HVGuard: Utilizing Multimodal Large Language Models for Hateful Video Detection

Anonymous ACL submission

Abstract

The rapid growth of video platforms has transformed information dissemination and led to an explosion of multimedia content. However, this widespread reach also introduces risks, as some users exploit these platforms to spread hate speech, which is often concealed through complex rhetoric, making hateful video detection a critical challenge. Existing detection methods rely heavily on unimodal analysis or simple feature fusion, struggling to capture cross-modal interactions and reason through 011 implicit hate in sarcasm and metaphor. To address these limitations, we propose HVGuard, the first reasoning-based hateful video detection framework with multimodal large language models (MLLMs). Our approach integrates Chain-of-Thought (CoT) reasoning to enhance multimodal interaction modeling and implicit 019 hate interpretation. Additionally, we design a Mixture-of-Experts (MoE) network for efficient multimodal fusion and final decisionmaking. The framework is modular and extensible, allowing flexible integration of different MLLMs and encoders. Experimental results demonstrate that HVGuard outperforms all existing advanced detection tools, achieving an improvement of 6.88% to 13.13% in accuracy and 9.21% to 34.37% in M-F1 on two public datasets covering both English and Chinese.

> Disclaimer: This paper contains harmful content, which has the potential to be offensive and may disturb readers.

1 Introduction

037

038

041

In recent years, video platforms like YouTube (Google, 2005), Bilibili (Kuanyu, 2009), and Tik-Tok (ByteDance, 2016) have transformed information dissemination and fueled multimedia growth. However, this also brings risks, as some users exploit these platforms to spread false information, extremist content, and hate speech (Ottoni et al., 2018). Hate speech, which demeans, attacks, or



Figure 1: A typical example of hateful video

042

043

044

045

051

052

055

060

061

062

063

marginalizes individuals or groups based on characteristics like race, religion, or gender (Hee et al., 2024b; Fortuna and Nunes, 2018). It may not only incite social conflicts but also cause real-world harm to individuals and groups. Thus, effectively detecting hate speech on video platforms (Alcântara et al., 2020; Das et al., 2023; Wu and Bhandary, 2020) has become an urgent challenge.

Compared with traditional text-based forms, the spread of hate speech in videos is more concealed and has a broader impact. Since video content typically includes multimodal information including text, audio, and visual elements, hate message is often embedded in a more subtle manner, making it difficult for single-modality detection methods to identify effectively. Figure 1 shows an offender joking with a bald victim: "Do you know why the man put a rabbit on his head?" "Because he wanted a hare on his head!" This uses the phonetic similarity between "hair" and "hare" to offend bald individuals. It highlights a challenge in detecting subtle hate speech and the need for inference. Fur-

102 104

105 107

108 109 110

111

112

113

114

115

thermore, content that appears harmless in a single modality may reveal its offensive nature when visual, auditory, and contextual cues are considered together. Effective hateful video detection thus requires an integrated understanding of multimodal interactions and rhetorical devices such as metaphor and wordplay.

Current hateful video detection methods mainly use modality encoders to extract features and then classify them (Yu et al., 2022; Wu and Bhandary, 2020; Wang et al., 2024; Das et al., 2023). However, these methods either use single modality or simply concatenate features from multiple modalities which is limited because it does not take into account the interaction between different modalities. At the same time, hateful videos often involves rhetorical devices such as metaphors, irony, and sarcasm (Xu et al., 2024; Ge et al., 2023), which cannot be addressed by simple modality feature extraction methods without some form of reasoning (Prystawski et al., 2022). Moreover, online hateful videos are increasing rapidly and are often related to specific cultural contexts (Ottoni et al., 2018). This requires the integration of rich world knowledge to enhance reasoning capabilities to address this issue. Therefore, research on hateful video detection has important practical significance and can provide more precise technical support for the content moderation mechanisms of social platforms.

Recent advancements in multimodal large language models (MLLMs) (Bai et al., 2023; Team et al., 2024; Liu et al., 2024; Wang et al., 2023) have demonstrated strong video understanding capabilities by leveraging extensive world knowledge and deep semantic comprehension (Tang et al., 2023). This makes them promising for tackling the challenges of hateful video detection. To enhance their effectiveness in this domain, we incorporate Chain-of-Thought (CoT) reasoning, which enables MLLMs to break down complex reasoning tasks into intermediate steps. This structured approach allows for a more systematic analysis of multimodal information—spanning audio, visual, and textual components-while capturing their interactions to form a coherent understanding of the video's overall semantics.

In this work, we first explore the effectiveness of MLLMs and CoT reasoning in understanding hateful videos, particularly their role in handling multimodal interactions and rhetorical devices such as metaphors. Building on these insights, we propose the first reasoning-based hateful video detection framework, HVGUARD¹. Our approach leverages MLLMs to generate multimodal rationales and incorporates a CoT strategy that explicitly models cross-modal interactions and rhetorical elements, addressing the challenges of implicit hate detection. Additionally, we design a Mixture-of-Experts (MoE) network (Jacobs et al., 1991) to effectively integrate diverse multimodal information. The MoE model fuses multimodal representations with MLLM-derived rationales, optimizing the decision-making process. This integration enables our framework to combine low-level feature extraction with high-level semantic reasoning, ultimately improving the accuracy and robustness of hateful video detection. Experimental results demonstrate that HVGUARD achieves a detection accuracy of up to 0.86, outperforming existing stateof-the-art methods.

