# Multimodal Understanding of Memes with Fair Explanations

Yang Zhong
Department of Computer Science
University of Pittsburgh, PA, USA
yaz118@pitt.edu

Bhiman Kumar Baghel
Department of Computer Science
University of Pittsburgh, PA, USA
bkb45@pitt.edu

## Abstract

*Digital Memes have been widely utilized in people's daily lives over social media platforms. Composed of images and descriptive texts, memes are often distributed with the flair of sarcasm or humor, yet can also spread harmful content or biases from social and cultural factors. Aside from mainstream tasks such as meme generation and classification, generating explanations for memes has become more vital and poses challenges in avoiding propagating already embedded biases. Our work studied whether recent advanced Vision Language Models (VL models) can fairly explain meme contents from different domains/topics, contributing to a unified benchmark for meme explanation. With the dataset, we semi-automatically and manually evaluate the quality of VL model-generated explanations, identifying the major categories of biases in meme explanations.*

## 1. Introduction

The concept of "Memes" was first introduced by Dawkins [6] as the idea/behavior that can spread between people within a culture. In the digital era, there is an explosion of memes on social media platforms such as Twitter, Instagram, and others. Those digital memes are often composed of images and descriptive texts, often distributed with the flair of sarcasm or humor [36]. As described by Tan [35], understanding the humor can be grouped into Proximal mechanisms, which " attempts to provide the mechanism behind the predicted label, i.e., how to infer the label from the text,". Memes can also be connected to the utilization of figurative languages to spread propaganda [7], as well as spreading harmful contents or biases [9, 13, 32]. Therefore, analyzing memes through a sociocultural lens and establishing well-informed regulations is imperative.

Mainstream tasks related to memes cover meme generation [23], meme classification [7, 9, 28], and meme caption/description generation [11, 31, 33]. The third task of explaining memes becomes more important in the current fast-evolving era. Hwang and Shwartz [11] curated 6.3K memes along with the title, meme caption, literal image caption, and annotated visual metaphors, which is a good test bed to study.

However, the generation and elucidation of memes present a multitude of challenges. On the one hand, the availability of training data can be limited: MemeCap [11][1] and MEMEX [33][2] cover a total of 10k image pairs, which are smaller and have domain differences. Other datasets are either annotated with only classification labels [7, 17] or covering semantic role labeling pairs [30] which do not provide natural language explanations of the meme contents. The domains of the memes often fall into the categories of political topics, where the meme composer's political standing plays an important role. On the other hand, memes depend on cultural factors related to both the author and the audience. While composing the explanations, one should know about the audience and the possible harmful/biases embedded within the meme. One viable approach is formulating explanations predicated on specified attributes, thereby mitigating the risk of harmful content being excessively or inaccurately interpreted. Our work is inspired by this idea, first trying to study and diagnose the toxic/biased contents embedded in memes and identify the representative biases in generated explanations. We further propose a list of taxonomies on the biases through manual annotations, finding that biases can have different origins and that more effort is needed to improve the AI models' capability to produce safe content. We foresee the future development of a better VL model on meme explanation because it provides an accurate interpretation of the message, extends a bit on the potential biases, and gives some justifications. Our datasets and curated meme explanations are publicly available at https://github.com/bhimanbaghel/FiME.

---

[1] https://github.com/eujhwang/meme-cap
[2] https://github.com/LCS2-IIITD/MEMEX_Meme_Evidence

| Dataset | Categories | Task | Data source | Caption/ Description | Labels |
|---|---|---|---|---|---|
| MEMEX [33] containing 3400 memes and related context, along with gold-standard human annotated evidence sentence-subset. | political, historical, English language memes | identify explanatory evidence for memes from their related contexts | Meme: Google Images, r/CoronavirusMemes, r/PoliticalHumor, r/PresidentialRace<br><br>Context:Wiki, Quora | Yes | No |
| FigMemes [17] 5141 memes dataset for figurative language classification, covering a wide range of topics and six different figurative language categories | refugees, racial minorities, U.S elections, Epstein, antisemitism, COVID, LGBTQ+, feminism | identify the type of (one or more) figurative language used in a meme. | 4chan /pol/ board<br><br>Similar datasets: HatefulMemes [12], HarMeme [24] | No | Yes |
| MemeCap [11] 6.3K memes along with the title of the post containing the meme, the meme captions, the literal image caption, and the visual metaphors. | text dominant, image dominant, complementary, had no metaphor<br><br>Removed offensive, sexual memes | our extensive experiments using state-of-the-art VL models show that they still struggle with visual metaphors, and perform substantially worse than humans. | r/memes<br><br>Similar datasets MultiMET [48], Met-Meme [42] | Yes | Yes |
| HVVMemes [30] 7K memes containing entities and their associated roles: hero, villain, victim, or other. | COVID-19, US Politics | Hero, Villain, and Victim: Dissecting Harmful Memes for Semantic Role Labeling of Entities. | reannotated the HarMeme [24] dataset | No | Yes |

