

# All Changes May Have Invariant Principles: Improving Ever-Shifting Harmful Meme Detection via Design Concept Reproduction

Ziyou Jiang<sup>1,2,3</sup>, Mingyang Li<sup>1,2,3\*</sup>, Junjie Wang<sup>1,2,3</sup>, Yuekai Huang<sup>1,2,3</sup>,  
Jie Huang<sup>1,2,3</sup>, Zhiyuan Chang<sup>1,2,3</sup>, Zhaoyang Li<sup>1,2,3</sup> and Qing Wang<sup>1,2,3\*</sup>

<sup>1</sup>State Key Laboratory of Complex System Modeling and Simulation Technology, Beijing, China <sup>2</sup>Science and Technology on Integrated Information System Laboratory Institute of Software Chinese Academy of Sciences, Beijing, China <sup>3</sup>University of Chinese Academy of Sciences

{ziyou2019, mingyang2017, junjie, yuekai2018, huangjie, zhiyuan2019, lizhaoyang2024, wq}@iscas.ac.cn,

## Abstract

Harmful memes are ever-shifting in the Internet communities, which are difficult to analyze due to their type-shifting and temporal-evolving nature. Although these memes are shifting, we find that different memes may share invariant principles, i.e., the underlying design concept of malicious users, which can help us analyze why these memes are harmful. In this paper, we propose REPMD, an ever-shifting harmful meme detection method based on the design concept reproduction. We first refer to the attack tree to define the Design Concept Graph (DCG), which describes steps that people may take to design a harmful meme. Then, we derive the DCG from historical memes with design step reproduction and graph pruning. Finally, we use DCG to guide the Multimodal Large Language Model (MLLM) to detect harmful memes. The evaluation results show that REPMD achieves the highest accuracy with 81.1% and has slight accuracy decreases when generalized to type-shifting and temporal-evolving memes. Human evaluation shows that REPMD can improve the efficiency of human discovery on harmful memes, with 15~30 seconds per meme.

**Disclaimer:** This paper may contain content that is disturbing to some readers.

## 1 Introduction

Nowadays, memes have emerged as important cultural symbols on the Internet. They use both visual and textual elements to convey the designer's opinion for organizations, regimes, and social events (Lin et al., 2024). However, malicious users may design harmful memes as a weapon to express their biased viewpoints (Sharma et al., 2022), which may pose threats to society and other users.

The rapid growth of the Internet community has led to the following ever-shifting characteristics

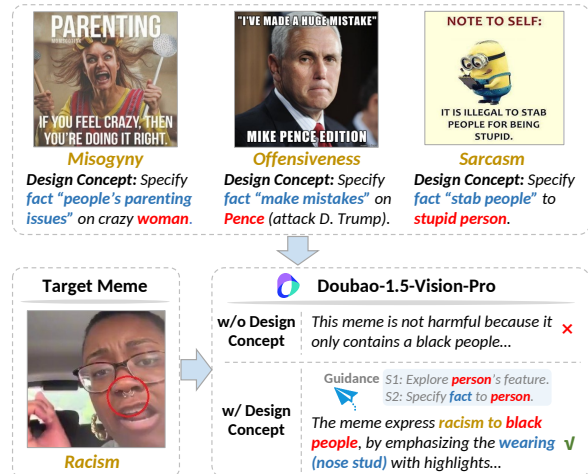


Figure 1: The motivation example of REPMD.

of harmful memes (Valensise et al., 2021): (1) **Type Shifting:** Memes come in many forms and attack different targets, and (2) **Temporal Evolving:** Memes are temporally related to some events. This ever-shifting nature makes the visual elements and expression of new memes quite different from the historical ones. For example, Figure 1's target meme uses a new and implicit way to express discrimination against black people, through the highlighted red circles on the people's accessories (i.e., nose stud). Existing harmful meme detection methods (Kiela et al., 2020; Suryawanshi et al., 2020) only learn the combination of harmful elements, and lack understanding of these implicit expressions, so they cannot accurately detect the harmfulness. Moreover, new slang abbreviations like "GOAT" and "Stan" are used in memes to express harmful viewpoints, and these rarely-seen expression improve the difficulty of detection.

Although harmful memes are ever-shifting, we find that they have "invariant principles". In Figure 1, we choose four example memes that have different harmful types, but share similar underlying ideas that MLLM can learn from each other.

\*Corresponding author.

For example, the idea of the leftmost meme is to specialize in fact (e.g., parenting issues) to a certain group of people (e.g., crazy women), thereby attacking these people. This is the **Design Concept** of memes, expressing the attacks through the *human specification* (Sălcudean, 2020). By analyzing the design concepts, the test MLLM, i.e., Doubao-1.5-Vision-Pro (Guo et al., 2025), can accurately identify the target meme’s racism: the malicious user specifies a fact "people who wear nose stud" onto "black people" to express the stereotype.

In the previous cases, we can see that the design concept is useful. However, previous works have not defined an explainable structure to formally describe this design concept. Therefore, we need to define its structure and derive the contents based on the visible information of memes. Moreover, we also need to propose the corresponding usage strategy of design concepts, which can effectively guide the MLLM to detect target harmful memes.

In this paper, we propose an automated approach that **Reproduce** memes’ design concept to improve the ever-shifting harmful **Meme Detection**, named **REPMD**. Inspired by the effectiveness of the attack tree (Schneier, 1999), we define the structure of DCG, a **heterogeneous graph** (Zhang et al., 2019; Gao et al., 2025) that describes steps that people may take to design a harmful meme, as well as the goal they want to achieve. This structure can explain the user’s design logic and guide MLLM to stepwise identify the harmfulness. Based on this structure, we derived a DCG from historical memes that MLLMs fail to predict accurately, through the design steps’ reproduction and graph pruning. For target memes, we retrieve similar reproduction steps from DCG, then form the stepwise guidance to help MLLM detect target harmful memes in the ever-shifting scenario.

To evaluate the performance of REPMD, we conduct two types of experiments on a dataset of 58,192 memes. For the type-shifting experiment on public GOAT-Bench memes, REPMD achieves the highest accuracy with 81.1% and has only a 2.1% accuracy decrease in out-of-domain evaluation. For the temporal-evolving experiment on our manually crawled Twitter memes, REPMD also outperforms baselines and has a 0.3% accuracy improvement in other quarters’ evaluation. Moreover, our human evaluation shows that REPMD can improve the efficiency of human discovery on harmful memes, with 15~30 seconds per meme.

This paper makes the following contributions:

- We propose REPMD, a harmful meme detection method with design concept reproduction, applicable for ever-shifting memes on the Internet.
- We evaluate REPMD on ever-shifting memes from two data sources, which outperforms baselines and can be generalized to type-shifting and temporal-evolving memes with the help of DCG.
- We conduct a human evaluation to illustrate that DCG has high explainability and can help evaluators manually identify harmful memes.
- We release the code and dataset<sup>1</sup> to facilitate REPMD’s reproducibility.

## 2 Definition of Design Concept

To describe the meme’s design concept, we refer to the idea of attack trees (Schneier, 1999), a threat model that contains STRIDE threat types, attack methods, goals, and logic of the attack reproduction. It can guide LLM-based security testing tools to warn of software vulnerabilities (Xiong and Lagerström, 2019), which aligns with MLLM’s logic for identifying memes’ harmful information.

Before extracting the DCG, we first analyze the reasons why harmful memes cannot be detected by MLLM, then form the fail reason tree. It contains the type and fail reason nodes, as well as the type and link edges. Then, we derive the fail reason node to the steps that the harmful users will take to design harmful memes, which contain the type, reproduction method, goal, and logic gate nodes, as well as the type, link, and achievement edges.

### 2.1 Definition of Fail Reason Tree

The structure of fail reason tree can be formed as  $\mathcal{G}_F = \{\{\mathcal{N}_T, \mathcal{N}_F\}, \{\mathcal{E}_T, \mathcal{E}_{Link}\}\}$ , defined as follows:

**Node: (1) Type Node  $\mathcal{N}_T$ :** It contains three node levels  $L_1 \sim L_3$ , where the chosen value is in Table 7, similar to STRIDE threat types in attack trees: **Macro Type Nodes  $\mathcal{N}_T^{L_1}$ :** The seven fixed categories in Table 7, defined in the previous works (Chen et al., 2025); **Subtype Nodes  $\mathcal{N}_T^{L_2,3}$ :** The detailed types that need to be specified, e.g., *Culture*→*Video Games*. These subtypes are **not limited to the table** and can be extended. **(2) Fail Reason Node  $\mathcal{N}_F$ :** It contains one-sentence MLLM’s fail reason and memes’ text description. **Edge: (1) Type Edge  $\mathcal{E}_T$ :** The edge is  $\mathcal{N}_T^{L_i} \rightarrow \mathcal{N}_T^{L_{i+1}}$ , indicating the division from the macro type

<sup>1</sup><https://github.com/jzySaber1996/RepMD>

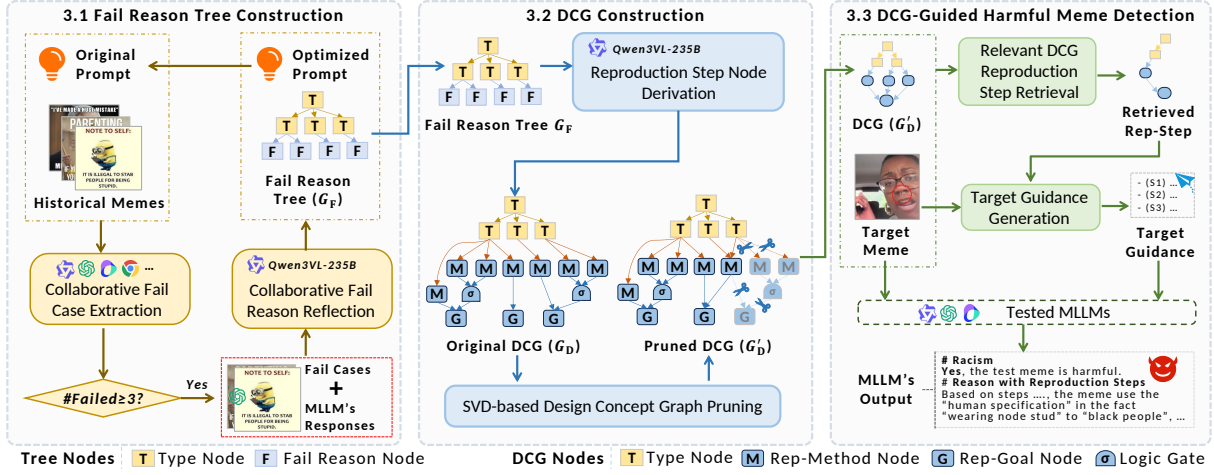


Figure 2: The overview of REPM D.

to subtypes; (2) **Link Edge**  $\mathcal{E}_{\text{Link}} = \mathcal{N}_T \rightarrow \mathcal{N}_F$ : It indicates that one type of incorrectly predicted harmful meme has corresponding fail reasons.

