

MemeIntel: Explainable Detection of Propagandistic and Hateful Memes

Anonymous ACL submission

Abstract

The proliferation of multimodal content on social media presents significant challenges in understanding and moderating complex, context-dependent issues such as misinformation, hate speech, and propaganda. While efforts have been made to develop resources and propose new methods for automatic detection, limited attention has been given to label detection and the generation of explanation-based rationales for predicted labels. To address this challenge, we introduce *MemeXplain*, an explanation-enhanced dataset for **propaganda memes in Arabic and hateful memes in English**, making it the *first* large-scale resource for these tasks. To solve these tasks, we propose a **multi-stage optimization approach** and train **Vision-Language Models (VLMs)**. Our results demonstrate that this approach significantly improves performance over the base model for both **label detection** and **explanation generation**, outperforming the current state-of-the-art with an **absolute improvement of $\sim 3\%$ on ArMeme and $\sim 7\%$ on Hateful Memes**. For reproducibility and future research, we aim to make the *MemeXplain* dataset and experimental resources publicly available.¹

1 Introduction

Despite the rapid growth of multimodal content—integrating images, text, and sometimes video—the automated detection of harmful and false information on online news and social media platforms has become increasingly critical. In particular, identifying propaganda and hate in memes is essential for combating misinformation and minimizing online harm. While most research has focused on textual analysis, multimodal approaches have received comparatively less attention. In propaganda detection, text-based methods have evolved from monolingual to multilin-

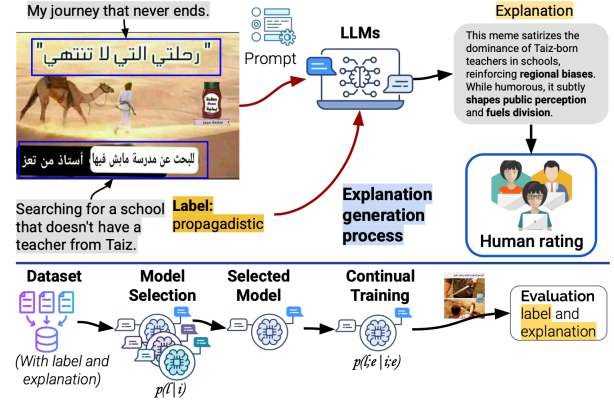


Figure 1: Experimental steps for explanation generation and training.

gual setups (Piskorski et al., 2023; Hasanain et al., 2023), initially through binary classification and later via multilabel and fine-grained span-level tasks (Barrón-Cedeno et al., 2019; Habernal et al., 2017, 2018; Da San Martino et al., 2019). Hate speech detection has similarly progressed from text-based to multimodal approaches that integrate both textual and visual elements. Recent methods have shifted from transformer-based text detection (Fortuna and Nunes, 2018) toward techniques that incorporate visual context (Kielbaso et al., 2020) by leveraging fusion strategies, attention mechanisms, and contrastive learning to boost accuracy, especially when hateful intent is conveyed through text-image interplay (Alam et al., 2022).

The emergence of LLMs has demonstrated significant capabilities across various disciplines. Consequently, efforts have been made to leverage Vision-Language Models (VLMs) (Zhang et al., 2024) and prompting techniques to enhance the detection and classification of harmful and propagandistic memes (Cao et al., 2023). LLM-based models utilize prompt-based learning (Cao et al., 2022), contrastive learning techniques such as CLIP (Kumar and Nandakumar, 2022), and cross-modal attention mechanisms to better capture implicit hate and propaganda.

¹anonymous.com

Despite significant progress, challenges remain in detecting implicit hatefulness, particularly when sarcasm or an ironic dissonance exists between text and images. Propagandistic memes further complicate detection by employing emotional appeals, humor, cultural references, manipulative language, and other rhetorical strategies. To address these nuances, it is crucial for a system to provide not only accurate predictions but also interpretable explanations that reveal the underlying reasoning behind its decisions (Hee et al., 2023; Yang et al., 2023; Huang et al., 2023; Sun et al., 2023). Such explanations enhance classifier reliability by helping end users understand the decision, provided they maintain a natural tone and relate closely to the visual elements of the meme.

An explanation-based approach offers numerous advantages and enhances performance across various tasks (Li et al., 2022; Magister et al., 2022; Nandi et al., 2024; Kumari et al., 2024). While most studies have focused on textual content (Li et al., 2022; Magister et al., 2022), a few recent approaches (Nandi et al., 2024; Kumari et al., 2024) have applied explainability to images. However, these methods rely on QA-based explanations that lack naturalness, use multiple inference calls with custom models—thereby increasing computational complexity—and employ explanations only during training rather than as an inference output. These limitations motivated us to explore a simplified procedure for meme classification and explanation generation. To overcome the above limitations, we propose a novel procedure that achieves state-of-the-art performance on the target tasks across two distinct datasets. Our contributions are briefly outlined below:

- We developed explanation-enhanced datasets, *MemeXplain*, using a rapid and low-cost annotation procedure;
- We investigated state-of-the-art VLMs to identify an appropriate model for meme classification and explanation generation;
- We proposed an efficient *multi-stage* optimization procedure that significantly improves performance;
- With the experiments we achieved state-of-the-art performance on two types of datasets related to propaganda and hateful content detection.

Our findings are as follows: (a) A higher human evaluation score suggests that explanations from

stronger models (e.g., GPT-4o) are reliable and can serve as gold-standard explanations for training smaller models. (b) Task-specific fine-tuning improves performance over the base model. (c) Our multi-stage optimization approach benefits both label detection and explanation generation. Overall, our work is the *first* to enhance VLMs for simultaneous propaganda and hateful content detection while providing natural reasoning to end users.

2 Related Work

The widespread use of social networks has become a major channel for spreading misinformation, propaganda, and harmful content. Significant research efforts have been directed toward addressing these challenges, particularly in multimodal disinformation detection (Alam et al., 2022), harmful memes (Sharma et al., 2022), and propagandistic content (Dimitrov et al., 2021a). However, most studies have focused on detection, while less attention has been given to generating natural explanations/reasons behind the predicted labels.

2.1 Multimodal Propagandistic Content

Following the previous research for propaganda detection using textual content (Da San Martino et al., 2019), Dimitrov et al. (2021b) introduced SemEval-2021 Task 6 focusing on persuasion techniques detection in both textual and visual memes. Subsequently, the focus has extended to the detection of multilingual and multimodal propagandistic memes (Dimitrov et al., 2024). Glenski et al. (2019) studied multimodal disinformation content on social media platforms in multilingual settings. Similar multimodal work on Arabic involves the development of datasets and shared task for propaganda detection (Alam et al., 2024b; Hasanain et al., 2024). For the detection problem, typical approaches include a fusion of textual and visual embedding and a classification head on top of them (Hasanain et al., 2024; Shah et al., 2024), graph attention network based approach for multimodal visual-textual objects Chen et al. (2024).