Table 1. Meme Dataset.

| Dataset | Labels |
|---|---|
| MEMEX [33] | NA |
| FigMemess [17] | Allusion<br>Exaggeration/Hyperbole<br>Irony/Sarcasm<br>Anthropomorphism/Zoomorphism<br>Metaphor/Smile<br>Contrasts |
| MemeCap [11] | text dominant<br>image dominant<br>complementary<br>no methaphore |
| HVVMemes [30] | Hero<br>Villain<br>Victim<br>Other |

Table 2. Label Analysis over the four Meme datasets.

## 2. Related Work

There are three broad tasks related to meme generation, meme classification, and meme caption/description generation. All these tasks have their particular set of datasets characterized distinctively based on the associated task. These work interests are most in line with meme caption/description generation datasets like MEMEX [33] and MemeCap [11]. Meme classification datasets like FigMemess [17][3] and HVVMemes [30][4] can also be utilized. However, to make them in line with this work, they will require support from the cation/description generation model. Table 1 shows the important details about these datasets, which will assist in identifying the correct dataset for this work. The table gives the dataset's general description. It then provides the dataset distribution in the form of categories and also mentions the data sources from which the dataset was curated. It then mentions the task for which the dataset was utilized and finally informs whether it contains captions/descriptions and labels. It can be observed from

---

[3]https://github.com/UKPLab/emnlp2022-FigMemess
[4]https://constraint-lcs2.github.io/

| VL Model | Language Model | Vision Model | Training data |
|---|---|---|---|
| OpenFlamingo-9B [3] | LLaMA 7B [37] | CLIP ViT/L-14 [25] | Multimodal C4 dataset [51] <br> LAION-2B [27] |
| MiniGPT-4 [50] | Vicuna [49] | BLIP-2 [15] | LAION [27], Conceptual Captions [29], <br> SBU [22] |
| LLaVA [18] | LLaMA [37] | CLIP [25] | [18] proposed Instructional vision-language data. |
| PaLI [5] | mT5 [43] | ViT [47] | [5] proposed WebLI,a multilingual image-language dataset |

Table 3. Vision Language models

the category column in Table 1 that MEMEX [33], Fig-Memess [17] and HVVMemes [30] do contain political topics. No such comment was made for MemeCap [11] from the initial analysis. So, an in-depth assessment is required for this dataset. MemeCap is important because it contains captions/descriptions, which is the primary requirement for this work. Further research of class labels present in the datasets (shown in Table 2) can help in filtering down the dataset for more precise targeted memes required for this work. Another noticeable thing is that FigMemess [17] and HVVMemes [30] don't contain a caption or description. However, this dataset can be essential as it contains political opinion data.

Memes, being images, are different from normal images because they contain visual and textual information. Both these pieces of information aid the meme's overall understanding and intent. So, to understand a meme and perform any task upon it, the underlying system should be able to comprehend vision and text modalities. This is where Vision Language models come into the picture. Vision language models (VL) can be utilized to overcome the shortcomings of datasets like FigMemess [17] and HVVMemes [30] and generate caption/description.

Vision language models generally combine two models, each handling one modality. According to [41], these combinations have four major flavors. First is jointly training image and text as a single feature vector [1, 16, 40]. Second is learning only image embedding for a frozen pre-trained language model [21, 38]. Third is employing a special mechanism to fuse visual context into layers of language model [2, 4, 5, 18, 20, 45, 46, 50]. Details of a few of these models are mentioned in Table 3. All previously mentioned categories required some level of training in the models. However, there are techniques [34, 44] that can combine the vision and language models without any training. This marks the fourth category. Although, chances are that these might not perform as well as their trained counterparts.

