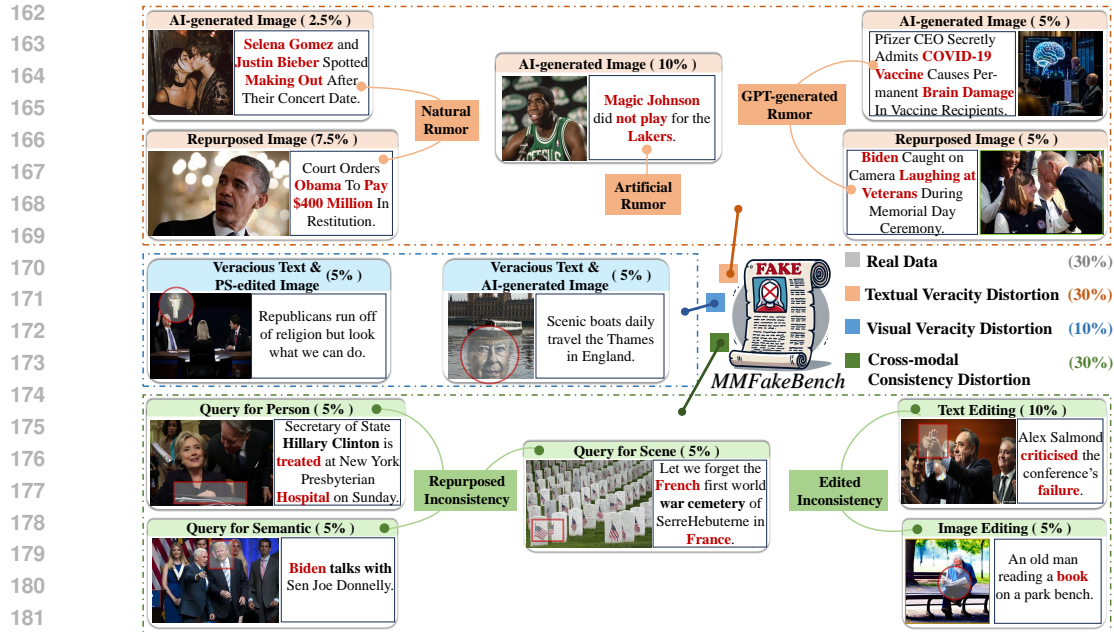


Figure 2: Statistics of the MMFakeBench Benchmark.

- *Textual Veracity Distortion*. It incorporates text-based rumors paired with supporting images to mimic real-world multimodal scenarios. An image that visually supports misleading text can make the misinformation appear more credible and persuasive to the users.
- *Visual Veracity Distortion*. Images in this category contain fact-conflicting misinformation through altered or fabricated elements, while the texts remain veracious. These visual manipulations often lead people to perceive falsified content as authentic, distorting their understanding of the information.
- *Cross-modal Consistency Distortion*. Even when the text and image are individually correct, their combination can generate potential misinterpretation if presented in a manner that introduces incorrect associations or semantic discrepancies between the two modalities, thereby misleading people.

### 3.1 THREE MISINFORMATION SOURCES

#### 3.1.1 TEXTUAL VERACITY DISTORTION

Textual veracity distortion is a critical misinformation category. Our dataset in this category comprises 3,300 samples. Previous works (Thorne et al., 2018; Shu et al., 2020; Hanselowski et al., 2019; Chen & Shu, 2024) focus on single-source and single-modal textual rumors. However, the types of real-world rumors are diverse, and those accompanied by images can have a significantly greater impact. To address this, MMFakeBench introduces a broader range of rumor types and augments them with highly relevant supporting images to enhance perceived credibility.

**Textual Rumor.** As shown in Table 1, unlike previous methods that consider only one type of textual rumor, we consider three types: (1) Natural Rumor. We select natural rumors from Politifact and Gossipcop (Shu et al., 2020), which provide political news and entertainment stories derived from fact-checking websites. (2) Artificial Rumor. We collect artificial rumors from the FEVER dataset (Thorne et al., 2018), which is curated by manually modifying Wikipedia sentences. (3) GPT-generated Rumor. We instruct ChatGPT (`gpt-3.5-turbo`) to produce rumors via three prompt approaches (Chen & Shu, 2024): arbitrary generation, rewriting generation, and information manipulation. Arbitrary generation is utilized to generate misinformation in specific domains. Rewriting generation addresses the concise and synthetic traces of artificial rumors, while information manipulation involves altering factual information in real claims from Politifact and Gossipcop.

**Supporting Image.** We use either AI-generated images or carefully selected real images to support the content presented in the rumor text. (1) AI-generated Image: For artificial rumors and their derived GPT-generated rumors, as well as some less harmful gossip, we utilize generative models to

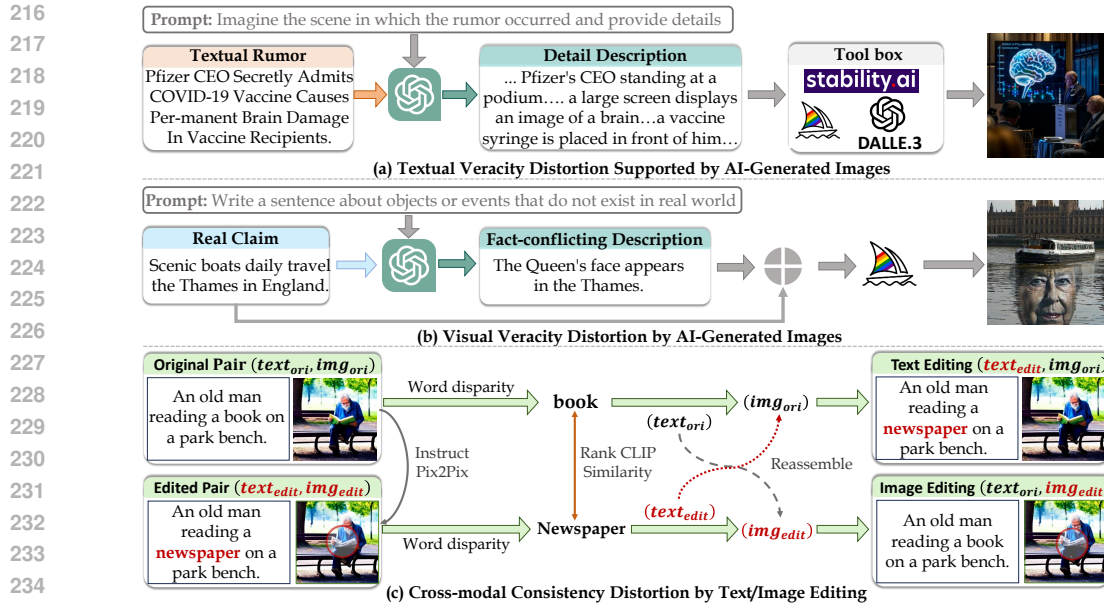


Figure 3: Illustrations of using collaborative generative models and AI tools to generate different sources of misinformation.

create supporting images. We utilize three advanced models: Stable Diffusion XL (Rombach et al., 2022), DALL-E3 (Ramesh et al., 2022), and Midjourney V6 (Midjourney, 2022) to enhance the diversity of the generated images. For each rumor, we randomly select a generative model to produce a corresponding supporting image. As many rumors are highly abstract and lack concrete descriptions of objects and scenes, directly using these texts as conditions often yields images that are neither realistic nor relevant. To address this, as shown in Fig. 3 (a), we instruct ChatGPT to enrich the rumors with more detailed descriptions. These enriched contexts are then used as input for generative models, ensuring alignment with the textual rumors. (2) Repurposed Image: To avoid creating new high-risk images, especially for sensitive topics like politics and gossip, we use repurposed images from the VisualNews (Liu et al., 2021) dataset, which contains numerous image-text pairs from real-world news sources. We select images with high semantic relevance to the textual rumors based on text-image CLIP similarity and text-text CLIP similarity. Images with the highest similarity scores are chosen as supporting images.

### 3.1.2 VISUAL VERACITY DISTORTION

The visual veracity distortion dataset comprises 1,100 samples where the text is real and the misinformation exists in the image. Previous datasets (Da et al., 2021; Nakamura et al., 2020) focus solely on PS-edited (Photoshop-edited) images, containing both misleading and non-misleading content. In this study, we manually select the misleading ones and include them in MMFakeBench. Besides, we incorporate AI-generated images with veracity distortion, which is increasingly harmful due to advancements in diffusion generators.

**PS-Edited Image.** The PS-edited images are derived from the "manipulated content" samples in the Fakeddit dataset (Nakamura et al., 2020), which is designed for multimodal fake news detection. These samples originate from the "Photoshop battles comments" on Reddit. PS-edited samples in the Fakeddit typically exhibit either aesthetic modifications or fact-conflicting manipulations. Since aesthetic modifications do not compromise the factuality of the visual content, they are excluded from our benchmark criteria. Consequently, ten of the volunteers participate in selecting 550 PS-edited images containing fact-conflicting content from the 7,693 samples of the "manipulated content" set.

**AI-generated Image.** We propose an automated pipeline that generates fact-conflicting descriptions from text captions and then creates high-quality images. Specifically, we first collect image-text pairs from the MS-COCO (Lin et al., 2014) and VisualNews datasets (Liu et al., 2021). Based on the original text captions, we use ChatGPT to generate corresponding fact-conflicting descriptions, which are manually verified, as depicted in Fig. 3 (b). For example, from the caption "Scenic boats daily travel the Thames in England", we generate the description "The Queen's face appears in

the Thames”. These descriptions, combined with the original captions, are used as prompts in the Midjourney V6 model (Midjourney, 2022) to create corresponding images. The resulting text-image pairs contain original factual text and generated images with additional fact-conflicting information.

### 3.1.3 CROSS-MODAL CONSISTENCY DISTORTION

In cross-modal consistency distortion, both the text and image with veracity, but either the text or image is replaced/manipulated to disrupt their overall consistency. Previous datasets (Sabir et al., 2018; Luo et al., 2021; Shao et al., 2023) focus on inconsistencies from a single source, either edit-based or repurposed-based. In contrast, our MMFakeBench integrates inconsistencies from both edited and repurposed perspectives into five distinct sub-categories, a total of 3,300 image-text pairs.

**Repurposed Inconsistency.** Our dataset contains three types of repurposed inconsistency, curated directly from the NewsCLIPings (Luo et al., 2021) dataset: semantic query, person query, and scene query. (1) Semantic query retrieves repurposed images based on specific semantic content. (2) Person query ensures the individual mentioned in the caption appears in the mismatched image. (3) Scene Query relies on spatial similarity to retrieve comparable scene information from repurposed images.

**Edited Inconsistency.** Our dataset contains two types of edited inconsistency: text editing and image editing. For text editing, we select samples from the DGM<sup>4</sup> (Shao et al., 2023) dataset, which modifies sentiment words with their antonyms. Notably, some samples in DGM<sup>4</sup> contain fact-conflicting content post-editing. To avoid redundancy with textual veracity distortion, we filter out these samples. For remaining text editing and all image editing inconsistencies, we build upon the COCO-Counterfactuals (Le et al., 2023) dataset. This dataset encompasses original image-text pairs ( $text_{ori}, img_{ori}$ ) and edited image-text pairs ( $text_{edit}, img_{edit}$ ) which are obtained via Instruct-Pix2Pix model (Brooks et al., 2023). As illustrated in Fig. 3 (c), we separately extract word disparities between  $text_{ori}$  and  $text_{edit}$  and select samples with significant semantic differences using CLIP similarity. Then, we reassemble the two pairs and obtain ( $text_{edit}, img_{ori}$ ) as text-edited consistency distortion samples and ( $text_{ori}, img_{edit}$ ) as image-edited consistency distortion samples.