## 2.2 Definition of DCG

We extract design concepts of historical memes and gather them into the DCG. The structure refers to the attack tree’s three-level structure, i.e., attack method, attack goal, and logic gate, which reflect how to reproduce a harmful meme. The structure of DCG can be formed as  $\mathcal{G}_D = \langle \{\mathcal{N}_T, \mathcal{N}_M, \mathcal{N}_G, \mathcal{N}_\sigma\}, \{\mathcal{E}_T, \mathcal{E}_A, \mathcal{E}_{\text{Link}}\} \rangle$ , where the reproduction nodes are derived from reason nodes  $\mathcal{N}_F \rightarrow (\mathcal{N}_M, \mathcal{N}_G, \mathcal{N}_\sigma)$ , which is the reproduction step for the DCG. defined as follows:

**Node:** (1) **Type Node**  $\mathcal{N}_T$ : They are the same as the fail meme tree. (2) **Reproduction Method**  $\mathcal{N}_M$ : These nodes describe the steps that malicious users take to design harmful memes; (3) **Logic Gate**  $\mathcal{N}_\sigma$ : There is a logical combination between the reproduction methods to design a meme, i.e., *And* ( $\wedge$ ), *Or* ( $\vee$ ), and *Not* ( $\neg$ ); (4) **Reproduction Goal**  $\mathcal{N}_G$ : A design goal that a malicious user aims to achieve, such as the *human specification*. Each node also contains a harmful indicator  $\{0, 1\}$  that identifies which node may pose a harmful.

**Edge:** (1) **Type Edge**  $\mathcal{E}_T$ : It is the same as the fail meme tree. (2) **Link Edge**  $\mathcal{E}_{\text{Link}} = \mathcal{N}_T \rightarrow \mathcal{N}_M$ : The structure is the same as the fail meme tree, but the meaning is slightly different. The link indicates that in historical memes, one type of harmful meme has the corresponding reproduction step. (3) **Achievement Edge**  $\mathcal{E}_A$ : Edges between  $(\mathcal{N}_M, \mathcal{N}_G, \mathcal{N}_\sigma)$ , which are the logic indicating the achievement between methods and the goal.

## 3 Overview of REPM D

The element changing in ever-shifting harmful memes makes them difficult to detect. As shown in Figure 1, the target meme expresses racism against people in an unknown and implicit way. With the help of the design concepts, we find that different types of memes may have invariant principles, such as the above-mentioned *human specification*. This invariant feature can help MLLMs infer why the unknown meme is harmful.

Figure 2 illustrates the overview of REPM D. Guided by the concept, we first explore the reasons in historical meme cases that MLLMs fail to detect harmful memes, formed as the fail reason tree defined in Section 2. Second, we derive reproduction steps from the fail reason tree and form DCG, then propose a fast Singular Value Decomposition (SVD)-based graph pruning method, proved effective in GNN’s dimensionality reduction (Cai et al., 2023). These steps form the design concept with the relation "*Fail Reason Tree*  $\rightarrow$  *DCG*  $\rightarrow$  *SVD*". Finally, REPM D retrieves relevant reproduction steps from DCG to help MLLMs detect whether the target meme is harmful.

### 3.1 Fail Reason Tree Construction

Before constructing the DCG  $\mathcal{G}_D$ , we aim to explore the reason why MLLMs fail to detect harmful memes. Given historical memes  $[M_1^H, \dots, M_n^H]$ , we propose an looped framework of "*extraction, reflection, and optimization*". After the loop, we can obtain a fail reason tree  $\mathcal{G}_F$  and the optimized prompt. Through prompt’s optimization, the fail reason tree will contain some harmful memes that cannot be accurately detected, no matter how the



### 3.2.2 SVD-based DCG Pruning

The DCG is the gathering of all historical memes, so some memes may have similar design concepts, introducing redundant nodes and edges to the DCG. Therefore, we propose the SVD-based pruning method, as shown in Algorithm 1. The main idea is how to design an **adjacency matrix** to describe the relations between nodes  $\mathcal{N}_i$  and  $\mathcal{N}_j$ , so we propose the reproduction score (line 2):

$$\text{Score}_{\text{rep}}(\mathcal{N}_{ij}) = \text{ReLU}(\text{sim}(\mathcal{N}_i, \mathcal{N}_{\text{root}}) - \text{sim}(\mathcal{N}_i, \mathcal{N}_j)) \quad (1)$$

where  $\text{sim}(\mathcal{N}_i, \mathcal{N}_j) = 0$  if  $\mathcal{N}_i/\mathcal{N}_j = \mathcal{N}_{\text{root}}$

where the function  $\text{sim}(\cdot, \cdot)$  is the cosine similarity between TF-IDF values of the node’s contents.  $\mathcal{N}_T$ ’s root node is  $\mathcal{N}_T^{L1}$ , and  $\mathcal{N}_M$ ’s root node is  $\mathcal{N}_G$ . The final scores are activated with the ReLU function (Glorot et al., 2011) to ensure the positive definiteness of the adjacent matrix. The equation implies that we want a node to have **higher correlation** to the root node and **lower similarity** to other nodes, which means that this node is not redundant and will not be pruned. Moreover, we find that redundancy mainly comes from reproduction steps, so we introduce scaled factors  $\alpha, \beta$  for the matrix of the reproduction step ( $\mathcal{N}_M, \mathcal{N}_G, \mathcal{N}_\sigma$ )  $\rightarrow \mathbf{A}_{\text{Repr}}$  and the edge matrix  $\mathcal{E}_{\text{Link}} \rightarrow \mathbf{A}_{\text{Link}}$  ( $1 > \beta > \alpha$ ).

From line 3 to 12, after we have initialized the base matrix  $\mathbf{A}$ , we calculate the  $t$ -hop matrix’s indirect feature value, i.e.,  $\mathbf{A}^t$  in the loop, which indicates how the nodes are correlated to other nodes through the graph’s relations. Then, we set a hyperparameter  $\theta$  to represent the retained nodes in  $\mathcal{G}'_D$ . We first introduce the Laplacian Normalization (Li et al., 2024b) to process the adjacent matrix to  $\mathbf{L}$ .

In line 5, during the SVD, we obtain singular values  $[\lambda_1, \dots, \lambda_n]$ , which can be used to represent the matrix’s feature with the low-rank approximation. Since the distribution of singular values may be long-tailed, we calculate the difference on a logarithmic scale and choose the cut-off value  $\lambda_{\text{cut}}$  with the steepest decline.

Finally, in line 11, if the proportion of retained nodes does not achieve  $\theta$ , we explore the impact of the  $t$ -hop matrix’s indirect relations.

### 3.3 DCG-Guided Harmful Meme Detection

Since REPM D has constructed DCG with multiple prompts, the detection is lightweight by reusing previous prompts. Given target memes  $[M_1^T, \dots, M_n^T]$  and the test MLLM, we retrieve the similar reproduction steps from the DCG  $\mathcal{G}'_R$  and guide the harmful meme detection. Since this step is executed

---

#### Algorithm 1: SVD-based DCG pruning.

---

**Input:** The original DCG  $\mathcal{G}_D$ , retained proportion  $\theta$ .

**Output:** The pruned DCG  $\mathcal{G}'_D$ .

```

1  $\mathbf{A} = \text{Score}_{\text{rep}}(\mathcal{G}_D)$  with Equation (1), where
    $\mathbf{A} = \begin{pmatrix} \mathbf{A}_{\text{Type}} & \beta \mathbf{E}_{\text{Link}} \\ \beta \mathbf{E}_{\text{Link}} & \alpha \mathbf{A}_{\text{Repr}} \end{pmatrix}$ ,  $\alpha, \beta$ : scaled factors;
2 Calculate Degree Matrix  $\mathbf{D}$ , Initialize  $\mathcal{G}'_D = \emptyset$ ;
3 while  $t \leq 5$  do
4    $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A}^t \mathbf{D}^{-1/2}$ ;
5   SVD:  $\mathbf{L} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T$  (See Appendix A.2.1);
    $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$ ;
6   Cut-off Determination (See Appendix A.2.2):
    $\text{cut} = \arg \max_i |\ln(\lambda_{i+1} - \lambda_i)|$ ;
7   Graph Pruning:  $\mathcal{G}'_D = \langle \{\mathcal{N}_1, \dots, \mathcal{N}_{\text{cut}}\}, \{\mathcal{E}'\} \rangle$ ,
   where  $\mathcal{N}_i := \lambda_i$ ,
    $\{\mathcal{E}'\} = \{\mathcal{N}_i \rightarrow \mathcal{N}_j | i, j \leq \text{cut}\}$ ;
8   if  $\text{cut}/n \geq \theta$  then
9     break;
10  end
11   $t$ -hop Indirect Features:  $t = t + 1$ ;
12 end
13 return  $\mathcal{G}'_D$ ;
```

---

online, the time cost of REPM D comes from the scale of the test dataset.

#### Step-1: Relevant DCG Reproduction Step Retrieval.

We extract the historical reproduction step as follows: **(1) DCG Node Retrieval:** We reuse the prompt  $P_F$  and  $P_D$  to select DCG’s nodes that may relate to  $M_i^T$ , aggregated as a set  $\{\mathcal{N}_i | \mathcal{N}_i \in \mathcal{G}'_D\}$ . **(2) Reproduction Step Formation:** We cluster the subgraphs based on the edges, i.e.,  $\mathcal{G}_{\text{sub}} = \langle \{\mathcal{N}_i\}, \mathcal{E}_{\text{sub}} \subset \mathcal{G}'_D \rangle$ .

#### Step-2: Target Guidance Generation.

We transform the retrieved step  $\mathcal{G}_{\text{sub}}$  into plain text with logic semantics (see Appendix A.1), then feed the target meme and this plain text into the test MLLM (same model for testing the follow-up harmful meme detection) to output the guidance as  $(S_1) \rightarrow (S_n)$ . Finally, REPM D inputs the generated guidance with the target meme  $M_i^T$  into the test MLLM, then uses the optimized prompt  $P'_{\text{Harm}}$  in Figure 9 to detect the target harmful memes.

## 4 Experimental Design

To evaluate the performance of REPM D, we introduce three Research Questions (RQs).

**RQ1:** *What are the performances of REPM D on detecting ever-shifting harmful memes, i.e., type-shifting or temporal-evolving memes.*

**RQ2:** *How do each component contribute to REPM D’s performances?*

**RQ3:** Can REPMD’s DCG help humans understand why the memes are harmful?

**Dataset Preparation.** We construct the dataset from two sources, i.e., GOAT-Bench (Lin et al., 2024) and our crawled memes from Twitter, and the steps are shown as follows.

**(1) Type-Shifting Dataset Preprocessing:** We first reuse all memes in GOAT-Bench. Then, we manually find that it has overlap categories (misogyny and hatefulness) and misclassified samples (classify misogyny memes as hatefulness), so we replace the type name, e.g., "Harmfulness" → "Toxicity" (violence and suicide, etc) and "Hatefulness" → "Racism" (mainly the racial hatefulness). Then, we reclassify 110 misogyny memes to the correct type. The revised data is used to detect type-shifting harmful memes.

**(2) Temporal-Evolving Dataset Annotation:** There are currently no usable benchmarks, so we manually crawl the memes on Twitter in **year: 2025** and divide them into four quarters. In each quarter, we randomly select 500 memes for DCG construction and evaluation. We have invited three researchers with  $\geq 3$  years of experience in harmful meme detection (not included in the authors), and asked them to annotate the dataset **independently** based on the criteria for determining harmful content (Pandiani et al., 2025). Each meme is labeled by these three annotators, and they check whether the labels are accurate. The average Cohen’s Kappa (Pérez et al., 2020) value is 0.86, which achieves a high agreement on the labels.

Table 1: The statistics of our constructed dataset.