2.2 Multimodal Hate speech

Similarly, there has been growing interest in detecting multimodal hate speech (Kiela et al., 2020; Veliloglu and Rose, 2020; Hee et al., 2022). Due to the lack of resources, Kiela et al. (2020) developed a large-scale dataset for multimodal hate identification. This study advanced research in this area

and emphasized the importance of integrating textual and visual features for effective detection. The issue has also been explored through a multi-task learning framework for identifying hate speech in memes using multimodal features (Sharma et al., 2022). To further progress in this field, efforts have been made to develop resources for multiple languages, including Arabic (Alam et al., 2024a), Bangla (Hossain et al., 2022), and English (Hee et al., 2023). A more detailed summary of these earlier efforts can be found in Sharma et al. (2022), which also highlights key challenges and outlines future research directions.

2.3 Training with Explanations

Integrating reasoning or explainability capabilities to enhance LLM/VLM performance has been shown to be highly beneficial for various tasks across multiple domains (Plaat et al., 2024). This approach has also proven effective for knowledge distillation and model compression (Li et al., 2022; Magister et al., 2022), where explanations generated by large LLMs improve the performance and capabilities of smaller LLMs. In the context of hateful speech, toxicity detection, and sentiment analysis, it has led to significant advancements (Yang et al., 2023; Huang et al., 2023; Sun et al., 2023). For example, in the hateful speech detection task, Hare (Yang et al., 2023) employs Chain-of-Thought (CoT) reasoning, while (Huang et al., 2023) utilizes Chain of Explanation (CoE). Their aim is to improve the effectiveness of LLM-based sentiment classifiers by leveraging reasoning capabilities. Likewise, Sun et al. (2023) introduced a technique called Clue and Reasoning Prompting (CARP), which incorporates both reasoning and keywords as clues to support the reasoning process. In the following subsections, we specifically examine approaches closely related to our task, focusing on those that have applied VLM-based methods for analyzing hateful or propagandistic memes.

CoT is a widely recognized prompting technique that generates a chain of reasoning to derive answers. A recent comprehensive CoT-based meme analysis study is presented in (Kumari et al., 2024), which proposed a framework based on text- and image-based entity-object relationships using a scene graph. They applied a hierarchical three-step CoT-based prompting strategy to guide the LLM in identifying Emotion, Target, and Context, using these elements to build a model for meme analysis. Another recent work, called SAFE-MEME (Nandi

et al., 2024), proposed two multimodal datasets and introduced a structured reasoning framework for hate speech detection in memes. They developed a CoT prompting-based framework that incorporates Q&A-style reasoning and hierarchical categorization to classify memes as hateful (explicit or implicit) or benign. However, they did not evaluate their approach using the popular Hateful Memes dataset, preventing a direct numerical comparison with their results.

One drawback of these CoT-based approaches is that they rely on multi-step reasoning, requiring multiple inferences with VLMs. Our approach differs from these CoT-based methods in the following ways: (a) we do not employ a complex multistep CoT approach, eliminating the need for multiple LLM inferences, which significantly improves computational efficiency and reduces costs. (b) we focus on providing explanations alongside classification, helping the end-users better understand the reasoning behind classification decisions, thereby increasing reliability.

(Hee et al., 2023) constructed a dataset providing explanations for hateful memes. However, unlike us, they focused solely on evaluating explanation generation and did not perform classification tasks. Despite the availability of their data, we do not use it due to the lack of *naturalness*. In particular, their explanations do not fully account for image content or image-centric contextual perspectives.

3 Dataset

3.1 ArMeme

The ArMeme dataset aimed to address the scarcity of Arabic-language datasets for multimodal propaganda detection. It comprises approximately ~6k Arabic memes collected from various social media platforms, each manually annotated to identify propagandistic content (Alam et al., 2024b). This dataset has been collected from different social media platforms, filtered, cleaned and manually annotated with four labels such as *Not propaganda*, *Propaganda*, *Not-meme* and *Other*. Table 1 provides the distribution of the data splits. The memes with “Not propaganda” category covers over half of the dataset (~66%), followed by “Propaganda” and the distribution of “Not-meme” and “Other” classes are significantly smaller. This distribution highlights a substantial class imbalance, particularly between “Not propaganda” and the other categories.

Class label	Train	Dev	Test	Total
Not propaganda	2,634	384	746	3,764
Propaganda	972	141	275	1,388
Not-meme	199	30	57	286
Other	202	29	56	287
Total	4,007	584	1,134	5,725

Table 1: Data splits for ArMeme datasets.

3.2 Hateful Meme

The Hateful Memes dataset (Kiel et al., 2020), is a benchmark designed to evaluate multimodal hate speech detection. It consists of $\sim 12k$ memes, combining both text and images, carefully curated to ensure that effective classification requires an understanding of both modalities. The dataset was created using a mix of synthetically generated memes and real-world examples, sourced from social media, while ensuring a balanced distribution of hateful and non-hateful content. A key feature of this dataset is the inclusion of benign confounders, where individual elements of a hateful meme—either the image or the text—are altered to make it non-hateful. This approach prevents unimodal models (which rely only on text or images) from achieving high performance, reinforcing the need for true multimodal understanding. In Table 2, we report the distribution of hateful meme dataset used for this study. Note that hateful meme dataset consists of two other splits (dev-seen and test-seen), here, we used unseen versions.

Class Label	Train	Dev-seen	Test-seen	Total
Not Hateful	5,481	253	510	6,244
Hateful	3,019	247	490	3,756
Total	8,500	500	1000	10,000

Table 2: Distribution of hateful meme dataset.

4 MemeXplain: Explanation Generation

The outcomes of an automatic system become more reliable for users if it provides decisions with adequate and interpretable natural explanations, which help users better understand the underlying reason behind the system’s decision (Hee et al., 2023; Yang et al., 2023; Huang et al., 2023; Sun et al., 2023). Technically, this approach provides numerous advantages in terms of knowledge distillation, model compression, and enhancing the performance of target tasks in different domains (Li et al.,

Data	Total Words	Avg. Words	Total Expl. Words		Avg. Expl. Words	
			Ar	En	Ar	En
ArMeme						
Train	58,688	15	280,341	375,843	70	94
Dev	8,583	15	40,756	55,336	70	95
Test	16,653	15	79,360	105,476	70	93
Total	83,924	15	400,457	536,655	70	94
Hateful Meme						
Train	99,812	12	–	740,624	–	87
Dev	4,904	9	–	43,956	–	81
Test	18,079	9	–	173,982	–	87
Total	122,795	10	–	958,562	–	85

Table 3: Descriptive statistics of the dataset. *Total Words* and *Avg.* refer to the total and average number of words in the text. The last two columns represent the corresponding values for the explanations.

2022; Magister et al., 2022; Nandi et al., 2024; Kumari et al., 2024). This motivates us to adopt the explanation-based approach in our research. However, we also aim to improve its efficiency, particularly with respect to dataset generation, model training, and system inference procedures.

