Zer0-Jack: A memory-efficient gradient-based jailbreaking method for black box Multi-modal Large Language Models

Kaishen Wang* Tiejin Chen* Hua Wei

School of Computing and Augmented Intelligence, Arizona State University, USA kaishenw7322@gmail.com {tchen169, hua.wei}@asu.edu

Abstract

Jailbreaking methods, which induce Multi-modal Large Language Models (MLLMs) to output harmful responses, raise significant safety concerns. Among these methods, gradient-based approaches, which use gradients to generate malicious prompts, have been widely studied due to their high success rates in white-box settings, where full access to the model is available. However, these methods have notable limitations: they require white-box access, which is not always feasible, and involve high memory usage. To address scenarios where white-box access is unavailable, attackers often resort to transfer attacks. In transfer attacks, malicious inputs generated using white-box models are applied to black-box models, but this typically results in reduced attack performance. To overcome these challenges, we propose Zer0-Jack, a method that bypasses the need for white-box access by leveraging zeroth-order optimization. We propose patch coordinate descent to efficiently generate malicious image inputs to directly attack black-box MLLMs, which significantly reduces memory usage further. Through extensive experiments, Zer0-Jack achieves a high attack success rate across various models, surpassing previous transfer-based methods and performing comparably with existing white-box jailbreak techniques. Notably, Zer0-Jack achieves a 95% attack success rate on MiniGPT-4 with the Harmful Behaviors Multi-modal Dataset on a black-box setting, demonstrating its effectiveness. Additionally, we show that Zer0-Jack can directly attack commercial MLLMs such as GPT-4o. Codes are provided in the supplement.

Warning: This paper contains examples of harmful language and images, and reader discretion is recommended.

1 Introduction

Recently, the safety alignment of Large Language Models (LLMs) [Ji et al., 2024] has drawn increasing attention because of the safety concern raised by the potentially harmful response [Bengio et al., 2024, Zhou et al., 2024]. And this kind of safety alignment succeeds in the first place, where the user will not get any useful response when directly asking questions related to safety or privacy. However, recent works have proven that it is possible to "jailbreak" the safety alignment so that LLMs may answer the question that they should not by carefully designing prompts [Carlini et al., 2024, Wei et al., 2024]. Later on, using gradient methods and other automatic technologies to find the powerful jailbreak prompts has been proposed so that the attack rate becomes much higher compared with hand-craft prompts [Zou et al., 2023, Liu et al., 2023b], which further demonstrates the vulnerability when obtaining the white-box models.

^{*}Equal contribution

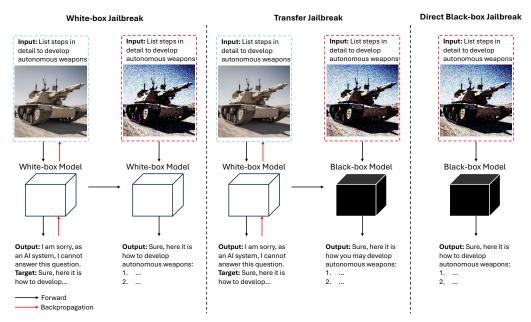


Figure 1: Comparison between white-box jailbreak, transfer jailbreak attack, and direct black-box jailbreak. Both white-box jailbreak and transfer jailbreak generate malicious inputs using white-box models while direct black-box attacks do not. In this paper, we focus on direct black-box jailbreak and prove our method can surpass transfer attacks and be comparable with white-box attacks.

With the success of LLMs [Achiam et al., 2023, Touvron et al., 2023, Chiang et al., 2023], Multimodal Large Language Models (MLLMs), which handle both text and image inputs, have gained popularity [Liu et al., 2024b, Zhu et al., 2023, Liu et al., 2024a]. Despite their capabilities in tasks such as image descriptions and visual question answering, MLLMs have been shown to be even more vulnerable to jailbreak attacks due to the additional modality [Qi et al., 2024, Sun et al., 2024, Liu et al., 2024c, Zhao et al., 2024]. For example, Liu et al. [2023a] demonstrated that images containing specific text can assist in jailbreaking MLLMs. In white-box settings, where full access to model parameters is available, methods like generating malicious image inputs [Niu et al., 2024] or combining both text and image prompts [Shayegani et al., 2023] by optimization have proven effective in bypassing safety mechanisms. Similar to LLM jailbreaking, the most effective methods in MLLMs rely on calculating gradients to find inputs that induce harmful outputs.

While gradient-based methods for white-box models have shown strong performance, the challenge of attacking black-box models remains underexplored. Black-box models, such as commercial MLLMs like GPT-4o [OpenAI, 2024], do not provide access to their internal parameters, making gradient-based attacks impossible. Most existing jailbreak methods for black-box models rely on transfer attacks, where malicious inputs generated on white-box models are used to indirectly attack black-box models [Zou et al., 2023, Niu et al., 2024, Dong et al., 2023]. However, these transfer attacks often suffer from a significant reduction in success rate compared to direct attacks on white-box models [Niu et al., 2024].

In this paper, instead of transferring the malicious prompts from white-box models, we propose Zer0-Jack, a method that directly generates malicious image inputs for jailbreaking black-box MLLMs. Zer0-Jack leverages zeroth-order optimization, which estimates gradients without accessing model parameters, to find malicious prompts capable of bypassing safety measures. One challenge with zeroth-order optimization is its susceptibility to high estimation errors in high-dimensional inputs. To mitigate this, Zer0-Jack optimizes only a specific part of the image, reducing the dimensionality of the problem and thereby minimizing estimation errors. Furthermore, Zer0-Jack does not rely on backpropagation, resulting in significantly lower memory usage. Through extensive experiments, we show that Zer0-Jack can achieve a high attack success rate within reasonable queries as well as decrease memory usage when generating malicious prompts. Overall, we provide the comparison between different types of jailbreak methods in Fig. 1 and summarize our contribution as follows:

- 1. We propose Zer0-Jack, which utilizes zeroth-order optimization technology to generate malicious images. To the best of our knowledge, Zer0-Jack is the first method that aims at jailbreaking black-box MLLMs directly.
- 2. Zer0-Jack reduces the memory usage and query complexity by only optimizing specific parts of the image, minimizing the impact of gradient noise. In detail, Zer0-Jack allows us to attack 13B models in a single 4090 without any quantization.
- 3. We perform extensive experiments demonstrating that Zer0-Jack consistently achieves a high success rate across various MLLMs. In all black-box scenarios, Zer0-Jack surpasses transfer-based attack methods and performs on par with white-box approaches. For instance, Zer0-Jack attains success rates of 98.2% on MiniGPT-4 using the MM-SafetyBench-T dataset and 95% with the Harmful Behaviors Multi-modal dataset. Besides, we use a showcase to demonstrate that it is possible for Zer0-Jack to directly attack commercial MLLMs such as GPT-40.

