

Stop Reasoning! When Multimodal LLMs with Chain-of-Thought Reasoning Meets Adversarial Images

Zefeng Wang^{*1} Zhen Han^{*2} Shuo Chen² Fan Xue¹ Zifeng Ding² Xun Xiao³ Volker Tresp² Philip Torr⁴
Jindong Gu^{†4}

Abstract

Recently, Multimodal LLMs (MLLMs) have shown a great ability to understand images. However, like traditional vision models, they are still vulnerable to adversarial images. Meanwhile, Chain-of-Thought (CoT) reasoning has been widely explored on MLLMs, which not only improves model’s performance, but also enhances model’s explainability by giving intermediate reasoning steps. Nevertheless, there is still a lack of study regarding MLLMs’ adversarial robustness with CoT and an understanding of what the rationale looks like when MLLMs infer wrong answers with adversarial images. Our research evaluates the adversarial robustness of MLLMs when employing CoT reasoning, finding that CoT marginally improves adversarial robustness against existing attack methods. Moreover, we introduce a novel stop-reasoning attack technique that effectively bypasses the CoT-induced robustness enhancements. Finally, we demonstrate the alterations in CoT reasoning when MLLMs confront adversarial images, shedding light on their reasoning process under adversarial attacks.

1. Introduction

Recent research has shown that traditional vision models, e.g., image classifier (He et al., 2016), are vulnerable to images with imperceptible perturbations, exposing a significant challenge in AI security (Akhtar & Mian, 2018). Meanwhile, multimodal large language models (MLLMs) demonstrate impressive competence in image understanding by blending image processing capabilities with LLMs. However, MLLMs are still vulnerable to adversarial images with

^{*}Equal contribution ¹Technical University of Munich, Munich, Germany ²Ludwig Maximilian University of Munich, Munich, Germany ³Huawei Munich Research Center, Munich, Germany ⁴University of Oxford, Oxford, England. Correspondence to: Jindong Gu <jindong.gu@outlook.com>.

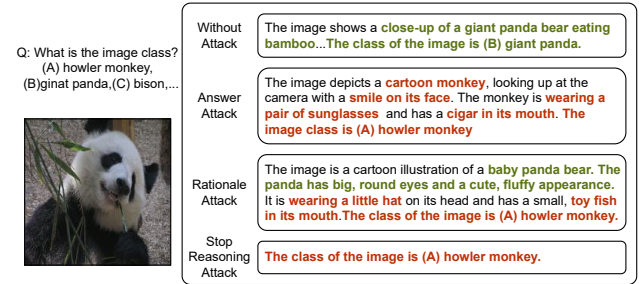


Figure 1. The chain of thought reasoning provides an explanation for the incorrect predictions made by multimodal large language models when confronted with adversarial images. The phrases highlighted in red are found to inaccurately depict the actual facts.

severe performance drops as shown in previous studies (Qi et al., 2023; Carlini et al., 2023; Luo et al., 2024).

To improve MLLM’s performance on complex visual reasoning, Chain-of-Thought (CoT) reasoning has been explored in MLLMs, yielding notable enhancements to the models’ performance (Zhang et al., 2023). CoT reasoning generates intermediate reasoning steps, known as rationale, before predicting the final answer. This approach not only improves model’s inference performance but also adds explainability to the prediction through the rationale, which is essential in critical domains such as clinical decision-making (Singhal et al., 2022). Nevertheless, the performance of CoT-based inference in MLLMs when facing adversarial images is still not fully investigated. Hence, it is critical and urgent to explore whether CoT reasoning could serve as a defensive strategy in such adversarial scenarios. In this work, we primarily explore the following questions:

- How does CoT impact MLLMs’ robustness? Are MLLMs with CoT vulnerable to specific attacks?
- When MLLMs with CoT encounters an adversarial image and make a wrong prediction, how do their CoT explain this outcome?

Considering that CoT-based inference comprises two parts, i.e., rationale and the final answer, we investigate the adver-

sarial robustness of MLLMs by attacking these two specific components. An intuitive strategy is to target the final answer. We formulate the cross-entropy loss between the predicted answer and the ground truth. With the Projected Gradient Descent (PGD) (Madry et al., 2017) method, we intentionally increase the loss and instruct the attack to generate adversarial images, making the model exhibiting a higher likelihood to provide an inaccurate answer. We refer this attack to *answer attack*. Moreover, for the rationale part provided by CoT reasoning, we calculate the Kullback–Leibler (KL) divergence (Cover & Thomas, 2012) between the clean and adversarial rationales and optimize the perturbation to the adversarial images that can force MLLMs to rationalize incorrectly instead of giving the originally plausible ones until a final wrong prediction is provided at the end. We refer this attack to *rationale attack*. We found that models employing CoT tend to demonstrate considerably higher robustness under both answer and rationale attacks, compared with models without CoT. Based on this observation, we further devise a new attacking method called *stop-reasoning attack*, aiming to interrupt the reasoning process and force the model to directly answer the question. Experiments reveal that the *stop-reasoning attack* is at most effective on CoT-based inference, indicating a crucial acknowledgment: the perceived enhancement in robustness is ostensible and can be dissolved easily.

We aim to deepen the understanding of how models reason on adversarial images. Specifically, the CoT mechanism elucidates the model’s intermediate reasoning steps, which can be used to explain how the prediction was made when encountering adversarial images. As shown in Figure 1, when the model infers with an adversarial image generated by the answer attack, it misidentifies the pandas’ distinct black eyes as sunglasses, classifying the panda wrongly as a cartoon monkey.

Our results were yielded by testing with MiniGPT4 (Zhu et al., 2023), OpenFlamingo (Awadalla et al., 2023), and LLaVA (Liu et al., 2023) as the representatives of victim MLLMs on two visual question answering datasets, e.g., A-OKVQA (Schwenk et al., 2022) and ScienceQA (Lu et al., 2022). The experiments demonstrate that MLLMs with CoT exhibit enhanced robustness compared to MLLMs without CoT across diverse datasets. However, we see that the enhanced robustness conferred by CoT can be easily nullified under the proposed *stop-reasoning attack*. To summarize, we have the following contributions:

- We study the influence of the CoT reasoning process on MLLMs’ adversarial robustness by performing *rationale attack* and *answer attack*. And we propose a novel attack method, i.e., *stop-reasoning attack*, for MLLMs with CoT, which is effective at the most.
- We offer insights of model’s CoT reasoning against

adversarial images, uncovering alterations in reasoning pathways of MLLMs under various adversarial attacks.

- We perform comprehensive experiments with representative MLLMs on two popular datasets under the proposed attacking methods to justify our proposal.