However, a general concern arises about the quality and accuracy of such generation. A recent study [11] observed that VL models struggle to understand visual metaphors. It would be interesting to study their performance when we add another lens of fairness. This leads to the research questions of this work, as discussed in the next section.



Figure 1. Bias in MEME Explanation

# 3. Problem Statement

We aimed to study these Research Questions:

**1. How can we enrich the dataset by generating pseudo explanations over data with meme images, the caption, and propaganda labels?** One inspiration is to align the annotated figurative labels and translate them by generating sentences from language models given the literal caption texts and the extracted text in the image through an OCR system. We plan to generate the explanatory captions using visual-caption models LLaVA [18], and MiniGPT-4 [50] and add to the model prompt (see Table 5). We will then append the task labels to the model to generate the pseudo-meme explanation. One challenge would be mapping and unifying the classification labels from different datasets into the same domain, thus producing explanations within a shared vocabulary of social factors.

**2. Can we identify the harmful ones and the targets for the bias for different explanations** Given the explanations, we want to apply machine-learning models to predict the harmful/toxicity of the generation or the social biases against a certain group. There are already publically available APIs such as PerspectiveAPI[5], as well as trained models [8, 26]. The goal is to identify and analyze the distribution of social biases in manually written and automatically generated explanations. We acknowledge that the models will make errors and conduct human verification in the middle to evaluate the results.

To answer these research questions, our work can be broken into the following sections: We start by discussing gathering and unifying meme datasets, which could span multiple genres and cover abundant annotated data related to figurative language or the employment of metaphor. We then prompt the VL models to generate explanations, providing the meme and text. Afterward, both automatic evaluation metrics and manual evaluations are applied to the generated text, and we conduct a systematic study to evaluate the biases. Lastly, we identify several ways to mitigate bias and conduct a preliminary study on the automatic ways to mitigate biases.

# 4. Meme Dataset Gathering and Unification

As aforementioned, the literature lacks datasets that evaluate the automatically generated meme explanations from the lens of fairness. To bridge this gap, we collected meme datasets for various tasks mentioned in Table 1. As they are from different sources and purposes, we first unified their input feature space as shown in Table 4. The OCR is generated using EasyOCR[6], and the caption is generated using

|   | Features | MEMEX | MEmeCap | FigMemes | HVVMemes |
|---|---|---|---|---|---|
| I N P U T | Meme image | Y | Y | Y | Y |
| | OCR text inside the Meme | Y | Y (generated) | Y | Y |
| | Title of the Meme | NA | Y | NA | NA |
| | Caption (Image Literal description) | Y (generated) | Y | Y (generated) | Y (generated) |
| | Labels (Metadata) | NA | Y | Y | Y |
| O U T | Explanation | Y (generated) | Y (generated) | Y (generated) | Y (generated) |

Table 4. Data Unification. 'Y' signifies this feature was already a part of dataset. 'Y (generated)' signifies a missing feature that is later generated, and NA signifies a missing feature that is not generated.

| Prompt | Data Point | Prompt |
|---|---|---|
| **raw** | **Image** | What is the meme poster trying to convey? |
| **p2** | **Image + OCR + Caption** | '''This is a meme. The image description is "{image_caption}". The following list of texts is written inside the meme: "{OCR_text}". \n\n What is the meme poster trying to convey?''' |
| **p3** | **Image + OCR + Caption + Metadata** | '''This is a meme. The image description is "{image_caption}". The following list of texts is written inside the meme: "{OCR_text}".{figurative_text}\n\n What is the meme poster trying to convey?''' |
| | **Image + Title** | This is a meme with the title <Title>. What is the meme poster trying to convey? (only applies to memecap) |

Table 5. Prompts for Explanation Generation.

VL models mentioned in Table 3. For datasets that missed image captions, We specifically prompt the MiniGPT-4 and LLaVA models to produce the image captions.

## 4.1. Explanation Generation

After unifying the input, we generated meme explanations using the VL model LLaVA 1.5 [19] and MiniGPT-4 [50][7]. We used three prompt variations inspired from [11] as shown in Table 5 to generate the explanations. This was done to monitor the behavior and change in the generated explanation of the VT model about the change in the input feature.

# 5. Fairness Evaluation

We performed two types of fairness evaluation: Automatic and Manual. The idea is to connect the notion of bias with toxicity/profanity evaluations, which have been long studied in the NLP area.