## 3.2 REAL DATA COLLECTION

In addition to the misinformation data, we collect 3,300 real data pairs, ensuring both textual and visual veracity and exhibiting strong image-text consistency. Given that our synthetic data is derived from multiple datasets, we construct the real dataset from the same corresponding sources, including MS-COCO, VisualNews, and real image-text pairs from Fakeddit. We further divide VisualNews into four distinct news sources: The Guardian, BBC, USA TODAY, and The Washington Post. Finally, we build the real dataset by equally selecting from six distinct sources.

## 3.3 MMFAKEBENCH ANALYSIS

MMFakeBench consists of 11,000 image-text pairs, which are divided into a validation set and a test set following (Yue et al., 2024). The validation set, comprising 1,000 image-text pairs, is intended for hyperparameter selection, while the test set contains 10,000 pairs. MMFakeBench encompasses one real category and three misinformation categories. Detailed statistics are shown in Fig. 2. MMFakeBench is partitioned into 30% for textual veracity distortion, 10% for visual veracity distortion, 30% for cross-modal consistency distortion, and 30% for real data. The three misinformation categories can be further subdivided into 12 detailed subcategories based on the sources of the text and images in Fig. 2. Such a comprehensive benchmark highlights the challenges of intertwining mixed-source and multiple-types multimodal misinformation in the real world.

## 4 MMD-AGENT FRAMEWORK

We present a simple yet effective framework, *MMD-Agent*, which integrates the reasoning, actions, and tool-use capabilities of LVLm agents. As shown in Fig. 4, MMD-Agent involves two main processes: (1) Hierarchical decomposition and (2) Integration of internal and external knowledge.

Specifically, we instruct LVLms  $\mathcal{M}$  to decompose the task of mixed-source multimodal misinformation detection into three smaller subtasks: textual veracity check, visual veracity check, and cross-modal consistency reasoning. During the intermediary phase, each subtask  $t$  is addressed

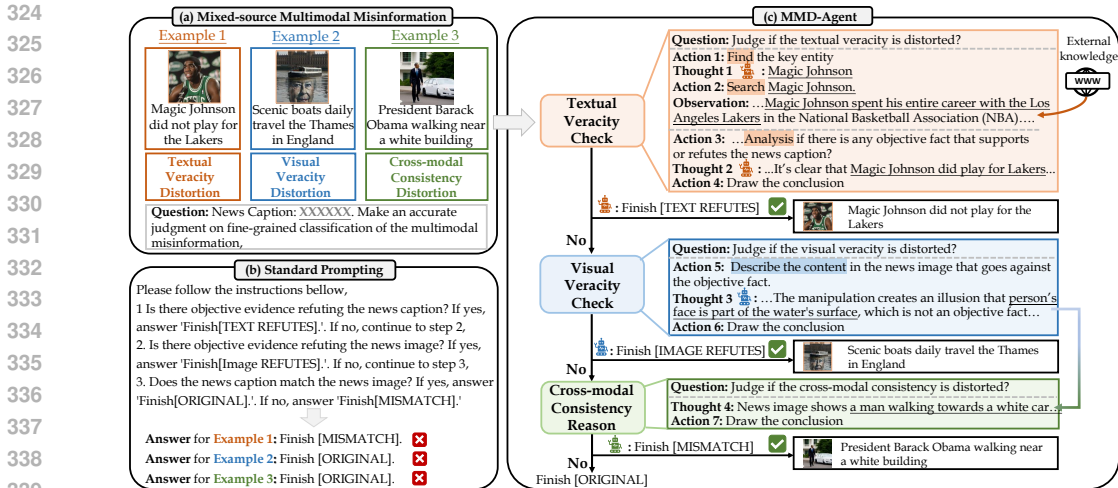


Figure 4: Comparison of standard prompting and proposed MMD-Agent. (a) Three examples of multimodal misinformation from distinct sources. (b) LVLMs with standard prompting methods fail to make correct judgments. (c) MMD-Agent instructs LVLMs to decompose mixed-source detection into smaller subtasks, which are solved by integrating model thoughts and environment observation.

through an interleaved sequence of reasoning and action. The LVLm is guided to reason and induce the needed action to solve the task. The actions  $a_t$  are then executed either by leveraging the model’s internal knowledge to generate “Thought”  $\mathcal{R}_t$  or by interacting with external sources to gather additional information (“Observation”)  $\mathcal{O}_t$ . These action outputs will be integrated into the sequence to facilitate subsequent decision-making  $d_t$ :

$$d_t = \mathcal{M}(\mathcal{R}_t, \mathcal{O}_t, a_t). \tag{1}$$

The LVLm’s internal capabilities and knowledge are utilized both to reason about which actions to take and to perform those actions from various perspectives, such as identifying key textual entities (Thought 1), conducting factual analysis (Thought 2), and applying commonsense reasoning (Thought 3). However, the model may experience hallucinations when relying solely on its internal knowledge. To address this, we enable the model to interact with external knowledge bases, such as the Wikipedia API, to retrieve reliable and up-to-date information for fact-checking. This approach ensures the accuracy and relevance of the knowledge used in verification.

## 5 EXPERIMENTS

We select 6 state-of-the-art misinformation detection methods and 15 large vision-language models (LVLms) for preliminary benchmarking using the MMFakeBench dataset, aiming to explore their applicability in the mixed-source MMD setting. Additionally, we evaluate the performance of our proposed framework, MMD-Agent. All evaluations are conducted under a zero-shot setting on our benchmark. All experiments are performed on eight NVIDIA GeForce 3090 GPUs with PyTorch.

### 5.1 EXPERIMENTAL SETTING

**Baseline Models.** We select LVLms of varying sizes as baseline models. (i) LVLms with 7B parameter including Otter (Li et al., 2023a), MiniGPT-4 (Zhu et al., 2023), InstructBLIP (Dai et al., 2023), Qwen-VL (Bai et al., 2023), VILA (Lin et al., 2024), PandaGPT (Su et al., 2023), mPLUG-Owl2 (Ye et al., 2024), BLIP-2 (Li et al., 2023b) and LLaVA-1.6 (Liu et al., 2024b). (ii) LVLms with 13B parameter including VILA, InstructBLIP, BLIP-2, and LLaVA-1.6. (iii) LVLms with 34B parameter including LLaVA-1.6. (iv) Closed-source model including GPT-4V (OpenAI, 2023).

**Evaluation Metrics.** We evaluate the performance of different baselines using multi-class classification, categorizing data into four distinct classes: textual veracity distortion, visual veracity distortion, cross-modal consistency distortion, and real class. Consistent with (Qian et al., 2021a; Zhang & Gao, 2023; Chen et al., 2023), we adopt the widely-used macro-F1 metric, which balances precision and recall through a harmonic mean. Beyond the F1 score, we also include macro-precision, macro-recall,

Table 2: Overall results (%) of different models on the MMFakeBench validation and test set with the comparison of standard prompting (Standard) and proposed MMD-Agent framework.

Model Name	Language Model	Prompt Method	Validation (1000)				Test (10000)			
			F1↑	Precision↑	Recall↑	ACC↑	F1↑	Precision↑	Recall↑	ACC↑
<i>Human Evaluation</i>			35.9	38.3	38.9	37.9	-	-	-	-
<b>LVLMs with 7B Parameter</b>										
Otter-Image	MPT-7B	Standard	5.2	10.5	3.4	4.1	4.9	9.3	3.3	4.0
MiniGPT4	Vicuna-7B	Standard	5.2	5.2	21.2	9.0	5.3	6.9	21.0	9.1
InstructBLIP	Vicuna-7B	Standard	7.1	7.9	6.5	7.8	8.1	16.4	7.2	8.5
Qwen-VL	Qwen-7B	Standard	7.5	10.3	24.3	11.0	8.0	35.9	25.5	11.6
VILA	LLaMA2-7B	Standard	11.5	7.5	25.0	30.0	11.5	7.5	25.0	30.0
PandaGPT	Vicuna-7B	Standard	11.8	9.8	25.0	30.0	11.6	8.6	25.0	30.0
mPLUG-Owl2	LLaMA2-7B	Standard	14.5	22.2	25.9	31.1	15.1	25.2	26.3	31.5
BLIP2	FlanT5-XL	Standard	16.4	20.1	27.5	33.0	16.7	17.3	27.7	33.2
LLaVA-1.6	Vicuna-7B	Standard	17.4	14.8	25.7	30.8	19.0	16.5	26.9	32.3
<b>LVLMs with 13B Parameter</b>										
VILA	LLaMA2-13B	Standard	11.5	7.5	<b>25.0</b>	<b>30.0</b>	11.6	<b>32.5</b>	25.0	<b>30.0</b>
		MMD-Agent	<b>22.7</b>	<b>27.3</b>	24.4	28.7	<b>24.0</b>	30.4	<b>25.5</b>	29.4
InstructBLIP	Vicuna-13B	Standard	13.7	13.2	24.0	28.8	13.9	25.5	24.3	<b>29.1</b>
		MMD-Agent	<b>26.0</b>	<b>33.3</b>	<b>30.1</b>	<b>29.5</b>	<b>24.5</b>	<b>32.1</b>	<b>28.8</b>	27.3
BLIP2	FlanT5-XXL	Standard	16.7	34.9	27.3	32.8	16.3	34.6	27.3	<b>32.8</b>
		MMD-Agent	<b>31.6</b>	<b>39.8</b>	<b>32.2</b>	<b>34.4</b>	<b>28.8</b>	<b>39.0</b>	<b>30.4</b>	32.1
LLaVA-1.6	Vicuna-13B	Standard	12.0	22.5	25.0	30.0	14.4	35.7	26.0	31.2
		MMD-Agent	<b>38.0</b>	<b>44.5</b>	<b>41.0</b>	<b>40.6</b>	<b>34.5</b>	<b>42.7</b>	<b>37.5</b>	<b>37.4</b>
<b>LVLMs with 34B Parameter</b>										
LLaVA-1.6	Nous-Hermes-2 -Yi-34B	Standard	25.7	44.5	33.7	40.4	25.4	44.1	33.8	40.5
		MMD-Agent	<b>49.9</b>	<b>54.4</b>	<b>52.9</b>	<b>48.7</b>	<b>47.7</b>	<b>52.1</b>	<b>49.6</b>	<b>46.6</b>
<b>Proprietary LVLMs</b>										
GPT-4V	ChatGPT	Standard	51.0	66.8	49.7	54.0	48.8	63.0	48.7	54.2
		MMD-Agent	<b>61.6</b>	<b>67.8</b>	<b>59.3</b>	<b>62.1</b>	<b>61.5</b>	<b>67.7</b>	<b>59.1</b>	<b>61.0</b>

and macro-accuracy as complementary evaluation metrics. Specifically, we construct robust regular expressions to extract key phrases from the long responses for accurate answer matching. Following (Liu et al., 2023b), if a model’s response lacks a valid answer, we classify it as a pseudo choice “Z” and consider the response incorrect.

