

Beyond Single-View Detection: A Dual-Space Reasoning Framework for Interpretable Harmful Meme Understanding

Anonymous ACL submission

Abstract

The identification of harmful memes extends beyond a mere classification task, encompassing challenges related to multi-perspective semantic comprehension and hierarchical reasoning. Prevailing approaches predominantly depend on modal alignment or black-box classifiers, which fail to capture implicit biases and lack interpretability. In this study, we propose BPDMoE-Hate, a novel framework grounded in dual-space mixture-of-experts, which innovatively conceptualizes harmful meme detection as an integrated process of “viewpoint decoupling and hierarchical fusion”. Our approach generates adversarial binary perspectives via Visual-Language Models (VLMs) and incorporates an adaptive viewpoint gating to facilitate viewpoint selection, thereby enabling the model to autonomously discern implicit semantic inclinations. Moreover, we propose the Hyperbolic-Euclidean space expert to effectively capture the hierarchical structural relationships and semantic correlations between multimodal and viewpoint features, thereby enabling interpretable reasoning at the geometric representation level. Empirical evaluations conducted on three mainstream datasets demonstrate that BPDMoE-Hate not only substantially surpasses existing methodologies in performance but also offers visual explanations for viewpoint selection and hierarchical structuring, thereby advancing the field of interpretable multimodal content analysis.

1 Introduction

The advancement of social media platforms has enhanced the capacity for individuals to express their emotions (Hermida and Santos, 2023; Liu et al., 2024a); however, it has concurrently contributed to the propagation of detrimental information. A notable example includes internet memes, which have gained widespread circulation in recent years. Typically, a meme comprises an image paired with

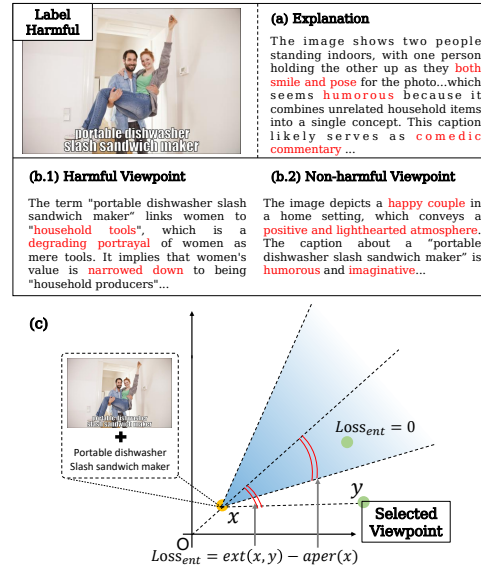


Figure 1: Binary viewpoints and entailment loss.

concise textual content conveying a particular viewpoint. Certain users exploit this format to disseminate hate speech, thereby facilitating and intensifying violent (Mei et al., 2024) conduct in real-world contexts.

Viewpoint decoupling. Previous hate speech meme detection methods either adopt a single classifier (Hebert et al., 2024; Lu et al., 2024) or cross-modal alignment (Yang et al., 2024) to narrow the inter-modality semantic gap, but these black-box models (Lin et al., 2024) exhibit limited interpretability. Recently, studies have leveraged VLMs with zero-shot prompting (Kojima et al., 2022) and multi-agent debate frameworks (Park et al., 2024; Zheng et al., 2024) to interpret harmful memes (Cao et al., 2023; Ji et al., 2024; Lin et al., 2025), and even incorporated explanatory text (Hee and Lee, 2025) or debate outcomes (Lin et al., 2024; Zhou et al., 2025) into model training for performance improvement. However, these interpretative approaches still cannot model multi-perspective semantic conflicts, and are thus prone to introducing

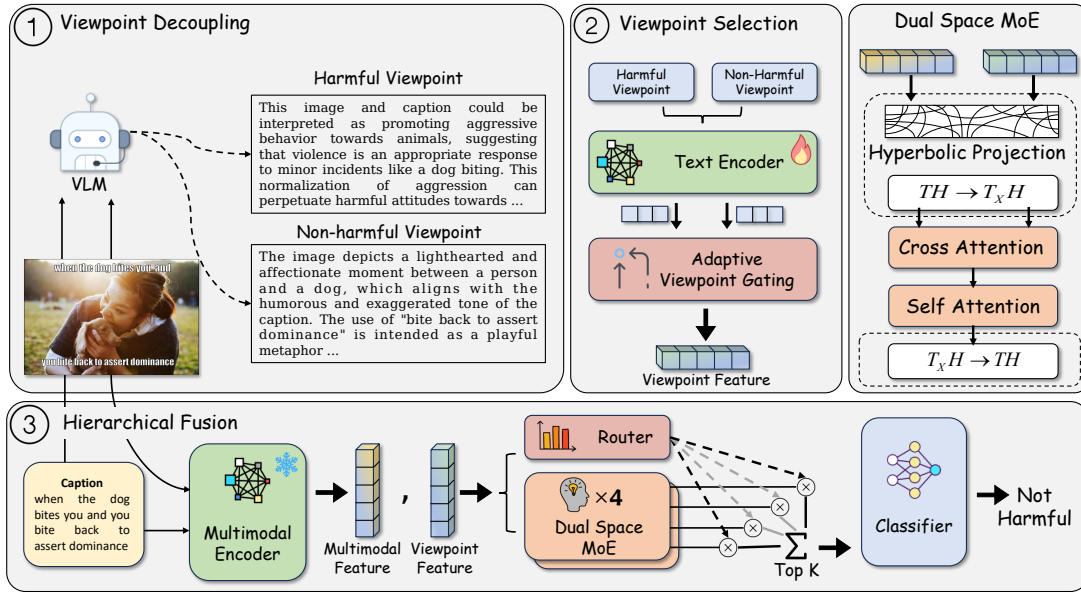


Figure 2: Our overall framework: generate decoupled binary viewpoints via the VLM, select viewpoints using AVG, and then leverage DSMoE to learn dual-space features for classifier-based prediction. The dashed box delineates the mapping and inverse mapping operations that are unique to hyperbolic space.

065 the model’s subjective bias towards memes.

066 As illustrated in Figure 1, the text within the
 067 meme diminishes the value of women by portray-
 068 ing them as mere “household tools”. Figure 1 (a)
 069 incorporates the model’s own subjective judgment,
 070 erroneously interpreting the content as humorous.
 071 Consequently, integrating such explanations into
 072 the training process risks propagating subjective
 073 biases, leading to misclassification. Panels (b.1)
 074 and (b.2) represent the explanatory approach we
 075 propose, which involves a more in-depth analysis
 076 of memes from two distinct perspectives. This de-
 077 coupled adversarial viewpoint partially mitigates
 078 the model’s inherent bias. Nonetheless, only one
 079 of these perspectives accurately reflects the true
 080 nature of the harmful meme. Consequently, this
 081 raises an important question: how can the model
 082 be trained to autonomously select the appropriate
 083 viewpoint?

084 **Hierarchical reasoning.** Beyond the viewpoint
 085 selection challenge, a more critical issue unad-
 086 dressed by existing methods is the lack of a system-
 087 atic framework to model the hierarchical seman-
 088 tic relationships between multimodal content and
 089 interpretative perspectives. Regarding the meme
 090 depicted in Figure 1, from the perspective of cog-
 091 nitive logic (Van Ditmarsch et al., 2008), human
 092 observers first synthesize information from both
 093 components (image and title) to achieve compre-
 094 hension, subsequently generating clear and specific

095 views on it. These perspectives represent a more
 096 profound understanding of the multimodal informa-
 097 tion—comprising both images and texts—thereby
 098 establishing a hierarchical structure (Vendrov et al.,
 099 2015) that progresses from “multimodal content” to
 100 “human viewpoints”. Given that hyperbolic space
 101 (Desai et al., 2023; Pal et al., 2024) near the origin
 102 encodes more general information, while regions
 103 closer to the boundary convey more specific at-
 104 tributes, it is appropriate to represent multimodal
 105 content and viewpoints as “roots” situated near the
 106 origin and “leaves” positioned closer to the periph-
 107 ery within the hyperbolic space.

108 This study revisits the problem of harmful meme
 109 detection from the perspective of geometric rep-
 110 resentation learning. We hypothesize that under-
 111 standing harmful memes entails a reasoning
 112 process of “viewpoint decoupling-hierarchical fu-
 113 sion”. Specifically, we propose BPDMoE-Hate,
 114 a **Binary Perspectives Dual-space Mixture-of-**
 115 **Experts** framework comprising two core compo-
 116 nents: an Adaptive Viewpoint Gating (AVG) mod-
 117 ule for viewpoint selection and a Dual-Space MoE
 118 (DSMoE) module for hierarchical feature fusion.
 119 In the viewpoint decoupling phase, dual viewpoints
 120 are generated via adversarial prompting, and the
 121 AVG module is designed to enable adaptive view-
 122 point selection, thereby empowering the model to
 123 discern semantic veracity. Subsequently, in the hi-
 124 erarchical fusion phase, the hierarchical structure

of multimodal and viewpoint features is modeled in hyperbolic space, while semantic relationships are captured in Euclidean space. Finally, the DSMoE module dynamically integrates these two types of representations to form a unified and discriminative feature space.

