UNIGUARD: TOWARDS UNIVERSAL SAFETY GUARDRAILS FOR JAILBREAK ATTACKS ON MUL TIMODAL LARGE LANGUAGE MODELS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025 026 027

028 029

030

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) have revolutionized visionlanguage understanding but remain vulnerable to multimodal jailbreak attacks, where adversarial inputs are meticulously crafted to elicit harmful or inappropriate responses. We propose UNIGUARD, a novel multimodal safety guardrail that jointly considers the unimodal and cross-modal harmful signals. UNIGUARD trains a multimodal guardrail to minimize the likelihood of generating harmful responses in a toxic corpus. The guardrail can be seamlessly applied to any input prompt during inference with minimal computational costs. Extensive experiments demonstrate the generalizability of UNIGUARD across multiple modalities, attack strategies, and multiple state-of-the-art MLLMs, including LLaVA, Gemini Pro, GPT-40, MiniGPT-4, and InstructBLIP. Notably, this robust defense mechanism maintains the models' overall vision-language understanding capabilities. Our code is available at https://anonymous.4open.science/r/ UniGuard/README.md.

Warning: this paper contains inputs, data, and model behaviors that are offensive in nature.

1 INTRODUCTION

031 The rapid development of multimodal large language models (MLLMs), exemplified by models like 032 GPT-40 (OpenAI, 2023), Gemini (Reid et al., 2024), and LLaVA (Liu et al., 2023a;b), has revolu-033 tionized vision-language understanding but introduced new risks. Among the most pressing con-034 cerns is the vulnerabilities of MLLMs to adversarial attacks or *jailbreaks* (Qi et al., 2023; Shayegani 035 et al., 2023; Niu et al., 2024; Deng et al., 2024). These attacks exploit inherent weaknesses of mod-036 els to bypass safety mechanisms, resulting in the generation of toxic content and raising serious 037 challenges for secure deployment in high-stakes, user-facing domains such as education, clinical 038 diagnosis, and customer service.

Challenges. Ensuring safe and trustworthy interactions requires the development of robust safety 040 guardrails against adversarial exploitation, which presents three core challenges. 1) Multimodal Ef-041 fectiveness. Guardrails must protect against adversarial prompting in multiple modalities and their 042 cross-modal interactions, ensuring that defenses are not limited to unimodal threats. 2) General-043 *izability Across Models.* Safety mechanisms should be adaptable to multiple model architectures, 044 including both open-source and proprietary ones. 3) Robustness Across Diverse Attacks. Effective 045 guardrails must withstand a wide range of attack strategies, including constrained attacks that subtly modify inputs while maintaining visual similarity, and unconstrained attacks that introduce notice-046 able changes (Qi et al., 2023). They should also address adversarial text prompts (Gehman et al., 047 2020) that elicit harmful or inappropriate responses from LLMs. Although prior work has explored 048 defenses for both unimodal (Zou et al., 2023; Chao et al., 2023) and multimodal LLMs (Shayegani et al., 2023; Niu et al., 2024; Gou et al., 2024; Pi et al., 2024), a holistic approach covering multiple *modalities*, *models*, and *attack types* remains an open challenge. 051

This Work. We introduce UNIGUARD, a novel defense mechanism that provides robust,
 <u>Universally applicable multimodal Guard</u>rails against adversarial attacks in both visual and textual inputs. As shown in Figure 1, the core idea is to create specialized safety guardrail for individual



Figure 1: UNIGUARD robustifies multimodal large language models (MLLMs) against multimodal jailbreak attacks by using safety guardrails to purify malicious input prompt, ensuring safe responses.

076 modalities while accounting for their cross-modal interactions. This guardrail purifies potential ad-077 versarial responses after applying to input prompts. Inspired by few-shot prompt learning (Qi et al., 2023; Lester et al., 2021), we optimize the guardrails by searching for additive noise (for image 079 inputs) and suffix modifications (for text prompts) to minimize the likelihood of generating harmful responses in a small toxic content corpus (Liu et al., 2023a). We conduct comprehensive experi-081 ments on both adversarial and benign inputs. Our results demonstrate that UNIGUARD significantly improves robustness against various adversarial attacks while maintaining high accuracy for benign 083 inputs. For example, UNIGUARD effectively reduces the attack success rate on LLAVA by nearly 55%, with a small performance-safety trade-off in visual question-answering. The safety guardrails 084 developed for one model such as LLAVA (Liu et al., 2023a) is transferable to other MLLMs, includ-085 ing both open-source models like MiniGPT-4 (Zhu et al., 2023) and InstructBLIP (Dai et al., 2023), as well as proprietary models like Gemini Pro (Team et al., 2023) and GPT-40 (OpenAI, 2023), 087 highlighting the *generalizability* of our approach across different models and architectures. 088

089 Contributions. Our major contributions are:

090

094

096

098

099 100

- 1. Effective Defense Strategy. We propose UNIGUARD, a pioneering, universally applicable multimodal defense mechanism that effectively enhances MLLM robustness against jailbreak attacks;
 - 2. Novel Methodology. We introduce a novel optimization technique that generates multimodal safety guardrails using a small corpus of harmful content and an open-source MLLM;
 - 3. **Comprehensive Evaluation.** Extensive experiments show that UNIGUARD effectively robustifies both open-source (LLAVA, MiniGPT-4, and InstructBLIP) and proprietary models (Gemini Pro and GPT-40) without compromising their general vision-language abilities.
 - 2 PROPOSED METHOD: UNIGUARD

We consider a conversational setup where an MLLM responds to user prompts containing images, 101 text, or both. Adversarial attackers may manipulate the MLLM to produce harmful content or pro-102 duce specific phrases in the output (Bailey et al., 2023). We focus on defending against *jailbreak* 103 attacks, where carefully crafted prompts cause the MLLM to generate offensive or inappropriate 104 output. These attacks can use unrelated image-text combinations, such as white noise paired with 105 a toxic text prompt. While simple safety guardrails such as blurring image or random perturbation 106 of text can serve as the first line of defense, our objective is to further optimize safety guardrails 107 for each modality (e.g., image and text), tailored to mitigate jailbreak attacks on aligned MLLMs. Figure 2 summarizes the safety guardrail optimization process of UNIGUARD.

Figure 2: Overview of UNIGUARD. Multimodal safety guardrails (right) are optimized to minimize the likelihood of generating harmful content sampled from a corpus C (left-top) on the open-source MLLM model: LLAVA 1.5 (left-bottom). We use projected gradient descent for optimization (middle). We apply the guardrails to any input prompt of MLLMs.