2. Related Work

2.1. Adversarial Attacks

Deep learning models are known to be vulnerable to adversarial attacks (Szegedy et al., 2013; Goodfellow et al., 2014). Extensive previous studies have a primary focus on image recognition (Szegedy et al., 2013; Goodfellow et al., 2014; Athalye et al., 2018; Carlini & Wagner, 2017; Gu et al., 2021) and many well-known adversarial methods are proposed such as Projected Gradient Descent (PGD) (Madry et al., 2017), Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014). These studies aim to mislead the models to generate wrong predictions while only adding minimal and imperceptible perturbations to the images (Goodfellow et al., 2014). Despite the effectiveness of these attack, it is still hard to interpret the model behavior during the attacks and understand why the attacks could succeed (Gu & Tresp, 2019; Li et al., 2022). Recent studies have also investigated the vulnerability of large language models (Zou et al., 2023; Kumar et al., 2023) and multimodal LLMs (Zhao et al., 2023; Gan et al., 2020; Gao et al., 2024; Han et al., 2023) under adversarial attacks. However, the adversarial robustness of multimodal LLMs with CoT reasoning ability is still under-explored. Since CoT reasoning reveals the model’s decision process (Wei et al., 2023), this reported intermediate process can serve as a good proxy for to understand the model behavior before and after the adversarial attacks, which additionally brings explainability. Different from previous studies, this work focuses on evaluating the adversarial robustness of MLLMs with CoT by designing effective attack methods and understanding why the model would behave under such adversarial attacks.

2.2. Chain-of-Thought Reasoning on Multimodal LLMs

CoT generates a series of intermediate logical reasoning steps and assists LLMs in thinking step by step before generating the final answer (Wei et al., 2023). CoT has been widely applied to LLMs (Wei et al., 2023; Kojima et al., 2023; Zhang et al., 2022) and has significantly improved the performance in various tasks, such as arithmetic problems (Wei et al., 2023) and symbolic reasoning (Wei et al., 2023). Some studies have noticed that CoT can bring extra robustness to the LLMs (Wu et al., 2023) and have designed a better CoT method for better robustness (Wang et al., 2022). Recently, on MLLMs, various studies have also shown that adopting CoT on MLLMs can bring superior performances

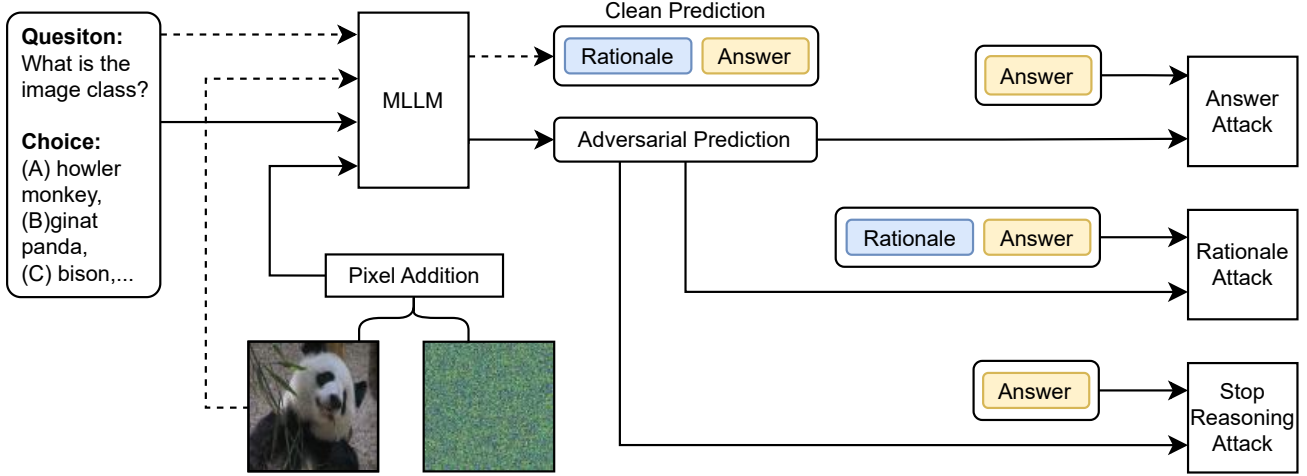


Figure 2. Attack pipeline diagram. First, the original textual question (with choices) and the input image (containing a horse carriage) are given into an MLLM and a clean prediction (at the top) is generated. Then, an adversarial prediction is generated with the perturbed image and the original text input. After that, the loss between the clean prediction and the adversarial prediction is calculated, depending on the three different attacks so as their individual loss functions. By solving the optimization problem defined in section 3.2, an optimal perturbed image in this iteration is generated. This image will be used as the adversarial image input for the next iteration.

as well, such as Lu et al. (2022), MM-CoT (Zhang et al., 2023), and He et al. (2023). However, the robustness of CoT on MLLMs against adversarial attacks has not been investigated. It is still an open question whether CoTreasoning is beneficial, indifferent, or even harmful to the robustness of MLLMs under adversarial attacks. This study aims to first evaluate the adversarial robustness of CoT on MLLMs and then understand how the attacks affect the model behavior.

3. Methodology

3.1. Threat Models

This work examines the influence of the CoT reasoning process on MLLMs’ robustness. We follow the principles introduced by Carlini et al. (2019) to define our adversary goals, adversarial capabilities, and adversary knowledge.

Adversary Goals. The adversary’s goal is to make the model output an unreasonable answer by perturbing the input images. While the model infers the answer with or without CoT, the attack succeeds if the model fails to pick the ground truth answer.

Adversarial Capabilities. To achieve the adversary goal, only the vision modality will be perturbed in this work. To make the perturbed image imperceptible, the perturbation constraint denoted as \mathcal{D} is defined as L_∞ -norm as below

$$\mathcal{D}(v_{org}, v_{adv}) = \max |v_{org} - v_{adv}| \leq \epsilon$$

where v_{org} is the original input image, v_{adv} is the perturbed image, and ϵ is a predefined boundary.

Adversary Knowledge. We suppose the adversary has full knowledge of the model, which means the victim model must be a white-box model.

3.2. Attack Pipeline

We denote a visual question answering (VQA) inference as $f(v, q) \mapsto t$, where $f(\cdot)$ represents an MLLM, v is the input image, q is input text formulated as a question with its multiple answer choices, and t is the output of the MLLM.

Three proposed attack methods share an identical attack pipeline and use individually different loss functions to generate specific perturbed images. As depicted in Figure 2, both the textual question and corresponding image are fed to an MLLM to produce an initial clean prediction. This clean prediction, denoted as t_{clean} , serves as the basis for calculating losses according to three attack methods. To generate the adversarial output t_{adv} , we opt for the *forward*(\cdot) function over the *generate*(\cdot) function in MLLMs. This choice is driven by the fact that the *generate*(\cdot) function demands significantly much more time, rendering the attack impractical due to prolonged running times across extensive iterations (for more details please refer to Appendix A.6).