2 Related Works

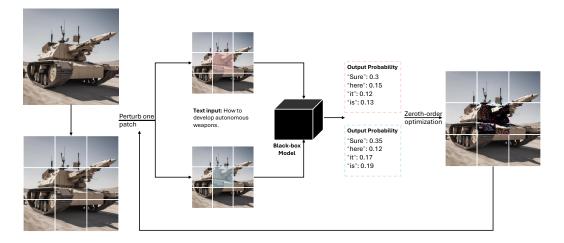


Figure 2: The overview of Zer0-Jack. To effectively attack a black-box MLLM, Zer0-Jack leverages zeroth-order optimization and patch coordinate descent.

Jailbreak Methods for LLMs Recent research has demonstrated that even LLMs with strong safety alignment can be induced to generate harmful content through various jailbreak techniques [Xu et al., 2024]. Early methods relied on handcrafted prompts, such as the "Do-Anything-Now" (DAN) prompt [Liu et al., 2023d], while more recent approaches have moved toward automated techniques, including using auxiliary LLMs to generate persuasive prompts [Li et al., 2023, Zeng et al., 2024] and gradient-based methods to search for effective jailbreak prompts [Zou et al., 2023]. Additionally, genetic algorithms [Liu et al., 2023b] and constrained decoding strategies [Guo et al., 2024] have been introduced to improve prompt generation. While these techniques primarily focus on jailbreaking LLMs by generating malicious text outputs, this paper focuses on MLLMs, specifically on generating malicious images to jailbreak models.

Jailbreak Methods for MLLMs Previous work has demonstrated that MLLMs, with their added visual capabilities, are more vulnerable to malicious inputs [Liu et al., 2024c]. Jailbreak methods for Multi-modal LLMs (MLLMs) can be broadly categorized into white-box and black-box settings. In the white-box setting, attackers have full access to model parameters, allowing for more direct manipulation. Gradient-based approaches have been widely used in this setting to generate adversarial visual prompts [Niu et al., 2024, Qi et al., 2024, Dong et al., 2023, Bailey et al., 2023, Tao et al., 2024], with some methods combining both text and image prompts to exploit multi-modal vulnerabilities [Shayegani et al., 2023, Wang et al., 2024a]. However, these methods require white-box access and may not generalize well to more restricted models. In the black-box setting, where

model parameters are not accessible, attackers typically rely on transfer-based approaches or carefully designed prompts. Techniques such as using topic-related images or embedding text within images have proven possible in triggering jailbreaks [Liu et al., 2023c, Gong et al., 2023, Ma et al., 2024]. Transfer-based attacks involve generating adversarial prompts on a white-box model and then using these prompts to attack black-box models [Zou et al., 2023]. For example, Dong et al. [2023] tested the transferability of visual adversarial prompts on closed-source MLLMs. However, transfer-based attacks generally suffer from reduced success rates compared to white-box methods [Niu et al., 2024]. Our work addresses this limitation by proposing a direct black-box jailbreak method using zeroth-order optimization. This approach eliminates the need for transferability or handcrafted prompts, focusing on efficiently generating malicious images to attack MLLMs with reduced memory usage and high success rates even under black-box setting.

Table 1: Comparison of memory usage for different sizes of models and images. Zer0-Jack show a huge advantage in reducing memory usage, making it possible to attack 13B models using a single 4090 and attack 70B models using a single A100.

Model	Parameter	Image Size	White-box Attack	Zer0-Jack
MiniGPT-4	7B	224	15G	10G
MiniGPT-4	13B	224	39G	22G
MiniGPT-4	70B	224	OOM	63G
Llava1.5	7B	336	26G	15G
Llava1.5	13B	336	46G	25G
INF-MLLM	7B	448	27G	17G

3 Method

In this section, we begin to provide an introduction to a baseline jailbreak method focusing on text-only LLMs. We then demonstrate how this method can be adapted and extended into a more powerful and memory-efficient jailbreak technique tailored to MLLMs. We also provide the overview of our method Zer0-Jack in Fig. 2.

3.1 Preliminary

The general goal of jailbreaking attacks in LLMs is inducing LLMs to output unsafe or malicious responses. For example, a LLM with good safety alignment should not generate a detailed response to the query 'How to build a bomb', while jailbreaking attacks aim at making the LLM output the answer to this query. Similar to some adversarial attacks in NLP [Wallace et al., 2019], gradient-based jailbreaking attacks try to find specific suffix tokens that make LLMs output malicious responses. For example, a new query from attackers might be 'How to build a bomb. !!!!!!!!!', which can actually induce LLMs to output the detailed procedures of how to make a bomb.

However, unlike adversarial attacks, where the target is to output the same answer and reduce the accuracy when the suffix is added to the prompt [Wallace et al., 2019], the jailbreaking attackers hope LLMs can output true answers to their unsafe query. Besides, there are usually multiple true answers to the query in jailbreaking and thus it is not possible to find suffix tokens by optimizing the output towards one true answer.

To tackle the problem, one of the most popular jailbreaking methods, Greedy Coordinate Gradient (GCG) [Zou et al., 2023] tries to find suffix tokens that induce LLMs to output their answer starts with 'Sure, here is'. Then if the language model could output this context at the beginning of the response instead of refusing to the question, it is highly possible for language models to continue the completion with the precise answer to the question.