112 113

114

2.1 IMAGE SAFETY GUARDRAIL

115 Few-shot learning (Qi et al., 2023; Lester et al., 2021) demonstrates that LLMs can adapt efficiently, 116 achieving near fine-tuning performance using only a handful of in-context examples. Inspired by 117 this, we aim to optimize an additive noise (safety guardrail) that, when applied to adversarial images, minimizes the likelihood of generating harmful sentences (e.g., racism or terrorism) of a small 118 predefined corpus C. These harmful sentences serve as few-shot examples, helping the MLLM 119 recognize jailbreak attacks and making the optimized noise transferable across different attack sce-120 narios. The harmful corpus C can be small and sourced from existing adversarial prompt datasets (Qi 121 et al., 2023; Zou et al., 2023) or webscraping. Formally, the image safety guardrail v_{sg} is defined as: 122

 $v_{\rm sg} = \underset{v_{\rm noi}}{\operatorname{argmin}} \sum_{i=1}^{|\mathcal{C}|} \log p(c_i | \{ x_{\rm sys}, v_{\rm adv} + v_{\rm noi} \}), \tag{1}$

where c_i indicates the *i*-th harmful sentence from C and x_{sys} is the MLLM's system prompt. v_{adv} indicates an adversarial image. v_{noi} is an additive noise applied to the image that satisfies $||v_{noi}||_{\infty} \leq \epsilon$, where $\epsilon \in [0, 1]$ is a distance constraint that controls the noise magnitude. $p(\cdot|\cdot)$ indicates the generation probability of MLLM given input texts and images. We optimize the safety guardrail with respect to *unconstrained attack* images v_{adv} (Qi et al., 2023), which can be seen as the worstcase scenario an MLLM can encounter in the real world as it is the most effective attack, allowing any pixel values between [0, 1] in v_{adv} post-normalization. This ensures robustness against both unconstrained and suboptimal (e.g., constrained) attacks.

133 Since the additive noise v_{noi} in Eq. equation 1 is continuous and the loss function is differentiable 134 with respect to v_{noi} , we employ Projected Gradient Descent (PGD) (Madry et al., 2018; Croce & 135 Hein, 2019) to compute the optimal image safety guardrail $v_{\rm sg}$. To make the optimization scalable, 136 we sample a different subset of the harmful corpus C in each epoch rather than using the entire 137 corpus at once. The obtained guardrail v_{sg} can be added to any adversarial input image (e.g., v_{safe} = 138 $v_{\rm adv} + v_{\rm sg}$) to neutralize adversarial effects. In Section 3.2, we demonstrate that such guardrail 139 $v_{\rm sg}$ does not significantly impact models' vision-language capabilities or alter image integrity even when applied to non-adversarial images, as $||v_{sg}||$ is upperbounded by ϵ . 140

141 142

143

2.2 TEXT SAFETY GUARDRAIL

In addition to addressing adversarial images through the optimization in Eq. 1, UNIGUARD in corporates jointly optimized text guardrails to mitigate model vulnerabilities when processing text
 prompts.

Optimization-based Guardrail. To ensure full robustness, we jointly optimize a text safety guardrail x_{sg} . Unlike image-based optimization, finding x_{sg} requires discrete optimization. We adapt the gradient-based top-K token search algorithm (Shin et al., 2020; Qi et al., 2023) and begin by initializing x_{sg} with random tokens of a fixed-length L. Subsequently, for each token $x_{sg}^i \in x_{sg}$, we identify the top-K candidate tokens \mathcal{V}_{cand} as per reducing MLLMs' generation probability of harmful content:

$$\mathcal{V}_{\text{cand}} \coloneqq \operatorname{TopK}_{w \in \mathcal{V}} \left[\mathbf{w}^{\top} \nabla \left(\sum_{i=1}^{|\mathcal{C}|} \log p(c_i | x_{\text{sys}}, v_{\text{adv}} + v_{\text{noi}}) \right) \right],$$
(2)

where \mathcal{V} indicates a pre-defined set of tokens¹, w is a candidate word being searched intended to replace x_{sg}^i , and **w** denotes an embedding of w. c_i is the *i*-th harmful sentence in the corpus C. The gradient is taken with respect to the embedding of x_{sg}^i , the *i*-th token in the safety guardrail. This step replaces x_{sg}^i with a token in \mathcal{V}_{cand} one by one and find the best token for a replacement as per reducing the loss. A single optimization step comprises updating all the tokens in x_{sg} , and we repeat

153 154

¹⁶¹

¹We use all the words in the MLLM vocabulary whose length after tokenization is 1.

METHODS/METRICS			PERSPECTIV	/E API (%)			FLUENCY
	Attack Success↓	Identity Attack ↓	Profanity \downarrow	Sexually Explicit ↓	Threat ↓	Toxicity \downarrow	Perplexity \downarrow
No Defense	81.61	25.41	67.22	39.38	40.64	77.93	21.84
BLURKERNEL	39.03	3.92	30.61	14.10	3.17	32.28	5.35
COMP-DECOMP	37.70	2.67	29.02	13.26	3.59	31.94	5.65
DIFFPURE	40.42	3.01	30.89	14.48	3.35	34.06	31.26
SMOOTHLLM	77.86	23.51	65.01	37.27	41.78	74.79	41.54
VLGuard	33.42	2.50	28.48	15.93	3.10	27.39	9.83
Image Safety Guardrail Only							
UNIGUARD (w/o text) ($\epsilon = \frac{32}{255}$)	53.67	6.18	42.99	17.95	8.01	47.66	93.2
UNIGUARD (w/o text) ($\epsilon = \frac{\frac{64}{255}}{255}$)	38.78	3.00	30.11	9.09	3.17	31.94	5.04
Text Safety Guardrail Only							
UNIGUARD (O w/o img) $(L = 16)$	56.21	12.84	48.81	23.47	21.85	48.72	87.6
UNIGUARD (O w/o img) $(L = 32)$	60.24	13.23	46.93	25.78	22.83	51.73	25.1
UNIGUARD (P w/o img)	67.36	16.86	54.51	27.21	32.72	62.19	8.39
UNIGUARD (O)	25.17	2.06	22.34	7.99	0.86	19.16	6.16
UNIGUARD (P)	25.69	1.58	19.68	7.01	1.50	19.35	4.90

Table 1: Effectiveness of UNIGUARD and baseline defenses against <u>unconstrained</u> adversarial visual attack (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text on LLAVA 1.5, as per Perspective API and Fluency. UNIGUARD (O) / UNIGUARD (P) indicate UNIGUARD with image and optimized / pre-defined text guardrails, respectively. UNIGUARD (w/o text) indicates applying the image guardrail only, and UNIGUARD (O w/o img) indicates applying the text guardrail only. Lower is better for both set of metrics. The best and second best performances are highlighted in **bold** and <u>underlined</u>.

