

# VEAttack: Downstream-agnostic Vision Encoder Attack against Large Vision Language Models

Paper ID: 363

Hefei Mei<sup>1</sup>, Zirui Wang<sup>1</sup>, Shen You<sup>1</sup>, Minjing Dong<sup>1</sup>, Chang Xu<sup>2</sup>

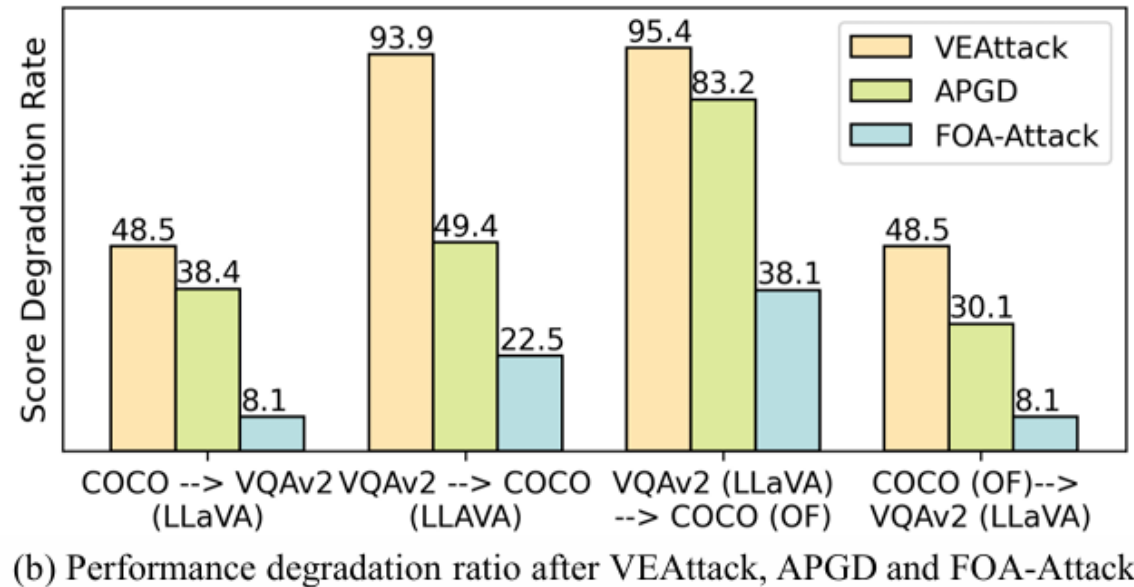
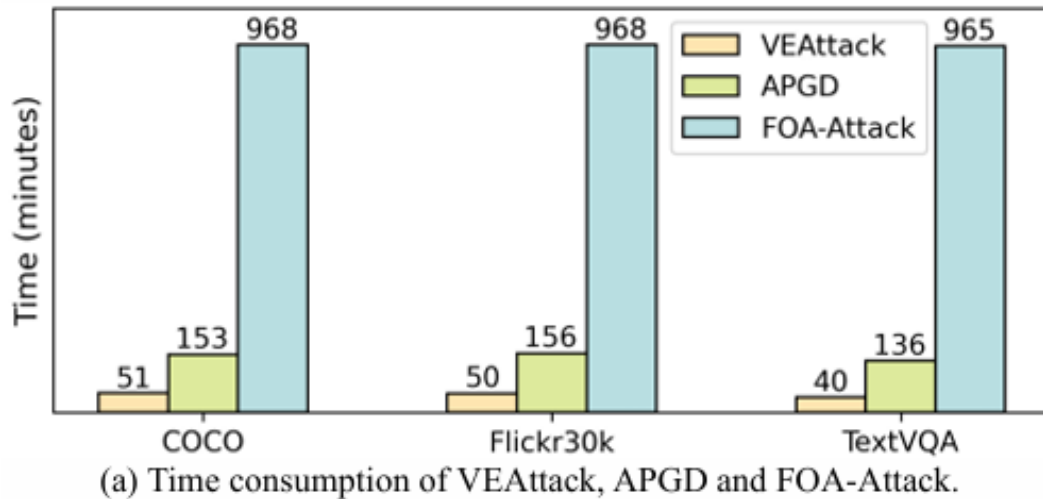
<sup>1</sup>City University of Hong Kong