Our contributions are as follows: 1) We propose a multi-view decoupled reasoning framework, which formalizes the task as a process of view generation, selection and fusion. 2) We design a DSMoE to provide an interpretable geometric foundation for multimodal reasoning. 3) Experimental results verify the superior performance of our approach.

2 Preliminaries

In this section, we provide a concise overview of the essential concepts underlying hyperbolic geometry; more details can be found in (Chami et al., 2019). Hyperbolic space is characterized as the unique complete, simply connected Riemannian manifold that is isotropic and exhibits constant negative sectional curvature (Wang et al., 2024). The curvature parameter quantifies the extent to which hyperbolic space diverges from Euclidean flatness. The models used to describe hyperbolic space mainly include the Lorentz model and the Poincaré model. This work focuses exclusively on the Poincaré model.

Consider the n -dimensional Poincaré sphere endowed with a constant negative curvature $-c$ ($c > 0$). The associated Riemannian manifold can be characterized as $\mathcal{H}^{n,c} = \{x \in \mathcal{H}^n \mid \|x\|^2 < \frac{1}{c}\}$. For any point $x_{\mathcal{H}} \in \mathcal{H}^{n,c}$ within this hyperbolic space, there exists a corresponding tangent space $\mathcal{T}_x \mathcal{H}^{n,c}$, which serves as a local, first-order approximation of the manifold at $x_{\mathcal{H}}$. The hyperbolic space and its tangent space at $x_{\mathcal{H}}$ are related through the following mapping:

$$\exp_x^c(v) = v \oplus_c \left(\tanh\left(\sqrt{c} \frac{\lambda_x^c \|v\|}{2}\right) \frac{v}{\sqrt{c} \|v\|} \right) \quad (1)$$

$$\log_x^c(y) = \frac{2 \tanh^{-1}(\sqrt{c} \| -x \oplus_c y \|) (-x \oplus_c y)}{\sqrt{c} \lambda_x^c \| -x \oplus_c y \|} \quad (2)$$

In this context, $v \in \mathcal{T}_x \mathcal{H}^{n,c}$, $x, y \in \mathcal{H}^{n,c}$. $\lambda_x^c = \frac{2}{1-c\|x\|^2}$ is the conformal factor, \oplus_c denotes Möbius addition (Ungar, 2001). The exponential operation maps $\mathcal{T}_x \mathcal{H}^{n,c}$ to $\mathcal{H}^{n,c}$, while the logarithmic operation maps $\mathcal{H}^{n,c}$ to $\mathcal{T}_x \mathcal{H}^{n,c}$. The distance

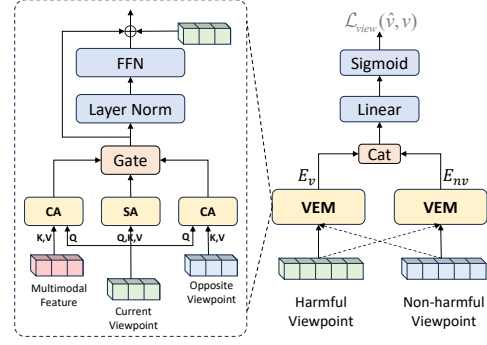


Figure 3: Adaptive viewpoint gating. VEM represents our viewpoint enhancement module, CA signifies Cross-Attention, while SA denotes Self-Attention.

within the Poincaré sphere is defined as the shortest path between two points x and y :

$$d_{xy}^c = \frac{2}{\sqrt{c}} \operatorname{arctanh}(\sqrt{c} \| -x \oplus_c y \|) \quad (3)$$

3 Method

3.1 Viewpoint Decoupling

The overall architecture is illustrated in Figure 2. To mitigate the subjective biases (Ye et al., 2024) present in models during the generation of VLM explanations, we conceptualize the generation of viewpoints as an adversarial semantic completion task. Concretely, we design a pair of adversarial prompt templates that direct the VLM to produce two opposing textual descriptions—one harmful and one harmless—for the identical meme, as illustrated in step 1 of Figure 2. This approach can be interpreted as a semantic decoupling of the original multimodal content, with the objective of disentangling implicit semantic dimensions that may contribute to ambiguity in the final evaluation. The corresponding prompt templates are presented in Appendix B.

3.2 Viewpoint Selection

We use the text encoder to encode the decoupled binary viewpoints. Nevertheless, the model must identify the viewpoint that more accurately corresponds to the true semantics between the two contrasting perspectives, which is the prerequisite for achieving reliable hierarchical fusion.

3.2.1 Adaptive Viewpoint Gating

Inspired by (Zhao et al., 2022), we design an AVG to conduct the selection of viewpoint features. As illustrated in Figure 3, the two viewpoint feature

vectors produced by the text encoder are denoted as $H_v, H_{nv} \in \mathcal{R}^{B \times L \times N}$. Where B represents the batch size, L denotes the length of the text sequence, N represents the output dimension of the hidden layer. These two vectors are individually enhanced through the core module VEM (Viewpoint Enhancement Module), allowing each viewpoint feature to interactively reference its complementary viewpoint as well as the original multimodal context. Consequently, this process facilitates the acquisition of more discriminative representations.

Subsequently, the enhanced feature representations E_v and E_{nv} are derived and subsequently fed into the classification layer to produce the final viewpoint selection weights $W_v \in \mathcal{R}^{B \times 1}$. This process is formally represented by the following equation:

$$W_v = \text{Sigmoid}(\text{Linear}(\text{Cat}(E_v, E_{nv}))) \quad (4)$$

Here, “ $\text{Cat}(\cdot)$ ” represents feature concatenation. After obtaining the weights for the viewpoint selection, for each sample $h_v^i \in \mathcal{R}^{L \times N}$, $h_{nv}^i \in \mathcal{R}^{L \times N}$ and $w_v^i \in \mathcal{R}^{1 \times 1}$, where $i \in B$, we perform the viewpoint selection through the following formula:

$$h_s^i = w_v^i \cdot h_v^i + (1 - w_v^i) \cdot h_{nv}^i \quad (5)$$

Moreover, we develop a viewpoint selection loss, wherein the true labels are consistent with the downstream harmful detection task, forcing the selection mechanism to be congruent with the final task objective.:

$$\mathcal{L}_{view} = \mathcal{L}_{CE}(\hat{v}, v) \quad (6)$$

Where \mathcal{L}_{CE} denotes the cross-entropy loss, \hat{v} signifies the prediction generated by AVG for the selection of viewpoints, and v corresponds to the true label.

3.2.2 Viewpoint Enhancement Module

Suppose the target sample currently input is h_v^i , that is “Current Viewpoint”, then h_{nv}^i is “Opposite Viewpoint”, while the features from the multimodal encoder are $m^i \in \mathcal{R}^{M \times N}$. Firstly, perform multi-head self-attention and cross-attention operations on the current features, that is, the SA and CA modules in Figure 3. Here, we use CA as an example. The Queries are derived from h_v^i , represented as $Q_t = I_t W_q$, while the Keys and Values come from m^i , represented as $K_s = I_s W_k$ and $V_s = I_s W_v$, where $W_q, W_k, W_v \in \mathcal{R}^{N \times N}$ are learnable parameters. The CA can then be formulated as:

$$\begin{aligned} o_{cm}^i &= \text{Softmax} \left(\frac{Q_t K_s^T}{\sqrt{d_k}} V_s \right) \\ &= \text{Softmax} \left(\frac{I_t W_q W_k^T I_s^T}{\sqrt{d_k}} \right) I_s W_v \end{aligned} \quad (7)$$

In the SA module, the Queries, Keys, and Values are all derived from the same feature. We represent the vectors produced by the SA and an additional CA as o_c^i and o_{co}^i , respectively. These are then combined using a “Gate” followed by a feed-forward network (FFN) to achieve an enhanced representation $e_v^i \subset E_v$ (Appendix E.8 demonstrates the effectiveness of VEM):

$$\begin{cases} e_v^i = h_v^i + e^i + \text{FFN}(\text{LN}(e^i)) \\ e^i = w_e^i \cdot o_c^i + (1 - w_e^i) \cdot (o_{cm}^i + o_{co}^i)/2 \\ w_e^i = \text{Sigmoid}(\text{Linear}(\text{Cat}(o_c^i, o_{cm}^i))) \end{cases} \quad (8)$$

