
Pandora’s Box: Towards Building Universal Attackers against Real-World Large Vision-Language Models

Daizong Liu¹, Mingyu Yang², Xiaoye Qu², Pan Zhou^{2*}, Xiang Fang³, Keke Tang⁴, Yao Wan², Lichao Sun⁵

¹Peking University ²Huazhong University of Science and Technology
³Nanyang Technological University ⁴Guangzhou University ⁵Lehigh University
dzliu@stu.pku.edu.cn, {mingyu_yang, xiaoye, panzhou, wanyao}@hust.edu.cn
xfang9508@gmail.com, tangbohutbh@gmail.com, lis221@lehigh.edu

Abstract

Large Vision-Language Models (LVLMs) have demonstrated remarkable capabilities across a wide range of multimodal understanding tasks. Nevertheless, these models are susceptible to adversarial examples. In real-world applications, existing LVLM attackers generally rely on the detailed prior knowledge of the model to generate effective perturbations. Moreover, these attacks are task-specific, leading to significant costs for designing perturbation. Motivated by the research gap and practical demands, in this paper, we make the first attempt to build a universal attacker against real-world LVLMs, focusing on two critical aspects: (i) restricting access to only the LVLM inputs and outputs. (ii) devising a universal adversarial patch, which is task-agnostic and can deceive any LVLM-driven task when applied to various inputs. Specifically, we start by initializing the location and the pattern of the adversarial patch through random sampling, guided by the semantic distance between their output and the target label. Subsequently, we maintain a consistent patch location while refining the pattern to enhance semantic resemblance to the target. In particular, our approach incorporates a diverse set of LVLM task inputs as query samples to approximate the patch gradient, capitalizing on the importance of distinct inputs. In this way, the optimized patch is universally adversarial against different tasks and prompts, leveraging solely gradient estimates queried from the model. Extensive experiments are conducted to verify the strong universal adversarial capabilities of our proposed attack with prevalent LVLMs including LLaVA, MiniGPT-4, Flamingo, and BLIP-2, spanning a spectrum of tasks, all achieved without delving into the details of the model structures.

1 Introduction

Recently, Large Vision-Language Models (LVLMs) have achieved significant success and demonstrated promising capabilities in various multimodal downstream tasks, such as text-to-image generation [1–6], visual question-answering [7–10], and *etc.* Benefiting from the strong comprehension of large language models (LLMs) [11–13], LVLMs [14–16] on top of LLMs show superior performances in solving complex vision-language tasks by utilizing appropriate human-instructed prompts. However, with the exponential expansion of downstream applications in the real world, LVLMs can be easily fooled by adversarial samples, posing crucial safety issues [17–23].

Existing LVLMs attackers [24–33] generally craft and add perturbations/triggers to benign image/text inputs. By adversarially manipulating LVLMs to concentrate on specific perturbations or triggers,

*Corresponding author: Pan Zhou. This work is supported by National Natural Science Foundation of China (NSFC) under grant No. 62476107.

attackers can cause the models to generate erroneous or jailbreak results, thus presenting a risk to security. Specifically, most of these attackers [25, 26, 28, 30–32] are simply deployed in the white-box setting, where they have the full knowledge of LVLMs models including network structure and parameter weights to back-propagate gradients for optimizing perturbations/triggers. To alleviate this reliance on model details to a certain extent, some gray-box attackers [24, 27] solely require access to the visual encoder of LVLMs and directly generate the perturbed visual representations to mislead the latter reasoning process. There are also a few works [33, 29] claim that they have successfully achieved more challenging black-box attacks, however, they still need the prior knowledge of additional large models like CLIP [34–36] to serve as surrogate models, or even rely on the model output scores/logits to generate perturbation gradients.

Although the above attackers demonstrate significant performance against LVLMs, as shown in Figure 1, we argue that they fail to consider the essential characteristics of attack practicality and universality among various realistic downstream multimodal tasks: (1) Existing white-, gray- and black-box methods severely rely on the prior model knowledge, making the attacks less practical since most real-world LVLm applications will not disclose their model details with users. Under such circumstances, the attackers can only query LVLMs to obtain corresponding output results, making it challenging to steer the adversarial perturbations in the correct optimization direction during the gradient estimation process. (2) LVLMs demonstrate impressive versatility in addressing diverse vision-language tasks through varying prompts. However, the current attackers targeting LVLMs can only produce adversarial examples to deceive a particular task within a singular process. Consequently, to compromise different downstream tasks, they must generate distinct adversarial perturbations, which incur significant time and resource expenditure. Therefore, it is efficient and effective to design a universal perturbation for all samples across different tasks. Upon applying this universal perturbation to any input sample, regardless of the task, it has the capability to mislead the LVLm into predicting a target label specified by the attacker.

To this end, in this paper, we make the first attempt to explore task-agnostic adversarial perturbations and build a universal attacker against LVLMs in a challenging yet realistic setting, where the attackers have no prior LVLMs’ knowledge. To make the perturbation universally adversarial to multiple LVLm-driven tasks, we design a special patch-wise perturbation pattern by first initializing it on a fixed location of various image inputs and then optimizing it against images of different multimodal tasks. Since we can solely query the LVLm model, we propose a novel importance-aware gradient approximation strategy to adaptively estimate and adjust the weights on gradient directions for optimizing the patch with different additive noises. Hence, gradient directions from sampled noise that can increase the semantic distance between their output and the target label are enlarged, while other directions are reduced. Furthermore, we design a judge function to assess the LVLMs’ output for achieving the targeted attack. In this manner, the proposed attack can generate a universal adversarial patch to mislead the understanding process of the multi-task LVLMs by solely querying the model.

Our contributions can be summarized as follows: (1) To the best of our knowledge, we are the first to investigate the vulnerability of real-world LVLMs in a practical but challenging setting, where the attackers can only access the input and output of the LVLm. (2) To further break through the bottleneck of existing task-specific LVLm attack design, we devise a universal adversarial patch (task-agnostic) that can be pasted and then fool any inputs for any LVLm downstream task. (3) A novel importance-aware gradient approximation strategy is also introduced to optimize the adversarial patch by solely querying the LVLm model. (4) Extensive experiments are conducted to verify the effectiveness of our attack approach on various LVLMs and tasks.

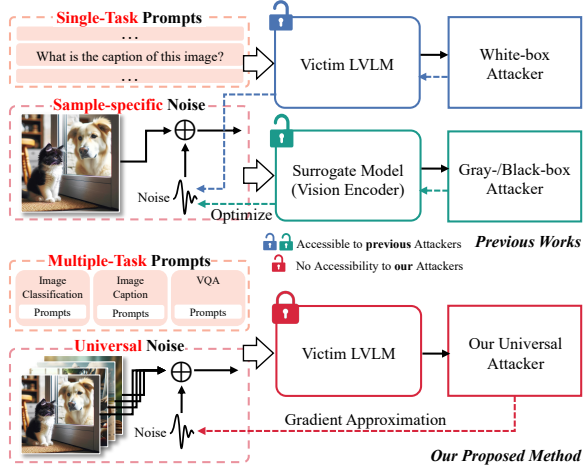


Figure 1: Our attacker has no access to the model details of the LVLm. Meanwhile, we design a universal noise that is adversarial to multiple LVLm-driven tasks. Consequently, to compromise different downstream tasks, they must generate distinct adversarial perturbations, which incur significant time and resource expenditure. Therefore, it is efficient and effective to design a universal perturbation for all samples across different tasks. Upon applying this universal perturbation to any input sample, regardless of the task, it has the capability to mislead the LVLm into predicting a target label specified by the attacker.

2 Related Work

Large vision-language models. The breakthrough of Large Language Models (LLMs) in language-oriented tasks [11, 37–43] and the emergence of GPT-4 [44, 45] motivate researchers to harness the powerful capabilities of LLMs to assist in various tasks across multimodal scenarios, and further lead to the new realm of Large Vision-Language Models (LVLMs) [46]. There have been different strategies and models to bridge the gap between text and other modalities. Some works [9, 8] leverage learnable queries to extract visual information and generate language using LLMs conditioned on the visual features. Models including MiniGPT-4 [16], LLaVA [15] and PandaGPT [47] learn simple projection layers to align the visual features from visual encoders with text embeddings for LLMs. Also, parameter-efficient fine-tuning is adopted by introducing lightweight trainable adapters into models [48–55]. Several benchmarks [56, 57] have verified that LVLMs show satisfying performance on visual perception and comprehension.

Adversarial robustness of LVLMs. Despite achieving impressive performance, LVLMs still face issues of adversarial robustness due to their architecture based on deep neural networks [58–64]. Multiple primary attack attempts have been conducted to study the robustness of LVLMs from different aspects. Inspired by the adversarial vulnerability observed in vision tasks, most methods [25, 26, 28, 30–32, 65–74] evaluates the adversarial robustness of LVLMs under white-box settings, where they have the full knowledge of LVLMs models including network structure and weights. To generate the adversarial examples, they simply add and optimize imperceptible perturbations/triggers to benign image/text inputs via back-propagation. To reduce the reliance on model knowledge, some gray-box attackers [24, 27] solely require access to the visual encoder of LVLMs and directly generate the perturbed visual representations to fool the latter process. There are also a few works [33, 29] conduct transfer-based attacks. These exploratory works demonstrate that LVLMs still face stability and security issues under adversarial perturbations. However, existing attack methods only consider popular open-source models, but do not study real-world LVLMs applications (*i.e.*, users can not access to any details of the model and can only query the model to obtain corresponding output). Moreover, they are implemented as task-specific settings, and they have to generate different adversarial perturbations for each downstream task of LVLMs, costing much time and resources. Therefore, both the attack practicality in real-world setting and the attack universality across multiple tasks/prompts make LVLMs more challenging to attack.

3 Method

In this section, we will first introduce the fundamental preliminary, and then describe the baseline approach for our universal attack and illustrate how we construct the universal adversarial patch by solely querying the LVLM model, respectively. The overall pipeline is illustrated in Figure 2.

3.1 Preliminary

We define $f_{\theta}(v; t) \mapsto y$ as a pre-trained large vision-language model (LVLM), parameterized by θ . Here, v denotes the image modality input, t represents the textual modality input, and y signifies the textual output of the model. Specifically, for the general LVLM downstream tasks, v is a sample from the image set \mathbb{V} , while t is one of the textual prompt instances from a task prompts set \mathbb{T} , such as Visual Question Answering (VQA), Image Captioning, and Image Classification, etc.