In this research, we generate explanations for two different stages: (a) during existing dataset enhancement, which leverages an expert VLM (such as GPT) to generate high-quality explanations and (b) during training/inference with a smaller VLM (such as Llama-3.2 11b). Figure 1 illustrates these different stages. Mathematically, these two stages can be described by the functions $f(i, l) = e$ and $g(i) = (l, e)$, where e denotes the explanation, l is the label, and i is the input image or meme. Specifically, $f(i, l)$ returns an explanation e given both i and l , whereas $g(i)$ generates both the label l and the explanation e from only the input i .

This research enhances two existing datasets with explanations, see Section 3 and Table 3 for the details and statistics. For the explanation generation task, it first uses a VLM for $f(i, l)$ and then involves human experts, which significantly accelerates high-quality explanation generation and lowers the overall cost and time. The following subsections provide step-by-step details.

4.1 VLMs for Explanation Generation

Figure 1 illustrates an example of an Arabic meme along with its explanation-generation process using a VLM. We leverage GPT-4o (version 2024-

11-20) for automated explanation generation. The choice of this model is motivated by prior studies Wang et al. (2023), which show that advanced GPT models can produce fluent, informative, persuasive, and logically sound explanations when properly prompted. In Listing 1, we present the *prompts* used for generating explanations for **ArMeme** and **Hateful Memes**. To refine the prompt, we iteratively tested several memes in both English and Arabic, selecting the one that produced the most reasonable explanations.

For Arabic memes, we generate two sets of explanations—one *in English* and one *in Arabic*. The motivation behind this approach is to assess the multilingual capability and quality of smaller VLMs, such as Llama-3.2 11b, in generating explanations and labels in both languages.

Size of the Explanation Determining the optimal length for explanations is important for balancing informativeness and cognitive load (Herm, 2023). Shen et al. (2022) explored the relationship between explanation length and human understanding, finding that the shortest rationales are often ineffective. Recently, Wang et al. (2023) also studied the effect of explanation size and found that human evaluators are reluctant to read longer explanations. To achieve an optimal balance, we iteratively tested various explanation lengths and ultimately set a limit of 100 words.

Model and Its Parameters To utilize GPT-4o (OpenAI, 2023), we accessed the OpenAI API via Azure services. Though recently released o1 models have shown promising directions for complex reasoning, they were not accessible to us. For explanation generation, we employed zero-shot learning. To ensure reproducibility, we set the temperature value to zero.

4.2 Human Evaluation

Given that our idea is to use the generated explanation as gold data for further training and evaluation, therefore, we intended to go through human evaluation process. Following the prior studies (Wang et al., 2023; Huang et al., 2023; Agarwal et al., 2024) we adopted four metrics discussed below. For each metric we use 5-point Likert scale.

Informativeness. Measures the extent to which the explanation provides relevant and meaningful information for understanding the reasoning behind the label. A highly informative explanation offers

detailed insights that directly contribute to the justification, while a low-informative explanation may be vague, incomplete, or lacking key details.

Clarity. Assesses how clearly the explanation conveys its meaning. A clear explanation is well-structured, concise, and easy to understand without requiring additional effort. It should be free from ambiguity, overly complex language, or poor phrasing that might hinder comprehension.

Plausibility. Refers to the extent to which an explanation logically supports the assigned label and appears reasonable given the meme’s content. A plausible explanation should be coherent, factually consistent, and align with the expected reasoning behind the label.

Faithfulness. Measures how accurately an explanation reflects the reasoning behind the assigned label. A faithful explanation correctly represents the key factors and logical steps that justify the label, without adding misleading or unrelated details.

For manual annotation, we first prepared an annotation guideline for the annotators. Additionally, we developed annotation guidelines and a platform (see Appendix B and A, respectively).

Evaluation Setting. For the Arabic meme task, we recruited annotators who are native Arabic speakers and fluent in English, all holding at least a bachelor’s degree. Because of their fluency, they also handled the hateful meme task. We provided necessary training and consultation, and all had prior experience with similar tasks.

A total of six annotators participated in the evaluation. In line with institutional requirements, each signed a Non-Disclosure Agreement (NDA), and a third-party company managed their compensation at standard hourly rates based on location.

Quality Assessment In Table 4, we summarize the quality assessment of the explanations. We used 5-point Likert scale for various human evaluation metrics, including informativeness, clarity, plausibility, and faithfulness. We compute the average of the Likert scale value for all evaluation metrics. We manually evaluated 359 and 202 random samples for ArMeme Arabic and English explanations while 200 random examples were evaluated for the Hateful meme dataset. The average agreement scores for the ArMeme dataset with Arabic explanations are 4.23, 4.38, 4.24, and 4.16 for faithfulness, clarity, plausibility, and informativeness, respectively, indicating high agreement across all evaluation metrics. However, for the English expla-

nations of ArMeme, the faithfulness and plausibility scores are relatively. To better understand this issue, we plan to conduct further evaluations on another set of explanations. For the Hateful Memes dataset, the average Likert scale agreement scores range from 4.562 to 4.682.

Dataset	Faithfulness	Clarity	Plausibility	Informative
ArMeme (Ar)	4.23	4.38	4.24	4.16
ArMeme (En)	3.91	4.50	3.81	4.13
Hateful meme	4.56	4.65	4.63	4.68

Table 4: Average Likert scale value for each human evaluation metric across different sets of explanations.

4.3 Basic Statistics

Table 3 presents the basic statistics for both datasets. The average explanation length is 94 words for Arabic and 85 words for English. Notably, we instructed GPT-4o to generate explanations with fewer than 100 words. Based on manual evaluation (Table 4), we conclude that both the quality and length of the explanations are appropriate.

5 Methodology

5.1 Instructions Dataset

Our approach follows the standard pipeline for aligning LLMs with user intentions and specific tasks through fine-tuning on representative data (Zhang et al., 2023; Kmainasi et al., 2024). This process typically involves curating and constructing instruction datasets that guide the model’s behavior, ensuring it generates responses that align with the desired objectives. For our study, the responses include label and explanation. Hence, we created instruction format for both datasets. For the ArMeme dataset, we replicated the experiments for both Arabic and English explanations.

5.2 Model Selection

As shown in Figure 1, our first experimental phase involves model selection among several recent VLMs, including Llama-3.2 (11b) (Dubey et al., 2024), Paligemma 2 (3b) (Steiner et al., 2024), Qwen2-vl (Wang et al., 2024), and Pixtral (12b) (Agrawal et al., 2024).

We evaluate the base models in a zero-shot setting and fine-tune them using an instruction-following paradigm. The instructions prompt the model to generate responses in the format “Label: (class_label)”. We use a regex-based function to extract the predicted labels.

Note that this stage fine-tunes the models to predict class labels only, allowing us to verify whether they can handle multilingual inputs—especially in understanding Arabic text, cultural nuances, and image context. We do not ask the model to generate explanations here, as that is a more complex task and could affect their performance.

