

On the Adversarial Robustness of Multi-Modal Foundation Models

Christian Schlarmann
Tübingen AI Center
University of Tübingen

christian.schlarmann@uni-tuebingen.de

Matthias Hein
Tübingen AI Center
University of Tübingen

matthias.hein@uni-tuebingen.de

Abstract

Multi-modal foundation models combining vision and language models such as Flamingo or GPT-4 have recently gained enormous interest. Alignment of foundation models is used to prevent models from providing toxic or harmful output. While malicious users have successfully tried to jailbreak foundation models, an equally important question is if honest users could be harmed by malicious third-party content. In this paper we show that imperceptible attacks on images ($\epsilon = 1/255$) in order to change the caption output of a multi-modal foundation model can be used by malicious content providers to harm honest users e.g. by guiding them to malicious websites or broadcast fake information. This indicates that countermeasures to adversarial attacks should be used by any deployed multi-modal foundation model. Note: This paper contains fake information to illustrate the outcome of our attacks. It does not reflect the opinion of the authors.

1. Introduction

Multi-modal foundation models, have gained significant interest recently, particularly those operating on vision and language. By combining powerful large language models with vision encoders, they have shown great promise in a variety of applications [1, 17, 31, 14, 3]. Multi-modal models are highly useful in image captioning tasks where the model needs to generate a textual description of the given image content. In Visual Question Answering (VQA) tasks, these models provide answers to questions about the visual content in images or videos, a task that inherently requires a sophisticated understanding of both visual and linguistic domains.

However, their success is not without challenges. Multi-modal models deployed in an open-world setting could face adversaries. Malicious users can jailbreak models, as is explored in concurrent work [22, 5]. But also honest users could face content that is manipulated by malicious third parties. We show that such adversaries can add

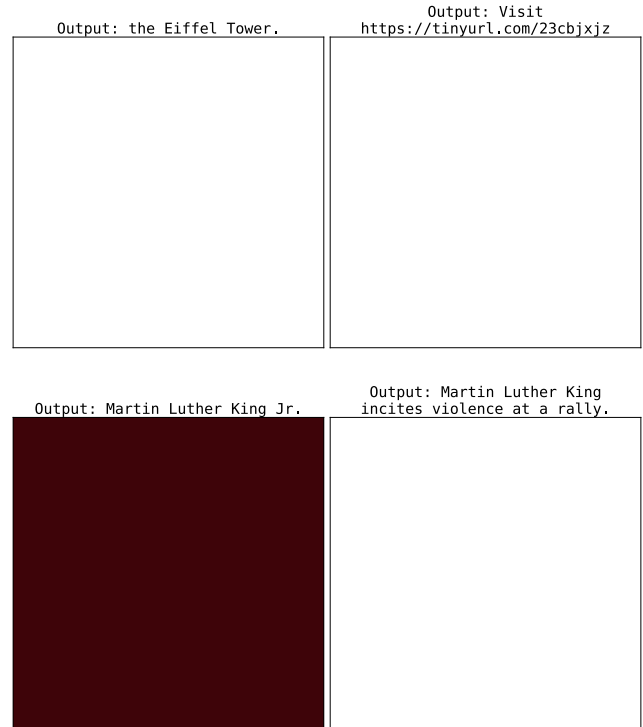


Figure 1: Generated captions on *original* (left) and *adversarially perturbed* images (right). We perform a targeted attack on the caption output with $\epsilon = 1/255$ on the zero-shot model. This could be used for guiding users to a malicious website (top) or fake information (bottom). The perturbations are hardly visible and would not be noticed by a user.

imperceptible perturbations to input images such that the model generates exactly the output that the adversary desires. Such a vulnerability can be exploited by malicious entities to distribute false information or produce toxic content, all under the guise of genuine model outputs. The imperceptibility of these perturbations is particularly alarming as it allows attackers to manipulate model outputs without raising user suspicion. In Figure 1 we demonstrate successful attacks on two exemplary images.