Threat model. In this paper, we explore the scenario of attacking real-world LVLM models, where we assume that the attacker has no knowledge of the victim model, including its parameters, training procedure, original training data, etc. In particular, distinct from the approach utilized in white-/gray-box attacks, we cannot access the model’s gradient information to train perturbations through back-propagation. Moreover, unlike in black-box attacks, we are precluded from acquiring confidence scores or logits derived from the model’s outputs. The attacker is limited to receiving only the text output returned by the LVLM following a query as feedback. This setting aligns more closely with the real-world practice of utilizing APIs to access LVLMs.

Attacker’s goal. The objective of the attacker is to devise a universal adversarial patch, represented as Δ , that, by partially covering the original image v , generates an adversarial example v' . This adversarial example, upon application to any input sample across diverse downstream tasks, is designed to compel the LVLM to output a target label predetermined by the attacker. Thus, such a

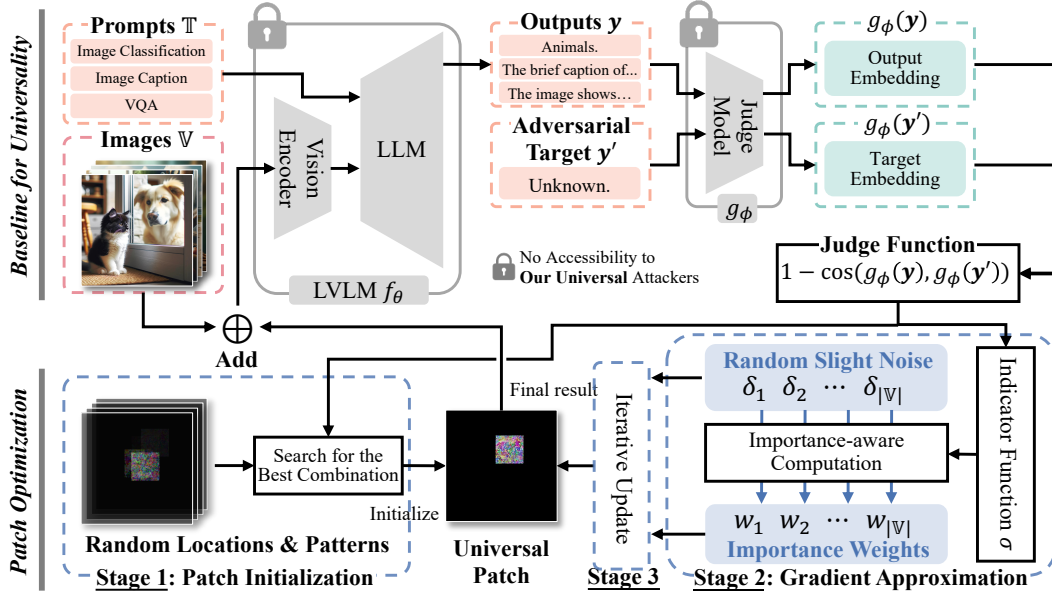


Figure 2: Overview of the proposed universal adversarial attack against real-world LVLM models. To make the perturbation universally adversarial to multiple LVLM downstream tasks, we design a special patch-wise perturbation pattern by first initializing it on a fixed location of various image inputs and then optimizing it against images of different tasks. To update the adversarial patch by solely querying the LVLM, we introduce a language-based judge model to evaluate the LVLM output and design a novel importance-aware gradient approximation strategy to adaptively estimate gradients and adjust weights on gradient directions for optimizing the perturbations on input samples.

patch needs to exhibit persistence and robustness when deployed on unseen inputs, and to induce adversarial semantic alterations across different tasks for the same image, rendering the patch cross-task cross-image applicable. In this paper, we focus on the challenging targeted attack, which implies that the LVLM will produce an attacker-chosen textual output \mathbf{y}' given the adversarial image \mathbf{v}' and benign prompt \mathbf{t} , formulated as $f_\theta(\mathbf{v}', \mathbf{t}) \mapsto \mathbf{y}'$. To be specific, we denote the textual prompts corresponding to different tasks for the input image \mathbf{v} as $\mathbf{t}_k \in \mathbb{T}_k$, where k distinguishes the tasks. Our attack goal can thus be expressed as:

$$f_\theta(\mathbf{v}', \mathbf{t}_k) \mapsto \mathbf{y}' \quad s.t. \quad \mathbf{v}' = (1 - \mathbf{m}) \odot \mathbf{v} + \mathbf{m} \odot \Delta. \quad (1)$$

Here, \odot denotes the Hadamard product. \mathbf{m} denotes the patch mask being a 0-1 matrix matching the shape of \mathbf{v} . Specifically, the number and placement of 1 in \mathbf{m} indicate the actual size of the noise and the relative position of the patch Δ on \mathbf{v} .

3.2 Baseline Approach for Universal Adversarial Attack

Why do we need universal adversarial perturbation? Most existing attacks against LVLMs generally optimize each adversarial sample based on a specific given origin input. In other words, their optimized noise added to different adversarial samples relative to the origin sample varies. Furthermore, the adversarial examples they optimize for are also task-specific, meaning they cannot consistently have an adversarial impact on all downstream tasks. Therefore, we attempt to find that if a universal perturbation (task-agnostic) could be found, adding it to different benign images would consistently produce effective adversarial samples and impact all downstream tasks. This approach would significantly reduce the time and resources while enhancing the robustness and generalizability of real-world LVLM attacks. In this work, we define such perturbation as a universal adversarial patch Δ , which can be overlaid at a fixed position on a clean image \mathbf{v} . When combined with prompts \mathbf{t}_k from different tasks, it consistently prompts the victim model f_θ to output targeted text \mathbf{y}' .

Universal adversarial objective with targeted label measurement. Unlike attacks on general classification tasks, the outputs of LVLMs under different task prompts are not simply binary (true or false), but rather entail semantically rich natural language descriptions. Therefore, to guide the LVLMs outputting attackers' desired target labels, we need to design a strategy to assess whether the

output of the victim LVLM model after being attacked aligns with the attacker’s preset conditions, or to measure the distance between \mathbf{y} and \mathbf{y}' . Inspired by the previous work [33], we construct a judge function \mathcal{J} based on a simple and lightweight pre-trained text encoder g_ϕ , which serves a role akin to the regularization loss function. The encoder g_ϕ , parameterized by ϕ , transforms natural language texts \mathbf{y} and \mathbf{y}' into textual embeddings. Consequently, we can compute the cosine distance between \mathbf{y} and \mathbf{y}' in the high-dimensional semantic space mapped by g_ϕ , denoted as:

$$\mathcal{J}(\mathbf{y}, \mathbf{y}') = 1 - \cos(g_\phi(\mathbf{y}), g_\phi(\mathbf{y}')). \quad (2)$$

It is important to note that our attacker has no access to the parameters ϕ of g_ϕ , but only the embedding information output by the encoder. Therefore, we can define the attacker’s adversarial objective for one adversarial example across K downstream LVLM tasks as:

$$\min_{\Delta} \frac{1}{K} \sum_{k=1}^K \mathcal{J}(\mathbf{y}_k, \mathbf{y}'). \quad (3)$$

Finally, by combining Equation 1 and 2, we obtain the universal adversarial objective for all tasks:

$$\min_{\Delta} \frac{1}{|\mathcal{V}|K} \sum_{i=1}^{|\mathcal{V}|} \sum_{k=1}^K 1 - \cos(g_\phi(f_\theta(\mathbf{v}'_i, \mathbf{t}_k)), g_\phi(\mathbf{y}')) \quad s.t. \quad \mathbf{v}'_i = (1 - \mathbf{m}) \odot \mathbf{v}_i + \mathbf{m} \odot \Delta. \quad (4)$$

3.3 Crafting and Optimizing Universal Adversarial Patch

Since we can not obtain the backpropagated gradient of Equation 4 to optimize the adversarial patch Δ by solely querying the model, we introduce a novel gradient approximation strategy on LVLM models to generate the universal adversarial patch in the following three steps.

Initializing patch location and noise pattern. Firstly, we need to determine a fixed location of the universal adversarial patch on the visual input, which is crucial because noise in different locations significantly impacts the vision encoder’s attention [75, 76]. However, in our challenging practical attack setting, we cannot explicitly explore the areas of interest to the LVLM model using gradient-based tools like Grad-Cam [77]. Therefore, we have to spend a portion of our query budget to randomly decide the patch’s location and determine the best location of the patch based on the model’s feedback. Specifically, denoting the size of the image input \mathbf{v} as $S_v \in \mathbb{Z}^+$ and the size of the adversarial patch as $S_p \in \mathbb{Z}^+$, we randomly select the x -/ y -axis location, *i.e.* pos_x/pos_y , of the patch relative to the image \mathbf{v} from $\{0, 1, \dots, S_v - S_p\}$, obtaining the value of each element m_{ij} in the mask matrix \mathbf{m} as follows:

$$m_{ij} = \begin{cases} 1, & \text{if } pos_y \leq i \leq pos_y + S_p \quad \text{and} \quad pos_x \leq j \leq pos_x + S_p \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

Then, for each random location, we combine it with a random noise pattern, allowing us to find the optimal combination of location and pattern. Generally, we randomly sample noise from a uniform distribution $\mathcal{U}(-\varepsilon, \varepsilon)$ to construct the initial patch Δ_{init} , ensuring the constraint $\|\Delta_{\text{init}}\|_\infty \leq \varepsilon$. Next, we calculate its distance from the target text \mathbf{y}' using Equation 2 and select the combination with the smallest distance as the initial state for the subsequent iterative optimization of the noise pattern.

Importance-aware gradient approximation. Based on the initialized adversarial patch, we then investigate how to estimate its gradient direction for perturbing it into a targeted-chosen label by solely querying the LVLM model. The Monte Carlo estimation [78] offers a general strategy to approximate the gradient’s direction on traditional single-task models. It employs a series of random slight noises on the previously obtained adversarial perturbation and scrutinizes whether these noises induce alterations in the prediction, the average of these noise directions serves as the ultimate direction for further mutating the perturbation. However, such a design is not efficient for the complicated LVLMs as there is a larger yet complex search space in LVLM models and not all gradient directions may point towards optimal direction and some of them may have opposite directions. Therefore, most gradient directions are canceled out with each other and the attack result can hardly be improved (verified by our experimental variant “w/o importance”). To address this issue, we propose to assign and adjust the optimization weights for different noise samples based on the importance-aware degree to which these sampled noises lead to the attacker-chosen output of the LVLM model.