## 5.2 MAIN RESULTS

**Comparison of Different LVLMs.** We present a comprehensive comparison of different LVLMs using the MMFakeBench, detailed in Table 2. Our key findings are summarized as follows:

**1) Challenges of MMFakeBench:** The benchmark poses substantial challenges to current models. Notably, GPT-4V, despite its advancement, achieves an F1-score of only 51.0% with the standard prompting. This indicates considerable room for improvement and highlights the rigorous standards of this benchmark.

**2) Disparity between Open-source Models and GPT-4V:** Although LLaVA-1.6-34b is the leading open-source model, it achieves an F1-score of just 25.7% with the standard prompting, significantly lower than GPT-4V. This highlights a pronounced disparity in detection capabilities between open-source and proprietary models.

**3) Impact of Parameter Quantity:** Comparing models within the same series, such as LLaVA-1.6-Vicuna-7b and LLaVA-1.6-34b, we observe that models with larger parameter counts exhibit better performance. Smaller LVLMs face constraints in instruction-following and high predicted consistency, as detailed in the Appendix A.1.1. These results indicate that 7B parameter models lack sufficient multimodal understanding to effectively combat misinformation.

**4) Effectiveness of MMD-Agent:** Due to the limited reasoning capability of small-scale models, we select moderately sized open-source and proprietary models as baselines to compare the proposed MMD-Agent with the standard prompting. MMD-Agent significantly improves the F1-score for both open-source models and GPT-4V. Notably, LLaVA-1.6-34B using MMD-Agent achieves an F1-score of 49.9%, approaching the 51% score of GPT-4V with the standard prompting. This suggests that MMD-Agent can serve as a general framework for future research on the MMFakeBench benchmark.



Table 3: (a) Comparison with single-source detectors for MMFakeBench. (b) Ablation studies on hierarchical (Hier.) decomposition and reasoning knowledge  $\mathcal{K}_t = \{\mathcal{R}_t, \mathcal{O}_t\}$  of each sub-task  $t$ . TVD, VVD, and CCD denote textual veracity distortion, visual veracity distortion, and cross-modal consistency distortion, respectively. Corpus refers to the general datasets used in LVLMs, not tailored for misinformation detection. \* denotes the chosen single-source detector applied for mixed detection.

(a)				(b)								
Existing Detector	Train Source	Binary Overall $\uparrow$	Multiclass Overall $\uparrow$	Hier.	$\mathcal{K}_1$	$\mathcal{K}_2$	$\mathcal{K}_3$	Real $\uparrow$	TVD $\uparrow$	VVD $\uparrow$	CCD $\uparrow$	Overall $\uparrow$
FakingFakeNews*	TVD	37.8	-					58.5	5.8	0.0	38.6	25.7
CNNSpot	VVD	23.8	-	✓				45.8 ( $\downarrow$ 12.7)	16.5 ( $\uparrow$ 10.7)	32.5 ( $\uparrow$ 32.5)	49.0 ( $\uparrow$ 10.4)	36.0
UnivFD	VVD	28.9	-	✓	✓			46.4 ( $\uparrow$ 0.6)	37.6 ( $\uparrow$ 21.1)	32.8 ( $\uparrow$ 0.3)	47.2 ( $\downarrow$ 1.8)	41.0
LNP*	VVD	33.0	-	✓	✓	✓		46.8 ( $\uparrow$ 0.4)	37.6 ( $\uparrow$ 0.0)	61.7 ( $\uparrow$ 28.9)	48.6 ( $\uparrow$ 1.4)	48.7
FakeNewsGPT4	CCD	41.7	-	✓	✓	✓	✓	51.1 ( $\uparrow$ 4.3)	37.6 ( $\uparrow$ 0.0)	61.7 ( $\uparrow$ 0.0)	49.2 ( $\uparrow$ 0.6)	49.9
HAMMER*	CCD	43.0	-	✓	✓	✓	✓					
Mixed Detection	-	47.6	22.5									
LLaVA-1.6-34B	Corpus	67.2	49.9					W/o Wiki Knowledge 49.7	18.0	61.0	46.3	43.8
GPT-4V	Corpus	<b>74.0</b>	<b>61.6</b>					W/ Google Knowledge 50.2	33.4	62.2	48.1	48.5

**Comparison with Single-source Detectors.** We compare LLaVA-1.6-34B and GPT-4V utilized in MMD-Agent, with several competitive single-source detectors including FakeingFakeNews (Huang et al., 2023), CNNSpot (Wang et al., 2020), UnivFD (Ojha et al., 2023b), LNP (Liu et al., 2022), FakeNewsGPT4 (Liu et al., 2024c), and HAMMER (Shao et al., 2023). The details of each detector are presented in the Appendix A.3.2. For a fair comparison, in addition to the single-source misinformation detection via existing detectors, we integrate the three most powerful detectors in distinct sources (i.e., FakeingFakeNews, LNP, and HAMMER) to assess the capability of mixed detection. Mixed detection utilizes our proposed hierarchical framework by replacing LVLMs with relevant detectors. The results in Table 3 (a) show that LLaVA-1.6-34B and GPT-4V perform better than single-source detectors for both binary and multiclass classification by a large margin. This demonstrates that LVLMs trained with a large general corpus achieve promising generalization performance in mixed-source MMD and can serve as potential baseline models for future study.

**Results of Human Evaluation.** We conduct a comprehensive user study using a validation set containing 1,000 samples as a question bank. For each questionnaire, 50 samples are randomly selected from this question bank. A total of 80 participants are asked to identify the source of misinformation for each news item. As shown in the first row of Table 2, the results reveal that 62.1% of the items are predicted incorrectly by users, highlighting the dataset’s high confusion and the realistic challenge it poses. This level of difficulty and confusion emphasizes the quality and challenge embedded in our dataset, making it an invaluable resource for pushing the boundaries of current understanding and capabilities in multimodal misinformation detection research.

5.3 EXPERIMENTAL ANALYSIS

**Ablation Study on Hierarchical Decomposition and Reasoning Knowledge.** We first investigate the effects of instructing LVLMs using only hierarchical decomposition compared to standard prompting. In Table 3 (b), the decomposition method performs better for solving multi-task interference. Additionally, we conduct an ablation study by sequentially generating multi-perspective knowledge for individual sub-tasks. Results in Table 3 (b) show that augmenting decisions with reasoning knowledge outperforms its ablation part, especially for checking content veracity. We further conduct ablations on external knowledge by removing it or integrating an alternative source (i.e., Google Knowledge Graph). Results show that models using external knowledge outperform those relying solely on internal reasoning, highlighting its critical role in validating textual veracity.

Table 4: Performance (F1 score (%)) of models on different sources of misinformation.

Model	Real $\uparrow$	TVD $\uparrow$	VVD $\uparrow$	CCD $\uparrow$	Overall $\uparrow$
VILA-13B	32.4	13.4	4.3	37.6	21.9
InstructBLIP-13B	41.9	18.8	19.6	23.8	26.0
BLIP2-FLAN-T5-XXL	41.5	39.2	13.1	32.6	31.6
LLaVA-1.6-34B	51.1	37.6	61.7	49.2	49.9
GPT-4V	<b>65.3</b>	<b>67.2</b>	<b>57.3</b>	<b>56.5</b>	<b>61.6</b>

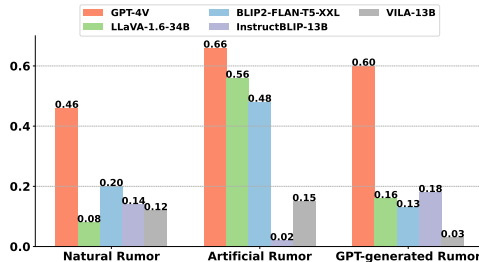


Figure 5: Performance (detection success rate) of models on different types of rumor.

**Analysis of Misinformation Sources and Types.** We compare the F1 scores of various LVLMS across misinformation sources in Table 4. Across all sources, the majority of open-source models perform worse than GPT-4V by a huge margin, particularly in terms of textual veracity distortion. This indicates that open-source models are considerably challenged by textual veracity distortion. Within textual veracity distortion, we further report the detection accuracy of selected models across three types of textual rumors. Fig. 5 shows that open-source models typically perform better in the artificial rumor. This might be attributed to the fact that artificial rumors are constructed by applying specific rules, which provide identifiable points of falsification enabling verification through external knowledge. In contrast, natural or GPT-generated rumors often leverage ambiguous language that makes detection challenging. Analysis of other misinformation sources is detailed in Appendix A.1.3.

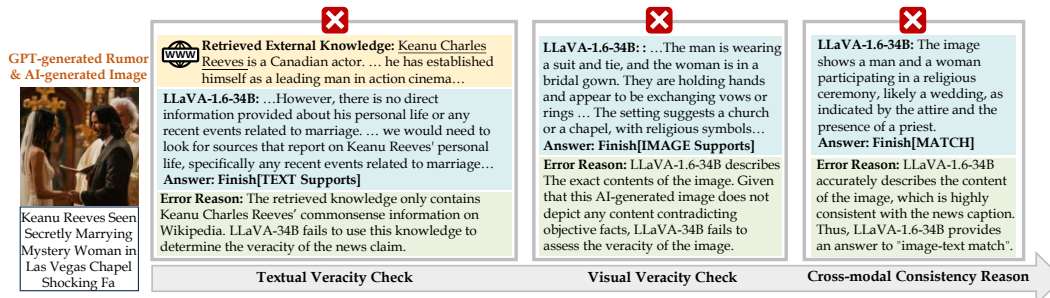


Figure 6: One of the most harmful examples involves a GPT-generated rumor supported by an AI-generated image, which is challenging for LLaVA-1.6-34B. More examples can be found in the Appendix A.1.7.

**Error Analysis.** A fundamental root cause of textual veracity checking errors in the LLaVA-1.6-34B is the lack of useful external knowledge. This deficiency is exemplified in Fig. 6, where the knowledge contains only commonsense information but fails to provide relevant events. Moreover, this AI-generated image in the image-text pair exhibits high fidelity and strong coherence, thus evading detection in visual veracity and cross-modal consistency. These instances underscore the dangers of using collaborative generative modes to automatically generate multimodal misinformation.

## 6 CONCLUSION

In this paper, we introduce MMFakeBench, the first comprehensive benchmark for detecting mixed-source multimodal misinformation. MMFakeBench contains three primary misinformation categories along with 12 sub-categories of forgery types. We conduct comprehensive evaluations of 6 prevalent detection methods and 15 LVLMS on the MMFakeBench dataset. Furthermore, we propose an innovative unified framework and perform extensive experiments to demonstrate its effectiveness.

## ETHICS STATEMENT AND LIMITATIONS

This paper contains examples of harmful texts or images, raising concerns about the manipulation of public safety. To reduce its social impact, we have implemented several safeguards: (1) Content Safeguards. Rigorous review sensitive content in data generation such as those related to politics and race etc. (2) Disable data generation code access. We open-source datasets and detection codes but do not release data generation codes for safety. (3) Data Access Restrictions. Data access is restricted to verified researchers by following a binding usage agreement. (4) Public Feedback Mechanism. We will offer a public feedback channel for ethical concerns and continuous dataset improvement. The licenses for the datasets contributed in this work are discussed in Appendix A.5.