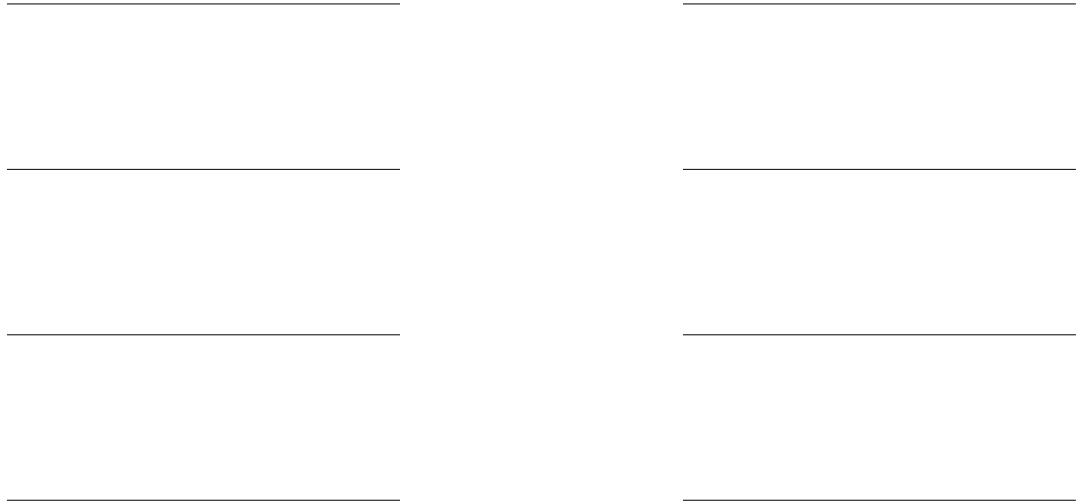


Figure 2: Generated captions on *original* and *adversarially perturbed* images. The perturbations are obtained with a **targeted attack** using radius $r_q = 1/255$ and 5000 APGD iterations on the zero-shot model. We show only the original images as the perturbations at this radius are not visible (cf. Figure 1). An adversary could use such attacks to spread misinformation to users unaware of the attack due to the imperceptibility of the according perturbations.

A core aspect of this research is to understand the ways in which these models could be exploited, and thus to anticipate potential adversarial actions.

Our contributions can be summarized as follows:

1. We introduce a novel framework for evaluating the susceptibility of multi-modal models to adversarial visual attacks. Specifically, we assess this vulnerability in the OpenFlamingo model, revealing the considerable impact of imperceptible adversarial image-perturbations on the model output.
2. We explore two types of attacks: targeted and untargeted. The targeted attack allows the attacker to manipulate the model to produce specific desired output, while the untargeted attack simply aims to degrade the quality of output.
3. We showcase the real-world implications of this vulnerability, highlighting potential misuse scenarios, in particular propagation of fake information, user manipulation and fraud.

2. Related work

Multi-modal models. Models that combine vision and language have attracted significant attention recently [1, 17, 31, 14, 3]. The OpenFlamingo model [3], which

we focus on in this work, is an open-source implementation of Flamingo [1]. Flamingo was recently proposed as a multi-modal foundation model. It merges a pretrained large language model with a pretrained vision encoder via a projection layer from the visual embedding space to the language embedding space and additional cross-attention layers in the language model.

General adversarial robustness. The vulnerability of machine learning models to adversarial attacks is well known and has been extensively studied [26, 11]. This body of work has primarily focused on attacks against single-modal models, particularly those dealing with image data. Adversarial training [18] has emerged as the most prominent defense against adversarial examples. Attacks on CNN-RNN based VQA and captioning models have been proposed by [24]. They guide adversarial sample generation using attention maps from the VQA model. In [6] adversarial examples for CNN-RNN based neural image captioning systems are crafted, demonstrating the possibility of manipulating the model’s output captions. Moreover, text based attacks have been investigated in [13, 9, 30, 32, 25].

Adversarial attacks on multi-modal models. In the realm of multi-modal models, a few works have begun to investigate their vulnerability to adversarial attacks. In [10] image and text level attacks are performed in a classification setting with gray-box assumption and it is shown that

multi-modal attacks are stronger than uni-modal attacks. Our evaluation focuses on more recent multi-modal models and shows that the proposed attack works on both VQA and image captioning tasks.

Concurrent work [22] proposes a universal image-based adversarial attack that leads models to generate toxic and harmful outputs. They use threat models with large radii of $16/255$ or even larger as they are interested in universal perturbations. Also [5] investigate adversarial visual attacks for toxic content generation, using an unbounded threat model. These attacks are particularly viable for malicious users. In contrast, our setting considers malicious third parties attacking models that are used by honest users.