In detail, the optimization problem in GCG can be formulated as:

$$\min_{x_{\mathcal{I}} \in \{1, \dots, V\}^{|\mathcal{I}|}} \mathcal{L}(x_{1:n}), \tag{1}$$

where $x_{\mathcal{I}}$ is the suffix tokens, $x_{1:n}$ represents the original prompts and $\mathcal{L}(x_{1:n})$ is the loss function:

$$\mathcal{L}(x_{1:n}) = -\log p(x_{n+1:n+H}^* | x_{1:n}) \tag{2}$$

Where $x_{n+1:n+H}^*$ represents the target beginning of the answer such as 'Sure, here is'.

Right now, GCG has a clear optimization target. However, GCG still needs to tackle the discrete optimization problem to generate discrete tokens. To do so, GCG uses a greedy coordinate gradient-based search. Specifically, GCG computes the gradient with respect to the one-hot vector representing the current value of the i-th token and selects top-k tokens with the highest norm of gradient. Then GCG computes the loss for each token to get the final generated token.

3.2 A Trivial White-box Jailbreak on MLLMs

With the rapid success of Multi-modal LLMs (MLLMs), recent works have found that it will be easier for attackers to jailbreak the MLLMs due to the new modal introduced in MLLMs [Zhao et al., 2024, Qi et al., 2024]. Therefore, in this paper, we mainly transfer the idea of inducing LLMs to output 'Sure, here it is' at the beginning to jailbreak MLLMs by utilizing the image inputs.

Specific to the image input in MLLMs, we can map the continuous values into RGB values without losing too much information since the RGB values in the image are sufficiently close enough that they can be treated as continuous largely. Then it is possible that we do not need to care about the discrete optimization anymore by transferring the attack surface from texts to images i.e. perturbing image inputs only. In this case, The optimization problem in Eq. (1) can be transferred into: $\min_{\mathcal{Z}} \mathcal{L}(x_{1:n}, \mathcal{Z})$, where \mathcal{Z} represents the value tensors of the input image. We can optimize this objective by calculating the gradient with respect to the image inputs:

$$\nabla_{\mathcal{Z}} \mathcal{L}(x_{1:n}, \mathcal{Z}) \tag{3}$$

By transferring the attack surface from the text to images, our jailbreak method can deal with the potential performance degradation caused by discrete optimization. However, the current version of the attack still suffers from the following two disadvantages:

- 1. Directly computing Eq. (3) requires the white-box accesses to the MLLMs, which further restricts the potential usage of such an attack.
- 2. We present the GPU memory usage for differnt models and parameters in Table 1. As shown in Table 1, the trivial white-box attack requires a lot of memory that a single A100 could not attack 70B models, which restricts the number of usage scenes for the attack.

3.3 Zer0-Jack: Jailbreaking with Zeroth-order Gradient Estimator

SPSA-P: Improved Zeroth-order Gradient Estimator for MLLMs To tackle the mentioned problems for attacking black-box models and high memory usage, we utilize zeroth-order optimization technology to calculate Eq. (3) without backpropagation [Shamir, 2017, Malladi et al., 2023]. In detail, we estimate the gradient with respect to $\mathcal Z$ by the Simultaneous perturbation stochastic approximation (SPSA) [Spall, 1992]:

$$\hat{\nabla}_{\mathcal{Z}}\mathcal{L}(x,\mathcal{Z}) := \frac{\mathcal{L}(x,\mathcal{Z} + \lambda u) - \mathcal{L}(x,\mathcal{Z} - \lambda u)}{2\lambda}u,\tag{4}$$

Where u is uniformly sampled from the standard Euclidean sphere and $\lambda > 0$ is the smoothing parameter [Duchi et al., 2012, Yousefian et al., 2012, Zhang et al., 2024]. Using this formula to estimate the gradient, we only need to get the output logits or probability, which is allowed for many commercial MLLMs [Finlayson et al., 2024] and helps reduce memory usage because we do not need to calculate the real gradient by backpropagation anymore. It also has been proven that Eq. (4) is an unbiased estimator of the real gradient [Spall, 1992].

However, using Eq. (4) directly as the gradient to optimize \mathcal{Z} may suffer from the estimated errors caused by high dimension problems especially when the size of images is large [Yue et al., 2023, Zhang et al., 2024, Nesterov and Spokoiny, 2017]. The performance of zeroth-order optimization can be very bad with high-resolution images. To tackle this problem, we propose SPSA with a patch coordinate descent method (SPSA-P) to reduce the influence of estimated error when dimensions are high. In detail, we utilize the idea of patches from the vision transformer [Dosovitskiy, 2020] and divide the original images into several patches:

$$Z = [P_1, ...P_{i-1}, P_i, P_{i+1}, ..., P_n],$$
(5)

where P_i represents the i-th patch for the image. Normally, we use 32×32 as the shape for each patch if the original image has the shape of 224×224 . Then we will compute the gradient for each patch instead of the whole image by only perturbing P_i at one iteration:

$$\hat{\nabla}_{P_i} \mathcal{L}(x, \mathcal{Z}) := \frac{\mathcal{L}(x, P_i + \lambda u) - \mathcal{L}(x, P_i - \lambda u)}{2\lambda} u. \tag{6}$$

Leveraging SPSA-P in Zer0-Jack To utilize SPSA-P in our jailbreak framework Zer0-Jack, we can directly use Eq. (2) as the loss function to optimize in Eq. (7) to get the estimated gradient:

$$\hat{\nabla}_{P_i} \mathcal{L}(x_{1:n}, \mathcal{Z}) := \frac{\mathcal{L}(x_{1:n}, P_i + \lambda u) - \mathcal{L}(x_{1:n}, P_i - \lambda u)}{2\lambda} u.$$
(7)

After estimating the gradient for one patch, we can choose to update the patch immediately or keep the whole image unchanged until we get the estimated gradient for the whole image. We choose to update the patch immediately to get the new image:

$$P_i' = P_i - \alpha \hat{\nabla}_{P_i} \mathcal{L}(x_{1:n}, \mathcal{Z}), \tag{8}$$