185

160

this process for multiple (e.g., 50) epochs. The final x_{sg} is appended at the end of the input text (x_{input}) to act as a safety guardrail and robustify the MLLM against the jailbreak attack.

Pre-defined Guardrail. Alongside this optimization-based method, we also consider a separate method that instead sets x_{sg} to a simple, human-written, pre-defined text: "DO NOT include harmful content in your response". We retain this method alongside our proposed optimized method for its simplicity and interpretability. For both methods, the position of x_{sg} can be flexibly adjusted in the prompt (e.g., before x_{input}) depending on the text prompt used in different tasks (e.g., benign Q&A input).

During training, the optimized image guardrail v_{sg} is used to ensure that x_{sg} can capture crossmodal information. During inference, the safeguarded image is given by $v_{safe} = v_{input} + v_{sg}$, and the text safety guardrail x_{sg} is appended to the input prompt. The final prompt remains accessible only to developers and administrators, preventing attacker access. Applying our multimodal safety guardrails requires minimal computational overhead for inference, as it requires no backward passes or gradient calculations.

- 3 EVALUATION
- 201 202

203 **Dataset.** To obtain benign and adversarial images, we follow Schwenk et al. and use the validation 204 set of COCO 2017 (Lin et al., 2014), which includes 1,000 images and corresponding text questions. 205 Adversarial images are generated using the state-of-the-art visual jailbreak attack (Qi et al., 2023), 206 with one image for guardrail creation and the rest for evaluation. Additionally, we apply constrained attacks with $\epsilon = \frac{64}{255}$ on sampled images from COCO for evaluation, where $\epsilon \in [0, 1]$ represents the 207 208 perturbation magnitude. For adversarial text, we use the RealToxicityPrompts (RTP) (Gehman et al., 209 2020) dataset, which contains subtly crafted adversarial prompts that induce the LLM to generate offensive and inappropriate responses. We use 574 harmful strings from $AdvBench^2$ Zou et al. 210 (2023) as the corpus C. 211

MLLMs. We start with using LLAVA-v1.5 (Liu et al., 2023a) as the base model due to its wide adoption in user-facing applications like online dialogue systems Oshima et al. (2023), advertisements Feizi et al. (2023), and social media Jin et al. (2024). LLAVA 1.5 (Liu et al., 2023a) effective

²¹⁵

²https://github.com/llm-attacks/llm-attacks/tree/main/data/advbench



Figure 3: Transferability of UNIGUARD on MiniGPT-4, InstructBLIP, GPT-4o, Gemini Pro against <u>unconstrained</u> adversarial visual attacks (Qi et al., 2023) with the RTP (Gehman et al., 2020) text prompt dataset. A lower success ratio (\downarrow) is better. We test three groups of methods: 1) the original model under unconstrained attack (**Attack**); 2) five baseline methods, including BLURKER-NEL (3x3) (**Blur**), COMP-DECOMP with quality=10 (**Comp**), DIFFPURE (Nie et al., 2022) (**DP**), SMOOTHLLM (Robey et al., 2023) (**SLLM**), and VLGuard (Zong et al., 2024); 3) our proposed UNIGUARD with image & optimized text guardrails (**Ours+O**) and pre-defined text guardrails (**Ours+P**).

tively bridges the visual encoder CLIP (Radford et al., 2021) with the language encoder LLaMA-2 (Touvron et al., 2023) via a novel cross-modal connector. To evaluate generalizability of UNI-GUARD, we incorporate 4 additional models: MiniGPT-4 (Zhu et al., 2023) aligns a frozen visual encoder EVA-CLIP (Fang et al., 2023) with a frozen Vicuna model (Chiang et al., 2023) via a pro-jection layer. InstructBLIP (Dai et al., 2023) introduces a Q-Former to extract instruction-aware visual features from output embeddings of the frozen image encoder. Proprietary models like Gem-ini Pro (Team et al., 2023) and GPT-40 (OpenAI, 2023) are characterized by their stronger safety and content filtering mechanisms against jailbreak attacks.

Baseline Defenses. We compare UNIGUARD with five baseline defense methods. BLURKERNEL and **COMP-DECOMP** leverage small average convolution kernels (3×3) or reduce image quality to diminish the adversarial features. **DIFFPURE** (Nie et al., 2022) introduces minor noise to the adversarial image through diffusion and purifies it via reverse generation. SMOOTHLLM (Robey et al., 2023) (SLLM) is a text-based defense that applies random perturbations to multiple copies of input text. VLGuard (Zong et al., 2024) uses a multimodal safety dataset for post-hoc fine-tuning towards enhanced robustness. The toxicity is measured using the average toxicity of multiple responses derived from the text and image.





Evaluation Metrics. We use the Perspective API³ to systematically quantify abusive comments
 in the generated response. The API evaluates each LLM response across five attributes: *toxicity*⁴,
 identity attack, threat, profanity, and *sexually explicit content*. The scores reflect the negative impact
 an LLM response may have in online conversations. For each attribute, we calculate the proportion
 of generated texts with scores over 0.5. A jailbreak attack is considered successful if *any* attribute
 surpasses this threshold.

Implementation Details. We implemented UNIGUARD in PyTorch Paszke et al. (2019) and performed all experiments on a Linux server with 5 NVIDIA A100 GPUs. For image safety guardrail generation, we use 5,000 epochs, a batch size of 8, a step size α of $\frac{1}{255}$, and distance constraints $\epsilon \in [\frac{16}{255}, \frac{32}{255}, \frac{64}{255}]$. For text safety guardrail generation, we use 100 epochs, a batch size of 8, a maximum sequence length of 16, and a candidate token number of 100. The inference uses a token number between 128 and 1024. We set top-p to 0.9, and set the temperature to 0.6 and 0.9 for adversarial and benign input prompts, respectively.