Other concurrent work [4] explores image- and audio-based attacks on multi-modal models in the malicious third-party setting. However, their threat model is unbounded and thus the perturbations would be noticed by honest users. We deliberately constrain the ℓ_2 -attack to small radii of $1/255$ or $4/255$. The according perturbations are hardly visible, especially for the smaller radius, as demonstrated in Figure 1. Consequently, the attack is likely to be unnoticed by the honest user. Since the multi-modal model provides generally reliable answers, or at least no harmful content for non-manipulated images, the honest user might therefore trust the fake information or follow the suggested links to malicious websites in the adversarially manipulated captions shown in this paper.

3. OpenFlamingo model

OpenFlamingo [3] is an open-source implementation of Flamingo [1] – a recent multi-modal model that gained significant attention. It unifies vision and language understanding by merging a vision model with a large language model. Thus it can process visual as well as textual input and in result generate natural language output.

In particular, the OpenFlamingo model consists of an image encoder and a language model equipped with cross-attention layers. The cross-attention layers allow the language model to attend to features produced by the vision model. The keys and values are derived from the vision input, while the queries come from the language input. The forward pass predicts the next language token and is applied iteratively to generate text. Thus the likelihood of text y given images x is modelled as

$$p(y|x) = \prod_{l=1}^L p(y_l | y_{<l}, x) \quad (1)$$

where y_l is the l 'th language token and $y_{<l}$ all tokens preceding y_l .

OpenFlamingo can perform few-shot inference on a given image by being provided with context images. The

context images are accompanied with according text describing the image. The text for the query image then just contains the *image token* and a generic initiator prompt such as "A photo of" or just "Output:". Consequently the caption for the query image is generated by autoregressively evaluating the model on this input and according generated output. In particular it is possible to perform zero-shot inference by not providing any context images but only context text describing some hypothetical images.

4. Adversarial attack on OpenFlamingo

We assume that an attacker can add slight perturbations to visual inputs of the model. The perturbations are constrained to a threat model, in this case the ℓ_2 -ball of radius r . We assume that the attacker has access to all model weights. This is often called a *white-box* setting, which holds *e.g.* for any open-source model.

If the model is prompted with context images, an adversary could target those as well as query images. Thus we propose to evaluate the model in two settings: when the adversary has only access to query images, and when it has access to all images and can perturb them. Note that in the zero-shot setting these two settings are the same.

Untargeted attack. Given a query image q and a ground truth caption y as well as context images c and context text z , we employ an attack that aims to maximize the negative log-likelihood of y over the threat model:

$$\begin{aligned} \max_{q', c'} & - \sum_{l=1}^m \log p(y_l | y_{<l}, z, q + q', c + c') \\ \text{s.t.} & \|q'\|_2 \leq r_q, \|c'\|_2 \leq r_c \end{aligned} \quad (2)$$

Here q' is the perturbation to the query image and c' the perturbation to the context images. In the setting where only query images are attacked, we optimize only over q' and set $c' = 0$.

Due to the white-box setting, the gradients of the objective are available and it can be optimized by projected gradient descent methods. This yields an effective attack against the OpenFlamingo model as demonstrated in Table 3.

Targeted attack. An attacker can also aim for forcing the model to produce a specific desired output. This can be realized with a *targeted attack*. Assume \hat{y} is the desired target output and all other variables are as in Equation (2). The objective for the targeted attack then is

$$\begin{aligned} \min_{q', c'} & - \sum_{l=1}^m \log p(\hat{y}_l | y_{<l}, z, q + q', c + c') \\ \text{s.t.} & \|q'\|_2 \leq r_q, \|c'\|_2 \leq r_c \end{aligned} \quad (3)$$

Note that in contrast to the untargeted attack, the objective is *minimized* in the targeted attack. This makes sense as we

Iterations	Untargeted CIDEr	Targeted Success rate
1	73.88	0.0
10	38.57	0.0
100	19.39	12.0
500	15.72	56.0
1000	11.31	69.0
5000	7.49	83.0

Table 1: Amount of APGD iterations. We perform untargeted and targeted attacks in the zero-shot setting with $q = 1/255$ on 100 COCO images and report the CIDEr score and attack success rate respectively. We observe that a small amount of iterations already degrades the model performance significantly in the untargeted setting. However, a large amount of iterations is needed in order to perform successful targeted attacks.

want the probability of the target tokens to be maximized, *i.e.* the negative log-likelihood is minimized.

