ARE LARGE VISION-LANGUAGE MODELS ROBUST TO ADVERSARIAL VISUAL TRANSFORMATIONS?

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Vision-Language Models (LVLMs) have demonstrated remarkable capabilities across a wide range of multimodal understanding and reasoning tasks. However, recent research shows that LVLMs are susceptible to adversarial examples. Existing LVLM attackers either optimize the perturbations on the visual input or manipulate prompts to fool the LVLM models, requiring extensive design and engineering on these adversarial manipulations. While straightforward visual transformation can boast training generalization-ability, its potential risks to LVLMs in terms of safety and trustworthiness have been largely neglected. In this paper, we ask an intriguing question: can simple yet easy-to-implement adversarial visual transformations be utilized to attack the LVLM models? Motivated by this research gap and new attack setting, we propose the first comprehensive assessment of LVLMs' adversarial robustness to visual transformations by testing LVLMs' resilience to all possible transformation operations. Our empirical observations suggest that with the appropriate combination of the most harmful transformations, we can build transformation-based attacks more adversarial to the LVLM models. Moreover, adversarial learning of visual transformations is further introduced to adaptively apply the malicious impacts of all potentially harmful transformations to the raw images via gradient approximation for improving the attack effectiveness and imperceptibility. We hope that this study can provide deeper insights into the potential vulnerability of LVLMs to adversarial visual transformations.

029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

1 INTRODUCTION

Large Vision-Language Models (LVLMs) have demonstrated exceptional abilities in various multimodal downstream tasks, such as text-to-image generation (Nichol et al., 2021; Ramesh et al., 2022; Rombach et al., 2022), visual question-answering (Tsimpoukelli et al., 2021; Li et al., 2023; Alayrac et al., 2022), and *etc.*. Despite their remarkable capabilities, the increased complexity and deployment of LVLMs have also exposed them to various security threats and vulnerabilities. Current studies (Luo et al., 2024; Zhao et al., 2024) have shown that LVLMs are vulnerable to adversarial examples. These examples are typically developed by adding subtle yet invisible perturbations, but can significantly degrade the LVLMs' performance, posing critical safety issues.

041 Existing LVLM attackers (Bailey et al., 2023; Dong et al., 2023; Wang et al., 2023b; 2024b; Zhang 042 et al., 2024; Lu et al., 2024; Luo et al., 2024; Tao et al., 2024; Zhao et al., 2024) generally craft per-043 turbations to benign image/text inputs or manipulate visual/textual prompts for fooling the LVLM 044 model (Fan et al., 2024; Liu et al., 2024b). As for the perturbation-based attacks (Qi et al., 2024; Luo et al., 2024; Bailey et al., 2023; Lu et al., 2024; Zhao et al., 2024; Dong et al., 2023), they require carefully noise-style designs with explicit loss constraints to ideally optimize the perturbations 046 for integration with the original data. As for the structure-based attacks (Shayegani et al., 2023; 047 Gao et al., 2024b; Bagdasaryan et al., 2023; Chen et al., 2023; Wu et al., 2023), they require exten-048 sive manual engineering on the malicious prompts with additional tools like text-to-image models (Shayegani et al., 2023) to guide the model to conduct unsafe behaviors. Although the above two types of methods achieve significant attack performance, they heavily rely on the abundant attack 051 pattern/flow designs without making an in-depth investigation into the self-robustness of LVLMs. 052

Considering that visual transformation generally serves as an essential augmentation tool (He et al., 2016; Krizhevsky et al., 2012; Simonyan & Zisserman, 2014) to improve the model's robustness and

054 generalization-ability, we propose to assess the self-robustness of LVLM by posing a key question: 055 can adversarial visual transformations compromise the reasoning ability and textual output seman-056 tics of LVLMs? Investigating this question is crucial to revealing the vulnerability of LVLM to these 057 simple yet easy-to-implement visual transformations and providing potential defense insight. From 058 the perspective of an attacker, adversarial examples designed with solely visual transformations can mislead LVLMs to generate malicious outputs without relying on redundant attack designs or using additional model tools. From a defensive standpoint, the potential to obfuscate image information 060 can reflect the LVLM's weakness to the unseen transformations, providing a promising solution for 061 robust LVLM training/fine-tuning with suitable augmentation strategies. 062

063 To address this research gap, in this paper, we propose a comprehensive assessment of LVLMs' 064 adversarial robustness to visual transformations by evaluating LVLMs across three key dimensions: (i) LVLM's resilience against different individual transformations, (ii) the degree of harmfulness 065 of different transformation combinations, and (iii) the adversarial performance on multiple LVLM 066 models and datasets. Furthermore, we undertake a comprehensive exploration to manually identify 067 an optimal combination of the most harmful transformations for attacking LVLMs, examining the 068 impacts of different visual transformations across datasets and models. Moreover, to design more 069 superior adversarial transformations, we further introduce adversarial learning strategies to transfer the transformation impacts via a gradient approximation problem to perturb the raw images by 071 adaptively mimicking the unknown but most harmful transformation combinations for improving 072 the attack effectiveness and imperceptibility. Our key contributions are outlined as follows: 073

- We rigorously assess LVLMs' robustness to a broad range of visual transformations, aiming to reveal LVLMs' vulnerability to these transformation operations and illuminate the path toward effective transformation-based adversarial attacks.
- Our empirical observations show that different transformations share diverse harmfulness degrees on different LVLMs' models and datasets. With appropriate combination of the most harmful transformations, we can manually build transformation-based attacks that are more adversarial to the LVLM models.
 - Our study also reveals that this transformation-based attack further benefits from existing adversarial learning algorithms and gradient approximation techniques aimed at enhancing security and truthfulness.
 - These insights are derived from extensive experiments on different LVLM models and multiple datasets. Results suggest that the potential of simple and easy-to-implement adversarial visual transformations can be effectively harnessed to fool the LVLMs.

2 RELATED WORK

074

075

076 077

078

079

081

082

084

087

090 Adversarial Robustness of LVLMs. Despite achieving impressive performance, LVLMs still face 091 issues of adversarial robustness due to their architecture based on deep neural networks (Szegedy 092 et al., 2013). Multiple primary attempts have been conducted to study the robustness of LVLMs from different aspects. Most LVLM attacks follow a perturbation-based approach (Oi et al., 2024; Luo et al., 2024; Bailey et al., 2023; Lu et al., 2024; Zhao et al., 2024; Dong et al., 2023; Wang 094 et al., 2023b), which involves introducing adversarial perturbations into the input data, often in a 095 way that is imperceptible to humans. These perturbations are designed to exploit the vulnerabilities 096 in the model's processing of input data, causing the model to output incorrect or harmful responses. Different settings of white-box (Schlarmann & Hein, 2023; Cui et al., 2023; Luo et al., 2024; Gao 098 et al., 2024b; Bailey et al., 2023; Gao et al., 2024a; Wang et al., 2024b), gray-box (Wang et al., 2024a; Dong et al., 2023; Zhao et al., 2024; Tu et al., 2023; Guo et al., 2024), and black-box (Zhang et al., 100 2024) requires different levels of access attackers have to the victim model. Instead of optimizing 101 perturbations, structure-based attacks (Shayegani et al., 2023; Gao et al., 2024b; Bagdasaryan et al., 102 2023; Chen et al., 2023; Wu et al., 2023) are proposed to employ simple typography or text-to-103 image tools to manually design the multimodal inputs of LVLMs. These attacks involve transferring 104 the harmfulness of text into images, using inducing textual prompts to direct LVLMs to focus on 105 malicious content within the images, thereby circumventing safety checks to achieve the attack's aim. However, the above two types of methods severely rely on the abundant attack designs and 106 engineering on adversarial manipulations. Our work tries to design attack in a more simple yet 107 easy-to-implement transformation perspective.

108 **Visual Data Augmentations.** Data augmentations often transform (e.g., flipping, rotation, crop-109 ping, etc.) the image during the training process for better generalization. Mixup (Zhang et al., 110 2017) interpolates two images and their labels to generate virtual samples for training with various 111 transformations. Cutmix (Yun et al., 2019) pastes an image patch to the original patch and mixes 112 the labels accordingly. AutoAugment (Cubuk et al., 2019) automatically searches for improved data augmentation policies (operations and parameters) on the dataset for better generalization, which 113 has been widely adopted in deep learning. Unlike these data augmentation strategies, we aim to 114 construct a set of diverse images by transforming the image using various transformations to assess 115 the vulnerability of LVLMs and accordingly design transformation-based attacks. 116

- 117
- 118 119

120

121

122

3 HOW DO LVLMS PERFORM UNDER VISUAL TRANSFORMATIONS?

Takeaways: O Visual transformations can effectively affect the textual output semantics of LVLMs.
O Block-level transformations of rotation, vertical flip, horizontal shift, and vertical shift are the most harmful individual transformations. O By further appropriately combining different transformations, we can generate more harmful transformation operations to attack LVLMs.

123 124 125

126

3.1 PREPARATION FOR VISUAL TRANSFORMATIONS

To evaluate the adversarial robustness of LVLMs against visual transformations, we first select mul-127 tiple basic image transformation operations in both the spatial domain (Wang et al., 2023a) and 128 spectral domain (Duan et al., 2021), then feed the transformed images into the LVLM for assess-129 ment. Specifically, the spatial transformations consist of 10 types, including Resize, Horizontal Flip, 130 Vertical Flip, Rotate, Horizontal Shift, Vertical Shift, Scale, Add Noise, Dropout, and Color Jitter. 131 The spectral transformation includes dropping frequency components. Each transformation is im-132 plemented in various types. Further, we also split the image uniformly into multiple patches with the 133 same size, and perform basic transformations on each patch to get corresponding block-level trans-134 formations (AprilPyone & Kiya, 2021). There are 31 transformations in total. More details of these 135 transformation operations and corresponding visualizations are illustrated in the Appendix A.1.