While our MMFakeBench marks a critical advancement in mixed-source multimodal misinformation detection, it is important to recognize certain limitations. Our proposed framework utilizes external knowledge retrieved from the Wikipedia API. While the integration of such external knowledge has enhanced the performance of our baseline models, it may not always provide useful information for particularly challenging natural rumors and GPT-generated rumors. Future research should explore the use of a more advanced retrieval augmentation generation (RAG), which could lead to further performance improvements.

## REFERENCES

- 540 Sahar Abdelnabi, Rakibul Hasan, and Mario Fritz. Open-domain, content-based, multi-modal  
541 fact-checking of out-of-context images via online resources. In *CVPR*, 2022.
- 542 Shivangi Aneja, Chris Bregler, and Matthias Nießner. Cosmos: Catching out-of-context misinforma-  
543 tion with self-supervised learning. In *AAAI*, 2023.
- 544 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang  
545 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.  
546 *arXiv preprint arXiv:2308.12966*, 2023.
- 547 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image  
548 editing instructions. In *CVPR*, 2023.
- 549 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
550 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
551 few-shot learners. In *NeurIPS*, 2020.
- 552 Canyu Chen and Kai Shu. Can llm-generated misinformation be detected? In *ICLR*, 2024.
- 553 Ziwei Chen, Linmei Hu, Weixin Li, Yingxia Shao, and Liqiang Nie. Causal intervention and  
554 counterfactual reasoning for multi-modal fake news detection. In *ACL*, 2023.
- 555 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
556 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot  
557 impressing gpt-4 with 90%\* chatgpt quality. <https://vicuna.lmsys.org>, 2023.
- 558 Jeff Da, Maxwell Forbes, Rowan Zellers, Anthony Zheng, Jena D Hwang, Antoine Bosselut, and  
559 Yejin Choi. Edited media understanding frames: Reasoning about the intent and implications of  
560 visual misinformation. In *ACL-IJCNLP*, 2021.
- 561 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,  
562 Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language  
563 models with instruction tuning. In *NeurIPS*, 2023.
- 564 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*,  
565 2021.
- 566 Bilal Ghanem, Simone Paolo Ponzetto, Paolo Rosso, and Francisco Rangel. Fakeflow: Fake news  
567 detection by modeling the flow of affective information. In *EACL*, 2021.
- 568 Andreas Hanselowski, Christian Stab, Claudia Schulz, Zile Li, and Iryna Gurevych. A richly  
569 annotated corpus for different tasks in automated fact-checking. In *CoNLL*, 2019.
- 570 Kung-Hsiang Huang, Kathleen McKeown, Preslav Nakov, Yejin Choi, and Heng Ji. Faking fake  
571 news for real fake news detection: Propaganda-loaded training data generation. In *ACL*, 2023.
- 572 Ayush Jaiswal, Ekraam Sabir, Wael AbdAlmageed, and Premkumar Natarajan. Multimedia semantic  
573 integrity assessment using joint embedding of images and text. In *ACM MM*, 2017.
- 574 Tiep Le, Vasudev Lal, and Phillip Howard. Coco-counterfactuals: Automatically constructed  
575 counterfactual examples for image-text pairs. In *NeurIPS*, 2023.
- 576 Nayeon Lee, Yejin Bang, Andrea Madotto, Madian Khabsa, and Pascale Fung. Towards few-shot  
577 fact-checking via perplexity. In *NACCL*, 2021.
- 578 Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A  
579 multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023a.
- 580 Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven  
581 Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum  
582 distillation. In *NeurIPS*, 2021.
- 583
- 584
- 585
- 586
- 587
- 588
- 589
- 590
- 591
- 592
- 593

- 594 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image  
595 pre-training with frozen image encoders and large language models. In *ICML*, 2023b.  
596
- 597 Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz,  
598 Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models. In *CVPR*,  
599 2024.
- 600 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
601 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.  
602
- 603 Bo Liu, Fan Yang, Xiuli Bi, Bin Xiao, Weisheng Li, and Xinbo Gao. Detecting generated images by  
604 real images. In *ECCV*, 2022.
- 605 Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visual news: Benchmark and  
606 challenges in news image captioning. In *EMNLP*, 2021.  
607
- 608 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*,  
609 2023a.
- 610 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction  
611 tuning. In *CVPR*, 2024a.  
612
- 613 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.  
614 Llava-next: Improved reasoning, ocr, and world knowledge. <https://llava-vl.github.io/blog/2024-01-30-llava-next>, 2024b.  
615
- 616 Xuannan Liu, Peipei Li, Huaibo Huang, Zekun Li, Xing Cui, Jiahao Liang, Lixiong Qin, Weihong  
617 Deng, and Zhaofeng He. Fakenews4: Advancing multimodal fake news detection through  
618 knowledge-augmented lvlms. *arXiv preprint arXiv:2403.01988*, 2024c.
- 619 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
620 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining  
621 approach. *arXiv preprint arXiv:1907.11692*, 2019.  
622
- 623 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi  
624 Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player?  
625 In *NeurIPS*, 2023b.
- 626 Grace Luo, Trevor Darrell, and Anna Rohrbach. Newsclippings: Automatic generation of out-of-  
627 context multimodal media. In *EMNLP*, 2021.  
628
- 629 Midjourney. <https://www.midjourney.com/home/>, 2022.
- 630 Erxue Min, Yu Rong, Yatao Bian, Tingyang Xu, Peilin Zhao, Junzhou Huang, and Sophia Ananiadou.  
631 Divide-and-conquer: Post-user interaction network for fake news detection on social media. In  
632 *WWW*, 2022.  
633
- 634 Kai Nakamura, Sharon Levy, and William Yang Wang. r/fakeddit: A new multimodal benchmark  
635 dataset for fine-grained fake news detection. In *LREC*, 2020.
- 636 Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize  
637 across generative models. In *CVPR*, 2023a.  
638
- 639 Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize  
640 across generative models. In *CVPR*, 2023b.
- 641 OpenAI. Gpt-4v(ision) system card. <https://api.semanticscholar.org/CorpusID:263218031>, 2023.  
642
- 643 Piotr Przybyla. Capturing the style of fake news. In *AAAI*, 2020.
- 644 Peng Qi, Zehong Yan, Wynne Hsu, and Mong Li Lee. Sniffer: Multimodal large language model for  
645 explainable out-of-context misinformation detection. In *CVPR*, 2024.  
646
- 647 Chen Qian, Fuli Feng, Lijie Wen, Chunping Ma, and Pengjun Xie. Counterfactual inference for text  
classification debiasing. In *ACL*, 2021a.

- 648 Shengsheng Qian, Jinguang Wang, Jun Hu, Quan Fang, and Changsheng Xu. Hierarchical multi-  
649 modal contextual attention network for fake news detection. In *SIGIR*, 2021b.
- 650  
651 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
652 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 653 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
654 resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- 655 Ekraam Sabir, Wael AbdAlmageed, Yue Wu, and Prem Natarajan. Deep multimodal image-  
656 repurposing detection. In *ACM MM*, 2018.
- 657  
658 Rui Shao, Tianxing Wu, and Ziwei Liu. Detecting and grounding multi-modal media manipulation.  
659 In *CVPR*, 2023.
- 660  
661 Rui Shao, Tianxing Wu, Jianlong Wu, Liqiang Nie, and Ziwei Liu. Detecting and grounding  
662 multi-modal media manipulation and beyond. *TPAMI*, 2024.
- 663 Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data  
664 repository with news content, social context, and spatiotemporal information for studying fake  
665 news on social media. *Big data*, 2020.
- 666  
667 Gabriela Ben Melech Stan, Raanan Yehezkel Rohekar, Yaniv Gurwicz, Matthew Lyle Olson,  
668 Anahita Bhiwandiwalla, Estelle Aflalo, Chenfei Wu, Nan Duan, Shao-Yen Tseng, and Vasudev  
669 Lal. Lvlm-intrepret: An interpretability tool for large vision-language models. *arXiv preprint*  
670 *arXiv:2404.03118*, 2024.
- 671  
672 Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to  
673 instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023.
- 674  
675 James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale  
676 dataset for fact extraction and verification. In *NAACL*, 2018.
- 677  
678 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
679 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
680 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 681  
682 Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated  
683 images are surprisingly easy to spot... for now. In *CVPR*, 2020.
- 684  
685 Yang Wu, Pengwei Zhan, Yunjian Zhang, Liming Wang, and Zhen Xu. Multimodal fusion with  
686 co-attention networks for fake news detection. In *Findings of ACL-IJCNLP*, 2021.
- 687  
688 Barry Menglong Yao, Aditya Shah, Lichao Sun, Jin-Hee Cho, and Lifu Huang. End-to-end multi-  
689 modal fact-checking and explanation generation: A challenging dataset and models. In *SIGIR*,  
690 2023.
- 691  
692 Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu,  
693 Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with  
694 multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- 695  
696 Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and  
697 Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality  
698 collaboration. In *CVPR*, 2024.
- 699  
700 Qichao Ying, Xiaoxiao Hu, Yangming Zhou, Zhenxing Qian, Dan Zeng, and Shiming Ge. Bootstrap-  
701 ping multi-view representations for fake news detection. In *AAAI*, 2023.
- 702  
703 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,  
704 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal  
705 understanding and reasoning benchmark for expert agi. In *CVPR*, 2024.
- 706  
707 Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot:  
708 Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs.  
709 In *CVPR*, 2024.

702 Xuan Zhang and Wei Gao. Towards llm-based fact verification on news claims with a hierarchical  
703 step-by-step prompting method. In *AAACL*, 2023.  
704

705 Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibeil Yang. Ddcot: Duty-distinct chain-of-  
706 thought prompting for multimodal reasoning in language models. In *NeurIPS*, 2023.

707 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: En-  
708 hancing vision-language understanding with advanced large language models. *arXiv preprint*  
709 *arXiv:2304.10592*, 2023.  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A APPENDIX

This appendix contains additional details for the ICLR 2025 submission, titled “MMFakeBench: A Mixed-Source Multimodal Misinformation Detection Benchmark for LVLMS”. The appendix is organized as follows:

- §A.1 Additional Experimental Result.
  - §A.1.1 Instruction Following and Prediction Consistency.
  - §A.1.2 Evaluation Results for Binary Classification.
  - §A.1.3 Analysis of More Misinformation Types.
  - §A.1.4 Comparison of Existing Datasets.
  - §A.1.5 Efficiency of MMD-Agent.
  - §A.1.6 Fine-tuning on Existing Detectors.
  - §A.1.7 [More Error Analysis](#).
  - §A.1.8 [Interpretable visualization](#).
- §A.2 Benchmark Analysis.
- §A.3 Implementation Details.
  - §A.3.1 LVLMS Configuration Details.
  - §A.3.2 Existing Single-source Detectors Details.
- §A.4 Instruct Prompts for ChatGPT.
- §A.5 Dataset Licenses.
- §A.6 More Visualization Examples.

### A.1 ADDITIONAL EXPERIMENTAL RESULTS

Table 5: Statistics of Instruction following capabilities and predicted consistency tendency of LVLMS.

