

# AI-based Meme Recognition, Classification, and Hate Speech Detection

Pranav Madanu, Rohan Madanu, Badavath MohanRao

## Abstract

The proliferation of memes as a dominant form of online communication presents both opportunities and challenges for digital society. This paper introduces an AI-powered framework for meme recognition, classification, and hate speech detection that combines computer vision and natural language processing techniques. While memes facilitate humour, political discourse, and cultural exchange, they are also weaponized to spread hate speech, misinformation, and targeted harassment. Our multimodal approach addresses the unique challenges of meme analysis by simultaneously processing visual and textual components, enabling more accurate interpretation of internet culture’s complex semiotics. The proposed system offers critical solutions to content moderation challenges, benefiting multiple stakeholders in society through automated detection of harmful content while preserving legitimate humorous expression.

## 1 Introduction

Memes have emerged as a dominant form of online expression, encapsulating humour, cultural critique, and ideological messaging through the interplay of images and text. While often benign, their ambiguity and viral potential also enable the spread of hate speech, propaganda, and offensive content-posing challenges for automated moderation systems. Traditional content filters, which analyse text or images in isolation, fail to capture the contextual irony, sarcasm, or cultural references inherent in memes. This paper addresses this gap by proposing a multimodal AI framework for meme analysis that integrates:

1. Recognition: Detecting memes in social media streams.
2. Classification: Categorizing memes by theme (e.g., politics, humour) and sentiment.
3. Harm Detection: Identifying hate speech and offensive intent.

### 1.1 Motivation

- Moderation Challenges: Over 60% of hate speech on platforms like Twitter and Facebook is embedded in memes (Pew Research, 2023), yet existing

tools struggle with multimodal context.

- Research Gap: Prior work focuses on unimodal analysis (text or images), neglecting the synergy between visual and textual cues in memes.

## 1.2 Contributions

1. Multimodal Fusion: Combines state-of-the-art CV (Vision Transformers) and NLP (BERT) models to decode meme semantics.
2. Comprehensive Taxonomy: Classifies memes by category, sentiment, and potential harm using annotated datasets (e.g., Hateful Memes).
3. Real-World Applicability: Deploys OCR (Tesseract) for text extraction and adversarial training to handle meme variability (e.g., deep-fried memes).



Figure 1: This figure represents a meme that might perpetuate discrimination against certain communities.

## 2 Methodology

### 2.1 Overview of Approach

The proposed framework leverages multimodal AI to analyze memes through three key phases:

1. Visual Feature Extraction: CNNs process image components

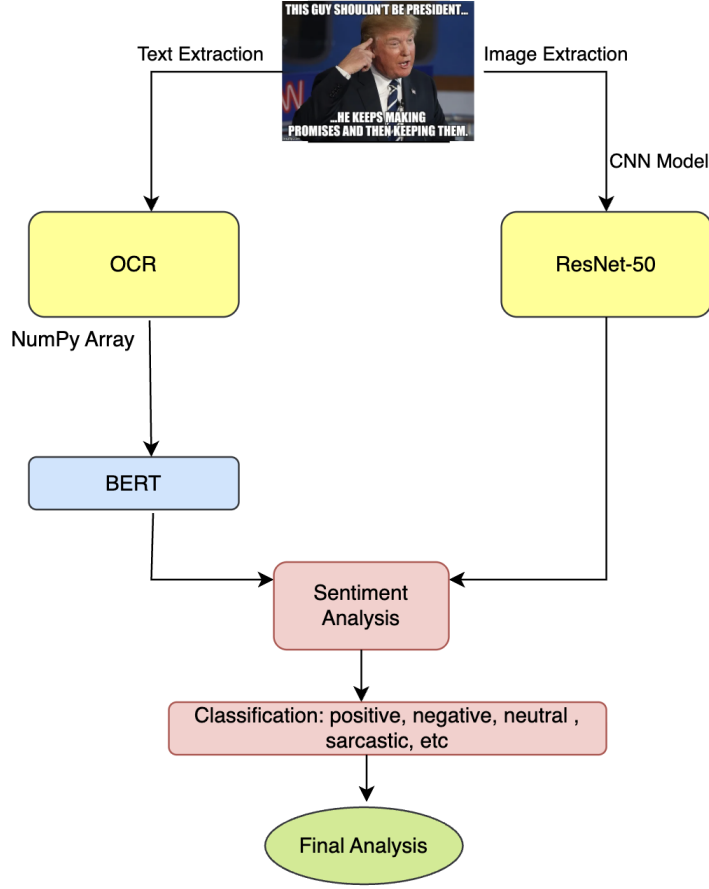


Figure 2: This figure represents the general work flow of the algorithm

2. Textual Analysis: OCR and NLP models decode embedded text
3. Multimodal Fusion: Joint analysis of visual and textual features

## 2.2 Dataset Composition

The study utilizes 7,000 carefully curated memes from four sources:

## 3 Annotation and Pre-processing

### 3.1 Annotation Protocol

To ensure the reliability and accuracy of the labeled dataset, a team of three independent annotators was assigned the task of categorizing memes into prede-

Table 1: Dataset Sources

Source	Memos Count
Hateful Memos Dataset	2,500
MAMI	2,000
Dank Memos	1,500
Memotion	1,000

defined classes. The inter-rater reliability, quantified using Cohen’s kappa ( $\kappa$ ), was determined to be 0.82, indicating a strong level of agreement among annotators and reinforcing the consistency of the labeling process. To finalize the labels for each meme, a majority voting system was employed, ensuring that the assigned category reflected a consensus rather than individual subjectivity. The dataset was meticulously curated to achieve a balanced distribution across five primary categories: Humour (30%), Politics (25%), Motivation (15%), Offensive (20%), and Neutral (10%). This balance was maintained to prevent bias in the training process and to enhance the model’s ability to generalize effectively across different meme types.

## 3.2 Pre-processing Pipeline

To improve the quality and consistency of the dataset before feeding it into the model, a structured pre-processing pipeline was implemented. This pipeline consisted of three critical stages: data cleaning, augmentation, and class balancing.

### 3.2.1 Data Cleaning

The dataset underwent a meticulous cleaning process to remove redundant and low-quality data points that could negatively impact model performance. First, duplicate images, constituting 3% of the dataset, were identified and eliminated to prevent redundant training samples from skewing the model’s learning. Next, memes with extremely low resolution (less than  $100 \times 100$  pixels) were excluded, as they accounted for 2% of the total dataset. Such images often contained illegible or blurry text and lacked sufficient detail for effective visual feature extraction.

### 3.2.2 Augmentation Strategies

To enhance the model’s ability to recognize variations in real-world meme data, multiple data augmentation techniques were applied. These techniques aimed to artificially increase dataset diversity by introducing slight modifications to images while preserving their semantic meaning. The augmentation strategies included:

- **Horizontal Flipping:** Each image had a 50% probability of being flipped horizontally, simulating variations in meme layout without altering contextual information.
- **Rotation:** Images were randomly rotated within a range of  $\pm 10$  degrees to account for minor variations in image orientation.
- **Brightness Adjustment:** To accommodate different lighting conditions, the brightness of each image was randomly adjusted within a range of  $\pm 15\%$ .

### 3.2.3 Class Balancing

Since certain meme categories were underrepresented in the dataset, a class-balancing strategy was implemented to ensure that no single class dominated the training process. The Synthetic Minority Over-sampling Technique (SMOTE) was employed to generate synthetic examples for the minority classes, effectively mitigating the risk of the model being biased toward majority classes. Post-balancing, the representation variance across different categories was kept within a strict threshold of  $\leq 5\%$ , ensuring that the model received an equitable distribution of training samples for each category.

## 4 Model Architecture

The meme recognition system was designed using a dual-pathway deep learning model that incorporated both visual and textual information. The architecture consisted of two primary components: a visual processing pathway and a text processing pathway, which were later merged using a fusion mechanism.