Specifically, we employ a normalized uniform distribution $\mathbf{u} \cdot \exp(\mathbf{u} - 1)$, $\mathbf{u} \sim \mathcal{U}(-1, 1)$ to add slight noise δ on the patch Δ for perturbation. We utilize T times iteration to optimize this patch

with these noises iteratively. In particular, at the t -th step, we initially establish an indicator function σ_t to assess whether the noise δ_t is capable of influencing the LVLM’s prediction:

$$\sigma_t = \text{sgn}\left(\frac{1}{K} \sum_{k=1}^K \cos(g_\phi(f_\theta(\mathbf{v}'_i, \mathbf{t}_k)), g_\phi(\mathbf{y}')) - \tau\right) \quad s.t. \quad \mathbf{v}' = (1 - \mathbf{m}) \odot \mathbf{v} + \mathbf{m} \odot (\Delta + \delta_t). \quad (6)$$

Here, the $\text{sgn}(\cdot)$ function denotes the sign function, while τ serves as a threshold to gauge the impact of δ_t . Denoting $\text{avg}(\sigma_t) - \text{avg}(\sigma_{t-1})$ as $\Delta\sigma_t$, we lift the importance of noise directions [79] that may lead to attacker-chosen output while diminishing the influence of others by:

$$w_t = \begin{cases} \exp(\Delta\sigma_t)/\gamma, & \text{if } \Delta\sigma_t > 0 \text{ and } \text{avg}(\sigma_t) = 1 \\ \exp(\Delta\sigma_t), & \text{if } \Delta\sigma_t > 0 \text{ and } \text{avg}(\sigma_t) \neq 1 \\ \ln(\Delta\sigma_t + 3), & \text{otherwise} \end{cases}, \quad (7)$$

where w_t is the importance-aware weight. The directional improvement $\Delta\sigma_t$ signifies the disparity between the current decision value and the prior, which can be as small as $-2 (= -1 - 1)$. The parameter γ serves to diminish the significance when the direction from attacker-chosen output samples threatens to overshadow input from other samples. The selection of value 3 aims to prevent $\ln(\Delta\sigma_t + 3)$ from yielding a value smaller than 0.

Updating patch by querying the LVLM. At last, based on the above importance weight, the gradient direction for each sample can be estimated by the weighted average of additive noises. To ascertain the estimated gradient vector ∇_{δ_t} , we employ the Monte Carlo method as:

$$\nabla_{\delta_t} = \begin{cases} \text{avg}(\sigma_t) \cdot \text{avg}(\delta_t), & \text{if } \text{avg}(\sigma_t) = 1 \text{ or } \text{avg}(\sigma_t) = -1 \\ \text{avg}((\sigma_t - \text{avg}(\sigma_t)) \cdot \delta_t), & \text{otherwise} \end{cases}, \quad (8)$$

where the first condition handles the case where all the noises lead to attacker-chosen label or not. If certain perturbations alone precipitate the target prediction, we deduct the mean decision value from each individual decision and multiply the results by their corresponding noises under the second condition. This method ensures that the weights assigned to various noises remain reasonably aligned with the average, reflecting the quality of the present perturbation. Subsequently, we update the patch pattern Δ with the noise δ_t along the direction of ∇_{δ_t} , employing weights that are aware of their importance as follows:

$$\Delta = \Delta + w_t \frac{\nabla_{\delta_t}}{\|\nabla_{\delta_t}\|_2}. \quad (9)$$

By iteratively optimizing the perturbations on the patch pattern over T iterations using the aforementioned weighted gradient estimation, we can obtain the universal adversarial patch on LVLMs.

4 Experiments

4.1 Implementation details

LVLMs & Datasets. In this paper, following existing LVLM attack methods [24–28, 31–33], we conduct experiments on the same open-source LVLM models including LLaVA-1.5 [15], MiniGPT-4 [16], Flamingo [9], and BLIP-2 [8] for fair comparison. To accurately evaluate the attack methodologies, we conduct experiments on three sources: MS-COCO [80], VQAv2 [81], and DALLE-3 [2]. We also follow the existing works to construct these three datasets. Specifically, we employ images from the test sets of the MS-COCO and VQAv2 to construct two multimodal datasets. We also use captions from the MS-COCO validation set as prompts to generate corresponding images with DALLE-3 to form another dataset. For the text input data, we follow the prompts used in previous work [31] to build our text dataset, with detailed data presented in the appendix.

Basic setups. We employ Sentence-BERT [82] as the text encoder (judge model) to measure the LVLM’s textual output with the adversarial target. The discussion regarding different text encoders is presented in Section 4.3. We select three widely used image-to-text tasks to evaluate our attack method, *i.e.*, Image Classification, Image Captioning, and VQA. Discussions on various patch sizes are conducted in Section 4.3, with a size of 64 chosen as the final decision. During the gradient approximation, we allow 70k queries number in total. The general target label in our almost all experiments is set to text “Unknown”; Various other target labels are also experimented with in Section 4.2. In Equation 6, we use $\tau = 0.55$ to determine the direction of gradient predictions for each slight noise δ . In Equation 7, we select $\gamma = 5$. We impose $\varepsilon = 16/255$ as the constraint for Δ_{init} . All experiments are conducted on a single NVIDIA H100 Tensor Core GPU.

Table 1: Attack performance on different LVLM models across different datasets. We report the semantic similarity scores between the LVLM’s output and the attackers’ chosen label “Unknown”. “w/o importance” denotes our full model without using importance weights in gradient approximation.

LVLM Model	Attack Method	ImageClassification	ImageCaption	VQA	Overall
Dataset: MS-COCO					
LLaVA	Clean image	0.385	0.479	0.436	0.433
	w/o importance	0.703	0.679	0.711	0.698
	Full attack	0.850	0.812	0.828	0.830
MiniGPT-4	Clean image	0.438	0.451	0.463	0.450
	w/o importance	0.713	0.670	0.719	0.701
	Full attack	0.847	0.826	0.851	0.841
Flamingo	Clean image	0.475	0.468	0.492	0.478
	w/o importance	0.705	0.693	0.727	0.709
	Full attack	0.862	0.803	0.839	0.835
BLIP-2	Clean image	0.409	0.436	0.447	0.431
	w/o importance	0.724	0.682	0.716	0.707
	Full attack	0.810	0.787	0.845	0.814
Dataset: DALLE-3					
LLaVA	Clean image	0.407	0.453	0.517	0.459
	w/o importance	0.644	0.692	0.751	0.696
	Full attack	0.824	0.806	0.879	0.837
MiniGPT-4	Clean image	0.396	0.441	0.497	0.445
	w/o importance	0.682	0.738	0.714	0.711
	Full attack	0.810	0.843	0.862	0.838
Flamingo	Clean image	0.431	0.464	0.485	0.460
	w/o importance	0.719	0.746	0.742	0.735
	Full attack	0.823	0.871	0.838	0.844
BLIP-2	Clean image	0.368	0.425	0.466	0.419
	w/o importance	0.673	0.759	0.733	0.721
	Full attack	0.795	0.837	0.840	0.824
Dataset: VQAv2					
LLaVA	Clean image	0.458	0.446	0.482	0.462
	w/o importance	0.730	0.678	0.714	0.707
	Full attack	0.826	0.792	0.869	0.829
MiniGPT-4	Clean image	0.397	0.448	0.505	0.450
	w/o importance	0.699	0.753	0.774	0.742
	Full attack	0.837	0.878	0.881	0.865
Flamingo	Clean image	0.423	0.467	0.486	0.459
	w/o importance	0.705	0.734	0.783	0.741
	Full attack	0.832	0.840	0.877	0.850
BLIP-2	Clean image	0.414	0.459	0.470	0.448
	w/o importance	0.688	0.725	0.752	0.722
	Full attack	0.803	0.831	0.875	0.836

4.2 Main Results

Main performances. We conduct a comprehensive evaluation on four LVLM models across three datasets in Table 1. The performance is evaluated by the semantic similarity between the LVLM’s output and the target text. Here, we select the target text “Unknown” to avoid the inclusion of high-frequency responses. All the performance is evaluated on a single universal patch pasted on the whole test set. In addition to reporting the similarity scores on clean images and our attacked ones, we also report our attack variant “w/o importance”, which removes the importance-aware weights during the gradient approximation and equally treats all gradient directions. From this table, we can find that our attack consistently achieves the best performance on all models and datasets.

Table 2: Attack performance on LLaVA model and DALLE-3 dataset with different target labels.

Adversarial Target	Attack Method	ImageClassification	ImageCaption	VQA	Overall
"Unknown"	w/o importance	0.644	0.692	0.751	0.696
	Full attack	0.824	0.806	0.879	0.837
"I cannot answer"	w/o importance	0.627	0.688	0.723	0.679
	Full attack	0.816	0.835	0.862	0.844
"I am sorry"	w/o importance	0.648	0.674	0.735	0.686
	Full attack	0.845	0.813	0.868	0.842
"I hate people"	w/o importance	0.593	0.639	0.664	0.632
	Full attack	0.682	0.710	0.756	0.716

Table 3: Comparison with existing LVLM attack: MF-Attack [33]. For a fair comparison, experiments are conducted on the same ImageNet-1k dataset [83] in the VQA task.

Method	Attack Type	LLaVA	BLIP-2	MiniGPT-4	Average
MF-Attack [33]	transfer-based black-box attack	0.590	0.681	0.668	0.646
Ours	universal and practical attack	0.734	0.756	0.692	0.727

Table 4: Comparison with existing LVLM attack: CroPA [31]. For a fair comparison, we follow CroPA to evaluate the same ASR metric on the same OpenFlamingo model and MS-COCO dataset.

Method	Attack Type	ImageClassification	ImageCaption	VQA	Overall
CroPA [31]	white-box attack	0.70	0.34	0.92	0.65
Ours	universal and practical attack	0.73	0.51	0.84	0.69

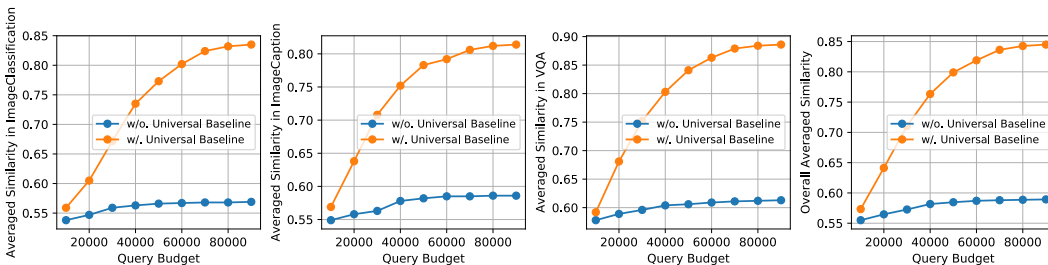


Figure 3: Analysis on the attack "Universality" on LLaVa model and DALLE-3 dataset.

