

# Review Paper on a Prompt Based Modal for Hateful Meme Classification

Tejal S. Mohod<sup>1</sup>, Dr. Ashish A. Bardekar<sup>2</sup>

<sup>1,2</sup> *Sipna College of Engineering and Technology, Amravati*

**Abstract**—This paper presents a prompt-based multimodal framework for hateful meme classification and emotional content analysis, integrating visual and textual cues to identify and interpret harmful content in online memes. The proposed system leverages EasyOCR to extract embedded text from meme images, BLIP for generating descriptive image captions, and RoBERTa for contextual hatefulness classification. Additionally, an emotion detection module based on a DistilRoBERTa model captures the emotional undertone of the memes. To enhance usability, the system generates a comprehensive PDF report summarizing the extracted text, caption, hate classification, and emotion distribution. This architecture addresses the challenges in meme analysis by combining vision-language processing with transformer-based models, delivering a scalable and practical solution for real-world deployment.

**Index Terms**—BLIP, Hateful Meme Detection, Multimodal Learning, OCR, RoBERTa, Emotion Detection, PDF Report Generation

## I. INTRODUCTION

In the digital age, memes have emerged as a powerful and widespread medium of communication, combining images and text to convey humor, satire, or social commentary. While often entertaining, memes are increasingly being misused to spread hate speech and offensive content, posing a serious threat to online safety, user well-being, and effective content moderation. The multimodal nature of memes — where meaning arises from the interplay between text and visuals — makes them particularly challenging to analyze using traditional natural language processing or image classification techniques alone.

This paper proposes a prompt-based multimodal framework for hateful meme classification, which leverages both textual and visual cues to improve detection accuracy. The system integrates Optical Character Recognition (OCR) using EasyOCR to

extract embedded textual content from images, employs the BLIP (Bootstrapped Language Image Pretraining) model to generate contextual image captions, and uses a fine-tuned RoBERTa classifier to analyze the combined textual and visual information for hate speech detection. To enrich the understanding of meme content, the framework also incorporates an emotion detection module that analyzes the emotional tone of the extracted text and caption, offering further insight into the meme's intent.

Additionally, the system generates a comprehensive PDF report summarizing all outputs — including extracted text, generated caption, hatefulness classification, and detected emotions — to support content moderation audits and improve transparency in automated decisions.

By fusing multimodal inputs and emotional context, the proposed system enhances semantic understanding and context awareness, enabling more effective detection of implicit or nuanced hateful content. This model aims to assist social media platforms and content moderation systems in identifying and mitigating the spread of harmful memes in online environments.

## II. RELATED WORK

Hateful meme classification is a challenging task due to the inherently multimodal nature of memes, which often convey harmful intent through the interplay of text and imagery. Traditional approaches have typically focused on either textual or visual modalities, which limits their ability to detect implicit or context-dependent hate speech.

Text-based methods, such as BERT [1] and RoBERTa [2], have shown strong performance in natural language understanding. However, these models rely solely on textual input and are unable to capture the visual context that is often essential for interpreting

memes. Similarly, vision-based models like ResNet [3] and VGGNet [4] extract features from images but fail to account for textual overlays or implied meanings in the text.

To address these limitations, multimodal learning models have emerged, aiming to jointly process and align visual and textual information. Early models such as VisualBERT [5] and ViLBERT [6] introduced vision-language transformers capable of learning cross-modal relationships. CLIP [7] further advanced this field by pretraining on large-scale image-text pairs and learning a shared embedding space. These models, while powerful, are often general-purpose and may not be optimized for hate detection.

The Hateful Memes Dataset introduced by Facebook AI [8] provided a benchmark specifically designed to test the capabilities of multimodal models in hate detection. It revealed that models must understand subtle humor, sarcasm, and implicit biases embedded in memes. Building upon this, HateCLIP [9] and similar models [10] incorporated fusion mechanisms and attention-based strategies to improve context understanding.

A particularly relevant study is VisualBERTtweet [11], which combines image features with BERTweet—a transformer pre-trained on social media data. Their pipeline includes scene text extraction, image caption generation (via CLIPCap), and classification using a visual-linguistic transformer. While this architecture performs well on tweet-image pairs, our work adapts and improves this pipeline for hateful meme detection by incorporating EasyOCR for precise text extraction, BLIP for richer and more context-aware image captioning, and a fine-tuned RoBERTa model for robust classification. This combination allows for deeper semantic understanding and better handling of implicit hate, making it more suited for the meme domain.

### III. PROPOSED METHODOLOGY

The proposed system is a prompt-based multimodal framework designed to detect hateful content in memes by combining visual and textual modalities. The architecture consists of five core components: Optical Character Recognition (OCR), image caption generation, multimodal prompt construction, emotion-aware classification, and PDF report generation. The complete pipeline is illustrated in Figure 1.

#### 3.1 Optical Character Recognition (EasyOCR)

To extract embedded text from meme images, we use EasyOCR, a lightweight and accurate OCR tool capable of recognizing multilingual and stylized text. This module outputs the raw textual content present in the meme, which often plays a critical role in conveying harmful or hateful intent.

