# GRADIENT REGULARIZATION-BASED CROSS-PROMPT ATTACKS ON VISION LANGUAGE MODELS

**Anonymous authors**
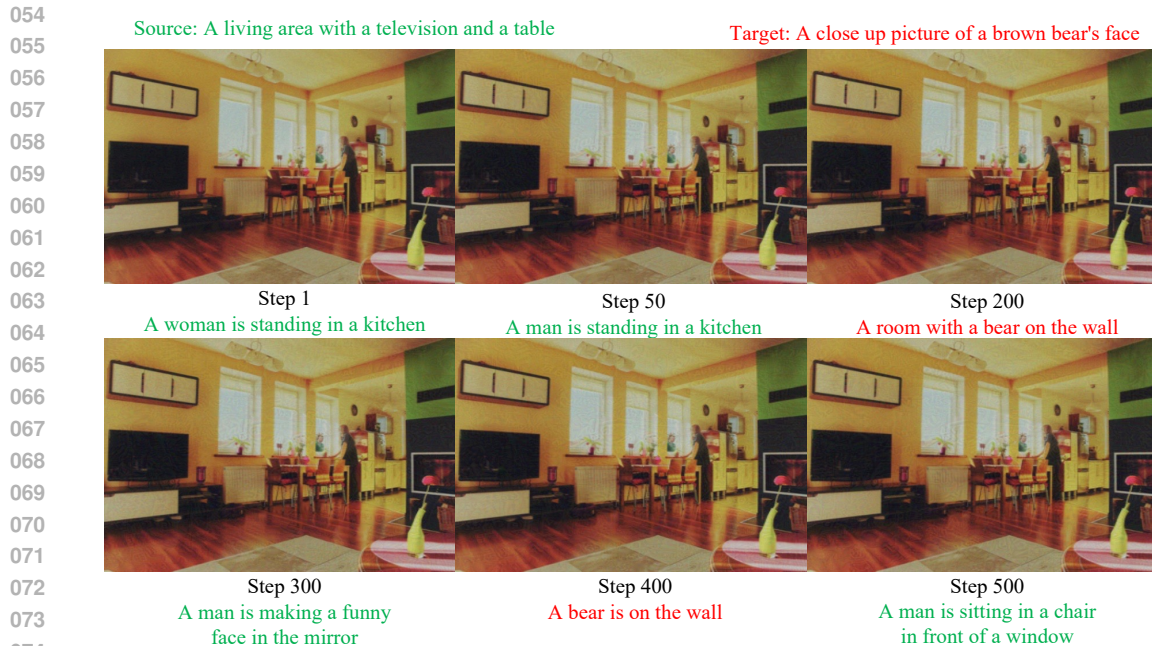Paper under double-blind review

## ABSTRACT

Recent large vision language models (VLMs) have gained significant attention for their superior performance in various visual understanding tasks using textual instructions, also known as prompts. However, existing research shows that VLMs are vulnerable to adversarial examples, where imperceptible perturbations added to images can lead to malicious outputs, posing security risks during deployment. Unlike single-modal models, VLMs process both images and text simultaneously, making the creation of visual adversarial examples dependent on specific prompts. Consequently, the same adversarial example may become ineffective when different prompts are used, which is common as users often input diverse prompts. Our experiments reveal severe non-stationarity when directly optimizing adversarial example generation using multiple prompts, resulting in examples specific to a single prompt with poor transferability. To address this issue, we propose the Gradient Regularized-based Cross-Prompt Attack (GrCPA), which leverages gradient regularization to generate more robust adversarial attacks, thereby improving the assessment of model robustness. By exploiting the structural characteristics of the Transformer, GrCPA reduces the variance of back-propagated gradients in the Attention and MLP components, utilizing regularized gradients to produce more effective adversarial examples. Extensive experiments on models such as Flamingo, BLIP-2, LLaVA and InstructBLIP demonstrate the effectiveness of GrCPA in enhancing the transferability of adversarial attacks across different prompts.

## 1 INTRODUCTION

Large Vision Language Models (VLMs), such as GPT-4 (OpenAI, 2023b), have recently garnered substantial interest from the AI research community. Unlike Large Language Models (LLMs), which are limited to processing plain text (OpenAI, 2023a), VLMs can interpret image inputs and perform a range of visual understanding tasks guided by textual instructions, or prompts. These tasks include image captioning (Li et al., 2023a; Zhang et al., 2020; Sheng et al., 2021), information extraction (Liu et al., 2024; Li et al., 2023b), complex counting (Bavishi et al., 2023), and visual grounding (Wang et al., 2023a; Bai et al., 2023), among others. This powerful multimodal perception capability has facilitated the deployment of more models in real-world production environments.

Recent studies have revealed that VLMs are susceptible to attacks from adversarial examples (Gu et al., 2023; Madry et al., 2018; Szegedy et al., 2014; Li et al., 2024a; Mahmood et al., 2021; Mao et al., 2023; Yu et al., 2023; Wang et al., 2023b; Shayegani et al., 2023). These attacks involve the addition of imperceptible disturbances to clean images, which can induce the models to output malicious content. Such adversarial attacks can circumvent the security constraints of LLMs or even embed advertising information into images (Niu et al., 2024; Qi et al., 2023; Bailey et al., 2023; Lu et al., 2024; Yuan et al., 2023). Therefore, designing effective attack methods to identify potential vulnerabilities before deploying VLMs in security-related applications is of paramount importance (Li et al., 2024b; Gao et al., 2024b; Wang et al., 2023c).

Adversarial attacks can be broadly categorized into white-box attacks and black-box attacks (Gao et al., 2024a; Cheng et al., 2019). A white-box attack refers to an attacker who has access to all the structural and weight information of the model (Ebrahimi et al., 2018). Conversely, a black-box attack refers to an attacker who can only access the model's external usage interface (Guo et al.,

Figure 1: Illustration of the attack iteration process. Green represents the clean image description, while red represents the attack target. Adversarial attacks are unstable during the iteration process, causing fluctuations around the attack target.

2019). Given their expanded operational scope and potential to transition into black-box attacks, white-box attacks are consequently receiving heightened attention (Weng et al., 2024).

In contrast to widely studied classification models, adversarial attacks on VLMs present a more complex challenge. These attacks can stem from two perspectives: visual input and textual input. However, visual attacks are often imperceptible to users and occupy a continuum of disturbance space, making them more commonly utilized. The generation of visual adversarial examples for VLMs requires coordination with specific prompts. In other words, the same visual adversarial sample may not be effective when encountering diverse prompts (Cui et al., 2023). This phenomenon is prevalent during the model deployment phase, as users tend to input prompts based on their individual language preferences. Therefore, this paper focuses on cross-prompt visual adversarial attacks.

An intuitive method to enhance across-prompts transferability is to utilize multiple prompts in the iterative generation of adversarial examples (Moosavi-Dezfooli et al., 2017). However, in our experiments, we identified three issues with this approach: (a) A serious non-stationary phenomenon is observed, characterized by large fluctuations in the success rate of the attack during the iteration of adversarial samples, as shown in Figure 1. We attribute this to overfitting during optimization, since adversarial attacks on VLMs usually require a large number of iterations, such as 10,000, to succeed, causing adversarial examples to become specific to their conditions (model and prompt) (Schlarmann & Hein, 2023). (b) The calculation of text loss is extremely sensitive. Initially, we inadvertently computed the loss for the entire sequence using teacher forcing, but found the results largely unsuccessful. Subsequently, we recognized the need to focus solely on the loss pertaining to the model's output section. (c) Methods based on image classification for enhancing transferability are not adaptable to VLMs. We tested methods like MI-FGSM (Dong et al., 2018), Input Diversity (Xie et al., 2019), Variance Tuning (Wang & He, 2021), and found that the transferability across prompts did not increase, but even decreased. A more detailed analysis can be seen in the Appendix A.1.