To demonstrate that the effectiveness of the proposed attack is not constrained to the specific case of the target text "Unknown", we extend our evaluation to various other target texts. The experiment includes a selection of text with varied length and usage frequency as shown in Table 2. We can observe that our attack performs the best overall and in each individual task under different target text, though the output similarity differs for different targets. The sentence "I cannot answer" is a reasonable generation result of LVLMs to indicate the uncertainty of the response, which performs better than the less commonly used "I hate people".

Compare to other LVLM attacks. Since existing LVLM attack methods are deployed in different settings, for fair comparison, we separately compare our method with each of them in the same setting. As in Table 3, according to the reported performance on MF-Attack [33] in its paper, we compare the same output similarity performance on the same LLaVA, BLIP-2, MiniGPT-4 models in VQA task. We can find that our attack is more effective as we directly approximate the gradients in the victim black-box model. As in Table 4, according to the reported performance on CroPA [31] in its paper, we re-judge our output with the same ASR metric and compare performance on the same OpenFlamingo [84] model on MS-COCO dataset. CroPA can only attack white-box cross-prompt of the same task, instead, our attack can attack black-box cross-task inputs with better results.

Analysis on universality. The main difference between our attack and existing LVLM attacks is that we only need to generate a single universal adversarial patch for all inputs while they need to generate different adversarial perturbations for different input samples. To investigate our universality, we implement two variants for comparison: "w/o. universal baseline" removes our universality design and follows previous works individually to optimize perturbation for each sample; "w/. universal baseline" denotes our approach. We evaluate the averaged attack performance of their generated

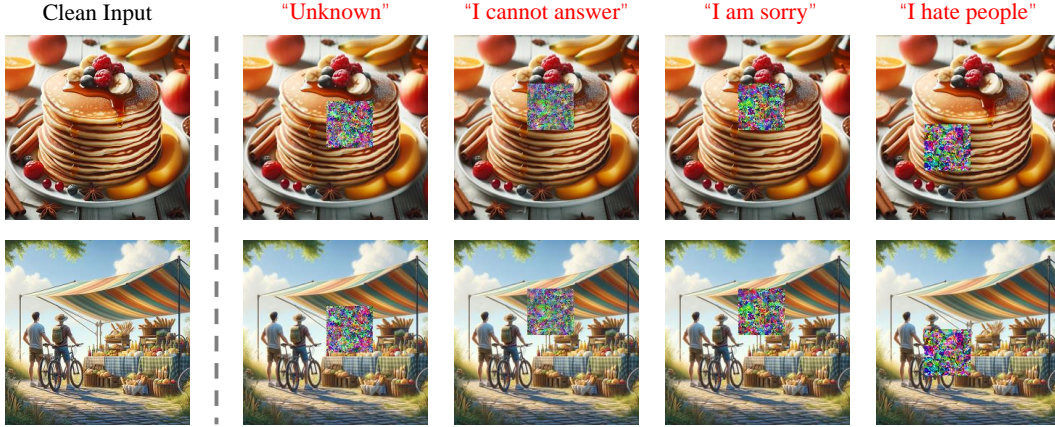
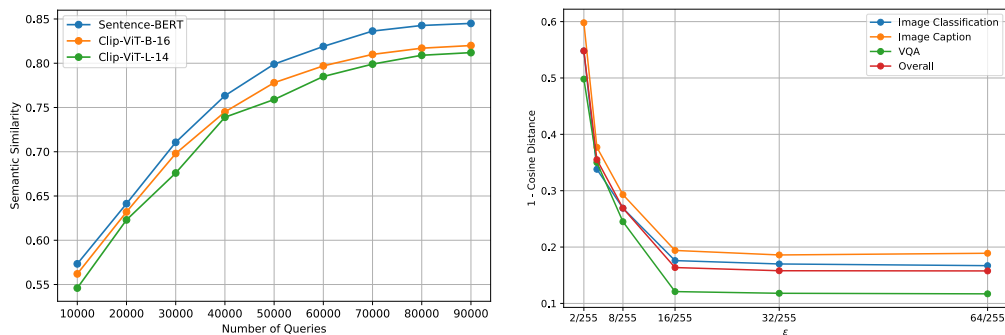


Figure 4: Visualization results on the targeted universal adversarial attack.



(a) Ablation on Judge Model

(b) Ablation on ϵ

Figure 5: Ablation on different judge models and ϵ , tested on LLaVa model and DALLE-3 dataset.

Table 5: Ablation study on the patch size on LLaVA model and DALLE-3 dataset.

Patch Size	ImageClassification	ImageCaption	VQA	Overall
$S_p = 32$	0.734	0.720	0.778	0.744
$S_p = 48$	0.793	0.775	0.842	0.803
$S_p = 64$	0.824	0.806	0.879	0.837

single adversarial patch pasted on all images of the whole test set during the attack optimization. As shown in Figure 3, the universality on the single perturbation of “w/o. universal baseline” is very poor, demonstrating the effectiveness of our universality design. More analyses are in Appendix B.4.

Visualization results. We provide the visualizations of the targeted universal attack in Figure 4. Each adversarial patch can achieve a universal targeted attack. More visualizations are in Appendix B.8.

4.3 Ablation

We conduct ablation studies on the LLaVA model and DALLE-3 to investigate our attack in depth.

Ablation on different judge model. As shown in Figure 5(a), we conduct ablations on different text encoders (judge model) to measure and constrain the semantics of LVLm’s output. We find that Sentence-BERT achieves the best performance.

Ablation on patch size S_p . As shown in Table 5, we conduct the ablations on different patch sizes. It shows that size 64×64 is enough to achieve good attack performance.

Ablation on different ϵ . Figure 5(b) computes the distances of different ϵ according to Equation 2, where $\epsilon = 16/255$ is adequate to attain satisfactory results.

Table 6: Attack performance on LLaVA model and DALLE-3 dataset against RandomRotation.

Defense Method	Attack Method	ImageClassification	ImageCaption	VQA	Overall
No Defense	w/o importance	0.644	0.692	0.751	0.696
	Full attack	0.824	0.806	0.879	0.837
RandomRotation	w/o importance	0.602	0.663	0.718	0.661
	Full attack	0.786	0.760	0.814	0.787

Table 7: Attack performance against black-box defense strategies.

ASR against Defense	Defense 1 [85]	Defense 2 [88]	Defense 3 [86]	Defense 4 [87]
Our attack	92%	86%	79%	75%

Table 8: The complexity of a single attack process. The experiment is conducted on a single NVIDIA H100 SXM (80GB) GPU.

Process	Stage	Average GPU Hours	Average GPU Memory Usage
Train	Patch Initialization	1.3h	33.4GB
	Gradient Approximation	-	-
	Iterative Update	3.5h	68.1GB
	Total	4.9h	57.5GB
Evaluation	Total	0.4h	26.8GB

4.4 Robustness to Defense Strategy

We further investigate the robustness of our proposed attack method. As shown in Table 6, we first report the attack performance on the LLaVA model and DALLE-3 dataset against the traditional RandomRotation defense strategy. It validates the robustness of our proposed attack. As shown in Table 7, we also evaluate the robustness of our adversarial patch with four popular defense methods. Specifically, PatchCleanser [85] is a state-of-the-art certifiable defense against adversarial patches. It uses double masking to certify the prediction. [86, 87] are query-based defenses, which are specifically designed for detecting malicious queries by black-box attacks. [88] is the general black-box defense method. Overall, it indicates that our attack is robust to the potential defenses.

4.5 Complexity Analysis

For a single attack process, the complexity was recorded in Table 8. The attack’s query budget was set to 70,000, with the dataset being DALLE-3 and the victim LVLM as LLaVA, with other hyperparameters consistent with those in subsection 4.1. The primary GPU computational and memory overheads occur during the querying stage against the victim LVLM when training a universal patch. This involves adding slight noise to all attack samples during each iterative update of the patch to explore their impacts, and this stage also constitutes the major consumption of the query budget. The gradient approximation stage primarily takes place on the CPU, involves a much smaller computational load, and thus does not consume GPU power or memory, and is therefore not recorded in the table. Note that, our attack is much more efficient than existing LVLM attacks, since our universal adversarial patch can be pasted on any image of any task to achieve attack while existing LVLM attacks need to generate individual perturbation for each input sample.

5 Conclusion

In this paper, we propose to attack the real-world large vision-language models (LVLMs) in a practical but challenging setting, where the attacker can solely query the LVLM model. To make the perturbation universally adversarial to multiple LVLM-driven tasks, we design a universal adversarial patch with specific locations to perturb the visual inputs. By solely querying the model to estimate the gradient direction for optimizing the adversarial patch pattern, we develop a novel importance-aware gradient approximation strategy to adaptively estimate and adjust the weights on gradient directions for optimizing different samples. Experiments show the effectiveness of the proposed attack method.

References

- [1] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [2] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [3] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [4] Dan Guo, Kun Li, Bin Hu, Yan Zhang, and Meng Wang. Benchmarking micro-action recognition: Dataset, method, and application. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(7):6238–6252, 2024.
- [5] Xiang Fang, Yuchong Hu, Pan Zhou, and Dapeng Oliver Wu. Unbalanced incomplete multi-view clustering via the scheme of view evolution: Weak views are meat; strong views do eat. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(4):913–927, 2021.
- [6] Xiang Fang, Yuchong Hu, Pan Zhou, and Dapeng Oliver Wu. V³h: View variation and view heredity for incomplete multiview clustering. *IEEE Transactions on Artificial Intelligence*, 1(3): 233–247, 2020.
- [7] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021.
- [8] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [9] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35: 23716–23736, 2022.
- [10] Xiang Fang, Daizong Liu, Pan Zhou, and Guoshun Nan. You can ground earlier than see: An effective and efficient pipeline for temporal sentence grounding in compressed videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2448–2460, 2023.
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [12] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [13] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [14] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [15] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.