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

153

154

155

156

158

159

160

161

162

163

The key contributions of this paper are as follows:

- First Exploration of MLLMs and CoT in Hateful Video Understanding. This is the first work to explore the potential of MLLMs and CoT reasoning for hateful video understanding, demonstrating their effectiveness in managing multimodal interactions and complex rhetorical devices, such as metaphors.
- · Novel Reasoning-Based Hateful Video Detection Framework. We propose the first reasoning-based hateful video detection framework, integrating MLLMs with CoT reasoning to enhance multimodal interaction modeling and implicit hate interpretation. Additionally, we introduce a MoE network to efficiently fuse multimodal representations and MLLM-generated rationales, optimizing the decision-making process.
- Extensive Evaluation of HVGUARD. Experimental results show that HVGuard achieves up to 0.86 accuracy, outperforming all existing detection tools with accuracy gains of 6.88% to 13.13% and M-F1 improvements of 9.21% to 34.37%. Extensive experiments on two public datasets, covering both English and Chinese, further validate its effectiveness in binary and ternary classification settings against five state-of-the-art baselines, including advanced MLLMs and existing detection tools.

¹We will open-source our framework for future research.



Figure 2: Visualization of features used by different methods. (a) Embedding of video titles, transcripts. (b) Embedding of MLLM rationale. (c) Embedding after incorporating the CoT prompts.

2 Preliminary Study

With the advancement of artificial intelligence, MLLMs have become the focal point of the latest developments. The complementarity of LLMs and VLMs has given rise to MLLMs, such as Gemini 1.5(Team et al., 2024) and GPT-4 series (Achiam et al., 2023). They can receive, reason, and output multi-modal information, showing impressive capabilities in various multi-modal tasks, including image reasoning and video understanding (Wu et al., 2023; Fu et al., 2024), thus opening up new ways to solve complex and novel challenges in the multi-modal field.

To more clearly demonstrate how the reason-177 ing capability of MLLMs aids in understanding of hateful content in videos, we conducted a visual analysis of embedding representations on the hateful video dataset Multihateclip (Wang et al., 2024). 181 Figure 2a visualizes the embeddings of pure textual 182 information (video title and transcript) extracted 183 using the pre-trained text encoder BERT (Devlin, 2018), which exhibit significant overlap with no 185 discernible class separability. This indicates the insufficiency of traditional approaches with single modality. However, when analyzing videos with 188 MLLMs (Figure 2b), a certain degree of class separability becomes observable. By further incorporat-190 ing the CoT prompting strategy (detailed in Section 191 3.4), we guide the MLLM to clarify rhetorical devices such as metaphors and puns in the videos, 194 ultimately achieving sharper classification boundaries (Figure 2c). Thus, MLLMs provide effective rationale for hateful video understanding, and the 196 CoT prompting strategy further amplifies this capability. 198

3 Method

3.1 Task Definition

The goal of hateful video detection is to extract features from videos and classify them based on these features. The video dataset is represented as $\mathcal{V} = \{v_1, \dots, v_i, \dots, v_{|\mathcal{V}|}\}$, where $|\mathcal{V}|$ is the number of videos. The task can be expressed as:

$$\arg \max_{c \in \{1,2,...,|C|\}} P(c|v_i)$$
(1)

200

201

202

203

204

206

208

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

227

228

where $c \in \{1, 2, ..., |C|\}$ represents the classification categories. Our work focuses on utilizing rationale generated by MLLM and multimodal information from the video itself for hateful video detection. Therefore, this task can be re-expressed as:

$$\arg\max_{c\in\{1,2,\dots,|C|\}} P(c|v_i^T, v_i^A, v_i^F, v_i^M) \quad (2)$$

where v_i^T represents the text information in the video (such as title, subtitles, or transcript), v_i^A represents the audio information of the video, v_i^F represents the frame information of the video, and v_i^M represents MLLM-derived rationales.

3.2 Overview

The overview of our framework, HVGUARD, is shown in Figure 3. Based on preliminary study, we design this novel framework for hateful video detection, leveraging MLLM-derived rationales to address challenges in multimodal interaction and the interpretation of metaphors and rhetorical devices. This framework extracts text, audio, and video frames from the input video, providing a comprehensive semantic representation of the video. A

165

166

168

170

171

172

173

174

175



Figure 3: Overview of the proposed framework.

CoT-based reasoning approach is then applied, progressively reasoning through the individual modalities and their interactions, to generate rationale from MLLM. In the final stage, these embeddings are ultimately integrated using a MoE network to yield the final classification results.

3.3 Multimodal Extraction Module

230

231

233

235

236

240

241

242

243

247

248

249

251

254

257

259

Considering that hateful videos encompass multiple modalities, we first extract feature information from the three main modalities of the video: text, audio, and video frames. We process audio signal v_i^A as a combination of semantic information and emotional information. We use FunASR (Gao et al., 2023), an open-source audio processing tool, to transcribe the audio into transcript v_i^{trans} and extract the emotion of the spoken content v_i^{emo} . Subsequently, following the approach of Vivit (Arnab et al., 2021), the video is uniformly sampled into 32 frames, with a fixed interval between consecutive frames to ensure equal temporal spacing throughout the video.

$$v_i^A, v_i^F = extract(v_i), v_i^{trans}, v_i^{emo} = trans(v_i^A)$$
(3)

where v_i^A represents the original audio signal, v_i^F represents the video frames, and v_i^{title} represents the video title.