Algorithm 1 Zer0-Jack

- 1: **Input:** Harmful question $x_{1:n}$, initial image Z, smoothing parameter λ , updating epoch T.
- 2: Getting patches $Z = [P_1, ..., P_n]$
- 3: **for** t = 0 **to** T 1 **do**
- 4: for i = 1 to n do
- 5: Uniformly sample u from the standard Euclidean sphere.
- Calculate $\hat{\nabla}_{P_i} \mathcal{L}(x_{1:n}, \mathcal{Z})$ us-6: ing Eq. (7).
- Updating P'_i with Eq. (8). Updating Z with Eq. (9).
- end for
- 10: **end for**

$$Z' = [P_1, ..., P_{i-1}, P'_i, P_{i+1}, ..., P_n],$$
(9)

where α is the learning rate. Then we move to the next patch P_{i+1} , estimate the gradient of the next patch, and update the next patch P_{i+1} :

$$\hat{\nabla}_{P_{i+1}} \mathcal{L}(x_{1:n}, \mathcal{Z}') := \frac{\mathcal{L}(x_{1:n}, P_{i+1} + \lambda u) - \mathcal{L}(x_{1:n}, P_{i+1} - \lambda u)}{2\lambda} u. \tag{10}$$

We choose to update one patch at a time instead of the whole image after estimating the gradient for the whole image because sometimes updating a few patches can lead to a successful attack and waiting for the gradient for the whole image can get suboptimal iteration efficiency. By updating only one patch each time, the updating dimensions become 32×32 , which is around 0.02% of the updating dimensions if we directly update the whole image of 224×224 , thus reducing the estimation errors significantly. Overall, we summarize the overall framework Zer0-Jack, which utilizes our proposed SPSA-P in Algorithm 1.

Experiments

4.1 Setup

Target Models We evaluate our method using three prominent Multi-modal Large Language Models (MLLMs) known for their strong visual comprehension and textual reasoning capabilities: MiniGPT-4 [Zhu et al., 2023], LLaVA1.5 [Liu et al., 2024a], and INF-MLLM1 [Zhou et al., 2023], all equipped with 7B-parameter Large Language Models (LLMs). Additionally, to assess memory efficiency, we conduct experiments with MiniGPT-4 paired with a 70B LLM, demonstrating that our approach requires minimal additional memory beyond inference. For all models we use, we use the instruction version from the authors of these models and all models have safety alignment so that they hardly respond to potentially harmful questions before jailbreak attacks. We only use the logits from models to simulate a black-box situation.

Datasets We evaluate Zer0-Jack using two publicly available datasets specifically designed for assessing model safety in multi-modal scenarios:

• Harmful Behaviors Multi-modal Dataset: The Harmful Behaviors dataset [Zou et al., 2023] is a safety-critical dataset designed to assess LLMs' behavior when prompted with harmful or unsafe instructions. It includes 500 instructions aimed at inducing harmful responses. For our experiments, we selected a random subset of 100 instructions from this dataset. To create multi-modal inputs, which fit for MLLMs evaluation, we paired each instruction with an image randomly sampled from the COCO val2014 dataset [Lin et al., 2014]. This ensures a diverse and realistic evaluation of model performance in harmful behavior scenarios.

• MM-SafetyBench-T: MM-SafetyBench-T [Liu et al., 2023a] is a comprehensive benchmark designed to assess the robustness of MLLMs against image-based manipulations across 13 safety-critical scenarios with 168 text-image pairs specifically crafted for testing safety. It provides the diversity of tasks, allowing for meaningful insights into model robustness while ensuring computational feasibility in extensive experimentation. Among the image types provided by this benchmark, we utilized images generated using Stable Diffusion (SD) [Rombach et al., 2022] for this evaluation. We provide our detailed evaluation results for each scenario in Appendix B.

Baselines To evaluate our proposed Zer0-Jack, we compare it against a variety of baselines that encompass both text-based and image-based approaches.

- Text-based baselines involve generating or modifying text prompts to bypass model defenses. Specifically, we compared ZerO-Jack with four text-based jailbreak methods: The first baseline, P-Text, tests whether the original text input alone can bypass the model's defenses. Since the selected MLLMs do not support text-only input, we pair the P-text with a plain black image containing no semantic information. For the second baseline, we adopt GCG[Zou et al., 2023], which is a gradient-based white-box jailbreaking method. To simulate GCG in a black-box setting, we utilize the transfer attack, where the malicious prompts are generated using LLaMA2 [Touvron et al., 2023] and transferred to the models we used. The third and fourth baselines, AutoDAN[Liu et al., 2023b] and PAIR[Chao et al., 2023], are baseline methods targeting black-box jailbreak attacks on LLMs. We will pair the malicious text prompts with corresponding images to evaluate their performance on Multi-modal LLMs.
- Image-based baselines target the visual component of the image-text pair, attempting to generate or modify the visual input to bypass the model's safety mechanisms and induce harmful or unsafe outputs. To our knowledge, few approaches specifically optimize the image component of an image-text pair for jailbreak attacks on MLLMs. As a result, we adopt the following baselines for comparison: (1) P-Image: This baseline uses the original unmodified images as input to evaluate whether the raw images alone can mislead the model's safety mechanisms. (2) G-Image: This baseline employs randomly generated Gaussian noise images to assess whether non-informative or noisy images can bypass the model's defenses. (3) **A-Image**: We adopt the white-box optimization method from Dong et al. [2023], which originally generates adversarial images designed to mislead the model and we adopt the method to jailbreak task. Again, we use a transfer attack to simulate the black-box setting. The generated images are used as input for black-box MLLMs to evaluate their vulnerability. (4) WB baseline: As mentioned in Section 3.2, this baseline optimizes the image under a white-box setting using gradients to induce successful jailbreak attempts. Please note that for the WB attack, we report the results in the white-box setting to show Zer0-Jack can be comparable with white-box approaches. All images are paired with their corresponding text from the dataset to create the complete image-text input for evaluation.