---

[5] https://perspectiveapi.com/
[6] https://github.com/JaidedAI/EasyOCR

[7] LLaMA-2 Chat 7B

| Dataset (size) | LLaVA / MiniGPT-4 | | |
|---|---|---|---|
| | Raw | P2 | P3 |
| FigMemes (1518) | 28 / 38 | 30 / 51 | 28 / 57 |
| MemeCap (559) | 1 / 8 | 2 / 7 | 3 / 4 |
| MEMEX (200) | 2 | 2 | N/A |
| HVV-Covid (300) | 1 | N/A | 3 |
| HVV-USPolitics (350) | 6 | 9 | 6 |

Table 6. Detection of biased explanation based on PerspectiveAPI-Toxicity for meme explanations, for FigMemes and MemeCap, we also report the MiniGPT-4 results.

## 5.1. Automatic Evaluation

For Automatic evaluation, we picked three models. We used the toxicity [10, 39][8] and the Profanity package [9] to predict the score for each explanation within the range of 0-1. Moreover, we leverage the PerspectiveAPI [14] to evaluate the explanations from six dimensions: *'INSULT', 'THREAT', 'TOXICITY', 'SEVERE_TOXICITY', 'IDENTITY_ATTACK', and 'PROFANITY'*. Following prior work, we treat an explanation scored 0.5 or higher as carrying the bias.

### 5.1.1 The Distribution of Scores

Overall, we observe that most explanations are tagged as unbiased based on metrics. We report the PerspectiveAPI-TOXICITY score in Table 6. For FigMemes and Memecap, we additionally experimented with MiniGPT-4 and found that more data points are tagged as biased. When evaluating the number of biased samples between the three types of prompts, we find that injecting text/OCR captions slightly enlarges the amount of biased data. We will conduct some analysis in later sections.

**Rejection to Response**  One additional model behavior was found for MiniGPT-4's results is the rejection behavior, that is, when the meme contains some offensive language or some internal problem in loading the image, the VL model outputs will say "I apologize, but I cannot that may be harmful or" or " I cannot access or show images". We thus apply the new category of "failures" into the aforementioned evaluation categories. We additionally analyze the overall biased data being labeled biased by at least one of the nine metrics, as denoted in Table 7. The numbers increased for both models.

---

[8] https : / / huggingface . co / spaces / evaluate - measurement/toxicity

[9] https : / / github . com / dimitrismistriotis / alt - profanity-check

| Dataset | LLaVA | | | MiniGPT-4 | | |
|---|---|---|---|---|---|---|
| | Raw | P2 | P3 | Raw | P2 | P3 |
| FigMemes | 74 | 78 | 74 | 66 | 85 | 79 |
| MemeCap | 9 | 9 | 8 | 9 | 13 | 12 |

Table 7. Biased explanation labeled by at least one model for FigMemes and MemeCap.

### 5.1.2 Agreement Between Metrics

While we have multiple metrics to detect the bias, it is interesting to understand how the different metrics agree with each other. We thus measure the pair-wise correlations between different models by computing the Spearman's rank correlation (range between -1 and 1) between two score lists. Figure 2 shows the agreement of scores across both models on FigMemes. We found that PerspectiveAPI scores have high agreements, while the two off-the-shelf models on toxicity and profanity have a lower agreement with each other. This unveils the limitations of automatic metrics, and we move on to the second section for manual evaluations.

### 5.1.3 Manual Evaluation

We sampled a small portion of distinct memes from the test dataset to perform the manual evaluation. We are working on 4 different datasets (see Table 1) to capture variety of memes. Additionally, a test portion of the datasets was chosen to capture a variety of memes within the dataset. For MEMEX and HVVMemes, samples were drawn in sequence, whereas for MemeCap and FigMemes, samples were selected from the ones marked as biased according to the automatic evaluation.

Bias evaluation on memes requires adequate familiarisation with the meme and language comprehension. Therefore, evaluation is performed by graduate-level students with high English proficiency. We also made sure the evaluators had some prior knowledge of bias evaluation. Since we are working on different types of memes with a high bar of technical background, it poses a challenge for the evaluators not to be familiar with all types of memes. In such scenarios, getting a high inter-annotator agreement is also difficult. To address this challenge, evaluators only evaluated the memes with which they were most familiar. We have classified whether an explanation is biased and identified the source of biases, as shown in Table 8. This level of detailing in evaluations makes high - meme familiarity, background knowledge, and proficiency critical aspects of the evaluator.