Model	Match	Consist.	Model	Match	Consist.	Model	Match	Consist.
<b>LVLMS with 7B Parameter</b>								
Otter-Image	9.8	100	Qwen-VL	88.6	92.9	mPLUG-Owl2	100	96.6
MiniGPT4	100	88.6	VILA	100	100	BLIP2	100	93.6
InstructBLIP	25.7	96.11	PandaGPT	100	98.9	LLaVA-1.6	100	76.9
<b>LVLMS with 13B Parameter</b>								
VILA	100	100	InstructBLIP	99.9	91.0	BLIP2-FlanT5-XXL	100	49.4
<b>LVLMS with 34B Parameter</b>								
LLaVA-1.6	100	52.6	-	-	-	-	-	-
<b>Proprietary LVLMS</b>								
GPT-4V	99.9	36.5	-	-	-	-	-	-

#### A.1.1 INSTRUCTION FOLLOWING AND PREDICTION CONSISTENCY.

We evaluated the instruction following capabilities and prediction consistency to further study the multimodal understanding of LVLMS on mixed-source multimodal misinformation detection (MMD). We report the success rate in heuristic matching (Match) with regular expressions and prediction consistency rate (Consist.). The results are shown in the Table 5. Among all LVLMS, small-scale models like Otter-Image and InstructBLIP, achieve the lower matching success rate. While there exist small-scale LVLMS that perfectly follow the format of the regular expressions and achieve high success rates (>99%) in matching, most small-scale models exhibit high predicted consistency rates. This indicates small-scale models may prefer to predict a certain category answer among all given choices. Additionally, the leading open-source model, LLaVA-1.6-34B, and the proprietary model GPT-4V demonstrate superior instruction-following capabilities and lower prediction consistency. This indicates their significant potential in addressing mixed-source MMD, positioning them as valuable baseline models.

Table 6: Binary overall results of different models on the MM-FakeBench validation and test set with the comparison of standard prompting (Standard) and proposed MMD-Agent framework.

Model Name	Language Model	Prompt Method	Validation (1000)				Test (10000)			
			F1	Precision	Recall	ACC	F1	Precision	Recall	ACC
<i>Human Evaluation</i>			54.9	56.6	57.8	56.8	-	-	-	-
<b>LVLMs with 7B Parameter</b>										
Otter-Image	MPT-7B	Standard	7.9	4.1	4.5	7.9	8.6	32.4	5.0	8.6
MiniGPT4	Vicuna-7B	Standard	40.4	38.2	45.7	63.1	41.7	41.0	47.4	65.2
InstructBLIP	Vicuna-7B	Standard	14.7	30.8	13.2	8.1	16.1	40.5	14.2	8.8
Qwen-VL	Qwen-7B	Standard	43.6	50.6	44.9	60.3	44.0	51.6	45.2	60.5
VILA	LLaMA2-7B	Standard	41.2	35.0	50.0	70.0	41.2	35.0	50.0	70.0
PandaGPT	Vicuna-7B	Standard	24.6	60.6	50.5	30.9	24.1	61.7	50.4	30.6
mPLUG-Owl2	LLaMA2-7B	Standard	47.2	64.9	52.3	70.6	48.7	71.1	53.3	71.4
BLIP2	FlanT5-XL	Standard	41.2	35.0	50.0	70.0	41.2	35.0	50.0	70.0
LLaVA-1.6	Vicuna-7B	Standard	48.1	48.2	48.5	59.5	52.5	53.0	52.6	62.5
<b>LVLMs with 13B Parameter</b>										
VILA	LLaMA2-13B	Standard	41.1	35.0	50.0	70.0	41.1	35.0	50.0	70.0
		MMD-Agent	<b>56.5</b>	<b>62.2</b>	<b>56.9</b>	<b>70.3</b>	<b>56.6</b>	<b>64.3</b>	<b>57.2</b>	<b>71.2</b>
InstructBLIP	Vicuna-13B	Standard	41.1	35.0	49.9	<b>69.9</b>	41.1	35.0	49.9	<b>69.8</b>
		MMD-Agent	<b>51.3</b>	<b>53.4</b>	<b>54.0</b>	53.1	<b>47.9</b>	<b>50.1</b>	<b>50.1</b>	49.9
BLIP2	FlanT5-XXL	Standard	31.6	<b>63.4</b>	53.6	35.5	30.6	<b>64.9</b>	53.4	34.9
		MMD-Agent	<b>51.5</b>	53.4	<b>54.0</b>	<b>53.6</b>	<b>51.8</b>	54.0	<b>54.7</b>	<b>53.5</b>
LLaVA-1.6	Vicuna-13B	Standard	41.1	35.0	50.0	69.7	42.3	57.3	50.1	69.5
		MMD-Agent	<b>51.8</b>	<b>66.7</b>	<b>54.6</b>	<b>71.4</b>	<b>50.2</b>	<b>67.3</b>	<b>53.9</b>	<b>71.3</b>
<b>LVLMs with 34B Parameter</b>										
LLaVA-1.6	Nous-Hermes-2	Standard	62.9	67.1	<b>70.0</b>	63.4	64.3	68.8	<b>71.7</b>	64.8
	-Yi-34B	MMD-Agent	<b>67.2</b>	<b>70.4</b>	66.0	<b>75.1</b>	<b>68.1</b>	<b>71.1</b>	67.0	<b>75.6</b>
<b>Proprietary LVLMs</b>										
GPT-4V	ChatGPT	Standard	72.3	72.1	72.8	75.6	<b>74.2</b>	<b>73.5</b>	<b>76.9</b>	<b>76.4</b>
		MMD-Agent	<b>74.0</b>	<b>73.4</b>	<b>75.5</b>	<b>76.8</b>	72.8	72.4	75.4	75.0

### A.1.2 EVALUATION RESULTS FOR BINARY CLASSIFICATION.

In addition to multi-class classification, we also provide binary classification performance to assess the overall detection capability of baseline models in mixed-source MMD. Based on the 4 categories in the mixed-source MMD settings, we develop binary evaluation metrics via mapping techniques. Specifically, we standardize the assignment of labels denoting “textual veracity distortion”, “visual veracity distortion”, and “cross-modal consistency distortion” to the classification of “Fake” while reserving the label “True” to denote real data. Similar to the multi-classification evaluation, we adopt the widely used F1 score. In addition to the F1 score, we also use precision, recall, and accuracy as supplementary evaluation metrics. The specific evaluation results are shown in Table 6. From the results, we make the following observations:

- Current models including open-source models and GPT-4V are challenged by the MMFakeBench dataset in binary classification detection. Despite being an advanced modal, GPT-4V attains a mere F1-score of 72.3% using the standard prompting on the validation set.
- The proposed framework MMD-Agent yields substantial improvement on recent LVLMs, especially on open-source models. For instance, for the BLIP2-FlanT5-XXL model, MMD-Agent achieves a 19.9% increase in F1 score on the validation set. This may be credited to the effective integration of reasoning, actions, and tool use in enhancing multimodal understanding in mixed-source MMD.

### A.1.3 ANALYSIS OF MORE MISINFORMATION TYPES.

Within visual veracity distortion, we report the detection accuracy of selected models across two types of fact-conflicting images. Fig. 7 (a) shows the challenging nature of both types of such fact-conflicting images for existing models. Notably, even the advanced GPT-4V achieves a detection success rate of less than 50% on both types. Additionally, we present an analysis of the detection accuracy of selected models across different types of image-text inconsistency. As shown in Fig. 7 (b),



edited inconsistency emerges as a more substantial challenge compared to repurposed inconsistency. This finding suggests that editing methods introduce minor alterations to images or text, necessitating enhanced multimodal reasoning capabilities.

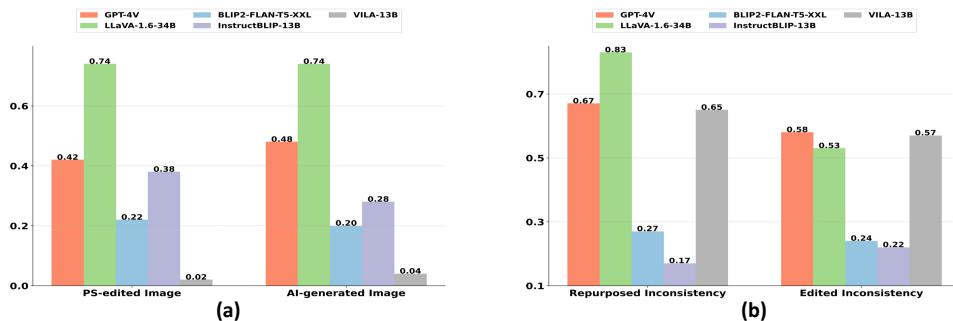


Figure 7: (a) Performance (detection success rate) of models on different types of fact-conflicting images. (b) Performance of models on different types of inconsistent image-text pairs.

Table 7: Performance (F1 score) comparison of selected models on the proposed benchmark and existing benchmarks (datasets).

Model Name	NewsClippings (Binary)	NewsClippings (Multiclass)	Fakeddit (Binary)	Fakeddit (Multiclass)	MMFakeBench (Binary)	MMFakeBench (Multiclass)
<b>Large Vision Language Models</b>						
BLIP2-FLanT5-XXL	33.6	-	40.5	-	31.6	16.7
LLaVA-1.6-34B	65.6	-	71.4	-	62.9	25.7
<b>Multimodal Specialized Detectors</b>						
FKA-Owl	52.0	-	55.0	-	41.7	-
HAMMER	55.0	-	52.6	-	43	-

#### A.1.4 COMPARISON OF EXISTING DATASETS

We conduct an experiment to compare the performance of selected models on the proposed benchmark against existing benchmarks (datasets) in Table 7. We employ four models: two powerful open-source LVM models, LLaVA-1.6-34B and BLIP2-FLAN-T5-XXL with two open-source multimodal specialized detectors, HAMMER (Shao et al., 2023) and FKA-Owl (Liu et al., 2024c). All these models are evaluated on two widely-used multimodal misinformation datasets, NewsClippings (Luo et al., 2021) and Fakeddit (Nakamura et al., 2020). The results indicate that our MMFakeBench poses a greater challenge in binary classification compared to the other two datasets. More importantly, the primary challenge of our benchmark lies in its capacity to perform fine-grained assessments of misinformation sources in scenarios where multiple forgery types coexist. This mirrors the complexity of real-world environments, where misinformation stems from diverse, overlapping sources. This capability is essential for addressing real-world challenges and underscores the importance of MMFakeBench in advancing multimodal misinformation detection.

#### A.1.5 EFFICIENCY OF MMD-AGENT

In Table 8, we compare the inference time and computational resource usage of MMD-Agent with standard prompting methods. Overall, while MMD-Agent shows an increase in inference time, it does not lead to additional GPU memory consumption and significantly enhances

both the performance and interpretability of multimodal misinformation detection. The experimental results are shown in the table, and the detailed analysis is as follows: (1) MMD-Agent introduces higher inference time compared to standard prompting methods, primarily due to the additional reasoning steps, such as extracting key entities from the text. (2) MMD-Agent does not lead to additional GPU memory consumption. This is because the Agent method does not affect the model parameter deployment and dataset storage.