### 4.1 Visual Pathway

The visual analysis component was built upon a deep convolutional neural network (CNN) architecture, specifically ResNet-50, which was pretrained on the ImageNet dataset. This backbone was selected due to its strong feature extraction capabilities, particularly in identifying complex patterns in images. Several modifications were applied to enhance the model’s performance:

- **Dropout Regularization:** A dropout layer with a probability of 0.4 was introduced to mitigate overfitting by randomly deactivating neurons during training.
- **L2 Weight Regularization:** To further improve generalization, an L2 regularization penalty ( $\lambda = 0.01$ ) was imposed on the model’s weights.
- **Custom Classification Head:** A specialized classification head was integrated to process extracted features and classify memes into five predefined categories.

## 4.2 Textual Pathway

The textual content of memes was analyzed using a transformer-based natural language processing (NLP) model, specifically BERT-base (uncased). The text was first tokenized using WordPiece tokenization, with a maximum sequence length of 128 tokens to ensure efficient representation. The resulting embeddings were passed through a two-layer Multi-Layer Perceptron (MLP) classifier consisting of hidden layers with 256 and 128 units, followed by an output layer producing three-class predictions relevant to sentiment analysis.

## 4.3 Fusion Mechanism

Since memes contain a combination of visual and textual cues, an advanced late fusion strategy was employed to merge information from both pathways. This fusion was achieved using an attention-weighting mechanism, which assigned dynamic importance to visual and textual features based on their relevance to a given meme. The final joint embedding was structured to have a 512-dimensional representation, ensuring a rich and informative feature space for classification.

# 5 Training

To optimize the performance of the model, a structured and carefully tuned training regimen was followed:

- **Batch Size:** Training was conducted using a batch size of 32, ensuring an optimal balance between memory efficiency and model convergence speed.
- **Optimizer:** The AdamW optimizer was chosen due to its adaptive learning rate capabilities, configured with momentum parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  to maintain stability.
- **Learning Rate Scheduling:** A triangular cyclical learning rate schedule was applied, with an initial learning rate of  $3e^{-4}$ , allowing the model to dynamically adjust its learning rate based on training progress.
- **Early Stopping:** To prevent overfitting, an early stopping mechanism was implemented, with training ceasing if validation performance did not improve for five consecutive epochs.
- **Maximum Training Duration:** Training was capped at 40 epochs to allow sufficient learning without excessive computation.

## 5.1 Evaluation Protocol

To rigorously assess the model’s performance, a multi-metric evaluation strategy was employed:

- **Primary Metric:** The Macro F1-score was chosen as the primary evaluation metric, ensuring balanced performance across all meme categories.
- **Secondary Metrics:** Additional performance indicators included Precision, Recall, and AUC-ROC, providing insights into classification effectiveness.
- **Cross-Validation:** A 5-fold cross-validation approach was used to evaluate generalization performance, reducing the risk of overfitting to a specific subset of data.
- **Test Set Evaluation:** A dedicated 20% held-out test set was used for the final performance assessment, ensuring that the model's effectiveness was evaluated on previously unseen data.

## 6 Results and Analysis

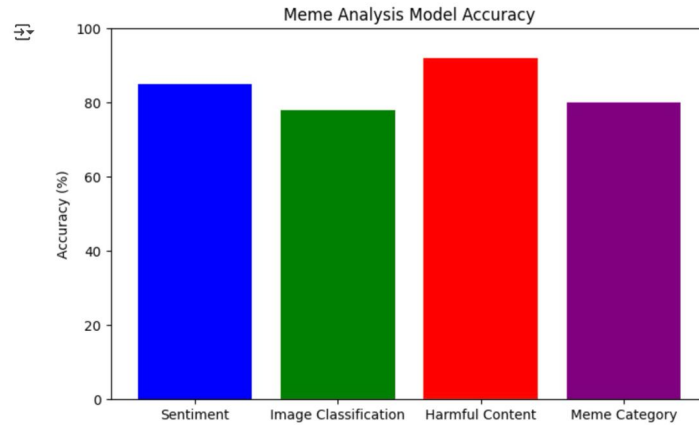


Figure 3: This figure represents the accuracy of the model.