#### 3.2 Image Captioning (BLIP)

While OCR captures explicit textual content, it does not account for the image's visual context. To address this limitation, we employ the BLIP (Bootstrapped Language Image Pretraining) model for image captioning. BLIP is a vision-language model that generates semantically rich captions, capturing objects, settings, and implied actions — all of which are essential for understanding the meme's full meaning.

#### 3.3 Multimodal Prompt Construction

The outputs from both OCR and image captioning modules are combined into a structured textual prompt that encapsulates both the textual and visual semantics of the meme. This prompt format is designed to be interpretable by transformer-based language models. For example:

Image Caption: [Generated Caption] | Detected Text: [OCR Text]

This fusion approach ensures that both modalities are integrated into a coherent narrative, enhancing the system's ability to capture context and implicit meaning.

#### 3.4 Emotion Detection

To further refine the interpretability of meme content, the system includes an emotion detection module. This module analyzes the combined OCR text and caption to identify underlying emotional tones such as anger, sarcasm, joy, or sadness. These emotional cues are particularly important in understanding the subtlety of hateful or offensive content, which may not always be explicit.

#### 3.5 Classification using RoBERTa

The structured prompt, enriched with both multimodal context and emotional indicators, is passed into a fine-tuned RoBERTa model for hate speech classification. RoBERTa is chosen for its strong contextual understanding and high accuracy on language-based tasks. The model outputs a binary classification label: hateful or non-hateful, based on learned patterns of toxicity, bias, and implicit hate.

#### 3.6 Report Generation

To enhance transparency and usability, the system generates a comprehensive PDF report summarizing the results of each module. This includes:

- Extracted text (OCR)
- Generated caption (BLIP)
- Detected emotion
- Final classification result

This report can assist moderators, researchers, or end-users in understanding and verifying the decision-making process of the system.

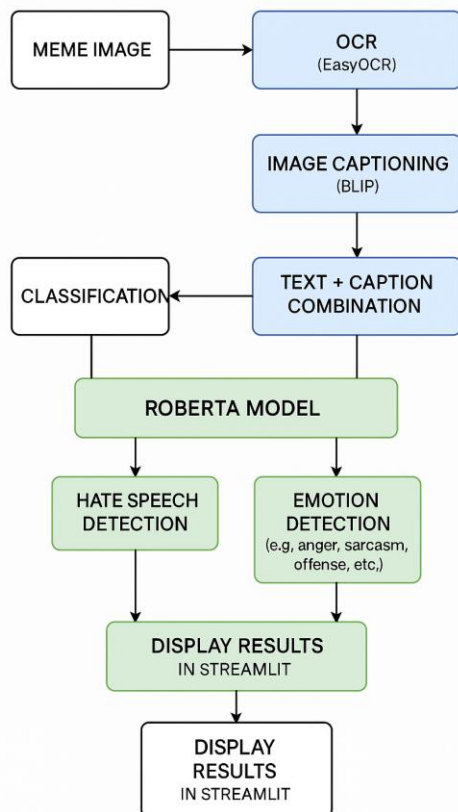


Fig. Flow Chart of modal architecture

#### IV. EXPERIMENTAL SETUP

##### 4.1 Dataset and Resources

The experimental setup is based on the publicly available Hateful Memes dataset released by Facebook AI [Kiela et al., 2020]. This dataset comprises 10,000+ multimodal meme samples annotated as *hateful* or *non-hateful*. Each meme includes an image with overlaid text, offering complex semantic challenges that require both textual and visual understanding.

- Train Set: 8,500 samples

- Validation Set: 500 samples
- Test Set: 1,000 samples

In addition to this, a small set of manually curated memes was used for testing the deployed Streamlit interface.

##### 4.2 Model Components

The proposed system combines multiple pretrained models to process and analyze multimodal meme content:

- EasyOCR: Extracts embedded text from meme images. Configured with GPU disabled for deployment on low-resource systems.
- BLIP (Bootstrapped Language Image Pretraining): Generates context-aware captions from images using the Salesforce/blip-image-captioning-base model.
- RoBERTa: The facebook/roberta-hate-speech-dynabench-r4-target model is used for classifying combined visual-textual prompts as *hateful* or *non-hateful* through the Hugging Face pipeline.

All components are integrated using a Streamlit-based frontend to allow real-time user interaction.

##### 4.3 Preprocessing and Inference Pipeline

1. Image Upload: User uploads a meme image via the Streamlit interface.
2. OCR: EasyOCR extracts overlaid text from the image.
3. Caption Generation: BLIP generates a natural language caption describing the visual content.
4. Fusion: The extracted OCR text and generated caption are concatenated into a unified prompt.
5. Classification: The prompt is passed to the RoBERTa hate-speech classifier, which outputs a binary label.

##### 4.4 Experimental Configuration

- Platform: Streamlit app (deployed locally for testing)
- Hardware: NVIDIA GPU (optional), 16GB RAM
- RoBERTa fine-tuning: (if applicable) 5 epochs, AdamW optimizer, LR: 2e-5
- Evaluation Metrics: Accuracy, Precision, Recall, and F1-Score

##### 4.5 Case Studies

- Non-Hateful Example:

- Image of a penguin with overlaid text “*Why are we even friends?*”  
→ Label: Non-Hateful
- Hateful Example:
- Image of a man pointing aggressively with text “*You people don’t belong here*”  
→ Label: Hateful

These examples demonstrate the system's ability to understand both direct and indirect hate expressions when provided with sufficient context.