Next, we construct the textual content v_i^T using the video title and transcript:

$$v_i^T = \{v_i^{title}, v_i^{trans}\}$$
(4)

3.4 MLLM Reasoning Module

To address the challenges in hateful video detection, such as metaphors, cultural contexts, and the complexity of multimodal interactions, it is necessary to leverage MLLMs to extract deep semantic information from the video. Based on preliminary study (Chapter 2), we find that hateful video detection is a complex process, requiring the extraction of key cues from multiple modalities, including text, visuals, and audio. Inspired by the works of (Xu et al., 2024; Vishwamitra et al., 2024), we employ carefully designed CoT prompts to decompose this complex task, thereby enabling the understanding of multimodal hateful content within the video. Specifically, our CoT prompt is as follows:

Adaption Prompt. In the field of hateful content detection, domain alignment, role description and task-specific adaptation is critical, as it equips MLLMs with essential cultural context and contextual comprehension. This focuses the model's capabilities on addressing the specific challenges of understanding both nuanced and overt hateful content, thereby improves its performance and reliability (Csurka, 2017; Qi et al., 2024). We employ the prompt:

This is a video that may contain harmful content, such as hate speech, explicit violence, discrimination, or other forms of harmful behavior. You are a content moderation specialist. Your task is to identify any instances of hate speech, violent imagery, discriminatory actions, or any other content that could be considered harmful, abusive, or offensive. Ensure the answer's accuracy while keeping it concise and avoiding overexplanation.

Visual Meaning Understanding. To guide the model to analyze the video progressively, starting

260

288

289 290

291

20

294 295 296

297 298

299

300

303

309 310

311

312

313 314

315

318

with the visual information while ignoring the subtitles in the video frames. The focus is placed on analyzing the characters and scenes in the frames. To achieve this, we employ the following prompt:

> Describe the video content based on {video frames}, ignoring subtitles in the frames. Pay attention to any special characters or scenes.

Given the video frames v_i^F and this prompt X_{prompt}^F , the output computation is as follows:

$$res1 = MLLM(v_i^F, X_{prompt}^F)$$
(5)

Textual Meaning Understanding. We guide the model to focus on textual information by analyzing the video titles and transcripts, paying special attention to the presence of rhetorical devices such as puns and homophonic wordplay used as promotional strategies. Based on this, we employ the following prompt:

> The video title is {video title}. The text in the video is {video transcript}. Please analyze the meaning of this text. Note that there may be homophonic memes and puns; distinguish and explain them.

Given the textual input v_i^T and the prompt X_{prompt}^T , the output computation is as follows:

$$res2 = MLLM(v_i^T, X_{prompt}^T)$$
(6)

Fusion Meaning Understanding. Given the complex relationships between semantics across different modalities, it is essential to comprehensively consider the meaning conveyed by the video after multimodal fusion. As illustrated by figure 1, some videos may contain no obvious offensive content in their text or visuals individually, yet their combination can give rise to new meanings. Therefore, we aim for the model to synthesize the results from the first two steps and further integrate the video's raw information, including video frames, text, and extracted emotions of spoken content. This approach seeks to uncover deeper cross-modal interactions and analyze potential new metaphors. We employ the following prompt:

Please combine the {video title}, {video transcript}, {video frames}, {voice emotion}, {response1}, {response2} and analyze both the visual, textual and audio elements of the video to detect and flag any hateful content. No need to describe the content of the video, only answer implicit meanings and whether this video expresses hateful content further.

The MLLM rationale is as follows:

 $v_i^M = MLLM(v_i^T, v_i^F, v_i^{emo}, res1, res2) \quad (7)$

319

320

321

322

323

324

325

326

327

329

330

331

332

333

334

335

336

337

338

340

341

343

346

347

348

350

351

3.5 Multimodal Fusion Module

After obtaining rationale generated by MLLM reasoning module, we designed a multimodal fusion module to fuse information from the aforementioned modalities. We employ modality-specific encoders for each type of modality to obtain their respective embedding representations:

$$E_i^T = f_T(v_i^T),$$

$$E_i^A = f_A(v_i^A),$$

$$E_i^F = f_F(v_i^F)$$
(8)

where f_T , f_A , and f_F represent the text, audio, and vision modality encoders, while E_i^T , E_i^A , and E_i^F represent corresponding embeddings. To reduce the inference burden, we designed an embedding cache, allowing the above process to be executed only once on the dataset.

The rationale v_i^M generated by the MLLM is presented in textual form. We treat it as additional textual input and feed it into the text modality encoder to obtain embeddings:

$$E_i^M = f_T(v_i^M) \tag{9}$$

To fuse the embeddings from different modalities, we designed a mixture of experts network. First, all embeddings are concatenated into a single long vector as the representation embedding E_i for the entire video:

$$E_i = concat(E_i^T, E_i^A, E_i^F, E_i^M)$$
(10)

Next, we constructed n identical expert networks and one gating network, where n is the number of experts. These experts and the gating network share the same input E_i . Each expert network extracts high-level information specific to certain feature.