3.3 Dual Space MoE

3.3.1 Hyperbolic Space Experts

To construct the hierarchical structure within hyperbolic space, we initially employ Eq.1 to project the selected viewpoint and multimodal features onto the hyperbolic manifold. As indicated in (Wang et al., 2024), it is necessary to first apply a linear transformation to the features to reduce their dimensionality, thereby enhancing the representational capacity of the hyperbolic manifold. Given a weight matrix $W \in \mathcal{R}^{m \times n}$ and a bias term $b \in \mathcal{R}^{1 \times 1}$ associated with a linear layer, the matrix multiplication operation within the Poincaré ball is performed via Möbius multiplication:

$$W \oplus_c x = \exp_0^c(W \cdot \log_0^c(x)) \quad (9)$$

In the context of biased addition, the vector b represents a vector located in the tangent space $\mathcal{T}_0 \mathcal{H}^{n,c}$. It is necessary to transport this vector to the tangent space $\mathcal{T}_x \mathcal{H}^{n,c}$ at the point x , after which it is projected onto the hyperbolic manifold:

$$T_{0 \rightarrow x}^c(b) = \log_x^c(x \oplus_c \exp_0^c(b)) \quad (10)$$

Where $x \mathcal{H} \oplus_c b = \exp_x^c(T_{0 \rightarrow x}^c(b)) = \frac{\lambda_x^c}{\lambda_x^c} b$. As illustrated in Figure 1(c), we further adopt the entailment loss (Desai et al., 2023) to constrain viewpoint features within the entailment cone centered on multimodal features, thus geometrically encoding the “multimodal entailment viewpoint” reasoning relationship explicitly. Let x and y represent

the multimodal and viewpoint vectors embedded in hyperbolic space, respectively. Here, x should lie closer to the hyperbolic origin, while y should reside within the entailment cone originating at x . The half-aperture angle $aper(x)$ of x is defined as follows:

$$aper(x) = \arcsin\left(\frac{r_{\min}\sqrt{c}}{\tanh\left(\frac{\sqrt{c}d_{ox}^c}{2} + \varepsilon\right)}\right) \quad (11)$$

Here $r_{\min} = 0.1$ denotes the minimum radius of the containment cone, $\varepsilon = 1 \times 10^{-8}$, c is a learnable curvature value, initialized to 1.0 and decreasing during training. The external angle $ext(x, y) = \pi - \angle Oxy$, thus $\angle Oxy$ can be obtained by the following formula:

$$\angle Oxy = \cos^{-1}\left(\frac{\cosh(D_{ox}^c)\cosh(D_{xy}^c) - \cosh(D_{oy}^c)}{\sinh(D_{ox}^c)\sinh(D_{xy}^c) + \varepsilon}\right) \quad (12)$$

Where $D_{AB}^c = \sqrt{c}d_{AB}^c$, while d_{AB}^c corresponds to the hyperbolic distance between points A and B , as determined by Eq.3. Then, the final entailment loss is expressed as:

$$\mathcal{L}_{ent} = \max(0, ext(x, y) - aper(x)) \quad (13)$$

Moreover, we enable the implementation of both cross-attention and self-attention modules within hyperbolic space (as shown in Figure 2). Subsequently, as described by Eq.2, the fused features are mapped back to Euclidean space.

3.3.2 Euclidean Space Expert

Euclidean space experts are responsible for modeling the semantic connections between multimodal and viewpoints. In this instance, the component enclosed by the dashed box in Figure 2 is excluded, retaining solely the attention mechanisms.

3.3.3 Hierarchical Fusion and Prediction

Our DSMoE innovatively integrates experts from different spaces into the MoE (Liu et al., 2025) layer, which is composed of a Router and n expert networks, where we set $n = 4$. The experts with odd indices operate in Euclidean space, while those with even indices function in hyperbolic space. As shown in Figure 2, the Router (a standard transformer block), dynamically integrates experts from

two distinct domains, enabling the model to adaptively integrate hierarchical reasoning and semantic correlation information, thereby facilitating hierarchical fusion subsequent to viewpoint decoupling. Subsequently, a classification layer is employed to derive the probabilities of n experts being assigned and top 2 experts are chosen to participate in the computation. The detailed formula is as follows:

$$O_{exp} = \sum_{i=1}^n w_{exp}^i \cdot o_{exp}^i \quad (14)$$

Where o_{exp}^i denotes the output from the i th expert. Ultimately, the harmful meme prediction is produced by the classification layer. We employ the cross-entropy loss between the model’s predicted category and the actual category, expressed as $\mathcal{L}_{task} = \mathcal{L}_{CE}(\hat{y}, y)$.

In addition, we also introduce contrastive learning loss $\mathcal{L}_{InfoICE}$ (Oord et al., 2018) to better distinguish the positive and negative samples in the sampling process, and load balancing loss $\mathcal{L}_{balance}$ (Fedus et al., 2022) to promote the balanced loading of experts. For detailed information, please refer to Appendix C and E.4. The overall training loss is as follows:

$$\mathcal{L}_{all} = \mathcal{L}_{task} + \alpha\mathcal{L}_{view} + \beta\mathcal{L}_{ent} + \gamma\mathcal{L}_{balance} + \eta\mathcal{L}_{InfoICE} \quad (15)$$

4 Experiments

4.1 Evaluation Dataset and Baselines

We evaluate BPDME-Hate on three widely used hate meme datasets, including Facebook’s Hateful Meme (FHM) (Kiela et al., 2020), Multimedia Automatic Misogyny Identification (MAMI) (Fersini et al., 2022), and Harmful Meme (HarMeme) (Praninck et al., 2021). Furthermore, we demonstrate the effectiveness of our method by comparing it with some state-of-the-art meme detection models (separated by a double line in Table 1), which include: 1) VLMs 2) pure text classification models and multimodal classification models 3) harmful meme detection framework. For detailed supplementary information on the datasets and references to the models, please refer to Appendix D.

4.2 Implementation Details

To guarantee the reliability of the decoupling of viewpoints, we employ Qwen2.5-VL-32B-Instruct (Qwen2.5-VL) as the viewpoint generation model.

Model	FHM		MAMI		HarMeme	
	AUC	ACC	AUC	ACC	AUC	ACC
Llama-3.2-V	-	63.38 \pm 0.83	-	69.14 \pm 0.69	-	69.21 \pm 1.61
Llava-1.5	-	52.34 \pm 0.81	-	53.72 \pm 1.45	-	55.37 \pm 1.03
Qwen2.5-VL	-	73.52 \pm 0.16	-	69.78 \pm 0.25	-	63.33 \pm 0.13
BERT-base	64.98 \pm 0.61	57.86 \pm 0.68	71.56 \pm 0.65	64.28 \pm 0.52	81.38 \pm 0.88	75.31 \pm 1.19
RoBERTa-large	64.40 \pm 1.29	58.56 \pm 0.43	72.46 \pm 0.65	66.08 \pm 1.15	81.72 \pm 0.83	76.38 \pm 0.90
FLAVA-full	78.60 \pm 0.74	70.02 \pm 1.39	81.24 \pm 0.55	70.10 \pm 0.61	85.45 \pm 1.06	79.72 \pm 1.78
VisualBERT*	68.71 \pm 1.02	61.48 \pm 1.19	78.71 \pm 0.59	71.06 \pm 0.94	80.46 \pm 1.04	75.31 \pm 1.44
ViLBERT*	73.05 \pm 0.62	64.70 \pm 1.12	77.71 \pm 1.20	69.48 \pm 1.00	84.11 \pm 0.88	78.70 \pm 1.17
BLIP2	63.52 \pm 0.62	58.18 \pm 0.96	82.05 \pm 0.20	65.50 \pm 3.53	89.94 \pm 0.14	80.62 \pm 1.84
ALBEF*	79.40 \pm 0.53	70.58 \pm 0.50	83.24 \pm 0.93	72.77 \pm 1.00	85.49 \pm 1.23	80.99 \pm 0.80
Mod-HATE	64.50 \pm 0.19	58.00 \pm 1.07	67.40 \pm 0.46	61.00 \pm 2.22	73.40 \pm 0.27	69.40 \pm 0.42
PromptHate	76.76 \pm 0.95	67.82 \pm 1.23	76.21 \pm 1.05	68.08 \pm 0.58	87.51 \pm 0.74	79.38 \pm 1.72
Pro-Cap	80.87 \pm 0.66	72.28 \pm 0.90	82.53 \pm 0.49	73.06 \pm 0.82	90.25 \pm 0.54	83.25 \pm 1.00
ExplainHM	82.32 \pm 1.12	72.22 \pm 1.62	79.07 \pm 2.13	71.03 \pm 0.88	75.58 \pm 4.92	77.16 \pm 1.98
IntMeme	81.50 \pm 1.11	71.52 \pm 1.49	81.89 \pm 1.15	72.30 \pm 1.79	89.35 \pm 1.22	81.92 \pm 2.47
BPDMoE-Hate	83.71 \pm0.39	75.18 \pm1.38	87.84 \pm0.54	76.70 \pm0.80	94.11 \pm0.28	86.10 \pm0.94

Table 1: The results are presented as the “mean \pm standard deviation” of the outcomes derived from five distinct random seeds. * indicates the results are from (Cao et al., 2023). We use the same VLM to generate the explanations.

The text encoder and multimodal encoder are instantiated using RoBERTa-large and BLIP2, respectively. During training, the multimodal encoder parameters are frozen, while the final two layers of the text encoder are fine-tuned over 4 epochs. The loss function incorporates weighted components with coefficients $\alpha, \beta, \gamma, \eta$ set to 0.5, 1×10^{-2} , 1×10^{-3} , and 0.1, respectively. A learning rate of 1×10^{-5} is adopted, and the entire training procedure is conducted for 10 epochs (batch size is set to 48) utilizing one NVIDIA A6000 GPU. Model performance is evaluated using Accuracy (ACC) and the Area Under the Receiver Operating Characteristic Curve (AUC) as metrics.