Based on the results reported in Tables 5 and 6, we selected Llama-3.2-vision-instruct (11b) for further training with explanations.

5.3 Multi-Stage (MS) Optimization Procedure

To emphasize our novel contribution, we introduce a dedicated optimization procedure to train VLM with *MemeXplain*, which decouples the classification and explanation generation tasks. This approach is designed to first endow the model with strong task-specific representations through classification-only fine-tuning, and then refine its ability to generate coherent, natural explanations.

Stage 1: Classification Fine-Tuning In this stage, the model is fine-tuned solely on the classification task. The training objective is restricted to predicting the correct class label. This focused objective encourages the model to develop robust, task-specific representations. We use the QLoRA setup described later with a learning rate of 2×10^{-4} to optimize the model.

Stage 2: Explanation Enhancement In this stage, the model is further fine-tuned on a combined label-with-explanation dataset. Here, we employ a reduced learning rate (1×10^{-5}) to gently adapt the model’s parameters for generating the explanations while preserving the classification performance achieved in Stage 1.

To validate the effectiveness of the multi-stage procedure, we compare it against a single-stage (SS) fine-tuning baseline where the model is directly trained on the label-with-explanation dataset. Our ablation studies (detailed in Section 6) demonstrate that the proposed multi-stage approach significantly outperforms the single-stage strategy.

5.4 Training Setup

Our fine-tuning experiments utilize QLoRA (Dettmers et al., 2023), which combines INT4 quantization with parameter-efficient fine-tuning through Low-Rank Adaptation (LoRA) (Hu et al., 2022). In our setup, the base model is quantized to 4-bit precision, with LoRA updates applied to a subset of the model

Model	Setup	Acc (%)	W-F1	M-F1
(Alam et al., 2024b)	Qarib	69.7	0.690	0.551
(Alam et al., 2024b)	mBERT	70.7	0.675	0.487
Llama-3.2 (11b)	Base	13.4	0.172	0.113
Llama-3.2 (11b)	FT	68.0	0.665	0.452
Paligemma2 (3b)	Base	15.3	0.090	0.080
Paligemma2 (3b)	FT	65.9	0.524	0.200
Qwen2 (7b)	Base	63.1	0.550	0.242
Qwen2 (7b)	FT	27.0	0.149	0.195
Pixtral (12b)	Base	14.6	0.177	0.133
Pixtral (12b)	FT	70.8	0.636	0.377

Table 5: Results for ArMeme. FT: Fine-tuned. Qarib (Abdelali et al., 2021) is a Arabic BERT (text only). mBERT - multilingual BERT (text only).

Model	Setup	Acc (%)	W-F1	M-F1
(Kiel et al., 2020)		69.47±2.06		
(Cao et al., 2022)		72.98±1.09		
Llama-3.2 (11b)	Base	66.1	0.650	0.618
Llama-3.2 (11b)	FT	77.7	0.770	0.748
Paligemma2 (3b)	Base	35.2	0.277	0.217
Paligemma2 (3b)	FT	69.2	0.664	0.623
Qwen2 (7b)	Base	66.4	0.669	0.442
Qwen2 (7b)	FT	77.9	0.773	0.753
Pixtral (12b)	Base	66.7	0.667	0.430
Pixtral (12b)	FT	77.2	0.766	0.746

Table 6: Results for Hateful meme. FT: Fine-tuned

parameters. This approach was selected to address computational resource constraints. Furthermore, deploying models for inference incurs significant costs. Therefore, we focus on quantized models and assessing their performance accordingly.

For all experiments, we fine-tuned the models using the QLoRA approach with 4-bit quantization. This approach was chosen due to its efficiency in reducing memory usage while maintaining model performance. We adapted all relevant submodules (vision, language, attention, and MLP layers) with a LoRA rank of 16, an alpha of 16, and no dropout. For training, we used a per-device batch size of 2 with gradient accumulation over 4 steps and optimized using AdamW with a learning rate of 2×10^{-4} , a weight decay of 0.01, and a linear scheduler with 5 warmup steps. For the second stage experiments (label-with-explanation), the learning rate was reduced to 1×10^{-5} .

5.5 Evaluation Setup and Metrics

We train the models using the training set, fine-tune the parameters with the development set, and evaluate their performance on the test set as reported in Tables 5 and 6. For performance measurement

across different experimental settings, we compute accuracy, weighted F_1 score, and macro- F_1 score. We evaluate the model’s explanation performance on the test set using semantic similarity-based metric, measured by the F_1 score within BERTScore (Zhang et al., 2020). This score is computed using contextual embeddings extracted from pre-trained BERT models. To enhance accuracy, we utilize language-specific transformer models for embedding extraction. For Arabic we use AraBERT (v2) (Antoun et al., 2020) model and for English we use bert-base-uncased model (Devlin et al., 2019). Although metrics such as BLEU and ROUGE are commonly used, studies have reported their limitations (Xu et al., 2023; Krishna et al., 2021). Therefore, we rely solely on BERTScore.

6 Experimental Results and Discussion

This section first presents competitive results among our proposed method and the state-of-the-art approaches. Next, it briefly analyzes and investigates the proposed method to validate and highlight the core contributions of this research.

Table 7 compares our proposed models with state-of-the-art approaches. On the ArMeme dataset, our method achieves the best accuracy at 72.1% and the best weighted F_1 at 0.699, with Qarib and mBERT following behind. Although Qarib attains the highest macro F_1 (0.551), our model remains competitive with a macro F_1 of 0.536. Importantly, our method stands out because it provides explanations that add significant value. On the Hateful Meme dataset, our approach clearly outperforms the state-of-the-art by achieving the best performance with an accuracy of 79.9%, a weighted F_1 of 0.802, and a macro F_1 of 0.792. These results clearly highlight the advantages of our explainability-enhanced dataset and the proposed multi-stage optimization procedure for both classification and explanation-generation tasks.

Table 8 provides classification and explanation-generation results on the *ArMeme* and *Hateful Meme* datasets. It briefly presents these results from several perspectives: (a) **Base vs. FT**: demonstrates the performance difference between the same model with and without fine-tuning (FT); (b) **Single-stage (SS) vs. Multi-stage (MS)**: highlights the necessity and benefits of the proposed optimization procedure and (c) **Eng-Exp vs. Ar-Exp**: showcases the multilingual capability of the selected VLM. Next, we provide a brief analysis of

Model	Setup	Acc(%)	W-F1	M-F1
ArMeme				
(Alam et al., 2024b)	Qarib	69.7	0.690	0.551
(Alam et al., 2024b)	mBERT	70.7	0.675	0.487
(Alam et al., 2024b)	ResNet50	66.0	0.637	0.434
Llama MS	FT	72.1	0.699	0.536
Llama (Ar-Exp) MS	FT	72.0	0.696	0.499
Hateful Meme				
(Kiela et al., 2020)		69.47±2.06		
(Cao et al., 2022)		72.98±1.09		
Llama MS	FT	79.9	0.802	0.792

Table 7: Comparison with SOTA and our results. ResNet50 (He et al., 2016) is an image only model. MS: Multi-stage.

the results based on these perspectives.