392

393

394

395

396

397

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

384

The output of the k-th expert is denoted as O_k and is computed as follows:

$$O_k = f_k(E_i; \theta_k), \quad k \in \{1, 2, \dots, n\}$$
 (11)

where f_k represents the mapping function of the *k*-th expert network, and θ_k denotes its parameters.

Simultaneously, the gating network $g(E_i; \phi)$ dynamically generates weights w_k to adjust the contribution of each expert's output. To prevent weight polarization, dropout is applied to the gating network's output weights. The gating network computes these weights as:

$$w_{k} = \text{Dropout}\left(\frac{\exp(g_{k}(E_{i};\phi))}{\sum_{j=1}^{n}\exp(g_{j}(E_{i};\phi))}\right), \quad (12)$$
$$k \in \{1, 2, \dots, n\}$$

where $g_k(E_i; \phi)$ is the unnormalized weight produced by the gating network, and ϕ represents the parameters of the gating network.

The final fused output O_{fusion} is obtained by combining the weighted outputs of all experts:

$$O_{fusion} = \sum_{k=1}^{n} w_k \cdot O_k \tag{13}$$

3.6 Final Decision

367

371

373

374

375

379

383

During training, we optimize the parameters of the expert and gating networks by minimizing a loss function. Assuming the ground truth labels are y and the final decision outputs are \hat{y} , we use a cross-entropy loss function:

$$\mathcal{L} = -\frac{1}{m} \sum_{i=1}^{m} \left[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right]$$
(14)

where m denotes the number of samples.

4 Experiments

Dataset	Language	Total	Η	0	Ν
HateMM	English	1,066	427	0	639
Multibataalin	English	891	72	218	601
Mutunatechp	Chinese	897	112	180	605

Table 1: Overview of datasets. H:hateful, O:offensive, N:normal

In this chapter, we first introduce the two datasets used to validate our method, the experimental setup, and the selection of baselines. Then, we present the experimental results and provide a detailed analysis.

4.1 Dataset

In our study, we employ two high-quality, up-todate public datasets for hateful video detection: the HateMM dataset and the MultiHateClip dataset.

HateMM(Das et al., 2023). The HateMM dataset consists of 1,083 videos sourced from BitChute, a platform with lenient content moderation, resulting in a higher prevalence of hateful content. Videos are labeled as either Hate or Non-Hate.

MultiHateClip(Wang et al., 2024). The MultiHateClip dataset is a multilingual benchmark dataset for hateful video detection, including 2,000 videos from YouTube and Bilibili, with 1,000 videos in English and 1,000 in Chinese. Each video is classified as Hateful, Offensive, or Normal.

To enhance data reliability, we filtered out corrupted and blurry videos. Additionally, to ensure high-quality textual information, we re-annotated the video transcripts using the speech transcription tool FunASR (Gao et al., 2023), improving the accuracy of multimodal analysis. The dataset we use is summarized in Table 1.

4.2 Experiment Settings

We randomly split all datasets into training, testing, and validation sets with a 7:2:1 ratio. For the ternary classification task on the MultiHateClip dataset, the labels used are Hateful, Offensive, and Normal. For binary classification on both the Multi-HateClip and HateMM datasets, we combine Hateful and Offensive into a single category, keeping the Normal label unchanged.

All models are trained with a learning rate of 1e-4, a batch size of 32, and early stopping after 100 epochs. Experiments are conducted on three Tesla V100-32G GPUs. Model performance is primarily evaluated using macro-averaged F1 score (M-F1) and accuracy (acc). We employ GPT-40(Achiam et al., 2023), XLM(Conneau, 2019), Vit(Dosovitskiy, 2020), and Wav2Vec(Baevski et al., 2020) as the fundamental MLLM and modality encoders.

4.3 Baseline Models

We evaluate HVGUARD with five baselines, including three advanced MLLMs and two state-of-theart methods in hateful video detection: (1) **GPT-40** (Achiam et al., 2023): An advanced MLLM by OpenAI, with high-level reasoning capabilities. (2) **Gemini-1.5-pro** (Team et al., 2024): A sophis-