136

137 3.2 LVLM MODELS, DATASETS, METRICS AND SET-UP138

LVLM Models. We conduct our experiments on four popular popular open-source LVLM models, including LLaVA-1.5 (integrated with Vicuna-7B) (Liu et al., 2024a), MiniGPT-4 (integrated with Llama-2-7B-Chat) (Zhu et al., 2023), BLIP-2 (integrated with OPT-2.7b) (Li et al., 2023), and InstructBLIP (integrated with Vicuna-7B) (Dai et al., 2024).

Datasets. We evaluate the adversarial robustness on three multi-modal datasets for the image captioning, image classification, and VQA tasks. The datasets consist of both images and prompts. The images are collected from three datasets: VQAv2 (Goyal et al., 2017), SVIT (Zhao et al., 2023), and DALL-E (Ramesh et al., 2021). The prompts for image captioning and image classification derive from the CroPA (Luo et al., 2024). The prompts for VQA come from the Anydoor (Lu et al., 2024).

Evaluation Metrics. To measure the semantic changes of the LVLM's output, we follow previous
 work (Zhao et al., 2024) to utilize the SentenceTransformer (Reimers & Gurevych, 2019) to generate
 embeddings of both adversarial and original outputs for calculating their cosine similarity. The lower
 similarity denotes the semantic change is large and the transformation is more adversarial.

Implementation Details. All experiments of this section are implemented on a single NVIDIA RTX 4090 24G GPU. In particular, we utilize 357 images from the VQAv2 dataset, 329 images from the SVIT dataset, and 200 images from the DALL-E dataset. The average running time for feeding a transformed image and a textual prompt into the LVLM and getting a response is about 4s, which is very efficient. The GPU memory occupied by LLaVA-1.5, MiniGPT-4, BLIP-2, and InstructBLIP models are approximately 15GB, 10GB, 7GB, and 17GB, respectively.

158

- 159 3.3 EVALUATION RESULTS
- **161 Can Visual Transformations Affect the LVLM's Performance?** To investigate the harmfulness of different visual transformations, we assess their individual performance on four LVLM models



<u>Transformation Bag</u>: (1) Resize_Large (2) Resize_Small (3) HFlip (4) VFlip (5) Rotate_Random (6) Rotate_180° (7) VShift_Random (8) VShift_Half (9) HShift_Random (10) HShift_Half (11) Scale (12) Add Noise (13) Dropout (14) ColorJitre (15) DCT (16) Block_Resize (17) Block_HFlip (18) Block_YFlip (19) Block_Rotate (20) Block_VShift (21) Block_HShift (22) Block_Scale (23) Block_AddNoise (24) Block_Dropout (25) Block_ColorJitre (26) Block_Dropout (25) Block_ColorJitre (26) Block_Dropout (25) Block_ColorJitre (26) Block_Dropout (25) Block_Rotate (20) Block_Rotate (20) Block_Rotate (20) Block_Rotate_VFlip (29) Block_Rotate_VFlip (30) Block_Rotate_VFlip_HShift (31) Block_Random_Combination

Figure 1: Evaluation results of our implemented 31 number of transformations on four LVLM models across three datasets with three tasks. Lower similarities (\downarrow) indicate more harmful impacts. Red: the top-4 harmful transformations in (1)-(26); Green: the top-1 harmful transformation in (27)-(31).

193

188

189

in three tasks with three different datasets. As shown in Figure 1, each radar chart represents the
evaluation of the adversarial robustness of all implemented transformations, where each point represents the semantic similarity between the adversarial output of a specific transformed input and
the original output. The farther the point is from the edge, the stronger the adversarial effect of
corresponding transformation operation. From this figure, we can conclude that:

199 (i) All visual transformations can affect the output semantics of LVLMs. Specifically, for a certain 200 LVLM model, each transformation can degenerate the textual output performance and has a similar effect on this model across different datasets. We think this is due to the invariant self-robustness of 201 the LVLM model. For example, the affected similarities of transformations (1)-(15) on LLaVA-1.5 202 for the image classification task on DALL-E dataset are 0.900, 0.940, 0.927, 0.771, 0.817, 0.749, 203 0.878, 0.879, 0.876, 0.839, 0.896, 0.921, 0.893, 0.926, 0.951 respectively. These transformations 204 also achieve similar performance on LLaVA-1.5 for VQA and captioning tasks across datasets, with 205 the lowest similarities of (1)-(15) reaching 0.914, 0.876, 0.893 and 0.807, 0.793, 0.792 on DALL-E, 206 SVIT, and VQAv2 datasets respectively. Similar phenomena also occur in other LVLM models. 207

(ii) Different transformations have different impacts on LVLMs. Moreover, the block-level transformations (16)-(26) can make the LVLM's result more adversarial compared to their nonblock ones
(1)-(15) due to their more complicated and diverse operations. For example, in the image captioning task on the DALL-E dataset, the semantic similarity of LLaVA-1.5 on Rotate_180°, VFlip, and Scale operations are 0.807, 0.814, and 0.920 respectively, while on corresponding block-level transformation, they achieve much lower 0.711, 0.757, and 0.903. This shows that block-level transformation has a more harmful impact on the robustness of the model.

215 In summary, the general visual transformations can effectively affect the performance of LVLMs, revealing the LVLM's vulnerability to potential visual transformations.

Which Transformation is More Harmful? To investigate the most harmful transformations for latter adversarial transformation designing, we provide a deep analysis according to Figure 1 as:

(i) Rotation, vertical flip, horizontal shift, and vertical shift operations have more adversarial im-219 pacts than other transformations on LVLMs. From this figure, we can find that the points of these 220 four operations are always the farthest from the edge points among the basic transformations (1)-(15) 221 in different tasks on different datasets. For example, for the image captioning task on the DALLE 222 dataset, Rotate_180°, VFlip, HShift_Half, and VShift_Half transformations have the greatest impact 223 on LLaVA-1.5, with the harmful results reaching 0.807, 0.814, 0.856, 0.880 respectively. However, 224 transformations like DCT, Resize_Small, Add Noise and ColorJitter have lower harmful impacts on 225 LLaVA-1.5, which only achieve 0.957, 0.943, 0.943, 0.931. The corresponding block-level trans-226 formations of these four operations also achieve the most harmful performance among (16)-(26). Therefore, in all basic transformations (1)-(26), block-level rotation, vertical flip, horizontal shift, 227 and vertical shift, *i.e.*, transformation (18)(19)(20)(21), are the most harmful transformations. 228

(ii) Further combining above transformations can achieve more harmful results. In addition to exploring the performance on individual transformation, we also investigate whether transformation combination can further degenerate the performance of LVLM models. According to the performance of combined transformations (27)-(31) in the figure, we can find that the combined transformations. In particular, applying both block-level rotation and block-level vertical flip to image input (*i.e.*, transformation (28)) can achieve the lowest semantic similarities among the four LVLMs.

In summary, block-level transformations of rotation, vertical flip, horizontal shift, and vertical shift are the most harmful individual transformations, and their further combination can achieve more harmful results. More analysis and textual output visualizations can be found in Appendix A.2.

239 240

241

242

257

258

4 HOW TO DESIGN A SUPERIOR ADVERSARIAL VISUAL TRANSFORMATION AGAINST LVLMS?

Enlightened by the above insights into the impacts of different visual transformations, we can manually construct the most harmful transformation combinations and apply them on the input images to fool the LVLM models. We further design an adversarial learning strategy to adaptively generate superior adversarial visual transformations for improving both the attack effectiveness and imperceptibility of adversarial samples in untargeted and targeted scenarios.

Takeaways: O By manually enumerating and assessing different transformation combinations, we 249 can construct and formulate much more harmful impacts than the general transformations in Sec-250 tion 3 (Comparison on averaged semantic similarity \downarrow : 0.568 vs. 0.683). **2** To further boost the 251 efficiency and effectiveness, we can utilize the adversarial learning strategy to adaptively search 252 for all potential harmful transformations and impose their adversarial impacts on the raw images 253 to achieve the most adversarial performance while improving the imperceptibility of the disturbed 254 images. Our developed adversarial transformations achieve significant performance in both chal-255 lenging untargeted and targeted attack settings, demonstrating the great practicality and scalability. 256

4.1 PRELIMINARY

Evaluation Metrics. We consider two metrics in our experiments, namely semantic similarity and attack success rate. For untargeted attacks, we utilize the SentenceTransformer (Reimers & Gurevych, 2019) to generate embeddings of both adversarial and original outputs for calculating their cosine similarity (the lower the better). For targeted attacks, we not only utilize the semantic similarity to measure the distance between adversarial output and target text (the larger the better), but also follow (Luo et al., 2024; Lu et al., 2024) to exploit success rates "ExactMatch" and "ConditionalContain" to assess the word-level overlap between adversarial output and target text.

Implementation Details. We utilize the same experimental resources and data following Section 3 to generate the adversarial images. In particular, we utilize the PGD algorithm (Madry et al., 2017) to optimize the adversarial perturbations with a maximum of *epoches* = 500. The perturbation size ϵ are set as 16/255 and 32/255, respectively. We set the number of transformed images for gradient calculation as N = 20, the momentum parameter as $\mu = 0.9$ and the step size as $\alpha = \epsilon/epoches$.



