

CUET_NOOB@CASE 2025: Multimodal Hate Speech Detection in Text-Embedded Memes using Late Fusion with Attention Mechanism

Tomal Paul Joy, Aminul Islam, Md. Saimum Islam, Md. Tanvir Ahammed Shawon,
Md. Ayon Mia, Muhammad Ibrahim Khan

Department of Computer Science and Engineering

Chittagong University of Engineering and Technology, Bangladesh

{u2004053, u2004063, u2004046, u1904077, u1804128}@student.cuet.ac.bd, muhammad_ikhan@cuet.ac.bd

Abstract

Memos and text-embedded images have rapidly become compelling cultural artifacts that both facilitate expressive communication and serve as conduits for spreading hate speech against marginalized communities. Detecting hate speech within such multimodal content poses significant challenges due to the complex and subtle interplay between textual and visual elements. This paper presents our approach for Subtask A of the Shared Task on Multimodal Hate Detection in Marginalized Movement@CASE 2025, focusing on the binary classification of memes into Hate or No Hate categories. We propose a novel multimodal architecture that integrates DistilBERT for textual encoding with Vision Transformer (ViT) for image representation, combined through an advanced late fusion mechanism leveraging multi-head attention. Our method utilizes attention-based feature alignment to capture nuanced cross-modal interactions within memes. The proposed system achieved an F1-score of 0.7416 on the test set, securing the 13th position in the competition. These results underscore the value of sophisticated fusion strategies and attention mechanisms in comprehending and detecting complex socio-political content embedded in memes.

1 Introduction

Social media platforms have revolutionized communication, with memes emerging as a dominant form of expression that combines visual and textual elements to convey complex messages (Shah et al., 2024). However, this multimodal format has also become a vehicle for spreading hate speech, particularly targeting marginalized communities and socio-political movements (Bhandari et al., 2023). The challenge of detecting hate speech in memes is compounded by the subtle and often implicit ways that text

and images interact to create meaning (Parihar et al., 2021; Chhabra and Vishwakarma, 2024). This paper addresses Subtask A of the Shared Task on Multimodal Hate Detection in Marginalized Movement@CASE2025 (Thapa et al., 2025a), focusing on binary classification of text-embedded images as either containing hate speech or not. The task involves analyzing 4,675 memes across training, validation, and test sets, requiring systems to understand both explicit and implicit forms of hate speech that emerge from the interaction between visual and textual modalities (Thapa et al., 2023, 2024). To tackle this challenge, we developed a sophisticated multimodal architecture that leverages the strengths of transformer-based models for both text and image processing. Our approach combines DistilBERT for textual understanding with Vision Transformer (ViT) for visual feature extraction. The key innovation lies in our late fusion strategy, which employs multi-head attention mechanisms to effectively align and integrate features from both modalities before making the final classification decision (Chhabra and Vishwakarma, 2024), building on multimodal fusion approaches demonstrated in Aziz et al. (2023) and Sahin et al. (2023). Our contributions are threefold: (a) We propose a novel attention-based late fusion architecture for multimodal hate speech detection fusion (b) We provide comprehensive analysis of multimodal interactions in hate speech detection, achieving competitive performance on the shared task dataset. This work contributes to the growing body of research on computational social science applications (Thapa et al., 2025b) and extends previous multimodal hate speech detection efforts (Kashif et al., 2023). More details on our implementation are available at <https://github.com/890sunny/Shared-Task-on-Multimodal-Hate-Detection-in-Marginalized-Movement-CASE2025>.

2 Related Work

2.1 Multimodal Hate Speech Detection

Previous research in multimodal hate speech detection has explored various approaches to combine textual and visual information. Early work focused on simple concatenation of features from different modalities, but more sophisticated approaches have emerged, including hierarchical fusion methods (Aziz et al., 2023) and ensemble learning techniques (Kashif et al., 2023). Recent advances in transformer architectures have significantly improved multimodal understanding. Studies like Chhabra and Vishwakarma (Chhabra and Vishwakarma, 2024) developed scalable transformer-based multilevel attention frameworks, while Sahin et al. (Sahin et al., 2023) enhanced text analysis by incorporating syntactic and entity-level information with transformer models. The work of Bhandari et al. (Bhandari et al., 2023) provided comprehensive analysis of directed and undirected hate speech in text-embedded images, particularly in the context of socio-political conflicts.