Dataset	Number of categories	Model	Acc	M-F1	F1(H)	R(H)	P(H)	F1(O)	R(O)	P(O)
Makih analia (Tanliah)	3	GPT-40	0.7326	0.3280	0.2957	0.2361	0.3953	0.4923	0.4486	0.5455
		Gemini-1.5-pro	0.6319	0.4458	0.2143	0.2000	0.2308	0.3409	0.3488	0.3333
		Qwen-VL	0.5618	0.4060	0.2051	0.6154	0.1231	0.2258	0.1556	0.4118
		HateMM	0.6966	0.4894	0.1333	0.1667	0.1111	0.5217	0.5516	0.5345
		Multihateclip	0.7079	<u>0.4946</u>	0.1667	0.1667	0.1667	0.4928	<u>0.5780</u>	0.4750
		HVGuard	0.8090	0.6646	0.4556	0.4722	0.5000	0.6488	0.6270	0.6994
Multinatechp(Elighsii)		GPT-40	0.7989	0.5019	/	/	/	0.6455	0.5699	0.7443
		Gemini-1.5-pro	0.7198	0.6020	/	/	/	0.3855	0.2759	0.6400
	2	Qwen-VL	0.6573	0.6549	/	/	/	0.6258	0.9273	0.4722
		HateMM	0.7191	0.6646	/	/	/	0.5421	0.4722	0.6548
		Multihateclip	0.7416	0.6806	/	/	/	0.5544	0.4861	0.7269
		HVGuard	0.8539	0.7714	/	/	/	0.6308	<u>0.5819</u>	0.7619
	3	GPT-40	0.6444	0.4460	0.2326	0.1852	0.3125	0.2941	0.3448	0.2564
		Gemini-1.5-pro	0.6648	0.4393	0.2069	0.1500	0.3333	0.2985	0.2703	0.3333
		Qwen-VL	0.5719	0.4472	<u>0.3333</u>	0.6875	0.2200	0.2491	0.1889	0.3656
		HateMM	0.6889	0.4163	0.0741	0.0476	0.1667	0.3667	0.3889	0.4722
		Multihateclip	<u>0.7111</u>	<u>0.4573</u>	0.1667	0.1111	<u>0.3333</u>	<u>0.3778</u>	0.3889	0.4167
Multihateclin(Chinese)		HVGuard	0.8045	0.5643	0.3563	0.2917	0.5278	0.4417	0.4190	0.6139
infantinateenp(chinese)		GPT-40	0.7389	0.6900	/	/	/	0.5766	0.5714	0.5818
		Gemini-1.5-pro	0.7443	0.6188	/	/	/	0.4000	0.2632	0.8333
		Qwen-VL	0.6704	0.6684	/	/	/	0.6424	0.9298	0.4907
		HateMM	0.7444	0.6908	/	/	/	0.5694	0.5694	0.5826
		Multihateclip	<u>0.7778</u>	0.6904	/	/	/	0.5299	0.4028	<u>0.7833</u>
		HVGuard	0.8603	0.8219	/	/	/	0.7408	<u>0.6905</u>	0.8274
HateMM	2	GPT-40	0.7308	0.7306	0.7238	0.8806	0.6144	/	/	/
		Gemini-1.5-pro	<u>0.7874</u>	0.7872	<u>0.7933</u>	0.8554	0.7396	/	/	/
		Qwen-VL	0.7089	0.7089	0.7075	0.8824	0.5906	/	/	/
		HateMM	0.7500	0.7454	0.7430	0.7259	0.7614	/	/	/
		Multihateclip	0.7614	0.7594	0.7611	0.7537	<u>0.7690</u>	/	/	/
		HVGuard	0.8563	0.8597	0.8479	0.8228	0.8809	/	/	/

Table 2: Results of different methods on the task of hateful video detection. H:hateful, O:offensive, Acc:accuracy, M-F1:macroF1, R:recall, P:precision.

ticated multimodal model by Google DeepMind, capable of handling diverse reasoning tasks and understanding multiple modalities, including audio, images, videos, and text. (3) Qwen-VL-7B (Bai et al., 2023): An open-source vision-language model by Alibaba Cloud, excelling in tasks like image captioning, question answering, and visual localization. (4) HateMM (Das et al., 2023): A multimodal hateful video detection model that combines text, audio, and visual pretrained models through a trainable fusion layer to make final predictions. (5) MultiHateClip (Wang et al., 2024): A model that processes each modality's features through independent fully connected layers, concatenates them, and performs final classification to determine whether the video contains hate speech.

For the MLLMs used, we employ a generalized prompt to detect hateful videos: "*Analyze whether the video contains hateful content*." To ensure test consistency, we reproduced all the baselines and conducted a unified evaluation.

4.4 Evaluation Results

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

To verify the effectiveness of HVGUARD, experiments were conducted on the dataset shown in Table 2. The Multihateclip dataset includes both English and Chinese videos, which are used to evaluate the generalization ability of the detection tools in cross-lingual environments. The binary and ternary classification tasks aim to assess the performance of the detection tools in tasks with varying levels of granularity. In real-world scenarios, the binary classification task aids platforms in quickly identifying and blocking hateful videos, while the ternary classification task enables more precise content moderation. The additional "Offensive" category in the ternary task allows for further differentiation, thereby reducing false positives. 458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

Overall, HVGUARD outperformed all other baselines, with an improvement of 6.88% to 13.13% in accuracy and 9.21% to 34.37% in M-F1 compared to existing SOTA detection tools. We then explored further conclusions through the following analysis.

HVGUARD achieved SOTA performance on both English and Chinese hateful video datasets, demonstrating its multilingual adaptability. Additionally, it outperformed other baselines in both ternary and binary classification tasks.