Types	Type-Shifting Dataset from GOAT-Bench				
	Racism	Misogyny	Offensiveness	Sarcasm	Toxity
DCG	8,310	10,190	7,000	19,816	4,250
Target	2,080	920	743	1,820	1,063

Quarters	Temporal-Evolving Dataset from Twitter				Total
	Jan~Apr	Apr~Jun	Jul~Sep	Oct~Dec	
DCG	400	400	400	400	51,166
Target	100	100	100	100	7,026

**(3) Dataset Size and Labeling Cost:** Table 1 shows that we collect 58,192 memes (54.0% are harmful) in the dataset, where 56,192 come from GOAT-Bench for type-shifting evaluation, and 2,000 are crawled memes for temporal-evolving experiments. Moreover, the manually-labeled dataset comes from GOAT-Bench’s adjustment and Twitter’s, so the cost of human labor is small.

**Baselines.** We select three novel and representative baselines. One method utilizes Retrieval-Augmented Generation (RAG) for harmful meme

detection, which meets the evaluation requirements of DCG construction and retrieval. The other two methods introduce few-shot Supervised Fine-Tuning (SFT), which meet the experimental requirements on ever-shifting memes: **ModeHate** (Cao et al., 2024) is a few-shot hateful meme detector fine-tuned with LoRA; **MIND** (Liu et al., 2025) is a zero-shot RAG-based framework with bidirectional insight-augmented inference; and **RAHMD** (Mei et al., 2025) is a fine-tuned adaptation framework for hateful meme detection. Moreover, we have also selected four representative MLLMs as baselines, i.e., **Qwen2.5VL-32B** (Bai et al., 2025), **GPT-4o**, **Doubao-1.5-Vision-Lite**, and **Doubao-1.5-Vision-Pro** (Guo et al., 2025).

**Metrics.** To ensure a fair and comprehensive evaluation, we measure model performance using two commonly-used metrics: **Accuracy** and **F1-Score**, following the setting of GOAT-Bench.

**Experimental Settings.** For implementation, we use MLLM’s APIs and set *temperature* = 0 to make the output fixed, *Top\_P*, and *max\_token* as default values in the vanilla MLLM. For hyperparameters, we set  $\alpha = 0.3$  and  $\beta = 0.6$  for SVD, and  $\theta = 75\%$  (see Appendix A.4). All experiments run on four GeForce RTX A6000 GPUs.

## 5 Results

### 5.1 Overall Performance of REPMD on Evolving Memes (RQ1)

**Type-Shifting Experiments.** We conduct this experiment in the following types: **(1) In-Domain (ID) Evaluation:** We construct DCG and detect target memes’ harmfulness within the same type. **(2) Out-Of-Domain (OOD) Evaluation:** We introduce the cross-type analysis, which means that we choose one specific type of memes in GOAT-Bench as the target, and use **the other four types of memes** as historical to generate DCG.

Table 2 illustrates the results for type-shifting scenario, where  $\Delta(\text{ID}, \text{Vanilla})$  and  $\Delta(\text{OOD}, \text{ID})$  indicate the absolute improvement of REPMD and the relative decrease when migrated to the OOD memes. We can see that, the absolute detection performance of REPMD achieves 81.1% accuracy, outperforming all baseline approaches. The improvement on vanilla MLLMs achieves +9.1% (F1) and +10.5% (Accuracy). Besides, the OOD experiment shows that the decrease of REPMD’s performance is smaller than baselines (which is -18.7%

Table 2: The performance of REPMd on In-Domain (ID) and Out-Of-Domain (OOD) memes(%).

Evaluation on Types Approaches	Racism		Misogyny		Offensiveness		Sarcasm		Toxicity		Average	
	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.
<i>Few-Shot SFT and RAG-based Baselines for Harmful Meme Detection</i>												
Mod-Hate <sub>ID</sub> (Cao et al., 2024)	58.3	58.1	61.9	61.5	63.0	64.1	61.3	63.3	71.2	69.3	63.1	63.3
Mod-Hate <sub>OOD</sub>	36.6	34.9	43.7	44.5	40.2	43.7	42.1	45.2	59.3	55.1	44.4	44.7
MIND <sub>ID</sub> (Liu et al., 2025)	69.1	66.2	69.2	68.7	63.7	62.5	64.0	63.2	71.6	71.2	67.5	66.4
MIND <sub>OOD</sub>	47.4	45.0	49.5	49.4	42.1	42.0	45.9	44.3	55.3	54.0	48.0	46.9
RA-HMD <sub>ID</sub> (Mei et al., 2025)	74.6	73.4	77.3	76.8	71.9	70.1	70.4	71.3	70.2	69.5	72.9	72.2
RA-HMD <sub>OOD</sub>	57.2	56.3	50.2	48.6	48.0	47.4	45.2	44.3	64.5	62.3	55.0	51.8
Average $\Delta$ (OOD, ID)	$\downarrow 20.3$	$\downarrow 20.5$	$\downarrow 21.7$	$\downarrow 21.5$	$\downarrow 22.8$	$\downarrow 21.2$	$\downarrow 20.8$	$\downarrow 21.3$	$\downarrow 11.3$	$\downarrow 12.9$	$\downarrow 18.7$	$\downarrow 19.5$
<i>Test MLLMs for Harmful Meme Detection</i>												
Qwen2.5VL-32B (Vanilla <sup>*</sup> )	72.6	68.3	65.2	60.9	63.2	55.3	71.0	68.7	62.3	53.9	66.9	61.4
+REPMd <sub>ID</sub>	78.2	77.4	74.9	74.5	71.0	70.6	80.4	79.6	79.1	77.0	76.7	75.8
+REPMd <sub>OOD</sub>	76.8	75.4	73.0	72.6	70.2	70.1	77.2	77.2	77.5	75.8	74.9	74.2
GPT-4o (Vanilla)	76.6	75.2	82.7	82.4	59.1	59.1	68.7	66.3	66.3	64.2	70.7	69.4
+REPMd <sub>ID</sub>	<b>88.5</b>	<b>87.6</b>	<b>90.5</b>	<b>90.4</b>	71.2	71.1	<b>84.2</b>	<b>83.8</b>	80.2	78.5	<b>82.9</b>	<b>82.3</b>
+REPMd <sub>OOD</sub>	87.1	87.0	88.7	86.2	70.4	70.2	80.5	79.6	77.6	76.9	80.9	80.0
Doubao-1.5-Vision-Lite (Vanilla)	70.5	69.1	67.3	68.0	66.4	65.5	68.6	68.7	70.0	70.2	68.6	68.3
+REPMd <sub>ID</sub>	77.4	77.2	74.5	74.1	76.3	75.4	75.0	74.9	78.6	78.5	76.4	76.0
+REPMd <sub>OOD</sub>	76.2	75.5	73.2	73.7	72.1	73.1	74.1	72.6	77.5	76.2	74.6	74.2
Doubao-1.5-Vision-Pro (Vanilla)	72.4	71.9	73.4	72.9	67.4	64.3	71.3	70.4	72.0	72.1	71.1	69.9
+REPMd <sub>ID</sub>	81.5	80.6	82.2	81.3	<b>80.2</b>	<b>81.3</b>	81.2	81.3	<b>83.0</b>	<b>82.8</b>	81.3	81.1
+REPMd <sub>OOD</sub>	80.3	79.4	80.4	79.3	78.2	77.6	78.5	78.4	81.7	80.9	79.4	78.7
Average $\Delta$ (ID, Vanilla)	$\uparrow 7.7$	$\uparrow 9.0$	$\uparrow 7.6$	$\uparrow 7.9$	$\uparrow 9.4$	$\uparrow 12.8$	$\uparrow 9.2$	$\uparrow 9.8$	$\uparrow 11.7$	$\uparrow 13.1$	$\uparrow 9.1$	$\uparrow 10.5$
Average $\Delta$ (OOD, ID)	$\downarrow 1.3$	$\downarrow 1.4$	$\downarrow 1.7$	$\downarrow 2.1$	$\downarrow 2.0$	$\downarrow 1.9$	$\downarrow 2.6$	$\downarrow 3.0$	$\downarrow 1.7$	$\downarrow 1.7$	$\downarrow 1.9$	$\downarrow 2.1$

\* "Vanilla" means using original MLLMs with the optimized detection prompt  $P'_{Harm}$  to detect harmful memes.

F1), with -1.9% (F1) and -2.1% (Accuracy). By analyzing the inference of MLLM, we find that they detect harmful memes by analyzing what elements are harmful, but may ignore the implicit expressions (e.g., Figure 1’s red circle) and seemingly harmless elements. From the experiment results, we can see that DCG guides MLLM to understand the idea of memes’ design, enabling REPMd to detect harmful memes across types. Even if there are not too many training samples labeled with fine-grained types, REPMd can achieve effective detection ability by  $\{0, 1\}$  harmful labels. Moreover, we have conducted a pairwise T-test to measure the significance between MLLMs and REPMd<sup>2</sup>. We can see that, in Table 3, REPMd can significantly outperform the vanilla MLLMs, and the column "OOD vs ID" shows that migrating REPMd to type-shifting memes will not reduce the performance significantly.

Table 3: The  $p$ -values of pairwise significant T-testing on type-shifting and temporal-evolving memes.

Tested MLLMs	REPMd vs Vanilla	OOD vs ID	TE vs TF
Qwen2.5VL-32B	$3.0 \times 10^{-6}$	0.007	0.017
GPT-4o	$7.0 \times 10^{-9}$	0.037	0.032
Doubao-1.5-Vision-Lite	$4.6 \times 10^{-10}$	0.008	0.016
Doubao-1.5-Vision-Pro	$7.1 \times 10^{-10}$	0.052	0.020

<sup>2</sup> $p < 0.01$  means high significance.

**Temporal-Evolving Experiments.** We conduct this experiment in the following types: **(1) Temporal Fixed (TF) Evaluation:** We construct DCG and evaluate REPMd on same quarter’s memes. **(2) Temporal Evolving (TE) Evaluation:** To ensure fairness of evaluation, we choose  $quarter_i$ ’s cases to construct DCG and detect harmful memes in  $quarter_{i+1}$ ’s cases.

Table 4: The performance of REPMd on Temporal Evolving (TE) and Temporal Fixed (TF) memes (%).

Evaluation on Quarters Models (Top-2 MLLMs)	Apr~Jun <sup>*</sup>		Jul~Sep		Oct~Dec	
	F1	Acc.	F1	Acc.	F1	Acc.
RA-HMD <sub>TF</sub>	45.9	44.9	52.7	51.5	56.4	55.2
RA-HMD <sub>TE</sub>	31.7	31.2	40.2	40.0	45.3	44.1
$\Delta$ (TE, TF)	$\downarrow 14.2$	$\downarrow 13.7$	$\downarrow 12.5$	$\downarrow 11.5$	$\downarrow 11.1$	$\downarrow 11.1$
GPT-4o (Vanilla)	51.3	50.0	56.3	55.0	67.4	65.0
+REPMd <sub>TF</sub>	60.2	60.0	63.0	66.0	83.5	80.0
+REPMd <sub>TE</sub>	67.4	65.0	61.6	65.0	82.4	78.0
$\Delta$ (TF, Vanilla)	$\uparrow 8.9$	$\uparrow 10.0$	$\uparrow 6.7$	$\uparrow 11.0$	$\uparrow 16.1$	$\uparrow 15.0$
$\Delta$ (TE, TF)	$\uparrow 7.2$	$\uparrow 5.0$	$\downarrow 1.4$	$\downarrow 1.0$	$\downarrow 1.1$	$\downarrow 2.0$
Doubao-1.5-V-Pro (Vanilla)	67.4	65.0	66.2	65.0	67.4	65.0
+REPMd <sub>TF</sub>	82.4	80.0	82.4	80.0	86.7	85.0
+REPMd <sub>TE</sub>	82.4	80.0	82.4	80.0	86.7	85.0
$\Delta$ (TF, Vanilla)	$\uparrow 15.0$	$\uparrow 15.0$	$\uparrow 16.2$	$\uparrow 15.0$	$\uparrow 19.3$	$\uparrow 20.0$
$\Delta$ (TE, TF)	0.0	0.0	0.0	0.0	0.0	0.0

\* The year: 2025 is divided into quarter-based intervals.