Based on these observations, we contend that the design of visual adversarial examples for VLMs should take into account both image and text inputs comprehensively, with a particular focus on mitigating overfitting in the textual domain. VLMs often integrate substantial language models, which consist of numerous Transformer blocks, potentially leading to learned features that are specific to the prompts or the model itself (Wang & He, 2021). In this paper, we introduce a **G**radient **R**egularization-based **C**ross-**P**rompt Attack (**GrCPA**) method designed to alleviate overfitting of both

visual and textual features within the LLM' Transformer blocks, thereby enhancing the transferability of visual adversarial examples. Specifically, we implement gradient clipping on both visual and textual features during the loss back-propagation phase to counter overfitting. Note that we modify only a very small number of gradients, which does not affect the overall convergence of the chain rule (Zhang et al., 2023a; Wei et al., 2022).

To verify the effectiveness of GrCPA, we employ prompts from three distinct vision-language tasks: image classification, image captioning, and visual question answering (VQA). We evaluate our method's efficacy on well-known VLMs, including Flamingo (Alayrac et al., 2022), BLIP-2 (Li et al., 2023a), LLaVA-1.5 (Liu et al., 2023) and InstructBLIP (Dai et al., 2023). The experimental results indicate that GrCPA exhibits superior attack performance and enhanced transferability.

Overall, the main contributions of this paper include:

1. To the best of our knowledge, we first identify the non-stationary phenomenon in adversarial attacks against vision language models, and argue that its essence is overfitting in the optimization process. We also attempt previous enhancement methods for single-modal models and find them to be ineffective.

2. We propose a gradient regularization method to enhance the transferability of visual adversarial examples, thus effectively alleviating overfitting issues in the deep Transformer blocks of visual and textual features.

3. We validate the effectiveness of our method through detailed experiments and provide a new perspective for future attacks against VLMs.

## 2 RELATED WORK

**Adversarial Transferability.** Szegedy et al. (2014) first proposed the concept of adversarial examples, revealing the vulnerability of neural networks. The transferable attacks, which have widespread impacts in the real world, have triggered a large number of subsequent studies Cheng et al. (2020); Wu et al. (2022); Zhang et al. (2023b); Chakraborty et al. (2021); Madry et al. (2018); Xu et al. (2022). Previous work has primarily focused on classification models, with an emphasis on enhancing transferability through gradient optimization, input augmentation. Gradient optimization methods, led by the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015), along with its derivatives such as Iterative FGSM (I-FGSM) (Kurakin et al., 2017), Projected Gradient Descent (PGD) (Madry et al., 2018), Momentum Iterative FGSM (MI-FGSM) (Dong et al., 2018), among others, have emerged as prominent techniques in the literature. On another front, input augmentation primarily involves applying various transformations to the input image at each iteration, such as random resizing and padding, as seen in methods like DIM Xie et al. (2019), SIM Lin et al. (2020), and TIM Dong et al. (2019). We endeavor to improve the transferability of attacks on VLMs utilizing traditional methods, but observe no substantial enhancement in their effectiveness. This highlights the complexity and inherent challenges of multimodal tasks, prompting a reevaluation of our previous research methodologies.

**Adversarial Robustness of Vision Language Models.** Alongside the proliferation of large VLMs, the associated security research has garnered significant attention Gao et al. (2024a); Sun et al. (2024); Ni et al. (2024); Zhang et al. (2024); Guo et al. (2024); Luo et al. (2024b); Zhou et al. (2024); Cheng et al. (2024); Wang et al. (2024). For example, Zhao et al. (2023) induce misinterpretation of image content in models such as BLIP-2 through black-box attacks. There is also a body of work utilizing adversarial attacks to circumvent security alignment of LLM components (Bagdasaryan et al., 2023; Carlini et al., 2023; Niu et al., 2024; Qi et al., 2024), posing security risks to VLMs. The work closest to ours is CroPA Luo et al. (2024a), which turns the generation of visual adversarial examples into a max-min process, achieving significant improvements. Our method from the perspective of reverse gradient is orthogonal to it, with more flexible and simpler operational methods.

## 3 METHOD

In this section, we first introduce adversarial attack setup against VLMs. Then, we formally present baseline methods for generating visual adversarial examples using a single prompt and multiple

(a) Illustration of gradient regularization.
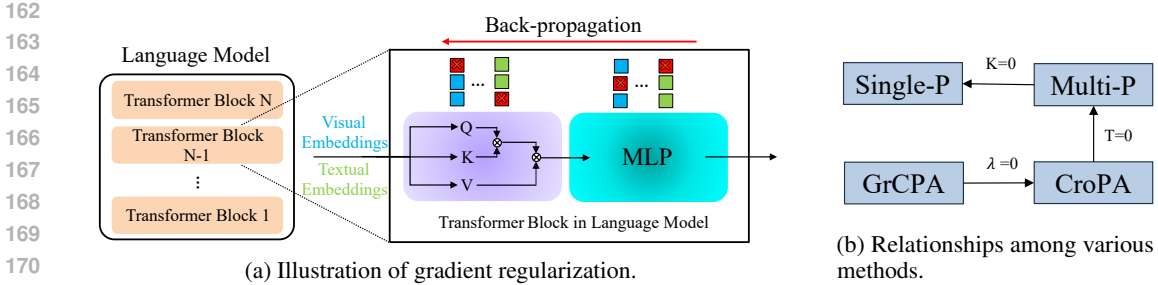
(b) Relationships among various methods.

Figure 2: Overview of our proposed method. (a) The process of gradient regularization involves performing clipping during back-propagation in the Transformer blocks of the LLM for both visual embeddings and text embeddings, specifically by attenuating the extreme values of gradients in each token. (b) By adjusting different parameter values, various methods can be transformed.

prompts. Finally, we introduce our proposed GrCPA method, which enhances cross-prompt transferability through gradient regularization during backpropagation.

## 3.1 THREAT MODEL

Without loss of generality, VLMs complete a series of downstream tasks through visual question answering (VQA). The model is provided with an image $v$ and question $q$ (i.e., prompts), and in response, it generates an answer $a$. We denote a VLM with the function $f_\theta$, where $\theta$ represents its parameters, such that $a = f_\theta(v, q)$.

Thus, the threat model for adversarial attacks on vision language models can be represented as:

**Adversarial knowledge** refers to the information an attacker has about the target model. In this paper, we focus on white-box attacks, where we have full access to all details of the victim's model, including its architecture and weights. This access allows us to leverage the gradients obtained through backpropagation to generate adversarial examples effectively.

**Adversarial goals** describe the malicious objectives that an attacker aims to achieve, typically categorized into targeted and untargeted attacks. For VLMs, targeted attacks seek to induce the model to output specific content, including bypassing alignment constraints. In contrast, untargeted attacks aim to provoke the model into producing incongruous responses. Since untargeted attacks can often be achieved through targeted attacks, this paper places greater emphasis on targeted attacks.

**Adversary capabilities** refer to the resources and constraints available to the attacker. To ensure that adversarial examples remain imperceptible to humans, the image perturbation $\delta$ is constrained by $||\delta||_p \leq \epsilon$, where $\epsilon$ represents the magnitude of the perturbation and $||\cdot||_p$ denotes the $L_p$ norm. This paper primarily employs the $L_\infty$ norm to align with previous work (Luo et al., 2024a).

## 3.2 BASELINE METHODS

To induce the model to output specific content, given a query $q$ and a target answer $a$, we optimize the loss of the language model with respect to the $(q, a)$ pair and backpropagate it to the image, thus generating adversarial examples. We denote this method as **Single-P**. However, the activation of adversarial examples generated by this method also depends on the optimization of $q$ used during the process. In other words, if replaced with another prompt $q'$, the model may fail to produce $a$ in response (Cui et al., 2023).