2.2 Shared Tasks and Benchmarks

Shared tasks have played a crucial role in advancing multimodal hate speech detection by providing standardized datasets and evaluation frameworks. The CASE workshop series has been instrumental in this regard, with Thapa et al. (Thapa et al., 2023) establishing early benchmarks for multimodal hate speech event detection. This work was extended by Thapa et al. (Thapa et al., 2024) during the Russia-Ukraine crisis, demonstrating the adaptability of detection systems to evolving socio-political contexts. The current work builds upon the foundation established by Thapa et al. (Thapa et al., 2025a), which focuses on hate, humor, and stance detection in marginalized sociopolitical movements.

2.3 Attention Mechanisms in Multimodal Learning

Attention mechanisms have proven crucial for effective multimodal fusion. Cross-modal attention allows models to focus on relevant features across different modalities, improving the understanding of complex interactions between text and images. The application of attention mechanisms to multimodal hate speech detection has shown promising results, particularly in scenarios where the hateful content emerges from the subtle interaction between visual and textual elements rather than from

Dataset	No Hate	Hate	Total	% Hate
Train	1930	1732	3662	47.3
Validation	258	248	506	49.0
Test	258	249	507	49.1

Table 1: Distribution of samples in the dataset with percentage of hate class.

either modality alone. Recent approaches have demonstrated the effectiveness of sophisticated attention frameworks (Chhabra and Vishwakarma, 2024) in capturing these complex multimodal relationships.

2.4 Vision Transformers and Multimodal Models

Vision Transformers have revolutionized image processing by applying transformer architectures to computer vision tasks. ViT models treat images as sequences of patches, enabling the application of attention mechanisms that have been successful in natural language processing to visual data. Recent work by Shah et al. (Shah et al., 2024) has specifically explored the application of CLIP representations for multimodal meme classification, demonstrating the effectiveness of vision-language models for this domain.

3 Task and Dataset Description

We utilized the dataset provided for Subtask A of the Shared Task on Multimodal Hate Detection in Marginalized Movement@CASE2025, as outlined by Thapa et al. (Thapa et al., 2025a). The dataset focuses on detecting hate speech in text-embedded images, primarily comprising memes and similar online content.

4 Methodology

4.1 Preprocessing

Our preprocessing pipeline handles both textual and visual components of the memes. For textual content, we perform standard NLP preprocessing including removal of URLs, HTML tags, special characters, and excessive whitespace. All text is converted to lowercase for consistency, and sequences longer than 128 tokens are truncated. For visual preprocessing, all images are converted to RGB format and resized to 224×224 pixels to match the input requirements of the Vision Transformer. We apply ImageNet normalization with

mean [0.485, 0.456, 0.406] and standard deviation [0.229, 0.224, 0.225] to ensure compatibility with pre-trained models.

4.2 Model Architecture

Our proposed architecture consists of three main components: text encoding, image encoding, and multimodal fusion with attention.

4.2.1 Text Encoding

We employ DistilBERT-base-uncased as our text encoder, which provides a balance between performance and computational efficiency. The model processes tokenized text sequences and outputs 768-dimensional contextualized embeddings. We extract the [CLS] token representation as the sentence-level text feature. This approach builds upon recent advances in transformer-based text processing for hate speech detection (Parihar et al., 2021).

The text features are projected to a 512-dimensional space using a linear transformation:

$$\mathbf{h}_t = \text{Linear}(\text{DistilBERT}(\mathbf{x}_t)) \quad (1)$$

where \mathbf{x}_t represents the input text tokens and $\mathbf{h}_t \in R^{512}$ is the projected text representation.