Table 4 shows the Top-2 MLLMs’ performances on Twitter’s temporal-evolving memes within the last three quarters. We can see that, REPMd improves the vanilla MLLMs with +13.7% (F1) and +14.3% (Accuracy) on average. Moreover, we calculate the difference  $\Delta$ (TE, TF) and find that introducing historical memes can even improve the

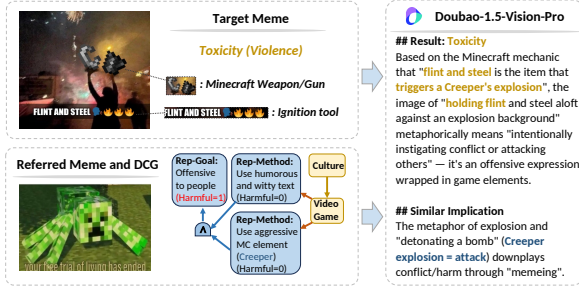


Figure 4: The case study of temporal evolving memes.

Table 5: The contribution of REPMD on detecting harmful memes with different design concepts.

Approach	Human Specification	Minecraft Replacement
MIND	21/50 (42.0%)	20/50 (40.0%)
RA-HMD	35/50 (79.0%)	36/50 (72.0%)
Vanilla GPT-4o	34/50 (68.0%)	31/50 (62.0%)
GPT-4o+REPMd	44/50 (88.0%)	47/50 (94.0%)

detection performance of memes, with +0.8% (F1) and +0.3% (Accuracy) on average. This advantage comes from the consistency with the design concepts of the same organization. Some subcultural groups (e.g., Bitcoin’s Dogwifhat, Minecraft, and K-pop Fancam, etc.) have opinion leaders, and group members will utilize a similar expression pattern to design memes, so the continuity of these patterns can help us identify the causes of harmful memes accurately. Finally, in Table 3’s column "TE vs TF", we can see that even when migrated to temporal evolving memes, REPMD’s performance reductions are not significant, which further indicates the generalizability of our approach.

**Case Study for Design Concepts.** Figure 4 shows that the Minecraft group’s users use elements like "Creepers" and "Ignition Tool" instead of regular weapons to express violence. This expression is successfully captured by REPMD in the DCG, so it can help the MLLMs understand ideas of the group’s opinion leaders and members, thus detecting harmful memes that are rapidly shifting.

We have conducted another experiment to illustrate the advantages of GPT-4o+REPMd with the help of DCG, where we manually select 200 examples with two design concepts, i.e., "Human Specification" and "Minecraft Replacement", where 150 memes are used for constructing DCG, and we detect harmful memes in the remaining 50 memes. The following results show the number of accurately predicted memes. We can see that, in Table 5, REPMD can not only outperform baselines, but

Table 6: The accuracy of REPMD’s ablation study (component removal or replacement) on three variants (%).

Variants	ID	OOD	TE
GPT-4o+REPMd	82.3	80.0	69.3
w/o Tree	81.1 (↓1.2)	79.2 (↓0.8)	68.1 (↓1.2)
w/o VoteMLLMs	81.0 (↓1.3)	79.2 (↓0.8)	65.3 (↓4.0)
w/o OptPrompt	79.0 (↓3.3)	78.6 (↓1.4)	64.9 (↓4.4)
w/o SVD	76.8 (↓5.5)	71.6 (↓8.4)	60.5 (↓8.8)
SVD→GPT-4o	79.9 (↓2.4)	78.2 (↓1.8)	68.1 (↓1.2)
SVD→Doubao-V-Pro	Failed	79.1 (↓0.9)	69.0 (↓0.3)
SVD→Qwen3VL	Failed	Failed	70.0 (↑0.7)
w/o DCG	76.3 (↓6.0)	70.9 (↓9.1)	60.1 (↓9.2)
Retrieval→ImgSim	78.2 (↓4.1)	77.1 (↓2.9)	64.5 (↓4.8)
Retrieval→GraphRAG	81.1 (↓1.2)	79.5 (↓0.5)	65.5 (↓3.8)

also detect 26 harmful memes with DCG, where vanilla GPT-4o cannot detect them.

## 5.2 Ablation Study of REPMD (RQ2)

To evaluate the contribution of model’s components to REPMD, we analyze three types of **variants** based on the structure in Section 3: **Sec 3.1:** Removing the failed reason tree and using all historical memes to construct DCG directly (**w/o Tree**); Removing the vote mechanism (**w/o VoteMLLMs**) and prompt optimization (**w/o OptPrompt**); **Sec 3.2:** Replacing DCG with fail reason tree (**w/o SVD**) and SVD with MLLMs (**SVD→MLLMs**); **Sec 3.3:** Directly generating reproduction steps without retrieval (**w/o DCG**), replacing retrieval with image similarity (**Retrieval→ImgSim**), and the most similar graph-based RAG method, i.e., (**Retrieval→GraphRAG** (Peng et al., 2024)).

Table 6 shows the results of the ablation study. We can see that, REPMD outperforms most variants (5/6) in terms of the experiments ID, OOD, and TE. The maximum decrease is in the components that are related to the DCG, i.e., *w/o SVD* and *w/o DCG*, which achieve over -8.0% and -9.0% Accuracy. Considering *SVD→MLLM*, we find that replacing SVD with Qwen3VL-235B fails on ID and OOD because of information truncation (too many input tokens) and incorrect output (pruned graph’s nodes are different from original DCG). However, MLLM successfully prunes the graph and slightly outperforms SVD on the TE task, with a time cost of over 14× (see Appendix A.4). This indicates that SVD offers advantages in balancing time cost and performance, but fine-tuning an LLM-based pruning method may improve DCG quality and detection performance.

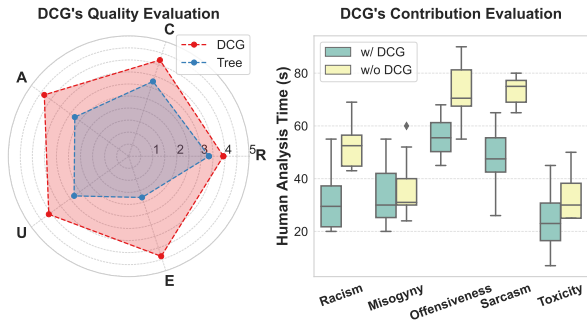


Figure 5: The results of human evaluation, i.e., scores of five evaluation criteria and analysis time cost.

### 5.3 Human Evaluation of REPMD (RQ3)

We conduct the human evaluation on the DCG in two aspects: (1) **DCG’s Quality**: We have invited three types of human evaluators categorized by research experiences on harmful memes (i.e.,  $\geq 3$  year veterans, 1  $\sim$  3 year beginners, and non-experienced social network users). Each category includes two evaluators. We randomly select 10 historical memes with their corresponding fail reason trees and the DCG’s design concepts, then rate these patterns on a 5-point Likert scale (from 1: poor to 5: excellent) based on five key criteria (Mazhar et al., 2025): *Relevance (R)*, *Correctness (C)*, *Actionability (A)*, *Uniqueness (U)*, and *Explainability (E, which is newly introduced based on DCG’s structure)*. (2) **DCG’s Contribution**: We randomly select another 10 target harmful memes with generated target guidance, then calculate the time costs that previous evaluators take to identify their harmfulness with/without DCG’s help. Note that we only count samples that can be successfully analyzed by humans. We have presented the questionnaire in Appendix A.6.

Figure 5 shows the results of human evaluation. We can see that, for the overall DCG’s quality evaluation, the DCG’s five-dimensional scores outperform the fail reason trees. All the scores are nearly or over 4, which means the DCG has high uniqueness and can be used to guide the detection of harmful memes. Moreover, the explainability of the DCG (score: 4.4) outperforms that of the trees (score: 1.8) by the largest margin (score: +2.6), which means the DCG is easier to understand by evaluators with different experiences. For the DCG’s contribution, the time costs of human evaluation show that DCG can help evaluators improve their efficiency in accurately analyzing why memes are harmful. Especially for the Offensiveness and Sarcasm memes that implicitly express the harmful-

ness, and humans cannot easily find their harmful elements, the DCG can reduce the average time costs from 15~30 seconds per meme. Therefore, REPMD can help humans identify harmful memes. Experiments in Appendix A.7 show that DCG can also help safeguard text-to-image models.

## 6 Related Works

Harmful meme detection has emerged as an important research direction in multimodal semantic analysis. Previous studies extract visual and textual feature embeddings and fuse them with attention mechanisms (Suryawanshi et al., 2020; Pramanick et al., 2021; Lee et al., 2021) and Transformer-based intermediate fusion (Kiela et al., 2019; Lu et al., 2019; Li et al., 2019; Chen et al., 2020), and contrastive reweighting approaches (Kiela et al., 2020; Muennighoff, 2020; Lippe et al., 2020). With the development of MLLMs, researchers proposed prompt-based (Cao et al., 2022; Ji et al., 2023; Cao et al., 2023; Ji et al., 2024) and multi-agent-based frameworks for deeper visual understanding (Lin et al., 2025). Considering the bias in MLLM’s inference steps, researchers proposed LLM debate frameworks (Hee et al., 2022) to improve the detection performance. To address challenges in low-resource generalization, researchers have explored few-shot enhanced methods with SFT and RAG methods (Cao et al., 2024; Huang et al., 2024; Liu et al., 2025). However, harmful memes are ever-shifting on the Internet, but few of these works concerns on improving their generalizability. Different from these works, REPMD introduces the explainable design concepts to facilitate the detection of ever-shifting harmful memes.

## 7 Conclusion

In this paper, we propose REPMD, an ever-shifting harmful meme detection method. We first define DCG, including the harmful types and reproduction steps of how users design harmful memes. Then, we derive the DCG from historical memes. Finally, we retrieve and utilize it to guide the MLLMs to detect whether the target meme is harmful. The evaluation results show that REPMD can detect harmful memes with 81.1% accuracy and can be generalized to type-shifting and temporal-evolving memes. Human evaluation shows that REPMD can improve the efficiency of human discovery on harmful memes, with 15~30 seconds per meme.

## Limitations

Although REPM is effective in detecting ever-shifting harmful memes, there are some cases that REPM cannot accurately identify their harmfulness. We manually investigate these bad cases and discuss the reasons for their failures.

First, 60% of the bad cases come from overly simple meme features. For example, the Minecraft group users always use the cubes to represent humans and objects, but MLLMs have a limited ability to learn these expressions. In the future, we plan to combine the MLLM fine-tuning and DCG to learn these features. Second, 35% of the cases come from the LLM hallucinations. MLLM may violate the guidance of DCG during the inference process and output guessed results, leading to errors in harmful meme detection. We plan to introduce the stepwise hallucination corrector to revise the correct biases in the responses. Finally, 5% of the cases come from newly designed memes with completely unseen design concepts. The proportion of these memes is small, which has relatively less impact on the detection performance. We also plan to continuously improve DCG for these cases.

## Ethical Statement

Our study aims to detect ever-shifting harmful memes on the Internet, but it may raise ethical considerations. We will discuss how these considerations are mitigated as follows. First, the dataset may have privacy and license issues because it comes from GOAT-Bench and online Twitter posts. We control permissions of visitors to the dataset through providing private dataset links and password access, and continuously track the whereabouts of the data to prevent data abuse. Second, the algorithms may have fairness issues because there are differences in the protection of individual freedom of speech in different countries. We publicly disclose the DCG and its related guidance, and plan to provide channels for misjudgment on our public website. Third, the annotators and evaluators may have mental health issues when labeling the dataset. We have provided regular psychological counseling to annotators, established a job rotation system, and provided reasonable salaries.