First, we compare the **Base vs. FT** setup, from which it is evident that the FT model significantly outperforms the baseline. For example, on the ArMeme dataset, while the baseline achieves an accuracy of 12.7%, the proposed fine-tuning boosts it to 72.1%. Similarly, on the Hateful Meme dataset, fine-tuning improves the base accuracy from 65.2% to 79.9%. We observe similar improvements in the F1 metrics for classification and BERTScore for explanation quality. These significant performance gains *validate our approach of fine-tuning the base models with the explainability enhanced dataset*, demonstrating its efficacy for the meme classification and explanation generation tasks.

Next, we compare the **SS vs. MS** setup, which reveals that multi-stage (MS) fine-tuning further enhances performance over the single-stage (SS) approach. For example, on the ArMeme dataset, the accuracy increased from 68.2% to 72.1%, the weighted F1 increased from 0.584 to 0.699, the macro F1 increased significantly from 0.257 to 0.536, and the BERTScore for Arabic explanation increased significantly from 0.58 to 0.72. A similar trend is observed on the Hateful Meme dataset, where additional fine-tuning iterations yield more robust classification (approximately 4% improvement) and enhanced explanation quality. These performance gains *validate our proposed multi-stage optimization procedure* to further refine the VLMs.

Finally, we assess the model’s multilingual capability by comparing the performance of **Llama MS - FT** with **Llama MS Ar-Exp**. The results show that fine-tuning using explanations generated in both languages yields comparable outcomes. This *validates the multilingual capability of our empirically chosen VLM* for the target task and enables users to

understand multilingual content even if they are not fluent in that language. For example, our model allows an English speaker to analyze Arabic memes and receive explanations in English.

Model	Setup	Acc (%)	W-F1	M-F1	BS
ArMeme					
Llama	Base	12.7	0.165	0.105	0.61
Llama SS	FT	68.2	0.584	0.257	0.70
Llama MS	FT	72.1	0.699	0.536	0.70
Llama Ar-Exp	Base	19.0	0.246	0.125	0.58
Llama MS Ar-Exp	FT	72.0	0.696	0.499	0.72
Hateful Meme					
Llama	Base	65.2	0.615	0.567	0.661
Llama SS	FT	75.9	0.760	0.745	0.767
Llama MS	FT	79.9	0.802	0.792	0.777

Table 8: Results with ArMeme and Hateful meme classification and explanation generation. Llama: Llama-3.2 (11b), BS: BERTScore. SS: Single-stage, MS: Multi-stage. Ar-Exp: Model trained with Arabic explanation.

7 Conclusions and Future Work

In this study, we introduce a *MemeXplain* dataset for propagandistic and hateful meme detection and natural explanation generation, making it the *first* resource of its kind. To address both detection and explanation generation tasks and ensure efficient VLMs model training on this dataset, we also propose a multi-stage optimization procedure. To evaluate the multilingual capability of the model, we developed Arabic and English explanations for Arabic memes. The inclusion of English explanations benefits non-Arabic speakers, whereas providing explanations in the native language ensures that cultural nuances are accurately conveyed. With our multi-stage training procedure, we demonstrate improved detection performance for both *ArMeme* and hateful memes. The higher performance of explanation generation further demonstrates the efficacy of our multi-stage training approach. We foresee several future directions to extend this research and explore the following: (a) training the model with additional data through data augmentation, which could help it become an instruction-generalized model and potentially enhance its performance further; (b) incorporating pseudo and self-labeled data using an active learning procedure to incrementally improve the model’s capabilities; and (c) developing a task-generalized model that addresses multiple tasks.

8 Limitations

Due to the complex nature of manual explanation creation, we have relied on GPT-4o for explanation generation. To ensure the reliability of the explanation we have manually evaluated in four criteria such as informativeness, clarity, plausibility, and faithfulness on a small sample for each set of explanation. The preliminary evaluation scores suggest that we can rely on the gold explanation as the reference. As a part of ongoing work we plan to conduct manual evaluation on a larger set. An important aspect of the ArMeme dataset is that it is highly imbalanced, which affects overall performance. One possible approach to address this issue is to increase the number of memes labeled as propaganda, other, and not-meme. This can be achieved through data augmentation or by collecting additional memes.

Ethics and Broader Impact

We extended existing datasets by adding explanations. To the best of our knowledge, the dataset does not contain any personally identifiable information, making privacy risks nonexistent. Regarding the explanations, we provided clear annotation instructions and cautioned annotators that some memes might be offensive. It is important to note that annotations are inherently subjective, which can introduce biases into the overall evaluation results. We encourage researchers and users of this dataset to remain critical when developing models or conducting further research. Models built using this dataset could be highly valuable for fact-checkers, journalists, and social media platforms.

References

Ahmed Abdelali, Sabit Hassan, Hamdy Mubarak, Kareem Darwish, and Younes Samih. 2021. [Pre-training bert on arabic tweets: Practical considerations](#).

Chirag Agarwal, Sree Harsha Tanneru, and Himabindu Lakkaraju. 2024. [Faithfulness vs. plausibility: On the \(un\)reliability of explanations from large language models](#). *arXiv preprint arXiv:2402.04614*, 2402.04614.

Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. 2024. Pixtral 12b. *arXiv preprint arXiv:2410.07073*.

Firoj Alam, Md Rafiul Biswas, Uzair Shah, Wajdi Zaghoulani, and Georgios Mikros. 2024a. Propaganda

to hate: A multimodal analysis of arabic memes with multi-agent llms. In *International Conference on Web Information Systems Engineering*, pages 380–390. Springer.

Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fabrizio Silvestri, Dimiter Dimitrov, Giovanni Da San Martino, Shaden Shaar, Hamed Firooz, and Preslav Nakov. 2022. A survey on multimodal disinformation detection. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6625–6643.

Firoj Alam, Abul Hasnat, Fatema Ahmad, Md. Arid Hasan, and Maram Hasanain. 2024b. [ArMeme: Propagandistic content in Arabic memes](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21071–21090, Miami, Florida, USA. Association for Computational Linguistics.

Wissam Antoun, Fady Baly, and Hazem Hajj. 2020. AraBERT: Transformer-based model for Arabic language understanding. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 9–15.

Alberto Barrón-Cedeno, Israa Jaradat, Giovanni Da San Martino, and Preslav Nakov. 2019. Propy: Organizing the news based on their propagandistic content. *Information Processing & Management*, 56(5):1849–1864.

Rui Cao, Ming Shan Hee, Adriel Kuek, Wen-Haw Chong, Roy Ka-Wei Lee, and Jing Jiang. 2023. [Procap: Leveraging a frozen vision-language model for hateful meme detection](#). In *Proceedings of the 31st ACM International Conference on Multimedia*, MM ’23, page 5244–5252, New York, NY, USA. Association for Computing Machinery.