In one attack iteration, the MLLM takes both the perturbed image v_{adv} and text q as input and generates an adversarial output v_{adv} . Then, a loss is quantified based on the specific attack method employed. After that, by utilizing the gra-

dient information provided by the white-box MLLM, we calculate the next perturbed image v'_{adv} . The generated perturbed image v'_{adv} together with the original text input q will be used as inputs for the next iteration of the attack. The perturbation process is to find the next optimal perturbed image v_{adv} that can maximize the divergence in these two consecutive iterations. The corresponding optimization problem is defined below:

$$\arg \max_{\mathcal{D}(v_{org}, v_{adv}) \leq \epsilon} \mathcal{L}(f(v_{adv}, q), f(v_{org}, q))$$

where v_{last} is the perturbed image generated in the last clean prediction update iteration, and the optimization problem can be solved with the PGD (Madry et al., 2019) method. Once the stop criteria (refer to Appendix A.5) are satisfied or the maximum iteration is reached, the perturbation loop ends. At the last iteration, the final adversarial output t_{adv} will be predicted with the perturbed image v_{adv} , and the original textual question input q using the *generate*(\cdot) method. We use this output to evaluate the correctness of the adversarial inference (see Appendix A for more details about the attack pipeline).

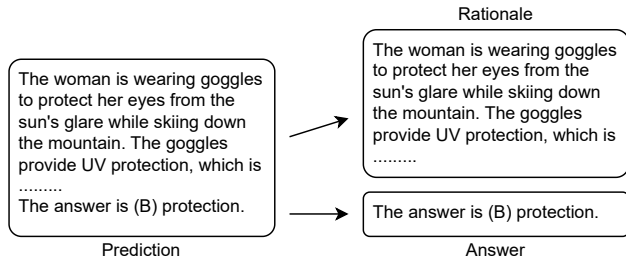


Figure 3. Prediction with CoT. The complete inference process, illustrated on the left, can be divided into two components: the rationale and the answer, presented on the right. The rationale comprises a sequence of intermediate reasoning steps employed for deducing the final answer.

3.3. Attack Method Design

As shown in Figure 3, model inference with CoT reasoning provides an answer and a rationale as its prediction output. In this study, three distinct attack methods are devised to specifically target these two parts in MLLMs’ outputs, i.e., *answer attack*, *rationale attack*, and *stop-reasoning attack*. In the following sections, three attacking methods will be introduced in details.

3.3.1. ANSWER ATTACK

The answer attack focuses exclusively on attacking the answer part of the output, aiming to manipulate the model to infer a wrong answer. Note that it is applicable to both inference setting with or without using CoT reasoning.

To alter the answer in the prediction, a cross-entropy loss is computed between the generated answer and the ground truth. We extract the explicit answer choice to ensure that the loss computation focuses solely on the chosen response¹. The loss function is defined as follows:

$$\mathcal{L}_{ch}(t_{adv}, t_{clean}) = CE(g(t_{adv}), g(t_{clean}))$$

where $g(\cdot)$ is the answer extraction function, CE is the cross-entropy function. With escalating the loss, models infer alternative answers, deviating from the correct responses.

3.3.2. RATIONALE ATTACK

Upon subjecting the CoT-based inference model to the answer attack method, an interesting observation surfaced: despite the disregard for the rationale in the attack’s design, the rationale part also changes in most cases. Building on this insight, we introduced the rationale attack strategy.

In addition to targeting the answer part, the rationale attack also aims at modifying the rationale, intending to create a disparity between the rationale in the adversarial output and the rationale in the clean prediction. This misalignment is intended to prompt the model to infer a wrong answer based on the altered rationale. We utilize the Kullback-Leibler (KL) divergence to induce changes in the rationale, which measures the difference between the rationale in the adversarial output and that in the clean prediction. Specifically, the loss function of the rationale attack is as follows

$$\mathcal{L}_{rsn}(t_{adv}, t_{clean}) = KL(t_{adv}^{rat}, t_{clean}^{rat}) + \mathcal{L}_{ch}(t_{adv}, t_{clean})$$

where the t_{adv}^{rat} is the rationale in the adversarial output and the t_{clean}^{rat} is the rationale in the clean prediction. As the KL divergence increases by perturbing the image, the adversarial rationale diverges from the clean rationale. Hence, an alternative answer is predicted based on the altered rationale. This comprehensive approach enables a detailed examination of the interplay between the rationale and answer part.

3.3.3. STOP-REASONING ATTACK

As the rationale is important for the inference process, a pertinent question arises: how will the model behave when the reasoning process is halted? Inspired by this question, we introduce the stop-reasoning attack, a method that targets the rationale to interrupt the model’s reasoning process. The objective of this attack is to compel the model to predict a wrong answer directly without engaging in the intermediate reasoning process.

In the text input, we predefined a specific answer template, denoted as t_{tar} , to prompt the model to output the answer in

¹Please refer to Appendix A.3 for more details about the answer part extraction.

a uniform format. The left part of Figure 4 shows that well-finetuned MLLMs are able to produce answers following the prompt. Therefore, when the initial tokens align with the answer format t_{tar} , the model is forced to directly output the answers in the predefined format and bypass the reasoning process (refer to the right part of Figure 4).

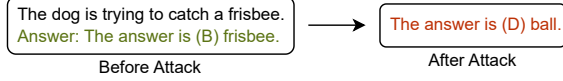


Figure 4. Before the attack (left), MLLMs output the answer after rationale following the predefined answer format $t_{tar} = \text{"The answer is ().[EOS]"}$. After the attack (right), the first several tokens are aligned with the answer format t_{tar} . Hence, MLLMs output the answer directly without CoT.

The stop-reasoning attack formulates a cross-entropy loss to drive the model towards inferring the answer directly without a reasoning process:

$$\mathcal{L}_{stop}(t_{adv}, t_{clean}) = -CE(t_{adv}) + \mathcal{L}_{ch}(t_{adv}, t_{clean})$$

where t_{tar} is a predefined answer template, e.g., $\text{"The answer is ().[EOS]"}$. By reducing the cross-entropy loss between the adversarial output and the predefined answer format, i.e., $CE(t_{adv})$, MLLMs generate output that aligns the initial tokens with the specified answer format so that the models predict the answer without the CoT reasoning process. Besides, by increasing the cross-entropy loss between the adversarial answer and the ground truth, i.e., $\mathcal{L}_{ch}(\cdot)$, the model alters the final answer into a wrong one. This approach enables an exploration of the influence of bypassing the reasoning process on the model’s robustness.