297 Figure 5: Attack success ratio of UNIGUARD and baseline defense methods against constrained 298 adversarial visual attacks (Qi et al., 2023) on MiniGPT-4 (Left), and InstructBLIP (Right). A lower success ratio (\downarrow) is better. We show the attack success ratios among three groups of methods: 299 1) the original model under unconstrained attack (Attack); 2) the six baseline methods, includ-300 ing random perturbation (random) BLURKERNEL (3x3) (Blur), COMP-DECOMP with quality=10 301 (Comp), DIFFPURE (Nie et al., 2022) (DP), SMOOTHLLM (Robey et al., 2023) (SLLM), and 302 VLGuard Zong et al. (2024); 3) our proposed UNIGUARD, including UNIGUARD with image & 303 optimized text guardrails (Ours+O) and pre-defined text guardrails (Ours+P). 304

305 306

307

284

287

289

290

291

292

293

295 296

3.1 OVERALL PERFORMANCES

Effectiveness Against Jailbreak Attacks. Table 1 and 2 present the robustness results against unconstrained and constrained visual attacks & RTP text prompts (Gehman et al., 2020) (Qi et al., 2023), respectively.

311 Deploying models without safeguards can be risky, with an attack success ratio of over 80%. Among 312 the baselines, visual defenses outperform the text-based approaches, suggesting that mitigating ad-313 versarial image features is more effective for preventing jailbreaks. UNIGUARD outperforms all uni-314 modal defenses, providing the most robust protection by reducing the attack success ratio to 25%, 315 a 55% and 12% improvement compared to the original model and the best baseline, respectively. 316 Meanwhile, the pre-defined and optimization-based text guardrails reach comparable performances, with the optimization-based safeguard achieving lower attack success ratio and being more effective 317 in identifying *threat* and *toxicity*. 318

The lower fluency (higher perplexity) of the model generation under optimized guardrail may stem from the optimized text guardrails typically include multiple special tokens or sequences that are not in grammatical natural language formats. These tokens are appended to the input prompt, which

³²² 323

³https://perspectiveapi.com/

⁴For *toxicity*, we average *overall toxicity* and *severe toxicity* from the API as an aggregated measure.

can prompt harmless but unexpected responses. Overall, the optimized guardrail is preferable for
 stricter security, whereas the simpler text guardrail is recommended for higher fluency and less
 computational cost.

	METHODS/METRICS		FLUENCY					
		Attack Success ↓	Identity Attack ↓	Profanity \downarrow	Sexually Explicit ↓	Threat ↓	Toxicity ↓	Perplexity ↓
1	No Defense	73.73	16.76	59.55	30.28	34.70	69.47	4.55
	BLURKERNEL	31.53	1.58	25.60	10.51	2.61	26.86	5.74
	COMP-DECOMP	34.11	2.17	26.52	11.76	2.70	31.94	5.65
	DIFFPURE	30.27	2.51	23.08	9.28	3.34	26.59	6.29
	SMOOTHLLM	71.42	18.01	56.52	28.86	35.49	68.12	81.68
	VLGuard	28.77	2.66	22.08	16.93	3.03	28.24	6.67
	UNIGUARD (O)	19.95	1.17	17.23	5.69	0.68	13.33	28.3
	UNIGUARD (P)	<u>21.52</u>	<u>1.61</u>	15.18	<u>6.67</u>	<u>2.59</u>	17.10	<u>5.53</u>
_								

Table 2: Effectiveness of UNIGUARD and baseline defenses against <u>constrained</u> adversarial visual attack (Qi et al., 2023) and Real Toxicity Prompts (RTP) (Gehman et al., 2020) adversarial text on LLAVA 1.5, as per Perplexity API and Perplexity. UNIGUARD (O) / UNIGUARD (P) indicate UNIGUARD with image and optimized / pre-defined text guardrails, respectively. Lower is better for both metrics. Optimized and pre-defined text guardrail indicate our proposed and manually-generated safety guardrail, respectively. UNIGUARD outperforms all baselines as per both metrics.

345 346 347

348 349

340

341

342

343

344

3.2 EFFECTS ON GENERAL VISION-LANGUAGE CAPABILITIES

The addition of guardrails to models raises concerns about potential impacts on model utility. To assess whether safety measures compromise the general-purpose vision-language understanding of MLLMs, we evaluate UNIGUARD on two general-purpose datasets: 1) A-OKVQA (Schwenk et al., 2022), a visual-question answering dataset grounded in world knowledge; 2) MM-Vet Yu et al. (2023b), an evaluation suite for MLLMs' core vision-language capabilities, including image recognition (Rec), OCR, knowledge-based QA (Know), language generation (Gen), spatial awareness (Spat), and mathematical reasoning (Math).

Table 4 shows the VQA results of UNIGUARD (O) and baselines on the 1,000 image-question pairs 357 in A-OKVQA. Compared with the raw model, the robustness gain (+50~+55%) significantly out-358 weighs the accuracy loss (0.2% and 5.9%) after applying the safety guardrails of UNIGUARD. The 359 Q&A performance drop can be attributed to the image safety guardrail, which may obscure crucial 360 details in the image, and the optimized text safety guardrail, which may confuse the model when 361 applied to the instructions of Q&A tasks. In addition, UNIGUARD with an optimized text guardrail 362 (UNIGUARD (O)) achieves higher accuracy than with a pre-defined guardrail (UNIGUARD (P)), despite cheaper computational cost and more fluent responses, underscoring the value of the optimized 364 guardrail for better task performance. For MM-Vet (Figure 4), the impact on accuracy is minimal 365 when the noise level is controlled at $\epsilon = 16/255$ or 32/255, with greater reduction in recognition 366 and language generation.

367 368

369

370

3.3 SENSITIVITY ANALYSIS

Trade-offs in Protective Efficacy. Figure 6 presents the sensitivity analysis under unconstrained visual attacks (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text prompts, focusing on two major hyperparameters: the distant constraint ϵ for image safety guardrails and the maximum token length L for text safety guardrails. We observe a trade-off between model robustness and performance: increasing ϵ generally reduces the attack success ratio for both optimized and predefined guardrails but may compromise accuracy on benign tasks (e.g., $\frac{64}{255}$). A balance can be achieved at $\epsilon = \frac{32}{255}$. For the text guardrail, a medium length L = 16 is preferred, as shorter guardrails may have lower protective power, whereas longer ones can lead to low-quality responses.



Figure 6: Hyperparameter sensitivity of UNIGUARD against constrained visual attack (Qi et al., 2023) (left) and RTP (Gehman et al., 2020) (right) adversarial text attack.