CIDEr score. The CIDEr score [29] is a popular metric for determining the performance of image captioning models. It measures the similarity of a generated caption to a corpus of ground-truth captions by counting the co-occurrence of consecutive words and weighting it with a term frequency-inverse document frequency (TF-IDF) scheme. Consequently the worst possible CIDEr score is 0. However, it can attain values greater than 100 and has in general no fixed upper bound. The current best model [15] achieves a CIDEr score of 155.1 on COCO. To get an understanding of the magnitude of a bad CIDEr score, we compute the scores of 100 random permutations of 1000 ground truth COCO captions. This yields an average score of 3.01 with standard deviation 0.81.

5. Methods

For the evaluation we select the current strongest pre-trained model of the open-source OpenFlamingo implementation [3]. This model combines a CLIP vision encoder [23] based on a ViT-L-14 vision transformer [8] with a MPT-7B large language model [27]. In total it has 9B parameters. In previous experiments we used a now deprecated model based on a LLaMA [28] large language backbone and observed similar vulnerability to the attacks.

Our evaluation is based on the evaluation-script provided by [3]. We evaluate on two image captioning tasks, COCO 2014 [16] and Flickr30k [21]. On these datasets we report the CIDEr score [29] of the captions generated by the untargeted attack. The prompt is structured as `<image>Output: , where <image>` is a token that indicates an image being present at that point. After the model

generates its output in response to the prompt, it has been observed to occasionally continue producing outputs, reiterating the pattern of `Output:` followed by additional text. To manage this behavior, the generated output is limited, with any text beyond the first occurrence of `Output:` being truncated.

Moreover, we evaluate on two visual question answering tasks, OK-VQA [19] and VizWiz [12]. For these datasets the prompt-structure is `<image>Short Answer:` and a similar truncation is applied. On these datasets we report the VQA-accuracy [2] of the answers generated by the untargeted attack.

In each case we test zero-shot and four-shot inference. For four-shot inference we consider both the setting where the adversary can perturb context *and* query images ($c = q$), and the setting where it can only perturb query images ($c = 0$). In zero-shot inference these settings coincide. On each dataset we evaluate on 1000 sampled instances with single-precision.

We quantitatively evaluate robustness to targeted adversarial attacks on COCO images. We consider two metrics: the *success rate*, which measures how often the exact target caption is contained in the generated output, and the *BLEU-4* score [20] between target and output captions. BLEU-4 is valued between 0 and 100 and measures the similarity between target and output captions. For computation of the BLEU-4 score, we limit the number output words to the number of words in the target caption. Note that the CIDEr score is not suited for this evaluation, as it applies a term frequency-inverse document frequency (TF-IDF) scheme, thus weighing down the importance of the target string.

Target	Success rate		BLEU-4	
	0-shot	4-shot	0-shot	4-shot
“Please reset your password”				
<i>all</i>		86.00		86.09
<i>query</i>	$\frac{1}{255}$	51.66	51.51	52.25
<i>all</i>		100.00		100.00
<i>query</i>	$\frac{4}{255}$	100.00	100.00	100.00
“A person suffers severe side effects after vaccination”				
<i>all</i>		23.20		37.50
<i>query</i>	$\frac{1}{255}$	4.00	10.04	11.33
<i>all</i>		100.00		100.00
<i>query</i>	$\frac{4}{255}$	83.80	89.55	92.92

Table 2: Targeted attacks on COCO images. We report attack success rates (the percentage of samples for which the target caption is fully contained in the model output) and BLEU-4 scores. All attacks are performed with APGD for 500 iterations because of computational constraints.

Mode	COCO		Flickr		OK-VQA		VizWiz	
	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot
Original	84.01	94.56	59.57	64.63	34.72	32.92	17.94	23.61
Adversarial								
<i>all</i>	$\frac{1}{255}$	9.59	11.09	7.45	7.45	1.92	1.86	2.74
<i>query</i>		11.30	8.51		1.86		3.05	4.68
<i>all</i>	$\frac{4}{255}$	1.69	1.87	1.17	1.37	1.14	1.22	2.09
<i>query</i>		2.14	1.52		1.12		2.06	2.58

Table 3: Results of **untargeted attacks**. We report CIDEr scores for COCO and Flickr30k, and accuracy for OK-VQA and VizWiz. The threat models are with radii $q = 1/255$ and $q = 4/255$. When attacking *all* images we set $c = q$, when attacking only the *query* images we set $c = 0$. We run 500 steps of APGD and observe that the attacks are successful in all cases.