**Results**  During our evaluation, we tried to categorize the type of bias in the explanations. An example of some of
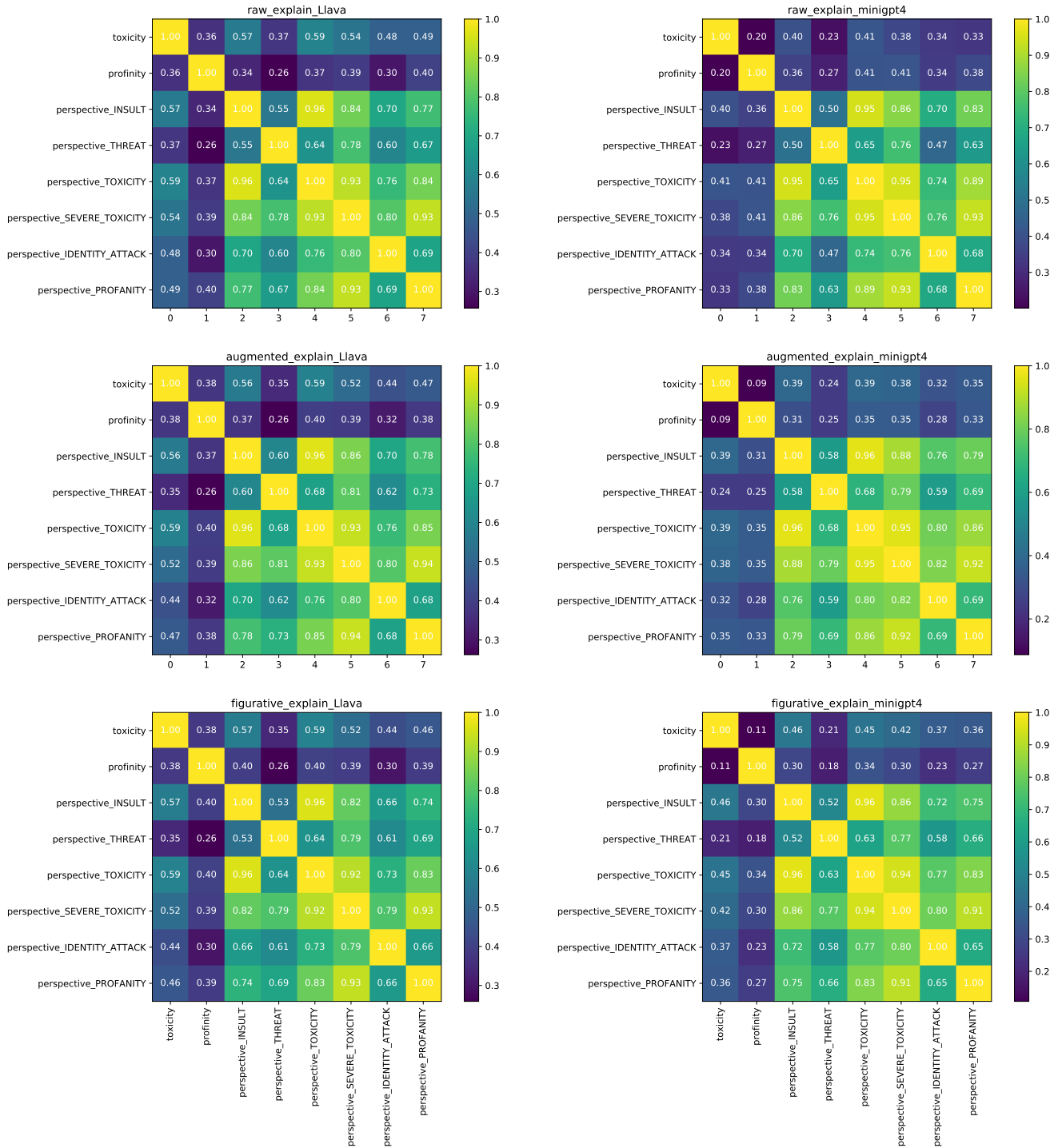
spearman_r



Figure 2. Correlations on metrics for FigMemes explanation.