We also achieved superior performance on most metrics for the crucial labels of "Hateful" and "Offensive," demonstrating the HVGUARD ability for hateful video detection. Notably, Qwen-VL

Model	Terr	nary	Binary		
Wouer	Acc	M-F1	Acc	M-F1	
w/o Vision encoder	0.7865	0.4760	0.8202	0.7397	
w/o Text encoder	0.7753	0.5633	0.8258	0.7090	
w/o Audio encoder	0.7697	0.5807	0.8258	0.7413	
w/o Modal features	0.7584	0.4816	0.8146	0.7126	
w/o CoT	0.7416	0.4715	0.7921	0.5512	
w/o MoE	0.7809	0.5936	0.8371	0.7466	
HVGuard	0.8090	0.6646	0.8539	0.7714	

Table 3: Ablation study for the components in HV-Guard.

achieved the highest recall rate for the "Hate" category, but performed poorly in accuracy and M-F1. This suggests that Qwen-VL tends to classify videos as "Hate", leading to the misclassification of some normal videos. In practical applications, an excessively high false positive rate may negatively impact the normal information flow within online communities.

To more clearly demonstrate the effectiveness of the proposed framework, we present a case study in Appendix A.

4.5 Effectiveness of Components in HVGuard

Table 3 summarizes the results of the ablation study on the MultiHateClip(English) dataset using HV-Guard. Removing the visual, text, or audio components individually resulted in performance declines, indicating that each modality plays a crucial role in hate detection. Furthermore, ablation of all modal features, relying solely on MLLM rationale—led to a noticeable decrease in performance. These findings underscore the importance of integrating comprehensive multimodal information for accurate detection.

Moreover, removing the CoT guidance for the MLLM and relying solely on generalized prompt templates resulted in a significant performance drop. This demonstrates that the CoT approach generates more informative supplementary features, enabling the multimodal fusion module to make more accurate predictions.

Furthermore, replacing the MoE in the model with a standard MLP also led to a performance decline. This indicates that MoE is crucial for the multimodal tasks in this context. MoE leverages information from different modalities, along with the rationale provided by the MLLM, to enhance hateful video detection.

In addition, we conducted comprehensive experiments on different combinations of MLLMs, Text encoders, Vision encoders, and Audio encoders, demonstrating the deployment flexibility of HV-GUARD. Details are shown in Appendix B. 524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

4.6 Hyper-parameter Study

To investigate the effects of the hyper-parameters in HVGUARD, we show the impact of hyper-parameters on the performance trend. Details can be found in Appendix C.

5 Related Work

Hate speech detection includes unimodal (text, image, audio) and multimodal methods. Most current methods focus only on a single modality. However, multimodal detection, especially in hateful video detection, integrates these modalities to achieve more comprehensive understanding. Our research not only explores the integration of modalities but also analyzes their interactions to enable a deeper analysis. For more detailed information, please refer to the Appendix D.

6 Conclusion

In this work, we propose a hateful video detection framework named HVGUARD, which is the first reasoning-based hateful video detection framework with MLLMs. This framework carefully designs a CoT reasoning strategy to fully leverage the reasoning ability of MLLMs and introduces a MoE network for the efficient utilization of rationale and multimodal features. Experiments demonstrate that the proposed framework achieves SOTA performance on two public datasets, containing both English and Chinese videos. In the future, we aim to improve the framework by incorporating larger, more diverse, and multilingual datasets to enhance its performance and adaptability across different contexts and languages. This expansion will help address the complexities of detecting hateful content in a broader range of scenarios.

Limitations

We only evaluated HVGUARD on the Chinese and English datasets and did not evaluate other languages. This limits our further exploration of the language generalizability of the framework.

Moreover, we believe that fine-grained detection of hateful videos is of great importance. Although we have considered both binary and ternary classification scenarios, more refined categorization may be more beneficial for the application of such research in real-world contexts.

518

519

522

523

484

571 References

572

573

576

577

578

582

583

585

586

587

588

589

590

592

593

596

597

599

604

606

607

610

611

612

613

614

615

616

617

618

619

620

624

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
 - Cleber Alcântara, Viviane Moreira, and Diego Feijo. 2020. Offensive video detection: dataset and baseline results. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4309– 4319.
 - Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. 2021.
 Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846.
 - Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
 - Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 1(2):3.
 - Aashish Bhandari, Siddhant B Shah, Surendrabikram Thapa, Usman Naseem, and Mehwish Nasim. 2023.
 Crisishatemm: 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.
 - ByteDance. 2016. tiktok. https://www.tiktok.com.
 - Rui Cao, Ming Shan Hee, Adriel Kuek, Wen-Haw Chong, Roy Ka-Wei Lee, and Jing Jiang. 2023. Procap: Leveraging a frozen vision-language model for hateful meme detection. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5244–5252.
 - A Conneau. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
 - Gabriela Csurka. 2017. A comprehensive survey on domain adaptation for visual applications. *Domain adaptation in computer vision applications*, pages 1–35.
- Mithun Das, Rohit Raj, Punyajoy Saha, Binny Mathew, Manish Gupta, and Animesh Mukherjee. 2023. Hatemm: A multi-modal dataset for hate video classification. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 17, pages 1014–1023.

- Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pages 512–515.
- Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Alexey Dosovitskiy. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Paula Fortuna and Sérgio Nunes. 2018. A survey on automatic detection of hate speech in text. ACM Computing Surveys (CSUR), 51(4):1–30.
- Antigoni Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. In *Proceedings of the international AAAI conference on web and social media*, volume 12.
- Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. 2024. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.
- Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao Du, Zhangyu Xiao, and Shiliang Zhang. 2023. Funasr: A fundamental end-to-end speech recognition toolkit. In *INTERSPEECH*.
- Francesca Gasparini, Giulia Rizzi, Aurora Saibene, and Elisabetta Fersini. 2022. Benchmark dataset of memes with text transcriptions for automatic detection of multi-modal misogynistic content. *Data in brief*, 44:108526.
- Mengshi Ge, Rui Mao, and Erik Cambria. 2023. A survey on computational metaphor processing techniques: From identification, interpretation, generation to application. *Artificial Intelligence Review*, 56(Suppl 2):1829–1895.