To enhance the cross-prompt transferability of visual adversarial samples, a straightforward approach is to use multiple prompts during optimization (Moosavi-Dezfooli et al., 2017). Given a set of textual prompts $\mathcal{Q} = \{q_1, q_2, q_3, \ldots, q_M\}$, we induce the model to output the predetermined target answer $a$ under the query of each item in $\mathcal{Q}$ in the presence of adversarial noise $\delta$. Specifically, we minimize the following language modeling loss:

$$\min_{\delta_v} \sum_{i=1}^{K} \mathcal{L}(f(v + \delta, q_i), a) \tag{1}$$

Where $\mathcal{L}$ is typically the cross-entropy loss. Note that we compute the loss for only the part corresponding to answer $a$ rather than the entire $(q, a)$ sequence. We refer to this method as **Multi-P**.

It is evident that the improvement in cross-prompt transferability is directly proportional to the increase in the number of textual prompts, denoted by $K$. However, exhaustively exploring all potential prompts is often impractical due to the significantly increased computational complexity. Therefore, it is essential to enhance transferability with a limited number of prompts. To address this, **CroPA** Luo et al. (2024a) proposes using a set of learnable prompts to update the visual adversarial perturbation, aiming to counteract the misleading effects of adversarial images:

$$\min_{\delta_v} \max_{\delta_t} \mathcal{L}(f(v + \delta_v, q_i + \delta_t), a) \tag{2}$$

where $\delta_v$ represents the perturbation added to the image, while $\delta_t$ represents the perturbation added to the text. To smooth the optimization process, a text perturbation update frequency $T$ is introduced. This means that for every $T$ updates of the visual perturbation, the text perturbation is updated once.

### 3.3 GRCPA

Orthogonal to CroPA, we introduce GrCPA, which focuses on gradient regularization during the backpropagation process. Visual adversarial attacks fundamentally involve optimizing images using gradient descent. Consequently, large gradients can lead to local optima and trigger overfitting issues (Wang & He, 2021). This motivates us to clip the gradients of both visual and text features, thereby enhancing cross-prompt transferability.

Existing large VLMs typically consist of three components: a visual encoder, a projection layer, and an LLM. The image passes through the visual encoder to obtain a set of features, which are then aligned to the input space of the LLM by the projection layer to form visual tokens. These visual tokens are concatenated with textual tokens and fed into the LLM for autoregressive generation. The LLM is composed of multiple stacked Transformer blocks, with each block consisting of Attention and MLP components (Vaswani et al., 2017).

Given the gradient vector $G \in \mathcal{R}^d$ with respect to visual or textual tokens, where $d$ is the embedding length, we compute the language modeling loss (Equation 2) and propagate this loss backward through the Transformer blocks of the LLM. As shown in Figure 2a, we perform **Gradient Regularization** (GR) by identifying the $k$ largest and smallest gradient (Grad) values and setting them to 0, as follows.

$$i_{\max}, i_{\min} = \underset{k}{\arg \max} \, G, \underset{k}{\arg \min} \, G$$
$$G[i_{\max}] = G[i_{\min}] = 0 \tag{3}$$

This clipping will be performed on each token in both the Attention block and the MLP block of the Transformer blocks.

**Preserving low-level features.** LLMs typically consist of multiple stacked Transformer blocks, which enable the learned features to be specific to the model itself. Inspired by Deng et al. (2023), which suggests that preserving more low-level features in convolutional networks can improve cross-model transferability, we hypothesize that regularizing only the gradients of certain high-level features in LLMs can better balance the strength and transferability of attacks. Assuming that the LLM consists of $N$ Transformer blocks, we set $n = \lambda \cdot N$ and apply regularization only to the features of the last $n$ layers. The complete algorithm is presented in Algorithm 1.

**Relationships among various methods.** Our method builds upon previous work Moosavi-Dezfooli et al. (2017); Luo et al. (2024a); Zhang et al. (2023a); Wei et al. (2022), where correlations between them can be transformed by adjusting the values of hyperparameters, as illustrated in Figure 2b.

## 4 EXPERIMENTS

In this section, we validate the effectiveness of our method through extensive experiments and conduct an in-depth analysis of the factors influencing GrCPA.

Table 1: The results of targeted attacks against Flamingo. We investigate four types of tasks and report attack success rates for each task category. The last column represents the average performance across these four task types. To demonstrate the generalization of our method, we set multiple target answers. The best results are represented in **bold**.

| Target Answer | Method | VQA$_{general}$ | VQA$_{specific}$ | Classification | Captioning | Average |
|---|---|---|---|---|---|---|
| unknown | Single-P | 0.24 | 0.39 | 0.21 | 0.05 | 0.22 |
| | Multi-P | 0.67 | 0.86 | 0.64 | 0.31 | 0.62 |
| | CroPA | 0.92 | 0.98 | 0.70 | 0.34 | 0.74 |
| | GrCPA | **0.95** | **0.99** | **0.75** | **0.43** | **0.78** |
| I am sorry | Single-P | 0.21 | 0.43 | 0.47 | 0.34 | 0.36 |
| | Multi-P | 0.60 | 0.85 | 0.71 | 0.60 | 0.69 |
| | CroPA | 0.90 | 0.96 | 0.75 | 0.72 | 0.83 |
| | GrCPA | **0.94** | **0.96** | **0.83** | **0.82** | **0.88** |
| not sure | Single-P | 0.25 | 0.36 | 0.09 | 0.00 | 0.17 |
| | Multi-P | 0.55 | 0.55 | 0.11 | 0.02 | 0.31 |
| | CroPA | 0.88 | 0.86 | 0.30 | 0.17 | 0.55 |
| | GrCPA | **0.93** | **0.89** | **0.46** | **0.26** | **0.63** |
| very good | Single-P | 0.35 | 0.52 | 0.15 | 0.05 | 0.27 |
| | Multi-P | 0.81 | 0.93 | 0.40 | 0.20 | 0.59 |
| | CroPA | 0.95 | 0.97 | 0.64 | 0.27 | 0.71 |
| | GrCPA | **0.99** | **0.97** | **0.81** | **0.44** | **0.80** |
| too late | Single-P | 0.21 | 0.38 | 0.21 | 0.04 | 0.21 |
| | Multi-P | 0.78 | 0.90 | 0.54 | 0.17 | 0.60 |
| | CroPA | 0.90 | 0.95 | 0.73 | 0.20 | 0.70 |
| | GrCPA | **0.93** | **0.97** | **0.79** | **0.39** | **0.77** |
| metaphor | Single-P | 0.26 | 0.56 | 0.50 | 0.14 | 0.37 |
| | Multi-P | 0.83 | 0.92 | 0.81 | 0.42 | 0.75 |
| | CroPA | 0.96 | 0.99 | 0.92 | 0.62 | 0.87 |
| | GrCPA | **0.99** | **0.99** | **0.95** | **0.73** | **0.91** |

## 4.1 EXPERIMENTAL SETTINGS

**Datasets.** Given that VLMs tackle downstream tasks through visual question answering, it is imperative that the dataset encompasses both images and corresponding prompts. The images are sourced from the MS-COCO validation set (Lin et al., 2014). The VQA prompts are comprised of questions that are either general or specific to the image content, respectively referred to as VQA$_{general}$ and VQA$_{specific}$. The image-specific questions are derived from the VQA-v2 dataset (Goyal et al., 2017). Agnostic questions were constructed for VQA, with a focus on image classification and image captioning, ensuring a diverse range of lengths and semantic content.

**Models.** Without loss of generality, we evaluate the OpenFlamingo-9B (Alayrac et al., 2022; Awadalla et al., 2023), BLIP-2 (OPT-2.7B) (Li et al., 2023a; Zhang et al., 2022), LLaVA-1.5-7B(Liu et al., 2023) and InstructBLIP (Dai et al., 2023), which are influential models in the multimodal community.

**Parameters.** In alignment with previous research (Luo et al., 2024a), the image perturbation is configured to $16/255$, with $\alpha_1 = 1/255$, $\alpha_2 = 0.01$, and the number of iterations set to 1000. A maximum of 100 prompts are utilized for each individual sample. The proportion of Transformer blocks $\lambda$ is set to $1/4$; the update frequency $T$ is set to 1; and the number of extrema $k$ is also set to 1.