| Dataset | Bias | | | Bias Towards some Common Visual Feature | Bias Towards some entity or group or gender | Bias Towards image than text | Bias Towards Usage of words in particular sentiment | No Bias but explanation wrong |
|---|---|---|---|---|---|---|---|---|
| | raw | p2 | p3 | | | | | |
| MEMEX (20) | 10 | 9 | NA | 5 | 10 | 6 | 7 | 6 |
| HVV-Covid (10) | 2 | 1 | 1 | 1 | 3 | 0 | 0 | 12 |
| HVV-USPolitics (10) | 3 | 1 | 1 | 0 | 1 | 0 | 4 | 12 |
| Figmeme (minigpt4) (10) | 6 | 1 | 2 | - | - | - | - | - |
| Memecap (10) | 3 | 7 | 2 | - | - | - | - | - |

Table 8. Manual Evaluation Results, for each dataset, we select 10/20 memes and annotate all three variations unless specifically noted.

the categories is shown in Fig 1. Specifically, it shows two types of bias. One 'Bias towards a common visual feature' is where people running in a group are identified as participating in some race. However, in reality, they are running in fear of the bomb. Another type of bias is 'Bias towards using the word in the particular sentiment.' Here, the 'democracy' word is seen in positive sentiment even if there is a bomb-tagged democracy that is about to kill people. We also identified some more categories mentioned in Table 8 along with their distribution across the sample space. We have also mentioned the count of explanations that showed no bias but were not coherent with the meme.



Figure 3. Proper OCR mitigates bias



Figure 4. Wrong Caption introduces bias

## 5.2. Bias Mitigation

Once we have identified the bias, mitigating it is also essential. Here are some of the findings we draw out from our evaluation, which helps mitigate bias:



Figure 5. An example of a good explanation with the help of context information.

1. Proper OCR mitigates bias: We found out that bias towards image over text was removed when text in the meme is provided in the prompt. An example is shown in Fig. 3, where with prompt with OCR information, VL model (LLaVA) generated explanation getting biased towards the image and text inside it. However, when OCR information is added to the prompt, the same model generates the correct explanation without any bias and is coherent with the meme.

2. Wrong Caption introduces bias: We found out that captions generated from the VL model can carry its bias and, when given a prompt to generate an explanation, influence the explanation to be biased. An example is shown in Fig. 4, where the VL model (LLaVA) without any caption in the prompt generated the correct explanation for the meme. However, when the caption was generated from the same model, it produced a caption biased towards a common visual feature. When this caption was provided in a prompt for explanation generation, influenced by the caption, the same model generated a biased explanation and amplified the bias by introducing new ones, i.e., bias towards some entity or group or gender.

However, scaling up the mitigation with the large models poses some challenges. Firstly, we notice that closed-source models such as GPT4 would have similar rejection behavior on offensive memes. This could be attributed to the alignment done in the model training and fine-tuning stages. This remains an unsure option whether the forced rejection

would be helpful to improve the meme explanation; instead, a good explanation with some reference to the biased source could be more helpful, as denoted by one good example in Figure 5. Compared to the raw explanation, which explains the texts but does not provide any justification, the latter part of the P2 explanation tried to provide a fairer view of the understanding and avoid propagating the biases introduced by the meme creator.

## 6. Conclusion and Future Work

In this work, we propose a generative task to produce an explanation of memes. We found that current VL models, such as LLaVA and MiniGPT-4, can have biases in generating meme explanations.

We contribute a unified dataset across four separate corpora and produce a diverse set of prompts for benchmark evaluation. We find that biases can have different origins through automatic and manual evaluations, and more effort is needed to improve the AI models' capability to produce safe content. We foresee the future development of a better VL model on meme explanation because it provides an accurate interpretation of the message, extends a bit on the potential biases, and gives some justifications. For future work, we plan to utilize large models to generate less harmful explanations with the original peers and fine-tune the large models on the neutralized data.

## Limitation and Social Impact

We acknowledge that the study in bias is complicated, and our analysis might be limited and focused only on the vehicles of the figurative languages used in memes. Moreover, we could not perform instruction tuning on the large models due to computing resource restrictions, and our goal was to test the off-the-shelf reasoning capability of those models. The meme explanation task involves employing background knowledge, which may vary between annotators. Meanwhile, more carefully selecting and instructing the annotators is crucial to alleviate misunderstandings or misrepresentations of different cultures or social groups. To further mitigate this limitation, a voluntary based[10] annotation strategy with importance on meme familiarity and adequate background can be a potential future direction. In addition, there is some level of subjectivity concerning the evaluation criteria for the meme explanation quality, as denoted by the inconsistency between automatic metrics such as PerspectiveAPI scores and manual judgments. Our study focused on benchmarking two open-sourced LVMs; while more powerful VLMs are being used, they lack sufficient benchmarking on their performances on bias-related tasks.