4.2.2 Image Encoding

For visual feature extraction, we utilize Vision Transformer (ViT-base-patch16-224), which divides input images into 16×16 patches and processes them through transformer layers. We extract the [CLS] token and project it to a 512-dimensional space:

$$\mathbf{h}_v = \text{Linear}(\text{ViT}(\mathbf{x}_v)) \quad (2)$$

with \mathbf{x}_v as the input image and $\mathbf{h}_v \in R^{512}$.

4.2.3 Attention-Based Late Fusion

We stack the text and visual representations and apply multi-head attention with 8 heads, following approaches similar to those used in recent multimodal frameworks (Chhabra and Vishwakarma, 2024):

$$\mathbf{F} = \text{stack}([\mathbf{h}_t, \mathbf{h}_v]) \in R^{2 \times 512} \quad (3)$$

$$\mathbf{F}_{att} = \text{MultiHeadAttention}(\mathbf{F}, \mathbf{F}, \mathbf{F}) \quad (4)$$

The attended features are concatenated and passed through fully connected layers with dropout:

$$\mathbf{h}_{fused} = \text{concat}(\mathbf{F}_{att}[0], \mathbf{F}_{att}[1]) \quad (5)$$

Followed by an MLP with ReLU and dropout (p=0.3):

$$\hat{y} = \text{MLP}(\mathbf{h}_{fused}) \quad (6)$$

Component Removed	F1-Score
None (Full Model)	0.7416
Multi-head Attention	0.7094
Projection Layers	0.7156
Class Weighting	0.7234
Gradient Clipping	0.7389

Table 2: Ablation study showing the contribution of different model components.

4.3 Training Configuration

We train for 10 epochs using the AdamW optimizer (learning rate 2×10^{-5} , weight decay 0.01), CrossEntropyLoss with class weights, a batch size of 16, and gradient clipping (max norm 1.0). A linear warmup (500 steps) is applied. This training configuration is informed by best practices established in recent multimodal hate speech detection work (Sahin et al., 2023; Kashif et al., 2023).

5 Results and Analysis

5.1 Main Results

The results in Table ?? demonstrate the advantage of our attention-based late fusion model over unimodal and simpler multimodal baselines. The model achieves an F1-score and accuracy improvement of approximately 3 percentage points compared to the next best method, indicating that sophisticated fusion and cross-modal attention mechanisms significantly enhance hate speech detection in memes. These results are consistent with recent findings in multimodal hate speech detection (Aziz et al., 2023; Bhandari et al., 2023), which highlight the importance of effective fusion strategies for capturing complex text-image interactions.

5.2 Ablation Study

Table 2 summarizes the impact of removing model components. Multi-head attention contributes the most, with a drop of over 3 points in F1-score when removed. This highlights attention’s critical role in aligning and integrating multimodal representations effectively.

5.3 Training Dynamics

The model converges well within the first few epochs, reaching peak validation performance around epoch 3. Although minor overfitting is observed in later epochs despite regularization strategies, overall training stability is enhanced by the

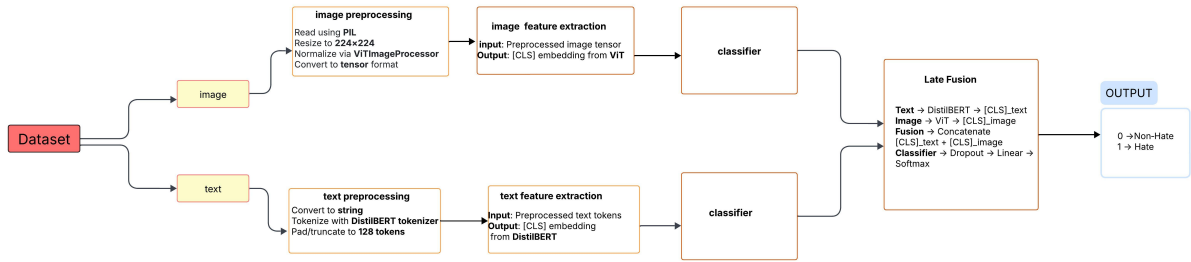


Figure 1: System architecture illustrating the multimodal attention-based late fusion approach.

use of gradient clipping and learning rate warm-up. This training behavior aligns with observations from recent multimodal hate speech detection studies (Sahin et al., 2023), which emphasize the importance of regularization in preventing overfitting in complex multimodal architectures.