4.3 Main Results

As presented in Table 1, the BPDMoE-Hate model demonstrated superior performance across all three datasets. Notably, while explanation-based approaches such as ExplainHM and IntMeme yielded competitive results—particularly with Pro-Cap attaining an AUC of 90.25 on the HarMeme dataset—our method, by effectively integrating two distinct perspectives, achieved a higher AUC of 94.11. This enhancement substantiates the efficacy of the dual-view decoupling and selection mechanism in identifying implicit harmful content. Furthermore, the consistent outperformance on the MAMI dataset, which focuses on explicitly hateful content directed towards women, as well as the FHM dataset, characterized by mixed content types,

suggests that the proposed dual-space hierarchical fusion framework possesses strong generalizability across diverse tasks.

4.4 Ablation Study

As shown in Figure 4, we demonstrate the contributions of different modules. The simultaneous removal of AVG and DSMoE resulted in a marked decline in performance, indicating that AVG plays a critical role in selection following perspective decoupling, whereas DSMoE is essential for hierarchical fusion. Both components are therefore indispensable. This highlights the importance of the two-stage architecture of our framework, termed “decoupling-fusion”. The elimination of AVG alone led to a more pronounced performance degradation, underscoring its critical role in mitigating model biases and autonomously identifying the accurate perspective. Interestingly, the removal of DSMoE alone on the MAMI dataset resulted in a slight increase in AUC. We believe that the misogynistic content within this dataset is predominantly expressed explicitly, thereby reducing the reliance on complex hierarchical reasoning.

4.5 The Significance of Viewpoints

Given the importance of uncovering the latent information within memes for enhancing model decision-making, we conduct a series of experiments to assess the influence of various interpretations. The results, presented in Table 2, indicate

Viewpoint Type	FHM		MAMI		HarMeme	
	AUC	ACC	AUC	ACC	AUC	ACC
BPDMoE-Hate	83.71 \pm 0.39	75.18 \pm 1.38	87.84 \pm 0.54	76.70 \pm 0.80	94.11 \pm 0.28	86.10 \pm 0.94
Random Select	79.78 \pm 1.21	69.88 \pm 1.43	86.54 \pm 0.87	73.04 \pm 1.29	89.95 \pm 0.90	83.33 \pm 1.47
Harmful View	82.06 \pm 0.56	72.78 \pm 1.88	87.12 \pm 0.50	74.04 \pm 1.31	90.29 \pm 1.07	83.79 \pm 1.11
Non-harmful View	79.35 \pm 0.55	69.04 \pm 1.09	86.60 \pm 1.05	73.90 \pm 1.71	89.98 \pm 1.45	83.45 \pm 1.33
No View	70.16 \pm 2.05	61.98 \pm 1.17	82.49 \pm 1.14	67.74 \pm 1.22	92.74 \pm 0.11	84.58 \pm 1.89
Only Explanation	82.80 \pm 0.16	73.18 \pm 1.81	87.10 \pm 0.61	75.38 \pm 1.09	92.75 \pm 0.24	85.76 \pm 0.96

Table 2: The influence of different viewpoints. We set up 5 types of viewpoints for training, including: 1) Randomly select one of the viewpoints. 2) Only use harmful viewpoints. 3) Only use harmless viewpoints. 4) Remove both types of viewpoints. 5) Replace both types of viewpoints with explanations for memes.

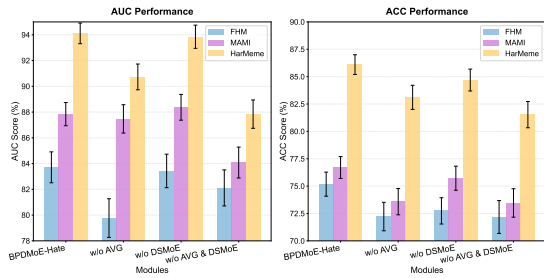


Figure 4: The impact of different modules on the overall performance of the model. We gradually removed AVG and DSMoE.

that employing a single viewpoint leads to a noticeable decline in performance. Furthermore, the approach of randomly selecting viewpoints, which disrupts the AVG’s selection strategy, yields even poorer outcomes compared to using only one viewpoint. The performance without using any viewpoint is the worst, underscoring the critical role of hidden information extraction in meme detection. Conversely, utilizing solely model explanations produces comparatively better results. These findings demonstrate that the concurrent integration of both viewpoints exerts a more beneficial impact on the model’s decision-making process.

4.6 The Significance of DSMoE

In this section, we conducted tests using the same random seed, and the experimental results are presented in Figure 5(a). It is evident that the AUC and ACC values for both “Only HS” and “Only ES” are inferior to those of BPDMoE-Hate, indicating that our model effectively captures complementary information from different spaces. The absence of either spatial feature representation adversely affects the model’s performance. Furthermore, the overall performance of “Only HS” surpasses “Only ES”, suggesting that hierarchical structure modeling in hyperbolic space contributes more positively

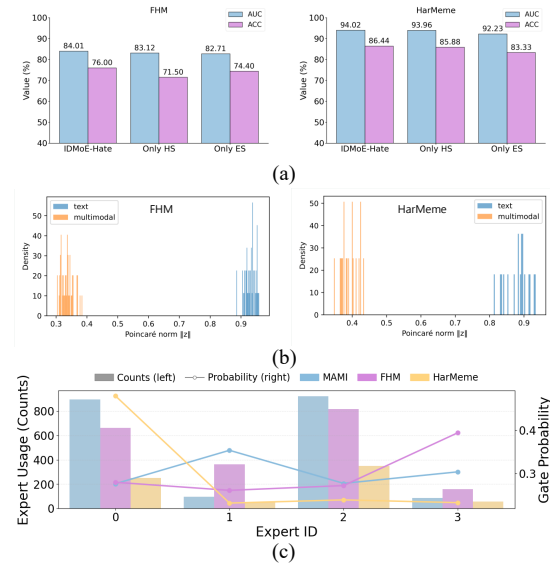


Figure 5: The importance of DSMoE. In (a), “Only ES” indicates replacing all experts with Euclidean experts, while “Only HS” represents “only hyperbolic space experts”. (b) represents the hierarchical structure in hyperbolic space. “text” represents the selected viewpoint feature. (c) denotes the frequency of activation for each expert alongside the corresponding mean probability.

to the model’s gains.

To validate the geometric assumptions underpinning the hierarchical structure modeling in hyperbolic space, we utilize the visualization method from (Pal et al., 2024) and present the learned hyperbolic space structure through low-dimensional visualization (showing the spatial norm distribution of the test set samples in the form of histograms). As shown in Figure 5(b), multimodal features consistently cluster near the origin of the hyperbolic space, while viewpoint features are located farther away. This notable geometric distinction offers quantitative evidence that the hierarchical structure of the “multimodal entailment viewpoint features” has been effectively represented.

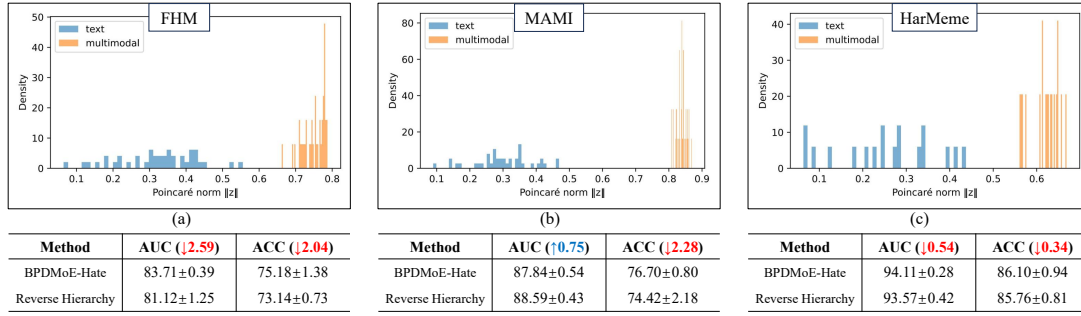


Figure 6: The influence of the reverse hierarchical structure. We forcibly reversed the hierarchical relationship between the viewpoint (“text”) and the multimodal information.

As shown in Figure 5(c), experts in the hyperbolic space are more frequently activated on most datasets. This finding suggests that the model predominantly depends on hyperbolic space representations, which are capable of encoding hierarchical semantic relationships, during decision-making processes. From the perspective of the activation probabilities of experts, except for the HarMeme dataset, which tends to assign greater weight to the hyperbolic space expert with ID 0, the other two datasets exhibit a preference for allocating higher probabilities to the Euclidean space experts. This observation clearly demonstrates the flexibility of our framework in selecting different experts.

4.7 Reverse Hierarchical Structure

We performed a causal intervention experiment by repositioning the viewpoint features to the origin of the hyperbolic space while situating the multimodal features at the periphery. The outcomes, presented in Figure 6, demonstrate a marked decrease in performance across the three datasets. These findings indicate that the observed positive hierarchical relationship is not merely a coincidental correlation within the data but constitutes a causal factor contributing to the model’s superior performance. Furthermore, we also verify the effect of removing this hierarchical relationship, detailed information is provided in Appendix E.6.