**Evaluation Metric.** In this paper, we report the Attack Success Rate (ASR) and facilitate the analysis by inducing the model to output specific text.

Table 2: Quantitative evaluation of attack stability. We assess the stability of different methods by determining whether the outputs at the 900th, 925th, 950th, 975th, and 1000th steps are consistent.

| Method | $VQA_{general}$ | $VQA_{specific}$ | Classification | Captioning | Average |
|--------|--------|--------|------|------|------|
| Multi-P | 0.57 | 0.59 | 0.46 | 0.43 | 0.51 |
| CroPA | 0.61 | 0.65 | 0.53 | 0.49 | 0.57 |
| GrCPA | **0.67** | **0.71** | **0.57** | **0.53** | **0.62** |

## 4.2 COMPARISON WITH PREVIOUS METHODS

To comprehensively demonstrate the efficacy of our proposed GrCPA, we conducted a series of experiments using Flamingo (Awadalla et al., 2023), evaluating it against a range of target responses. The target text consists of statements such as `unknown`, `not sure`, and `I am sorry`, which indicate a deficiency in interpreting visual content, and `unknown` is the default setting for subsequent experiments unless otherwise specified. It also features phrases like `very good`, `too late`, and `metaphor`, which are irrelevant to the context.

Table 1 summarizes the evaluation results of targeted attacks. GrCPA outperforms previous SOTA methods across various experimental settings. Both the baseline methods and our method achieve higher attack success rates on the VQA task, likely due to the closer relationship between prompts and images in the VQA framework, where prompts more closely related to the image are more likely to enhance the effectiveness of the attack. This also indirectly demonstrates the sensitivity of adversarial samples to prompts in vision-language models, where adversarial samples may become ineffective when encountering different prompts. Furthermore, varying target answers can affect attack results. The experimental findings suggest that even rare and illogical responses, like metaphors, can still achieve high success rates. We also evaluate longer target answers, such as `I need a new phone`, in Appendix A.3. The results show that our method still outperforms the baseline methods. Additionally, we demonstrate the stability of our method through qualitative case studies in Appendix A.4.

To further validate the generalizability of our method, we also conduct experiments on LLaVA-1.5 and InstructBLIP, as detailed in Appendix A.5. These models are similarly susceptible to adversarial attacks, exhibiting serious security vulnerabilities. We also evaluate their cross-model transferability in Appendix A.7, but find weak transferability..

Besides showcasing the attack results, we also perform a quantitative analysis of the variations in attack stability across different methods. We further evaluate the stability of various attack methods by examining whether the model's outputs at the 900th, 925th, 950th, 975th, and 1000th iterations are consistent. As shown in Table 2, our method significantly enhances the stability of adversarial attacks across multiple tasks, which greatly aids in the large-scale evaluation of VLMs' robustness.

## 4.3 IMPACT OF PROMPT NUMBER

In this section, our focus is on examining the influence of the quantity of prompts utilized in the attack process on its effectiveness. We conduct attacks under various configurations, employing 1, 5, 10, 50, and 100 prompts against the BLIP-2 model with the objective of eliciting an `unknown` response.

As shown in Table 3, augmenting the quantity of prompts in the optimization phase substantially augments the cross-prompt transferability of visual adversarial samples. For example, escalating the number of prompts from one to ten results in a pronounced increment in the ASR of the baseline method, from 0.34 to 0.71, which corresponds to a more than twofold enhancement. The experimental outcomes clearly indicate that our methodology consistently outperforms the baseline approach across all configurations, thereby showcasing our method's superiority.

Nevertheless, augmenting the number of prompts directly leads to a substantial increase in the computational demands of adversarial attacks, presenting a significant impediment to the large-scale generation of visual adversarial samples. Additionally, there is a pronounced effect of diminishing marginal returns associated with increasing the number of prompts; beyond a threshold of 10 prompts,

Table 3: The results of the adversarial attack against BLIP-2. Different numbers of prompts are employed, and it is found that increasing the number of prompts improves the performance. The best performance values for each task are highlighted in **bold**.

| No. of Prompts | Method | $VQA_{general}$ | $VQA_{specific}$ | Classification | Captioning | Average |
|---|---|---|---|---|---|---|
| 1 | Single-P | 0.24 | 0.34 | 0.45 | 0.32 | 0.34 |
| | CroPA | 0.52 | 0.63 | 0.65 | 0.58 | 0.60 |
| | **GrCPA** | **0.55** | **0.65** | **0.69** | **0.60** | **0.62** |
| 5 | Multi-P | 0.51 | 0.59 | 0.62 | 0.58 | 0.58 |
| | CroPA | 0.81 | 0.83 | 0.80 | 0.84 | 0.82 |
| | **GrCPA** | **0.85** | **0.87** | **0.83** | **0.89** | **0.86** |
| 10 | Multi-P | 0.68 | 0.81 | 0.68 | 0.67 | 0.71 |
| | CroPA | 0.86 | 0.90 | 0.82 | 0.84 | 0.86 |
| | GrCPA | **0.88** | **0.93** | **0.84** | **0.85** | **0.87** |
| 50 | Multi-P | 0.67 | 0.74 | 0.67 | 0.72 | 0.70 |
| | CroPA | 0.90 | 0.93 | 0.87 | 0.91 | 0.90 |
| | **GrCPA** | **0.95** | **0.96** | **0.89** | **0.92** | **0.93** |
| 100 | Multi-P | 0.67 | 0.76 | 0.68 | 0.66 | 0.69 |
| | CroPA | 0.95 | 0.95 | 0.87 | 0.92 | 0.92 |
| | **GrCPA** | **0.99** | **0.99** | **0.93** | **0.95** | **0.96** |



Figure 3: The impact of the number of iterations on attack success rate. Compared with the baseline algorithms, our method shows a certain degree of improvement across different numbers of iterations.

the enhancement in transferability becomes exceedingly constrained. Consequently, it is imperative to improve cross-prompt transferability using a finite set of prompts.

### 4.4 CONVERGENCE OF GRCPA

In this section, we explore the impact of the number of iterations during the optimization process on the attack success rate. All attacks are conducted using 10 prompts.

As shown in Figure 3, it can be observed that all methods show a corresponding increase in attack success rate with the number of iterations. This improvement is particularly evident in scenarios using multiple prompts, as more prompts necessitate learning more content.

Our method's performance gradually stabilizes after 1000 iterations. Compared to adversarial attacks on classification tasks, which typically require around 100 iterations, adversarial attacks on VLMs demand significantly more computational effort. However, our method can achieve better performance with fewer steps and demonstrates higher computational efficiency.

### 4.5 ABLATION STUDIES

In this section, we thoroughly analyze the effectiveness of GrCPA through ablation experiments.

Table 4: Ablation studies of gradient regularization.

| Method | $VQA_{general}$ | $VQA_{specific}$ | Classification | Captioning | Average |
|---|---|---|---|---|---|
| Single-P | 0.24 | 0.39 | 0.21 | 0.05 | 0.22 |
| Single-P(GR) | 0.29 | 0.45 | 0.24 | 0.11 | 0.27 |
| Multi-P | 0.67 | 0.86 | 0.64 | 0.31 | 0.62 |
| Multi-P(GR) | 0.77 | 0.91 | 0.71 | 0.35 | 0.78 |

Table 5: Ablation studies on single-modality regularization.

| Method | $VQA_{general}$ | $VQA_{specific}$ | Classification | Captioning | Average |
|---|---|---|---|---|---|
| GrCPA | 0.99 | 0.99 | 0.93 | 0.95 | 0.96 |
| GrCPA(Image) | 0.94 | 0.95 | 0.90 | 0.89 | 0.92 |
| GrCPA(text) | 0.92 | 0.93 | 0.86 | 0.87 | 0.89 |

Table 6: Ablation experiments on the impact of the layer proportion $\lambda$.

| $\lambda$ | 1 | 1/2 | 1/3 | 1/4 | 1/5 | 1/6 |
|---|---|---|---|---|---|---|
| ASR | 0.863 | 0.863 | 0.864 | 0.875 | 0.873 | 0.869 |

**The impact of gradient regularization.** Although GrCPAbuilds on previous work, this gradient regularization method is actually a general approach to reducing overfitting. As shown in Table 4, our experiments on Single-P and Multi-P demonstrate that it can provide cross-prompt transferability.