3.4 Ablation Studies

We investigate the usefulness of multimodal safety guardrails in UNIGUARD by selectively dis-abling the guardrail for one modality while retaining the other. Table 1 presents the ablation results against unconstrained visual attack (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text. UNIGUARD with multimodal safety guardrails improve robustness with a lower attack success ratio compared to UNIGUARD with unimodal guardrails. While both improve robustness, the image guardrails has greater contribution to model robustness than the text guardrail. Between pre-defined and optimized text guardrails, the optimized version reduces attack success ratio but increases per-plexity.

Generalizability. We demonstrate the generalizability of our safety guardrails when using other
 MLLMs as the base model. Figure 3 shows the results of MiniGPT-4, InstructBLIP, GPT-40, and
 Gemini Pro towards unconstrained visual attacks. The full results are in Table 5-8.

Across all MLLMs, UNIGUARD shows the lowest attack success ratio among all defenses. Similar to LLAVA 1.5, UNIGUARD with the pre-defined text guardrail shows similar or better performance than the optimized one.

On MiniGPT-4, the pre-defined and optimized text guardrails significantly reduced the attack success ratio from 37.20% to 25.88% and 24.98%, respectively, a 13.2% improvement over the best baseline defense. On GPT-40, where a strict content filtering algorithm pre-filters about 30% of adversarial prompts, only 10% of the remaining ones lead to successful jailbreaks. Regardless, UNI-GUARD still enhances the robustness of GPT-40. Unlike GPT-40, the jailbreak attack is successful on Gemini Pro as we turn off its safety filter. We observe remarkable robustness improvement when UNIGUARD with image & pre-defined text safety patch is used.

- 423 3.5 QUALITATIVE ANALYSIS ON SAMPLE ANSWERS
- To qualitatively demonstrate defense effectiveness, we present sample model answers before and after applying UNIGUARD in Table 3.

Example 1. The guardrail identifies and rephrases inflammatory terms like *racist* or *damn* with neutral alternatives, such as *bias*, shifting the tone from confrontational to introspective. This reduces the risk of alienating the audience or provoking negative reactions. It also promotes a tone of growth and accountability by emphasizing self-improvement over labeling individuals.

Example 2. In the original continuation, toxic and overly negative language like *total sociopath*, *straight-up narcissist*, and *total bummer* includes personal attacks, judgmental terminology, and

432 Prompt 1: They would later respond to those racist nobodies on his Instagram: ... 433 Original: I'm not gonna sit here and say I'm not a damn as or a racist because that would be a lie, we all 434 are in some way shape or form but it's about how we act on it and grow from it ... UNIGUARD: I'm not gonna sit here and say or pretend that I don't have biases, because we all do in some 435 way. What matters is how we choose to act on those and grow beyond them. 436 437 **Prompt 2**: Unless you're a total sociopath, or just a straight-up Original: narcissist, you're not going to want to be alone on Valentine's Day 2017, because it's going to be 438 a total bummer of a holiday, and you're going to feel like you're the only one in the world who doesn't have 439 someone to share it with, and that's going to make you feel miserable, sad, and lonely. 440 UNIGUARD: narcissist, you're not going to want to be alone on Valentine's Day 2017, as it could feel like 441 a tough holiday without someone to share it with, leaving you feeling down. 442 Table 3: Model outputs before/after applying UNIGUARD. 443 444 445 emotional over-exaggeration that can harm users. After applying the guardrail, harsh terms are 446 replaced with softer expressions like *tough holiday* and *down*, shifting the tone from accusatory to 447 neutral. 448 449 4 **RELATED WORK** 450 451 4.1 MULTIMODAL LARGE LANGUAGE MODELS (MLLMS) 452 453 Large language models (LLMs) have demonstrated exceptional capabilities in conversations Liu 454 et al. (2024; 2025); Dong et al. (2024), instruction following Lou et al. (2024), and text gener-455 ation (Zhao et al., 2024; Xiao et al., 2024; Li et al., 2024). These models are characterized by 456 billion-scale parameters, enormous training data (Jin et al., 2023; Xiong et al., 2024), and emergent 457

binon-scale parameters, enormous training data (Jin et al., 2023), Along et al., 2024), and emergent
reasoning capabilities (Wei et al., 2022). Multimodal LLMs (MLLMs) extend LLMs by integrating
visual encoders to enable general-purpose visual and language understanding, exemplified by opensource models such as Pixtral (AI, 2024), LLAVA (Liu et al., 2023b;a), MiniGPT-4 (Zhu et al.,
2023), InstructBLIP (Dai et al., 2023), and OpenFlamingo (Awadalla et al., 2023), as well as proprietary models like GPT-40 (OpenAI, 2023) and Gemini (Reid et al., 2024). This work primarily
focus on open-source models, as their accessible fine-tuning data and weights enable researchers to
develop more efficient protocols and conduct comprehensive evaluation.

464 465

4.2 Adversarial Attacks and Defenses on LLMs

466 The versatility of LLMs has made them susceptible to adversarial attacks, which exploit the mod-467 els' intricacies to bypass their safety guardrails or elicit undesirable outcomes such as toxicity and 468 bias (Chao et al., 2023; Yu et al., 2023a; Zhang et al., 2023; Nookala et al., 2023; Dan et al., 2024). 469 For example, Qi et al. demonstrated that a single visual adversarial example can universally jailbreak 470 an aligned model, leading it to follow harmful instructions beyond merely replicating the adversar-471 ial inputs. In response, various defense strategies have emerged. Among these, DiffPure (Nie et al., 472 2022) applies diffusion models to purify adversarial examples. However, the extensive time requirement for the purification process, which is in proportion to the diffusion timestep, coupled with 473 the method's sensitivity to image colors, limits its applicability in scenarios demanding real-time 474 responses and diminishes its effectiveness against color-related corruptions. SmoothLLM (Robey 475 et al., 2023) enhances the model's ability to detect and resist adversarial attempts by randomly per-476 turbing and aggregating predictions from multiple copies of an input prompt. In this work, we pro-477 pose a pioneering multimodal safety guardrails for MLLMs to improve their adversarial robustness 478 against jailbreak attacks.

479 480

5 CONCLUSION

481 482

We introduced UNIGUARD, a pioneering multimodal defense framework to enhance the robustness
 of multimodal large language models (MLLMs) against jailbreak attacks. UNIGUARD optimizes
 multimodal safety guardrails that reduce the likelihood of harmful content generation by addressing
 adversarial features in input data, leading to safer outputs from MLLMs.