## Acknowledgement

We sincerely appreciate all the reviewers for their constructive suggestions. This work was supported by the National Key Research and Develop-

ment Program of China (No.2024YFF0618800), National Natural Science Foundation of China Grant No.62402484, No.62232016, Youth Innovation Promotion Association Chinese Academy of Sciences, and Basic Research Program of ISCAS Grant No.ISCAS-JCZD-202405.

## References

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Ming-Hsuan Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. [Qwen2.5-vl technical report](#). *CoRR*, abs/2502.13923.
- Xuheng Cai, Chao Huang, Lianghao Xia, and Xubin Ren. 2023. [Lightgcl: Simple yet effective graph contrastive learning for recommendation](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- 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.
- Rui Cao, Roy Ka-Wei Lee, Wen-Haw Chong, and Jing Jiang. 2022. Prompting for multimodal hateful meme classification. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 321–332.
- Rui Cao, Roy Ka-Wei Lee, and Jing Jiang. 2024. [Modularized networks for few-shot hateful meme detection](#). In *Proceedings of the ACM Web Conference 2024*, WWW '24, page 4575–4584, New York, NY, USA. Association for Computing Machinery.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer.
- Zixin Chen, Hongzhan Lin, Kaixin Li, Ziyang Luo, Zhen Ye, Guang Chen, Zhiyong Huang, and Jing Ma. 2025. Adammeme: Adaptively probe the reasoning capacity of multimodal large language models on harmfulness. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, page 4234–4253.
- Hang Gao, Chenhao Zhang, Fengge Wu, Changwen Zheng, Junsuo Zhao, and Huaping Liu. 2025. Bootstrapping heterogeneous graph representation learning via large language models: A generalized approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 16717–16726.

- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323. JMLR Workshop and Conference Proceedings.
- Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, Jianyu Jiang, Jiawei Wang, Jingji Chen, Jingjia Huang, Kang Lei, Liping Yuan, Lishu Luo, Pengfei Liu, Qinghao Ye, Rui Qian, Shen Yan, and 81 others. 2025. [Seed1.5-v1 technical report](#). *CoRR*, abs/2505.07062.
- Ming Shan Hee, Roy Ka-Wei Lee, and Wen-Haw Chong. 2022. On explaining multimodal hateful meme detection models. In *Proceedings of the ACM web conference 2022*, pages 3651–3655.
- Jianzhao Huang, Hongzhan Lin, Liu Ziyang, Ziyang Luo, Guang Chen, and Jing Ma. 2024. Towards low-resource harmful meme detection with lmm agents. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2269–2293.
- Junhui Ji, Xuanrui Lin, and Usman Naseem. 2024. Capalign: Improving cross modal alignment via informative captioning for harmful meme detection. In *Proceedings of the ACM Web Conference 2024*, pages 4585–4594.
- Junhui Ji, Wei Ren, and Usman Naseem. 2023. [Identifying creative harmful memes via prompt based approach](#). In *Proceedings of the ACM Web Conference 2023*, WWW '23, page 3868–3872. Association for Computing Machinery.
- Douwe Kiela, Suvrat Bhooshan, Hamed Firooz, Ethan Perez, and Davide Testuggine. 2019. Supervised multimodal bitransformers for classifying images and text. *arXiv preprint arXiv:1909.02950*.
- 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](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*.
- Roy Ka-Wei Lee, Rui Cao, Ziqing Fan, Jing Jiang, and Wen-Haw Chong. 2021. Disentangling hate in online memes. In *Proceedings of the 29th ACM international conference on multimedia*, pages 5138–5147.
- Bohao Li, Yuying Ge, Yi Chen, Yixiao Ge, Ruimao Zhang, and Ying Shan. 2024a. Seed-bench-2-plus: Benchmarking multimodal large language models with text-rich visual comprehension. *arXiv preprint arXiv:2404.16790*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Zhengpin Li, Mengzhe Jia, Zheng Wei, and Jian Wang. 2024b. Beyond smoothness: A general optimization framework for graph neural networks with negative laplacian regularization. *Neural Networks*, 180:106704.
- Hongzhan Lin, Ziyang Luo, Bo Wang, Ruichao Yang, and Jing Ma. 2024. [Goat-bench: Safety insights to large multimodal models through meme-based social abuse](#). *CoRR*, abs/2401.01523.
- 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 *Proceedings of the ACM on Web Conference 2025*, pages 4734–4744.
- Phillip Lippe, Nithin Holla, Shantanu Chandra, Santhosh Rajamanickam, Georgios Antoniou, Ekaterina Shutova, and Helen Yannakoudakis. 2020. A multimodal framework for the detection of hateful memes. *arXiv preprint arXiv:2012.12871*.
- Ziyang Liu, Chunxiao Fan, Haoran Lou, Yuexin Wu, and Kaiwei Deng. 2025. [MIND: A multi-agent framework for zero-shot harmful meme detection](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 923–947. Association for Computational Linguistics.
- Chaochao Lu, Chen Qian, Guodong Zheng, Hongxing Fan, Hongzhi Gao, Jie Zhang, Jing Shao, Jingyi Deng, Jinlan Fu, Kexin Huang, Kunchang Li, Lijun Li, Limin Wang, Lu Sheng, Meiqi Chen, Ming Zhang, Qibing Ren, Sirui Chen, Tao Gui, and 17 others. 2024. [From GPT-4 to gemini and beyond: Assessing the landscape of mllms on generalizability, trustworthiness and causality through four modalities](#). *CoRR*, abs/2401.15071.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32.
- Abdullah Mazhar, Zuhair Hasan Shaik, Aseem Srivastava, Polly Ruhnke, Lavanya Vaddavalli, Sri Keshav Katragadda, Shweta Yadav, and Md Shad Akhtar. 2025. Figurative-cum-commonsense knowledge infusion for multimodal mental health meme classification. In *Proceedings of the ACM on Web Conference 2025*, pages 637–648.
- Jingbiao Mei, Jinghong Chen, Guangyu Yang, Weizhe Lin, and Bill Byrne. 2025. [Robust adaptation of large multimodal models for retrieval augmented hateful meme detection](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 23817–23839. Association for Computational Linguistics.
- Niklas Muennighoff. 2020. Vilio: State-of-the-art visiolinguistic models applied to hateful memes. *arXiv preprint arXiv:2012.07788*.

Delfina Sol Martinez Pandiani, Erik Tjong Kim Sang, and Davide Ceolin. 2025. ‘toxic’ memes: A survey of computational perspectives on the detection and explanation of meme toxicities. *Online Soc. Networks Media*, 47:100317.

Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. 2024. Graph retrieval-augmented generation: A survey. *ACM Transactions on Information Systems*.

Jorge E. Pérez, Jessica Díaz, Javier García Martín, and Bernardo Tabuenca. 2020. Systematic Literature Reviews in Software Engineering - Enhancement of the Study Selection Process Using Cohen’s Kappa Statistic. *J. Syst. Softw.*, 168:110657.

Shraman Pramanick, Shivam Sharma, Dimitar Dimitrov, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021. MOMENTA: A multimodal framework for detecting harmful memes and their targets. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4439–4455. Association for Computational Linguistics.

Minodora Sălcudean. 2020. Visual humor through internet memes (ii) from harmless humour to the discriminatory potential of (anti) memes. case study: “the transgender bathroom debate”. *Revista transilvania*.

Bruce Schneier. 1999. Attack trees. *Dr. Dobbs’s journal*, 24(12):21–29.

Shivam Sharma, Firoj Alam, Md Shad Akhtar, Dimitar Dimitrov, Giovanni Da San Martino, Hamed Firooz, Alon Halevy, Fabrizio Silvestri, Preslav Nakov, and Tanmoy Chakraborty. 2022. Detecting and understanding harmful memes: A survey. *arXiv preprint arXiv:2205.04274*.

Shardul Suryawanshi, Bharathi Raja Chakravarthi, Michael Arcan, and Paul Buitelaar. 2020. Multimodal meme dataset (multioff) for identifying offensive content in image and text. In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying, TRAC@LREC 2020*, pages 32–41. European Language Resources Association (ELRA).

Carlo M Valensise, Alessandra Serra, Alessandro Galeazzi, Gabriele Etta, Matteo Cinelli, and Walter Quattrocchi. 2021. Entropy and complexity unveil the landscape of memes evolution. *Scientific Reports*, 11(1):20022.

Wenjun Xiong and Robert Lagerström. 2019. Threat Modeling - A Systematic Literature Review. *Comput. Secur.*, 84:53–69.

Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V. Chawla. 2019. Heterogeneous graph neural network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019*, pages 793–803. ACM.

Table 7: The macro type and example subtypes.

Macro-Type	Example Subtypes in DCG’s Type Part
Nationality	Countryhumans Historical Event Parody National Stereotype
Gender	Gender Role Reversal Implicit Misogyny Slangs Performative-Male Symbols Transgender Symbols
Religion	Holiday & Ritual Meme Templates Scripture & Figure Parody Islamic Muslim Symbols Buddhist Symbols
Human	Racism Stereotype Disability Stereotype
Animal	Personified Animal Memes Animal Behavior as Metaphor Mimicry in Nature as Meme Template
Culture	Nonsense Literature Abbreviation Culture Versailles Literature Video Game
Political	Politician as Meme Template Political Figure’s Biological Reduction Policy & Event Satire

## A Appendix

### A.1 Detailed Definition of DCG

**Type Node.** In this part, we predefined seven macro types for the meme classification, i.e., *Nationality, Gender, Religion, Human, Animal, Culture, and Political*. Based on this classification system, we can classify the Internet memes into those that represent what the users want to express.

For each macro type, we ask the MLLM to divide it into subtypes to specifically illustrate the details in these memes. Table 7 shows the corresponding relationship between macro types and the subtypes. In this table, we only present the example subtypes in DCG’s type nodes, where the subtypes will be extended when new memes occur.

With the previous classification system, the macro types and subtypes will be mapped to the type node  $\mathcal{N}_{\text{type}}$ ’s three-level information. The prompts  $P_T$  for type node generation (Figure 10) will control the subtype extension in the meme reason tree and DCG construction.

**Reproduction Nodes.** In this part, we have defined the design concept, which is the idea and steps that users convey their harmfulness. We need to reproduce their ideas by observing the visible meme’s elements, and the graph represents the MLLM’s inner understanding of this meme.

To illustrate the details of the reproduction nodes in the DCG, we have provided an example in Figure

3. In this example, the DCG incorporates four reproduction step nodes and one reproduction goal. The macro type is the "Gender" and the subtype is the "Implicit Misogyny Slangs".

**Representation of Reproduction Steps.** To represent the graph in the MLLM’s input with plain text together with logic-based semantics, we use the logic expression, which is concatenated by logic symbols ( $\wedge, \vee, \neg, \rightarrow$ ). The example’s reproduction steps can be formulated as follows:

$$\mathcal{N}_{M_1} \wedge (\mathcal{N}_{M_1} \rightarrow \mathcal{N}_{M_2} \rightarrow \mathcal{N}_{M_3}) \rightarrow \mathcal{N}_G \quad (2)$$

where the  $\mathcal{N}_{\text{rep}}^{M_1}$  to  $\mathcal{N}_{\text{rep}}^{M_3}$  means the content in the reproduction nodes, and the  $\mathcal{N}_{\text{rep}}^G$  means the content in the reproduction goal. In the guidance of MLLM’s inference of Section 3.3, the input logic expression utilizes the "then", "and", "or", and "on the contrary" to replace the logic symbols. Based on this criterion, the DCG is parsed into text according to the following rules.