Google. 2005. Youtube. https://www.youtube.com.

- Ming Shan Hee, Rui Cao, Tanmoy Chakraborty, and Roy Ka-Wei Lee. 2024a. Understanding (dark) humour with internet meme analysis. In *Companion Proceedings of the ACM on Web Conference 2024*, pages 1276–1279.
- Ming Shan Hee, Shivam Sharma, Rui Cao, Palash Nandi, Preslav Nakov, Tanmoy Chakraborty, and Roy Lee. 2024b. Recent advances in online hate speech moderation: Multimodality and the role of large models. *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4407–4419.

625

626

652

653

654

655

656

668 669

670

671

672

673 674

675

676

677

678

782

783

784

785

786

Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. 1991. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87.

679

699

703

704

705

706

707

710

711

712 713

714

715

716

717

718

721

722

723

724

725

726

727 728

729

731

732

- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026.
- Kuanyu. 2009. Bilibili. https://www.bilibili.com.
 - Hongzhan Lin, Ziyang Luo, Jing Ma, and Long Chen. 2023. Beneath the surface: Unveiling harmful memes with multimodal reasoning distilled from large language models. *arXiv preprint arXiv:2312.05434*.
 - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.
 - Robin Matthew Medina, Judith Nkechinyere Njoku, and Dong-Seong Kim. 2022. Audio-based hate speech detection for the metaverse using cnn. , pages 667– 668.
 - Raphael Ottoni, Evandro Cunha, Gabriel Magno, Pedro Bernardina, Wagner Meira Jr, and Virgílio Almeida. 2018. Analyzing right-wing youtube channels: Hate, violence and discrimination. In *Proceedings of the 10th ACM conference on web science*, pages 323–332.
 - Ben Prystawski, Paul Thibodeau, Christopher Potts, and Noah D Goodman. 2022. Psychologicallyinformed chain-of-thought prompts for metaphor understanding in large language models. *arXiv preprint arXiv:2209.08141*.
 - Peng Qi, Zehong Yan, Wynne Hsu, and Mong Li Lee. 2024. Sniffer: Multimodal large language model for explainable out-of-context misinformation detection. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 13052–13062.
 - Yunhang Shen, Chaoyou Fu, Peixian Chen, Mengdan Zhang, Ke Li, Xing Sun, Yunsheng Wu, Shaohui Lin, and Rongrong Ji. 2024. Aligning and prompting everything all at once for universal visual perception. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13193–13203.
 - Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. 2023. Video understanding with large language models: A survey. *arXiv preprint arXiv:2312.17432*.
 - Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al.

2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.

- Nishant Vishwamitra, Keyan Guo, Farhan Tajwar Romit, Isabelle Ondracek, Long Cheng, Ziming Zhao, and Hongxin Hu. 2024. Moderating new waves of online hate with chain-of-thought reasoning in large language models. In 2024 IEEE Symposium on Security and Privacy (SP), pages 788–806. IEEE.
- Han Wang, Tan Rui Yang, Usman Naseem, and Roy Ka-Wei Lee. 2024. Multihateclip: A multilingual benchmark dataset for hateful video detection on youtube and bilibili. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7493–7502.
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*.
- Ching Seh Wu and Unnathi Bhandary. 2020. Detection of hate speech in videos using machine learning. In 2020 International Conference on Computational Science and Computational Intelligence (CSCI), pages 585–590. IEEE.
- Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. 2023. Multimodal large language models: A survey. In 2023 IEEE International Conference on Big Data (BigData), pages 2247–2256. IEEE.
- Yanzhi Xu, Yueying Hua, Shichen Li, and Zhongqing Wang. 2024. Exploring chain-of-thought for multimodal metaphor detection. In *Proceedings of the* 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 91–101.
- Midia Yousefi and Dimitra Emmanouilidou. 2021. Audio-based toxic language classification using selfattentive convolutional neural network. In 2021 29th European Signal Processing Conference (EUSIPCO), pages 11–15. IEEE.
- Xinchen Yu, Eduardo Blanco, and Lingzi Hong. 2022. Hate speech and counter speech detection: Conversational context does matter. *arXiv preprint arXiv:2206.06423*.

A Case Study

To provide a more comprehensive demonstration of HVGuard's effectiveness, we present a detailed case study in Figure 4. In this example, a video titled "*When Find Out a Gay Friend Nearby.mp4*" is processed, where understanding the reactions of different gender groups to homosexuality requires analyzing both visual and textual modalities. In HVGUARD, MLLM leverages CoT prompts to guide



Figure 4: Example of case study.

reasoning from both video frames and transcripts, with the analysis from these modalities integrated to accurately interpret the video content. In contrast, baseline methods lacking MLLM reasoning fail to capture the complementary information between the visuals and the text, leading to incomplete analysis and misclassification.