4.8 Case Study

Figure 7(a) is the malicious pun on the Asian surname “Wong” by deliberately spelling “something wrong” as “sum ting wong” to imitate a stereotypical Asian accent is offensive and defamatory to the Asian community. Our harmful viewpoint correctly explains this point. Additionally, our Router utilizes hyperbolic space experts to accurately identify such harmful memes, and AVG successfully



Figure 7: Case illustration. “Router” represents the activated experts and their assigned probabilities, while “View_Select_Pred” indicates the probability of each viewpoint being selected.

select the correct viewpoint. Figure 7(b) depicts the actor portraying Iron Man exhibiting relief upon confronting the masked criminal. Since typical robberies do not necessitate superhero intervention, this scenario conveys a humorous and teasing tone. These observations further substantiate the efficacy of our proposed framework.

5 Conclusion

This paper presents BPDMoE-Hate, a dual-space viewpoint decoupled reasoning framework for detecting harmful memes. Our framework utilizes the generated dual viewpoints as the input for semantic decoupling, and employs an AVG to autonomously determine the semantic authenticity, effectively alleviating the model’s subjective bias. Furthermore, we designed a dual-space MoE, which explicitly models the hierarchical entailment relationship between multimodal and perspective features in the hyperbolic space, and learns semantic associations in the Euclidean space, achieving the synergy of structured reasoning and semantic matching.

529 Limitations

530 The limitations of this study are as follows. Our
531 proposed framework depends on a robust VLM
532 to produce high-quality binary perspectives. Al-
533 though the implementation of a perspective selec-
534 tion mechanism mitigates bias to some extent, the
535 diversity and comprehensiveness of the generated
536 perspectives remain constrained by the inherent
537 cognitive limitations of the VLM. In Appendix
538 E.3, we provide a detailed argument supporting
539 the decoupling perspective generated by the small
540 model. Future research will investigate retrieval-
541 augmented generation (RAG) techniques to en-
542 hance the quality and controllability of perspective
543 generation in smaller models.

544 Ethical Considerations

545 **Data Privacy and Compliance:** The datasets used
546 in this study (FHM, MAMI, HarMeme) are pub-
547 licly available and used in strict accordance with
548 the original authors’ terms of use. We have not
549 collected any additional personal data, and all mul-
550 timodal content (images, captions) has undergone
551 anonymization processing (e.g., removing identi-
552 fiable personal information, desensitizing faces or
553 private logos) to avoid infringing on user privacy
554 (e.g., age, location, gender identity).

555 **Bias Mitigation and Fairness:** To mitigate
556 potential biases in harmful meme detection, our
557 dual-perspective generation mechanism—anchored
558 in adversarial harmful/non-harmful viewpoint
559 prompts—aims to eliminate implicit semantic bi-
560 ases within multimodal data, such as gendered or
561 racial stereotypes. Qualitative analyses and ab-
562 lation studies validate that our hierarchical fusion
563 framework, by balancing semantic modeling across
564 these dual perspectives, reduces disproportionate
565 impacts on marginalized groups. Despite these
566 safeguards, the model may still encounter gener-
567 alization challenges in cross-cultural contexts or
568 low-resource language settings, as the underlying
569 perspective generation mechanism struggles to cap-
570 ture cultural nuances and language-specific char-
571 acteristics in such scenarios. We acknowledge this as
572 a significant limitation of our approach.

573 **Expected Use and Misuse Prevention:** This
574 framework is exclusively designed for detecting
575 and mitigating harmful memes, with the goal of
576 promoting a safer online environment. We explic-
577 itly prohibit its use for any purpose that violates
578 ethical or legal norms, including but not limited

579 to: (1) propagating hate speech or discriminatory
580 content; (2) abusive surveillance or inappropriate
581 censorship of legitimate speech; (3) targeting vul-
582 nerable groups with biased detection results. We
583 will provide a detailed user guide to clarify the
584 model’s applicable scenarios and limitations, and
585 encourage users to report misuse cases.

586 **Social Benefits and Limitations:** This research
587 contributes to reducing the spread of harmful con-
588 tent online, protecting vulnerable groups (e.g.,
589 women targeted by misogynistic memes, racial mi-
590 norities) from discrimination, and fostering inclu-
591 sive digital spaces. However, we acknowledge that
592 no detection model can achieve 100% accuracy:
593 false positives may restrict legitimate speech, and
594 false negatives may allow harmful content to evade
595 detection. Future work will focus on improving
596 cross-cultural adaptability and reducing such risks.

597 **Code Release:** The publicly available dataset
598 we utilized contains harmful language and attacks
599 directed at minority groups; therefore, we have cho-
600 sen not to release the binary opinion dataset gen-
601 erated by the VLM. However, the relevant dataset
602 can be reasonably obtained by requesting it from
603 the original authors, and our prompt templates can
604 be employed to create new binary opinion datasets.
605 Additionally, our code is included in the supple-
606 mentary file “code.zip” and will be made publicly
607 accessible upon acceptance of the paper.

608 References

- 609 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
610 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie
611 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl
612 technical report. *arXiv preprint arXiv:2502.13923*.
- 613 Rui Cao, Ming Shan Hee, Adriel Kuek, Wen-Haw
614 Chong, Roy Ka-Wei Lee, and Jing Jiang. 2023. Pro-
615 cap: Leveraging a frozen vision-language model for
616 hateful meme detection. In *Proceedings of the 31st
617 ACM international conference on multimedia*, pages
618 5244–5252.
- 619 Rui Cao, Roy Ka-Wei Lee, Wen-Haw Chong, and Jing
620 Jiang. 2022. Prompting for multimodal hateful meme
621 classification. In *Proceedings of the 2022 conference
622 on empirical methods in natural language processing*,
623 pages 321–332.
- 624 Rui Cao, Roy Ka-Wei Lee, and Jing Jiang. 2024. Modu-
625 larized networks for few-shot hateful meme detection.
626 In *Proceedings of the ACM Web Conference 2024*,
627 pages 4575–4584.
- 628 Ines Chami, Zhitao Ying, Christopher Ré, and Jure
629 Leskovec. 2019. Hyperbolic graph convolutional

630	neural networks. <i>Advances in neural information processing systems</i> , 32.	Jinfa Huang, Jinsheng Pan, Zhongwei Wan, Hanjia Lyu, and Jiebo Luo. 2025. Evolver: Chain-of-evolution prompting to boost large multimodal models for hateful meme detection. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 7321–7330.	684
631			685
632	Anusha Chhabra and Dinesh Kumar Vishwakarma. 2023. Multimodal hate speech detection via multi-scale visual kernels and knowledge distillation architecture. <i>Engineering Applications of Artificial Intelligence</i> , 126:106991.	Junhui Ji, Xuanrui Lin, and Usman Naseem. 2024. Capalign: Improving cross modal alignment via informative captioning for harmful meme detection. In <i>Proceedings of the ACM Web Conference 2024</i> , pages 4585–4594.	686
633			687
634			688
635			689
636			690
637	Karan Desai, Maximilian Nickel, Tanmay Rajpurohit, Justin Johnson, and Shanmukha Ramakrishna Vedantam. 2023. Hyperbolic image-text representations. In <i>International Conference on Machine Learning</i> , pages 7694–7731. PMLR.	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes. <i>Advances in neural information processing systems</i> , 33:2611–2624.	691
638			692
639			693
640			694
641			695
642	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, volume 1 (long and short papers)</i> , pages 4171–4186.	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	696
643			697
644			698
645			699
646			700
647			701
648			702
649	Alexey Dosovitskiy. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. <i>arXiv preprint arXiv:2010.11929</i> .	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International conference on machine learning</i> , pages 19730–19742. PMLR.	703
650			704
651			705
652	William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. <i>Journal of Machine Learning Research</i> , 23(120):1–39.	Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. <i>Advances in neural information processing systems</i> , 34:9694–9705.	706
653			707
654			708
655			709
656	Elisabetta Fersini, Francesca Gasparini, Giulia Rizzi, Aurora Saibene, Berta Chulvi, Paolo Rosso, Alyssa Lees, and Jeffrey Sorensen. 2022. Semeval-2022 task 5: Multimedia automatic misogyny identification. In <i>Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)</i> , pages 533–549.	Liunian Harold Li, Mark Yatskar, D Yin, CJ Hsieh, and KW Chang. 2019. Visualbert: A simple and performant baseline for vision and language. <i>arXiv preprint arXiv:1908.03557</i> , 10.	710
657			711
658			712
659			713
660			714
661			715
662			716
663	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	Hongzhan Lin, Ziyang Luo, Wei Gao, Jing Ma, Bo Wang, and Ruichao Yang. 2024. Towards explainable harmful meme detection through multimodal debate between large language models. In <i>Proceedings of the ACM Web Conference 2024</i> , pages 2359–2370.	717
664			718
665			719
666			720
667			721
668	Liam Hebert, Gaurav Sahu, Yuxuan Guo, Nanda Kishore Sreenivas, Lukasz Golab, and Robin Cohen. 2024. Multi-modal discussion transformer: Integrating text, images and graph transformers to detect hate speech on social media. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 22096–22104.	Xuanrui Lin, Chao Jia, Junhui Ji, Hui Han, and Usman Naseem. 2025. Ask, acquire, understand: A multimodal agent-based framework for social abuse detection in memes. In <i>Proceedings of the ACM on Web Conference 2025</i> , pages 4734–4744.	722
669			723
670			724
671			725
672			726
673			727
674			728
675	Ming Shan Hee and Roy Ka-Wei Lee. 2025. Demystifying hateful content: Leveraging large multimodal models for hateful meme detection with explainable decisions. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 19, pages 774–785.	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36:34892–34916.	729
676			730
677			731
678			732
679			733
680			734
681	Paulo Cezar de Q Hermida and Eulanda M dos Santos. 2023. Detecting hate speech in memes: a review. <i>Artificial Intelligence Review</i> , 56(11):12833–12851.	Junxi Liu, Yanyan Feng, Jiehai Chen, Yun Xue, and Fenghuan Li. 2024a. Prompt-enhanced network for hateful meme classification. <i>arXiv preprint arXiv:2411.07527</i> .	735
682			736
683			737
			738