On the other hand, memes keep evolving and become obsolete quickly as online social trends change quickly. While

our dataset collected memes spanning different periods, we admit that a more comprehensive benchmark should be frequently updated. Given the increasing use of Large Language / Vision Language models in understanding and generating culturally and contextually nuanced content, it is crucial to study the potential biases of those models carefully. However, the contents of currently available meme datasets may be limited to their specific domains and the designing goals of the original dataset. While we propose a first step towards unifying different sources, covering political, historical, and more recent pandemic-related memes, we find that the major sources of memes may still be biased toward the Western world in English. We advocate for a multi-lingual, multi-domain study on the memes study. It is also important to protect the private information of real people from the publically available benchmarks. The use of publicly available memes does not automatically negate potential privacy violations or the ethical implications of analyzing potentially sensitive content.

## References

[1] Armen Aghajanyan, Bernie Huang, Candace Ross, Vladimir Karpukhin, Hu Xu, Naman Goyal, Dmytro Okhonko, Mandar Joshi, Gargi Ghosh, Mike Lewis, and Luke Zettlemoyer. Cm3: A causal masked multimodal model of the internet, 2022. 3

[2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. 3

[3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 3

[4] Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. Visualgpt: Data-efficient adaptation of pretrained language models for image captioning, 2022. 3

[5] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov,

---

[10] https://www.labinthewild.org/

Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali: A jointly-scaled multilingual language-image model, 2023. 3

[6] Richard Dawkins. *The Selfish Gene*. Oxford University Press, London, England, 2006. 1

[7] Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. SemEval-2021 task 6: Detection of persuasion techniques in texts and images. In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 70–98, Online, 2021. Association for Computational Linguistics. 1

[8] Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. Latent hatred: A benchmark for understanding implicit hate speech. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 345–363, Online and Punta Cana, Dominican Republic, 2021. Association for Computational Linguistics. 4

[9] Elisabetta Fersini, Francesca Gasparini, Giulia Rizzi, Aurora Saibene, Berta Chulvi, Paolo Rosso, Alyssa Lees, and Jeffrey Sorensen. SemEval-2022 task 5: Multimedia automatic misogyny identification. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 533–549, Seattle, United States, 2022. Association for Computational Linguistics. 1

[10] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*, 2020. 5

[11] EunJeong Hwang and Vered Shwartz. Memecap: A dataset for captioning and interpreting memes, 2023. 1, 2, 3, 4

[12] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. In *Advances in Neural Information Processing Systems*, pages 2611–2624. Curran Associates, Inc., 2020. 2

[13] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. *Advances in neural information processing systems*, 33:2611–2624, 2020. 1

[14] Alyssa Lees, Vinh Q Tran, Yi Tay, Jeffrey Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. A new generation of perspective api: Efficient multilingual character-level transformers. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3197–3207, 2022. 5

[15] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023. 3

[16] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language, 2019. 3

[17] Chen Liu, Gregor Geigle, Robin Krebs, and Iryna Gurevych. FigMemes: A dataset for figurative language identification in politically-opinionated memes. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7069–7086, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. 1, 2, 3

[18] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 3, 4

[19] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 4

[20] Ziyang Luo, Yadong Xi, Rongsheng Zhang, and Jing Ma. A frustratingly simple approach for end-to-end image captioning, 2022. 3

[21] Ron Mokady, Amir Hertz, and Amit H. Bermano. Clipcap: Clip prefix for image captioning, 2021. 3

[22] Vicente Ordonez, Girish Kulkarni, and Tamara L Berg. Im2text: Describing images using 1 million captioned photographs. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, page 1143–1151, Red Hook, NY, USA, 2011. Curran Associates Inc. 3

[23] Abel L Peirson V and E Meltem Tolunay. Dank learning: Generating memes using deep neural networks. *arXiv preprint arXiv:1806.04510*, 2018. 1

[24] Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. Detecting harmful memes and their targets. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2783–2796, Online, 2021. Association for Computational Linguistics. 2