- **Parsed Text of Reproduction Method:** Directly use the reproduction method’s text.
- **Parsed Text of Harmful Indicator:** Append the parentheses with labels (i.e., "harmful" and "no harmful") after reproduction methods.
- **Parsed Text of Logic Gate and Achievement Edge:** Append the conjunctions (i.e., "then", "and", "or", and "on the contrary") between texts of reproduction methods.

Therefore, the input design concept with the logic expression can be illustrated as follows:

**## Harmful Type**  
 Macro Type: Gender; Subtype: Implicit Misogyny Slangs.  
**## Reproduction Step**  
*The malicious user writes a fact text (no harmful), then search the related images based on the fact (no harmful). Then, it specify the image to the biased people (woman, harmful). The fact and people specification simultaneously cause the attack goal (harmful): Specify fact to person, express misogyny to people emotion.*

## A.2 Mathematical Rationality for Algorithm 1

### A.2.1 Mathematical Rationality of SVD

We consider a symmetric matrix  $\mathbf{A}$  of the form:

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{\text{Type}} & \beta \mathbf{E}_{\text{Link}} \\ \beta \mathbf{E}_{\text{Link}} & \alpha \mathbf{A}_{\text{Repr}} \end{pmatrix} \quad (3)$$

where  $\mathbf{A}_{\text{Type}}$  and  $\mathbf{A}_{\text{Repr}}$  are symmetric matrices representing the adjacency or feature similarity within

the type nodes and the reproduction nodes, respectively;  $\mathbf{E}_{\text{Link}}$  is the connection matrix between the two parts; and  $\alpha$  and  $\beta$  are scaling factors satisfying  $1 > \alpha > \beta > 0$ .

**Effect of Scaling Factors.** To show that the scaling factors  $\alpha$  and  $\beta$  lead to smaller singular values for the reproduction nodes, we analyze the Rayleigh quotient of  $\mathbf{A}$ . Let  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  be a vector partitioned according to the two parts. We can derive the equation:

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}_1^\top \mathbf{A}_{\text{Type}} \mathbf{x}_1 + 2\beta \mathbf{x}_1^\top \mathbf{E}_{\text{Link}} \mathbf{x}_2 + \alpha \mathbf{x}_2^\top \mathbf{A}_{\text{Repr}} \mathbf{x}_2 \quad (4)$$

Consider a vector  $\mathbf{x}$  that is nonzero only in the reproduction nodes, i.e.,  $\mathbf{x} = (\mathbf{0}, \mathbf{x}_2)$  with  $|\mathbf{x}_2| = 1$ , we can derive the equation:

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \alpha \mathbf{x}_2^\top \mathbf{A}_{\text{Repr}} \mathbf{x}_2 \leq \alpha \lambda_{\max}(\mathbf{A}_{\text{Repr}}) \quad (5)$$

where  $\lambda_{\max}(\mathbf{A}_{\text{Repr}})$  is the largest eigenvalue of  $\mathbf{A}_{\text{Repr}}$ . By the Courant-Fischer minimax theorem, the largest singular value of  $\mathbf{A}$  (which equals the largest eigenvalue in absolute value, since  $\mathbf{A}$  is symmetric) satisfies:

$$\max_{\|\mathbf{x}\|=1} |\mathbf{x}^\top \mathbf{A} \mathbf{x}| \geq \alpha \lambda_{\max}(\mathbf{A}_{\text{Repr}}) \quad (6)$$

For vectors concentrated on the reproduction nodes, the Rayleigh quotient scales with  $\alpha$ , i.e., the singular values associated with the reproduction nodes are roughly proportional to  $\alpha$ . Similarly, for vectors that involve both parts, the Rayleigh quotient scales with  $\beta$ . Since  $\alpha < 1$  and  $\beta < 1$ , the scaling reduces the Rayleigh quotient for vectors involving the reproduction nodes, leading to **smaller singular values** for those nodes.

**Low-rank Approximation and Pruning.** In the SVD-based pruning, we approximate  $\mathbf{A}$  by a low-rank matrix  $\tilde{\mathbf{A}}_k = \mathbf{U}_k \Sigma_k \mathbf{V}_k^\top$ , where only the  $k$  largest singular values are retained. Since the scaling factors  $\alpha$  and  $\beta$  decrease the singular values corresponding to the reproduction nodes, these values are more likely to be truncated in the low-rank approximation. Consequently, the edges within the parts are identified as redundant and can be pruned.

Thus, the scaling strategy with  $\alpha$  and  $\beta$  is mathematically justified as it systematically reduces the influence of the reproduction nodes in the singular value spectrum, enabling reproduction nodes’ pruning while preserving the type nodes.

## A.2.2 Mathematical Rationality for Logarithmic Cut-off Determination

**Properties of Long-Tail Distributions.** Let  $\{x_i\}_{i=1}^n$  be a set of  $n$  positive real-valued observations assumed to have a long-tail distribution. After sorting in descending order, we denote:  $x_1 \geq x_2 \geq \dots \geq x_n > 0$ . Our target is to find a threshold  $x_c$  such that:

- $x_i \geq x_c$  belongs to the "head" region.
- $x_i < x_c$  belongs to the "tail" region.

The long-tail distributions exhibit the fundamental property that their generating processes are *multiplicative* rather than additive. This leads to the following characterization:

**Theorem 1** (Logarithmic Normalization). *For a random variable  $X$  following a long-tail distribution, the logarithmic transform  $Y = \ln(X)$  typically follows a distribution with exponentially decaying tails (e.g., normal distribution for log-normal, exponential for power-law).*

This theorem converts multiplicative differences in the original scale to **additive** differences in the **logarithmic** scale.

**Extreme Value Behavior in Logarithmic Space.** Let  $Y_1 \geq Y_2 \geq \dots \geq Y_n$  be the order statistics of the log-transformed data with a long-tail distribution. The spacings between consecutive order statistics satisfy:

**Lemma 1** (Spacing Distribution). *For observations from a distribution with continuous density  $f_Y(y)$ , the normalized spacings  $n \cdot (Y_i - Y_{i+1}) \cdot f_Y(Y_i)$  converge to independent standard exponential random variables as  $n \rightarrow \infty$ .*

This lemma suggests that in regions where the density  $f_Y(y)$  is relatively constant, the gaps  $\Delta_i$  should be approximately equally sized. At the boundary between the head and tail regions, we expect a **discontinuity** in the density, leading to an **anomalously large gap**.

**Optimality of Maximum Gap Selection.** We provide the theorem for maximum gap optimality:

**Theorem 2** (Maximum Gap Optimality). *Assume the data generation follows two regimes:*

- **Head regime:**  $Y \sim F_1$  with support on  $[a, \infty)$
- **Tail regime:**  $Y \sim F_2$  with support on  $(-\infty, b]$
- With  $a > b$  and minimal overlap between regimes

As  $n \rightarrow \infty$ , the probability that the maximum gap occurs at the true boundary converges to 1.

Then, we provide the proof for this theorem:

*Proof.* Let  $Y_{k^*}$  and  $Y_{k^*+1}$  be the order statistics straddling the true boundary. The expected gap at this position is:

$$\mathbb{E}[\Delta_{k^*}] = \mathbb{E}[Y_{k^*} - Y_{k^*+1}] \geq a - b + o(1) \quad (7)$$

whereas for position  $i \neq k^*$ , the expected gap is:

$$\mathbb{E}[\Delta_i] = O\left(\frac{1}{n \cdot f_Y(Y_i)}\right)$$

Since  $a - b > 0$  represents a fixed separation between regimes, and the other gaps shrink as  $O(1/n)$ , we have:

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\Delta_{k^*} > \max_{i \neq k^*} \Delta_i\right) = 1$$

Thus, the maximum gap asymptotically identifies the true boundary.  $\square$

## A.3 Case Study of DCG Construction and Harmful Meme Detection

In this section, we illustrate the case of how the DCG is pruned in the REPMD's second step, including the node extension, score calculation in the SVD, and the details are shown in Figure 6.

In this case, we can see that REPMD first analyzes the surface-level elements of the two historical memes from Figure 1's motivation example. The malicious users incorporate the "crazy woman", "parenting issues", "nose stud", and "black people" in these two memes, which are apparently shown in the elements.

Then, REPMD extends the elements to the inner ideas, which are the reproduction steps of the design concept that reflect what users have done to design such harmful memes. The two memes are extended to the DCG with 13 nodes (nine reproduction step nodes and four type nodes). From bottom to top, the DCG represents what the users have done to express misogyny and racism. However, the graph contains too many nodes and edges, which makes it difficult to understand.

To reduce the DCG's scale, we prune it based on the SVD, which is calculated based on the node's similarity and root correlation. For the node's similarity, we find that the reproduction methods like "select woman's image" and "select black people" have high TF-IDF similarity, where "circle

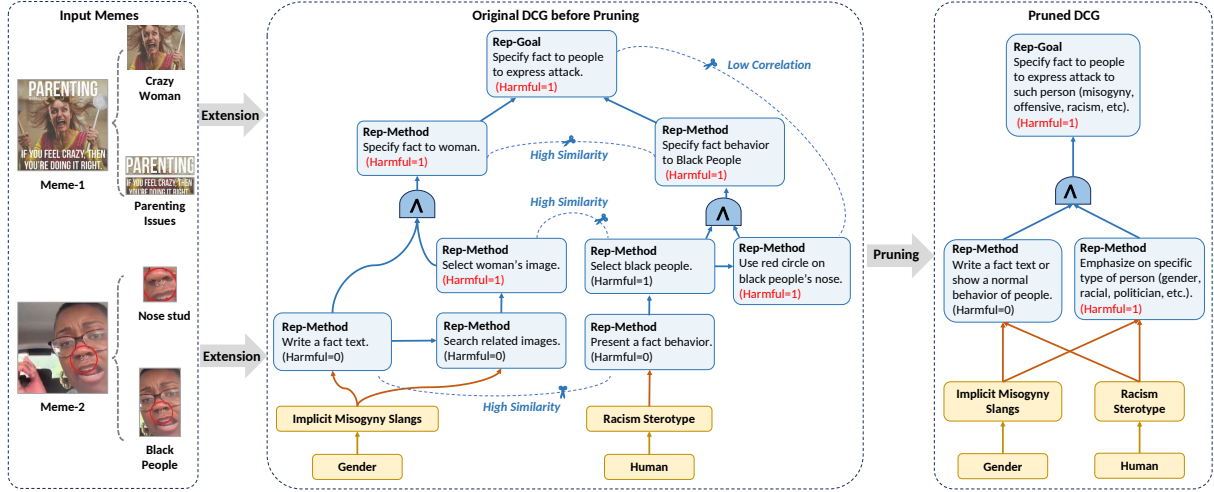


Figure 6: The case of REPMD’s DCG pruning.

the black people’s nose" and the "specify fact to people" seem to have low correlation that need to be calibrated and pruned. After SVD pruning and merging, these edges and nodes are removed in the DCG, which only keeps the core part that has seven nodes to represent the user’s design concept.