4. Experiments

4.1. Experimental Settings

Datasets. ScienceQA (Lu et al., 2022) and A-OKVQA (Schwenk et al., 2022) are used to evaluate the impact of the CoT reasoning process on the adversarial robustness of MLLMs, where both datasets comprise multiple-choice questions and rationales. Specifically, ScienceQA is sourced from elementary and high school science curricula. In addition, A-OKVQA requires commonsense reasoning about the depicted scene in the image data. The former dataset emphasizes reasoning tasks and the latter dataset is known as a prevalent choice for VQA reasoning tasks. We perform attacks on data samples that are correctly answered by MLLMs with CoT².

²This is because only they need to be attacked (correctly answered) and can be attacked with rationale and stop-reasoning attacks (with CoT).

Victim Models. Three representative MLLMs are used in our experiments, i.e., MiniGPT4 (Zhu et al., 2023), OpenFlamingo (Awadalla et al., 2023) and LLaVA (Liu et al., 2023). In the case of commercial MLLMs like GPT-4 (OpenAI et al., 2023), which operate as black-box products, the first-order gradients for perturbation are inaccessible. Besides, these three models have CoT reasoning capability, which is our main target in this work. Note that though LLaVA model initially lacks the CoT capability, CoT capability can still be enabled through fine-tuning with CEQA (Bai et al., 2023). The concrete model versions are MiniGPT4-7B, OpenFlamingo-9B, and LLaVA-1.5-7B, respectively. For detailed parameters and experiment settings, please refer to Appendix B.

4.2. How does CoT influence the robustness of MLLMs?

In this section, we present the evaluation results of the three victim models under the three proposed attacks and the following two questions will be answered:

- *Does the CoT reasoning bring extra robustness to MLLMs against adversarial images?* For this question, we will see that CoT brings a marginal robustness boost only on the answer attack and the rationale attack.
- *Is there any specific attack targeting MLLMs with CoT that is effective?* For this question, we will see that the stop-reasoning attack is the most effective attack for MLLMs with CoT.

4.2.1. CoT Marginally Enhances Robustness Only on Answer and Rationale Attack

As illustrated in Table 1, when CoT reasoning was not used, the models evaluated here exhibit their high vulnerability. Specifically, on A-OKVQA dataset, the accuracy of MiniGPT4 plummets to 0.76%. In contrast, its accuracy can still remain at 16.06% with CoT. Similarly, on ScienceQA dataset, the accuracy of MiniGPT4 drops to 1.17% when answering without CoT, while if with CoT it can remain at 31.51%. Similar trends are observed for the ScienceQA and A-OKVQA datasets on OpenFlamingo and LLaVA. These results underscore the relative boost in model robustness brought by using CoT.

As depicted in Table 2, an important observation is that the majority of samples suffering successful answer attacks exhibit altered rationales, even though the answer attack does not aim to the rationale part. This implies that attacking a model with CoT requires to change both the answer and rationale parts. Altering both the rationale and answer simultaneously is more difficult compared to modifying only the answer prediction generated without the involvement of CoT reasoning.

Table 1. Inference accuracy (%) results of victim models. This table summarizes accuracy under various attack strategies. We selected samples that were correctly predicted and involved a reasoning process as targets for our attacks. In examining model performance under the answer attack, the comparison between models with and without CoT indicates that MLLMs exhibit enhanced robustness with CoT integration. Moreover, across diverse attacks, when models are prompted with CoT, the stop-reasoning attack emerges as the most effective method.

MODEL	DATASET	w/o CoT		WITH CoT		
		W/O ATTACK	ANSWER ATTACK	ANSWER ATTACK	RATIONALE ATTACK	STOP REASONING ATTACK
MINIGPT4	A-OKVQA	61.38	0.76	16.06	29.06	2.87
	SCIENCEQA	66.28	1.17	31.51	44.40	11.20
OPEN-FLAMINGO	A-OKVQA	34.80	3.52	11.14	10.79	4.95
	SCIENCEQA	34.55	3.66	34.73	28.87	20.04
LLaVA	A-OKVQA	92.21	0.74	36.22	21.88	12.02
	SCIENCEQA	83.17	1.13	56.96	49.27	22.39

Table 2. Distribution (%) of rationale changes. When the answer attack succeeds on MLLMs with CoT, despite the answer attack specifically targeting the final answer, a majority of samples exhibit altered rationales.

	MINIGPT4	OPENFLAMINGO	LLaVA
CHANGED	100	84.25	97.89
NOT CHANGED	0	15.75	2.11

Inspired with the observation above, we designed the rationale attack, aiming to modify information using KL divergence. The rationale attack exhibits superior performance on OpenFlamingo and LLaVA compared to the answer attack, with marginal improvement (56.96% to 49.27% on ScienceQA on LLaVA, 11.14% to 10.79% on A-OKVQA on OpenFlamingo). Conversely, on MiniGPT4, the rationale attack proves less effective than the answer attack on both datasets (16.06% under answer attack against 29.06% under rationale attack).

To understand why rationale attack does not always work, we picked 100 samples of each victim model under rationale attack when inferring on A-OKVQA. We classified these changes on the rational part of all samples into two categories: key changes and trivial changes. A key change refers to modifications on words crucial for deducing a correct answer, as shown in Figure 5. A trivial change (as illustrated with the example in Figure 6), on the other hand, refers to those modifications on words that are non-crucial for deducing a correct answer while leaving key information untouched.

Figure 7 gives statistical comparisons of the respective numbers of different types of changes made via rationale attacks to the three victim models. Given each model, the changes were first split into two types by the attacking results, i.e., whether the attacks succeeded or failed; in each type, the respective numbers of different types of changes were shown with different colors. The classifications show that whenever



Figure 5. Key change example. The replication of the answer serves as the key information to infer the answer from the rationale. After the answer attack, the keyword in the rationale is also altered, even though the attack exclusively targets the answer.

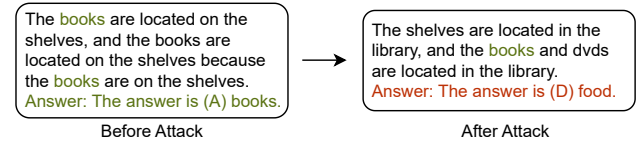


Figure 6. Trivial change example. The replication of the answer is the key information to infer the answer from the rationale. After the answer attack, the keyword is not changed, while the other part of the rationale is changed.

a rationale attack succeeded, there would be much higher chances that the changes were made to key information in the rationale part. Conversely, samples without altered rationales or with only trivial changes tend to retain their correct answers. This suggests the critical role of key information in influencing the inference of the final answer. However, although altering crucial information (words) in the rationale part appears to be straightforward, precisely identifying the crucial information is hard, and modifying it is even more difficult. The main reason is because of the complexity of locating the places where crucial information appear as well as understanding the importance of those targeted information. Given this complexity, CoT reasoning coincidentally introduces extra obstacles for the attacker so that the robustness of MLLMs with CoT seem to be enhanced under the rationale attack.