[25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 3

[26] Nihar Sahoo, Himanshu Gupta, and Pushpak Bhattacharyya. Detecting unintended social bias in toxic language datasets. In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 132–143, Abu Dhabi, United Arab Emirates (Hybrid), 2022. Association for Computational Linguistics. 4

[27] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. 3

[28] Chhavi Sharma, Deepesh Bhageria, William Scott, Srinivas Pykl, Amitava Das, Tanmoy Chakraborty, Viswanath Pulabaigari, and Bjorn Gamback. Semeval-2020 task 8: Memotion analysis–the visuo-lingual metaphor! *arXiv preprint arXiv:2008.03781*, 2020. 1

[29] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia, 2018. Association for Computational Linguistics. 3

[30] Shivam Sharma, Tharun Suresh, Atharva Kulkarni, Himanshi Mathur, Preslav Nakov, Md. Shad Akhtar, and Tanmoy Chakraborty. Findings of the CONSTRAINT 2022 shared task on detecting the hero, the villain, and the victim in memes. In *Proceedings of the Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situations*, pages 1–11, Dublin, Ireland, 2022. Association for Computational Linguistics. 1, 2, 3

[31] Shivam Sharma, Siddhant Agarwal, Tharun Suresh, Preslav Nakov, Md Shad Akhtar, and Tanmoy Chakraborty. What do you meme? generating explanations for visual semantic role labelling in memes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9763–9771, 2023. 1

[32] Shivam Sharma, Atharva Kulkarni, Tharun Suresh, Himanshi Mathur, Preslav Nakov, Md. Shad Akhtar, and Tanmoy Chakraborty. Characterizing the entities in harmful memes: Who is the hero, the villain, the victim? In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 2141–2155. Association for Computational Linguistics, 2023. 1

[33] Shivam Sharma, Ramaneswaran S, Udit Arora, Md. Shad Akhtar, and Tanmoy Chakraborty. MEMEX: Detecting explanatory evidence for memes via knowledge-enriched contextualization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5272–5290, Toronto, Canada, 2023. Association for Computational Linguistics. 1, 2, 3

[34] Yixuan Su, Tian Lan, Yahui Liu, Fangyu Liu, Dani Yogatama, Yan Wang, Lingpeng Kong, and Nigel Collier. Language models can see: Plugging visual controls in text generation, 2022. 3

[35] Chenhao Tan. On the diversity and limits of human explanations. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2173–2188, Seattle, United States, 2022. Association for Computational Linguistics. 1

[36] Kohtaro Tanaka, Hiroaki Yamane, Yusuke Mori, Yusuke Mukuta, and Tatsuya Harada. Learning to evaluate humor in memes based on the incongruity theory. In *Proceedings of the Second Workshop on When Creative AI Meets Conversational AI*, pages 81–93, Gyeongju, Republic of Korea, 2022. Association for Computational Linguistics. 1

[37] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. 3

[38] Maria Tsimpoukelli, Jacob Menick, Serkan Cabi, S. M. Ali Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models, 2021. 3

[39] Bertie Vidgen, Dong Nguyen, Helen Margetts, Patricia Rossini, and Rebekah Tromble. Introducing CAD: the contextual abuse dataset. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2289–2303, Online, 2021. Association for Computational Linguistics. 5

[40] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision, 2022. 3

[41] Lilian Weng. Generalized visual language models. *Lil'Log*, 2022. 3

[42] Bo Xu, Ting Li, Junzhe Zheng, Mehdi Naseriparsa, Zhehuan Zhao, Hongfei Lin, and Feng Xia. Met-meme: A multimodal meme dataset rich in metaphors. *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2022. 2

[43] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online, 2021. Association for Computational Linguistics. 3

[44] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. An empirical study of gpt-3 for few-shot knowledge-based vqa, 2022. 3

[45] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models, 2022. 3

[46] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models, 2021. 3

[47] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers, 2022. 3

[48] Dongyu Zhang, Minghao Zhang, Heting Zhang, Liang Yang, and Hongfei Lin. MultiMET: A multimodal dataset for metaphor understanding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3214–3225, Online, 2021. Association for Computational Linguistics. 2

[49] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. 3

[50] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. 2023. 3, 4

[51] Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig

Schmidt, William Yang Wang, and Yejin Choi. Multimodal C4: An open, billion-scale corpus of images interleaved with text. *arXiv preprint arXiv:2304.06939*, 2023. 3