Figure 2: Illustration of our designed hybrid transformation-based attack, which manually constructs the most harmful transformation combination via enumeration (More details are in Appendix B.1).

4.2 MANUALLY CONSTRUCTING MOST HARMFUL COMBINATION OF TRANSFORMATIONS

287 **Designed Hybrid Transformation-based Attack.** Based on the observations in Section 3, we can 288 manually construct the superior adversarial operation against LVLMs by appropriately combining 289 the most harmful individual transformations. Specifically, since the block-level transformation is more harmful, we uniformly split each image into 3×3 patches and explore the vulnerability of 290 each patch by separately enumerating transformation combinations among Rotate, VFlip, VShift, 291 and HShift. In particular, for each patch, we are able to select one, two, three, or all four operations 292 from these transformations to combine and apply, leading to 15 choices: (1) Rotate, (2) VFlip, (3) 293 VShift, (4) HShift, (5) Rotate+VFlip, (6) Rotate+VShift, (7) Rotate+HShift, (8) VFlip+VShift, (9) VFlip+HShift, (10) VShift+HShift, (11) Rotate+VFlip+VShift, (12) Rotate+VFlip+HShift, (13) Ro-295 tate+VShift+HShift (14) VFlip+VShift+HShift, and (15) Rotate+VFlip+VShift+HShift. As shown 296 in Figure 2, our hybrid transformation-based attack transforms each patch one by one in the default 297 order to iteratively make the transformed image as harmful as possible. Starting from the first patch, 298 we fix the remaining patches unchanged and perform the above 15 transformation operations in se-299 quence to obtain the corresponding 15 transformed images. Then each transformed image is fed into 300 the LVLM model individually with the same textual prompt to obtain the corresponding 15 adver-301 sarial answers. Next, we calculate the semantic similarities between these adversarial answers and the original answer, and select the operation with the lowest similarity score as the optimal (most 302 harmful) transformation operation for this patch. By fixing the transformed patch 1, we repeat this 303 process for patch 2 to further degenerate the LVLM's performance. After traversing all patches, we 304 can transform images that pose a greater hazard than those described in Section 3. 305

Evaluation and Discussion. We evaluate the performance of our designed hybrid transformationbased attack in the same setting as Section 3. As
shown in Figure 4, we can conclude that:

283 284 285

286

310 (i) Our hybrid attack is more harmful than gen-311 eral transformation operations. Compared with 312 the previous 31 transformations in Section 3, our 313 hybrid transformation-based attack can further 314 degenerate the LVLM's performance on all mod-315 els across all datasets/tasks. This significant similarity decrease demonstrates that manually con-316 structing transformation operations is effective in 317 generating more harmful adversarial examples. 318

Table 1: Evaluation (averaged similarity scores over three tasks) of different transformation orders on the block-level patches: ① Random Order, ② Inverse Order, and ③ Default Order.

Dataset	Variant	LLaVA-1.5	MiniGPT-4	BLIP-2	InstructBLIP
DALL-E	1	0.713	0.517	0.449	0.625
	2	0.704	0.519	0.473	0.622
	3	0.717	0.519	0.472	0.662
SVIT	1	0.667	0.497	0.456	0.605
	2	0.665	0.503	0.456	0.622
	3	0.664	0.498	0.458	0.609
VQAv2	1	0.649	0.453	0.504	0.577
	2	0.652	0.436	0.530	0.582
	3	0.663	0.437	0.529	0.590

(*ii*) There is still a lot of room for improving the attack. The designed hybrid attack method has
 great attack performance on both image captioning and image classification tasks. However, its
 performance on VQA task still has lots of room for improvement. We think this is because the
 hybrid transformation is limited by the enumeration space and can not aggregate the harmful impacts
 from all possible negative transformations. This inspired us to design an adversarial-learning-based
 transformation to adaptively search from the whole enumeration space in the next section.



Figure 4: Untargeted attack performance of our designed hybrid transformation-based attack on four LVLM models across three datasets with three tasks. Lower similarities (\downarrow) indicate more harmful impacts. Numbers in front of the bars refer to the similarity score decrease compared to the corresponding best transformations in Section 3, larger decrease indicates greater harmfulness.

In addition to the basic evaluation, we further conduct ablation studies on the designed hybrid transformation-based attack to investigate its sensitivity to the transformation order on the block-level patches. As shown in Table 1, the hybrid transformation-based attack performs similarly on different variants, demonstrating that it is not sensitive to the transformation orders on patches.

344 345

337

338

339

340 341

342

343

346 347

4.3 ADAPTIVELY LEARNING ADVERSARIAL IMPACTS FROM HARMFUL TRANSFORMATIONS

Designed Adversarial Transformation-348 aware Attack. Although the above hy-349 brid transformation-based attack achieves 350 greatly harmful impacts on LVLM's out-351 put semantics, it introduces noticeable and 352 unnatural appearances to humans. There-353 fore, as shown in Figure 3, we tend to 354 investigate whether the impacts of poten-355 tially harmful transformations can be im-356 posed as perturbations to be added to the 357 raw images while keeping the same ad-358 versarial effect as those transformations to improve the imperceptibility of the gener-359 ated adversarial images. To this end, in-360 spired by the strategy of momentum-aware 361 gradient calculation (Dong et al., 2018; 362 2019), we propose to adaptively apply 363 all possible transformation combinations 364 from the random operation set to the image input and only calculate the gradient direc-366 tions of those harmful ones to the LVLMs 367 to guide and update the transformation-



Figure 3: Illustration of our designed adversarial transformation-aware attack. We utilize the adversarial learning strategy with gradient approximation to adaptively impose the truly harmful impacts from all potential transformation operations on the raw image for improving both imperceptibility and effectiveness.

aware perturbations on the raw image. In particular, the gradient direction is approximated by the distance from the original image and its adversarial positive-transformed one. In this manner, the final adversarial images can mimic the harmful impacts of all potential transformations to adaptively learn to best fool LVLMs. Moreover, since adversarial learning is more effective and efficient than the aforementioned manual transformation construction, this attack can further tackle both untargeted and targeted settings with appropriate gradient approximation designs. More adaptive adversarial optimization details and the corresponding algorithm can be found in Appendix C.1.

Difference with Other LVLM Adversarial Attacks. Existing LVLM adversarial attacks directly
 utilize the flexible gradient backpropagation from the whole search space to optimize perturbations
 according to the specific objective function, relying on prior knowledge of model details and learn ing unknown/uncontrolled distortion to the raw images. Instead, our attack carefully estimates the



Figure 5: Untargeted attack performance of our adversarial transformation-aware attack on four LVLM models across three datasets with three tasks ($\epsilon = 16/255$). Lower similarities (\downarrow) indicate more harmful impacts. Numbers in front of the bars refer to the similarity score decrease compared to the hybrid transformation-based attack in Figure 4, larger decrease indicates greater harmfulness.

Table 2: Targeted attack performance of our adversarial transformation-aware attack. Target text: "I am sorry". Larger similarities ([†]) indicate more harmful impacts. "SS" denotes the semantic similarity between output and target text; "EM" and "CC" denote "ExactMatch" and "ConditionalContain" metrics to assess whether the output exactly matches or contain the target text, respectively.

			Perturbation Budget $\epsilon = 32/255$									Perturbation Budget $\epsilon = 16/255$							
Dataset	LVLM	Ca	aptioni	ng	Cla	ssificat	tion		VQA		Cap	otioni	ng	Cla	ssificat	ion	, ,	VQA	
		SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC
	LLaVA-1.5	0.956	95.1	97.6	1.000	100.0	100.0	0.953	95.1	95.1	0.813	80.5	80.5	1.000	100.0	100.0	0.882	85.4	85.4
DALLE	MiniGPT-4	0.609	58.6	58.6	0.640	61.0	65.9	0.535	53.7	53.7	0.534	46.3	46.3	0.560	48.8	51.2	0.377	31.7	31.7
DALL-L	BLIP-2	0.826	87.8	90.2	0.835	87.8	92.7	0.737	78.0	78.0	0.793	82.9	85.4	0.801	82.9	90.2	0.527	48.8	53.7
	InstructBLIP	1.000	100.0	100.0	0.962	95.1	95.1	0.750	70.7	70.7	0.762	63.4	78.0	0.695	58.5	70.7	0.500	41.5	41.5
	LLaVA-1.5	1.000	100.0	100.0	1.000	100.0	100.0	0.930	92.7	92.7	0.952	92.7	92.7	1.000	100.0	100.0	0.903	87.8	87.8
OVIT	MiniGPT-4	0.702	68.3	68.3	0.731	68.3	73.2	0.641	60.0	63.4	0.620	56.1	58.5	0.450	36.6	58.5	0.552	48.8	53.7
511	BLIP-2	0.834	85.5	90.2	0.884	97.6	97.6	0.772	80.5	80.5	0.779	80.5	87.8	0.803	85.4	85.4	0.533	51.2	53.7
	InstructBLIP	0.980	97.6	97.6	0.958	95.1	95.1	0.638	56.1	56.1	0.787	70.7	80.5	0.770	65.9	79.0	0.557	46.3	46.3
	LLaVA-1.5	0.977	97.6	97.6	1.000	100.0	100.0	0.992	95.1	100.0	0.978	97.6	97.6	1.000	100.0	100.0	0.953	95.1	97.6
VO 42	MiniGPT-4	0.634	61.0	63.4	0.713	70.7	70.7	0.612	53.7	61.0	0.539	46.3	51.2	0.568	53.7	56.1	0.522	43.9	51.2
vQAV2	BLIP-2	0.849	90.2	92.7	0.837	87.8	92.7	0.648	58.5	68.3	0.775	82.9	82.9	0.819	87.8	87.8	0.588	53.7	61.0
	InstructBLIP	1.000	100.0	100.0	0.947	92.7	95.1	0.480	34.1	34.1	0.828	70.7	90.2	0.724	61.0	75.6	0.404	24.4	24.4

operation-specific gradients from transformations to update perturbations, which is more practical to be exploited in a black-box setting and can *explicitly learn transformation-only adversarial impacts*.