739	Kangzheng Liu, Feng Zhao, Yu Yang, and Guandong Xu. 2024b. Dysarl: Dynamic structure-aware representation learning for multimodal knowledge graph reasoning. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 8247–8256.	794
740		795
741		796
742		797
743		798
744	Yifan Liu, Yaokun Liu, Zelin Li, Ruichen Yao, Yang Zhang, and Dong Wang. 2025. Modality interactive mixture-of-experts for fake news detection. In <i>Proceedings of the ACM on Web Conference 2025</i> , pages 5139–5150.	799
745		800
746		
747		801
748		802
749	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	803
750		804
751		805
752		
753		806
754	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. <i>Advances in neural information processing systems</i> , 32.	807
755		808
756		809
757		810
758		811
759	Junyu Lu, Bo Xu, Xiaokun Zhang, Hongbo Wang, Hao-hao Zhu, Dongyu Zhang, Liang Yang, and Hongfei Lin. 2024. Towards comprehensive detection of chinese harmful memes. <i>Advances in Neural Information Processing Systems</i> , 37:13302–13320.	812
760		813
761		814
762		815
763	Paolo Mandica, Luca Franco, Konstantinos Kallidromitis, Suzanne Petryk, and Fabio Galasso. 2024. Hyperbolic learning with multimodal large language models. In <i>European Conference on Computer Vision</i> , pages 382–398. Springer.	816
764		817
765		
766		818
767		819
768	Jingbiao Mei, Jinghong Chen, Weizhe Lin, Bill Byrne, and Marcus Tomalin. 2024. Improving hateful meme detection through retrieval-guided contrastive learning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5333–5347.	820
769		821
770		822
771		823
772		
773		824
774	Gabriel Moreira, Manuel Marques, João Paulo Costeira, and Alexander Hauptmann. 2024. Hyperbolic vs euclidean embeddings in few-shot learning: Two sides of the same coin. In <i>Proceedings of the IEEE/CVF Winter conference on applications of computer vision</i> , pages 2082–2090.	825
775		826
776		827
777		
778		828
779		829
780	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. <i>arXiv preprint arXiv:1807.03748</i> .	830
781		831
782		832
783	Avik Pal, Max van Spengler, Guido Maria D’Amely di Melendugno, Alessandro Flaborea, Fabio Galasso, and Pascal Mettes. 2024. Compositional entailment learning for hyperbolic vision-language models. <i>arXiv preprint arXiv:2410.06912</i> .	833
784		834
785		835
786		836
787		837
788		838
789	Someen Park, Jaehoon Kim, Seungwan Jin, Sohyun Park, and Kyungsik Han. 2024. Predict: multi-agent-based debate simulation for generalized hate speech detection. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 20963–20987.	839
790		840
791		841
792		842
793		843
		844
		845
		846
		847
		848
		849
	Zelin Peng, Zhengqin Xu, Zhilin Zeng, Changsong Wen, Yu Huang, Menglin Yang, Feilong Tang, and Wei Shen. 2025. Understanding fine-tuning clip for open-vocabulary semantic segmentation in hyperbolic space. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 4562–4572.	
	Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021. Detecting harmful memes and their targets. <i>arXiv preprint arXiv:2110.00413</i> .	
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	
	Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. Flava: A foundational language and vision alignment model. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 15638–15650.	
	Ilya Stepanov, Junaid Rashid, Jong Weon Lee, Salman Naseem, and Imran Razzak. 2025. Cbrcl: A clip-bert with retrieval-guided contrastive learning multimodal approach for crisis-driven hate speech detection. In <i>Companion Proceedings of the ACM on Web Conference 2025</i> , pages 1993–1999.	
	Abraham A Ungar. 2001. Hyperbolic trigonometry and its application in the poincaré ball model of hyperbolic geometry. <i>Computers & Mathematics with Applications</i> , 41(1-2):135–147.	
	Hans Van Ditmarsch, Wiebe van Der Hoek, and Barteld Kooi. 2008. <i>Dynamic epistemic logic</i> . Springer.	
	Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2015. Order-embeddings of images and language. <i>arXiv preprint arXiv:1511.06361</i> .	
	Bin Wang, Fuyong Xu, Peiyu Liu, and Zhenfang Zhu. 2024. Hypermr: Hyperbolic hypergraph multi-hop reasoning for knowledge-based visual question answering. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 8505–8515.	
	Bo Xu, Erchen Yu, Jiahui Zhou, Hongfei Lin, and Linlin Zong. 2025. Hyperhateprompt: A hypergraph-based prompting fusion model for multimodal hate detection. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 3825–3835.	
	Chuanpeng Yang, Fuqing Zhu, Yaxin Liu, Jizhong Han, and Songlin Hu. 2024. Uncertainty-aware cross-modal alignment for hate speech detection. In <i>Proceedings of the 2024 Joint International Conference</i>	

- 850 *on Computational Linguistics, Language Resources*
851 *and Evaluation (LREC-COLING 2024)*, pages 16973–
852 16983.
- 853 Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen,
854 Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer,
855 Chao Huang, Pin-Yu Chen, and 1 others. 2024. Jus-
856 tice or prejudice? quantifying biases in llm-as-a-
857 judge. *arXiv preprint arXiv:2410.02736*.
- 858 Xianbing Zhao, Yixin Chen, Wanting Li, Lei Gao, and
859 Buzhou Tang. 2022. Mag+: An extended multimodal
860 adaptation gate for multimodal sentiment analysis.
861 In *ICASSP 2022-2022 IEEE International Confer-*
862 *ence on Acoustics, Speech and Signal Processing*
863 *(ICASSP)*, pages 4753–4757. IEEE.
- 864 Changmeng Zheng, Dayong Liang, Wengyu Zhang,
865 Xiao-Yong Wei, Tat-Seng Chua, and Qing Li. 2024.
866 A picture is worth a graph: A blueprint debate
867 paradigm for multimodal reasoning. In *Proceedings*
868 *of the 32nd ACM International Conference on Multi-*
869 *media*, pages 419–428.
- 870 Hengyang Zhou, Jinwu Yan, Yaqing Chen, Rongman
871 Hong, Wenbo Zuo, and Keyan Jin. 2025. Ldgnet:
872 Llm debate-guided network for multimodal sarcasm
873 detection. In *ICASSP 2025-2025 IEEE International*
874 *Conference on Acoustics, Speech and Signal Process-*
875 *ing (ICASSP)*, pages 1–5. IEEE.

A Related Work

A.1 Detection of Harmful Memes

Harmful memes have proliferated extensively across social media platforms, inflicting harm on vulnerable populations and contributing to increased social fragmentation to some degree. To address the challenge of effective detection, (Stepanov et al., 2025) employs retrieval-guided contrastive learning to improve hate speech identification. (Yang et al., 2024) introduces a cross-modal alignment framework to model multimodal error alignment and uncertainty perception. (Chhabra and Vishwakarma, 2023) utilizes a multi-scale adaptive receptive field to emphasize salient spatial regions, while (Xu et al., 2025) leverages hypergraphs to capture the hateful content arising from cross-modal information. Additionally, (Lin et al., 2024) implements modular networks for the detection of hate memes. Although these methods primarily focused on directly training classifiers for hate meme detection, they often overlooked the supplementary role of implicit information embedded within the memes. In contrast, studies such as (Ji et al., 2024; Park et al., 2024) directly apply VLMs to extract indicative features from memes. Furthermore, (Lin et al., 2024; Huang et al., 2025; Hee and Lee, 2025; Liu et al., 2024a) propose that VLMs be employed to extract implicit information from memes, which is then integrated into detection models during training. Distinct from the aforementioned approaches, our method introduces a novel concept termed “Viewpoint Decoupling”. This approach has the potential to mitigate the inherent biases of models toward harmful content to a certain extent.

A.2 Hyperbolic Space

Hyperbolic space can effectively model data with potential hierarchical structures, thereby enhancing the generalization ability of models. (Liu et al., 2024b) utilizes hyperbolic space to solve the task of completing multimodal knowledge graphs and designs a dual-space multi-hop structure learning module. (Mandica et al., 2024) proposes a BLIP-2 hyperbolic version training strategy. (Moreira et al., 2024) demonstrates that hyperbolic embeddings achieve the best few-shot classification performance. (Peng et al., 2025) proposes HyperCLIP, which fine-tunes text embeddings by adjusting their hyperbolic radius through scaling transformations. (Desai et al., 2023; Pal et al., 2024) prove the ex-

Harmful Prompt
Given an image and its corresponding title, please explain, from the perspective of hatred or harm, why this image and caption have been labeled as hateful or harmful. **Please note that you need to provide an explanation from a hateful and harmful perspective, and keep it within no more than two sentences.** Here is the title:{title} Your answer:
Non-harmful Prompt
Given an image and its corresponding title, please explain, from a non-hateful or harmless perspective, why this image and caption have been labeled as non-hateful or harmless. **Please note that you need to provide the explanation from a non-hateful and harmless perspective, and keep it within no more than two sentences.** Here is the title:{title} Your answer:

Figure 8: A prompt used to generate different perspectives.

istence of a hierarchical structure between image-text pairs. We posit the existence of a hierarchical relationship between multimodal information and interpretative perspectives, and draw upon the aforementioned concepts to model this hierarchical structure within hyperbolic space.