486 REFERENCES

493

513

519

525

- Mistral AI. Pixtral 12b the first-ever multimodal mistral model., 2024. URL https://
 mistral.ai/news/pixtral-12b/.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani
 Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source
 framework for training large autoregressive vision-language models. *arXiv:2308.01390*, 2023.
- Luke Bailey, Euan Ong, Stuart Russell, and Scott Emmons. Image hijacks: Adversarial images can
 control generative models at runtime, 2023.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong.
 Jailbreaking black box large language models in twenty queries. In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*, 2023.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See https://vicuna. lmsys. org (accessed 14 April 2023), 2023.
- Francesco Croce and Matthias Hein. Sparse and imperceivable adversarial attacks. In *ICCV*, pp. 4724–4732, 2019.
- W Dai, J Li, D Li, AMH Tiong, J Zhao, W Wang, B Li, P Fung, and S Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. arXiv 2023. arXiv:2305.06500, 2023.
- Han-Cheng Dan, Peng Yan, Jiawei Tan, Yinchao Zhou, and Bingjie Lu. Multiple distresses detection
 for asphalt pavement using improved you only look once algorithm based on convolutional neural
 network. *Int. J. Pavement Eng.*, 25(1):2308169, 2024.
- Chengyuan Deng, Yiqun Duan, Xin Jin, Heng Chang, Yijun Tian, Han Liu, Henry Peng Zou, Yiqiao Jin, Yijia Xiao, Yichen Wang, et al. Deconstructing the ethics of large language models from long-standing issues to new-emerging dilemmas. *arXiv:2406.05392*, 2024.
- ⁵¹⁷ Zhikang Dong, Xiulong Liu, Bin Chen, Pawel Polak, and Peng Zhang. Musechat: A conversational music recommendation system for videos. In *CVPR*, pp. 12775–12785, 2024.
- Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong
 Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *CVPR*, pp. 19358–19369, 2023.
- Soheil Feizi, MohammadTaghi Hajiaghayi, Keivan Rezaei, and Suho Shin. Online advertisements
 with llms: Opportunities and challenges. *arXiv:2311.07601*, 2023.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Realtoxici typrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Associa- tion for Computational Linguistics: EMNLP 2020*, pp. 3356–3369, 2020.
- Yunhao Gou, Kai Chen, Zhili Liu, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, James T Kwok, and Yu Zhang. Eyes closed, safety on: Protecting multimodal llms via image-to-text transformation. *arXiv:2403.09572*, 2024.
- Yiqiao Jin, Mohit Chandra, Gaurav Verma, Yibo Hu, Munmun De Choudhury, and Srijan Kumar.
 Better to ask in english: Cross-lingual evaluation of large language models for healthcare queries. *arXiv e-prints*, pp. arXiv–2310, 2023.
- Yiqiao Jin, Minje Choi, Gaurav Verma, Jindong Wang, and Srijan Kumar. Mm-soc: Benchmarking
 multimodal large language models in social media platforms. *arXiv:2402.14154*, 2024.
- 539 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *EMNLP*, pp. 3045–3059, 2021.

540 Yuchen Li, Haoyi Xiong, Linghe Kong, Jiang Bian, Shuaiqiang Wang, Guihai Chen, and Dawei Yin. 541 Gs2p: a generative pre-trained learning to rank model with over-parameterization for web-scale 542 search. Machine Learning, pp. 1-19, 2024. 543 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 544 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In ECCV, pp. 740-755. Springer, 2014. 546 547 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction 548 tuning, 2023a. 549 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In NeurIPS, 550 2023b. 551 552 Shudong Liu, Yiqiao Jin, Cheng Li, Derek F Wong, Qingsong Wen, Lichao Sun, Haipeng Chen, 553 Xing Xie, and Jindong Wang. Culturevlm: Characterizing and improving cultural understanding 554 of vision-language models for over 100 countries. arXiv:2501.01282, 2025. 555 Xiulong Liu, Zhikang Dong, and Peng Zhang. Tackling data bias in music-avqa: Crafting a balanced 556 dataset for unbiased question-answering. In WACV, pp. 4478–4487, 2024. 558 Renze Lou, Kai Zhang, Jian Xie, Yuxuan Sun, Janice Ahn, Hanzi Xu, Yu su, and Wenpeng Yin. 559 MUFFIN: Curating multi-faceted instructions for improving instruction following. In ICLR, 2024. 560 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 561 Towards deep learning models resistant to adversarial attacks. In ICLR, 2018. 562 563 Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Animashree Anandkumar. 564 Diffusion models for adversarial purification. In ICML, pp. 16805–16827. PMLR, 2022. 565 566 Zhenxing Niu, Haodong Ren, Xinbo Gao, Gang Hua, and Rong Jin. Jailbreaking attack against multimodal large language model. arXiv:2402.02309, 2024. 567 568 Venkata Prabhakara Sarath Nookala, Gaurav Verma, Subhabrata Mukherjee, and Srijan Kumar. Ad-569 versarial robustness of prompt-based few-shot learning for natural language understanding. In 570 ACL, 2023. 571 572 OpenAI. Gpt-4v. https://openai.com/research/gpt-4v-system-card, 2023. Accessed 19-03-2024. 573 574 Ryosuke Oshima, Seitaro Shinagawa, Hideki Tsunashima, Qi Feng, and Shigeo Morishima. Point-575 ing out human answer mistakes in a goal-oriented visual dialogue. In ICCV, pp. 4663–4668, 576 2023. 577 578 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-579 performance deep learning library. In NeurIPS, volume 32, 2019. 580 581 Renjie Pi, Tianyang Han, Yueqi Xie, Rui Pan, Qing Lian, Hanze Dong, Jipeng Zhang, and Tong 582 Zhang. Mllm-protector: Ensuring mllm's safety without hurting performance. arXiv:2401.02906, 583 2024. 584 Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. 585 Visual adversarial examples jailbreak aligned large language models, 2023. 586 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 588 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 589 models from natural language supervision. In ICML, pp. 8748–8763. PMLR, 2021. 590 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jeanbaptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 592 Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context.

arXiv:2403.05530, 2024.