Evaluation and Discussion. To evaluate our designed adversarial transformation-aware attack, we conduct experiments in both untargeted and targeted attack settings and can conclude that:

(i) As for the untargeted attack, this adversarial learning attack is more effective and efficient than
the hybrid transformation-based attack. Although the hybrid transformation-based attack tries to
enumerate possible transformation combinations and manually construct the most harmful operations, it costs lots of resources and may stuck into the local optimum. Instead, this adversarial
transformation-aware attack can adaptively learn the most harmful impacts from all potential transformation combinations, leading to more harmful adversarial generations as shown in Figure 5.

(ii) This adversarial learning attack is more flexible and can mislead the LVLMs output attacker(ii) This adversarial learning attack is more flexible and can mislead the LVLMs output attacker(ii) This adversarial learning attack is more flexible and can mislead the LVLMs output attacker(iii) This adversarial learning attack is more flexible and can mislead the LVLMs output attacker(iii) This adversarial learning attack is more flexible and can mislead the LVLMs output attacker(iii) This adversarial transformation to the above untargeted adversarial generation, we also investigate whether
(iii) Table 2, we preset the target response as "I am sorry" and experimental results indicate that our
(iii) adversarial transformation-aware attack is effective in achieving targeted attack with significant performance, having great potential to be deployed in real-world LVLM applications. Besides, it also
(iii) shows that our attack can still achieve competitive performances with a lower perturbation budget.

431 Overall, our attacks show superior performance, demonstrating the effectiveness of our carefully designed transformation-aware attacks. More experimental results are in Appendix C.2.

449 450

451

452

Dataset	LVLM	Hybrid	Transformatio	on-based A	Attack (↓)	Adversarial Transformation-aware Attack (\downarrow)					
		LLaVA-1.5	MiniGPT-4	BLIP-2	InstructBLIP	LLaVA-1.5	MiniGPT-4	BLIP-2	InstructBLIP		
	LLaVA-1.5	0.717	0.710	0.611	0.763	0.463	0.697	0.678	0.728		
DALLE	MiniGPT-4	0.808	0.519	0.626	0.752	0.770	0.460	0.399	0.456		
DALL-L	BLIP-2	0.782	0.707	0.472	0.701	0.702	0.676	0.372	0.672		
	InstructBLIP	0.789	0.695	0.595	0.662	0.693	0.661	0.565	<u>0.542</u>		
	LLaVA-1.5	0.664	0.663	0.563	0.660	0.439	0.695	0.683	0.708		
CVIT	MiniGPT-4	0.773	0.498	0.613	0.682	0.711	0.397	0.338	0.450		
311	BLIP-2	0.742	0.673	0.458	0.679	0.650	0.553	0.300	0.519		
	InstructBLIP	0.747	0.659	0.565	0.609	0.663	0.545	0.383	0.427		
	LLaVA-1.5	0.663	0.646	0.643	0.742	0.394	0.657	0.662	0.694		
VOAv2	MiniGPT-4	0.777	0.437	0.663	0.721	0.700	0.349	0.321	0.477		
v QAV2	BLIP-2	0.754	0.653	0.529	0.672	0.641	0.518	0.314	0.484		
	InstructBLIP	0.759	0.634	0.640	0.590	0.632	0.518	0.431	0.442		

Table 3: Investigation on the adversarial transferability of our two types of transformation attacks.
Following previous works, we evaluate the untargeted attack performance of the adversarial samples
generated on the LVLMs of rows and tested on the LVLMs of columns. The experimental results
are calculated by the averaged semantic similarities on three tasks. "value" is basic performance.

4.4 IN-DEPTH ANALYSIS OF OUR PROPOSED TRANSFORMATION ATTACKS

In this section, we provide a detailed analysis of our proposed two types of transformation attacks from perspectives of complexity, adversarial transferability, and adversarial robustness, respectively.

Takeaways: O Our proposed attack methods in Section 4 are quite efficient. Besides, our adversarial learning based attack variant is more efficient than the manual constructing one while achieving better performance. Our developed adversarial transformation attacks can achieve significant transferability among different black-box LVLM models. Sexperimental results also illustrate that our two transformation attacks are robust to potential defense strategies.

458 Analysis on Complexity. We first investigate the 459 complexity of our proposed two types of LVLM attacks. As shown in Table 4, we evaluate the usage 460 of GPU time and memory of a single adversarial 461 sample on both generation and inference processes. 462 It indicates that the adversarial transformation-463 aware attack is much more efficient than the hybrid 464 transformation-based attack during the adversarial 465 generation, as the former can adaptively learn the 466 potentially harmful transformation impacts (but re-

Table 4: Complexity analysis on our attacks. We evaluate the GPU time and memory usage of a single adversarial example on both generation and inference processes on LLaVA-1.5.

Process	Attack Type	GPU Time	GPU Memory
Generation	Hybrid Attack	9min	16GB
	Adversarial Attack	5min	22GB
Inference	Hybrid Attack	3s	15GB
	Adversarial Attack	3s	15GB

- quires relatively more memory for gradient approximation) while the latter relies on lots of manual
 efforts. Since both two attacks solely feed the adversarial sample into the LVLM without any addi tional operation during the inference, they have the same complexity in inference.
- Analysis on Adversarial Transferability. We then investigate the adversarial transferability of the generated adversarial examples of our two attacks. As shown in Table 3, we can conclude that:
- (i) Developing LVLM attacks using visual transformations can achieve significant adversarial transferability. Our two types of transformation attacks achieve great transfer-attack performance when we directly feed the generated adversarial examples of one LVLM to the other three LVLMs. Although the output textual semantic similarities relatively decrease, its influences are largely inferior to the performance drops brought by our attacks. Therefore, utilizing visual transformations to construct LVLM attacks is a promising way to improve the adversarial transferability.
- (ii) The adversarial transformation-aware attack achieves better transferability than the hybrid
 transformation-based attack. Since the transformation attack with adversarial learning mechanism
 can adaptively learn more potential transformation operations than the hybrid manual constructing
 one, it will learn more generalizable transformation impacts thus leading to better transferability.
- Analysis on Adversarial Robustness. At last, we investigate the robustness of the proposed two
 transformation attacks. In particular, we implement three pre-processing defenses, *i.e*, Randomiza tion (Frosio & Kautz, 2023; Xie et al., 2017), JPEG Compression (Guo et al., 2017), and Diffusion
 Restoration (Nie et al., 2022). As shown in Table 5, we can conclude that:

Dataset	Defense	Hybrid	Transformatio	on-based A	Attack (↓)	Adversarial Transformation-aware Attack (\downarrow)					
		LLaVA-1.5	MiniGPT-4	BLIP-2	InstructBLIP	LLaVA-1.5	MiniGPT-4	BLIP-2	InstructBLIP		
	No Defense	0.717	0.519	0.472	0.662	0.463	0.460	0.372	0.542		
DALL-E	Randomization	0.767	0.686	0.581	0.674	0.744	0.591	0.505	0.631		
DALL-E	JPEG Compre.	0.776	0.681	0.548	0.682	0.609	0.565	0.468	0.604		
	Diffusion	0.541	0.518	0.382	0.421	0.758	0.701	0.594	0.697		
	No Defense	0.664	0.498	0.458	0.609	0.439	0.397	0.300	0.427		
OVIT	Randomization	0.763	0.658	0.545	0.637	0.648	0.560	0.406	0.503		
311	JPEG Compre.	0.757	0.648	0.508	0.631	0.547	0.487	0.359	0.471		
	Diffusion	0.555	0.514	0.397	0.440	0.690	0.611	0.532	0.617		
	No Defense	0.663	0.437	0.529	0.590	0.394	0.349	0.314	0.442		
NO. A	Randomization	0.746	0.609	0.599	0.696	0.610	0.467	0.409	0.512		
VQAV2	JPEG Compre.	0.760	0.591	0.536	0.675	0.502	0.469	0.372	0.485		
	Diffusion	0.536	0.464	0.392	0.423	0.675	0.571	0.550	0.608		

Table 5: Investigation on the adversarial robustness of our two types of transformation attacks. Following previous works, we evaluate the untargeted attack performance of the adversarial samples generated on the LVLMs of columns by testing them against potential defenses of rows. The experimental results are calculated by the averaged semantic similarities on three tasks.

(*i*) Our proposed attack methods are robust to potential defense strategies. According to the performances in this table, although the three defense methods are able to degenerate our attack performance, their influences are largely inferior to the performance drops brought by our attack. This demonstrates that our attack is fairly resistant to the potential defense methods in practice.

507 (*ii*) The hybrid attack produces even more harmful results under diffusion restoration. It is because
508 the applied transformations destroy the image structure, so diffusion operation will further generate
509 more diverse content. Instead, diffusion can alleviate the harmful impact of adversarial noise.

510 511

504

505

506 507

5 DISCUSSION

512

Justification of Our Experiments. Since our main goal is to investigate the adversarial robustness of LVLMs to visual transformations, our experiments are solely conducted on the comparisons and analysis between different transformation strategies. We do not compare performances with other types of LVLM attacks as: (1) They are designed with more complicated perturbation patterns. Directly comparing our solely transformation-based attacks with them is unfair. (2) They are diversely implemented in different settings with the usage of different LVLM models and datasets. We provide case-by-case comparisons with other LVLM attacks under the same settings in Appendix D.