B Viewpoint Generation Template

To facilitate the model’s ability to produce viewpoints from two distinct perspectives, we meticulously crafted a prompt that includes instructions representing both a harmful and a harmless viewpoint, as illustrated in Figure 8. It is important to note that the generated viewpoints were restricted to a maximum of two sentences, taking into account the computational constraints of the visual language model.

C Loss Calculation

To enhance the discrimination between positive and negative samples within our framework, we incorporate a contrastive learning loss. Let the total number of samples be denoted by n , and suppose that among these, there are m positive samples sharing the same label. The feature representations of all samples are given by the set $\{f_1, f_2, \dots, f_n\}$. Under these conditions, the following formulation applies:

$$\mathcal{L}_{InfoICE} = \frac{1}{n} \sum_{i=1}^n \frac{1}{m} \sum_{j=1}^m \log \frac{\exp(s_{i,j})}{\sum_{k:k \neq i}^n \exp(s_{i,k})} \quad (16)$$

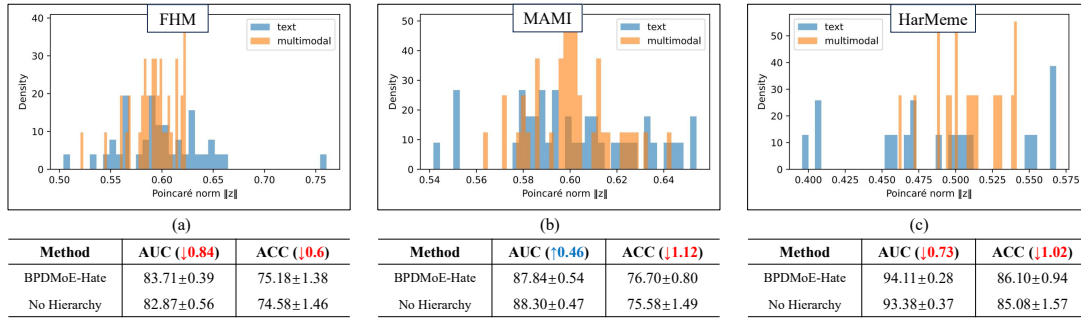


Figure 9: Eliminate the hierarchical relationship between viewpoints and multimodal features in hyperbolic space.

Here, $s_{i,j} = f_i^T f_j / \tau$, where $\tau = 0.07$ is the temperature coefficient. Moreover, to ensure the balanced loading of experts, our load balancing loss is as follows:

$$\mathcal{L}_{balance} = \frac{1}{n} \sum_{i=1}^n (w_{exp}^i - \bar{w}_{exp})^2 + \frac{1}{n} \sum_{i=1}^n (c_{exp}^i - \bar{c}_{exp})^2 \quad (17)$$

Where n represents the number of experts, w_{exp}^i represents the weight assigned to the i th expert, and c_{exp}^i represents the loading frequency of the i th expert in a sampling set.

Dataset	FHM	MAMI	HarMeme
Train	8,500	10,000	3,013
Valid	500	100	177
Test	1,000	1,000	354

Table 3: Dataset Distribution

D Detailed Supplement of the Dataset and Baselines

We conducted an evaluation of BPDME-Hate utilizing three extensively recognized hate meme datasets. The FHM dataset comprises a diverse collection of harmful memes sourced from the internet, consisting of 8,500 samples designated for training and 1,000 samples reserved for testing. The MAMI dataset focuses specifically on misogynistic memes gathered from prominent social media platforms. The HarMeme dataset includes harmful memes obtained from social media websites as well as through crowdsourcing efforts, notably encompassing a substantial subset of memes associated with the 2019 novel coronavirus. The size of the relevant dataset is shown in table 3.

The models we used for comparison include: 1) VLMs, such as Qwen2.5-vl-Instruct-32B (Bai et al., 2025), Llama-3.2-11B-Vision (Grattafiori et al., 2024) and Llava-1.5 (Liu et al., 2023). 2) Pure text classification models, including Bert-base (Devlin et al., 2019) and RoBERTa-large (Liu et al., 2019). 3) Multi-modal classification models, including FLAVA-full (Singh et al., 2022), Visual-BERT (Li et al., 2019), ViLBERT (Lu et al., 2019), BLIP2 (Li et al., 2023) and ALBEF (Li et al., 2021). 4) Harmful meme detection frameworks, including Mod-HATE (Cao et al., 2024), PromptHate (Cao et al., 2022), Pro-Cap (Cao et al., 2023), Explain-HM (Lin et al., 2024) and IntMeme (Hee and Lee, 2025).

E Additional Experiments

E.1 The Influence of the Viewpoint Encoder

For the encoding of viewpoints, we employed the more advanced RoBERTa-large model as the text encoder. This section examines the influence of various viewpoint encoders on the overall model performance, as shown in Table 4. Specifically, we substituted the original encoder with four alternative models: RoBERTa-base, T5-base (Raffel et al., 2020), BERT-base, and BERT-large. Experimental results indicate that an increase in the number of model parameters does not necessarily correspond to improved performance in meme detection. Rather, the effectiveness appears to depend on the intrinsic capabilities of the encoder itself. We hypothesize that this outcome may be attributed to the extensive freezing of parameters during the training phase. Consequently, future work should focus on selecting a more effective text encoder to further enhance the framework’s performance.

View Encoder	FHM		MAMI		HarMeme	
	AUC	ACC	AUC	ACC	AUC	ACC
Roberta-large	83.71 \pm 0.39	75.18 \pm 1.38	87.84 \pm 0.54	76.70 \pm 0.80	94.11 \pm 0.28	86.10 \pm 0.94
Roberta-base	83.89 \pm 0.49	75.06 \pm 1.20	88.91 \pm 0.48	75.40 \pm 1.06	93.85 \pm 0.43	86.05 \pm 1.48
T5-base	83.11 \pm 0.74	73.92 \pm 0.72	88.21 \pm 0.36	75.02 \pm 1.53	93.23 \pm 0.32	84.40 \pm 1.75
Bert-large	83.14 \pm 0.92	73.28 \pm 2.12	88.63 \pm 0.32	76.90 \pm 2.21	92.73 \pm 0.26	84.92 \pm 1.14
Bert-base	83.49 \pm 0.58	73.72 \pm 0.76	89.05 \pm 0.52	76.46 \pm 0.97	93.78 \pm 0.30	85.54 \pm 0.70

Table 4: The impact of the viewpoint encoder (text encoder) on the model’s performance.

Multimodal Encoder	FHM		MAMI		HarMeme	
	AUC	ACC	AUC	ACC	AUC	ACC
BLIP2	83.71 \pm 0.39	75.18 \pm 1.38	87.84 \pm 0.54	76.70 \pm 0.80	94.11 \pm 0.28	86.10 \pm 0.94
ViT	80.83 \pm 0.37	71.94 \pm 0.86	85.98 \pm 0.42	74.20 \pm 0.63	85.95 \pm 1.49	79.15 \pm 1.12
Flava-full	83.34 \pm 0.67	74.50 \pm 0.87	87.40 \pm 0.70	75.74 \pm 1.51	89.62 \pm 0.77	83.33 \pm 0.67

Table 5: The impact of the multimodal encoder on the model’s performance.

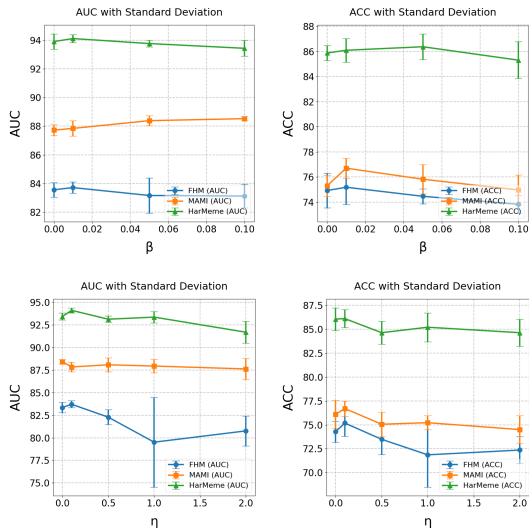


Figure 10: Compare the influence of the coefficients of the contrastive loss and the entailment loss. “ β ” represents the coefficient of the entailment loss, and “ η ” indicates the coefficient of the contrastive loss.