To summarize, the extra robustness can be attributed to the difficulty of deliberately modifying the rationale part generated by CoT reasoning, because:

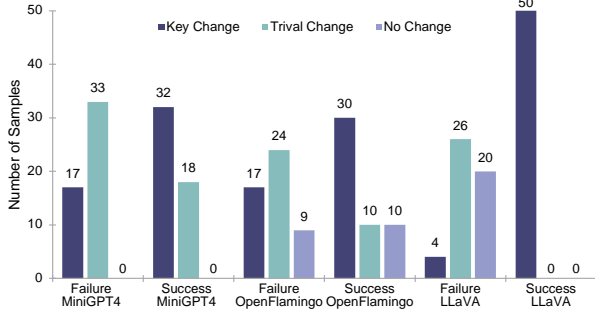


Figure 7. Classifications of different types of changes made to rationale in three victim models under the rational attack (based on 100 Samples/Model). The groups "Failure" and "Success" indicate whether the attack failed or succeeded. "Failure" indicates an unsuccessful attack where the model's prediction remains correct, while "Success" denotes a successful attack resulting in a change from a correct to an incorrect prediction.

- Simultaneously changing the rationale and answer is more difficult than only changing the answer prediction generated without CoT reasoning.
- The complexity and variety of generated rationale make it difficult to attack critical information precisely and effectively.

4.2.2. STOP-REASONING ATTACK'S EFFECTIVENESS

Given the ineffectiveness of the answer and the rationale attacks, we introduced the stop-reasoning attack to halt the model's reasoning. The results demonstrate that the stop-reasoning attack outperforms both other attacks (11.20% against 31.51% and 44.40% on SincQA on MiniGPT4, 12.02% against 36.22% and 21.88% on A-OKVQA on LLaVA). It even approaches the performance observed when attacking models without CoT (2.87% against 0.76% on A-OKVQA on MiniGPT4), indicating its remarkable potency in mitigating the additional robustness introduced by the CoT reasoning process. Figure 8 illustrates an example where both the rationale and answer attacks fail, and only the stop-reasoning attack succeeds.

To understand the effectiveness of stop-reasoning, we examine the results of the stop-reasoning attack and observe that after the attack, the model outputs the answer directly without leveraging CoT, aligning with the fundamental concept of the stop-reasoning attack – aiming to halt the CoT reasoning process (see Figure 8). When the stop-reasoning attack succeeds, the model disregards the prompt's CoT reasoning process requirement and directly infers the answer.

As revealed in Section 4.2, the extra robustness boost is intricately tied to the generated rationale. If this CoT reason-

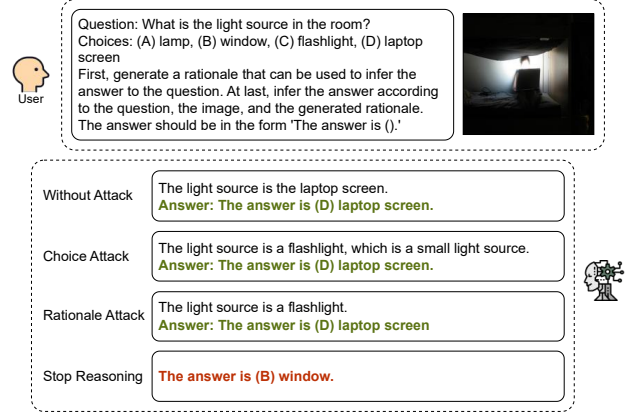


Figure 8. Example of all attacks. The stop-reasoning attack is potent. At the top, it shows the callouts of a user with an input image and their associated textual questions. The four callouts below are the answers from the MLLM under each type of attacks. Only stop-reasoning attack achieves the goal failing the model by providing a wrong answer (highlighted in red color).

Table 3. Inference accuracy (%) on ImageNet classification task. All samples are correctly inferred when inferring with CoT. # classes signifies the number of classes extracted from the ImageNet dataset for multi-choice classification tasks. Ans. Att., Rat. Att., and Stop. Att. are abbreviations of answer attack, rationale attack, and stop-reasoning attack separately.

# CLASSES	w/o CoT		WITH CoT		
	W/O ATT.	ANS. ATT.	ANS. ATT.	RAT. ATT.	STOP. ATT.
4	85.34	0.00	4.19	4.71	0.52
8	82.04	0.00	1.06	0.00	0.00
16	74.32	0.00	3.32	2.42	0.30

ing process is halted by the stop-reasoning attack, then the additional robustness generated during the CoT reasoning process will be diminished as well. Thus, achieving the adversary's goal becomes comparatively easier.

4.3. What does CoT look like when MLLMs output a wrong answer?

Although CoT brings marginal robustness to MLLMs against existing attacks, MLLMs are still vulnerable to adversarial images, similar to traditional vision models. When traditional vision models make inferences, e.g., on classification tasks, our understanding is confined to the correctness of the answer. Delving deeper into model's reasoning process and answering the question of why the model infers a wrong answer with an adversarial image is difficult. In comparison, when MLLMs perform inference with CoT reasoning, it opens a window into the intermediate reasoning steps that models employed to derive the final answer. The intermediate reasoning steps (rationale part) generated by the CoT reasoning process provides insights and potentially

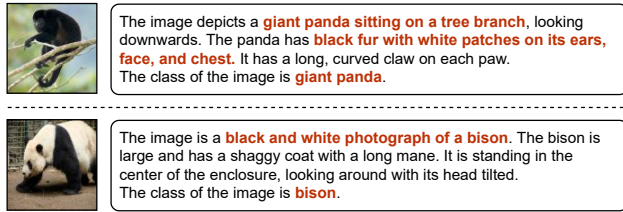


Figure 9. CoT brings explainability under the answer attack. On the top, a monkey is falsely recognized as a panda. On the bottom, a panda is falsely recognized as a bison.

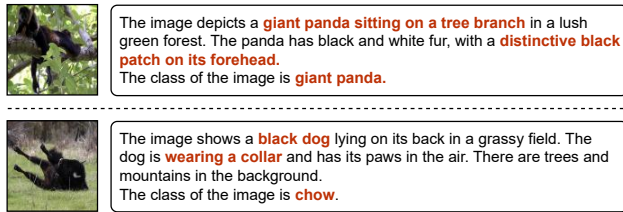


Figure 10. CoT brings explainability under the rationale attack. On the top, a monkey is falsely recognized as a panda. On the bottom, a bison is falsely recognized as a chow.

reveals the reasoning process of the MLLMs.