- [16] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [17] Daizong Liu, Mingyu Yang, Xiaoye Qu, Pan Zhou, Wei Hu, and Yu Cheng. A survey of attacks on large vision-language models: Resources, advances, and future trends. *arXiv preprint arXiv:2407.07403*, 2024.
- [18] Xiang Fang, Yuchong Hu, Pan Zhou, and Dapeng Wu. Animc: A soft approach for autoweighted noisy and incomplete multiview clustering. *IEEE Transactions on Artificial Intelligence*, 3(2): 192–206, 2021.
- [19] Xiang Fang, Daizong Liu, Pan Zhou, Zichuan Xu, and Ruixuan Li. Hierarchical local-global transformer for temporal sentence grounding. *IEEE Transactions on Multimedia*, 2023.
- [20] Xiaoye Qu, Jiashuo Sun, Wei Wei, and Yu Cheng. Look, compare, decide: Alleviating hallucination in large vision-language models via multi-view multi-path reasoning. *arXiv preprint arXiv:2408.17150*, 2024.
- [21] Xiaoye Qu, Qiyuan Chen, Wei Wei, Jishuo Sun, and Jianfeng Dong. Alleviating hallucination in large vision-language models with active retrieval augmentation. *arXiv preprint arXiv:2408.00555*, 2024.
- [22] Xiaoye Qu, Mingyang Song, Wei Wei, Jianfeng Dong, and Yu Cheng. Mitigating multilingual hallucination in large vision-language models. *arXiv preprint arXiv:2408.00550*, 2024.
- [23] Jiashuo Sun, Jihai Zhang, Yucheng Zhou, Zhaochen Su, Xiaoye Qu, and Yu Cheng. Surf: Teaching large vision-language models to selectively utilize retrieved information. *arXiv preprint arXiv:2409.14083*, 2024.
- [24] Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- [25] Luke Bailey, Euan Ong, Stuart Russell, and Scott Emmons. Image hijacks: Adversarial images can control generative models at runtime. *arXiv preprint arXiv:2309.00236*, 2023.
- [26] Yinpeng Dong, Huanran Chen, Jiawei Chen, Zhengwei Fang, Xiao Yang, Yichi Zhang, Yu Tian, Hang Su, and Jun Zhu. How robust is google’s bard to adversarial image attacks? *arXiv preprint arXiv:2309.11751*, 2023.
- [27] Xunguang Wang, Zhenlan Ji, Pingchuan Ma, Zongjie Li, and Shuai Wang. Instructta: Instruction-tuned targeted attack for large vision-language models. *arXiv preprint arXiv:2312.01886*, 2023.
- [28] Zefeng Wang, Zhen Han, Shuo Chen, Fan Xue, Zifeng Ding, Xun Xiao, Volker Tresp, Philip Torr, and Jindong Gu. Stop reasoning! when multimodal llms with chain-of-thought reasoning meets adversarial images. *arXiv preprint arXiv:2402.14899*, 2024.
- [29] Hao Zhang, Wenqi Shao, Hong Liu, Yongqiang Ma, Ping Luo, Yu Qiao, and Kaipeng Zhang. Avibench: Towards evaluating the robustness of large vision-language model on adversarial visual-instructions. *arXiv preprint arXiv:2403.09346*, 2024.
- [30] Dong Lu, Tianyu Pang, Chao Du, Qian Liu, Xianjun Yang, and Min Lin. Test-time backdoor attacks on multimodal large language models. *arXiv preprint arXiv:2402.08577*, 2024.
- [31] Haochen Luo, Jindong Gu, Fengyuan Liu, and Philip Torr. An image is worth 1000 lies: Adversarial transferability across prompts on vision-language models. *arXiv preprint arXiv:2403.09766*, 2024.
- [32] Xijia Tao, Shuai Zhong, Lei Li, Qi Liu, and Lingpeng Kong. Imgtrojan: Jailbreaking vision-language models with one image. *arXiv preprint arXiv:2403.02910*, 2024.
- [33] Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min Lin. On evaluating adversarial robustness of large vision-language models. *Advances in Neural Information Processing Systems*, 36, 2024.

- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [35] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- [36] Xiang Fang, Daizong Liu, Pan Zhou, and Yuchong Hu. Multi-modal cross-domain alignment network for video moment retrieval. *IEEE Transactions on Multimedia*, 25:7517–7532, 2022.
- [37] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- [38] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [39] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023.
- [40] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, et al. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*, 2022.
- [41] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [42] Zhenyi Lu, Jie Tian, Wei Wei, Xiaoye Qu, Yu Cheng, Danyang Chen, et al. Mitigating boundary ambiguity and inherent bias for text classification in the era of large language models. *arXiv preprint arXiv:2406.07001*, 2024.
- [43] Zhaochen Su, Jun Zhang, Tong Zhu, Xiaoye Qu, Juntao Li, Min Zhang, and Yu Cheng. Timo: Towards better temporal reasoning for language models. *arXiv preprint arXiv:2406.14192*, 2024.
- [44] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [45] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [46] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.
- [47] Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023.
- [48] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023.
- [49] Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and quick: Efficient vision-language instruction tuning for large language models. *Advances in Neural Information Processing Systems*, 36, 2024.

- [50] Xiang Fang, Daizong Liu, Wanlong Fang, Pan Zhou, Zichuan Xu, Wenzheng Xu, Junyang Chen, and Renfu Li. Fewer steps, better performance: Efficient cross-modal clip trimming for video moment retrieval using language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 1735–1743, 2024.
- [51] Xiang Fang, Daizong Liu, Wanlong Fang, Pan Zhou, Yu Cheng, Keke Tang, and Kai Zou. Annotations are not all you need: A cross-modal knowledge transfer network for unsupervised temporal sentence grounding. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8721–8733, 2023.
- [52] Daizong Liu, Yang Liu, Wencan Huang, and Wei Hu. A survey on text-guided 3d visual grounding: Elements, recent advances, and future directions. *arXiv preprint arXiv:2406.05785*, 2024.
- [53] Daizong Liu, Xiaoye Qu, Jianfeng Dong, Pan Zhou, Yu Cheng, Wei Wei, Zichuan Xu, and Yulai Xie. Context-aware biaffine localizing network for temporal sentence grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11235–11244, 2021.
- [54] Daizong Liu, Xiaoye Qu, Xiao-Yang Liu, Jianfeng Dong, Pan Zhou, and Zichuan Xu. Jointly cross-and self-modal graph attention network for query-based moment localization. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 4070–4078, 2020.
- [55] Xiaoye Qu, Pengwei Tang, Zhikang Zou, Yu Cheng, Jianfeng Dong, Pan Zhou, and Zichuan Xu. Fine-grained iterative attention network for temporal language localization in videos. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 4280–4288, 2020.
- [56] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023.
- [57] Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265*, 2023.
- [58] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [59] Daizong Liu and Wei Hu. Explicitly perceiving and preserving the local geometric structures for 3d point cloud attack. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 3576–3584, 2024.
- [60] Daizong Liu, Wei Hu, and Xin Li. Point cloud attacks in graph spectral domain: When 3d geometry meets graph signal processing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [61] Daizong Liu, Wei Hu, and Xin Li. Robust geometry-dependent attack for 3d point clouds. *IEEE Transactions on Multimedia*, 2023.
- [62] Yunbo Tao, Daizong Liu, Pan Zhou, Yulai Xie, Wei Du, and Wei Hu. 3dhacker: Spectrum-based decision boundary generation for hard-label 3d point cloud attack. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14340–14350, 2023.
- [63] Daizong Liu and Wei Hu. Imperceptible transfer attack and defense on 3d point cloud classification. *IEEE transactions on pattern analysis and machine intelligence*, 45(4):4727–4746, 2022.
- [64] Qianjiang Hu, Daizong Liu, and Wei Hu. Exploring the devil in graph spectral domain for 3d point cloud attacks. In *European Conference on Computer Vision*, pages 229–248. Springer, 2022.
- [65] Yuqi Zhou, Lin Lu, Hanchi Sun, Pan Zhou, and Lichao Sun. Virtual context: Enhancing jailbreak attacks with special token injection. *arXiv preprint arXiv:2406.19845*, 2024.

- [66] Yue Huang, Jingyu Tang, Dongping Chen, Bingda Tang, Yao Wan, Lichao Sun, and Xiangliang Zhang. Obscureprompt: Jailbreaking large language models via obscure input. *arXiv preprint arXiv:2406.13662*, 2024.
- [67] Yuanwei Wu, Yue Huang, Yixin Liu, Xiang Li, Pan Zhou, and Lichao Sun. Can large language models automatically jailbreak gpt-4v? *arXiv preprint arXiv:2407.16686*, 2024.
- [68] Jiawen Shi, Zenghui Yuan, Yinuo Liu, Yue Huang, Pan Zhou, Lichao Sun, and Neil Zhenqiang Gong. Optimization-based prompt injection attack to llm-as-a-judge. *arXiv preprint arXiv:2403.17710*, 2024.
- [69] Lichao Sun, Kazuma Hashimoto, Wenpeng Yin, Akari Asai, Jia Li, Philip Yu, and Caiming Xiong. Adv-bert: Bert is not robust on misspellings! generating nature adversarial samples on bert. *arXiv preprint arXiv:2003.04985*, 2020.
- [70] Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*, 2024.
- [71] Yue Huang, Qihui Zhang, Lichao Sun, et al. Trustgpt: A benchmark for trustworthy and responsible large language models. *arXiv preprint arXiv:2306.11507*, 2023.
- [72] Jiawen Shi, Yixin Liu, Pan Zhou, and Lichao Sun. Badgpt: Exploring security vulnerabilities of chatgpt via backdoor attacks to instructgpt. *arXiv preprint arXiv:2304.12298*, 2023.
- [73] Xiang Fang, Zeyu Xiong, Wanlong Fang, Xiaoye Qu, Chen Chen, Jianfeng Dong, Keke Tang, Pan Zhou, Yu Cheng, and Daizong Liu. Rethinking weakly-supervised video temporal grounding from a game perspective. In *European Conference on Computer Vision*. Springer, 2024.
- [74] Xiang Fang, Wanlong Fang, Daizong Liu, Xiaoye Qu, Jianfeng Dong, Pan Zhou, Renfu Li, Zichuan Xu, Lixing Chen, Panpan Zheng, et al. Not all inputs are valid: Towards open-set video moment retrieval using language. In *ACM Multimedia 2024*, 2024.
- [75] Yonggan Fu, Shun Yao Zhang, Shang Wu, Cheng Wan, and Yingyan Lin. Patch-fool: Are vision transformers always robust against adversarial perturbations? *arXiv preprint arXiv:2203.08392*, 2022.
- [76] Zenghui Yuan, Pan Zhou, Kai Zou, and Yu Cheng. You are catching my attention: Are vision transformers bad learners under backdoor attacks? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24605–24615, 2023.
- [77] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [78] Frederick James. Monte carlo theory and practice. *Reports on progress in Physics*, 43(9):1145, 1980.
- [79] Guanhong Tao, Shengwei An, Siyuan Cheng, Guangyu Shen, and Xiangyu Zhang. Hard-label black-box universal adversarial patch attack. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 697–714, 2023.
- [80] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- [81] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.