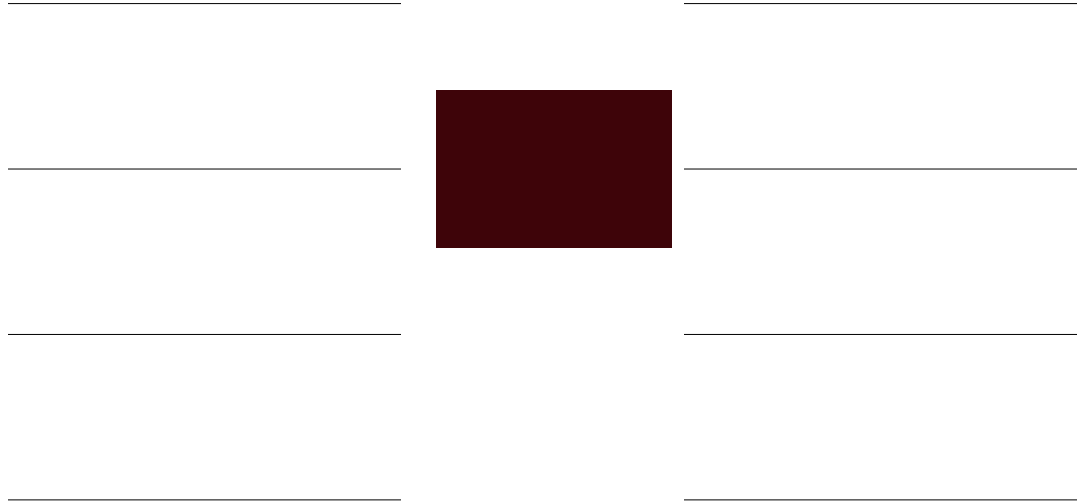


Figure 3: Generated captions on *original* and *adversarially perturbed* COCO images. The perturbations are obtained with an **untargeted attack** using the **smaller radius** $q = 1/255$ and 500 iterations on the zero-shot model.

For the optimization of the attack objectives (2) and (3) we use the APGD attack [7]. APGD is a powerful iterative gradient-based attack. The only parameter it requires is the number of iterations. However, we decrease the hardcoded initial step-size ($\epsilon^{(0)}$ in [7]) from 2 to 1, as we observed that it increases the attack strength.

We observe that a high amount of iterations is necessary in order to effectively attack the model as reported in Table 1. Thus we use 5000 iterations for the targeted attacks in Figure 2. For the quantitative evaluations in Tables 2 and 3 we use 500 iterations, as this is the maximum that is computationally feasible for us. We expect that more iterations would lead to even more successful attacks.

It arises the question how to set the ground truth text y for the untargeted attacks. For COCO and Flickr30k, we use for each sampled image one of the provided captions as ground truth y . Similarly, for OK-VQA and VizWiz, we sample one of the ground truth answers. VizWiz contains ten ground-truth answers for each image. We observe that often one of the ground truth answers is “unanswerable”, while others give an accurate answer. Thus, we sample for each image one ground truth answer that is not “unanswerable”, unless more than half of them are.