E.2 The Influence of the Multimodal Encoder

The performance of the multimodal encoder directly influences the feature representations during the hierarchical fusion stage. Since both images and their corresponding captions are simultaneously input into the multimodal encoder, the features output by this encoder have a greater impact on the BPDME-Hate prediction. As shown in Table 5, compared to Table 1, employing the lower-performing Flava-full model results in a significant decline in overall model performance. This finding underscores the critical role of multimodal information as the “root” within the hierarchical struc-

ture; inadequate representation of this information severely impairs hierarchical modeling. Furthermore, we replaced the multimodal encoder with a ViT (Dosovitskiy, 2020) model capable of encoding only images to investigate whether a more effective hierarchical structure could be formed between “image” and “viewpoint”. Experimental results indicate that this hierarchical structure contributes less to our model than the “multimodal-viewpoint” structure, thereby further validating the feasibility of the proposed hierarchical framework.

E.3 Different Viewpoint Generation Models

Our BPDME-Hate framework employs the advanced Qwen2.5-VL-32B-Instruct model for viewpoint generation. In this section, we investigate the influence of utilizing various VLMs to produce binary viewpoints on the overall system performance. The evaluation was conducted on the FHM and HarMeme datasets, with the corresponding results presented in table 6. Our findings indicate that VLMs possessing stronger self-inference capabilities exert a more beneficial impact on the framework. Conversely, VLMs with comparatively weaker performance, constrained by limited internal knowledge, tend to generate binary viewpoints with reduced informational content, thereby impeding the model’s evaluative accuracy. Consequently, we infer that binary viewpoints derived from models with enhanced reasoning abilities more accurately capture the authentic expression of memes, leading to improved judgment within the framework.

VLM	FHM		HarMeme	
	AUC	ACC	AUC	ACC
Qwen2.5-VL-32B-Instruct	83.71 \pm 0.39	75.18 \pm 1.38	94.11 \pm 0.28	86.10 \pm 0.94
Gemma3-12B	84.37 \pm 0.55	75.50 \pm 1.28	93.88 \pm 0.59	83.95 \pm 1.84
Qwen3-VL-8B	82.65 \pm 0.59	72.44 \pm 1.14	94.16 \pm 0.29	84.41 \pm 1.20
Qwen2.5-VL-7B-Instruct	80.48 \pm 0.41	69.66 \pm 2.65	92.74 \pm 0.51	84.01 \pm 1.91
Qwen2-VL-2B-Instruct	74.47 \pm 0.34	64.82 \pm 1.98	92.59 \pm 1.22	83.90 \pm 1.83

Table 6: The influence of different viewpoint generation models on the results.

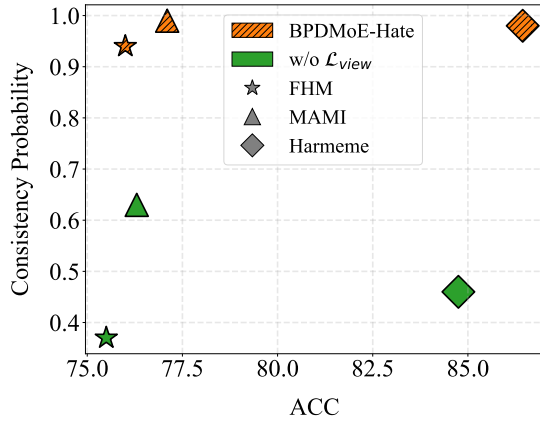


Figure 11: Removing the viewpoint selection loss. The term ‘‘Consistency Probability’’ denotes the probability that the model’s predicted viewpoint selection aligns with the harmful meme prediction within the test dataset. We fix the random seed for the verification.

E.4 The Proportion of Different Losses

This section examines the effects of the entailment loss and the contrastive learning loss on model performance. Since the viewpoint selection loss and the task loss are components of the same optimization objective, and no discernible pattern was observed regarding the impact of the corresponding loss coefficient on the outcomes during experimentation, further discussion on this aspect is omitted. The experimental results are presented in Figure 10. Our findings indicate that optimal performance was achieved at parameter values of $\beta = 0.01$ and $\eta = 0.1$. Although increasing β led to a slight improvement in the AUC metric for FHM, the ACC of HarMeme peaked at $\eta = 0.05$. Furthermore, we hypothesize that the proportion of entailment loss should not be excessively large because it is applied exclusively within the hyperbolic space. An overly high weighting of this loss may inhibit the model’s ability to select experts effectively in the hyperbolic space, thereby diminishing overall performance.

E.5 Consistency Between Viewpoint Selection and Prediction

In this section, we remove the viewpoint selection loss to examine the consistency between the model’s viewpoint selection predictions and its harmful meme predictions. The results are presented in the figure 11. We observe that after training, BPDME-Hate demonstrates a high degree of alignment between viewpoint selection and the prediction of whether a meme is harmful, underscoring the significance of the viewpoint selection loss for our model. Upon removal of this loss, the ‘‘Consistency Probability’’ experiences a substantial decline across all three datasets, and the ACC of the model is inferior to that of BPDME-Hate. This indicates that accurate viewpoint selection positively influences model performance, whereas incorrect viewpoint choices tend to mislead harmful meme predictions. These findings further validate the efficacy of our model in mitigating bias.

Type	Parameter Size
Trainable params	138 M
Non-trainable params	4.1 B
Total Params	4.2 B

Table 7: The total number of parameters.

E.6 Eliminate Hierarchical Relationships

In this section, we verify the impact of eliminating the hierarchical relationship in hyperbolic space on the model’s performance. We artificially bring the norm values of the projected viewpoint features and multimodal features in hyperbolic space closer and observe the experimental results as shown in Figure 9. We find that the absence of hierarchical relationship constraints leads to a certain decline in model performance. Moreover, the performance metrics in Figure 6 decline more significantly compared to the elimination of hierarchical relationships, demonstrating the importance of correctly

setting the hierarchical relationship order.

Hyperbolic Distance	FHM	MAMI	HarMeme
Multimodal	0.63	0.60	0.67
Viewpoint	7.82	8.57	7.17
Difference Value	7.19	7.97	6.50

Table 8: The average hyperbolic distance of different features from the origin. Here, we fixed the random seed for the test set.

E.7 Quantitative Analysis of Hierarchical Structure

The preceding observations have been intuitively illustrated through the distribution histograms of the norms of various vectors, indicating the existence of a hierarchical relationship between multimodal features and viewpoint features. Furthermore, multimodal features represent a more generalized form compared to viewpoint features. In this section, we quantitatively validate this hierarchical relationship by analyzing the average hyperbolic distance from the origin for both multimodal and viewpoint features within different test sets. The results, presented in table 8, demonstrate that across different test sets, multimodal features consistently exhibit shorter hyperbolic distances to the origin, whereas the selected viewpoint features are positioned farther away. These findings substantiate that the proposed model effectively captures the hierarchical structure inherent between these two feature types. It is observed that the range of values for the hyperbolic distance is $[0, +\infty)$.

Dataset	BPDMoE-Hate	w/o VEM
FHM	75.18 \pm 1.38	71.74 \pm 0.72
MAMI	76.70 \pm 0.80	74.04 \pm 0.55
HarMeme	86.10 \pm 0.94	83.05 \pm 0.89

Table 9: The significance of the viewpoint enhancement module. The evaluation metric used is ACC.

E.8 Removal of The Viewpoint Enhancement Module

The design objective of the VEM is to comprehensively incorporate the essential characteristics of contrasting perspectives and multimodal data into the target viewpoint, thereby better capturing the harmful information hidden in memes. Consequently, this module constitutes a critical component of our overall framework. To evaluate its

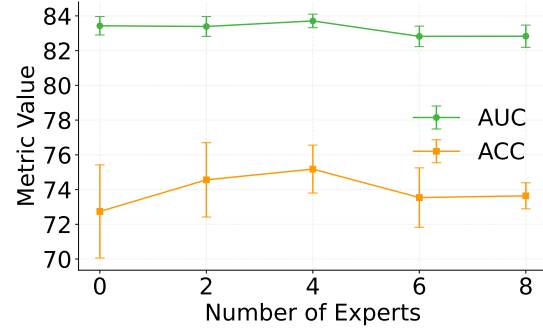


Figure 12: The impact of the number of experts. This assessment is conducted using the FHM dataset.

significance, we conducted an ablation study by removing the VEM and assessing the resultant performance variations. The findings, as presented in table 9, indicate a performance decline following the module’s removal. Notably, the VEM exerted the most pronounced effect on the FHM dataset and the least on the MAMI dataset. This observation further substantiates our assertion that the MAMI dataset, which focuses on misogyny, possesses relatively straightforward discriminative features, wherein the integration of images and titles sufficiently conveys explicit information.

E.9 Number of Experts

The number of selected experts is fixed at 2, while the total number of experts is varied. The corresponding results are presented in the Figure 12. It is important to note that a total expert count of zero indicates the direct removal of experts from the dual space. The findings demonstrate that the performance of BPDMoE-Hate is influenced by the total number of experts. Specifically, optimal performance is observed when the number of experts was set to 4. Consequently, the total number of experts is established at 4 for subsequent experiments.

F The Number of Parameters

The total number of parameters within the proposed framework was computed and is presented in table 7. It is important to note that the parameters associated with the VLM were excluded from this calculation. The framework comprises 138 million trainable parameters, whereas the majority of non-trainable parameters originate from the multimodal encoder component. Substituting the multimodal encoder with a more lightweight alternative, such as FLAVA-full (as shown in table 5), enables de-

1181 ployment on a wider range of devices; however,
1182 this modification will incur a slight reduction in
1183 performance.