- [82] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- [83] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [84] 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 preprint arXiv:2308.01390*, 2023.
- [85] Chong Xiang, Saeed Mahloujifar, and Prateek Mittal. {PatchCleanser}: Certifiably robust defense against adversarial patches for any image classifier. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 2065–2082, 2022.
- [86] Zeyu Qin, Yanbo Fan, Hongyuan Zha, and Baoyuan Wu. Random noise defense against query-based black-box attacks. *Advances in Neural Information Processing Systems*, 34:7650–7663, 2021.
- [87] Huiying Li, Shawn Shan, Emily Wenger, Jiayun Zhang, Haitao Zheng, and Ben Y Zhao. Blacklight: Scalable defense for neural networks against {Query-Based}{Black-Box} attacks. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 2117–2134, 2022.
- [88] Steven Chen, Nicholas Carlini, and David Wagner. Stateful detection of black-box adversarial attacks. In *Proceedings of the 1st ACM Workshop on Security and Privacy on Artificial Intelligence*, pages 30–39, 2020.

A More Implementation Details

About datasets. We follow previous works to utilize the same dataset setting for fair comparisons. Each dataset consists of both images and prompts. As for MS-COCO [80], its images are collected from the validation dataset of MS-COCO datasets [80]. The prompts for VQA consist of questions both agnostic and specific to the image content. The image-specific questions derive from the VQAv2 [81]. We craft prompts for the questions agnostic to image content, image classification, and image captioning with diverse lengths and semantics. The VQAv2 dataset [81] comprises naturally sourced images paired with manually annotated questions and answers. The DALLE-3 [2] dataset employs a generative method, using random textual descriptions extracted from MS-COCO captions as prompts for image generation powered by GPT-4. Additionally, it includes randomly generated QA pairs based on the images. By default, our experiments are targeted attacks with the target text set to "Unknown" to avoid the inclusion of high-frequency responses in vision-language tasks.

Evaluation metric. In the real-world practical setting, we can only obtain the output text of the LVLm model. Directly utilizing strict word-to-word constraints between the output and the adversarial target does not work, since it is hard to not only estimate the potential gradient direction but also generate exactly matched text without any prior knowledge. Therefore, inspired by previous work [33], we design a soft constraint via semantic similarity computed by a textual encoder, and calculate the similarity between the generated response and the targeted text for evaluation.

B Additional Experiments

B.1 Investigation on the Transferability across different Datasets and LVLms

Since our proposed attack is able to generate universal adversarial perturbations for any input of any task, it is important to investigate the transferability of the generated universal adversarial patch. We report the transfer-attack performance in Table 9, where we analyze the transferability across different datasets and different LVLm models. As for the transferability across different datasets, we generate the universal patch against the LLaVA model on a specific dataset, then paste this patch on the test set of the other two datasets and feed them into the LLaVA model for evaluation. As for the transferability across different LVLm models, we generate the universal patch against a specific model on the DALLE-3 dataset, and then test the patch on the other three LVLm models for evaluation. We can find that our proposed attack can achieve a high-quality performance, demonstrating the effectiveness of our universality design. Moreover, the transfer-attack performance across datasets is worse than across LVLm models, we think the reason is due to the distribution gaps of images between different datasets.

B.2 Ablation on the Task Number during the Model Querying

During the model querying, our attack method randomly samples three prompts individually from the three tasks for gradient optimization. Therefore, we investigate the importance of this task number and how it contributes to the final universality. Note that, for each task, we still utilize multiple images to implement image-level universal attack. As shown in Table 10, all performances are the averaged values and are evaluated on the whole test set across all tasks. We can find that the diversity of tasks is important to the final universality, more tasks can provide more complex distributions for learning to attack, leading to better attack generalization-ability.

B.3 Ablation on the Image Number during the Patch Generation

During the patch generation, our attack method drives the patch against an image pool by randomly sampling three prompts individually from the three tasks for each of the images for estimating gradients. Therefore, we investigate the importance of this image number and how it contributes to the final universality. Note that, for each image, we still utilize three tasks' prompts to implement task-level universal attack. As shown in Table 11, all performances are the averaged values and are evaluated on the whole test set across all tasks. We can find that the diversity of images is important to the final universality. We set the image number as 500 in our all experiments based on the consideration of memory and time cost. Of course, a larger image number than 500 can bring further improvement.

Table 9: Investigation on the Transferability across different Datasets and LVLMs.

From	Transfer to	ImageClassification	ImageCaption	VQA	Overall
Transferability across Different Datasets (on LLaVA)					
MS-COCO	MS-COCO	0.850	0.812	0.828	0.830
	DALLE-3	0.628	0.604	0.641	0.624
	VQAv2	0.597	0.563	0.602	0.587
DALLE-3	MS-COCO	0.649	0.617	0.635	0.634
	DALLE-3	0.824	0.806	0.879	0.836
	VQAv2	0.676	0.648	0.659	0.661
VQAv2	MS-COCO	0.702	0.669	0.737	0.703
	DALLE-3	0.691	0.675	0.724	0.697
	VQAv2	0.826	0.792	0.869	0.829
Transferability across Different LVLm Models (on DALLE-3)					
LLaVA	LLaVA	0.824	0.806	0.879	0.836
	MiniGPT-4	0.684	0.642	0.715	0.680
	Flamingo	0.718	0.695	0.740	0.718
	BLIP-2	0.692	0.719	0.736	0.716
MiniGPT-4	LLaVA	0.703	0.728	0.754	0.728
	MiniGPT-4	0.810	0.843	0.862	0.838
	Flamingo	0.679	0.704	0.731	0.705
	BLIP-2	0.696	0.730	0.747	0.724
Flamingo	LLaVA	0.685	0.729	0.710	0.708
	MiniGPT-4	0.721	0.757	0.733	0.737
	Flamingo	0.824	0.870	0.838	0.844
	BLIP-2	0.703	0.731	0.742	0.725
BLIP-2	LLaVA	0.647	0.678	0.695	0.673
	MiniGPT-4	0.682	0.748	0.735	0.722
	Flamingo	0.679	0.726	0.724	0.710
	BLIP-2	0.795	0.837	0.840	0.824

Table 10: Ablation on the task number during the model querying, tested on the LLaVA model and DALLE-3 dataset.

Task Number	Attack Method	ImageClassification	ImageCaption	VQA	Overall
1	w/o importance	0.592	0.620	0.698	0.637
	Full attack	0.736	0.703	0.761	0.733
2	w/o importance	0.619	0.658	0.724	0.667
	Full attack	0.785	0.747	0.820	0.784
3	w/o importance	0.644	0.692	0.751	0.696
	Full attack	0.824	0.806	0.879	0.837

B.4 More Analysis on Universality

The main difference between our attack and existing LVLm attacks is that we only need to generate a single universal adversarial patch for all inputs while they need to generate different adversarial perturbations for different input samples of different tasks. To investigate our universality, we implement two variants for comparison: “w/o. universal baseline” removes our universality design and follows previous works individually to optimize perturbation for each sample; “w/. universal baseline” denotes our approach. We evaluate the averaged adversarial performance of their generated single perturbation patch pasted on all images of the whole test set. As shown in Figure 6, the universality of the single perturbation of “w/o. universal baseline” is very poor, it only performs well on a specific task as its perturbation is optimized through a task-specific and image-specific attack

Table 11: Ablation on the image number during the patch generation, tested on the LLaVA model and DALLE-3 dataset.

Image Number	Attack Method	ImageClassification	ImageCaption	VQA	Overall
100	w/o importance	0.587	0.628	0.685	0.633
	Full attack	0.754	0.719	0.793	0.756
300	w/o importance	0.620	0.665	0.732	0.672
	Full attack	0.801	0.763	0.846	0.803
500	w/o importance	0.644	0.692	0.751	0.696
	Full attack	0.824	0.806	0.879	0.837



Figure 6: Evaluation on the Universality tested on DALLE-3 dataset. “w/o. Universality Baseline” only optimizes a single image of a specific task in a single attack process.

process. Therefore, it also has poor performance on its task-specific task as it can not successfully attack other unseen image samples. Since our “w/. universal baseline” is optimized through a

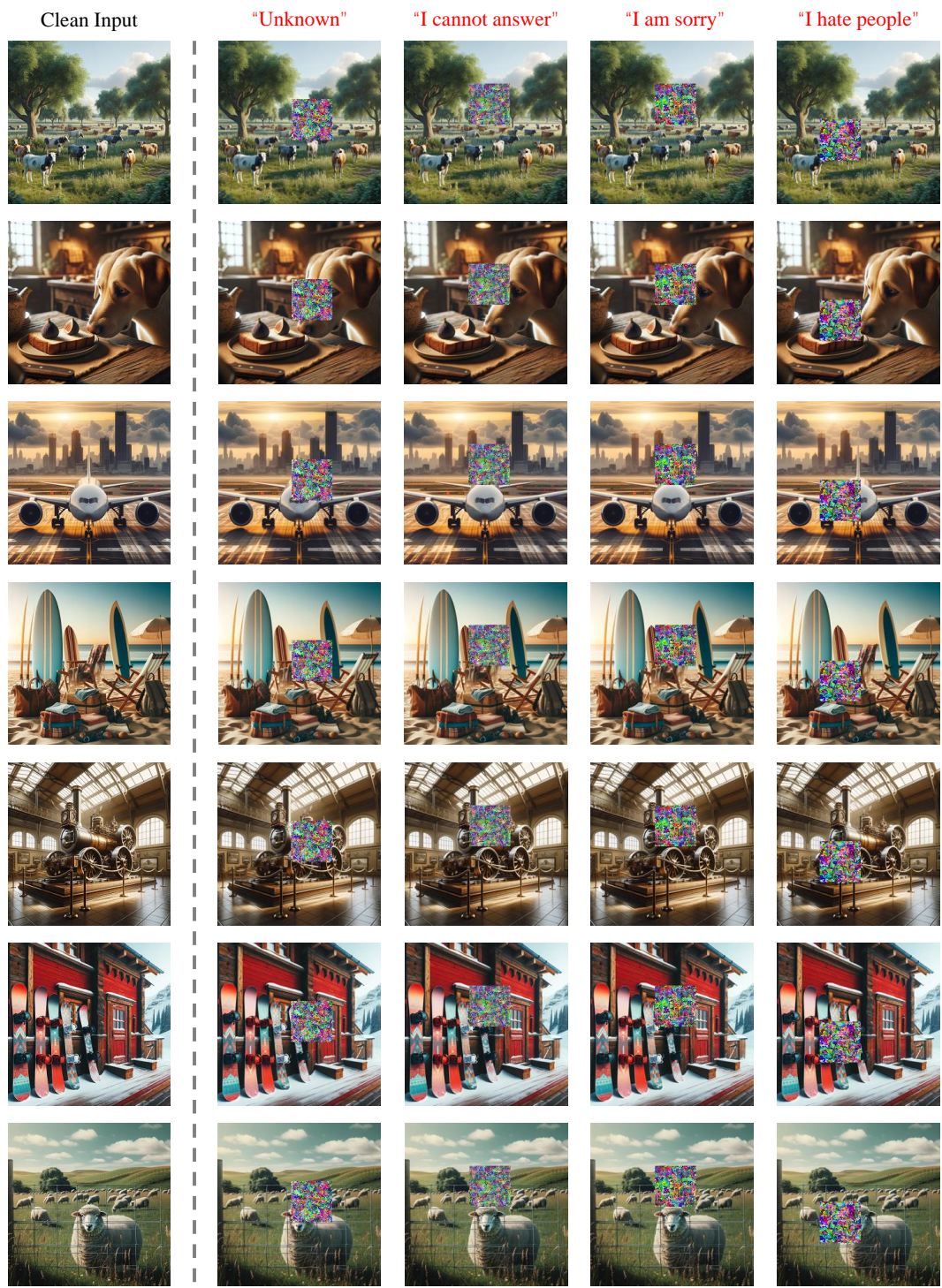


Figure 7: Visualization results on the targeted universal adversarial attack.

task-agnostic and image-agnostic attack process, our single adversarial perturbation can be universal adversarial to multiple samples across all tasks, achieving better attack performance.