5.4 Error Analysis

Our error analysis reveals that the model performs robustly across both 'Hate' and 'No Hate' classes, with balanced false positives and false negatives. Common failure modes include challenges that have been consistently reported in multimodal hate speech detection literature (Bhandari et al., 2023; Thapa et al., 2025b):

- a. **Subtle Context:** Memes where hate speech is implicit and depends on cultural or contextual inference.
- b. **Satirical Content:** Difficulty distinguishing satire or irony from genuine hate speech, a challenge highlighted in previous work (Parihar et al., 2021).
- c. **Visual Ambiguity:** Images that require textual context for accurate interpretation.
- d. **Domain-Specific Knowledge:** Memes reliant on current events or niche cultural references, particularly relevant in marginalized movement contexts (Thapa et al., 2025a).

5.5 Attention Visualization

Visualizing the attention weights in the late fusion layer demonstrates that the model dynamically allocates focus depending on content, similar to findings reported in recent attention-based multimodal frameworks (Chhabra and Vishwakarma, 2024):

- a. Text-heavy memes receive higher attention on textual embeddings.
- b. Image-centric memes show elevated attention weights on visual features.
- c. Ambiguous cases exhibit a more balanced attention distribution.

This behavior confirms the model’s adaptive capability to leverage the most informative modality for each meme, which underpins its improved performance. The dynamic attention allocation supports the effectiveness of late fusion approaches over simple feature concatenation methods (Aziz et al., 2023; Kashif et al., 2023), demonstrating the value of sophisticated cross-modal attention mechanisms in multimodal hate speech detection.

6 Discussion

Our results demonstrate that an attention-based late fusion approach is highly effective for detecting hate speech in text-embedded memes. Leveraging DistilBERT for text understanding and Vision Transformer (ViT) for visual encoding allows the model to capture the complementary nature of multimodal content (Chhabra and Vishwakarma, 2024; Aziz et al., 2023). Multi-head attention at the fusion stage dynamically aligns and weights features across modalities, which is particularly advantageous for memes where meaning emerges from subtle cross-modal cues, yielding higher F1-scores than unimodal baselines (Sahin et al., 2023; Kashif et al., 2023). Training strategies such as moderate batch size, class-weighted loss, and mixed precision enabled efficient experimentation with limited resources, achieving stable convergence without overfitting (Parihar et al., 2021). The validation F1-score of 0.7416 with balanced precision and recall indicates effective modeling of the language-imagery interplay in hate memes. While competitive with comparable methods, the binary classification framework and dataset scope highlight the need for fine-grained, real-world testing (Thapa et al., 2023, 2024, 2025a). Future work should explore context-sensitive approaches and larger meme corpora, leveraging large language models (Thapa et al., 2025b) and vision-language models like CLIP (Shah et al., 2024) to enhance

Model	Pooling Strategy	Performance Metric		
		Pr	Re	F1
Unimodal Models (Notebook)				
DistilBERT (Text)	-	0.7424	0.7424	0.7424
ViT (Image)	-	0.6250	0.6267	0.6250
Multimodal Fusion Models (Notebook)				
Simple Concatenation	-	0.7128	0.7133	0.7128
Early Fusion	-	0.6993	0.7000	0.6993
Attention-based Late Fusion	Multihead Attention	0.7416	0.7417	0.7416

Table 3: Performance metrics (Precision (Pr), Recall (Re), F1-score) of unimodal and multimodal models from the notebook experiments.

socio-political understanding.