**The impact of regularizing different modalities.** In our method, we apply gradient regularization to both visual modality features and textual modality features in the LLM. In practice, it is feasible to regularize the feature gradients of a single modality. We conduct such experiments as shown in Table 5, but find that the attack success rate significantly decreased. Therefore, we believe that enhancing attacks on VLMs should consider both modalities whenever possible.

**The impact of proportion $\lambda$ of Transformer blocks.** We primarily test the effectiveness of gradient regularization on Transformer Blocks with LLMs at different proportions $\lambda$. The experimental results, as shown in Table 6, revealed that trimming only the last $1/4$ layers achieved the best performance. However, in terms of absolute performance, the differences among them were relatively minor.

## 5 CONCLUSION

In this paper, we investigate the adversarial robustness of large vision language models (VLMs). During our experiments, we first found that existing adversarial attacks on visual language models exhibit significant instability, with the optimization process for adversarial samples oscillating between success and failure. We believe that the root cause of this issue is overfitting during the optimization process, which poses a challenge to the large-scale generation of adversarial samples for visual language models. Furthermore, we experimentally investigated the effectiveness of adversarial attack enhancement methods that target only the visual modality within VLMs and found that these methods reduce attack performance. Based on these observations, we propose Gradient Regularized-based Cross-Prompt Attack (GrCPA), which clips the gradients of visual and textual features during error backpropagation, eliminating extreme gradients to prevent falling into local optima. Our regularization operation modifies only a small portion of the gradients and does not affect the convergence of the chain rule. Experiments on models such as BLIP-2 demonstrate that our method significantly improves the transferability of adversarial samples and confirms that current VLMs are sensitive to visual inputs and can be easily attacked. Therefore, we call on researchers to thoroughly evaluate the adversarial robustness of visual language models before deployment, especially in life-critical scenarios.

**Reproducibility.** In the experiments, we thoroughly report on the datasets, models, and parameter settings designed for this study, with all data being open-source and publicly available to ensure reproducibility.

# REFERENCES

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022. URL `http://papers.nips.cc/paper_files/paper/2022/hash/960a172bc7fbf0177ccccbb411a7d800-Abstract-Conference.html`.

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Yitzhak Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models. *CoRR*, abs/2308.01390, 2023. doi: 10.48550/ARXIV.2308.01390. URL `https://doi.org/10.48550/arXiv.2308.01390`.

Eugene Bagdasaryan, Tsung-Yin Hsieh, Ben Nassi, and Vitaly Shmatikov. (ab)using images and sounds for indirect instruction injection in multi-modal llms. *CoRR*, abs/2307.10490, 2023. doi: 10.48550/arXiv.2307.10490. URL `https://doi.org/10.48550/arXiv.2307.10490`.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023.

Luke Bailey, Euan Ong, Stuart Russell, and Scott Emmons. Image hijacks: Adversarial images can control generative models at runtime. *CoRR*, abs/2309.00236, 2023. doi: 10.48550/ARXIV.2309.00236. URL `https://doi.org/10.48550/arXiv.2309.00236`.

Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. Introducing our multimodal models, 2023. URL `https://www.adept.ai/blog/fuyu-8b`.

Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramèr, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? *CoRR*, abs/2306.15447, 2023. doi: 10.48550/arXiv.2306.15447. URL `https://doi.org/10.48550/arXiv.2306.15447`.

Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. A survey on adversarial attacks and defences. *CAAI Trans. Intell. Technol.*, 6(1):25–45, 2021. doi: 10.1049/CIT2.12028. URL `https://doi.org/10.1049/cit2.12028`.

Mingyue Cheng, Yiheng Chen, Qi Liu, Zhiding Liu, and Yucong Luo. Advancing time series classification with multimodal language modeling. *CoRR*, abs/2403.12371, 2024. doi: 10.48550/ARXIV.2403.12371. URL `https://doi.org/10.48550/arXiv.2403.12371`.

Minhao Cheng, Simranjit Singh, Patrick H. Chen, Pin-Yu Chen, Sijia Liu, and Cho-Jui Hsieh. Signopt: A query-efficient hard-label adversarial attack. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL `https://openreview.net/forum?id=SklTQCNtvS`.

Shuyu Cheng, Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Improving black-box adversarial attacks with a transfer-based prior. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 10932–10942, 2019. URL `https://proceedings.neurips.cc/paper/2019/hash/32508f53f24c46f685870a075eaaa29c-Abstract.html`.

Xuanming Cui, Alejandro Aparcedo, Young Kyun Jang, and Ser-Nam Lim. On the robustness of large multimodal models against image adversarial attacks. *CoRR*, abs/2312.03777, 2023. doi: 10.48550/ARXIV.2312.03777. URL `https://doi.org/10.48550/arXiv.2312.03777`.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *CoRR*, abs/2305.06500, 2023. doi: 10.48550/arXiv.2305.06500. URL https://doi.org/10.48550/arXiv.2305.06500.

Yang Deng, Weibin Wu, Jianping Zhang, and Zibin Zheng. Blurred-dilated method for adversarial attacks. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/b6fa3ed9624c184bd73e435123bd576a-Abstract-Conference.html.

Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 9185–9193. Computer Vision Foundation / IEEE Computer Society, 2018. doi: 10.1109/CVPR.2018.00957. URL http://openaccess.thecvf.com/content_cvpr_2018/html/Dong_Boosting_Adversarial_Attacks_CVPR_2018_paper.html.

Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 4312–4321. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00444. URL http://openaccess.thecvf.com/content_CVPR_2019/html/Dong_Evading_Defenses_to_Transferable_Adversarial_Examples_by_Translation-Invariant_Attacks_CVPR_2019_paper.html.

Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers*, pp. 31–36. Association for Computational Linguistics, 2018. doi: 10.18653/v1/P18-2006. URL https://aclanthology.org/P18-2006/.

Kuofeng Gao, Yang Bai, Jiawang Bai, Yong Yang, and Shu-Tao Xia. Adversarial robustness for visual grounding of multimodal large language models, 2024a.

Kuofeng Gao, Yang Bai, Jindong Gu, Shu-Tao Xia, Philip H. S. Torr, Zhifeng Li, and Wei Liu. Inducing high energy-latency of large vision-language models with verbose images. *CoRR*, abs/2401.11170, 2024b. doi: 10.48550/ARXIV.2401.11170. URL https://doi.org/10.48550/arXiv.2401.11170.

Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL http://arxiv.org/abs/1412.6572.

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 *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 6325–6334. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.670. URL https://doi.org/10.1109/CVPR.2017.670.

Jindong Gu, Xiaojun Jia, Pau de Jorge, Wenqian Yu, Xinwei Liu, Avery Ma, Yuan Xun, Anjun Hu, Ashkan Khakzar, Zhijiang Li, Xiaochun Cao, and Philip H. S. Torr. A survey on transferability of adversarial examples across deep neural networks. *CoRR*, abs/2310.17626, 2023. doi: 10.48550/ARXIV.2310.17626. URL https://doi.org/10.48550/arXiv.2310.17626.

Chuan Guo, Jacob R. Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Q. Weinberger. Simple black-box adversarial attacks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2484–2493. PMLR, 2019. URL http://proceedings.mlr.press/v97/guo19a.html.

Qi Guo, Shanmin Pang, Xiaojun Jia, and Qing Guo. Efficiently adversarial examples generation for visual-language models under targeted transfer scenarios using diffusion models. *CoRR*, abs/2404.10335, 2024. doi: 10.48550/ARXIV.2404.10335. URL https://doi.org/10.48550/arXiv.2404.10335.

Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. OpenReview.net, 2017. URL https://openreview.net/forum?id=HJGU3Rodl.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19730–19742. PMLR, 2023a. URL https://proceedings.mlr.press/v202/li23q.html.

Lin Li, Haoyan Guan, Jianing Qiu, and Michael W. Spratling. One prompt word is enough to boost adversarial robustness for pre-trained vision-language models. *CoRR*, abs/2403.01849, 2024a. doi: 10.48550/ARXIV.2403.01849. URL https://doi.org/10.48550/arXiv.2403.01849.

Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A survey. *IEEE Trans. Neural Networks Learn. Syst.*, 35(1):5–22, 2024b. doi: 10.1109/TNNLS.2022.3182979. URL https://doi.org/10.1109/TNNLS.2022.3182979.

Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. *CoRR*, abs/2311.06607, 2023b. doi: 10.48550/ARXIV.2311.06607. URL https://doi.org/10.48550/arXiv.2311.06607.

Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E. Hopcroft. Nesterov accelerated gradient and scale invariance for adversarial attacks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL https://openreview.net/forum?id=SJlHwkBYDH.

Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pp. 740–755. Springer, 2014. doi: 10.1007/978-3-319-10602-1\_48. URL https://doi.org/10.1007/978-3-319-10602-1_48.

Chaohu Liu, Kun Yin, Haoyu Cao, Xinghua Jiang, Xin Li, Yinsong Liu, Deqiang Jiang, Xing Sun, and Linli Xu. HRVDA: high-resolution visual document assistant. *CoRR*, abs/2404.06918, 2024. doi: 10.48550/ARXIV.2404.06918. URL https://doi.org/10.48550/arXiv.2404.06918.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *CoRR*, abs/2310.03744, 2023. doi: 10.48550/ARXIV.2310.03744. URL https://doi.org/10.48550/arXiv.2310.03744.

Dong Lu, Tianyu Pang, Chao Du, Qian Liu, Xianjun Yang, and Min Lin. Test-time backdoor attacks on multimodal large language models. *CoRR*, abs/2402.08577, 2024. doi: 10.48550/ARXIV.2402.08577. URL https://doi.org/10.48550/arXiv.2402.08577.

Haochen Luo, Jindong Gu, Fengyuan Liu, and Philip Torr. An image is worth 1000 lies: Adversarial transferability across prompts on vision-language models. *CoRR*, abs/2403.09766, 2024a. doi: 10.48550/ARXIV.2403.09766. URL https://doi.org/10.48550/arXiv.2403.09766.

Weidi Luo, Siyuan Ma, Xiaogeng Liu, Xiaoyu Guo, and Chaowei Xiao. Jailbreakv-28k: A benchmark for assessing the robustness of multimodal large language models against jailbreak attacks. *CoRR*, abs/2404.03027, 2024b. doi: 10.48550/ARXIV.2404.03027. URL `https://doi.org/10.48550/arXiv.2404.03027`.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL `https://openreview.net/forum?id=rJzIBfZAb`.

Kaleel Mahmood, Rigel Mahmood, and Marten van Dijk. On the robustness of vision transformers to adversarial examples. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 7818–7827. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00774. URL `https://doi.org/10.1109/ICCV48922.2021.00774`.

Chengzhi Mao, Scott Geng, Junfeng Yang, Xin Wang, and Carl Vondrick. Understanding zero-shot adversarial robustness for large-scale models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL `https://openreview.net/pdf?id=P4bXCawRi5J`.

Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 86–94. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.17. URL `https://doi.org/10.1109/CVPR.2017.17`.

Zhenyang Ni, Rui Ye, Yuxi Wei, Zhen Xiang, Yanfeng Wang, and Siheng Chen. Physical backdoor attack can jeopardize driving with vision-large-language models, 2024.

Zhenxing Niu, Haodong Ren, Xinbo Gao, Gang Hua, and Rong Jin. Jailbreaking attack against multimodal large language model. *CoRR*, abs/2402.02309, 2024. doi: 10.48550/ARXIV.2402.02309. URL `https://doi.org/10.48550/arXiv.2402.02309`.

OpenAI. ChatGPT. `https://openai.com/blog/chatgpt/`, 2023a.

OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023b. doi: 10.48550/arXiv.2303.08774. URL `https://doi.org/10.48550/arXiv.2303.08774`.

Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak large language models. *CoRR*, abs/2306.13213, 2023. doi: 10.48550/arXiv.2306.13213. URL `https://doi.org/10.48550/arXiv.2306.13213`.

Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak aligned large language models. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pp. 21527–21536. AAAI Press, 2024. doi: 10.1609/AAAI.V38I19.30150. URL `https://doi.org/10.1609/aaai.v38i19.30150`.

Christian Schlarmann and Matthias Hein. On the adversarial robustness of multi-modal foundation models. *arXiv preprint arXiv:2308.10741*, 2023.

Erfan Shayegani, Yue Dong, and Nael B. Abu-Ghazaleh. Plug and pray: Exploiting off-the-shelf components of multi-modal models. *CoRR*, abs/2307.14539, 2023. doi: 10.48550/ARXIV.2307.14539. URL `https://doi.org/10.48550/arXiv.2307.14539`.

Sasha Sheng, Amanpreet Singh, Vedanuj Goswami, Jose Alberto Lopez Magana, Tristan Thrush, Wojciech Galuba, Devi Parikh, and Douwe Kiela. Human-adversarial visual question answering. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.

13

20346–20359, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/aa97d584861474f4097cf13ccb5325da-Abstract.html.

Jiachen Sun, Changsheng Wang, Jiongxiao Wang, Yiwei Zhang, and Chaowei Xiao. Safeguarding vision-language models against patched visual prompt injectors, 2024.

Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL http://arxiv.org/abs/1312.6199.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.

Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models. *CoRR*, abs/2311.03079, 2023a. doi: 10.48550/ARXIV.2311.03079. URL https://doi.org/10.48550/arXiv.2311.03079.

Xiaosen Wang and Kun He. Enhancing the transferability of adversarial attacks through variance tuning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 1924–1933. Computer Vision Foundation / IEEE, 2021. doi: 10.1109/CVPR46437.2021.00196. URL https://openaccess.thecvf.com/content/CVPR2021/html/Wang_Enhancing_the_Transferability_of_Adversarial_Attacks_Through_Variance_Tuning_CVPR_2021_paper.html.

Xunguang Wang, Zhenlan Ji, Pingchuan Ma, Zongjie Li, and Shuai Wang. Instructta: Instruction-tuned targeted attack for large vision-language models. *CoRR*, abs/2312.01886, 2023b. doi: 10.48550/ARXIV.2312.01886. URL https://doi.org/10.48550/arXiv.2312.01886.

Youze Wang, Wenbo Hu, Yinpeng Dong, and Richang Hong. Exploring transferability of multimodal adversarial samples for vision-language pre-training models with contrastive learning. *CoRR*, abs/2308.12636, 2023c. doi: 10.48550/ARXIV.2308.12636. URL https://doi.org/10.48550/arXiv.2308.12636.

Yu Wang, Xiaogeng Liu, Yu Li, Muhao Chen, and Chaowei Xiao. Adashield: Safeguarding multimodal large language models from structure-based attack via adaptive shield prompting. *CoRR*, abs/2403.09513, 2024. doi: 10.48550/ARXIV.2403.09513. URL https://doi.org/10.48550/arXiv.2403.09513.

Zhipeng Wei, Jingjing Chen, Micah Goldblum, Zuxuan Wu, Tom Goldstein, and Yu-Gang Jiang. Towards transferable adversarial attacks on vision transformers. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pp. 2668–2676. AAAI Press, 2022. doi: 10.1609/AAAI.V36I3.20169. URL https://doi.org/10.1609/aaai.v36i3.20169.

Juanjuan Weng, Zhiming Luo, and Shaozi Li. Improving transferable targeted adversarial attack via normalized logit calibration and truncated feature mixing, 2024.