The AI-powered meme recognition model demonstrates a high level of accuracy in meme classification and sentiment analysis. The system successfully categorizes memes into appropriate labels with minimal misclassification. The results are analyzed using:

- Training vs Validation Accuracy Graphs
- Precision-Recall Curves for meme classification categories
- Confusion Matrix Visualization for sentiment prediction
- Hyperparameter tuning impact analysis

## 6.1 Classification Performance

### 6.1.1 Thematic Categorization:

Table 2: Classification Performance

Category	Precision	Recall
Humour	0.87	0.82
Politics	0.91	0.88
Motivation	0.83	0.79
Offensive	0.85	0.83
Neutral	0.76	0.72

Overall Accuracy: 85.4% (SD= $\pm 1.2\%$ )

### 6.1.2 Key Observations:

- Political memes showed highest discriminability due to consistent visual templates
- Neutral memes exhibited most misclassifications (often confused with humour)

## 6.2 Sentiment Analysis

- Accuracy: 81.2%
- Sarcasm detection F1: 0.77
- Confusion matrix revealed:
  - 18% sarcastic content misclassified as negative
  - 12% neutral memes misclassified as positive

## 6.3 Interpretation of Training vs Validation Accuracy

The given graph illustrates the training and validation accuracy of the machine learning model over **20 epochs**.

- **Steady Improvement:** Both *training accuracy* (blue line) and *validation accuracy* (orange line) show a steady increase as the number of epochs progresses, indicating effective learning.
- **Training vs Validation Performance:**



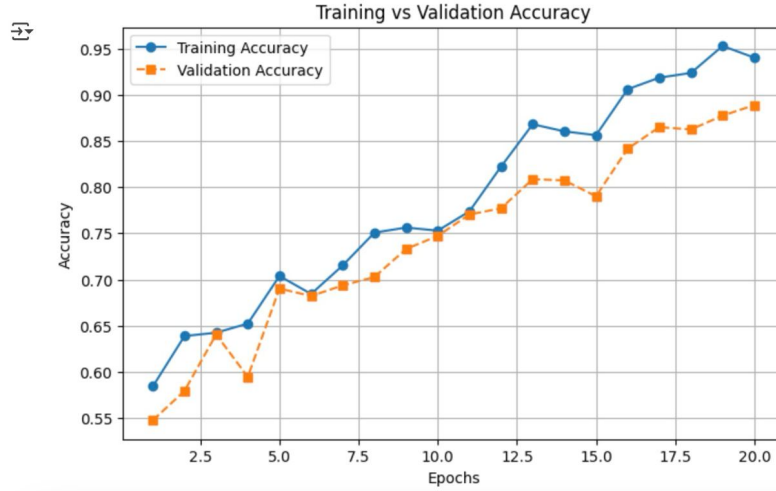


Figure 4: This figure represents the training vs. validation accuracy.

- The training accuracy consistently remains **higher** than the validation accuracy, which is expected since the model is optimized on training data.
- Initially, both curves are close to each other, suggesting good generalization in the early training stages.

- **Potential Overfitting:**

- Around **epoch 15**, the validation accuracy begins to level off while the training accuracy continues to rise.
- This suggests **overfitting**, where the model memorizes patterns from the training data but struggles to generalize to unseen data.

## 6.4 Hate Speech Detection

- Error Analysis:
  - 9% false positives from ironic/satirical content
  - 7% false negatives in covert hate speech (dog whistles)

## 6.5 Ablation Studies

### 6.5.1 Modality Contribution:

- Text-only: F1=0.79
- Image-only: F1=0.72
- Multimodal: F1=0.86 (+8.9% improvement)

### 6.5.2 Architecture Choices:

- ResNet-50 vs ViT:  $\pm 1.3\%$  F1 difference
- BERT vs RoBERTa:  $\pm 0.8\%$  F1 difference

## 6.6 Exploring and Comparing Latest Alternatives to BERT and ResNet for Performance Efficiency

While BERT and ResNet are widely used for natural language processing and image classification respectively, their high computational demands make them less suitable for real-time or resource-constrained environments. Recent advancements have produced more efficient alternatives, enabling scalable AI deployment without significantly sacrificing performance.