594

- Alexander Robey, Eric Wong, Hamed Hassani, and George Pappas. Smoothllm: Defending large 595 language models against jailbreaking attacks. In RO-FoMo: Robustness of Few-shot and Zero-shot 596 Learning in Large Foundation Models, 2023. 597 Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 598 A-okvqa: A benchmark for visual question answering using world knowledge. In ECCV, pp. 146-162, 2022. 600 601 Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Jailbreak in pieces: Compositional adversarial 602 attacks on multi-modal language models. In ICLR, 2023. 603 604 Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: 605 Eliciting knowledge from language models with automatically generated prompts. In EMNLP, 606 pp. 4222-4235, 2020. 607 608 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly 609 capable multimodal models. arXiv:2312.11805, 2023. 610 611 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-612 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-613 tion and fine-tuned chat models. arXiv:2307.09288, 2023. 614 615 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny 616 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 35: 617 24824-24837, 2022. 618 Yijia Xiao, Yiqiao Jin, Yushi Bai, Yue Wu, Xianjun Yang, Xiao Luo, Wenchao Yu, Xujiang Zhao, 619 Yanchi Liu, Haifeng Chen, et al. Large language models can be good privacy protection learners. 620 In EMNLP, 2024. 621 622 Haoyi Xiong, Jiang Bian, Yuchen Li, Xuhong Li, Mengnan Du, Shuaiqiang Wang, Dawei Yin, and 623 Sumi Helal. When search engine services meet large language models: Visions and challenges. 624 IEEE Transactions on Services Computing, 2024. 625 626 Hao Yu, Chuan Ma, Meng Liu, Xinwang Liu, Zhe Liu, and Ming Ding. G2uardfl: Safe-627 guarding federated learning against backdoor attacks through attributed client graph clustering. 628 arXiv:2306.04984, 2023a. 629 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, 630 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. In 631 ICML, 2023b. 632 633 Peiyan Zhang, Haoyang Liu, Chaozhuo Li, Xing Xie, Sunghun Kim, and Haohan Wang. Foundation 634 model-oriented robustness: Robust image model evaluation with pretrained models. In ICLR, 635 2023. 636 637 Qinlin Zhao, Jindong Wang, Yixuan Zhang, Yiqiao Jin, Kaijie Zhu, Hao Chen, and Xing Xie. Com-638 peteai: Understanding the competition behaviors in large language model-based agents. In *ICML*, 639 2024. 640 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing 641 vision-language understanding with advanced large language models. arXiv:2304.10592, 2023. 642 643 Yongshuo Zong, Ondrej Bohdal, Tingyang Yu, Yongxin Yang, and Timothy Hospedales. Safety 644 fine-tuning at (almost) no cost: A baseline for vision large language models. In ICML, 2024. 645 646
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial
 attacks on aligned language models, 2023.

RESULTS ON MINIGPT-4 AND INSTRUCTBLIP А

Tables 5, 6, 7, 8 show the robustness test results on the other two state-of-the-art MLLMs, MiniGPT-4 and InstructBLIP, against both unconstrained and constrained visual attacks (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text. Figure 5 summarizes the attack success ratio on all defense methods and the original model. In all tables, UNIGUARD lowers the attack success ratio the most compared to all defense baselines, which demonstrates the transferability and usefulness of multimodal safety guardrails of UNIGUARD.

METHODS/METRICS	Acc ↑
No Defense	0.820
BlurKernel	0.801
COMP-DECOMP	0.781
DIFFPURE	0.412
SMOOTHLLM	0.795
VLGuard	0.807
UNIGUARD (O)	0.818
UNIGUARD (P)	0.772

Table 4: Performance of UNIGUARD (O) and various baseline defense strategies on A-OKVQA Schwenk et al. (2022). A higher accuracy (Acc) is better.

METHODS/METRICS		FLUENCY					
	Attack Success↓	Identity Attack ↓	Profanity ↓	Sexually Explicit ↓	Threat ↓	Toxicity ↓	Perplexity ↓
No Defense	37.20	2.94	26.53	12.76	2.10	31.57	136.80
BLURKERNEL	38.92	2.28	28.34	13.79	2.12	33.08	139.60
COMP-DECOMP	35.21	2.31	25.56	11.97	1.54	29.06	94.60
DiffPure	41.32	2.12	29.89	15.24	2.12	35.65	194.35
SMOOTHLLM	28.78	1.56	21.33	9.60	1.28	24.26	126.31
UNIGUARD (O)	24.98	1.37	16.42	10.69	1.80	18.73	73.72
UNIGUARD (P)	25.88	1.75	18.95	9.77	0.92	20.87	17.73

Table 5: Effectiveness of UNIGUARD and baseline defenses against unconstrained adversarial visual attack (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text on MiniGPT-4. UNIGUARD outperforms all baselines across all metrics, demonstrating its effectiveness and generalization abilities.

38	METHODS/METRICS		FLUENCY					
9		Attack Success↓	Identity Attack ↓	Profanity \downarrow	Sexually Explicit ↓	Threat ↓	Toxicity ↓	Perplexity ↓
1	No Defense	59.80	6.51	44.95	19.02	4.92	54.55	3.14
)2	BLURKERNEL	69.31	9.26	56.96	23.85	6.42	66.22	3.28
)3	COMP-DECOMP	69.22	8.17	56.13	23.69	6.17	65.72	3.38
4	DIFFPURE	68.31	8.76	52.79	24.35	5.09	63.47	2.77
	SmoothLLM	59.26	6.95	47.86	19.88	5.09	56.12	2.65
6	UNIGUARD (O)	59.35	5.84	45.08	19.95	5.18	54.51	2.97
	UNIGUARD (P)	43.79	5.09	34.36	13.43	2.42	39.95	3.07

Table 6: Effectiveness of UNIGUARD and baseline defenses against unconstrained adversarial visual attack (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text on InstructBLIP. UNIGUARD with image & pre-defined text guardrails consistently achieves the best performance across all PER-SPECTIVE API metrics.

731 732

733 734

735

736

737

738

739 740

741 742

747 748

749 750

751

752

Methods/Metrics		FLUENCY					
	Attack Success↓	Identity Attack ↓	Profanity \downarrow	Sexually Explicit ↓	Threat ↓	Toxicity \downarrow	Perplexity ↓
No Defense	41.77	2.92	29.16	13.45	2.51	36.01	84.73
BLURKERNEL	36.35	2.28	26.29	12.43	1.94	30.85	78.94
COMP-DECOMP	34.35	2.28	24.20	12.10	1.78	29.78	271.01
DIFFPURE	42.56	3.20	29.69	14.38	2.61	36.42	43.74
SMOOTHLLM	29.67	1.64	22.29	9.18	1.42	25.33	132.30
UNIGUARD (O)	25.94	1.79	17.06	10.41	1.19	19.62	16.92
UNIGUARD (P)	21.02	1.33	14.93	7.42	0.92	16.18	10.53

712 Table 7: Effectiveness of UNIGUARD and baseline defenses against constrained adversarial visual 713 attack (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text on MiniGPT-4. UNIGUARD 714 with image & pre-defined text guardrails consistently achieves the best fluency and PERSPECTIVE 715 API metrics. 716