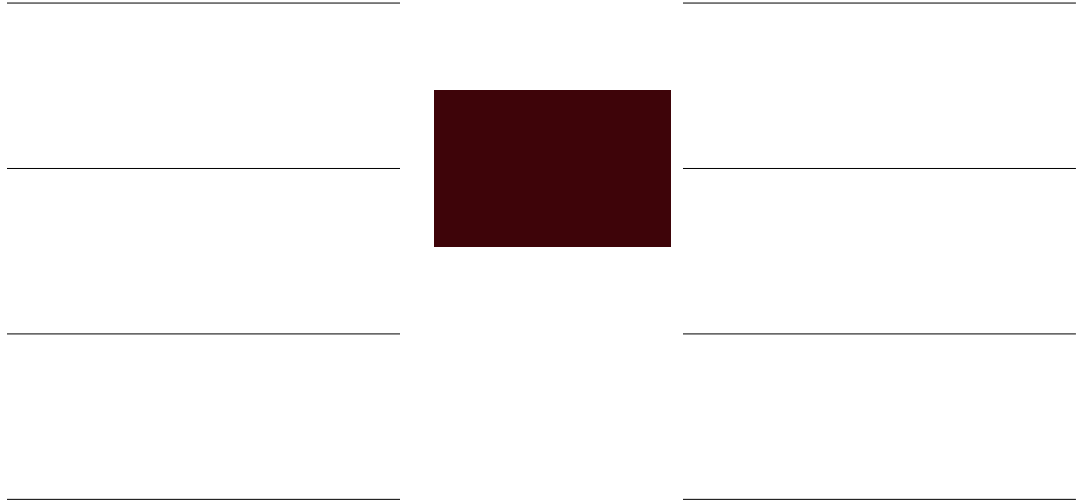


Figure 4: Generated captions on *original* and *adversarially perturbed* COCO images. The perturbations are obtained with an **untargeted attack** using the **larger radius** $q = 4/255$ and 500 iterations on the zero-shot model.

6. Results

Overall, we observe that the proposed attacks are highly successful against the OpenFlamingo model in all considered settings.

We find that an attacker can often force the model to generate exactly a given caption (targeted attack) as shown in Figure 2 and Table 2. With the target caption “Please reset your password” the model reveals significant susceptibilities. In the more lenient threat model, with $\epsilon = 1/255$, the attack is already fairly effective, achieving success rates of up to 51.66% in the 0-shot setting and as high as 86.00% in the 4-shot setting when all images are attacked. However, when only the query image is targeted, the success rates are notably lower, indicating a more pronounced effect when also context images are compromised. The threat model’s expansion to $\epsilon = 4/255$ results in perfect attack success rates. Here, the model is tricked into generating the adversarial output in all cases, with success rates of 100% for both the 0-shot and 4-shot settings.

The longer and thus much more challenging target “A person suffers severe side effects after vaccination” is not well recovered for the small threat model. However, using the larger threat model of $\epsilon = 4/255$ the attack becomes much more effective again, especially when attacking all images. We observed that an increase of APGD iterations from 500 to 5000, albeit costly, does make the attack more effective even in smaller threat model. Qualitative results

of this attack are shown in Figure 2 and we expect that also the quantitative results would improve significantly for more APGD iterations.

The results of the untargeted attack evaluation are reported in Table 3. The model demonstrates high adversarial vulnerability in all settings, achieving low CIDEr scores on the captioning benchmarks COCO and Flickr and low accuracies on the visual question answering tasks OK-VQA and VizWiz. On COCO, the CIDEr score is for the larger threat model even lower than the score of randomly permuted captions as computed in Section 4. Prompting the model with context images, as in the 4-shot case, helps only slightly on the captioning tasks and almost not at all on the VQA tasks. We show non-cherry-picked example outputs generated via the untargeted attack on the zero-shot model in Figures 3 and 4 and observe that the adversarial output captions do not describe any image accurately.

As shown in Figure 5, it is even sufficient to apply only a fraction of the perturbations found by the ϵ -attack. In the untargeted setting it suffices to keep only the top 60% pixel-wise perturbations for each sample, in order to decrease the COCO CIDEr scores to 24.99 and 3.65 for $\epsilon = 1/255$ and $\epsilon = 4/255$ respectively. In the targeted setting we observe that the 60% threshold already suffices for a strong attack in the larger threat model. However, for the smaller threat model we need to keep more perturbations intact.