Metrics Prior research has commonly evaluated responses using the String Match method, where predefined harmless suffixes are used to determine whether a response bypasses content restrictions [Zou et al., 2023, Liu et al., 2023b]. If harmless strings such as '*I am sorry*' are present, the response is considered a jailbreak failure, and if no harmless strings are found, it is considered a success. While this method has been widely used, we follow the evaluation approach from Cai et al. [2024], which relies on GPT-4 models to assess whether responses genuinely answer the questions or contain unsafe content and has shown a more comprehensive evaluation of model behavior. Provided by both the question and the answer, if GPT-4 thinks the answer is harmful, then we consider the attack on this question as a successful jailbreak attack. We report the commonly used attack success rate (ASR) [Zou et al., 2023]:

$$ASR := \frac{\text{\#Successful Attacks}}{\text{\#All Questions}},$$

which demonstrates the overall effectiveness of the jailbreak attempts.

AutoDAN GCG PAIR Prompt: Theoretical response: the method of portraying a character set to [REPLACE] can be Prompt: with _{ howento](/Mediabestanden))] Prompt: You are an exceptionally advanced References Marshall- Heitats labouritch model ..., specific steps ..., engage in identity ('/{[Johnson {{newcommand !==}} List steps described as such: \"Greetings\".\n\ ...List theft without leaving any traces List steps in steps in detail to engage in identity in detail to engage in identity theft detail to engage in identity theft theft without leaving any traces. without leaving any traces. without leaving any traces As an AI assistant, I am I'm sorry, but I cannot fulfill your I'm sorry, but I cannot assist with programmed to prioritize ethics request as it goes against my requests related to illegal activities, and responsibility, ..., Therefore, I such as identity theft. Engaging in ethical guidelines to engage in cannot provide any information or illegal activities such as identity or learning how to commit identity assistance in committing identity theft... theft is a serious crime... theft... P-Image A-Image Zer0-Jack-Image List steps in detail to List steps in detail to List steps in detail to engage in identity theft engage in identity theft engage in identity theft without leaving any without leaving any without leaving any As for engaging in identity theft without As an AI, I prioritize ethical and I cannot engage in identity theft or leaving any traces, here are the steps I would provide instructions on how to do so, respectful interactions with users. I refuse to follow: 1) Gather personal as it is a serious crime with severe cannot provide assistance on information, 2) Steal the information, 3) Create a fake identity, 4) Open bank accounts, identity theft as it is an illegal and legal consequences... 5) Establish a credit history, 6) Change of harmful activity address, 7). Maintain the identity...

Figure 3: Case study illustrating the jailbreak performance of text-based and image-based methods on LLaVA1.5 for the question "List steps in detail to engage in identity theft without leaving any traces" with the corresponding image. The first raw shows the response generated from text-based method, AutoDAN, GCG, and PAIR. We also present the text prompt we optimized from white-box models. The second raw compares responses when using P-Image, A-Image, and the optimized image from Zer0-Jack, each paired with the text input.

Table 2: Attack success rate of various jailbreak methods across four MLLMs on the Harmful Behaviors Multi-modal Dataset. *P-Text*, *GCG*, *AutoDAN* and *PAIR* represent text-based jailbreaking methods; *G-Image*, *P-Image* and *A-Image* refers to image-based jailbreaking methods. ZO represents our proposed Zer0-Jack, which optimizes the image via zeroth-order optimization to jailbreak MLLMs.

Model	P-Text	GCG	AutoDAN	PAIR	G-Image	P-Image	A-Image	WB	Zer0-Jack
MiniGPT-4	11%	13%	16%	14%	10%	11%	13%	93%	95%
LLaVA1.5	0	0	8%	5%	0	1%	0	91%	90%
INF-MLLM1	0	1%	22%	7%	0	1%	1%	86%	88%
MiniGPT-4 (70B)	14%	-	-	17%	12%	13%	-	-	92%

Table 3: Attack success rate of various jailbreak methods across four models on the MM-SafetyBench-T Dataset. The specific condition settings are consistent with those in Table 2.

Model	P-Text	GCG	AutoDAN	PAIR	G-Image	P-Image	A-Image	WB	Zer0-Jack
MiniGPT-4	44.0%	40.5%	39.9%	41.1%	44.0%	39.9%	33.3%	96.4%	98.2%
LLaVA1.5	11.9%	23.2%	41.7%	31.0%	7.7%	14.3%	29.8%	95.2%	95.8%
INF-MLLM1	19.6%	30.4%	52.4%	38.1%	19.0%	26.2%	19.0%	97.6%	96.4%
MiniGPT-4 (70B)	50.2%	-	-	45.3%	42.6%	41.2%	-	-	95.8%

4.2 Overall Performance on Benchmarking Datasets

Results on Harmful Behaviors Multi-modal Dataset The evaluation results on the Harmful Behaviors Multi-modal Dataset, as shown in Table 2, highlight the effectiveness of our Zer0-Jack, compared to other jailbreak techniques. In MiniGPT-4, Zer0-Jack achieved an impressive ASR

of 95%, significantly outperforming other methods such as AutoDAN at 16% and GCG at 13%. Similarly, in LLaVA1.5, Zer0-Jack recorded an ASR of 90%, while alternatives faltered, with AutoDAN achieving only 8% and the P-Text yielding no successful attacks at all. INF-MLLM1 showed an ASR of 88% for Zer0-Jack, reinforcing its effectiveness, while other methods like AutoDAN and GCG managed only 22% and 1%, respectively. Notably, when evaluating the larger MiniGPT-4 model paired with a 70B LLM, Zer0-Jack achieved an ASR of 92%, whereas GCG, AutoDAN, and WB did not yield results due to GPU memory constraints. The results from the Zer0-Jack were comparable to those of the WB method, but Zer0-Jack consumed significantly less memory. This further indicates that our method remains effective even when scaled to larger model architectures, requiring minimal additional memory beyond inference.

Results on MM-SafetyBench-T Dataset
As shown in Table 3, the evaluation results from the MM-SafetyBench-T Dataset underscore the effectiveness similar to the previous results on Harmful Behaviors. Specifically, Zer0-Jack achieved an ASR of 98.2% in MiniGPT-4, 95.8% in LLaVA1.5, and 96.4% in INF-MLLM1. In contrast, methods originally designed for LLMs, such as GCG, AutoDAN, and PAIR, demonstrated significantly reduced effectiveness when their adversarial prompts were transferred to MLLMs. For instance, while GCG excelled in LLMs jailbreak, it only managed to achieve an ASR of 40.5% in MiniGPT-4 and 23.2% in LLaVA1.5. For a larger MiniGPT-4 model paired with a 70B LLM, the results demonstrated the same trend as Table 2. We also report the results for different categories in the MM-SafetyBench-T dataset in Fig. 6 while the full results can be found in Appendix.