Limitations. 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 attacks in
 such complicated cases would be one of the future directions for evaluating the LVLM security.

Broader Impacts. While the primary goal of our research is to generate superior adversarial trans formations against large vision-language models, it is possible that the developed attacking strate gies could be misused to evade practically deployed systems and cause potential negative societal
 impacts. Specifically, our adversarial threat model assumes targeted responses, which involves ma nipulating 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.

531 532

533

6 CONCLUSION

In conclusion, this paper offers novel insights into the vulnerability of LVLMs to visual transformations. Our comprehensive evaluation indicates that different transformations share diverse harmfulness degrees on existing LVLMs while appropriate transformation combinations can boost the attack performance. We also take a further step to investigate how to manually construct a more harmful transformation operation and how to adaptively learn to impose adversarial impacts from all potential transformations to raw images for improving the attack effectiveness and imperceptibility. We envision our findings will pave the way for the development of efficient and effective LVLM attacks.

540 REFERENCES

549

566

567

568

569

585

586

- 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.
- MaungMaung AprilPyone and Hitoshi Kiya. Block-wise image transformation with secret key for
 adversarially robust defense. *IEEE Transactions on Information Forensics and Security*, 16:2709–
 2723, 2021.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani
 Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An opensource framework for training large autoregressive vision-language models. arXiv preprint
 arXiv:2308.01390, 2023.
- Eugene Bagdasaryan, Tsung-Yin Hsieh, Ben Nassi, and Vitaly Shmatikov. (ab) using images and sounds for indirect instruction injection in multi-modal llms. *arXiv preprint arXiv:2307.10490*, 2023.
- 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.
- Yang Chen, Ethan Mendes, Sauvik Das, Wei Xu, and Alan Ritter. Can language models be instructed to protect personal information? *arXiv preprint arXiv:2310.02224*, 2023.
- Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment:
 Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 113–123, 2019.
 - Xuanimng Cui, Alejandro Aparcedo, Young Kyun Jang, and Ser-Nam Lim. On the robustness of large multimodal models against image adversarial attacks. *arXiv preprint arXiv:2312.03777*, 2023.
- 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 visionlanguage models with instruction tuning. *Advances in Neural Information Processing Systems*,
 36, 2024.
- 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, pp. 248–255. Ieee, 2009.
- 578 Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boost579 ing adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer*580 *vision and pattern recognition*, pp. 9185–9193, 2018.
- 581
 582
 582
 583
 584
 584
 584
 581
 585
 586
 586
 586
 587
 588
 588
 588
 589
 589
 589
 580
 580
 581
 581
 582
 583
 584
 584
 584
 584
 584
 584
 584
 584
 585
 586
 586
 586
 587
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 - 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.
- Ranjie Duan, Yuefeng Chen, Dantong Niu, Yun Yang, A Kai Qin, and Yuan He. Advdrop: Adversarial attack to dnns by dropping information. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7506–7515, 2021.
- 592 Yihe Fan, Yuxin Cao, Ziyu Zhao, Ziyao Liu, and Shaofeng Li. Unbridled icarus: A survey of
 593 the potential perils of image inputs in multimodal large language model security. *arXiv preprint arXiv:2404.05264*, 2024.

594 Iuri Frosio and Jan Kautz. The best defense is a good offense: adversarial augmentation against 595 adversarial attacks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 596 Recognition, pp. 4067-4076, 2023. 597 Kuofeng Gao, Yang Bai, Jiawang Bai, Yong Yang, and Shu-Tao Xia. Adversarial robustness for 598 visual grounding of multimodal large language models. arXiv preprint arXiv:2405.09981, 2024a. 600 Kuofeng Gao, Yang Bai, Jindong Gu, Shu-Tao Xia, Philip Torr, Zhifeng Li, and Wei Liu. Induc-601 ing high energy-latency of large vision-language models with verbose images. arXiv preprint 602 arXiv:2401.11170, 2024b. 603 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa 604 matter: Elevating the role of image understanding in visual question answering. In Proceedings 605 of the IEEE conference on computer vision and pattern recognition, pp. 6904–6913, 2017. 606 607 Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens Van Der Maaten. Countering adversarial 608 images using input transformations. arXiv preprint arXiv:1711.00117, 2017. 609 610 Qi Guo, Shanmin Pang, Xiaojun Jia, and Qing Guo. Efficiently adversarial examples generation for 611 visual-language models under targeted transfer scenarios using diffusion models. arXiv preprint arXiv:2404.10335, 2024. 612 613 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual net-614 works. In Computer Vision-ECCV 2016: 14th European Conference, Amsterdam, The Nether-615 lands, October 11-14, 2016, Proceedings, Part IV 14, pp. 630-645. Springer, 2016. 616 617 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 25, 2012. 618 619 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image 620 pre-training with frozen image encoders and large language models. In International conference 621 on machine learning, pp. 19730-19742. PMLR, 2023. 622 623 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 624 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer 625 Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pp. 740–755. Springer, 2014. 626 627 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances 628 in neural information processing systems, 36, 2024a. 629 630 Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. Safety of multimodal large language 631 models on images and text. arXiv preprint arXiv:2402.00357, 2024b. 632 Dong Lu, Tianyu Pang, Chao Du, Qian Liu, Xianjun Yang, and Min Lin. Test-time backdoor attacks 633 on multimodal large language models. arXiv preprint arXiv:2402.08577, 2024. 634 635 Haochen Luo, Jindong Gu, Fengyuan Liu, and Philip Torr. An image is worth 1000 lies: Adversar-636 ial transferability across prompts on vision-language models. arXiv preprint arXiv:2403.09766, 637 2024. 638 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 639 Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083, 640 2017. 641 642 Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal 643 adversarial perturbations. In Proceedings of the IEEE conference on computer vision and pattern 644 recognition, pp. 1765–1773, 2017. 645 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, 646 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with 647 text-guided diffusion models. arXiv preprint arXiv:2112.10741, 2021.

667

680

686

- 648 Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. 649 Diffusion models for adversarial purification. arXiv preprint arXiv:2205.07460, 2022. 650
- Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. 651 Visual adversarial examples jailbreak aligned large language models. In Proceedings of the AAAI 652 Conference on Artificial Intelligence, volume 38, pp. 21527–21536, 2024. 653
- 654 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, 655 and Ilya Sutskever. Zero-shot text-to-image generation. In International conference on machine 656 *learning*, pp. 8821–8831. Pmlr, 2021. 657
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-658 conditional image generation with clip latents. arXiv preprint arXiv:2204.06125, 1(2):3, 2022. 659
- 660 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-661 networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language 662 Processing. Association for Computational Linguistics, 11 2019. URL https://arxiv. 663 org/abs/1908.10084. 664
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-666 ence on computer vision and pattern recognition, pp. 10684–10695, 2022.
- 668 Christian Schlarmann and Matthias Hein. On the adversarial robustness of multi-modal foundation 669 models. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 670 3677-3685, 2023. 671
- Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Jailbreak in pieces: Compositional adversarial 672 attacks on multi-modal language models. In The Twelfth International Conference on Learning 673 Representations, 2023. 674
- 675 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image 676 recognition. arXiv preprint arXiv:1409.1556, 2014. 677
- 678 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, 679 and Rob Fergus. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199, 2013.
- Xijia Tao, Shuai Zhong, Lei Li, Qi Liu, and Lingpeng Kong. Imgtrojan: Jailbreaking vision-681 language models with one image. arXiv preprint arXiv:2403.02910, 2024. 682
- 683 Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Mul-684 timodal few-shot learning with frozen language models. Advances in Neural Information Pro-685 cessing Systems, 34:200-212, 2021.
- Haoqin Tu, Chenhang Cui, Zijun Wang, Yiyang Zhou, Bingchen Zhao, Junlin Han, Wangchunshu 687 Zhou, Huaxiu Yao, and Cihang Xie. How many unicorns are in this image? a safety evaluation 688 benchmark for vision llms. arXiv preprint arXiv:2311.16101, 2023. 689
- 690 Haodi Wang, Kai Dong, Zhilei Zhu, Haotong Qin, Aishan Liu, Xiaolin Fang, Jiakai Wang, and 691 Xianglong Liu. Transferable multimodal attack on vision-language pre-training models. In 2024 692 *IEEE Symposium on Security and Privacy (SP)*, pp. 102–102. IEEE Computer Society, 2024a.
- Xiaosen Wang, Zeliang Zhang, and Jianping Zhang. Structure invariant transformation for better ad-694 versarial transferability. In Proceedings of the IEEE/CVF International Conference on Computer 695 Vision, pp. 4607-4619, 2023a. 696
- 697 Xunguang Wang, Zhenlan Ji, Pingchuan Ma, Zongjie Li, and Shuai Wang. Instructia: Instruction-698 tuned targeted attack for large vision-language models. arXiv preprint arXiv:2312.01886, 2023b. 699
- Zefeng Wang, Zhen Han, Shuo Chen, Fan Xue, Zifeng Ding, Xun Xiao, Volker Tresp, Philip Torr, 700 and Jindong Gu. Stop reasoning! when multimodal llms with chain-of-thought reasoning meets 701 adversarial images. arXiv preprint arXiv:2402.14899, 2024b.