#### A.4 Other Experimental Results

**Pruning Time Cost of SVD Replacement.** To evaluate the time cost between SVD and MLLMs for graph pruning, we separately evaluate the time the REPMD takes when pruning the historical memes  $\{M_1^H, \dots, M_n^H\}$ . The time cost is calculated with the following equation:

$$Cost_{\text{time}} = Sec(\mathcal{G}_R \rightarrow \mathcal{G}'_R) / N(\{M_i^H\}) \quad (8)$$

where the DCG  $\mathcal{G}_R$  is prepared based on the historical meme  $M_i^H$ ’s analysis results. With this equation, we do not consider the time cost of graph generation, but only analyze how effectively the graph can be pruned. The tradeoff between the model performance and the efficiency illustrates the benefit of REPMD.

Table 8: The pruning time cost of variants in SVD replacement (seconds per meme for DCG pruning).

Variants	ID	OOD	TE
RepMD w/ SVD	4.9	6.1	5.3
SVD→GPT-4o	62.3 (13×)	61.3 (10×)	52.9 (10×)
SVD→Doubao-V-Pro	Failed	71.2 (12×)	69.7 (13×)
SVD→Qwen3VL	Failed	Failed	75.5 (14×)

Table 8 illustrates the pruning time cost of the SVD and the other MLLMs. We can see that, MLLMs can prune the DCG, but they have too much time cost, which is over 10× to the SVD. In practice, these methods are not applicable, which

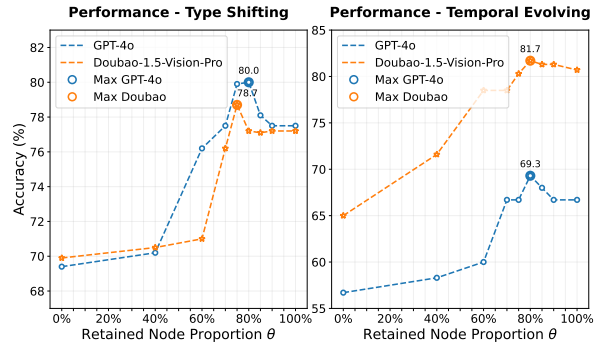


Figure 7: The effect of the retained node’s proportion  $\theta$ .

may come from the fact that the input graph has too many tokens, as well as the long CoT and inference steps that reduce their efficiency. However, we can also see that Qwen3VL-235B can slightly improve the detection performances in the TE task, which means if we can use SVD’s results as ground-truth and fine-tune these MLLMs to an acceptable time range (around 2× to 5× to SVD), the MLLMs may be useful to improve REPMD’s performance.

**Hyperparameter Analysis.** We also evaluate the performance of REPMD when we tune the Retained Nodes’ Proportion  $\theta$  from 0% ~ 100% with a 5% interval. We evaluate the performance of GPT-4o and Doubao-1.5-Vision-Pro on these intervals, where we calculate the **average results** of the OOD and TE experiment.

Figure 7 shows the results of hyperparameter tuning. We can see that, most of the highest value comes from  $\theta = 75\%$  and 80%, where the difference between these two parameters is slightly different. We also manually check the values across different types and temporal quarters, and find that



Target Meme Case 1	DCG-based Guidance
 <p>Misogyny</p>	<p>S1: The image is divided into top panel depicting a robotic figure holding a weapon and saying "men's rights are only equal to men's rights," which implies a self-centered or exclusive view of rights that excludes others.</p> <p>S2: The bottom panel shows the same robot now facing a yellow burst of light labeled "EQUAL RIGHTS EQUAL FIGHTS," suggesting a shift toward equality but using aggressive imagery to represent this change.</p> <p>S3: The robot's design and posture—militant, dominant, and imposing—visually associate gender equality with dominance and confrontation.</p> <p>S4: The use of the word "men" in the top panel without context or balance, combined with the absence of female representation, subtly positions male perspectives as default or primary, normalizing male-centric views.</p> <p>S5: The contrast between the two panels suggests that equality is achieved through force or struggle, implying that women's rights must be won through conflict rather than mutual recognition.</p>
<b>Output</b>	
<p><b>Output Result:</b> By centering a masculine figure as both the source and enforcer of rights discourse, the image implicitly marginalizes women's agency and promotes a competitive, adversarial view of gender equity. This reinforces harmful stereotypes that equate women's rights with threats to men's status, thereby expressing misogynistic undertones under the guise of advocating fairness. The rhetoric of "equal fights" further reduces equality to conflict, suggesting that women must fight aggressively to gain rights, which perpetuates a culture of hostility rather than cooperation. Thus, the image conveys an <b>misogynistic</b>.</p>	
Target Meme Case 2	DCG-based Guidance
 <p>Sarcasm</p>	<p>S1: The image uses a panda bear as the central visual subject, a creature that is biologically black and white in coloration, but which is also culturally associated with Asia due to its native habitat in China.</p> <p>S2: The text "He's Black / He's White / He's Asian" directly references the panda's physical appearance (black and white fur) while simultaneously assigning racial categories to it, creating a juxtaposition.</p> <p>S3: The phrase "Be like panda" at the top suggests an invitation to emulate the panda.</p> <p>S4: The inclusion of a cartoon man with glasses and a surprised expression on the right side visually contrasts with the panda, suggesting a human observer reacting to the absurdity or contradiction presented. This reinforces the humorous or ironic tone.</p> <p>S5: The title "Destoy Racism" (likely a misspelling of "Destroy Racism") sets up an apparent moral message against racism, but the content undermines this by using racial categorization in a trivial, reductive way.</p>
<b>Output</b>	
<p><b>Output Result:</b> The image employs sarcasm by presenting a seemingly anti-racist message ("Destoy Racism") while simultaneously using a non-human animal to mock the concept of race through oversimplification and absurdity. The design leverages irony to critique the rigid and often arbitrary nature of racial classification by applying it to a panda—an animal whose colors are natural and not socially constructed. The humor arises from the contrast between the serious intent of the title and the trivial, illogical application of race to a panda. This satirical approach exposes how easily race can be misused or misunderstood, but it does so in a way that may also trivialize real struggles related to racial identity. Therefore, the image functions as <b>sarcasm</b>—yet potentially risking the minimization of actual racial issues through humor.</p>	

Figure 8: The case study of REPMO of other two meme’s Qwen2.5VL’s outputs.

the highest results all distribute in the proportion within 70% ~ 80%. Therefore, we choose the balance value for all the experiments, i.e., 75%, for the REPMO’s SVD pruning.

### A.5 Case Study of Other Harmful Memes

In Figure 8, we provide two more memes’ Qwen2.5VL’s output based on our DCG guidance. The output can correct the misleading analysis (e.g., the original racism prediction in case 2) from these memes and output the correct harmful types (e.g., sarcasm of case 2), illustrating its usefulness in harmful meme detection.

### A.6 Questionnaire for Human Evaluation

We have provided the content of the questionnaire provided to evaluators, including the introduction of the 5-point Likert scale, the detailed definition for the five criteria, as well as the core dimensions for evaluators to rate our DCG and memes. We also suggest evaluators provide open feedback to our DCG’s improvement, so REPMO can be continuously optimized.

### Questionnaire: Evaluation of Design Concept Graph (DCG) for Harmful Meme Detection

**Instructions:** Please evaluate the provided **Design Concept Graph (DCG)** in relation to the corresponding **meme image** based on the following criteria. For each statement, select the response that best matches your assessment using the **5-point Likert scale**, where:

- 1 = Strongly Disagree
- 2 = Disagree
- 3 = Neutral (neither good nor bad / borderline)
- 4 = Agree
- 5 = Strongly Agree

**Note:** A score of 3 represents the threshold between positive and negative evaluation. Scores above 3 indicate a favorable assessment, while scores below 3 indicate an unfavorable assessment.

#### Part 1: Criteria Definition

- (1) **Relevance (R):** The DCG is closely related to the content and meaning of the meme image.
- (2) **Correctness (C):** The DCG accurately represents the key elements and message of the meme image.
- (3) **Actionability (A):** The DCG provides clear and practical insights that could guide further actions (e.g., content moderation, design adjustments, or analysis).
- (4) **Uniqueness (U):** The DCG is distinct and differs meaningfully from graphs of other similar meme images.
- (5) **Explainability (E):** The DCG effectively explains how the harmful meme image was designed and reflects the creator’s intent or thought process.

#### Part 2: Core Dimensions

No.	Dimension	Rating (1–5)
1	Relevance (R)	1 2 3 4 5
2	Correctness (C)	1 2 3 4 5
3	Actionability (A)	1 2 3 4 5
4	Uniqueness (U)	1 2 3 4 5
5	Explainability (E)	1 2 3 4 5

**Part 3: Open Feedback (Optional)** Please provide any additional comments or suggestions regarding the DCG or its evaluation:

Thank you for your participation!

### A.7 Safeguarding the Text-to-Image Models

After the previous human evaluation that illustrates the benefits of REPMO, we think further on the usefulness of REPMO, i.e., whether the DCG can cheat the Text-to-Image models and output harmful information. Theoretically, if one MLLM cannot correctly detect the harmfulness of the meme, it may be misled by these design concepts without security alignment. Inspired by it, if we can successfully jailbreak the MLLMs to generate such harmful memes, we may have the chance to protect them based on the design concepts:

**Step-1: Attack Prompt Construction:** We obtain the target harmful meme  $M_i^T$  and its stepwise guidance, then we directly form it as the prompt that will be input to the MLLM:

Table 9: The safeguarding of harmful image generation.

Metric	Generated Img. / Total Guidance		$\Delta$
	w/o DCG-Guard	w/ DCG-Guard	
Racism	7/12 (58.3%)	1/12 (8.3%)	50.0%
Misogyny	8/10 (80.0%)	1/10 (10.0%)	70.0%
Offensiveness	7/10 (70.0%)	1/10 (10.0%)	60.0%
Sarcasm	10/10 (100.0%)	0/10 (0.0%)	100.0%
Toxicity	1/10 (10.0%)	0/10 (0.0%)	10.0%

Metric	Average SSIM		$\Delta$
	w/o DCG-Guard	w/ DCG-Guard	
Racism	0.76	0.12	0.64
Misogyny	0.90	0.33	0.57
Offensiveness	0.54	0.19	0.35
Sarcasm	0.82	-	-
Toxicity	0.24	-	-

**## Task of Qwen-Image**

Please draw a meme graph with humor.

**## Reproduction Step**

The malicious user writes a fact text sentence like "parenting is too difficult and makes people crazy", then draws a crazy lady in the image.

**Step-2: DCG-Guard Safeguarding:** We ask another Qwen3VL-235B model to refer to the DCG, then safeguard the prompt and remove the biased information in the memes, and add the forbidden sentences. The sentences like "crazy lady in the image" will be replaced with commands like "stop generating such image because of misogyny", or replace them with other benign elements like "parents that are raising children".

We select the Qwen-Image as the target MLLM, which is a model that has not been successfully aligned and may generate harmful meme content, and randomly choose 10~12 cases as the target memes. We introduce two parts of the evaluation metrics. First, we evaluate how many images are successfully generated based on our attack prompt, and calculate the success rates of the attack prompts. Then, we evaluate the average SSIM, which is the metric that represents the similarity between AI-reconstructed images and original images. We evaluate the result with/without DCG-Guard on the five harmful types of memes.

Table 9 shows the results of safeguarding the harmful images through DCG-Guard. We can see that, before safeguarding the Qwen-Image, some attack prompts can achieve nearly 100% success rates (except for the Toxicity memes that contain apparent harmful features in the prompt). With the help of DCG-Guard, the success rates reduce by over 50%, where most of the harmful information in memes is detected and refused to generate with the stop commands. Moreover, even if some images are successfully generated, they may be

different from the previous images based on the benign element replacement, with over 0.3 SSIM decrease. These results show that DCG can also help safeguard text-to-image models.