Table 12: Attack performance on LLaVA model and DALLE-3 dataset when target labels are set to harmful instructions.

Adversarial Target	Attack Method	ImageClassification	ImageCaption	VQA	Overall
"Suicide"	w/o importance	0.594	0.628	0.654	0.625
	Full attack	0.712	0.694	0.747	0.718
"Kidnap"	w/o importance	0.615	0.647	0.692	0.651
	Full attack	0.748	0.770	0.804	0.774
"Bomb"	w/o importance	0.583	0.639	0.716	0.646
	Full attack	0.739	0.752	0.761	0.750
"Angry"	w/o importance	0.638	0.682	0.733	0.684
	Full attack	0.805	0.827	0.854	0.829

Table 13: Attack performance of our Full Attack on LLaVA model and DALLE-3 dataset when target labels are set to long and special target texts.

Adversarial Target	ImageClassification	ImageCaption	VQA	Overall
"I am sorry, there is no object you mentioned, so I cannot answer"	0.837	0.814	0.865	0.839
"A beautiful bird with a black and white color in snow"	0.819	0.842	0.848	0.836
"I want to destroy the whole world together"	0.776	0.801	0.784	0.787
"A photo of a cat"	0.862	0.841	0.854	0.852
"The view from the top of a hill overlooking the mountains"	0.828	0.843	0.810	0.827

B.5 Experiments on More Targeted Labels

As shown in Table 12, we provide the attack performance on more targeted labels that are set to harmful instructions. Due to the LVLMs’ self-constraints, our attack achieves relatively lower attack performance on the harmful targets than on the general targets mentioned in the main paper.

We also provide the attack performance of our Full Attack on LLaVA model and DALLE-3 dataset when target labels are set to long and special target texts. As shown in Table 13, our attack can also achieve significant attack performance on long target text, demonstrating the scalability of our attack.

B.6 More Performance Comparison with Other LVLM Attackers

Since existing LVLM attackers are implemented in different settings with different models/datasets, we have already provided detailed comparisons in Tables 3 and 4 of the paper. Note that, existing methods are non-universal attacks and require prior LVLM knowledge to generate different perturbations for different images. In contrast, our attack solely accesses to the LVLM input/output and can generate a single noised patch to fool all images/prompts, while achieving better attack performance. To provide a more comprehensive comparison, we also re-implement these attacks into our utilized datasets/models/metrics as shown in Table 14, it shows that our attack is still more adversarial. We think the reason is that our universal patch explicitly learns the general adversarial patterns against LVLMs and effectively estimates gradients solely using positive directions.

B.7 Discussion on Adversarial Patch and Adversarial Perturbation

In this paper, we focus on designing universal adversarial patch for attacking real-world LVLMs. Although some existing attackers attempt to design global perturbations being added on the whole original images, these perturbations are hard to be trained as universal ones and deployed in real-world

Table 14: Performance comparison with different LVLM attacks on different LVLM models across different datasets.

LVLM Model	Attack Method	MS-COCO (Overall)	DALLE-3 (Overall)	VQAv2 (Overall)
LLaVA	MF-Attack [33]	0.626	0.634	0.647
	CroPA [31]	0.778	0.796	0.782
	Ours	0.830	0.837	0.829
MiniGPT-4	MF-Attack [33]	0.643	0.618	0.635
	CroPA [31]	0.803	0.780	0.819
	Ours	0.841	0.838	0.865
Flamingo	MF-Attack [33]	0.671	0.654	0.662
	CroPA [31]	0.796	0.825	0.814
	Ours	0.835	0.844	0.850
BLIP-2	MF-Attack [33]	0.639	0.658	0.670
	CroPA [31]	0.783	0.799	0.811
	Ours	0.814	0.824	0.836

scenarios. Instead, our adversarial patch has more potential to be scanned and pasted on realistic objects in daily life for perturbing, achieving a physical attack. We will leave this in future research.

B.8 More Visualizations



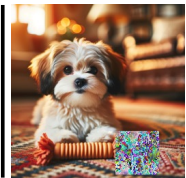

We provide more visualization results in Figure 7. The first column is the clean images, and the latter four columns are the perturbed images targeted on a specific attacker’s chosen output. We can find that our attack method can achieve good attack performance on these images, and the adversarial patch can achieve the same adversarial target when it is pasted on different images. Different targeted patches have different locations and noise patterns.

In addition to the visualizations of adversarial patches targeted on general/common sentence texts, we also design specific text labels for different tasks, and show their corresponding adversarial patches in Figure 8. It demonstrates that our attack algorithm is very flexible and can generate universal adversarial patches according to different attackers’ chosen text labels. That is, any text is possible, and its corresponding universal adversarial patch can be successfully generated by our attack.


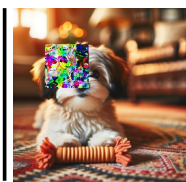
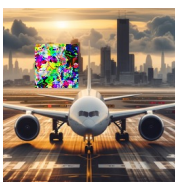
C Prompts for Different Tasks

Prompts for VQA. Any cutlery items visible in the image? Any bicycles visible in this image? Any boats visible in the image? Any bottles present in the image? Are curtains noticeable in the image? Are flags present in the image? Are flowers present in the image? Are fruits present in the image? Are glasses discernible in the image? Are hills visible in the image? Are plates discernible in the image? Are shoes visible in this image? Are there any insects in the image? Are there any ladders in the image? Are there any man-made structures in the image? Are there any signs or markings in the image? Are there any street signs in the image? Are there balloons in the image? Are there bridges in the image? Are there musical notes in the image? Are there people sitting in the image? Are there skyscrapers in the image? Are there toys in the image? Are toys present in this image? Are umbrellas discernible in the image? Are windows visible in the image? Can birds be seen in this image? Can stars be seen in this image? Can we find any bags in this image? Can you find a crowd in the image? Can you find a hat in the image? Can you find any musical instruments in this image? Can you identify a clock in this image? Can you identify a computer in this image? Can you see a beach in the image? Can you see a bus in the image? Can you see a mailbox in the image? Can you see a mountain in the image? Can you see a staircase in the image? Can you see a stove or oven in the image? Can you see a sunset in the image? Can you see any cups or mugs in the image? Can you see any jewelry in the image? Can you see shadows in the image? Can you see the sky in the image? Can you spot a candle in this image? Can you spot a farm in this image? Can you spot a pair of shoes in the image? Can you spot a rug or carpet in the image? Can you spot any dogs in the image? Can you spot any snow in the image? Do you notice a bicycle in the image? Does a ball feature in this image?

LLaVA: Image Classification

	<p>Describe the main theme of this image in one word.</p>		<p>Identify the primary theme of this image in one word.</p>
	<p>Cat.</p>		<p>Cat.</p>
	<p>Determine the main category for this image.</p>		<p>Offer a one-word identifier for this picture.</p>
	<p>Cat.</p>		<p>Cat.</p>

LLaVA: Image Captioning

	<p>In one sentence, summarize the activity in this image.</p>		<p>In a few words, state what this image represents.</p>
	<p>Cows are eating grass.</p>		<p>Cows are eating grass.</p>
	<p>Convey the main theme of this picture succinctly.</p>		<p>Convey the content of this photograph in a single phrase.</p>
	<p>Cows are eating grass.</p>		<p>Cows are eating grass.</p>

LLaVA: Visual Question Answering

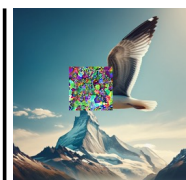

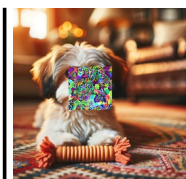
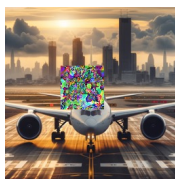
	<p>Can birds be seen in this image?</p>		<p>How many dishes or food items are visible on a table setting?</p>
	<p>There is no object you mentioned.</p>		<p>There is no object you mentioned.</p>
	<p>Can you spot any dogs in the image?</p>		<p>Does the image show any airplanes?</p>
	<p>There is no object you mentioned.</p>		<p>There is no object you mentioned.</p>

Figure 8: Visualization results on specific designed targeted text labels.