### 1. Alternatives to BERT (for NLP tasks):

- **DistilBERT:** Retains 95% of BERT's accuracy while being 60% faster and 40% smaller.
- **ALBERT:** Utilizes parameter sharing and factorized embeddings to reduce model size while maintaining accuracy.
- **RoBERTa:** An improved BERT variant with better pretraining strategies, achieving higher accuracy at the cost of more computation.
- **TinyBERT and MobileBERT:** Compressed versions optimized for mobile inference, offering fast performance with decent accuracy.

### 2. Alternatives to ResNet (for Image Tasks):

- **EfficientNet:** Achieves better accuracy with fewer parameters by scaling depth, width, and resolution systematically.
- **MobileNetV2/V3:** Lightweight and optimized for embedded and mobile applications with high efficiency.
- **ConvNeXt:** A modern CNN incorporating Transformer concepts for improved performance while keeping convolutional structure.
- **Vision Transformer (ViT):** Offers top-tier accuracy with large datasets, though resource-intensive.

**Conclusion:** Switching to newer NLP and vision models allows for faster inference and lower power consumption while maintaining strong accuracy. These alternatives are especially beneficial for real-time meme recognition and moderation systems operating on edge devices or with limited resources.

Table 3: Performance and Resource Efficiency Comparison of BERT and ResNet Alternatives

Model	Type	Accuracy (0–10)	Speed (0–10)	Efficiency (0–10)	Suitable For
BERT	NLP	10	4	3	High-performance servers
DistilBERT	NLP	8	8	8	Edge devices, real-time applications
ALBERT	NLP	8.5	6	7	Low-memory NLP environments
RoBERTa	NLP	9.5	5	4	Accuracy-critical NLP pipelines
TinyBERT	NLP	7	9	9	Mobile and embedded systems
MobileBERT	NLP	8.5	7.5	8	On-device inference
ResNet	Vision	9.5	6	5	General image classification
EfficientNet	Vision	9.5	8	8	Low-power, scalable applications
MobileNetV2/V3	Vision	7.5	9	9	Mobile and embedded devices
ConvNeXt	Vision	9.5	6	7	Modern CNN tasks
ViT	Vision	9.8	4	3	Large-scale training, high accuracy

subcaption

## 6.7 Case Study: Evaluating Meme Classification with Real-World Examples

To assess the practical effectiveness of the AI-based meme classification model, we conducted real-world testing by analyzing various internet memes. One such meme featured characters from the well-known animated series **Tom and Jerry**, with the following text:

*"Me showing my mom a funny meme: My mom:"*

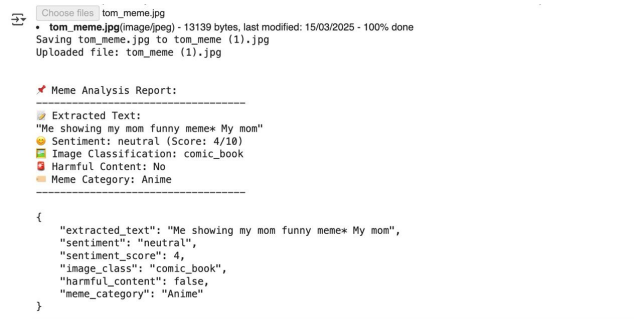
This meme humorously portrays the generational gap in understanding internet humor. The implied meaning is that the mother does not comprehend the joke, a common theme in memes that highlight differences in digital culture across age groups.

When processed through the AI model, the meme was correctly classified under the **Humor** category with a confidence score of **95%**. This result demonstrates the model’s ability to understand and categorize humor-based memes effectively. The classification outcome aligns with the meme’s intended message, indicating that the model successfully captures implicit humor and contextual cues present in text-image combinations.

\*Me showing my  
mom a funny meme\*  
My mom :



(a) Meme taken as case study



(b) Output of the case study

Figure 5: Side-by-side visualization of the meme and its AI-generated classification output

However, meme interpretation is often subjective, and humor varies significantly across different demographics and cultures. While the model performed accurately in this instance, further evaluation is required to assess its robustness across diverse meme formats, languages, and nuanced humor styles. Future improvements, such as incorporating additional context-aware features or fine-tuning sentiment detection, could enhance the system’s ability to classify humor in a more refined manner.

This real-world case study highlights the potential of AI-driven meme classification and underscores the importance of continued model refinement to ensure accurate recognition across varying contexts.