Table 8: Efficiency of MMD-Agent on LLaVA-1.6-34B.

Metric	Standard	MMD-Agent
Average Inference Time (s)	5.97s	49.04s
Memory (GB)	82G	82G
Performance (F1 score)	25.7	49.9

While the increased inference time is a consideration, the substantial gains brought by the MMD-Agent should not be overlooked. On the one hand, it greatly enhances the detection performance, with the F1 score improving from 25.7 to 49.9. On the other hand, MMD-Agent offers stronger interpretability because it offers a detailed analysis of the misinformation rather than only providing classification labels as used in standard prompting methods. Specifically, as shown in Fig. 4, the content of the intermediate reasoning traces such as Thought 2, Thought 3, and Thought 4 accurately analyzes the specific reasons for misinformation from different sources.

Table 9: (a) Results of the specialized detector, HAMMER, fine-tuning on the MMFakeBench dataset. (b) Results of the specialized detector, FKA-Owl, fine-tuning on the MMFakeBench dataset.

(a) Fine-tune on HAMMER Model.			(b) Fine-tune on FKA-Owl Model.		
Model	DGM4	MMFakeBench	Model	DGM4	MMFakeBench
HAMMER before fine-tuning	83.2	44.1	FKA-Owl before fine-tuning	78.5	44.6
HAMMER after fine-tuning (10)	78.7	46.8	FKA-Owl after fine-tuning (10)	78.3	46.9
HAMMER after fine-tuning (100)	70.1	60.8	FKA-Owl after fine-tuning (100)	66.4	61.1
HAMMER after fine-tuning (1000)	63.2	72.4	FKA-Owl after fine-tuning (1000)	64.8	76.2

### A.1.6 FINE-TUNING ON EXISTING DETECTORS

We conducted experiments in Table 9 to investigate the adaptation of existing models to the MM-FakeBench dataset. Specifically, we selected two existing powerful specialized detectors, HAMMER Shao et al. (2023) and LVLm-based detector, FKA-Owl Liu et al. (2024c). Then we fine-tuned them incrementally with 10, 100, and 1000 examples. We reported the results on both the DGM4 Shao et al. (2023) dataset (used for training and evaluating HAMMER and FKA-Owl) and our proposed dataset. The results indicate:

**1) Tuning-based models can be hard to generalize to unseen forgery data.** For off-the-shelf dedicated detectors without fine-tuning, their performance on the proposed benchmark dataset is notably poor. With the rapid development of generative models, new forgery techniques and synthesized data continue to emerge. Models trained on limited samples face significant challenges in generalizing to unseen types of forgery data, highlighting a critical issue in the field of fake detection Liu et al. (2024c); Ojha et al. (2023b).

**2) Catastrophic forgetting issues are inevitable.** Fine-tuning on new data improves performance on MMFakeBench but simultaneously leads to a decline in performance on the original dataset. This degradation underscores the phenomenon of catastrophic forgetting, where previously acquired knowledge, such as that from the DGM4 dataset, is progressively lost.

### A.1.7 MORE ERROR ANALYSIS

In Fig. 8, We have presented more error analysis that span four models of different scales (i.e., GPT-4V, LLaVA-1.6-34B, BLIP2-FLAN-T5-XXL, and VILA-13B) and on three distinct forgery sources of textual veracity distortions, visual veracity distortion and cross-modal consistency distortion. Based on these cases, we provide a deep analysis of the shortcomings of both largely and moderately sized models when encountering different sources of multimodal misinformation.

**1) For cases of textual veracity distortion,** our analysis is summarized as follows:

- **Reliance on external knowledge.** While Largely sized models like LLaVA-1.6-34B and GPT-4V exhibit strong reasoning capabilities, they are challenging to infer the factual correctness of a statement without access to a broader context. For instance, when faced with factually incorrect statements that appear linguistically accurate, such as the GPT-generated rumor, all these four models lack the inherent capability to question the information without retrieving corroborative evidence.

**2) For cases of visual veracity distortion,** our analysis is summarized as follows:

- **Limited Sensitivity to Abnormal Physical Features.** Models including GPT-4V, LLaVA-1.6-34B, BLIP2-FLAN-T5-XXL, and VILA-13B, face challenges in discerning abnormal physical characteristics. For instance, in the PS-edited examples, these four models fail to detect subtle yet unrealistic manipulations, such as swollen necks or distorted facial features in images of Donald

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

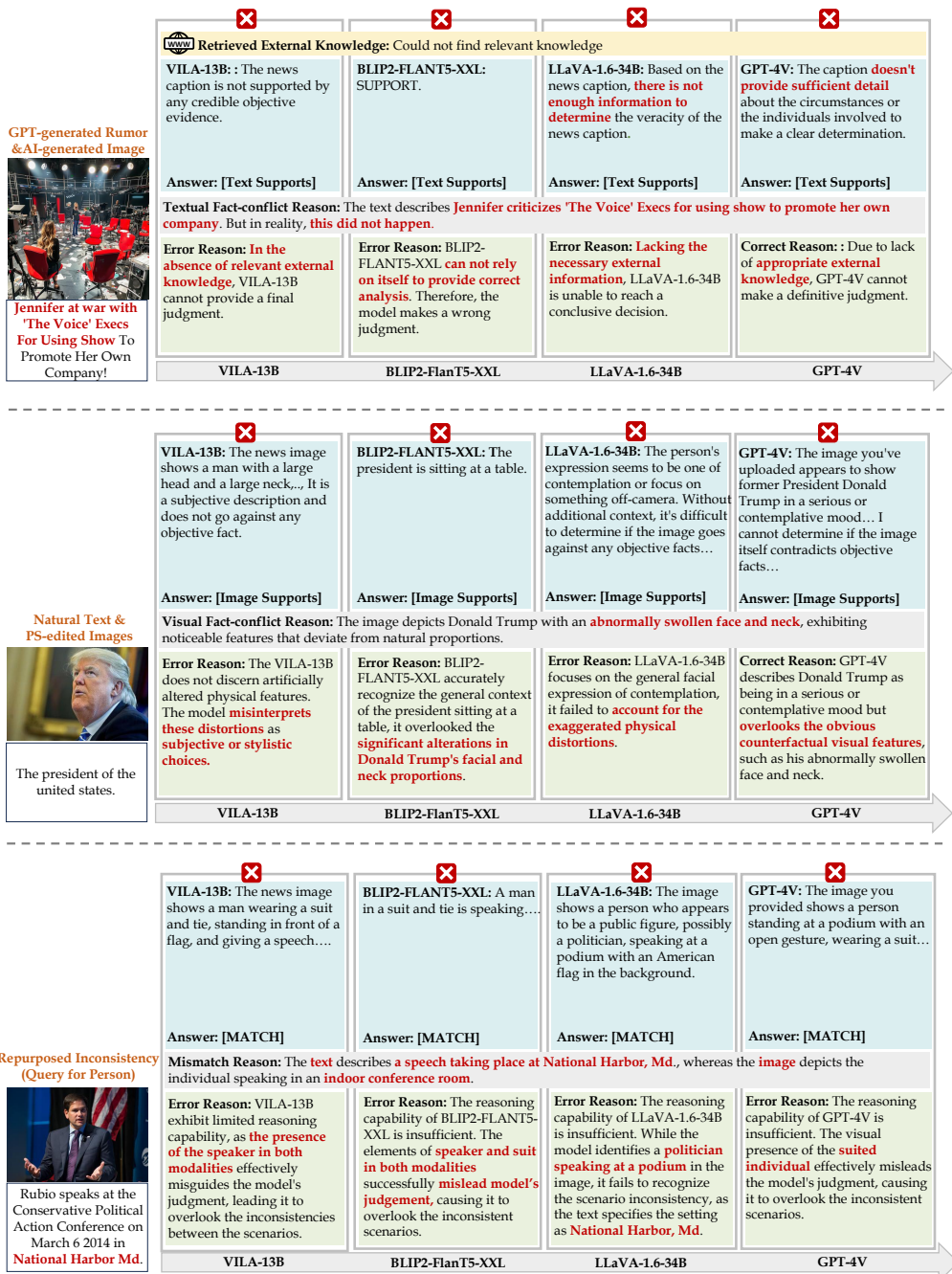


Figure 8: More Error analysis from three different forgery sources.

Trump. Instead, they are frequently distracted by more prominent visual elements, such as facial expressions or gestures, resulting in misjudgments when identifying these manipulations.

3) For cases of cross-modal consistency distortion, our analysis is summarized as follows:

- Distraction by global consistent semantics. When confronted with scenarios where images and text exhibit only subtle inconsistencies while other aspects remain largely consistent, large-scale models such as LLaVA-1.6-34B and GPT-4V often struggle to detect these discrepancies. These models can be distracted by the dominant presence of consistent content, which obscures the subtle mismatches critical for accurate misinformation detection. For instance, in the case of repurposed inconsistency, all four models focus on global semantics, such as a person delivering a speech, which diverts attention from deeper, subtle inconsistencies in the scenario.

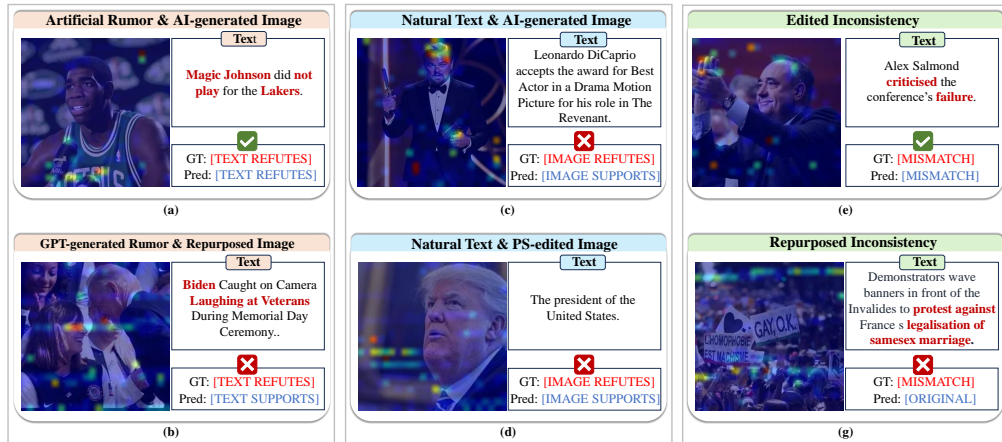


Figure 9: Illustration of relevancy maps to showcase interpretability for the predicted output of the LLaVA-13B model.

### A.1.8 INTERPRETABLE VISUALIZATION

We have incorporated the recent LVLm-interpret Stan et al. (2024) as a visualization tool to provide a more transparent understanding of the model’s decision-making process. This method adapts the calculation of relevancy maps to LVLm, thus providing a detailed representation of the regions of the image most relevant to each generated token. Specifically, as shown in Fig. 9 (e) for cross-modal consistency distortion, the relevancy map focuses on the clapping action of the characters in the image. Clapping typically signifies recognition or approval of an event, conveying semantics that differ from the meaning of the associated text. This observation aligns with the model’s detection output, providing intuitive explanations for the identified inconsistency. In contrast, Fig. 9 (g) shows that the model predominantly focuses on the person depicted in the image while neglecting the slogan displayed at its center. This oversight prevents the model from identifying the semantic inconsistencies between the textual and visual modalities, resulting in a failure to capture the underlying mismatch.