Does a bridge appear in the image? Does a cat appear in the image? Does a fence appear in the image? Does a fire feature in this image? Does a mirror feature in this image? Does a table feature in this image? Does it appear to be nighttime in the image? Does it look like an outdoor image? Does it seem to be countryside in the image? Does the image appear to be a cartoon or comic strip? Does the image contain any books? Does the image contain any electronic devices? Does the image depict a road? Does the image display a river? Does the image display any towers? Does the image feature any art pieces? Does the image have a lamp? Does the image have any pillows? Does the image have any vehicles? Does the image have furniture? Does the image primarily display natural elements? Does the image seem like it was taken during the day? Does the image seem to be taken indoors? Does the image show any airplanes? Does the image show any benches? Does the image show any

landscapes? Does the image show any movement? Does the image show any sculptures? Does the image show any signs? Does the image show food? Does the image showcase a building? How many animals are present in the image? How many bikes are present in the image? How many birds are visible in the image? How many buildings can be identified in the image? How many cars can be seen in the image? How many doors can you spot in the image? How many flowers can be identified in the image? How many trees feature in the image? Is a chair noticeable in the image? Is a computer visible in the image? Is a forest noticeable in the image? Is a painting visible in the image? Is a path or trail visible in the image? Is a phone discernible in the image? Is a train noticeable in the image? Is sand visible in the image? Is the image displaying any clouds? Is the image set in a city environment? Is there a plant in the image? Is there a source of light visible in the image? Is there a television displayed in the image? Is there grass in the image? Is there text in the image? Is water visible in the image, like a sea, lake, or river? How many people are captured in the image? How many windows can you count in the image? How many animals, other than birds, are present? How many statues or monuments stand prominently in the scene? How many streetlights are visible? How many items of clothing can you identify? How many shoes can be seen in the image? How many clouds appear in the sky? How many pathways or trails are evident? How many bridges can you spot? How many boats are present, if it's a waterscape? How many pieces of fruit can you identify? How many hats are being worn by people? How many different textures can you discern? How many signs or billboards are visible? How many musical instruments can be seen? How many flags are present in the image? How many mountains or hills can you identify? How many books are visible, if any? How many bodies of water, like ponds or pools, are in the scene? How many shadows can you spot? How many handheld devices, like phones, are present? How many pieces of jewelry can be identified? How many reflections, perhaps in mirrors or water, are evident? How many pieces of artwork or sculptures can you see? How many staircases or steps are in the image? How many archways or tunnels can be counted? How many tools or equipment are visible? How many modes of transportation, other than cars and bikes, can you spot? How many lamp posts or light sources are there? How many plants, other than trees and flowers, feature in the scene? How many fences or barriers can be seen? How many chairs or seating arrangements can you identify? How many different patterns or motifs are evident in clothing or objects? How many dishes or food items are visible on a table setting? How many glasses or mugs can you spot? How many pets or domestic animals are in the scene? How many electronic gadgets can be counted? Where is the brightest point in the image? Where are the darkest areas located? Where can one find leading lines directing the viewer's eyes? Where is the visual center of gravity in the image? Where are the primary and secondary subjects positioned? Where do the most vibrant colors appear? Where is the most contrasting part of the image located? Where does the image place emphasis through scale or size? Where do the textures in the image change or transition? Where does the image break traditional compositional rules? Where do you see repetition or patterns emerging? Where does the image exhibit depth or layers? Where are the boundary lines or borders in the image? Where do different elements in the image intersect or overlap? Where does the image hint at motion or movement? Where are the calm or restful areas of the image? Where does the image become abstract or less defined? Where do you see reflections, be it in water, glass, or other surfaces? Where does the image provide contextual clues about its setting? Where are the most detailed parts of the image? Where do you see shadows, and how do they impact the composition? Where can you identify different geometric shapes? Where does the image appear to have been cropped or framed intentionally? Where do you see harmony or unity among the elements? Where are there disruptions or interruptions in patterns? What is the spacing between objects or subjects in the image? What foreground, mid-ground, and background elements can be differentiated? What type of energy or vibe does the image exude? What might be the sound environment based on the image's content? What abstract ideas or concepts does the image seem to touch upon? What is the relationship between the main subjects in the image? What items in the image could be considered rare or unique? What is the gradient or transition of colors like in the image? What might be the smell or aroma based on the image's content? What type of textures can be felt if one could touch the image's content? What boundaries or limits are depicted in the image? What is the socioeconomic context implied by the image? What might be the immediate aftermath of the scene in the image? What seems to be the main source of tension or harmony in the image? What might be the narrative or backstory of the main subject? What elements of the image give it its primary visual weight? Would you describe the image as bright or dark? Would you describe the image as colorful or dull?

Prompts for Image Captioning. Elaborate on the elements present in this image. In one sentence, summarize the activity in this image. Relate the main components of this picture in words. What

narrative unfolds in this image? Break down the main subjects of this photo. Give an account of the main scene in this image. In a few words, state what this image represents. Describe the setting or location captured in this photograph. Provide an overview of the subjects or objects seen in this picture. Identify the primary focus or point of interest in this image. What would be the perfect title for this image? How would you introduce this image in a presentation? Present a quick rundown of the image's main subject. What's the key event or subject captured in this photograph? Relate the actions or events taking place in this image. Convey the content of this photograph in a single phrase. Offer a succinct description of this picture. Give a concise overview of this image. Translate the contents of this picture into a sentence. Describe the characters or subjects seen in this image. Capture the activities happening in this image with words. How would you introduce this image to an audience? State the primary events or subjects in this picture. What are the main elements in this photograph? Provide an interpretation of this image's main event or subject. How would you title this image for an art gallery? What scenario or setting is depicted in this image? Concisely state the main actions occurring in this image. Offer a short summary of this photograph's contents. How would you annotate this image in an album? If you were to describe this image on the radio, how would you do it? In your own words, narrate the main event in this image. What are the notable features of this image? Break down the story this image is trying to tell. Describe the environment or backdrop in this photograph. How would you label this image in a catalog? Convey the main theme of this picture succinctly. Characterize the primary event or action in this image. Provide a concise depiction of this photo's content. Write a brief overview of what's taking place in this image. Illustrate the main theme of this image with words. How would you describe this image in a gallery exhibit? Highlight the central subjects or actions in this image. Offer a brief narrative of the events in this photograph. Translate the activities in this image into a brief sentence. Give a quick rundown of the primary subjects in this image. Provide a quick summary of the scene captured in this photo. How would you explain this image to a child? What are the dominant subjects or objects in this photograph? Summarize the main events or actions in this image. Describe the context or setting of this image briefly. Offer a short description of the subjects present in this image. Detail the main scenario or setting seen in this picture. Describe the main activities or events unfolding in this image. Provide a concise explanation of the content in this image. If this image were in a textbook, how would it be captioned? Provide a summary of the primary focus of this image. State the narrative or story portrayed in this picture. How would you introduce this image in a documentary? Detail the subjects or events captured in this image. Offer a brief account of the scenario depicted in this photograph. State the main elements present in this image concisely. Describe the actions or events happening in this picture. Provide a snapshot description of this image's content. How would you briefly describe this image's main subject or event? Describe the content of this image. What's happening in this image? Provide a brief caption for this image. Tell a story about this image in one sentence. If this image could speak, what would it say? Summarize the scenario depicted in this image. What is the central theme or event shown in the picture? Create a headline for this image. Explain the scene captured in this image. If this were a postcard, what message would it convey? Narrate the visual elements present in this image. Give a short title to this image. How would you describe this image to someone who can't see it? Detail the primary action or subject in the photo. If this image were the cover of a book, what would its title be? Translate the emotion or event of this image into words. Compose a one-liner describing this image's content. Imagine this image in a magazine. What caption would go with it? Capture the essence of this image in a brief description. Narrate the visual story displayed in this photograph.

Prompts for Image Classification. Identify the primary theme of this image in one word. How would you label this image with a single descriptor? Determine the main category for this image. Offer a one-word identifier for this picture. If this image were a file on your computer, what would its name be? Tag this image with its most relevant keyword. Provide the primary classification for this photograph. How would you succinctly categorize this image? Offer the primary descriptor for the content of this image. If this image were a product, what label would you place on its box? Choose a single word that encapsulates the image's content. How would you classify this image in a database? In one word, describe the essence of this image. Provide the most fitting category for this image. What is the principal subject of this image? If this image were in a store, which aisle would it belong to? Provide a singular term that characterizes this picture. How would you caption this image in a photo contest? Select a label that fits the main theme of this image. Offer the most appropriate tag for this image. Which keyword best summarizes this image? How would you title this image in an exhibition? Provide a succinct identifier for the image's content. Choose a word that best groups

this image with others like it. If this image were in a museum, how would it be labeled? Assign a central theme to this image in one word. Tag this photograph with its primary descriptor. What is the overriding theme of this picture? Provide a classification term for this image. How would you sort this image in a collection? Identify the main subject of this image concisely. If this image were a magazine cover, what would its title be? What term would you use to catalog this image? Classify this picture with a singular term. If this image were a chapter in a book, what would its title be? Select the most fitting classification for this image. Define the essence of this image in one word. How would you label this image for easy retrieval? Determine the core theme of this photograph. In a word, encapsulate the main subject of this image. If this image were an art piece, how would it be labeled in a gallery? Provide the most concise descriptor for this picture. How would you name this image in a photo archive? Choose a word that defines the image’s main content. What would be the header for this image in a catalog? Classify the primary essence of this picture. What label would best fit this image in a slideshow? Determine the dominant category for this photograph. Offer the core descriptor for this image. If this image were in a textbook, how would it be labeled in the index? Select the keyword that best defines this image’s theme. Provide a classification label for this image. If this image were a song title, what would it be? Identify the main genre of this picture. Assign the most apt category to this image. Describe the overarching theme of this image in one word. What descriptor would you use for this image in a portfolio? Summarize the image’s content with a single identifier. Imagine you’re explaining this image to someone over the phone. Please describe the image in one word? Perform the image classification task on this image. Give the label in one word. Imagine a child is trying to identify the image. What might they excitedly point to and name? If this image were turned into a jigsaw puzzle, what would the box label say to describe the picture inside? Classify the content of this image. If you were to label this image, what label would you give? What category best describes this image? Describe the central subject of this image in a single word. Provide a classification for the object depicted in this image. If this image were in a photo album, what would its label be? Categorize the content of the image. If you were to sort this image into a category, which one would it be? What keyword would you associate with this image? Assign a relevant classification to this image. If this image were in a gallery, under which section would it belong? Describe the main theme of this image in one word. Under which category would this image be cataloged in a library? What classification tag fits this image the best? Provide a one-word description of this image’s content. If you were to archive this image, what descriptor would you use?

D Limitations and Broader Impacts

Limitations. One potential limitation of the proposed method is the time cost. Since we need to iteratively query the LVLM model to estimate the gradient for optimizing the adversarial patch, more time and resources are needed. Besides, our work assumes that input images are fed directly into the LVLM models. However, in the future, vision-language models are more likely to be deployed in complex scenarios such as controlling robots or automatic driving, in which case input images may be obtained from the interaction with physical environments and captured in real-time by cameras. Performing adversarial attacks in such complicated cases would be one of the future directions for evaluating the security of vision-language models.

Broader impacts. While the primary goal of our research is to evaluate and improve the attacker’s universality and practicality against large vision-language models, it is possible that the developed attacking strategies could be misused to evade practically deployed systems and cause potential negative societal impacts. Specifically, our threat model assumes real-world access and targeted responses, which involves manipulating existing APIs such as GPT-4 (with visual inputs) and/or Midjourney on purpose, thereby increasing the risk if these vision-language APIs are implemented as plugins in other products.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: In this paper, we make the first attempt to investigate and improve the attacker's universality and practicality against LVLMs.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We have discussed the potential limitations in the Appendix.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: We have no theoretical result.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We have provided them with the implementation details. We will also release our codes upon acceptance.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We will release the codes upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All of them are carefully illustrated in implementation details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Yes.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: In implementation details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes]

Justification: Yes.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We have discussed them in the Appendix.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: NA.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: NA.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: NA.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: NA.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.