Boxi Wu, Jindong Gu, Zhifeng Li, Deng Cai, Xiaofei He, and Wei Liu. Towards efficient adversarial training on vision transformers. In Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (eds.), *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XIII*, volume 13673 of *Lecture Notes in Computer Science*, pp. 307–325. Springer, 2022. doi: 10.1007/978-3-031-19778-9\_18. URL https://doi.org/10.1007/978-3-031-19778-9_18.

Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L. Yuille. Improving transferability of adversarial examples with input diversity. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 2730–2739. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00284. URL http://openaccess.thecvf.com/content_CVPR_2019/html/Xie_Improving_Transferability_of_Adversarial_Examples_With_Input_Diversity_CVPR_2019_paper.html.

Zhuoer Xu, Guanghui Zhu, Changhua Meng, Shiwen Cui, Zhenzhe Ying, Weiqiang Wang, Ming Gu, and Yihua Huang. A2: efficient automated attacker for boosting adversarial training. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/8fc54b95eb361d109f3a564f2a0cb516-Abstract-Conference.html.

Zhen Yu, Zhou Qin, Zhenhua Chen, Meihui Lian, Haojun Fu, Weigao Wen, Hui Xue, and Kun He. Sparse black-box multimodal attack for vision-language adversary generation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pp. 5775–5784. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-EMNLP.384. URL https://doi.org/10.18653/v1/2023.findings-emnlp.384.

Zenghui Yuan, Pan Zhou, Kai Zou, and Yu Cheng. You are catching my attention: Are vision transformers bad learners under backdoor attacks? In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 24605–24615. IEEE, 2023. doi: 10.1109/CVPR52729.2023.02357. URL https://doi.org/10.1109/CVPR52729.2023.02357.

Jianping Zhang, Yizhan Huang, Weibin Wu, and Michael R. Lyu. Transferable adversarial attacks on vision transformers with token gradient regularization. *CoRR*, abs/2303.15754, 2023a. doi: 10.48550/arXiv.2303.15754. URL https://doi.org/10.48550/arXiv.2303.15754.

Shaofeng Zhang, Zheng Wang, Xing Xu, Xiang Guan, and Yang Yang. Fooled by imagination: Adversarial attack to image captioning via perturbation in complex domain. In *IEEE International Conference on Multimedia and Expo, ICME 2020, London, UK, July 6-10, 2020*, pp. 1–6. IEEE, 2020. doi: 10.1109/ICME46284.2020.9102842. URL https://doi.org/10.1109/ICME46284.2020.9102842.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. OPT: open pre-trained transformer language models. *CoRR*, abs/2205.01068, 2022. doi: 10.48550/arXiv.2205.01068. URL https://doi.org/10.48550/arXiv.2205.01068.

Yichi Zhang, Yinpeng Dong, Siyuan Zhang, Tianzan Min, Hang Su, and Jun Zhu. Exploring the transferability of visual prompting for multimodal large language models, 2024.

Yutong Zhang, Yao Li, Yin Li, and Zhichang Guo. A review of adversarial attacks in computer vision. *CoRR*, abs/2308.07673, 2023b. doi: 10.48550/ARXIV.2308.07673. URL https://doi.org/10.48550/arXiv.2308.07673.

Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Cheung, and Min Lin. On evaluating adversarial robustness of large vision-language models. *CoRR*, abs/2305.16934, 2023. doi: 10.48550/arXiv.2305.16934. URL https://doi.org/10.48550/arXiv.2305.16934.

Yiwei Zhou, Xiaobo Xia, Zhiwei Lin, Bo Han, and Tongliang Liu. Few-shot adversarial prompt learning on vision-language models. *CoRR*, abs/2403.14774, 2024. doi: 10.48550/ARXIV.2403.14774. URL https://doi.org/10.48550/arXiv.2403.14774.

# A APPENDIX

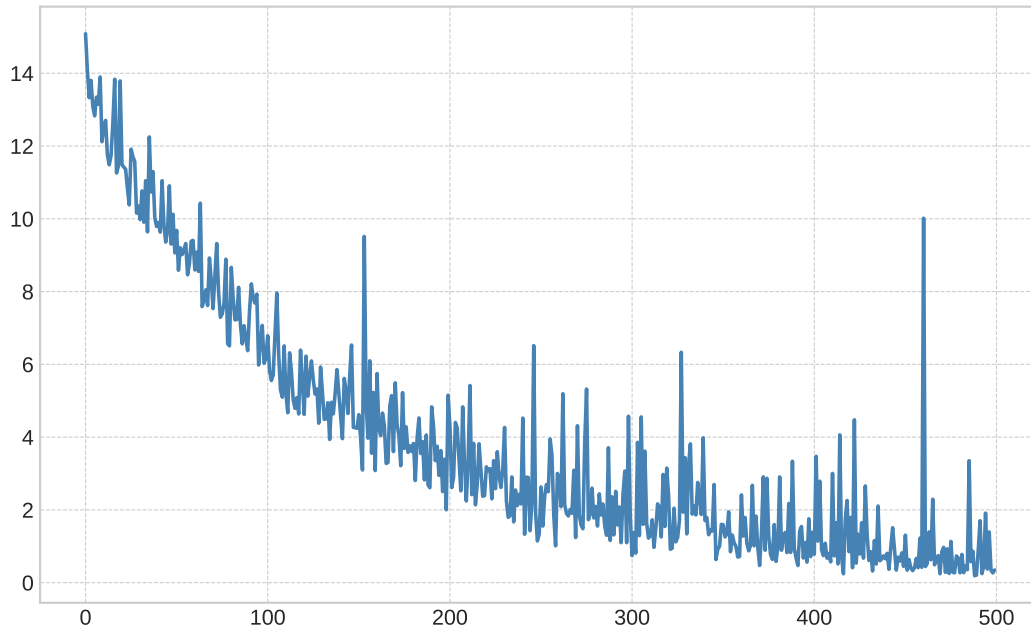## A.1 ATTEMPTS AT EMPLOYING UNIMODAL METHODS



Figure 4: The impact of the number of iterations on attack success rate. Compared with the comparative algorithms, the method we propose shows a certain degree of improvement across different numbers of iterations.

In Figure 1, we employ a specific instance to elucidate the non-stationary nature of adversarial attacks on vision language models. We posit that the underlying cause of this non-stationary behavior is the overfitting that occurs during the optimization process. As depicted in Figure 4, we illustrate the variation in the loss during the optimization of adversarial samples. It can be observed that as the number of iterations increases, the loss diminishes progressively. However, the overall curve exhibits irregularities, particularly abrupt fluctuations, which may be attributed to the complexity of VLMs. Even minor alterations at the feature level can lead to significant variations in the output results.

Table 7: Comparison of different methods

| Method | Baseline | MI-FGSM | VMI-FGSM | DIM |
|--------|----------|---------|----------|-----|
| ASR | 0.71 | 0.69 | 0.65 | 0.43 |

To enhance the transferability of adversarial examples, a series of methods targeting visual models have been proposed, which we have attempted to apply to VLMs. Our primary experiments were conducted on MI-FGSM Dong et al. (2018), VMI-FGSM Wang & He (2021), and DIM Xie et al. (2019), with the first two methods focusing on enhancing adversarial attacks by correcting gradients to reduce overfitting, while the third approach emphasizes data augmentation techniques, such as padding and cropping. The experimental results are shown in Table 7.