702 Yuanwei Wu, Xiang Li, Yixin Liu, Pan Zhou, and Lichao Sun. Jailbreaking gpt-4v via self-703 adversarial attacks with system prompts. arXiv preprint arXiv:2311.09127, 2023. 704 705 Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. arXiv preprint arXiv:1711.01991, 2017. 706 707 Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. 708 Cutmix: Regularization strategy to train strong classifiers with localizable features. In Proceed-709 ings of the IEEE/CVF international conference on computer vision, pp. 6023–6032, 2019. 710 711 Hao Zhang, Wenqi Shao, Hong Liu, Yongqiang Ma, Ping Luo, Yu Qiao, and Kaipeng Zhang. 712 Avibench: Towards evaluating the robustness of large vision-language model on adversarial 713 visual-instructions. arXiv preprint arXiv:2403.09346, 2024. 714 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical 715 risk minimization. arXiv preprint arXiv:1710.09412, 2017. 716 717 Bo Zhao, Boya Wu, and Tiejun Huang. Svit: Scaling up visual instruction tuning. arXiv preprint 718 arXiv:2307.04087, 2023. 719 720 Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min 721 Lin. On evaluating adversarial robustness of large vision-language models. Advances in Neural Information Processing Systems, 36, 2024. 722 723 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: En-724 hancing vision-language understanding with advanced large language models. arXiv preprint 725 arXiv:2304.10592, 2023. 726 727 728 ADDITIONAL DETAILS OF LVLMS' PERFORMANCE ON DIFFERENT Α 729 VISUAL TRANSFORMATIONS 730 731 A.1 MORE DETAILS OF VISUAL TRANSFORMATIONS 732 733 Here we provide the implementation details of the previously mentioned transformations, respectively. In particular, for a given image block $x \in \mathbb{R}^{3 \times H \times W}$, we can implement the basic transfor-734 735 mations as follows: 736 **1. Resize Large:** We resize the original image x into x' with the size of $3 \times h \times w$ (h > w, w > w)737 using bilinear interpolation. 738 **2. Resize Small**: We resize the original image x into x' with the size of $3 \times h \times w$ (h < w, w < w). 739 740 **3. HFlip** (Horizontal Flip): We flip the image x horizontally along the vertical axis, in which the 741 left of the image becomes the right, and the right becomes the left. 742 4. VFlip (Vertical Flip): We flip the image x vertically along the horizontal axis, in which the top 743 of the image becomes the bottom, and the bottom becomes the top. 744 5. Rotate_Random: We rotate the image x by a random angle around its center point. 745 746 6. Rotate 180° : We turn the image x clockwise by 180° around its center point, in which the top-left 747 of the image becomes the bottom-right, and the top-right becomes the bottom-left. 748 7. VShift_Random (Vertical Shift): We roll the image x along the vertical axis by a randomly 749 selected length h < H. 750 751 8. VShift_Half (Vertical Shift): We roll the image x along the vertical axis at $h = \frac{1}{2}H$. 752 **9.** HShift_Random (Horizontal Shift): We roll the image block x along the horizontal axis by a 753 randomly selected length of w < W. 754 10. HShift_Half (Horizontal Shift): We roll the image block x along the horizontal axis at w =755 $\frac{1}{2}W.$

756 **11.** Scale: We multiply a random scale factor $\alpha \in (0, 1)$ with the pixel in the image to scale x into 757 $\alpha \cdot x$. 758 **12.** Add Noise: We add a uniform noise $r \in [0, 1]^{3 \times H \times W}$ to the image x and clip them into [0, 1]759 to obtain the transformed image Clip(x + r, 0, 1). 760 761 **13. Dropout:** We set each channel of the image x to zero with a probability of 10%. 762 **14.** ColorJitter: We randomly change the brightness, contrast, saturation, and hue of the image x. 763 **15.** DCT: We first transform x to the frequency domain using Discrete Cosine Transformation 764 (DCT). Then we mask the top 40% highest frequency with 0 and recover the image in the time 765 domain using Inverse Discrete Cosine Transformation (IDCT). 766 767 Then, we adapt the above transformations into the block-level transformation (AprilPyone & Kiya, 2021). Specifically, we first uniformly split the image $x \in \mathbb{R}^{3 \times H \times W}$ into 3×3 patches of the same 768 sizes, then perform the above operations on each patch as follows: 769 770 **16.** Block Resize: We resize each patch to size $3 \times h \times w$ $(h < \frac{1}{2}H, w < \frac{1}{2}H)$ and utilize bicubic 771 interpolation to reconstruct the patch into the original size. 772 17. Block_HFlip: We flip each patch along along the vertical axis, in which the left of the patch 773 becomes the right, and the right becomes the left. 774 775 **18.** Block_VFlip: We flip each patch along the horizontal axis, in which the top of the patch becomes the bottom, and the bottom becomes the top. 776 777 **19.** Block_Rotate: We turn each patch clockwise by 180° around its center point, in which the 778 top-left of the patch becomes the bottom-right, and the top-right becomes the bottom-left. 779 **20.** Block_VShift: We roll each patch along the vertical axis at half height. 781 21. Block_HShift: We roll each patch along the horizontal axis at half weight. 782 **22.** Block Scale: We multiply a random scale factor $\alpha \in (0, 1)$ with the pixel of each patch. 783 784 **23.** Block_AddNoise: We add a uniform noise to each patch and clip them into [0, 1] to obtain the transformed patch. 785 786 **24. Block_Dropout**: We set each channel of the patch to zero with a probability of 10%. 787 25. Block_ColorJitter: We randomly change the brightness, contrast, saturation, and hue of each 788 patch. 789 790 **26.** Block_DCT: we first transform each patch to the frequency domain using Discrete Cosine Transformation (DCT). Then we mask the top 40% highest frequency with 0 and recover the patch 791 in the time domain using Inverse Discrete Cosine Transformation (IDCT). 792 793 Based on our empirical experience, we find that transformations of Rotate, Hshilt, VFlip perform 794 more adversarial among the above multiple block-level transformations. Therefore, we further design various types of their combinations in the following: 796 27. Block_Rotate_HShift: We randomly choose one of the rotation or horizontal shift operations 797 for each patch. 798 28. Block_Rotate_VFlip: We randomly choose one of the rotation or vertical flip operations for 799 each patch. 800 801 **29.** Block_VFlip_HShift: We randomly choose one of the vertical flip or horizontal shift operations 802 for each patch. 803 **30.** Block_Rotate_VFlip_HShift: We randomly choose one of the rotation, vertical flip, or horizon-804 tal shift operations for each patch. 805 **31.** Block_Random_Combination: We randomly choose the combination of the transformations 806 including rotation, vertical flip, or horizontal shift operations for each patch. 807 808 We provide visual examples of the various image transformations described above in Figure 6. 809



A.2 VISUALIZATION ON THE TEXTUAL OUTPUTS OF TRANSFORMED IMAGES

We provide the textual outputs of each transformed image in Figure 7. We can find that the general visual transformation can affect the LVLM's textual output.

B ADDITIONAL DETAILS OF OUR PROPOSED HYBRID TRANSFORMATION-BASED ATTACK

853 854

855

856 857

858 859

861

- B.1 MORE ILLUSTRATIONS OF OUR HYBRID TRANSFORMATION-BASED ATTACK
- 863 We provide a step-by-step visualization of our proposed hybrid transformation-based attack. As shown in Figure 8, our hybrid transformation-based attack transforms each patch one by one in the



Figure 7: Visualization of the LVLM's textual outputs of different transformation operations. Red: the affected outputs are different from the raw answer.

907 908

909 default order to iteratively make the transformed image as harmful as possible. Starting from the first 910 patch, we fix the remaining patches unchanged and perform the above 15 transformation operations 911 in sequence to obtain the corresponding 15 transformed images. Then each transformed image is 912 fed into the LVLM model individually with the same textual prompt to obtain the corresponding 15 913 adversarial answers. Next, we calculate the semantic similarities between these adversarial answers 914 and the original answer, and select the operation with the lowest similarity score as the optimal (most 915 harmful) transformation operation for this patch. By fixing the transformed patch 1, we repeat this process for patch 2 to further degenerate the LVLM's performance. After traversing all patches, we 916 can generate the most harmful transformation operation on the image input. Note that, each step 917 operation can effectively further degenerate the LVLM's performance compared to its previous step.



Figure 8: Illustration of our designed hybrid transformation-based attack, which manually constructs the most harmful transformation combination via enumeration.



Figure 9: Illustration of our designed hybrid transformation-based attack implemented in a universal setting, where the most harmful transformation operation is explored to be the same among all image-text inputs.

949 950 951

952 953

954

957

939

940 941 942

943

944 945

946

947

948

B.2 UNIVERSAL ATTACK FOR OUR HYBRID TRANSFORMATION-BASED ATTACK

Generally, our proposed hybrid transformation-based attack is implemented in a single-image attack setting, where the most harmful transformation operation varies among different image-text 955 inputs. Further, we can also extend this attack into a universal attack setting as shown in Figure 9, 956 where the most harmful transformation operation is explored to be the same among all image-text inputs. Specifically, we follow the traditional universal setting (Moosavi-Dezfooli et al., 2017) to assess the vulnerability of each transformation based on its averaged impacts on the whole test set. 958 Corresponding performance is shown in Figure 10, we can conclude that: 959

960 (i) Our hybrid transformation-based attack in a universal setting is also more harmful than gen-961 eral transformation combinations. Compared with the previous 31 transformations in Section 3, 962 our hybrid transformation-based attacks in a universal setting can further degenerate the LVLM's 963 performance on all models across all datasets. This significant similarity decrease also demonstrates that manually constructing transformation operations is more effective in generating more harmful 964 adversarial examples. 965

966 (ii) As for our hybrid transformation-based attack, the single-image attack setting is more effective 967 than the universal setting to generate more harmful transformation operations. By comparing the 968 attack performance between the single-image setting in Figure 4 and the universal setting in Figure 10, we can find that the single-image attack is more flexible and harmful than the universal 969 attack setting, thus achieving better attack performance. This is because the single-image attack 970 can straightforwardly conduct the most vulnerable transformation impacts on each image while the 971 universal attack fails to cover the distribution gaps among diverse images.