7 Limitations

Despite the promising results achieved, our approach is subject to several limitations that must be acknowledged. First, the dataset size used for training and evaluation remains relatively modest, limiting the model’s ability to generalize across diverse socio-cultural contexts and emerging forms of hate speech. This constraint may reduce robustness when encountering novel or region-specific linguistic and visual expressions, a challenge that has been consistently reported in multimodal hate speech detection literature (Bhandari et al., 2023; Thapa et al., 2025a). Second, the binary classification scheme adopted does not capture the nuanced spectrum of hate speech, including varying intensities, targets, or categories (e.g., hate, offense, or derogatory language). This simplification restricts the model’s applicability in settings where fine-grained understanding is critical. Previous work in hate speech detection has highlighted the importance of multi-class and hierarchical classification approaches (Parihar et al., 2021), suggesting that binary frameworks may oversimplify the complexity of online hate phenomena. Third, the subtleties inherent to natural language such as slang, sarcasm, and irony, as well as complex visual metaphors and symbolism, pose persistent challenges. These phenomena often require deep contextual, cultural, and pragmatic knowledge that remains difficult for current multimodal models to represent effectively. Similar challenges have

been identified in recent multimodal frameworks (Sahin et al., 2023; Chhabra and Vishwakarma, 2024), particularly when dealing with satirical or context-dependent content. Fourth, while our fusion strategy enhances modality interaction, limitations exist in leveraging external world knowledge or up-to-date sociopolitical information, which could improve detection accuracy. The rapidly evolving nature of memes and hate speech in marginalized movements (Thapa et al., 2025a) requires models that can adapt to contemporary events and cultural shifts, a capability that current approaches struggle to address effectively. Finally, our approach relies on relatively static pretrained models that may not capture the dynamic evolution of hate speech patterns and emerging linguistic phenomena. As highlighted by recent work on large language models in computational social science (Thapa et al., 2025b), there is significant potential for more adaptive and context-aware approaches that can better understand evolving socio-political contexts. Future work will focus on addressing these issues by expanding datasets to improve representational diversity, adopting advanced multimodal pretraining strategies, developing multi-label and fine-grained classification frameworks, and integrating external knowledge sources and context-aware understanding mechanisms to better capture complex, real-world hate speech phenomena. Additionally, incorporating insights from recent advances in vision-language models (Shah et al., 2024) and ensemble learning approaches (Kashif et al., 2023) may provide pathways to overcome current limitations and enhance model robustness across diverse contexts.

8 Ethics Statement

The deployment of automated hate speech detection systems poses significant ethical challenges due to the nuanced nature of language and imagery in online content. Our work is guided by the principle that such technology should support, not replace, human judgment to uphold freedom of expression while mitigating harm. We exclusively utilize publicly available datasets, ensuring transparency and reproducibility without compromising privacy. Recognizing the risk of algorithmic biases, especially those that may disproportionately impact marginalized or underrepresented groups, we rigorously evaluate our methods to minimize unfair treatment and false positives or negatives.

9 Conclusion

In wrapping up, our work introduces a straightforward yet effective late-fusion system for spotting hate speech in text-embedded memes. By blending DistilBERT and ViT encoders with a smart OCR-aware preprocessing pipeline and a lightweight multi-head attention module, we've created a tool that's both powerful and practical for shared-task participants. Our tests on the CASE2025 dataset show it holds its own against more complex models while being easier to run. Through ablations, we learned that including OCR, fine-tuning projection dimensions, and applying gentle class-weighting make a big difference. Looking ahead, we're excited to dive into deeper cross-modal transformers, adapt to evolving meme trends, and weave in external knowledge to better grasp cultural and topical nuances.

Our model incorporates mechanisms to adaptively balance precision and recall, reducing unwarranted censorship of legitimate content and limiting the proliferation of harmful speech. We emphasize the importance of continuous monitoring and updating of hate speech detection systems in response to evolving language, culture, and societal contexts. Furthermore, we advocate for inclusive stakeholder engagement, involving domain experts and affected communities, to guide responsible design and deployment. We acknowledge the limitations of automated approaches and the ethical imperative for human oversight, transparent reporting, and accountability to foster safer and fairer online environments. Our work aspires to contribute positively to the broader efforts combating hate speech, promoting respect, dignity, and inclusivity in digital

spaces without undermining fundamental rights.