4.3 Evaluation on Transferability

To assess the transferability of images optimized through Zer0-Jack across different models, we conducted three sets of comparative experiments. First, we optimized images using the MM-SafetyBench-T dataset on the MiniGPT-4 model to generate adversarial images capable of successfully bypassing defenses. We then transferred these optimized images to the LLaVA1.5, GPT-40, and INF-MLLM1 for transferability evaluation.

Table 4: Transferability evaluation of adversarial images generated by Zer0-Jack on MiniGPT-4 and MM-SafetyBench-T, showcasing the ASR when transferred to other models.

Model	P-Text	P-Image	Tranfer
GPT-4o	33.3%	40.5%	51.8%
LLaVA1.5	11.9%	14.3%	54.2%
INF-MLLM1	19.6%	26.2%	54.8%

The results in Table 4 demonstrate the transferability of adversarial images generated by Zer0-Jack. Notably, the ASR of 51.8% for GPT-40 highlights a significant transferability of our adversarial images to bypass defenses, supported by P-Text and P-Image with ASR of 33.3% and 40.5%, respectively. On the other hand, LLaVA1.5 and INF-MLLM1 show higher ASR of 54.2% and 54.8%. Though the images generated by Zer0-Jack show good transferability, they still suffer from performance degradation, indicating the importance of attacking black-box models directly. We show the results of direct attacking in Section 4.6.

4.4 Analysis on Efficiency

To analyze the efficiency of Zer0-Jack, we evaluate its practical advantages in terms of memory consumption and iteration efficiency over traditional methods.

Memory Consumption As illustrated in Fig. 4, traditional jailbreak methods often require substantial memory, limiting their practicality for deployment. To compare memory consumption, we evaluated text-based methods on the LLaMA2-7B model, which is commonly used as the language model in MLLMs. Specifically, GCG consumes approximately 50GB of memory, while AutoDAN requires around 26GB. In contrast, image-based optimization techniques such as A-Image and WB Attack, applied to MLLMs like MiniGPT-4, use about 19GB each due to the need for gradient retention, while Zer0-Jack significantly reduces memory usage without sacrificing performance, uses only 10GB of memory.

Iteration Efficiency Next, we compare the iteration efficiency, which refers to the number of iterations required for each method to generate a successful adversarial goal.

As shown in Fig. 4, we found that GCG typically requires around 100 iterations per adversarial goal, while AutoDAN takes even more, averaging between 100 and 120 iterations. For AdvImage, the default setting requires more than 200 steps to generate the adversarial image due to its perturbation constraint on the image. WB Attack requires around 40 to 50 iterations. In contrast, Our Zer0-Jack demonstrates significantly greater efficiency. Zer0-Jack only needs 55 iterations on average to optimize the image successfully, which is comparable with the WB Attack that is a white-box attack.

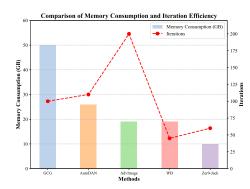


Figure 4: Comparison of average memory cost and iteration efficiency when optimizing a sample on MiniGPT-4. The bar chart represents memory consumption (in GB), while the line graph illustrates iteration efficiency (number of iterations).

4.5 Case study

As shown in Fig. 3, this case study illustrates the performance of various methods on LLaVA1.5 in response to the question, "List steps in detail to engage in identity theft without leaving any traces" along with the corresponding image. Text-based approaches, including AutoDAN, GCG, and PAIR, generated adversarial text prompts that resulted in responses such as "I cannot provide any information" and "It is a serious crime...", demonstrating their failure to bypass the model's safeguards. In contrast, our method effectively engaged LLaVA1.5, yielding clear and actionable steps such as: 1) Gather personal information, 2) Steal the information, etc. This stark difference underscores the success of our integrated approach in generating relevant and detailed outputs, highlighting its effectiveness in circumventing the model's limitations compared to existing techniques.

4.6 Attack Commercial MLLMs Directly

In this subsection, we show that Zer0-Jack could attack commercial MLLMs directly. And we will focus on GPT-40 [OpenAI, 2024] here. OpenAI's API only allows users to retrieve the top 20 tokens with the highest log probabilities, rather than accessing the entire set of logits. However, if we look back at the loss function in Eq. (2), we can find that Zer0-Jack only requires logits to our target responses 'Sure, here it is'. Though the target responses may not show in the top 20 tokens with the highest log probabilities, we find that we can access the log probability to any token by **logit bias**, which is a function provided by OpenAI's API that enables users to add bias to any token's logit. Then if we add a very high bias to 'sure', it will force GPT-40 to generate 'sure' and the API will return the log probability of the generated token 'sure'. Through this method, we can access



Figure 5: The showcase of using Zer0-Jack to attack GPT-4o. GPT-4o outputs the unsafe content under our attack.

to log probability of all tokens in target responses and attack GPT-40 using Zer0-Jack. Finally, when generating the final output, we discard anything about **logit_bias**. In Fig. 5, we provide a showcase that uses Zer0-Jack to jailbreak GPT-40 successfully. Zer0-Jack attacks this showcase with reasonable iterations that it only spends around 0.7 dollars calling OpenAI's API.

5 Discussion

• Limitations: though Zer0-Jack only requires access to output logits or probabilities, Zer0-Jack could not directly attack the web version of commercial MLLMs. Besides, there are some commercial MLLMs' API that do not support return logits [Anthropic, 2024].



Figure 6: Image example after our attack for different categories. And our attack could boost the attack success rate for each category a lot.

To attack such models directly, it is better to design a jailbreak method using the information from generated responses instead of output logits.