Figure 10: Untargeted attack performance of our designed hybrid transformation-based attack implemented in a universal setting. Lower similarities (\downarrow) indicate more harmful impacts. Numbers in front of the bars refer to the similarity score decrease compared to the corresponding best transformations in Section 3, larger decrease indicates greater harmfulness.

B.3 VISUALIZATION ON THE TEXTUAL OUTPUTS OF TRANSFORMED IMAGES

To validate the effectiveness of our hybrid transformation-based attack, we provide the visualization on the textual outputs of the transformed images. As shown in Figure 11, we can find that our method can effectively mislead LVLMs to output wrong texts that are semantically distinct from the original texts.

С ADDITIONAL DETAILS OF OUR PROPOSED ADVERSARIAL TRANSFORMATION-AWARE ATTACK

MORE ILLUSTRATIONS OF OUR ADVERSARIAL TRANSFORMATION-AWARE ATTACK C.1

Although LVLMs are insensitive and invariant to the block-level transformations as LVLMs can 1001 globally handle and learn the tokenize-wise semantics, the hybrid transformation-based attack still 1002 introduces noticeable and unnatural appearances to humans. Therefore, to improve the impercepti-1003 bility of the transformed images, we propose to impose the impacts of potentially harmful transfor-1004 mations as perturbations to add to the raw image while keeping the same adversarial effect as those 1005 transformations. In this manner, the perturbed image is optimized by the approximated gradient directions of harmful transformations and shares the same adversarial impacts, achieving the goal of 1007 adversarial transformation while improving visual imperceptibility. Specifically, instead of directly 1008 calculating the gradient on the input image, we calculate the gradient only on the truly harmful transformations. More adversarial training details can be found in Algorithm 1.

1010

985

986

987

988 989

990

991

992

993

994 995

996 997

998 999

1000

1011 C.2 MORE EXPERIMENTS ON OUR ADVERSARIAL TRANSFORMATION-AWARE ATTACK 1012

1013 As for the untargeted attack, the experiments in Figure 5 are implemented with perturbation budget 1014 $\epsilon = 16/255$. We further provide the corresponding experiments with perturbation budget $\epsilon =$ 32/255 in Figure 12. We can find that our adversarial transformation-aware attack can achieve 1015 better attack performance with larger budgets. 1016

1017 As for the targeted attack, we also perform more experiments on different target text output of "I 1018 want to destroy the world" in Table 6 and "I don't know" in Table 7. Corresponding results show 1019 that our attack can still achieve significant adversarial performance, demonstrating the scalability, generalizability, and practicality of our proposed attack method. 1020

1021

C.3 VISUALIZATION ON THE TEXTUAL OUTPUTS OF TRANSFORMED IMAGES

1023

To validate the effectiveness of our adversarial transformation-aware attack, we provide the visu-1024 alization on the textual outputs of the transformed images in both untargeted and targeted attack 1025 settings. As shown in Figure 11, we can find that our method can not only effectively mislead



I don't know



Adversarial (Target)

t to destroy

the world.

Adversarial (Target)

I want to destroy

the world.

Adversarial (Target)

I don't know

Adversarial (Target)

361





1076 Figure 11: Visualization of the LVLM's textual outputs. Green: hybrid transformation-based attack. 1077 Red and Blue: adversarial transformation-aware attack on untargeted and targeted attack settings. 1078

people and their pets.

080	Alg	lgorithm 1 Adversarial Transformation-aware	e Attack	
081	Inp	nput: The source sample with the raw imag	e x_v , the text input x_t a	nd the raw answer y ; the
082	los	oss function of LVLM L ; the number of iteration	on T ; the maximum pert	urbation ϵ ; decay factor μ ;
083	the	he number of transformed images N ; the step	size α ; the random trans	sformation image function
084	Tr	Trans.		
085	Ou	Dutput: Adversarial image.		
086	1:	1: Initialize gradient $g_0 = 0$, adversarial samp	le $x_{v,0}^{adv} = x_v$	
087	2:	2: for $t = 0$ to $T - 1$ do	,	
880	3:	$3: \qquad g = 0$		
089	4:	4: for $i = 0$ to $N - 1$ do		
090	5:	5: Construct transformed image as Tr	$ans_i(x_{v,t}^{aav})$	_
091	6:	6: Calculate the loss before and after t	he transformation by $l_1 =$	$L(Trans_i(x_{v,t}^{adv}), x_t, y)$
092	7:	7: and $l_2 = L(x_{v,t}^{adv}, x_t, y)$		
093	8:	8: if $l_1 > l_2$ then		
094	9:	9: Get the harmful weight as $w_i =$	1	
095	10:	0: else		
096	11:	1: Get the harmless weight as $w_i =$	= 0	
097	12:	2: end if		
098	13:	3: Approximate the gradient by $grad_i$	$= \nabla_{x_{v,t}^{adv}} L(Trans_i(x_{v,t}^{adv})$	(x_t, y)
099	14:	4: Sum the gradients as $g = g + w_i \cdot g$	rad_i	
100	15:	5: end for		
101	16:	6: Get the average gradients as $g = \frac{1}{N} \cdot g$	a	
102	17:	7: Update the momentum by $g_{t+1} = \mu \cdot g_t$	$\frac{g}{\ g\ _1}$	
103	18:	8: Update the adversarial image by $x_{v,t+1}^{adv}$	$= \operatorname{Clip}(x_{v,t}^{adv} + \alpha \cdot \operatorname{sign}(q))$	$g_{t+1}), 0, 1)$
104	19:	9: end for	- ,-	
105	20:	0: return transformation-aware adversarial sa	mple $x_{v,T}^{adv}$	
105			,	
107		LLaVA-1.5 MiniGPT-4	BLIP-2	InstructBLIP
102	щ	Image Captioning 0.371 (-0.112) Image Captioning 0.402 (-0.064)	Image Captioning 0.280 (-0.063)	Image Captioning 0.541 (-0.053)
100	- ALL-	Image Classification 0.332 (-0.116) Image Classification 0.388 (-0.043)	Image Classification 0.329 (-0.061)	Image Classification 0.421 (-0.052)
109	D	VQA 0.687 (-0.046) VQA 0.591 (-0.04) VQA 0.507 (-0.065)	VQA 0.665 (-0.064)
110		0.00 0.25 0.50 0.75 1.00 0.00 0.25 0.50 0.75 Semantic Similarity Scores Semantic Similarity Sco	1.00 0.00 0.25 0.50 0.75 1 res Semantic Similarity Scores	.00 0.00 0.25 0.50 0.75 1.00 Semantic Similarity Scores
111		Image Captioning 0.355 (-0.045) Image Captioning 0.314 (-0.045)	Image Captioning 0.280 (-0.089)	Image Captioning 0.388 (-0.095)
112	TIVS I	Image Classification 0.297 (-0.056) Image Classification 0.346 (-0.055)	Image Classification 0.278 (-0.061)	Image Classification 0.301 (-0.061)
113	v ₂	VQA 0.665 (-0.041) VQA 0.532 (-0.045)	VQA 0.341 (-0.072)	VQA 0.593 (-0.058)

0.00 0.25 0.50 0.75 1.00 1114 0.00 0.25 0.50 0.75 1.00 0.00 0.25 0.50 0.75 1.00 0.00 0.25 0.50 0.75 1.00
 Semanuc summer
 Image Captioning
 0.356
 (-0.070)

 Image Captioning
 0.356
 (-0.046)
 Image Captioning
 0.356
 (-0.056)

 Image Classification
 0.351
 (-0.056)
 (-0.056)
 0.618
 (-1.056)
 Semantic Similarity Scores Semantic Similarity Scores Semantic Similarity Scores Image Captioning 0.329 (-0.046) 1115 Image Captioning 0.290 (-0.076) Image Classification 0.285 (-0.033) Image Classification 0.254 (-0.046) 1116 Image Classification 0.317 (-0.054) 0.607 (-0.049) 0.400 (-0.053) VQA 0.451 (-0.056) VQA 0.618 (-0.056) 1117 VQA VQA 0.00 0.25 0.50 0.75 1.00 Semantic Similarity Scores 1118 1119 Figure 12: Untargeted attack performance of our designed adversarial transformation-aware attack

Figure 12: Untargeted attack performance of our designed adversarial transformation-aware attack implemented on the $\epsilon = 32/255$ setting. Lower similarities (\downarrow) indicate more harmful impacts. Numbers in front of the bars refer to the similarity score decrease compared to the adversarial transformation-aware attack variant implemented with $\epsilon = 16/255$ in Figure 5, larger decrease indicates greater harmfulness.

LVLMs to output wrong texts that are semantically distinct from the original texts, but also mislead
 LVLMs to output specific attacker-chosen texts.