Rui Cao, Roy Ka-Wei Lee, Wen-Haw Chong, and Jing Jiang. 2022. [Prompting for multimodal hateful meme classification](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 321–332, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Pengyuan Chen, Lei Zhao, Yangheran Piao, Hongwei Ding, and Xiaohui Cui. 2024. Multimodal visual-textual object graph attention network for propaganda detection in memes. *Multimedia Tools and Applications*, 83(12):36629–36644.

Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeño, Rostislav Petrov, and Preslav Nakov. 2019. Fine-grained analysis of propaganda in news article. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Hong Kong, China.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning of quantized llms](#). *Preprint*, arXiv:2305.14314.

772	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , NAACL-HLT '19, Minneapolis, Minnesota, USA.	828
773		829
774		830
775		
776		831
777		832
778		833
779		834
780	Dimitar Dimitrov, Firoj Alam, Maram Hasanain, Abul Hasnat, Fabrizio Silvestri, Preslav Nakov, and Giovanni Da San Martino. 2024. Semeval-2024 task 4: Multilingual detection of persuasion techniques in memes. In <i>Proceedings of the 2024 Annual Conference of the North American Chapter of the Association for Computational Linguistics</i> .	835
781		836
782		837
783		838
784		
785		839
786		840
787	Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021a. Detecting propaganda techniques in memes. In <i>ACL-IJCNLP</i> .	841
788		842
789		843
790		
791		844
792	Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021b. SemEval-2021 task 6: Detection of persuasion techniques in texts and images . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 70–98, Online. Association for Computational Linguistics.	845
793		846
794		847
795		848
796		
797		849
798		850
799		851
800	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	852
801		853
802		854
803		855
804		
805	Paula Fortuna and Sérgio Nunes. 2018. A survey on automatic detection of hate speech in text. <i>ACM Computing Surveys (CSUR)</i> , 51(4):1–30.	856
806		857
807		858
808	Maria Glenski, E. Ayton, J. Mendoza, and Svitlana Volkova. 2019. Multilingual multimodal digital deception detection and disinformation spread across social platforms. <i>ArXiv</i> , abs/1909.05838.	859
809		860
810		861
811		862
812	Ivan Habernal, Raffael Hannemann, Christian Pol-lak, Christopher Klammer, Patrick Pauli, and Iryna Gurevych. 2017. Argotario: Computational argumentation meets serious games. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , EMNLP '17, pages 7–12, Copenhagen, Denmark.	863
813		864
814		865
815		866
816		867
817		868
818		869
819	Ivan Habernal, Patrick Pauli, and Iryna Gurevych. 2018. Adapting serious game for fallacious argumentation to German: Pitfalls, insights, and best practices. In <i>LREC</i> . European Language Resources Association (ELRA).	870
820		
821		871
822		872
823		873
824	Maram Hasanain, Firoj Alam, Hamdy Mubarak, Samir Abdaljalil, Wajdi Zaghouani, Preslav Nakov, Giovanni Da San Martino, and Abed Freihat. 2023. ArAIEval shared task: Persuasion techniques and	874
825		875
826		
827		876
		877
		878
		879
		880
		881
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881

882	Mohamed Bayan Kmainasi, Ali Ezzat Shahroor, Maram	Shivam Sharma, Firoj Alam, Md. Shad Akhtar, Dimitar	939
883	Hasanain, Sahinur Rahman Laskar, Naeemul Has-	Dimitrov, Giovanni Da San Martino, Hamed Firooz,	940
884	san, and Firoj Alam. 2024. LlamaLens: Specialized	Alon Halevy, Fabrizio Silvestri, Preslav Nakov, and	941
885	multilingual llm for analyzing news and social media	Tanmoy Chakraborty. 2022. Detecting and under-	942
886	content. <i>arXiv preprint arXiv:2410.15308</i> .	standing harmful memes: A survey . In <i>Proceedings</i>	943
887	Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021.	<i>of the Thirty-First International Joint Conference on</i>	944
888	Hurdles to progress in long-form question answering .	<i>Artificial Intelligence, IJCAI '22</i> , pages 5597–5606,	945
889	In <i>Proceedings of the 2021 Conference of the North</i>	Vienna, Austria. International Joint Conferences on	946
890	<i>American Chapter of the Association for Computa-</i>	Artificial Intelligence Organization. Survey Track.	947
891	<i>tional Linguistics: Human Language Technologies</i> ,		
892	pages 4940–4957, Online. Association for Computa-	Hua Shen, Tongshuang Wu, Wenbo Guo, and Ting-Hao	948
893	tional Linguistics.	Huang. 2022. Are shortest rationales the best expla-	949
894	Gokul Karthik Kumar and Karthik Nandakumar. 2022.	nations for human understanding? In <i>Proceedings</i>	950
895	Hate-CLIPper: Multimodal hateful meme classifica-	<i>of the 60th Annual Meeting of the Association for</i>	951
896	tion based on cross-modal interaction of clip features.	<i>Computational Linguistics (Volume 2: Short Papers)</i> ,	952
897	In <i>Proceedings of the Second Workshop on NLP for</i>	pages 10–19.	953
898	<i>Positive Impact (NLP4PI)</i> , pages 171–183.	Andreas Steiner, André Susano Pinto, Michael Tschan-	954
899	Gitanjali Kumari, Kirtan Jain, and Asif Ekbal. 2024.	nen, Daniel Keysers, Xiao Wang, Yonatan Bitton,	955
900	M3hop-cot: Misogynous meme identification with	Alexey Gritsenko, Matthias Minderer, Anthony Sher-	956
901	multimodal multi-hop chain-of-thought. <i>arXiv</i>	bondy, Shangbang Long, et al. 2024. Paligemma 2:	957
902	<i>preprint arXiv:2410.09220</i> .	A family of versatile vlms for transfer. <i>arXiv preprint</i>	958
903	Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen,	<i>arXiv:2412.03555</i> .	959
904	Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian,	Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei	960
905	Baolin Peng, Yi Mao, et al. 2022. Explanations from	Guo, Tianwei Zhang, and Guoyin Wang. 2023. Text	961
906	large language models make small reasoners better.	classification via large language models . In <i>Find-</i>	962
907	<i>arXiv preprint arXiv:2210.06726</i> .	<i>ings of the Association for Computational Linguis-</i>	963
908	Lucie Charlotte Magister, Jonathan Mallinson, Jakub	<i>tics: EMNLP 2023</i> , pages 8990–9005. Association	964
909	Adamek, Eric Malmi, and Aliaksei Severyn. 2022.	for Computational Linguistics.	965
910	Teaching small language models to reason. <i>arXiv</i>	Riza Velioglu and Jewgeni Rose. 2020. Detecting hate	966
911	<i>preprint arXiv:2212.08410</i> .	speech in memes using multimodal deep learning	967
912	Palash Nandi, Shivam Sharma, and Tanmoy	approaches: Prize-winning solution to hateful memes	968
913	Chakraborty. 2024. SAFE-MEME: Structured	challenge. <i>arXiv preprint arXiv:2012.12975</i> .	969
914	reasoning framework for robust hate speech detec-	Han Wang, Ming Shan Hee, Md Rabiul Awal, Kenny	970
915	tion in memes. <i>arXiv preprint arXiv:2412.20541</i> .	Tsu Wei Choo, and Roy Ka-Wei Lee. 2023. Evaluat-	971
916	OpenAI. 2023. GPT-4 technical report . Technical re-	ing GPT-3 generated explanations for hateful content	972
917	port, OpenAI.	moderation. In <i>Proceedings of the Thirty-Second</i>	973
918	Jakub Piskorski, Nicolas Stefanovitch, Nikolaos Niko-	<i>International Joint Conference on Artificial Intelli-</i>	974
919	laidis, Giovanni Da San Martino, and Preslav Nakov.	<i>gence</i> , pages 6255–6263.	975
920	2023. Multilingual multifaceted understanding of on-	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	976
921	line news in terms of genre, framing, and persuasion	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin	977
922	techniques . In <i>Proceedings of the 61st Annual Meet-</i>	Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhanc-	978
923	<i>ing of the Association for Computational Linguistics</i>	ing vision-language model’s perception of the world	979
924	<i>(Volume 1: Long Papers)</i> , pages 3001–3022, Toronto,	at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	980
925	Canada. Association for Computational Linguistics.	Fangyuan Xu, Yixiao Song, Mohit Iyyer, and Eunsol	981
926	Aske Plaat, Annie Wong, Suzan Verberne, Joost	Choi. 2023. A critical evaluation of evaluations for	982
927	Broekens, Niki van Stein, and Thomas Back. 2024.	long-form question answering. In <i>Proceedings of the</i>	983
928	Reasoning with large language models, a survey.	<i>61st Annual Meeting of the Association for Compu-</i>	984
929	<i>arXiv e-prints</i> , pages arXiv–2407.	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	985
930	Uzair Shah, Md. Rafiul Biswas, Marco Agus, Mowafa	3225–3245.	986
931	Househ, and Wajdi Zaghouni. 2024. MemeMind at	Yongjin Yang, Joonkee Kim, Yujin Kim, Namgyu Ho,	987
932	ArAIEval shared task: Generative augmentation and	James Thorne, and Se-Young Yun. 2023. HARE:	988
933	feature fusion for multimodal propaganda detection	Explainable hate speech detection with step-by-step	989
934	in Arabic memes through advanced language and	reasoning. In <i>Findings of the Association for Com-</i>	990
935	vision models . In <i>Proceedings of The Second Ara-</i>	<i>putational Linguistics: EMNLP 2023</i> , pages 5490–	991
936	<i>bic Natural Language Processing Conference</i> , pages	5505.	992
937	467–472, Bangkok, Thailand. Association for Com-	Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu.	993
938	putational Linguistics.	2024. Vision-language models for vision tasks: A	994
		survey. <i>IEEE Transactions on Pattern Analysis and</i>	995
		<i>Machine Intelligence</i> .	996

Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with BERT. In *International Conference on Learning Representations*.

A Annotation Guideline

You will be shown a meme, a label assigned to it, and an explanation for the assigned label. As an annotator, your task is to carefully examine each meme, label, and explanation. Then assess the quality of the explanation provided for the assigned label. Follow the steps below to ensure a thorough evaluation:

Analyze the Meme

- Observe the image and read the accompanying text.
- Understand the overall message and the potential implications of the meme.

Check the Assigned Label

- Check the given label. The label is the result of annotation done by multiple human annotators.

Evaluate the Explanation

- Read the explanation provided for why the meme has been assigned its label.
- Assess the explanation based on the metrics below. Each metric is scored on a Likert scale from 1-5.

Kindly note that to evaluate the explanation, you do not have to agree or disagree with the given label.

A.1 Metrics

A.1.1 Informativeness

Measures the extent to which the explanation provides relevant and meaningful information for understanding the reasoning behind the label. A highly informative explanation offers detailed insights that directly contribute to the justification, while a low-informative explanation may be vague, incomplete, or lacking key details.

As an annotator, you are judging if the explanation provides enough information to explain the label assigned to the meme.

- 1 = Not informative: The explanation lacks relevant details and does not help understand why the meme is labeled as such.
- 2 = Slightly informative: The explanation provides minimal information, but key details are missing or unclear.
- 3 = Moderately informative: The explanation contains some useful details but lacks depth or supporting reasoning.
- 4 = Informative: The explanation is well-detailed, providing a clear and meaningful justification for the label.
- 5 = Very informative: The explanation is thorough, insightful, and fully justifies the label with strong supporting details.

A.1.2 Clarity

Assesses how clearly the explanation conveys its meaning. A clear explanation is well-structured, concise, and easy to understand without requiring additional effort. It should be free from ambiguity, overly complex language, or poor phrasing that might hinder comprehension.

As an annotator, you are judging the language and structure of the explanation. Spelling mistakes, awkward use of language, and incorrect translations will negatively impact this metric.

- 1 = Very unclear: The explanation is confusing, vague, or difficult to understand.
- 2 = Somewhat unclear: The explanation has some clarity but includes ambiguous or poorly structured statements.
- 3 = Neutral: The explanation is somewhat clear but may require effort to fully grasp.
- 4 = Clear: The explanation is well-structured and easy to understand with minimal ambiguity.
- 5 = Very clear: The explanation is highly readable, precise, and effortlessly understandable.

A.1.3 Plausibility

Refers to the extent to which an explanation logically supports the assigned label and appears reasonable given the meme's content. A plausible explanation should be coherent, factually consistent, and align with the expected reasoning behind the label. While it does not require absolute correctness, it should not contain obvious contradictions or illogical claims.

As an annotator, you are judging if the explanation actually supports the label assigned to the meme. For example, if a meme is labeled as

Not Propaganda, the explanation given should justify that label.

- 1 = Not plausible at all: The explanation does not align with the label and seems completely incorrect.
- 2 = Weakly plausible: The explanation has some relevance but lacks strong justification or contains logical inconsistencies.
- 3 = Moderately plausible: The explanation somewhat supports the label but may be incomplete or partially flawed.
- 4 = Plausible: The explanation logically supports the label and is mostly reasonable.
- 5 = Highly plausible: The explanation is fully aligned with the label and presents a strong, logical justification.

A.1.4 Faithfulness

Measures how accurately an explanation reflects the reasoning behind the assigned label. A faithful explanation correctly represents the key factors and logical steps that justify the label, without adding misleading or unrelated details. High faithfulness means the explanation stays true to the actual reasoning used for classification, ensuring reliability and consistency.

As an annotator, you are judging how well the explanation reflects the logic behind the label. For example, if the explanation claims an implication of the meme, it should also present the logical reasoning behind it.