789

790

795 796

797

802

807

810

811

812

814

815

B Flexibility of framework component

Table 4 shows the impact of different combinations of MLLMs and Encoders. We conducted tests on the ternary classification scenario of Multihateclip(English). The combination of GPT-4o(Achiam et al., 2023), XLM(Conneau, 2019), Vit(Dosovitskiy, 2020), and Wav2Vec(Baevski et al., 2020) achieved the highest M-f1 value, while the combination of Qwen-VL(Bai et al., 2023), Bert(Devlin, 2018), ViViT(Arnab et al., 2021), and Wav2Vec achieved the highest accuracy. MFCC as an Audio Encoder significantly lowered the results, indicating that excellent modality encoders are necessary.

We found that different combinations have varying impacts on performance, with the capabilities of the MLLM being the most significant factor. However, even the least effective combination significantly outperformed the baseline, demonstrating the flexibility and generalizability of our proposed HVGuard framework.

C Hyper-parameter Analysis

Figure 5 illustrates the impact of varying numbers of experts, learning rate and batch size on
the performance through a line chart, showing that
the model achieves optimal performance when the

MIIM	Text	Vision Audio		1 00	M F1	
MILLINI	Encoder	Encoder	Encoder	Acc	IVI-F I	
GPT-40	XLM	Vit	Wav2Vec	0.8090	0.6646	
		vit	MFCC	0.7809	0.4762	
		ViViT	Wav2Vec	0.7921	0.5881	
		VIVII	MFCC	0.7865	0.5604	
		Vit	Wav2Vec	0.8202	0.5562	
	Bert	vit	MFCC	0.7978	0.5590	
		ViViT	Wav2Vec	0.8034	0.6175	
			MFCC	0.8146	0.5384	
Qwen-VL	XLM	Vit	Wav2Vec	0.7865	0.6276	
		vit	MFCC	0.7640	0.4759	
		VAVET	Wav2Vec	0.7809	0.5744	
		VIVII	MFCC	0.7697	0.5637	
	Bert	Vit	Wav2Vec	0.7921	0.5652	
		MFCC 0.775			0.5022	
		ViViT	Wav2Vec	0.7978	0.5282	
		v1 v11	MFCC	0.7809	0.4835	

Table 4: Results of different model combinations onMultihateclip(English)

number of experts is eight, and the learning rate and batch size have little to no impact on the performance. Despite experimenting with different values for these hyperparameters, the model's performance remained relatively stable across the variations, indicating that the performance is primarily influenced by the number of experts rather than the learning rate or batch size.

D Related Work

D.1 Hate Speech Detection

Hate speech detection can be divided into unimodal hate speech detection and multimodal hate speech detection, based on the type of data used. Unimodal hate speech detection is further categorized into text-based, image-based, and audio-based approaches.

Text hate speech detection: This primarily ad-

820

821

822



Figure 5: (a) Number of experts hyper-parameter study. (b) Learning rate hyper-parameter study. (c) Batch size hyper-parameter study.

dresses binary classification tasks, with some studies expanding to three categories: hate speech, offensive speech, and normal speech. Notable studies, such as those by (Davidson et al., 2017) and (Founta et al., 2018), have made significant contributions in classifying hate speech from text. More recently, researchers have explored the subtleties of black humor (Hee et al., 2024a) and discourse context (Yu et al., 2022) to better understand the complexities of text-based hate speech.

Image-based hate speech detection: This area focuses on visual content, such as memes, with studies investigating methods to detect hate speech in images and build appropriate datasets (Gasparini et al., 2022; Bhandari et al., 2023). Approaches like Pro-Cap (Cao et al., 2023) and MR.HARM (Lin et al., 2023) attempt to address challenges in implicit hate speech detection.

Audio-based hate speech detection: Techniques in this domain often involve the use of CNNs to process audio features, such as spectrograms. Works like (Medina et al., 2022) and (Yousefi and Emmanouilidou, 2021) explore methods to enhance audio feature extraction for better detection.

Multimodal hate speech detection: This approach integrates text, image, and audio modalities to enhance hate speech detection, particularly in video content. Studies such as (Das et al., 2023) and (Wang et al., 2024) demonstrate the potential of multimodal techniques in capturing complex context, thereby improving detection performance.

In our study, we focus on multimodal hate speech detection, particularly in videos. While existing research typically concatenates modality information, our approach delves deeper into the interactions between different modalities, improving the understanding of hate speech videos and their intricate patterns.

D.2 Multimodal Large Language Models (MLLMs)

876

877

878

879

880

881

882

883

884

885

886

887

888

889

891

892

893

894

895

896

The emergence of large language models (LLMs) has led to significant advances in natural language processing, enabling models like Gemini (Team et al., 2024) to handle multimodal inputs, such as images and text. While LLMs excel at reasoning and world knowledge, they lack the ability to "see" images, making them less effective at understanding multimodal data. Conversely, large visual models (VLMs) excel in image recognition but are limited in reasoning and world knowledge (Kirillov et al., 2023; Shen et al., 2024). The combination of LLMs and VLMs in MLLMs allows for more robust multimodal understanding, making them highly effective in tasks like image reasoning and video understanding (Wu et al., 2023). In our research, we leverage MLLMs to analyze and understand the complex interaction patterns in hate speech videos, providing valuable insights for reasoning models.

874

875