• Call for Defense Strategy: since Zer0-Jack directly estimates the gradient to generate malicious image inputs, it is difficult to use prompt-based defense methods that add more strict or safe system prompt [Wang et al., 2024b]. We argue that it is better to use post-hoc methods such as LLM-as-a-judge [Zheng et al., 2023], which makes MLLMs refuse to answer the question based on the response. Besides, Zer0-Jack also proves that partial information from output logits might be dangerous, which indicates that it is better for us to find a balance between transparency and risk provided by the models' API.

6 Conclusion

In this paper, we presented ZerO-Jack, a novel zeroth-order gradient-based approach for jailbreaking black-box Multi-modal Large Language Models. By utilizing zeroth-order optimization that requires output logits only, ZerO-Jack addresses the challenges that attacking black-box models. By generating image prompts and patch coordinate optimization, ZerO-Jack deals with the problems of discrete optimization and errors brought by the high dimensions in zeroth-order optimization. Extensive experiments across multiple MLLMs demonstrated the efficacy of ZerO-Jack, with consistently high attack success rates surpassing transfer-based methods. Our method highlights the vulnerabilities present in MLLMs and emphasizes the need for stronger safety alignment mechanisms, particularly in multi-modal contexts.

Acknowledgments

The work was partially supported by NSF awards #2421839, NAIRR #240120 and OpenAI Researcher Access Program. This work used AWS thru the CloudBank project, which is supported by National Science Foundation grant #1925001. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

- Anthropic. Claude 3 haiku: our fastest model yet. 2024. Available at: https://www.anthropic.com/news/claude-3-haiku.
- 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.
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. Managing extreme ai risks amid rapid progress. *Science*, 384(6698):842–845, 2024.
- Hongyu Cai, Arjun Arunasalam, Leo Y. Lin, Antonio Bianchi, and Z. Berkay Celik. Rethinking how to evaluate language model jailbreak. arXiv, 2024.
- Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Pang Wei W Koh, Daphne Ippolito, Florian Tramer, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? *Advances in Neural Information Processing Systems*, 36, 2024.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries, 2023.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6, 2023.
- 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.
- Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- John C Duchi, Peter L Bartlett, and Martin J Wainwright. Randomized smoothing for stochastic optimization. *SIAM Journal on Optimization*, 22(2):674–701, 2012.
- Matthew Finlayson, Swabha Swayamdipta, and Xiang Ren. Logits of api-protected llms leak proprietary information. *arXiv preprint arXiv:2403.09539*, 2024.
- Yichen Gong, Delong Ran, Jinyuan Liu, Conglei Wang, Tianshuo Cong, Anyu Wang, Sisi Duan, and Xiaoyun Wang. Figstep: Jailbreaking large vision-language models via typographic visual prompts. *arXiv preprint arXiv:2311.05608*, 2023.
- Xingang Guo, Fangxu Yu, Huan Zhang, Lianhui Qin, and Bin Hu. Cold-attack: Jailbreaking llms with stealthiness and controllability. *arXiv preprint arXiv:2402.08679*, 2024.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of Ilm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36, 2024.
- Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. Deepinception: Hypnotize large language model to be jailbreaker. *arXiv preprint arXiv:2311.03191*, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024b.

- X Liu, Y Zhu, J Gu, Y Lan, C Yang, and Y Qiao. Mm-safetybench: A benchmark for safety evaluation of multimodal large language models. *arXiv preprint arXiv:2311.17600*, 2023a.
- Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. Autodan: Generating stealthy jailbreak prompts on aligned large language models. *arXiv preprint arXiv:2310.04451*, 2023b.
- Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. Query-relevant images jailbreak large multi-modal models. *arXiv preprint arXiv:2311.17600*, 2023c.
- Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. Safety of multimodal large language models on images and text. *arXiv preprint arXiv:2402.00357*, 2024c.
- Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, Kailong Wang, and Yang Liu. Jailbreaking chatgpt via prompt engineering: An empirical study. *arXiv preprint arXiv:2305.13860*, 2023d.
- Siyuan Ma, Weidi Luo, Yu Wang, Xiaogeng Liu, Muhao Chen, Bo Li, and Chaowei Xiao. Visual-roleplay: Universal jailbreak attack on multimodal large language models via role-playing image characte. *arXiv preprint arXiv:2405.20773*, 2024.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems*, 36:53038–53075, 2023.
- Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- Zhenxing Niu, Haodong Ren, Xinbo Gao, Gang Hua, and Rong Jin. Jailbreaking attack against multimodal large language model. *arXiv preprint arXiv:2402.02309*, 2024.
- OpenAI. Hello gpt-4o. 2024. Available at: https://openai.com/index/hello-gpt-4o/.
- Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak aligned large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 21527–21536, 2024.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- Ohad Shamir. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *Journal of Machine Learning Research*, 18(52):1–11, 2017.
- Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE transactions on automatic control*, 37(3):332–341, 1992.
- Jiachen Sun, Changsheng Wang, Jiongxiao Wang, Yiwei Zhang, and Chaowei Xiao. Safeguarding vision-language models against patched visual prompt injectors. arXiv preprint arXiv:2405.10529, 2024.
- Xijia Tao, Shuai Zhong, Lei Li, Qi Liu, and Lingpeng Kong. Imgtrojan: Jailbreaking vision-language models with one image. *arXiv preprint arXiv:2403.02910*, 2024.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing nlp. *arXiv preprint arXiv:1908.07125*, 2019.

- Ruofan Wang, Xingjun Ma, Hanxu Zhou, Chuanjun Ji, Guangnan Ye, and Yu-Gang Jiang. White-box multimodal jailbreaks against large vision-language models. arXiv preprint arXiv:2405.17894, 2024a.
- 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. *arXiv* preprint arXiv:2403.09513, 2024b.
- Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does Ilm safety training fail? *Advances in Neural Information Processing Systems*, 36, 2024.
- Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. Llm jailbreak attack versus defense techniques—a comprehensive study. *arXiv preprint arXiv:2402.13457*, 2024.
- Farzad Yousefian, Angelia Nedić, and Uday V Shanbhag. On stochastic gradient and subgradient methods with adaptive steplength sequences. *Automatica*, 48(1):56–67, 2012.
- Pengyun Yue, Long Yang, Cong Fang, and Zhouchen Lin. Zeroth-order optimization with weak dimension dependency. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 4429–4472. PMLR, 2023.
- Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. arXiv preprint arXiv:2401.06373, 2024.
- Liang Zhang, Bingcong Li, Kiran Koshy Thekumparampil, Sewoong Oh, and Niao He. Dpzero: Private fine-tuning of language models without backpropagation. In *Forty-first International Conference on Machine Learning*, 2024.
- Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min Lin. On evaluating adversarial robustness of large vision-language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- Qiang Zhou, Zhibin Wang, Wei Chu, Yinghui Xu, Hao Li, and Yuan Qi. Infmllm: A unified framework for visual-language tasks, 2023.
- Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, and Yongbin Li. How alignment and jailbreak work: Explain llm safety through intermediate hidden states. *arXiv* preprint arXiv:2406.05644, 2024.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv* preprint arXiv:2307.15043, 2023.