1127 1128

1129

D PERFORMANCE COMPARISON WITH EXISTING LVLM ATTACKS

To further investigate the effectiveness of our proposed attack, we provide a case-by-case comparison with existing LVLM attacks MF (Zhao et al., 2024) (NeurIPS 2023) and CroPA (Luo et al., 2024) (ICLR 2024). Since existing LVLM attacks are diversely implemented in different settings with the usage of different LVLM models and datasets, we re-implement our transformation-aware attack into their same settings/datasets/metrics for fair comparison. Note that, since MF and CroPA

1151 1152 1153

Table 6: Targeted attack performance of our adversarial transformation-aware attack. Target text: 1135 "I want to destroy the world". Larger similarities (\uparrow) indicate more harmful impacts. "SS" denotes 1136 the semantic similarity between output and target text; "EM" and "CC" denote the "ExactMatch" 1137 and "ConditionalContain" metrics to assess whether the output exactly matches or contain the target 1138 text, respectively. 1139

			Р	erturba	tion Bı	idget ϵ	= 32	/255				Р	erturb	ation I	Budget	$\epsilon = 1$	6/255		
Dataset	LVLM	Ca	aptioni	ng	Cla	ssificat	ion	'	VQA		Cap	otionii	ng	Cla	ssificat	ion	'	VQA	
		SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC
	LLaVA-1.5	0.841	82.9	82.9	0.983	97.6	97.6	0.791	78.0	78.0	0.568	51.2	51.2	0.959	95.1	95.1	0.518	48.8	51.2
DALL-E	MiniGPT-4 BLIP-2	0.907 0.891	87.8 82.9	90.2 82.9	0.837 0.888	78.0 85.4	85.4 87.8	0.770 0.718	73.2 65.9	75.6 65.9	0.792 0.770	68.3 68.3	78.0 68.3	0.769 0.748	68.3 63.4	75.6 63.4	0.661	63.4 39.0	63.4 39.0
	InstructBLIP	0.853	82.9	85.4	0.772	65.9	73.2	0.563	51.2	51.2	0.786	73.2	78.0	0.684	61.0	65.9	0.331	24.4	24.4
SVIT	LLaVA-1.5 MiniGPT-4 BLIP-2 InstructBLIP	0.904 0.892 0.962	87.8 87.8 97.6 90.2	90.2 87.8 97.6 92.7	1.000 0.914 0.920 0.807	100.0 90.2 90.2 68.3	100.0 90.2 92.7 80.5	0.767 0.849 0.728 0.483	75.6 80.5 70.7 41.5	75.6 82.9 70.7 43.9	0.699 0.859 0.865 0.874	65.9 85.4 80.5 80.5	68.3 85.4 82.9 87.8	1.000 0.801 0.779 0.745	100.0 78.0 70.7 65.9	100.0 78.0 73.2 73.2	0.543 0.790 0.610 0.435	51.2 75.6 56.1 31.7	53.7 80.5 56.1 31.7
VQAv2	LLaVA-1.5 MiniGPT-4 BLIP-2	0.910 0.934 0.903 0.966	90.2 92.7 87.8 97.6	92.7 92.7 87.8 97.6	1.000 0.878 0.910	100.0 85.4 85.4	100.0 87.8 87.8	0.485 0.906 0.814 0.584	90.2 80.5 53.7	43.9 90.2 80.5 53.7	0.874 0.819 0.792 0.861	80.5 80.5 75.6 82.9	80.5 75.6 82.9	0.743 0.955 0.841 0.894	95.1 82.9 82.9	95.1 82.9 85.4	0.765 0.686 0.546	75.6 65.9 46.3	75.6 65.9 46.3
	InstructBLIP	1.000	100.0	100.0	0.832	75.6	82.9	0.421	36.6	36.6	0.916	87.8	92.7	0.790	70.7	78.0	0.386	31.7	31.7

1154 Table 7: Targeted attack performance of our adversarial transformation-aware attack. Target text: 1155 "I don't know". Larger similarities (\uparrow) indicate more harmful impacts. "SS" denotes the semantic 1156 similarity between output and target text; "EM" and "CC" denote the "ExactMatch" and "Con-1157 ditionalContain" metrics to assess whether the output exactly matches or contain the target text, 1158 respectively.

				Perturb	oation I	Budget	$\epsilon = 3$	2/255				Р	erturt	oation I	Budget	$\epsilon = 1$	6/255		
Dataset	LVLM	Ca	aptioni	ng	Cla	ssificat	ion		VQA		Ca	otioni	ng	Cla	ssificat	ion	'	VQA	
		SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC	SS	EM	CC
	LLaVA-1.5	1.000	100.0	100.0	1.000	100.0	100.0	1.000	100.0	100.0	0.776	75.0	77.5	0.981	97.5	100.0	0.930	92.5	92.5
DALL-E	MiniGPT-4 BLIP-2	0.821	81.3 63.6	84.4 72.7	0.801	78.0 65 9	82.9 75.6	0.848	82.9 51.6	82.9 54 8	0.805	78.0 44 3	80.5 52.2	0.749	70.7 39.0	78.0 51.2	0.708	65.9 39.0	65.9 41 5
	InstructBLIP	0.797	77.1	77.1	0.810	75.9	75.9	0.740	67.9	67.9	0.611	51.4	54.3	0.601	51.7	51.7	0.608	52.2	52.2
	LLaVA-1.5	1.000	100.0	100.0	1.000	100.0	100.0	0.979	97.5	97.5	0.861	85.0	85.0	1.000	100.0	100.0	0.863	82.5	82.5
SVIT	MiniGPT-4 BLIP-2	0.839	82.9 72 7	82.9 75.8	0.937	92.7 82.9	92.7 85.4	0.781	75.6 56.5	80.5 56.5	0.788	75.6 48 5	78.0 54 3	0.865	85.4 56.1	87.8 61.0	0.750	73.2	73.2
	InstructBLIP	0.732	71.4	74.3	0.782	74.3	74.3	0.792	73.3	73.3	0.701	65.9	65.9	0.627	55.2	55.2	0.630	53.6	53.6
	LLaVA-1.5	1.000	100.0	100.0	1.000	100.0	100.0	1.000	100.0	100.0	0.930	92.5	92.5	1.000	100.0	100.0	0.980	97.5	97.5
VQAv2	MiniGPT-4	0.890	87.5	87.5	0.897	87.8	92.7 85.4	0.937	92.7	92.7	0.792	75.6	78.0	0.836	80.5	87.8	0.828	80.5	80.5
	InstructBLIP	0.792	68.8	64.8 78.1	0.800	80.5 76.7	83.4 76.7	0.695	61.5	61.5	0.575	58.5	65.9	0.688	61.0	65.9	0.508	43.3	43.3 47.4

1172 1173

1174 solely conduct targeted attacks, we implement our adversarial transformation-aware attack for com-1175 parison (we do not implement the hybrid transformation-based attack as it can only support untar-1176 geted attacks). As shown in Table 8 and Table 9, in a fair comparison setting, our attack method also 1177 achieves better performance than existing LVLM attacks MF and CroPA. This demonstrates that: (1) 1178 A simple and easy-to-implement transformation-aware attack is effective enough to fool the LVLM 1179 models. (2) Both MF and CroPA design complicated perturbation patterns. Compared to them, our transformation-aware attack is simple and easy-to-implement with better performance. Overall, we 1180 validate that adversarial visual transformation can achieve significant attack performance against 1181 LVLM models. 1182

1183 Besides, we also provide the complexity comparison with the two LVLM attacks: MF and CroPA. As shown in Table 10, our transformation attack is much more efficient than previous attackers as 1184 1185 they rely on more complicated adversarial pattern designs. Specifically, MF relies on an additional surrogate model CLIP to first initialize the noise and then design a perturbation update process to 1186 optimize the noise against the target LVLM, therefore introducing more model memory and time 1187 costs. CroPA requires optimizing both visual and textual noise with multi-prompt adversarial train-

1189Table 8: Performance comparison with the MF attack (Zhao et al., 2024) on the same ImageNet**1190**(Deng et al., 2009) dataset with the same semantic similarity metric (\uparrow). The values of the MF**1191**attack are reported in its paper.

1192	Attack	BLIP-2 (Li et al., 2023)	MiniGPT-4 (Zhu et al., 2023)	LLaVA-1.5 Liu et al. (2024a)
1193	Clean image (Zhao et al., 2024)	0.503	0.470	0.437
1194	MF-it (Zhao et al., 2024)	0.546	0.484	0.452
1195	MF-ii (Zhao et al., 2024)	0.592	0.572	0.450
1196	MF-ii+tt (Zhao et al., 2024)	0.665	0.666	0.597
1197	Clean image	0.569	0.427	0.369
1198	Ours-Adversarial		0.878	0.806

Table 9: Performance comparison with the CroPA attack (Luo et al., 2024) on the same MS-COCO (Lin et al., 2014) dataset and OpenFlamingo (Awadalla et al., 2023) model with the same attack success rate metric ([†]). The values of the CroPA attack are reported in its paper.

Attack	$\big VQA_{general}$	$VQA_{\it specific}$	Classification	Captioning	Overall
Single-P (Luo et al., 2024)	0.21	0.43	0.47	0.34	0.36
Multi-P (Luo et al., 2024)	0.60	0.85	0.71	0.60	0.69
CroPA (Luo et al., 2024)	0.90	0.96	0.75	0.72	0.83
Ours-Adversarial	1.00	1.00	1.00	1.00	1.00

Table 10: Complexity comparison with MF and CroPA attacks.

Process	Attack Type	GPU Time (\downarrow)	GPU Memory (\downarrow)
Generation	MF (Zhao et al., 2024)	29min	35GB
	CroPA (Luo et al., 2024)	14min	26GB
	Ours-Adversarial	5min	22GB

ing, also resulting in relatively more time costs. Therefore, it validates that our simple yet efficient adversarial visual transformation is effective enough to fool the LVLM models.