10 References

References

- Abdul Aziz, Md. Akram Hossain, and Abu Nowshed Chy. Csecu-dsg@multimodal hate speech event detection 2023: Transformer-based multimodal hierarchical fusion model for multimodal hate speech detection. In *Proceedings of the 6th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2023)*, pages 101–107, Varna, Bulgaria, 2023. INCOMA Ltd.
- Aashish Bhandari, Siddhant B Shah, Surendrabikram Thapa, Usman Naseem, and Mehwish Nasim. Crisshatem: Multimodal analysis of directed and undirected hate speech in text-embedded images from russia-ukraine conflict. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1994–2003, 2023.
- Anusha Chhabra and Dinesh Kumar Vishwakarma. Mhs-stma: Multimodal hate speech detection via scalable transformer-based multilevel attention framework. *arXiv preprint arXiv:2409.05136*, 2024.
- Ali Hürriyetoğlu, Surendrabikram Thapa, and Hristo Tanev. Findings and insights from the 8th workshop on challenges and applications of automated extraction of socio-political events from text. In *Proceedings of the 8th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2025)*, 2025.
- Mohammad Kashif, Mohammad Zohair, and Saquib Ali. Lexical squad@multimodal hate speech event detection 2023: Multimodal hate speech detection using fused ensemble approach. In *Proceedings of the 6th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2023)*, pages 84–91, Varna, Bulgaria, 2023. INCOMA Ltd.
- Anil Singh Parihar, Surendrabikram Thapa, and Sushruti Mishra. Hate speech detection using natural language processing: Applications and challenges. In *2021 5th International Conference on Trends in Electronics and Informatics (ICOEI)*, pages 1302–1308. IEEE, 2021.
- Umitcan Sahin, Izzet Emre Kucukkaya, Oguzhan Ozcelik, and Cagri Toraman. Arc-nlp at multimodal hate speech event detection 2023: Multimodal methods boosted by ensemble learning, syntactical and entity features. In *Proceedings of the 6th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2023)*, pages 71–78, Varna, Bulgaria, 2023. INCOMA Ltd.
- Siddhant Bikram Shah, Shuvam Shiwakoti, Maheep Chaudhary, and Haohan Wang. Meme-clip: Leveraging clip representations for multimodal meme classification. pages 17320–17332, 2024.

- Surendrabikram Thapa, Farhan Ahmad Jafri, Ali Hürriyetoğlu, Francielle Vargas, Roy Ka Wei Lee, and Usman Naseem. Multimodal hate speech event detection-shared task 4. In *CASE 2023-Proceedings of the 6th Workshop on Challenges and Applications of Automated Extraction of Socio-Political Events from Text, associated with 14th International Conference on Recent Advances in Natural Language Processing, RANLP 2023*, pages 151–159. Association for Computational Linguistics, 2023.
- Surendrabikram Thapa, Kritesh Rauniyar, Farhan Ahmad Jafri, Hariram Veeramani, Raghav Jain, Sandesh Jain, Francielle Vargas, Ali Hürriyetoğlu, and Usman Naseem. Extended multimodal hate speech event detection during russia-ukraine crisis-shared task at case 2024. In *7th Workshop on Challenges and Applications of Automated Extraction of Socio-Political Events from Text, CASE 2024*, pages 221–228. Association for Computational Linguistics, 2024.
- Surendrabikram Thapa, Siddhant Bikram Shah, Kritesh Rauniyar, Shuvam Shiwakoti, Surabhi Adhikari, Hariram Veeramani, Kristina T. Johnson, Ali Hürriyetoğlu, Hristo Tanev, and Usman Naseem. Multimodal hate, humor, and stance event detection in marginalized sociopolitical movements. In *Proceedings of the 8th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2025)*, 2025.
- Surendrabikram Thapa, Shuvam Shiwakoti, Siddhant Bikram Shah, Surabhi Adhikari, Hariram Veeramani, Mehwish Nasim, and Usman Naseem. Large language models (llm) in computational social science: prospects, current state, and challenges. *Social Network Analysis and Mining*, 15(1):1–30, 2025.