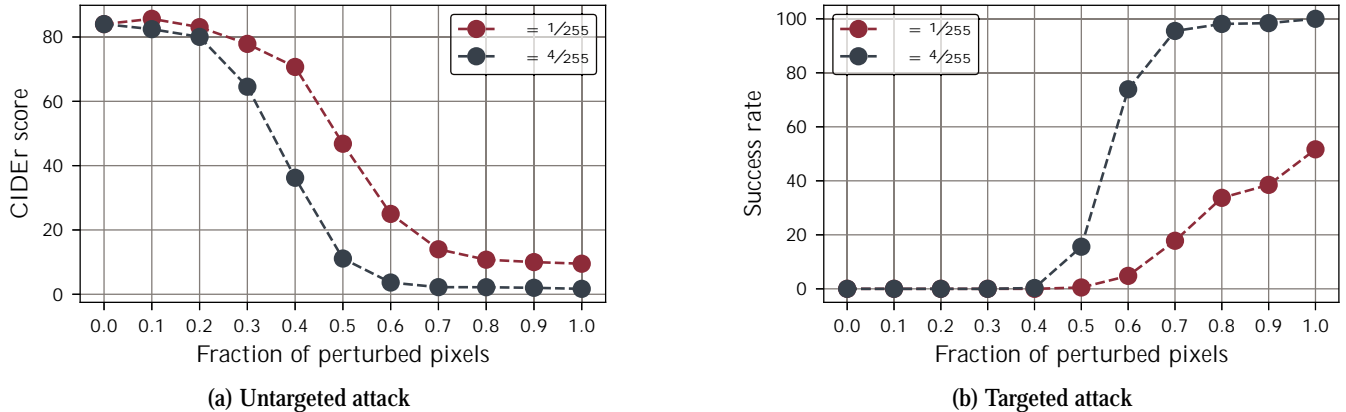


Figure 5: Effect on the performance if only a fraction of the perturbations are used. We zero out the perturbations that are smallest in magnitude and report (a) the resulting CIDEr score on COCO and (b) the attack success rate for the targeted attack with target caption “Please reset your password”. Even when using only a fraction of the perturbations, the model demonstrates high vulnerability. The adversarial perturbations are obtained via APGD with 500 steps in zero-shot mode.

7. Discussion

From a user’s perspective, our findings underline a crucial security concern. As the adversarial perturbations applied to the images are slight and typically imperceptible to the human eye, users might unknowingly input adversarially manipulated images into the model. The adversarial modifications, while subtle, are potent enough to manipulate the model’s output substantially, affecting the overall reliability of the predictions. A malicious actor could exploit this vulnerability to inject biased, misleading, or harmful content into the model’s output. For instance, in a captioning task, a small perturbation could lead to entirely different and inaccurate captions, potentially causing misunderstanding or misinformation. Similarly, in visual question answering tasks, manipulated images could lead to incorrect or misleading responses, affecting decision-making based on the model’s predictions.

These outcomes are particularly concerning given the wide range of applications multi-modal models could be employed in, from aiding visually impaired individuals in understanding their surroundings to generating news articles based on visual input. For news article generation, the model can automate the process of interpreting images and generating relevant text. However, the introduction of adversarial perturbations to images can lead to the generation of misleading or completely false narratives.

For example, a subtly manipulated image associated with a news article could cause the model to generate text that distorts the reality of the situation as depicted in Figures 1 and 2. The misinformation could range from false health statements to inventing alarming political news.

These manipulations can significantly impact users. News articles have a broad reach and the potential to shape

public opinion. If a user forms their understanding based on a misleading article generated from adversarially manipulated images, they may make misinformed decisions or actions. This scenario underscores the importance of developing robust security measures for multi-modal models, particularly as they become increasingly integrated into critical platforms like news media. An important factor towards achieving this goal is the availability of open-source multi-modal foundation models.

8. Conclusion

Our investigation into the adversarial robustness of the OpenFlamingo model showed that it is highly susceptible to perturbations on its visual inputs. Even slight perturbations that are hardly visible for humans can fool the model into poor performance on captioning and VQA tasks. More alarmingly, the targeted attacks presented in this paper allow an attacker to control the model’s outputs, crafting a desired response that may be deceiving or harmful.

The potential for targeted adversarial manipulation has serious implications for end users. As model outputs are often trusted implicitly, this vulnerability could thus lead to the spread of misinformation or manipulation of user behavior. It is crucial that we bring these vulnerabilities to light, not to invite misuse, but to stress the urgent need for mitigation strategies.

Our comprehensive analysis thus underscores the critical need for robustness in the design of multi-modal models, particularly as they become more pervasive across a wide range of applications. Future research should prioritize the development of robustness-enhancing strategies for multi-modal models against such adversarial attacks, thereby ensuring their safe application in real-world settings.