**A.8 Prompts in REPM**

In this section, we show all the details of the prompts that we have used in REPM

- **Harmful meme detection:**  $P_{\text{Harm}}, P'_{\text{Harm}}$  in Figure 9, where  $P_{\text{Harm}}$  is the original prompt before generating the target guidance (the Vanilla MLLMs use this prompt), and  $P'_{\text{Harm}}$  is the prompt after generating the target guidance.
- **Fail reason tree construction:**  $P_{\text{F}}$  in Figure 10, which indicates the steps of how to find the fail cases and generate the reasons.
- **DCG construction:**  $P_{\text{D}}$  in Figure 11, which indicates how to generate and prune the nodes/edges in DCGs based on the basic structure.

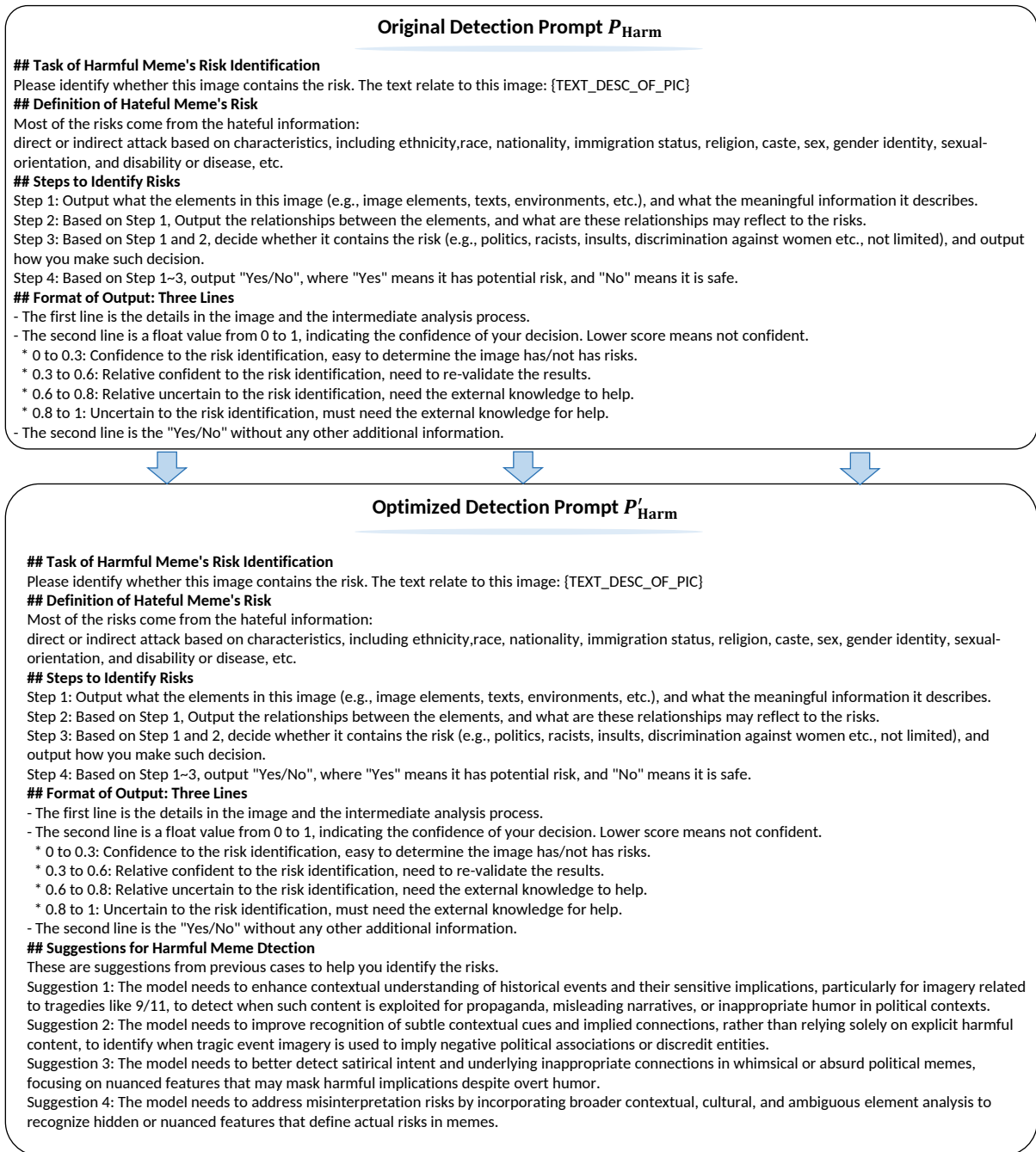


Figure 9: The detection prompt  $P_{\text{Harm}}$  and  $P'_{\text{Harm}}$ .

### Prompt of Fail Reason Tree Construction $P_F$

#### ## Objective:

Given a set of failed meme detection cases (Mfail), you are required to: Generate Reason Nodes (NF): Analyze and explain why the MLLMs failed to detect the harmful content in the meme. Generate Type Nodes (NT): Classify the meme into a macro-type (L1) and further stratify it into subtypes (L2, L3) based on the predefined typology.

#### ## Output Format:

Provide clear, structured reasoning for both the failure explanation and the typological classification.

#### ### Step 1: Reason Nodes Generation (Mfail → NF)

**Input:** Concatenated responses from all failed MLLMs for a given meme.

**Task:** Describe the content of the meme in detail. Analyze and explain why the MLLMs failed to identify the harmful elements. Focus on aspects such as: Subtlety or implicit nature of the harmful content. Cultural or contextual nuances that were overlooked. Limitations in the MLLMs' training data or detection algorithms. Visual or textual ambiguities that led to misinterpretation. Output a succinct, coherent reason node summarizing the failure analysis.

#### ### Step 2: Type Nodes Generation ((Mfail, NF) → NT)

**Input:** The meme description, the generated reason node (NF), and the predefined typology hierarchy.

**Task:** First, classify the meme into one macro-type (L1) from the following list: **Nationality, Gender, Religion, Human, Animal, Culture, Political.**

**Then, stratify the classification into subtypes:** L2: A subcategory under the chosen macro-type. L3: A further refined subcategory under L2 (if applicable).

Ensure the classification is consistent with the content description and the failure analysis. Provide a brief justification for each level of classification.

#### ### Predefined Typology Hierarchy (for reference):

##### Nationality:

L2: Stereotyping, Discrimination, Xenophobia, National Conflict. L3: (e.g., under Stereotyping: Generalization, Caricature, etc.)

##### Gender:

L2: Sexism, Gender Stereotyping, LGBTQ+ Discrimination, Sexual Harassment. L3: (e.g., under Sexism: Objectification, Role Enforcement, etc.)

##### Religion:

L2: Blasphemy, Religious Stereotyping, Inter-religious Conflict, Persecution. L3: (e.g., under Blasphemy: Sacred Symbol Misuse, Doctrine Mockery, etc.)

##### Human:

L2: Dehumanization, Violence, Marginalized Groups, Disability Mockery. L3: (e.g., under Dehumanization: Comparison to Objects, Animalization, etc.)

##### Animal:

L2: Animal Cruelty, Speciesism, Wildlife Exploitation, Pet Abuse. L3: (e.g., under Animal Cruelty: Graphic Harm, Neglect, etc.)

##### Culture:

L2: Cultural Appropriation, Tradition Mockery, Ethnic Stereotyping, Heritage Disrespect. L3: (e.g., under Cultural Appropriation: Symbol Misuse, Ritual Trivialization, etc.)

##### Political:

L2: Propaganda, Hate Speech, Conspiracy Theories, Leader Mockery.

L3: (e.g., under Propaganda: Misinformation, Manipulative Imagery, etc.)

#### ## Output Structure:

**Meme Description:** [Clear description of the meme's content]

**Failure Reason (NF):** [Concise analysis of why detection failed]

**Typology Classification (NT):** L1 (Macro-type): [Selected from the predefined list] L2 (Subtype): [Justified subcategory] L3 (Refined Subtype): [Further refinement if applicable]

**Justification:** [Brief explanation for the classification]

#### ### Note

Ensure the analysis is objective, detailed, and grounded in the provided inputs. The classification should align with the harmful elements identified (or missed) in the meme. Use the generated reason node (NF) to inform the typology classification (NT).

Figure 10: The fail reason tree's construction prompt  $P_F$ .

### Prompt of DCG Node Extension $P_D$

#### ## Task Overview

You are Qwen3VL-235B, tasked with extending failure reason nodes (NF) into comprehensive diagnostic graph nodes (GD). This process involves two main phases: Reproduction Method and Logic Gate Extension (analyzing meme construction) Graph Calibration and Goal Extension (validating and contextualizing the analysis)

**Input:** Failure reason node NF (containing analysis of why detection failed for a specific meme) and the original detection prompt PD

**Output:** A structured diagnostic graph node GD ready for integration into the overall Diagnostic Causal Graph (DCG)

#### ## Phase 1: Reproduction Method and Logic Gate Extension

##### Step 1.1: Surface-Level Element Extraction

Analyze the meme associated with the failure case NF and identify ALL surface-level elements: Visual elements: Characters, objects, symbols, colors, composition, facial expressions, gestures Textual elements: Captions, labels, speech bubbles, hashtags Structural elements: Layout, arrangement, visual hierarchy. For each element, describe: Its explicit appearance/representation; Its immediate semantic meaning (literal interpretation)

Combination Logic No Analysis: Describe how these surface elements combine to create meaning. Answer:

logical relationship between elements: and, or, not

##### Step 1.2: Inner-Level Design Method Analysis (NM)

For EACH identified surface element, conduct a three-question drill-down analysis:

**QUESTION A:** "Is there a replacement method for this element?"

Consider alternative elements that could serve similar structural or functional roles. Evaluate whether replacement would maintain the meme's format/template; Identify if multiple replacement options exist

**QUESTION B:** "Why is that element chosen (and not alternatives)?"

Analyze the specific cultural, emotional, or symbolic value of the chosen element. Consider why alternative elements would be less effective for the intended purpose. Examine how this element contributes to the meme's persuasive/rhetorical impact

**QUESTION C:** "Is the replaced element harmful?"

For each replacement option from Question A, assess whether it would still convey harmful content. If harm persists, describe how the harm mechanism changes or remains.

Identify which element properties are essential for conveying harm

**Output for Step 1:** Complete element inventory with surface descriptions; Combination logic No explanation; For each element: A/B/C analysis results

#### ## Phase 2: Graph Calibration and Goal Extension

##### Step 2.1: Design Concept Validation

Act as a validator (simulating "another Qwen3VL") to assess:

Does the design method analysis (NM) from Phase 1 accurately reflect the actual visual elements? Are there any misinterpretations or overinterpretations of elements? Does the combination logic No correctly describe how the meme functions?

**Calibration Checklist:** All identified elements actually exist in the meme. Element functions/meanings are not exaggerated or understated. Combination logic reflects actual viewer perception patterns. No significant elements were missed in Phase 1

##### Step 2.2: User Goal Inference (NG)

Based on the validated analysis, infer what the meme creator likely aimed to achieve:

**Primary Goals:** What immediate response was sought from viewers? (laughter, anger, agreement, etc.) What belief or attitude was being reinforced or challenged?

**Secondary Goals:** What social or cultural commentary was being made? What group was being targeted or elevated?

**Strategic Goals:** Why was this particular format/template chosen? How does the design maximize spread or impact?

**Output for Step 2:** Calibration report noting any discrepancies or confirmations. Comprehensive user goal analysis NG

Figure 11: The DCG construction's prompt  $P_D$ .