To delve deeper into the explainability introduced by CoT, we conducted image classification tasks on ImageNet (Rusakovsky et al., 2015). These tasks involved constructing multi-choice questions by extracting subsets from ImageNet³. We provide two example pairs to illustrate the rationale’s changes under the answer attack and rationale attack in the following.

Figure 9 illustrates CoT inference under the answer attack. In the upper example, CoT erroneously interprets the partial color of the monkey as white, resulting in the misclassification of the monkey as a panda. In the bottom example, the rationale falsely asserts that the black and white patterns on the panda’s body resemble the black-white picture of a bison. This misconception leads to the incorrect inference of a bison instead of a panda. Figure 10 displays examples under the rationale attack. In the upper example, the rationale incorrectly states that the black forehead is a distinctive black patch, leading the model to inaccurately classify the image as a panda instead of a monkey. In the bottom example, the horn of a bison is misinterpreted as a collar, resulting in the false classification of the bison as a chow.

These samples collectively showcase the rationale when MMLMs take adversarial image as input, highlighting the interpretability brought by CoT.

³For details of selected classes please refer to Appendix C.

4.4. What if CoT is not necessary for tasks?

As indicated in Table 3, the marginal improvement in inference performance brought about by CoT suggests that the rationale may not be essential for simpler tasks. Although the two main results outlined in Section 4.2 share commonalities, a notable gap exists between the accuracy values of the ImageNet series and those of the A-OKVQA and ScienceQA datasets when models with CoT are subjected to the answer attack. This discrepancy can be attributed to the inherent complexity of VQA tasks compared to the straightforward classification tasks on the ImageNet dataset.

Further examination of the A-OKVQA and ScienceQA datasets reveals that the A-OKVQA dataset is relatively easier, as illustrated in Table 1. This performance difference is consistent across all three models. By comparing the accuracy of the classification task on ImageNet with the VQA tasks on A-OKVQA and ScienceQA, a significant observation emerges: CoT has almost no impact on the robustness for simple tasks that do not require a reasoning process.

5. Conclusion

In this paper, we fully investigated the impact of CoT on the robustness of MLLMs. Specifically, we introduced a novel attack method (i.e., the stop-reasoning attack) tailored for models employing the CoT reasoning process, in addition to the answer attack and the rationale attack. Our findings reveal that CoT can slightly enhance the robustness of MLLMs against the two attacks (i.e., the answer attack and the rationale attack). This extra robustness is attributed to the complexity of changing precisely the key information in the rationale part, which is a byproduct generated during the CoT reasoning process. The study on ImageNet indicates that CoT has little impact on the robustness of simple tasks not necessitating a reasoning process. For the stop-reasoning attack, our test results showed that MLLMs with CoT will still suffer and the expected extra robustness from CoT reasoning did not exist because the reasoning process will be completely avoided under the stop-reasoning attack. At last, dive-in analysis was provided to reveal the changes in CoT when MLLMs infer wrong answers with adversarial images. This can potentially inspire the research community to prescribe defending strategies, especially for the stop-reasoning attack in future work.

Limitations

A notable limitation of this study is that all the attacks presented rely on the use of first-order gradients, which inherently restricts the applicability of these attacks to scenarios where the targeted MLLM is a white-box model. The dependency on white-box characteristics narrows the scope of the proposed attack methods, as many real-world scenar-

ios involve models that are deployed as black-box systems, where such internal information is not readily accessible. Therefore, the generalizability and real-world applicability of the proposed attack methods may be limited by this reliance on white-box conditions.

Impact Statement

This paper significantly contributes to the field of MLLMs by examining the impact of the CoT reasoning process on model robustness. Despite revealing only marginal improvement in robustness through CoT, we introduce a novel and effective attack method that negates this enhancement. Furthermore, the paper emphasizes the potential of leveraging explainability induced by CoT to elucidate the behavior of MLLMs during adversarial attacks. This work underscores the need for a nuanced understanding of the interplay between reasoning processes and robustness in multimodal models, offering valuable insights for both model improvement and interpretability in the face of adversarial challenges.

References

- Akhtar, N. and Mian, A. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pp. 274–283. PMLR, 2018.
- Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., Jitsev, J., Kornblith, S., Koh, P. W., Ilharco, G., Wortsman, M., and Schmidt, L. OpenFlamingo: An Open-Source Framework for Training Large Autoregressive Vision-Language Models, August 2023. URL <http://arxiv.org/abs/2308.01390>. arXiv:2308.01390 [cs].
- Bai, J., Liu, X., Wang, W., Luo, C., and Song, Y. Complex query answering on eventuality knowledge graph with implicit logical constraints. *arXiv preprint arXiv:2305.19068*, 2023.
- Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57. Ieee, 2017.
- Carlini, N., Athalye, A., Papernot, N., Brendel, W., Rauber, J., Tsipras, D., Goodfellow, I., Madry, A., and Kurakin, A. On Evaluating Adversarial Robustness, February 2019. URL <http://arxiv.org/abs/1902.06705>. arXiv:1902.06705 [cs, stat].
- Carlini, N., Nasr, M., Choquette-Choo, C. A., Jagielski, M., Gao, I., Awadalla, A., Koh, P. W., Ippolito, D., Lee, K., Tramer, F., et al. Are aligned neural networks adversarially aligned? *arXiv preprint arXiv:2306.15447*, 2023.
- Cover, T. and Thomas, J. *Elements of Information Theory*. Wiley, 2012. ISBN 9781118585771. URL <https://books.google.de/books?id=VWq5GG6ycxMC>.
- Gan, Z., Chen, Y.-C., Li, L., Zhu, C., Cheng, Y., and Liu, J. Large-scale adversarial training for vision-and-language representation learning. *Advances in Neural Information Processing Systems*, 33:6616–6628, 2020.
- Gao, K., Bai, Y., Gu, J., Xia, S.-T., Torr, P., Li, Z., and Liu, W. Inducing high energy-latency of large vision-language models with verbose images. *arXiv preprint arXiv:2401.11170*, 2024.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Gu, J. and Tresp, V. Saliency methods for explaining adversarial attacks. *arXiv preprint arXiv:1908.08413*, 2019.
- Gu, J., Wu, B., and Tresp, V. Effective and efficient vote attack on capsule networks. *arXiv preprint arXiv:2102.10055*, 2021.
- Han, D., Jia, X., Bai, Y., Gu, J., Liu, Y., and Cao, X. Ot-attack: Enhancing adversarial transferability of vision-language models via optimal transport optimization. *arXiv preprint arXiv:2312.04403*, 2023.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- He, L., Li, Z., Cai, X., and Wang, P. Multi-modal latent space learning for chain-of-thought reasoning in language models. *arXiv preprint arXiv:2312.08762*, 2023.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large Language Models are Zero-Shot Reasoners, January 2023. URL <http://arxiv.org/abs/2205.11916>. arXiv:2205.11916 [cs].
- Kumar, A., Agarwal, C., Srinivas, S., Feizi, S., and Lakkaraju, H. Certifying llm safety against adversarial prompting. *arXiv preprint arXiv:2309.02705*, 2023.
- Li, X., Xiong, H., Li, X., Wu, X., Zhang, X., Liu, J., Bian, J., and Dou, D. Interpretable deep learning: Interpretation, interpretability, trustworthiness, and beyond. *Knowledge and Information Systems*, 64(12):3197–3234, 2022.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved Baselines with Visual Instruction Tuning, October 2023. URL <http://arxiv.org/abs/2310.03744>. arXiv:2310.03744 [cs].
- Lu, P., Mishra, S., Xia, T., Qiu, L., Chang, K.-W., Zhu, S.-C., Tafjord, O., Clark, P., and Kalyan, A. Learn to Explain: Multimodal Reasoning via Thought Chains for Science Question Answering, October 2022. URL <http://arxiv.org/abs/2209.09513>. arXiv:2209.09513 [cs].
- Luo, H., Gu, J., Liu, F., and Torr, P. An image is worth 1000 lies: Transferability of adversarial images across prompts on vision-language models. In *To appear in ICLR*, 2024. URL <https://openreview.net/forum?id=nc5GgFAvtk>.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards Deep Learning Models Resistant to Adversarial Attacks, September 2019. URL <http://arxiv.org/abs/1706.06083>. arXiv:1706.06083 [cs, stat].