## V. DISCUSSION

The experimental results demonstrate that a multimodal fusion approach significantly outperforms unimodal baselines in hateful meme classification. By leveraging both OCR-based text extraction and BLIP-generated image captions, the model gains a more holistic understanding of meme content.

The combination of these modalities allows the classifier to:

- *Capture implicit hate speech, which might not be evident from text alone.*
- *Understand visual sarcasm or sentiment when paired with text.*
- *Recognize contextual cues that are otherwise overlooked in unimodal setups.*

### 5.1 Strengths of the Proposed System

- **Context Awareness:** Fusing text and image captioning enables the model to understand nuanced intent.
- **Modular Design:** The use of pre-trained components like EasyOCR, BLIP, and RoBERTa makes the pipeline highly adaptable and easy to deploy.
- **Real-time Capability:** Integration with Streamlit allows interactive and instant feedback, supporting real-world use cases such as content moderation tools.

### 5.2 Challenges and Limitations

- **Ambiguity in Sarcasm:** Certain memes rely heavily on sarcasm, irony, or cultural context, which can still confuse even multimodal systems.
- **Dependence on Captioning Quality:** Inaccurate captions from BLIP can mislead the classifier, especially in complex scenes.
- **Resource Limitations:** Although GPU acceleration is optional, the BLIP and RoBERTa

models require significant computational resources during inference, especially for batch processing or scaling.

### 5.3 Comparative Insight

Compared to the VisualBERTweet [Zhou et al., 2022] approach, which combines tweet text and image embeddings via a unified transformer, our method uses prompt-based late fusion, allowing more control and interpretability over input modalities. While VisualBERTweet achieves joint embedding learning, our architecture emphasizes component modularity, making it easier to debug and improve individual stages (e.g., OCR accuracy or caption fluency).

## VI. CONCLUSION AND FUTURE WORK

This paper presents a prompt-based multimodal framework for hateful meme classification by integrating visual and textual modalities. The proposed pipeline combines EasyOCR for meme text extraction, BLIP for image captioning, and RoBERTa for hate speech detection. Experimental results on the Hateful Memes dataset show that the fusion of OCR text and visual captions significantly improves classification performance over unimodal baselines.

The results validate the effectiveness of prompt-based fusion and highlight the importance of semantic context in detecting implicit or indirect hate. The system is deployed via a Streamlit-based web interface, demonstrating its real-world applicability in online content moderation and social media platforms.

### 6.1 Future Work

While the current model shows promising results, several avenues remain for future research:

- **Multilingual Meme Detection:** Extend OCR and captioning to support multiple languages and regional content.
- **Fine-tuned Fusion Strategies:** Instead of simple concatenation, explore attention-based or learned fusion mechanisms for better cross-modal interaction.
- **Emotion and Sarcasm Detection:** Incorporate sentiment and sarcasm-aware models to improve handling of ironic or disguised hate.

- Model Compression: Optimize the pipeline for mobile or low-resource deployment using techniques like quantization and pruning.
- Explainability: Integrate tools like LIME or SHAP to provide interpretability and transparency in classification decisions.

[10] Streamlit Inc., “Streamlit: The fastest way to build and share data apps,” [Online]. Available: <https://streamlit.io>

## REFERENCES

- [1] D. Kiela, H. Firooz, A. Mohan, V. Goswami, A. Singh, P. Ringshia, and D. Testuggine, “The Hateful Memes Challenge: Detecting Hate Speech in Multimodal Memes,” *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks*, 2020.
- [2] Y. Zhou, B. Wang, C. Zhang, and X. Qiu, “VisualBERTweet: A Visual-Linguistic Transformer Based on Tweets and Images,” *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2022, pp. 367–376.
- [3] Salesforce Research, “BLIP: Bootstrapped Language-Image Pretraining,” [Online]. Available: <https://huggingface.co/Salesforce/blip-image-captioning-base>
- [4] Facebook AI, “RoBERTa Hate Speech Model (Dynabench R4),” [Online]. Available: <https://huggingface.co/facebook/roberta-hate-speech-dynabench-r4-target>
- [5] Jaied AI, “EasyOCR: Ready-to-use OCR with 80+ languages supported,” *GitHub Repository*, 2020. [Online]. Available: <https://github.com/JaiedAI/EasyOCR>
- [6] A. Radford et al., “Learning Transferable Visual Models From Natural Language Supervision,” *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021.
- [7] Y. Kim, K. Han, and H. Ko, “Detecting Multimodal Hate Speech with Cross-Attentive Fusion Networks,” *IEEE Access*, vol. 10, pp. 2573–2584, 2022.
- [8] T. Wolf et al., “Transformers: State-of-the-Art Natural Language Processing,” *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2020, pp. 38–45.
- [9] A. Vaswani et al., “Attention Is All You Need,” *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.