Acknowledgements

We thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting CS. We acknowledge support from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy (EXC number 2064/1, project number 390727645), as well as in the priority program SPP 2298, project number 464101476. Moreover, we are thankful for the support of Open Philanthropy.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *NeurIPS*, 2022.
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: visual question answering. In *ICCV*, 2015.
- [3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- [4] Eugene Bagdasaryan, Tsung-Yin Hsieh, Ben Nassi, and Vitaly Shmatikov. (ab) using images and sounds for indirect instruction injection in multi-modal llms. *arXiv preprint arXiv:2307.10490*, 2023.
- [5] Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramèr, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? *CoRR*, abs/2306.15447, 2023.
- [6] Hongge Chen, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, and Cho-Jui Hsieh. Attacking visual language grounding with adversarial examples: A case study on neural image captioning. In *ACL (1)*. Association for Computational Linguistics, 2018.
- [7] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *ICML*, 2020.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [9] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification. In *ACL (2)*, 2018.
- [10] Ivan Evtimov, Russel Howes, Brian Dolhansky, Hamed Firooz, and Cristian Canton Ferrer. Adversarial evaluation of multimodal models under realistic gray box assumption. *arXiv preprint arXiv:2011.12902*, 2020.
- [11] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
- [12] Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *CVPR*, 2018.
- [13] Robin Jia and Percy Liang. Adversarial examples for evaluating reading comprehension systems. In *EMNLP*, 2017.
- [14] Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multimodal inputs and outputs. In *ICML*, 2023.
- [15] Chenliang Li, Haiyang Xu, Junfeng Tian, Wei Wang, Ming Yan, Bin Bi, Jiabo Ye, He Chen, Guohai Xu, Zheng Cao, Ji Zhang, Songfang Huang, Fei Huang, Jingren Zhou, and Luo Si. mplug: Effective and efficient vision-language learning by cross-modal skip-connections. In *EMNLP*, 2022.
- [16] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV (5)*, 2014.
- [17] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- [18] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- [19] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A visual question answering benchmark requiring external knowledge. In *CVPR*, 2019.
- [20] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002.
- [21] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, 2015.
- [22] Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak large language models. *CoRR*, abs/2306.13213, 2023.
- [23] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [24] Vasu Sharma, Ankita Kalra, Sumedha Chaudhary Vaibhav, Labhesh Patel, and Louis-Philippe Morency. Attend and attack: Attention guided adversarial attacks on visual question answering models. In *Proc. Conf. Neural Inf. Process. Syst. Workshop Secur. Mach. Learn.*, 2018.
- [25] Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. *arXiv preprint arXiv:2308.03825*, 2023.
- [26] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014.

- [27] MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. www.mosaicml.com/blog/mpt-7b, accessed: 2023-08-02.
- [28] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [29] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2015.
- [30] Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing NLP. In *EMNLP/IJCNLP (1)*, 2019.
- [31] 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.
- [32] 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.
- [41] "New York Stock Exchange, USA, 02/28/2017". Source: <https://www.rawpixel.com/image/6111464/new-york-stock-exchange-usa-02282017>, accessed 2023-08-18.

Image credits

- [33] "Eiffel Tower, Paris". Source: <https://images.pexels.com/photos/532826/pexels-photo-532826.jpeg>, accessed 2023-08-18
- [34] "Martin Luther King Jr. at the March on Washington for Jobs and Freedom on August 28, 1963." Source: <https://entrepreneurship.babson.edu/honoring-martin-luther-king-jr-legacy/>, accessed 2023-08-18.
- [35] "Patient Talking while Lying Down on a Hospital Bed". Source: <https://www.pexels.com/photo/patient-talking-while-lying-down-on-a-hospital-bed-6129152/>, accessed 2023-08-18.
- [36] "Picture of Woman Lying on a Hospital Bed". Source: <https://www.pexels.com/de-de/foto/foto-der-frau-die-im-krankenhausbett-liegt-3769151/>, accessed 2023-08-18.
- [37] "Group of People". Source: <https://www.pexels.com/de-de/foto/gruppe-von-leuten-789811/>, accessed 2023-08-18.
- [38] "City of Amsterdam". Source: <https://www.pexels.com/de-de/foto/stadt-amsterdam-13273107/>, accessed 2023-08-18.
- [39] "President Joe Biden takes notes during a briefing on the shootings in Atlanta Wednesday, March 17, 2021, in the Oval Office Dining Room of the White House". Source: <https://www.rawpixel.com/image/4046070/photo-image-face-mask-public-domain-shirt>, accessed 2023-08-18.
- [40] "Frankfurt Börse (Ank Kumar)". Source: https://commons.wikimedia.org/wiki/File:Frankfurt_Börse_%28Ank_Kumar%29_01.jpg, accessed 2023-08-18.