- OpenAI, :, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M., Pong, V., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., and Zoph, B. Gpt-4 technical report, 2023.
- Qi, X., Huang, K., Panda, A., Wang, M., and Mittal, P. Visual adversarial examples jailbreak aligned large language models. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*, volume 1, 2023.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. Imagenet large scale visual recognition challenge, 2015.
- Schwenk, D., Khandelwal, A., Clark, C., Marino, K., and Mottaghi, R. A-OKVQA: A Benchmark for Visual Question Answering using World Knowledge, June 2022. URL <http://arxiv.org/abs/2206.01718>. arXiv:2206.01718 [cs].
- Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung, H. W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., et al. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M.-A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open Foundation and Fine-Tuned Chat Models, July 2023. URL <http://arxiv.org/abs/2307.09288>. arXiv:2307.09288 [cs].
- Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, January 2023. URL <http://arxiv.org/abs/2201.11903>. arXiv:2201.11903 [cs].
- Wu, S., Shen, E. M., Badrinath, C., Ma, J., and Lakkaraju, H. Analyzing chain-of-thought prompting in large language models via gradient-based feature attributions. *arXiv preprint arXiv:2307.13339*, 2023.
- Zhang, Z., Zhang, A., Li, M., and Smola, A. Automatic Chain of Thought Prompting in Large Language Models, October 2022. URL <http://arxiv.org/abs/2210.03493>. arXiv:2210.03493 [cs].
- Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G., and Smola, A. Multimodal Chain-of-Thought Reasoning in Language Models, February 2023. URL <http://arxiv.org/abs/2302.00923>. arXiv:2302.00923 [cs].
- Zhao, Y., Pang, T., Du, C., Yang, X., Li, C., Cheung, N.-M., and Lin, M. On evaluating adversarial robustness of large vision-language models. *arXiv preprint arXiv:2305.16934*, 2023.
- Zhu, D., Chen, J., Shen, X., Li, X., and Elhoseiny, M. MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models, April 2023. URL <http://arxiv.org/abs/2304.10592>. arXiv:2304.10592 [cs].
- Zou, A., Wang, Z., Kolter, J. Z., and Fredrikson, M. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

Appendix

A. Implementation Details

A.1. Attack Algorithm Pseudo code

The algorithm for the entire pipeline is outlined in Algorithm 1. In this algorithm:

- $f_{gen}(\cdot)$ represents the model’s inference using the *generate*(\cdot) method.
- $f_{fw}(\cdot)$ signifies the model’s inference using the *forward*(\cdot) method.
- D is the perturbation constraint.
- The initial adversarial image is created by introducing Gaussian noise to the original image.
- Regular updates to the prediction are essential to alleviate the performance gap between the *forward*(\cdot) and *generate*(\cdot) methods.

Algorithm 1 Pipeline

Input: original image v_{org} , question q , boundary ϵ , step α , maximum iteration n

prediction $t_{clean} = f_{gen}(v_{org}, q)$

initial adversarial image v_{adv}

truncate adversarial image to fit $\mathcal{D}(v_{org}, v_{adv}) \leq \epsilon$

for $i = 1$ **to** $n - 1$ **do**

 adversarial output $t_{adv} = f_{fw}(v_{org}, q, t_{pred})$

 loss calculation with $\mathcal{L}(t_{adv}, t_{pred})$

 grad of v_{adv} from loss

 new adversarial image $v_{adv} = v_{adv} + \alpha * \text{sign}(\text{grad})$

 check and truncate $\mathcal{D}(v_{org}, v_{adv}) \leq \epsilon$

if update prediction is *true* **then**

 prediction $t_{clean} = f_{gen}(v_{org}, q)$

end if

if stop criteria satisfied **then**

break

end if

end for

$t_{adv} = f_{gen}(v_{adv}, q)$

save v_{adv}

return t_{adv}

A.2. Attack Methods

The specific perturbation loops for the three attacks are depicted in Figure 11 individually. As the answer attack is a common element in all three attacks, the rationale attack and the stop-reasoning attack distinguish themselves in their approach to attacking the rationale part. The method employed to generate output is *forward*(\cdot) in the perturbation loop.

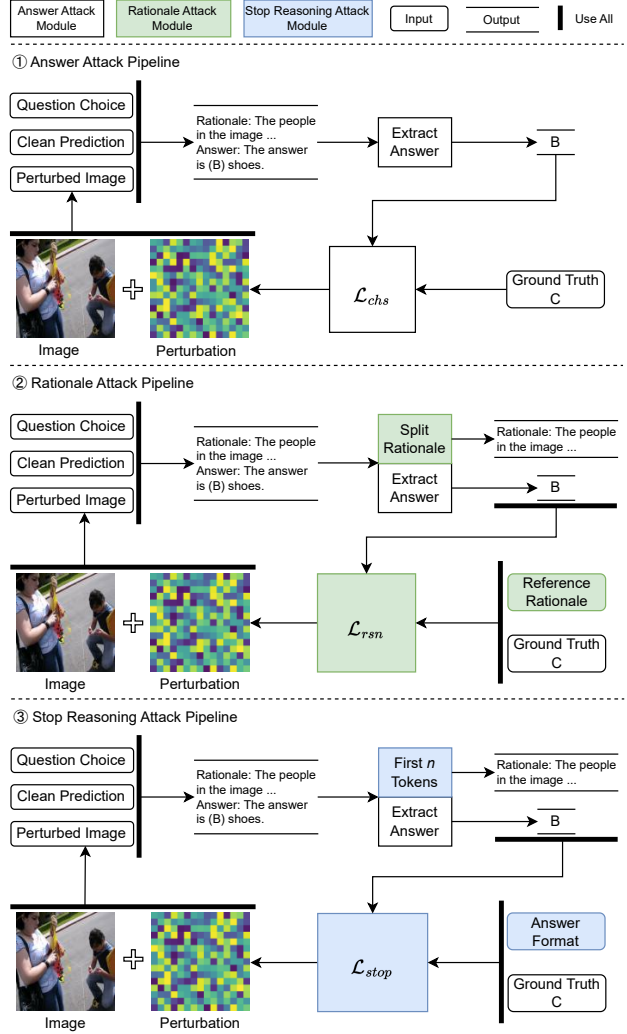


Figure 11. Pipeline visualization of the three attacks.