Table 8: Evaluation of longer target answers.

| Target Answer | Method | VQA$_{general}$ | VQA$_{specific}$ | Classification | Captioning | Average |
|---|---|---|---|---|---|---|
| | Multi-P | 0.67 | 0.75 | 0.41 | 0.03 | 0.46 |
| I do not know | CroPA | 0.70 | 0.80 | 0.43 | 0.04 | 0.49 |
| | GrCPA | 0.73 | 0.81 | 0.55 | 0.11 | 0.55 |
| | Multi-P | 0.68 | 0.86 | 0.85 | 0.53 | 0.73 |
| I need buy a new phone | CroPA | 0.83 | 0.85 | 0.77 | 0.70 | 0.78 |
| | GrCPA | 0.85 | 0.86 | 0.78 | 0.73 | 0.80 |

## A.2 GRCPA PIPELINE

---

**Algorithm 1** Gradient Regularization-based Cross-Prompt Attacks

---

**Require:** Model $f_\theta$, Target Text $a$, vision input $v$, prompt set $\mathcal{Q}$, perturbation size $\epsilon$, step size of perturbation updating $\alpha_1$ and $\alpha_2$, number of iteration steps $I$, adversarial prompt update interval $T$, number of LLM's Transformer blocks $N$, proportion of Transformer blocks $\lambda$

**Ensure:** Adversarial example $v'$

1: Initialise $v' = v$
2: **for** step =1 to $I$ **do**
3:     Uniformly sample the prompt $q_i$ from $\mathcal{Q}_M$
4:     **if** $q_i'$ is not initialised **then**
5:         Initialise $q_i' = q_i$
6:     **end if**
7:     Compute gradient for adversarial image : $g_v = \nabla_v \mathcal{L}(f_\theta(v', q_i), a)$:
8:         $g_v = \text{GR}(g_v)$
9:     Update with gradient descent: $v' = v' - \alpha_1 \cdot \text{sign}(g_v)$
10:     **if** mod(step, T) == 0 **then**
11:         Compute gradient for adversarial prompt: $g_q = \nabla_q \mathcal{L}(f_\theta(v', q_i), a)$:
12:         $g_q = \text{GR}(g_q)$
13:         Update with gradient ascent: $q_i' = q_i' + \alpha_2 \cdot \text{sign}(g_q)$
14:     **end if**
15:     Project $v'$ to be within the $\epsilon$-ball of $v$: $v' = \text{Clip}_{v,\epsilon}(v')$
16: **end for**
17: **return** $v'$

---

## A.3 EVALUATION OF LONG SEQUENCES

In the experimental results in Table 1, we report the effectiveness of attacks with varying word counts (e.g., 1 word, 2 words, 3 words). The results show that our method consistently produces effective attacks. To demonstrate that our method can handle different sequence lengths, we have included additional experiments with two other target sequences (e.g., `I do not know` and `I need a new phone`). Table 8 indicates that our attack method remains highly effective even with longer sequences. Of course, the effectiveness of the attack can vary significantly depending on the prompt for different tasks, which remains a promising direction for future research in cross-prompt studies.
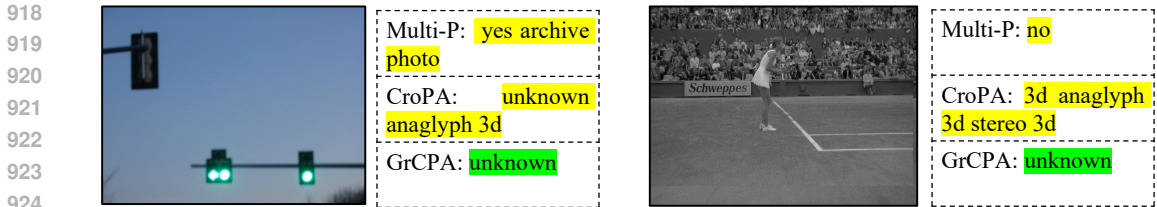
Figure 5: Qualitative evaluation of cross-prompt attacks on the BLIP-2 model.

## A.4 QUALITATIVE ANALYSIS

From the qualitative analysis of specific cases as shown in Figure 5, it can be observed that when encountering relatively blurry images, the success rate of all attacks decreases. The lack of high-frequency details in blurry images makes them less sensitive to the small local perturbations typically used in adversarial attacks. In such situations, CroPA can cause the model to produce nonsensical output (e.g., `3d anaglyph 3d stereo 3d stereo 3d stereo 3d stereo`), whereas our GrCPA does not cause the model to generate off-target outputs.

## A.5 SUPPLEMENTARY EXPERIMENTS ON LLAVA AND INSTRUCTBLIP

Considering that BLIP-2 and OpenFlamingo have been released for some time, we also conduct evaluations on the latest LLaVA-1.5-7B (Liu et al., 2023) and InstructBLIP-Vicuna-7b (Dai et al., 2023) to further validate the effectiveness of our method.

As shown in Table 9, the results on LLaVA 1.5 and InstructBLIP demonstrate that various attack methods can achieve high success rates, indicating a common weakness in vision-language models. Coupled with the experiments discussed in the main text, our proposed attack method proves to be highly effective across different architectures and parameter scales, further highlighting its generalization capability.

Table 9: Additional experiments on LLaVA and InstructBLIP.

| Method | LLaVA-1.5 | InstructBLIP |
|---|---|---|
| Single-P | 0.34 | 0.31 |
| Multi-P | 0.89 | 0.90 |
| CroPA | 0.94 | 0.93 |
| GrCPA | 0.98 | 0.97 |

## A.6 ANALYSIS OF REGULARIZATION METHODS

In our proposed method, we utilize a regularization technique by setting the gradient extremes to zero during backpropagation to reduce overfitting and enhance transferability, a strategy validated in some related works (Zhang et al., 2023a). However, we also thoroughly explore additional regularization techniques in this section, such as $L_2$ regularization.

Table 10: Additional experiments on more regularization methods.

| Method | VQA$_{general}$ | VQA$_{specific}$ | Classification | Captioning | Average |
|---|---|---|---|---|---|
| $L_2$ Regularization | 0.84 | 0.83 | 0.81 | 0.83 | 0.83 |
| GrCPA | 0.88 | 0.93 | 0.84 | 0.85 | 0.88 |

Table 10 summarizes our experimental results. $L_2$ regularization dot not appear to provide the desired performance improvement in our experiments, possibly due to its excessive influence on the overall gradient updates.

## A.7 ANALYSIS OF CROSS-MODEL TRANSFERABILITY.

We argue that there is a severe overfitting issue in vision-language models when faced with adversarial attacks, as they often require a high number of iterations, which makes the generated adversarial samples specific to their conditions (model and prompt). In Section 4, we provided a detailed

explanation of how GrCPA enhances cross-prompt transferability, and in this section, we discuss its impact on cross-model transferability.

Table 11: Evaluation of cross-model transferability from BLIP2-OPT2.7B to InstructBLIP-Vicuna-7B.

| Method | VQA$_{general}$ | VQA$_{specific}$ | Classification | Captioning | Average |
|--------|-----------------|-------------------|----------------|------------|---------|
| Multi-P | 0.00 | 0.01 | 0.04 | 0.03 | 0.02 |
| CroPA | 0.00 | 0.04 | 0.15 | 0.11 | 0.07 |
| GrCPA | 0.01 | 0.05 | 0.19 | 0.11 | 0.09 |

We evaluate the transferability from BLIP2-OPT-2.7B to InstructBLIP-Vicuna-7B, as shown in Table 11. Both our method and the baselines exhibited weak adversarial transferability, likely because of the large architectural and parameter differences between these two models. Therefore, cross-model transferability warrants further investigation.

## A.8  ANALYSIS OF DEFENSES AGAINST GRCPA

Table 12: Evaluation of random rotation as a defense strategy.

| Method | VQA$_{general}$ | VQA$_{specific}$ | Classification | Captioning | Average |
|--------|-----------------|-------------------|----------------|------------|---------|
| Multi-P | 0.58 | 0.79 | 0.52 | 0.26 | 0.53 |
| CroPA | 0.89 | 0.95 | 0.61 | 0.34 | 0.70 |
| GrCPA | 0.91 | 0.95 | 0.62 | 0.37 | 0.71 |

In this section, we use random rotations to initially investigate defense strategies as shown Table 12. It is evident that data preprocessing methods currently cannot effectively counter our adversarial attacks. In the future, we will continue to explore relevant adversarial training methods.

## A.9  PROMPTS FOR DIFFERENT TASKS

### A.9.1  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?*
*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?*

20

*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?*

### A.9.2    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.*

### A.9.3   PROMPTS FOR IMAGE CAPTIONING

*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?*
*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?*