- 1 = Not faithful at all: The explanation is completely unrelated to the given label and does not reflect a valid reasoning process.
- 2 = Weakly faithful: Some elements of the explanation are relevant, but much of it is misleading, inconsistent, or lacks proper justification.
- 3 = Moderately faithful: The explanation captures parts of the reasoning but includes unrelated, unclear, or unnecessary justifications.
- 4 = Faithful: The explanation aligns well with the reasoning behind the label and includes relevant, logical details.
- 5 = Highly faithful: The explanation fully and accurately reflects the correct reasoning, without any misleading or irrelevant information.

B Annotation Platform

In Figure 2, we present the screenshot of the interface designed for the explanation evaluation of

hateful meme, which consisted of an image, respective label, and explanation for the label, annotation guidelines, and four different evaluation metrics. We used 5-point Likert scale for each evaluation metric. Annotators select one of the Likert scale value following the annotation guideline for each metric and submit.

C Prompt for Explanation Generation


In Listings 1 and 2, we provide the prompts used to generate explanations for ArMeme and Hateful Meme. The prompt in Listing 1 is specifically for generating Arabic explanations for ArMeme. To generate English explanations, the same prompt was used, except it was adapted for English.

```
You are a Propaganda Image Detection Expert. A human expert has already classified the image as {class_label}. Do not change or re-identify the classified label of the image. First, analyze the image's visual elements (objects, symbols, color usage, composition) and provide a concise description. Next, read and paraphrase any text in the image-especially if it is in non-Arabic-so that a Arabic speaker could understand its literal or intended meaning. Explain how that text, in conjunction with the visual elements, supports or illustrates the known propaganda technique. Think carefully about emotional appeals, humor, cultural references, manipulative language, or other rhetorical strategies. Finally, provide a thoughtful explanation in Arabic (up to {explanation_length} words) of why these elements led the human to classify the image as propaganda. Be specific about how the text and visuals interact, referencing the Arabic context or cultural nuances if relevant. Your answer must be valid JSON with a single field:
{{
  "explanation": "Your explanation here (up to {explanation_length} words) in Arabic."
}}
```

Listing 1: Prompt for generating Arabic explanation for ArMeme dataset. The place-holders refers to label, and number of words, respectively.

```
You are a Hateful Image Detection Expert. A human expert has already classified the image as {}. Do not change or re-identify the classified label of the image. First, analyze the image's visual elements (objects, symbols, color usage, composition) and provide a concise description. Next, read and paraphrase any text in the image so
```

English Hateful Memes Explanation - Verification



Label:
not-hateful

Explanation:
The image features a lion lying in the grass with text that reads, 'hmm today i will commit 50% of all violent crimes despite being only 13% percent of the population' and 'being oppressed time.' The text references a controversial and often misused statistic about crime rates and demographics, which can perpetuate harmful stereotypes. The juxtaposition of the lion, a symbol of strength, with the text may imply a critique or mockery of these stereotypes. The human expert likely found it not hateful due to the satirical or ironic tone, though it touches on sensitive topics.

Informativeness
Select Informativeness

Clarity
Select Clarity

Plausibility
Select Plausibility

Faithfulness
Select Faithfulness

Annotation Guidelines

Figure 2: A screenshot of the annotation platform for the explanation evaluation of hateful meme.

anyone speaking English can understand its literal or intended meaning. Explain how that text, in conjunction with the visual elements, supports or illustrates the known hateful content. Finally, provide a thoughtful explanation in English (up to {} words) of why these elements led the human to classify the image as hateful. Be specific about how the text and visuals interact, referencing the context or cultural nuances if relevant. Your answer must be valid JSON with a single field:

```
{{
  "explanation": "Your explanation here (up to {} words) in English."
}}
```

Listing 2: Prompt for generating explanation. The place-holders refers to label, and number of words, respectively.

D Error Analysis

In Figure 3, we present examples from the Hateful Memes dataset, showcasing cases where the model made both correct and incorrect predictions.

In figure 3a, the Gold explanation describes the image as reinforcing a harmful racial stereotype by juxtaposing a joyful scene of Asian individuals eating with offensive text. The Predicted explanation correctly identifies the derogatory language and its racist implications, aligning with the gold annotation. The model’s BERT-F1 score of 0.873 shows the high confidence in associating textual and visual elements to detect hate speech effectively.


In figure 3, the Gold explanation interprets the image as a humorous juxtaposition, using wordplay

between nationality and species without targeting any group. However, the Predicted explanation classifies it as hateful. This missclassification suggests that the model struggled to distinguish linguistic humor from implicit hate speech, as reflected in its BERT-F1 score of 0.6259. This highlights the challenge of detecting context-dependent content, where intent and interpretation play a crucial role in classification.

E Data Release

The *MemeXplain* dataset² will be released under the CC BY-NC-SA 4.0 – Creative Commons Attribution 4.0 International License: <https://creativecommons.org/licenses/by-nc-sa/4.0/>.

²anonymous.com



Gold: Hateful **Predicted:** Hateful
BERT-F1 : 0.873


Explanation

Gold: The image shows a group of Asian individuals enjoying food, with text saying 'no little asian don't eat them raw.' The text uses a derogatory and infantilizing term 'little asian' and implies a stereotype about eating habits, which is offensive and dehumanizing. The juxtaposition of the joyful scene with the derogatory text creates a contrast that highlights the stereotype. This combination of visual and textual elements perpetuates harmful racial stereotypes, leading to its classification as hateful content.

Predicted: The image shows a group of people eating, with text saying 'no little asian don't eat them raw.' The text uses a derogatory stereotype about Asian people, implying they eat raw meat, which is offensive and perpetuates harmful racial stereotypes. The juxtaposition of the text with the image of people enjoying food suggests a dehumanizing and racist message. This combination of visual and textual elements likely led to the classification of the image as hateful, as it targets a specific racial group with a demeaning stereotype.

Correct prediction

(a) Correct prediction.



Gold: not-hateful,
Predicted: hateful
BERT-F1: 0.6259

Explanation

Gold: The image consists of two parts: the top shows a man in a suit with text stating he was born in Sweden and is Swedish. The bottom shows a goldfish with text saying it was born in an aquarium and is a 'fish ken.' The juxtaposition humorously contrasts nationality with species, using wordplay on 'Swedish' and 'fish ken.' The human expert likely found it non-hateful as it plays on linguistic humor without targeting or demeaning any group, focusing instead on a playful comparison between human and fish identities.

Predicted: The image juxtaposes two individuals, one labeled 'Hasan' and another labeled 'Misho,' with a fish. The text implies a comparison between Hasan, who is described as Swedish, and Misho, a fish, suggesting a derogatory equivalence. This comparison is likely intended to dehumanize Hasan by equating him to an animal, which is a common tactic in hate speech. The use of nationality and personal details in the text further targets Hasan based on his background, reinforcing xenophobic or racist sentiments. The combination of these elements likely led to the classification of the image as hateful.

Incorrect prediction

(b) Incorrect prediction.

Figure 3: Example of correct and incorrect label prediction with explanation.