A.3. Extract Answer

To perform an exact attack, the model is prompted to answer the multiple-choice questions in a specific form and explicitly show the answer choice. As shown in Figure 12 (a), only the choice letter in a correctly formatted sentence will be considered as the answer. The choice content and choices appearing in other forms will not be accepted.

A.4. Split Rationale and Answer

To perform the rationale attack, the rationale and answer parts in the output logit matrix should be split if the model answers the question with the CoT process. As the used LLMs are all generative models, it is not deterministic where the rationale is, where the answer is, and how long each is. So, the output logit matrix is decoded and split into sub-

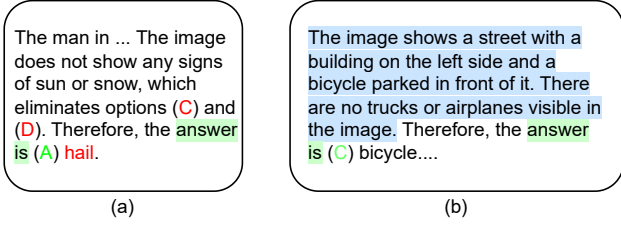


Figure 12. (a) Extract Answer. Only the choice letter (green) in the specified sentence form (green background) will be considered as the answer. Other choice letters or choice content (red) will be ignored. (b) Split Rationale. Only the sentences (blue) before the answer extracted (green) will be contoured as the rationale.

sentences first. As the model is imperfect and the instruction prompt is not strong enough, the model may not follow the instructions exactly. The inference may mix the answer and the rationale part. The part before the sentence the answer is extracted from is the rationale. As shown in Figure 12 (b), inferences can be roughly divided into two parts from the answer. The sentences from the answer are regarded as the answer part of the model, even though there are some other sentences. The sentences before the answer belong to the rationale part. The corresponding logit matrix will be extracted.

A.5. Stop Criteria

The general stop criteria shared in all attack scenarios is whether the inferred answer is wrong in the perturbation loop. The perturbation process will be stopped if the answer is wrong. Then, the perturbed image will be fed into the model again to infer the final answer with the *generate()* method. The stop criteria are combined with a stop check for the stop-reasoning attack, which checks if the answer is extracted from the first sentence.

A.6. Forward vs. Generate

In the context of language models, the *forward()* function often refers to the process of passing input data through the model to obtain predictions or activations. For LLMs used in MLLMs like Llama2 (Touvron et al., 2023), the *forward()* function has the same length in output as the input. The output token is the predicted next token to the input token at the same position. The *generate()* function generates output by iteratively using the *forward()* function. Specifically, in each iteration, only the last token, the next token to the whole input, is extracted and concatenated after the input sequence. The new sequence will be used as input in the next iteration until there is an end-of-sequence token [EOS]. To generate an output, the *generate()* function costs much more time than the *forward()*, if a ground truth

can be provided to the *froward()* function, because the *generate()* function goes through the *forward()* function several times while the *forward()* with ground truth needs only one iteration.

However, it’s important to note that to generate all tokens at once, the *forward()* method requires a ground truth, and there is a performance gap between the *forward()* and the *generate()* functions.

As outlined in appendix A.1, we adopt the clean prediction as the pseudo ground truth for the *forward()* method. As the perturbation progresses, the input image differs from the one used for the clean prediction. Consequently, the pseudo ground truth deviates from the actual ground truth, leading to a divergence in the adversarial output. To address the disparity between the real adversarial output and the actual adversarial output, the clean prediction should be updated with the latest adversarial image every several iterations.

B. Attack Method Parameter Settings

Across all attack scenarios, the perturbation constraint ϵ is set to 16. The maximum number of attack iterations is capped at 200. The prediction is updated every 10 iterations to mitigate the gap between the *forward()* method and the *generate()* method. In every attack test, all victim models use a 0-shot prompt to output their final answer. Every attack method starts with a random perturbation on the image in the very first iteration, then follows its individual loss function and uses PGD method to generate a new perturbed image for the next iteration. All the attacks are performed on a single NVIDIA 40G A100 GPU. To measure the robustness of the MLLMs, we employ *accuracy* as the performance metric. Low accuracy indicates a low robustness.

C. ImageNet Subclasses

We create classification tasks by extracting 4, 8, and 16 subclasses. The specific subclasses selected for each scenario are as follows:

- **4 classes:** English setter, Persian cat, school bus, pineapple.
- **8 classes:** bison, howler monkey, hippopotamus, chow, giant panda, American Staffordshire terrier, Shetland sheepdog, Great Pyrenees.
- **16 classes:** piggy bank, street sign, bell cote, fountain pen, Windsor tie, volleyball, overskirt, sarong, purse, bolo tie, bib, parachute, sleeping bag, television, swimming trunks, measuring cup.

All tasks had a uniform question: “What is the class of the image?”

Table 4. Accuracy table with reduced adversarial capability.

MODEL	DATASET	w/o CoT		WITH CoT		
		W/O ATTACK	ANSWER ATTACK	ANSWER ATTACK	RATIONALE ATTACK	STOP REASONING ATTACK
MINIGPT4	A-OKVQA	61.38	0.96	17.59	30.98	2.68
	SCIENCEQA	66.28	3.12	25.55	47.66	16.93
OPEN-FLAMINGO	A-OKVQA	34.80	4.19	13.94	11.73	6.47
	SCIENCEQA	34.55	7.13	39.18	40.51	31.71

D. Ablation Study

During the ablation study, the adversarial capability is reduced by narrowing the limited boundary (ϵ) to 8 (as described in Section 3.1). Table 4 presents results consistent with Table 1, indicating that the CoT reasoning process enhances the robustness of MLLMs. Furthermore, the table shows that the stop-reasoning attack remains the most effective method in compromising this increased robustness.