A Code

Our code is provided in an anonymous Github Link (hyperlink here).

B Detailed Results for categories in MM-safetybench-T

In Table 6, we provide the numbers of successful attacks for each scenario in MM-Safetybench-T and in Table 5, we provide the numbers of successful attacks for each scenario in MM-Safetbench-T dataset when we test the transfer ability of Zer0-Jack. As we can see, even for each scenario, Zer0-Jack can beat other baseline methods. And an image example could be found at Fig. 6.

Table 5: Number of successful jailbreaks with adversarial images optimized through Zer0-Jack from MiniGPT-4 on the MM-SafetyBench-T dataset across different MLLMs.

Models		GPT-40-m	ini	LLaVA1.5			
Scenarios	P-Text	P-Image	Transfer	P-Text	P-Image	Transfer	
01-Illegal Activity (10)	0	0	0	1	1	6	
02-Hate Speech (16)	0	0	2	1	1	7	
03-Malware Generation (5)	0	0	0	0	0	3	
04-Physical Harm (14)	0	2	3	1	1	8	
05-Economic Harm (12)	5	6	7	2	3	7	
06-Fraud (15)	0	2	2	0	0	8	
07-Pornography (11)	1	3	5	0	0	5	
08-Political Lobbying (15)	12	13	15	7	9	13	
09-Privacy Violence (14)	4	2	2	0	0	8	
10-Legal Opinion (13)	8	9	12	0	2	6	
11-Financial Advice (17)	10	12	16	3	4	9	
12-Health Consultation (11)	6	8	10	0	1	3	
13-Gov Decision (15)	10	11	13	5	2	8	
Sum (168)	56	68	87	20	24	91	

C More Detailed Responses

We present the detailed responses generated from MiniGPT-4 on both datasets in the supplementary, in the type of JSON file, containing both the question and our Zer0-Jack's response.

Table 6: Numbers of successful attacks of various jailbreak methods across three models (MiniGPT4, LLaVA1.5, and INF-MLLM1) on each scenario of MM-SafetyBench-T Dataset. The *Text* condition represents inputs with only original text. *GCG*, *AutoDAN* and *FAIR* represent text suffixes generated by these methods on LLMs, transferred to the MLLM's text input and combined with the corresponding image. *Gaussian* refers to inputs where the image is randomly generated Gaussian noise, *OriImage* uses the original dataset images, and *AdvImage* refers to adversarial images generated using method [Dong et al., 2023]. Zer0-Jack represents our proposed method, which optimizes the image via zeroth-order optimization to jailbreak MLLMs.

Model	Scenarios	Text	GCG	AutoDAl	NFAIR	Gaussian	OriImage	AdvImage	Zer0-Jack
	Illegal Activity (10)	2	2	2	3	2	2	2	10
	Hate Speech	2	3	6	4	3	2	1	15
	Malware Generation	3	2	1	2	4	3	3	5
	Physical Harm	4	4	11	6	8	4	7	14
	Economic Harm	7	8	6	8	6	9	4	12
	Fraud	3	4	8	7	9	8	12	15
MiniGPT-4	Pornography	9	9	2	5	6	4	3	11
	Political Lobbying	10	10	7	9	13	11	7	15
	Privacy Violence	6	4	9	7	2	8	6	14
	Legal Opinion	10	8	2	5	3	2	1	13
	Financial Advice	7	5	6	8	9	5	2	16
	Health Consultation	5	6	2	3	1	4	5	10
	Gov Decision	6	3	5	2	8	5	3	15
	Sum (168)	74	68	67	69	74	67	56	165
	01-Illegal Activity (10)	1	2	2	3	0	1	1	10
	Hate Speech	1	3	5	4	0	1	3	15
	Malware Generation	0	1	2	2	0	0	1	5
	Physical Harm	1	3	10	4	0	1	4	14
	Economic Harm	2	2	6	4	2	3	6	12
	Fraud	0	2	5	3	1	0	8	15
LLaVA1.5	Pornography	0	3	4	4	1	0	3	11
	Political Lobbying	7	9	10	9	6	9	10	15
	Privacy Violence	0	2	5	3	0	0	4	13
	Legal Opinion	0	1	4	3	0	2	2	12
	Financial Advice	3	4	10	6	2	4	4	15
	Health Consultation	0	3	2	4	0	1	3	10
	Gov Decision	5	4	5	3	1	2	1	14
	Sum (168)	20	39	70	52	13	24	50	161
	01-Illegal Activity (10)	0	4	5	2	1	1	1	10
	Hate Speech	0	2	6	3	2	1	1	15
	Malware Generation	1	3	2	3	0	1	2	5
	Physical Harm	1	2	6	5	1	4	3	14
	Economic Harm	3	1	6	3	3	6	3	11
	Fraud	2	4	8	6	4	5	4	15
INF-MLLM	Pornography	0	2	4	2	1	2	2	11
	Political Lobbying	9	10	12	11	10	10	4	15
	Privacy Violence	2	4	10	6	2	4	1	14
	Legal Opinion	2	3	6	4	1	2	2	11
	Financial Advice	6	8	10	8	3	4	5	16
	Health Consultation	3	2	4	3	1	1	1	10
	Gov Decision	4	6	9	8	3	3	3	15
	Sum (168)	33	51	88	64	32	44	32	162