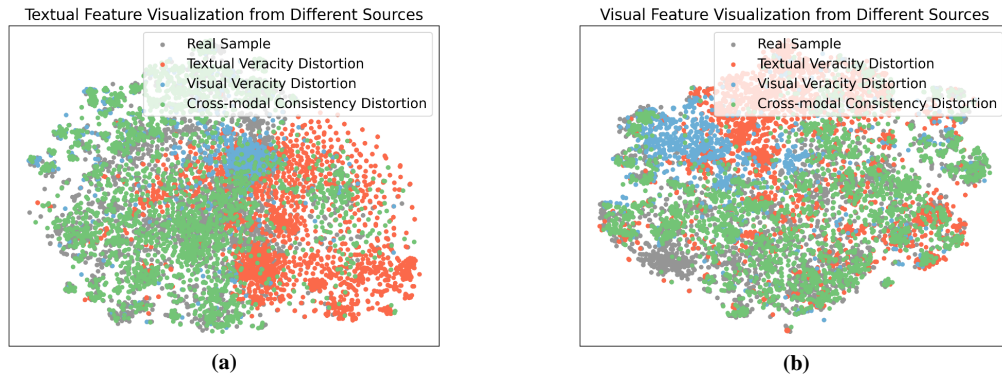


Figure 10: (a) TSNE visualization of textual features from different sources. (b) TSNE visualization of visual features from different sources.

## A.2 BENCHMARK ANALYSIS

In Table 10, We present benchmark analysis based on count statistics (i.e., the average number of words) and diversity (i.e., feature distribution and word frequency) and semantic similarity metrics for different forgery sources. The results indicate: 1) The average number of words is over 10 words, which suggests that the samples are sufficiently informative. 2) We have provided an analysis of the dataset’s diversity through t-SNE visualizations of textual and visual features presented in Fig. 10, as well as

Table 10: Benchmark statistics analysis.

Metric	Real Unit	Fake Unit
Avg. Number of Words	15.8	12.6
Word Frequency Entropy	10.9	11.6
Semantic Similarity (Text-Text)	0.41	0.47
Semantic Similarity (Image-Image)	0.39	0.41

1080 through the entropy of word frequencies detailed in Table 10. This feature visualization provides a  
1081 clear depiction of the dataset’s diversity where each forgery sources are distinctly dispersed across  
1082 the feature space. Moreover, we calculate the overall entropy of word frequency, where the entropy  
1083 for real units ranges from 0 to 15.5 (calculated by  $\log_2(3000 \times 15.8)$ ), 0 to 16.4 (calculated by  
1084  $\log_2(7000 \times 12.6)$ ) for fake units where this dataset’s entropy reached 10.9 and 11.6, demonstrating  
1085 significant diversity. (3) Additionally, we computed the pairwise semantic similarity between samples,  
1086 with results showing an average similarity below 0.5, further confirming the dataset’s rich diversity.

### 1087 A.3 IMPLEMENTATION DETAILS

#### 1088 A.3.1 LVLMS CONFIGURATION DETAILS

1089 **Model Version.** As for ChatGPT model, we use GPT-3.5 (gpt-3.5-turbo) or GPT-4 (gpt-4-vision-  
1090 preview) as generators or detectors. As for text-to-image models, we use DALLE (DALLE-E3),  
1091 Stable-Diffusion (Stable Diffusion XL), and Midjourney (Midjourney V6).  
1092

1093 **Inference Hyperparameters.** To achieve the justified evaluation, we have set the sampling hy-  
1094 perparameter of the off-the-shelf LVLMS, “do\_sample = False” or “Temperature = 0”, to guarantee  
1095 consistency in the predicted outputs. We adopt the default setting of other hyperparameters such as  
1096 “max\_new\_tokens = 512”.  
1097

#### 1098 A.3.2 EXISTING SINGLE-SOURCE DETECTORS DETAILS

1099 **FakingFakeNews.** The FakingFakeNews Huang et al. (2023) is designed for the detection of  
1100 textual fake news, particularly for those natural human-written misinformation. It proposes an  
1101 innovative approach for generating training instances, leveraging established styles and strategies  
1102 commonly employed in human-authored propaganda. FakingFakeNews employs the ROBERTA Liu  
1103 et al. (2019) model as the backbone and trains it on its own proposed PROPANEWS dataset. In our  
1104 experiments, we utilized the default configuration of the ROBERTA detector provided within the  
1105 FakingFakeNews framework, retaining its default hyperparameters.  
1106

1107 **CNNSpot.** CNNSpot Wang et al. (2020) is an artificial image detector designed specifically for  
1108 identifying images produced by generative models. It employs the ResNet-50 model as the classifier  
1109 backbone. Notably, CNNSpot recognizes that data augmentation, including JPEG compression and  
1110 Gaussian blur, can enhance the generalization capabilities of the detector. In our study, we utilize the  
1111 pre-trained CNNSpot model with default hyperparameters to perform the detection of visual veracity  
1112 distortion.  
1113

1114 **UnivFD.** UnivFD Ojha et al. (2023a) is a general-purpose fake image detector that uses a feature  
1115 space not explicitly trained to distinguish between real from fake images. When given access to  
1116 the feature space of a pre-trained vision-language model, UnivFD employs the nearest neighbor to  
1117 identify fake images originating from various sources. The utilization of the large pre-trained model  
1118 results in a smooth decision boundary, thereby enhancing the generalization capability of the detector.  
1119 In our work, we use the pre-trained detector of UnivFD with default hyperparameters to conduct the  
1120 visual veracity distortion detection task.  
1121

1122 **LNP.** LNP Liu et al. (2022) utilizes a well-trained denoising model to extract noise patterns from  
1123 spatial images. Subsequently, it discerns fake images by analyzing the frequency domain of these  
1124 noise patterns. Additionally, LNP employs the ResNet-50 model as the classifier backbone. In our  
1125 study, we utilize the pre-trained LNP detector with default hyperparameters to conduct the visual  
1126 veracity distortion detection task.  
1127

1128 **HAMMER.** HAMMER Shao et al. (2023) is a multimodal detector designed to identify multimedia  
1129 manipulation. It is built upon the pre-trained vision-language model, ALBEF Li et al. (2021) which  
1130 comprises two unimodal encoders and a multimodal Aggregator. To accomplish the multimodal  
1131 manipulation detection task, HAMMER employs hierarchical manipulation reasoning consisting  
1132 of shallow and deep manipulation reasoning. Shallow manipulation reasoning involves semantic  
1133 alignment between image and text embeddings, while deep manipulation reasoning performs deep

cross-modal fusion for forgery detection. In our paper, we employ the off-the-shelf HAMMER detector with default hyperparameters to detect cross-modal consistency distortion.

**FakeNewsGPT4.** FakeNewsGPT4 Liu et al. (2024c) is developed based on large vision language models (LVLMs) to detect multimodal fake news. It identifies two types of forgery-specific knowledge: semantic correlation and artifact tract, and augments LVLMs with these two knowledge. Specifically, it extracts semantic correlations using a multi-level cross-modal reasoning module and comprehends unimodal localized details through a dual-branch fine-grained verification model. In our study, we employed the off-the-shelf FakeNewsGPT4 with default hyperparameters to accomplish the cross-modal consistency distortion detection task.

**Mixed Detection** We combine the three most powerful models on each single-source detection task (i.e., FakingFakeNews for textual veracity distortion, LNP for visual veracity distortion, and HAMMER for cross-modal consistency distortion) to perform mixed detection. Specifically, utilizing our proposed hierarchical decomposition framework, we sequentially assess textual veracity, visual veracity, and cross-modal consistency with three single-source detectors and assign corresponding multi-class labels.

#### A.4 INSTRUCT PROMPT FOR CHATGPT

The construction of MMFakeBench employs the advanced ChatGPT to assist us in generating textual rumors, expanding detailed descriptions, and generating fact-conflicting descriptions. The specific prompts provided in this work are summarized as follows.

**Instruct Prompts to Ask ChatGPT to Generate Textual Rumors.** Fig. 11 illustrates the prompt utilized for asking ChatGPT to generate textual rumors with different prompt methods. These methods include arbitrary generation, rewriting generation, and information manipulation.

##### (1) Arbitrary Generation Prompt

Please write a piece of misinformation title. The domain should be one of gossip, science, health and politics. The time period should be within the past ten years. The type should be fake news/rumors/misleading claims. Avoid answering words like fake, rumor, confusion, disbelief, misinformation, etc.

#query

Misinformation title is:

##### (2) Rewriting Generation Prompt

Given a sentence, please write a piece of misinformation title. The content should be the same. The writing style should be serious, informative and convincing. Avoid answering words like fake, rumor, confusion, disbelief, misinformation, etc.

#query

Sentence:\_\_\_\_\_.

Misinformation title is:

##### (3) Information Manipulation Prompt

Given a true claim, please write a piece of misinformation. It should be long enough, convincing and detailed. The error type should be fake news/rumors/misleading claims. Avoid answering words like fake, rumor, confusion, disbelief, misinformation, etc.

#query

The true claim is:\_\_\_\_\_.

The answer is:

Figure 11: Prompt used to ask ChatGPT to generate textual rumors with arbitrary generation, rewriting generation, and information manipulation methods.

**Instruct Prompts to ask ChatGPT to Expand Detailed Descriptions.** Fig. 12 illustrates the prompt utilized for asking ChatGPT to expand detailed descriptions for textual rumors. Based on

1188 Given a news rumor that is not consistent with reality, I need you to imagine the scene in which the  
 1189 rumor occurred and provide details about the characters, famous buildings, physical objects, etc. in  
 1190 the scene. Avoid sentences that state opinions and only describe physical objects. Avoid words like  
 1191 false, rumor, confusion, disbelief, misinformation, etc.

1192 # in-context examples

1193 Rumor: Peking University is in Thailand.  
 1194 The answer is: A Thai university in the real world, with the Thai flag flying above, including the  
 1195 landmark building of Peking University, the gate of Peking University, a plaque with the name of  
 1196 Peking University, the Boya Tower of Peking University, and Weiming Lake of Peking University.

1197 Rumor: Carlos Santana is a US president.  
 1198 The answer is: This a realistic photo of Carlos Santana in the white house in the real world. Carlos  
 1199 Santana stands at a podium adorned with the presidential seal. Behind him, an American flag hangs  
 1200 proudly. A row of microphones and a table are set before him, as cameras flash around the room.

1201 Rumor: The Chrysler Building has yet to be surpassed in height.  
 1202 The answer is: This is a realistic photo of the center of New York City in the real world. The Chrysler  
 1203 Building stands tall. The slope of the Chrysler Building forms a crown. There are many high-rise  
 1204 buildings on the side.

1205 #query

1206 Rumor: \_\_\_\_\_.

1207 The answer is:

1209 Figure 12: Prompt used to ask ChatGPT to expand detailed descriptions of rumors.

1211 Imagine you as a science fiction writer. Given a true claim, please write a piece of misinformation.  
 1212 The error type should be contrary to objective facts, that is, objects or events that do not exist in the  
 1213 real scene, etc. Avoid answering words like fake, rumor, confusion, disbelief, misinformation, etc.