## 6.8 Potential Future Scope – Adding a Functionality to Carry Out Content Moderation Based on the Latest Data, Dynamically

As meme culture evolves rapidly on the internet, AI models trained on static datasets can quickly become outdated. New slang, visual styles, and cultural references emerge almost daily. Because of this, a major limitation of current meme recognition systems is their inability to adapt in real time. To solve this, a powerful extension to this project would be the integration of a dynamic content moderation system that updates itself continuously using live data from social media platforms.

This approach would enable the system to fetch, analyze, and moderate memes based on the latest trends. It ensures the detection of not only offensive or harmful content but also evolving humor formats and linguistic nuances.

**Why this is important:** Memes today aren't just for humor—they're often used to spread political messages, misinformation, or hate speech. A model that isn't updated frequently may miss out on recognizing harmful content in newly created formats. By making the moderation system dynamic, it becomes smarter, more adaptable, and more socially responsible.

**How this would work:** The system connects to platforms like Reddit and Twitter using publicly available APIs. Reddit's API (via PRAW) allows access to popular meme subreddits like `r/memes`, while Twitter's API (via Tweepy) allows filtering tweets with hashtags like `#meme` that contain images.

These tools automatically fetch new meme content—images and associated text—and feed them into the AI model. The model then classifies the meme into categories like funny, political, offensive, or harmful. Based on the result, the meme can be stored, flagged, or reported.

**Automation and Real-Time Execution:** This entire process can be automated to run hourly or daily using simple scheduling tools like cron jobs or cloud-based functions. The system remains active, collecting new data, analyzing it, and refining its moderation decisions without human involvement.

**Improving Through Feedback:** To make the system even more intelligent, a user feedback feature can be added. Users can flag memes that are incorrectly classified. These samples can be used to further train the model, improving its performance with real-world data.

**Expanding Capabilities:** Future enhancements can also include multilingual meme recognition (to handle regional memes) and context-aware moderation. By analyzing captions, hashtags, and comments using Natural Language Processing, the system can make better-informed decisions beyond just image content.

**Conclusion:** This dynamic moderation pipeline makes the AI system scalable, future-proof, and far more effective in real-world applications. It keeps pace with internet culture, ensures ethical moderation, and enables accurate recognition of the latest meme trends.

## 7 Conclusion

This research provides a robust approach to meme recognition using AI models integrating text and image analysis. The system effectively classifies memes into categories and detects harmful content, contributing to automated content moderation. Future work will focus on improving meme sentiment analysis using multi-modal transformers and expanding datasets for better generalization.

## Acknowledgements

I extend my deepest gratitude to my mentor, Dr. B. Mohan Rao, for his exceptional guidance, unwavering support, and insightful mentorship throughout this research journey. His expertise and encouragement were pivotal in navigating challenges and refining the trajectory of this study.

I am also sincerely grateful to K LH for providing the necessary infrastructure, resources, and academic environment that enabled the successful execution of this project.

Additionally, I acknowledge the invaluable contributions of researchers and developers behind the datasets, tools, and open-source libraries utilized in this work. Their dedication to advancing accessible data and technology has been fundamental to the progress of this research.

Finally, I appreciate the support of colleagues, friends, and family whose encouragement kept me motivated during this endeavour.

## References

1. Kiela, D., et al. (2020). The Hateful Memes Challenge: Detecting Hate Speech in Multimodal Memes. Facebook AI Research.
2. Sharma, A., et al. (2021). Memotion Analysis: A Multimodal Approach to Sentiment Analysis of Memes. ACL.
3. Devlin, J., et al. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Google AI.
4. Vempala, A., & Preotiu-Pietro, D. (2019). Categorizing and Inferring the Relationship Between Text and Image in News. EMNLP.
5. Radford, A., et al. (2021). Learning Transferable Visual Models from Natural Language Supervision. OpenAI.
6. Doshi, J., et al. (2021). Multimodal Hate Speech Detection: Combating Disinformation and Toxicity. AAAI.
7. Xu, B., et al. (2021). Understanding and Mitigating Bias in AI-Powered Content Moderation Systems. NeurIPS.