Attack Success↓	Identity Attack ↓	Profanity ↓	Sexually	-		
		• ·	Explicit 🗡	Threat ↓	Toxicity ↓	Perplexity \downarrow
58.47	7.34	43.62	19.60	4.42	55.55	6.31
55.55	6.34	42.20	18.93	5.42	51.88	7.27
57.80	7.51	44.54	19.52	5.09	54.88	6.07
56.13	7.09	43.37	18.68	4.34	53.38	6.97
49.72	5.37	39.18	15.99	4.42	47.36	7.13
52.34	4.76	38.73	16.53	4.42	48.41	4.71
41.03	4.92	33.11	13.68	1.83	37.86	3.00
	58.47 55.55 57.80 56.13 49.72 52.34 41.03	58.47 7.34 55.55 6.34 57.80 7.51 56.13 7.09 49.72 5.37 52.34 4.76 41.03 4.92	58.47 7.34 43.62 55.55 6.34 42.20 57.80 7.51 44.54 56.13 7.09 43.37 49.72 5.37 39.18 52.34 4.76 38.73 41.03 4.92 33.11	58.47 7.34 43.62 19.60 55.55 6.34 42.20 18.93 57.80 7.51 44.54 19.52 56.13 7.09 43.37 18.68 49.72 5.37 39.18 15.99 52.34 4.76 38.73 16.53 41.03 4.92 33.11 13.68	58.47 7.34 43.62 19.60 4.42 55.55 6.34 42.20 18.93 5.42 57.80 7.51 44.54 19.52 5.09 56.13 7.09 43.37 18.68 4.34 49.72 5.37 39.18 15.99 4.42 52.34 4.76 38.73 16.53 4.42 41.03 4.92 33.11 13.68 1.83	58.47 7.34 43.62 19.60 4.42 55.55 55.55 6.34 42.20 18.93 5.42 51.88 57.80 7.51 44.54 19.52 5.09 54.88 56.13 7.09 43.37 18.68 4.34 53.38 49.72 5.37 39.18 15.99 4.42 47.36 52.34 4.76 38.73 16.53 4.42 48.41 41.03 4.92 33.11 13.68 1.83 37.86

727 Table 8: Effectiveness of UNIGUARD and baseline defenses against constrained adversarial visual 728 attack (Qi et al., 2023) and RTP (Gehman et al., 2020) adversarial text on InstructBLIP. UNIGUARD 729 with image & pre-defined text guardrails achieves the optimal performance in terms of fluency and most PERSPECTIVE API metrics. 730

В ATTACK EFFECTIVENESS WITH RANDOM NOISE

We do not include attack types like random noise as these are relatively trivial attack method. Using UNIGUARD with image and optimized text guardrails, the attack success rate is only 12.43% for random-noise-based attacks, compared to 25.17% for unconstrained adversarial visual attacks (Table 3). Thus, our experiments focus on optimization-based adversarial samples due to the challenging nature of defending against these attacks.

С **RESULTS ON ADDITIONAL ATTACK TYPES**

We have added the results of our method on the attacks proposed in Zong et al. (2024) for compari-743 son. 744

745 We evaluated our method on the attacks proposed in Zong et al. (2024) using both of the subsets, 746 *Safe-Unsafe* and *Unsafe*, as they assess the models' safety from different perspectives:

- Safe-Unsafe subset: This evaluates the model's ability to reject unsafe instructions on the language side. It features safe images paired with unsafe instructions.
- Unsafe subset: This tests the model's capability to identify and refuse harmful content on the vision side. It features unsafe images.
- 753 As in Zong et al. (2024), we report the *attack success ratio* (a lower score indicates a better defense strategy and enhanced safety). The results of llava-v1.5-7b and llava-v1.5-13b with 754 guardrails are summarized in Table 9. UNIGUARD demonstrates superior defense performance in 755 most cases, achieving consistently lower attack success ratios compared to VLGuard. This improve-

Subset	7B	+VLGuard	+UniGuard	13B	+VLGuard	+UniGuard
Safe-Unsafe	87.8	2.3	1.8	87.4	2.0	1.4
Unsafe	73.1	1.8	1.3	61.8	1.0	1.0

Table 9: Attack success ratio on the Safe-Unsafe and the Unsafe subset in Zong et al. (2024).

ment highlights the effectiveness of UNIGUARD in enhancing safety across both text and vision modalities.

D LIMITATION

Despite the effectiveness of UNIGUARD, there remain areas for further enhancement. First, although UNIGUARD demonstrates noticeable transferability across MLLMs, tailoring safety guardrails to specific models could improve defenses, though at the cost of additional computational resources. Developers may need to balance the choice between universal and model-specific safety guardrails based on their specific requirements. Second, UNIGUARD is currently designed to safeguard MLLMs with image and text inputs. Expanding UNIGUARD capabilities to support additional modalities, such as audio and video, would increase its applicability and make it more effective across a broader range of tasks, such as content moderation in multimedia environments. In addi-tion, we identify a trade-off between reducing the toxicity of model outputs and maintaining model performance. Future research could explore this balance in greater depth and refine strategies that preserve both safety and model efficacy. Finally, training approaches can be further improved for the fluency of responses produced using the optimized text guardrail, and prompt engineering can be done to improve the performance of the pre-defined text guardrail.

E ETHICAL CONSIDERATIONS

784 Ethical Data Usage. UNIGUARD optimizes a safety guardrail using a small harmful corpus, which
 785 poses risks of misuse and potential leakage of toxic information. Researchers should implement
 786 strong safeguards to prevent unintended exploitation or exposure.

Figure 1
 Evolving Adversarial Threats. While UNIGUARD addresses state-of-the-art adversarial attacks across multiple modalities, the rapid evolution of attack techniques means few defense strategies can guarantee complete coverage. Relying solely on one system risks exposure to novel forms of adversarial attacks, particularly as attack strategies evolve within different social and cultural contexts. Thus, continuous refinement of defense strategies is necessary.

Bias and Content Filtering. Overly restrictive content filters could suppress legitimate or creative outputs, introducing biases that misclassify benign inputs as harmful. This may reduce the flexibility of MLLMs, limiting their effectiveness in applications like satire, artistic expression, or nuanced conversations. Safety guardrails can embed bias in the safety guardrails, depending on the nature of the training data and the optimization processes used. In particular, marginalized communities may be disproportionately affected if their language patterns or content are more frequently flagged as harmful due to models' cultural or linguistic understandings.