1214 # in-context examples

1215 The true claim is: People are standing on top of a snowy mountain.  
 1216 The answer is: There are angels with wings welcoming these people.

1217 The true claim is: A person sailing in the air on a snow board.  
 1218 The answer is: A person is snowboarding in the clouds.

1219 The true claim is: A man wearing a blue tie with the ten commandments on it.  
 1220 The answer is: This man is holding a huge fireball in his hand.

1221 #query

1222 The true claim is: \_\_\_\_\_.

1223 The answer is:

1224 Figure 13: Prompt used to ask ChatGPT to generate fact-conflicting descriptions.

1225 the responses, we prompt the stage-of-art diffusion generators to generate realistic and relevant  
 1226 supporting images.

1227 **Instruct Prompts to ask ChatGPT to Generate Fact-conflicting Descriptions.** Fig. 13 illustrates  
 1228 the prompt utilized for asking ChatGPT to generate fact-conflicting descriptions. Then, we combine  
 1229 these descriptions with original captions as prompts in the Midjourney V6 model Midjourney (2022)  
 1230 to create corresponding images with additional fact-conflicting information.

#### 1231 A.5 DATASET LICENSES

1232 The licenses of the existing datasets used in this work is as follows:

- 1233 • **FakeNewsNet**: free to use by all.
- 1234 • **FEVER**: Apache License 2.0.
- 1235 • **Fakeddit**: free to use by all.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

- **NewsClippings**: free to use by all.
- **DGM4**: S-Lab License 1.0
- **COCO-Counterfactuals**: Attribution 4.0 International (CC BY 4.0).
- **VisualNews**: free to use by all.
- **MSCOCO**: free to use by all.

All datasets provided in this work are licensed under the Attribution Non-Commercial ShareAlike 4.0 International (CC BY-NC-SA 4.0) license. We chose this license because some of the original datasets have this license and we provide our datasets with the same level of access.

### A.6 MORE VISUALIZATION EXAMPLES

We have provided more visualization examples of multimodal misinformation from different sources in Fig. 14, Fig. 15 and Fig. 16.







Textual Rumor & AI-generated Image			
			
<b>Charles de Gaulle</b> lost all <b>elections</b> for President of the Fifth French Republic.	The look of love: <b>Katie Holmes</b> and <b>Jamie Foxx</b> pack on the PDA in rare public appearance at pre-Grammys gala.	<b>Pregnant Kylie Jenner</b> 'Looked Great' at Family's Christmas Eve Party, Source Says.	<b>Jennifer Lawrence</b> and <b>Leonardo DiCaprio</b> Spotted <b>Kissing</b> at Exclusive Hollywood Party.
			
In 2015, among Americans, <b>30% of adults</b> had consumed <b>alcoholic drink</b> in the last year.	<b>Northwestern University</b> is a founding member of the Big Ten Conference starting in <b>1893</b>	Scientist <b>Uncovers</b> Previously Unknown <b>Galactic Cluster</b> Near <b>Milky Way Galaxy</b> .	<b>NASA</b> Confirms <b>Existence of Alien Life</b> Form Found on <b>Mars' Surface</b>
Textual Rumor & Repurposed Image			
			
Democrat <b>Maxine Waters</b> Has Shown Up To Only <b>10% Of Congressional Meetings</b> For 35 years.	<b>Obama</b> Admits Secret Deal With China to <b>Hand Over U.S. Territory</b> to Settle Debt	<b>Trump</b> Votes For <b>Death Penalty</b> For <b>Being Gay</b>	<b>Kim Kardashian</b> and <b>Kanye West's</b> Wedding: All of the Best Photos from <b>Paris and Florence</b> Photos.

Figure 14: Visualization examples of multimodal misinformation from the textual veracity distortion.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

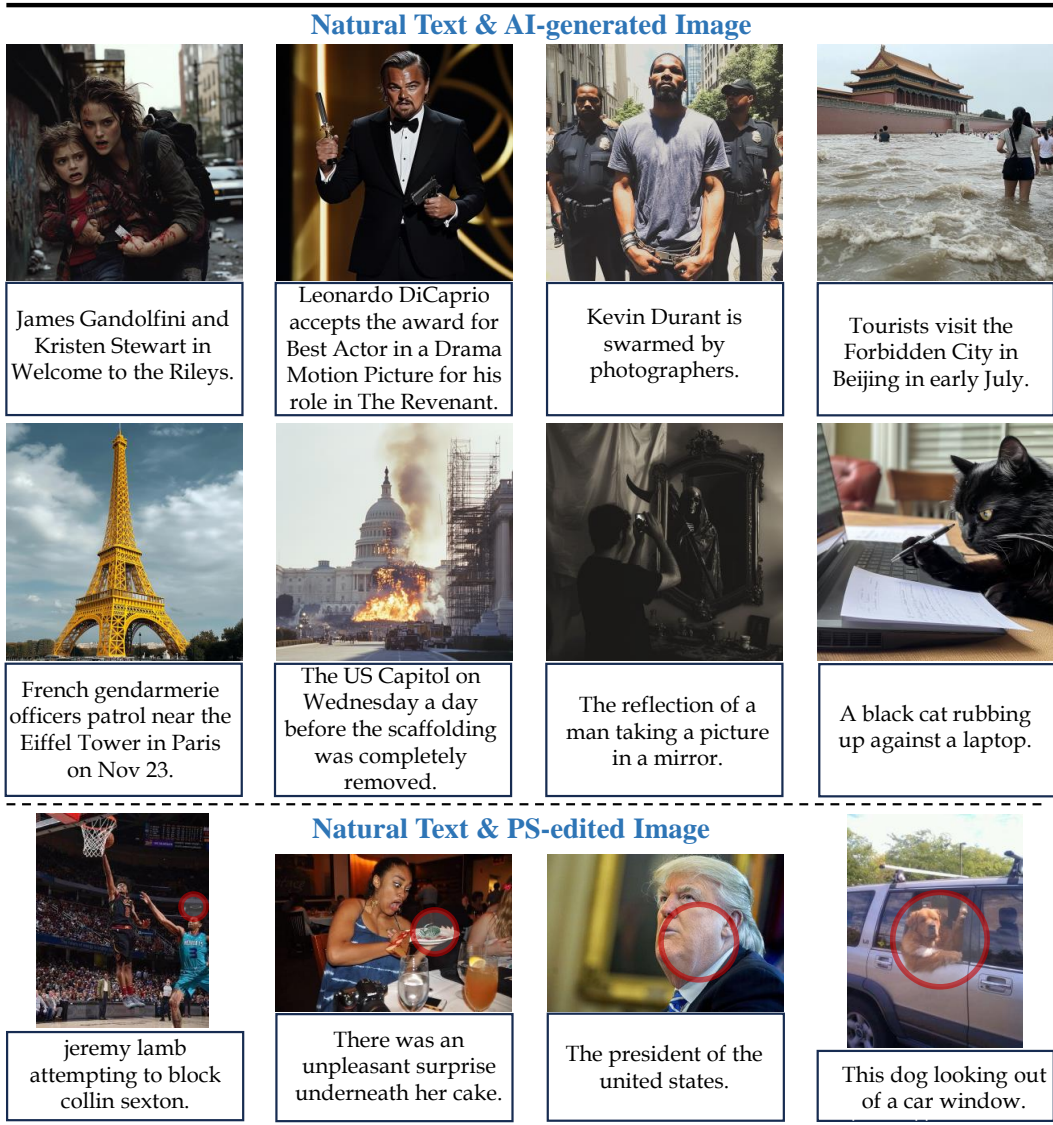


Figure 15: Visualization examples of multimodal misinformation from the visual veracity distortion.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

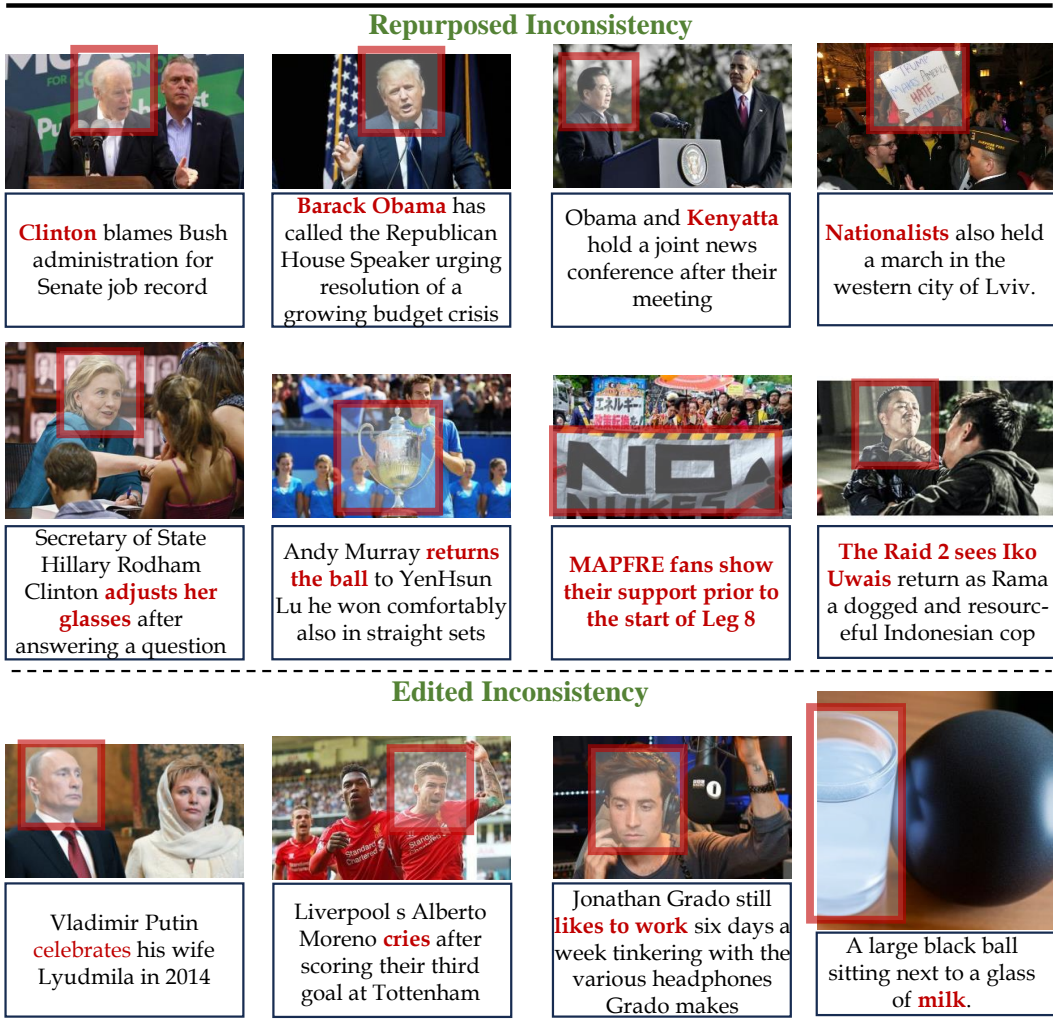


Figure 16: Visualization examples of multimodal misinformation from the cross-modal consistency distortion.