<sup>2</sup>University of Sydney

hefeimei2-c@my.cityu.edu.hk

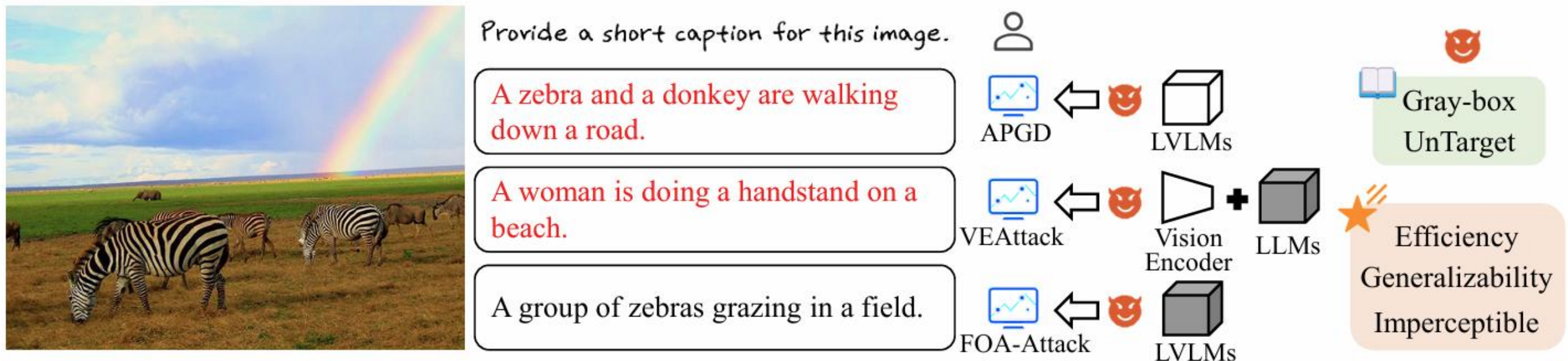
# Background

- LVLMs typically serve as foundation models for diverse multimodal tasks.
- Two stubborn challenges for attacking LVLMs lie in **efficiency** and **task generalization**.
- Given that traditional adversarial attacks focus on **specific tasks** and the **increasing parameter scale** of LVLMs, VEAttack aims to alleviate **task-specific dependency** and **reduce the computational overhead**.
- Black-box attack FOA-Attack **struggles** to achieve a **high** score reduction rate (SRR) of LVLM performance, despite employing **larger** perturbations.



# Motivation

- A **gray-box** attack targeting a foundational LVLM **module** provides better trade-offs among efficiency, transferability, and practicality since only **partial parameters are required** and overfitting to a specific LVLM task is relaxed.
- Drawing from traditional vision tasks, where a strong vision backbone enhances downstream task success, and recognizing the critical role of the **vision encoder** in LVLMs for overall performance, we target the vision encoder of LVLMs for attack.



# Redefine adversarial objective

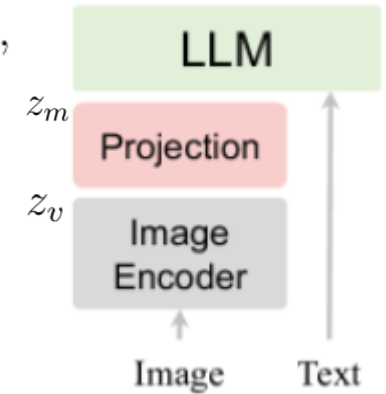
## ➤ White-box Attack:

$$\max_{\|\delta\| \leq \epsilon} \mathcal{L}(v + \delta, t; \theta) = \max_{\|\delta\| \leq \epsilon} -\log p(y | f_A(f_V(v + \delta)); \text{tokenizer}(t); \theta),$$

multi-task  $T = \{T^1, \dots, T^m\}$

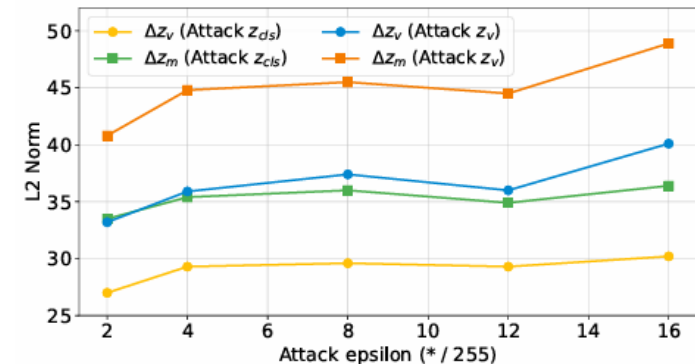
$$\max_{\|\delta^{T^i}\| \leq \epsilon} \mathcal{L}(\tilde{v}^{T^i}, t^{T^i}; \theta) = \max_{\|\delta^{T^i}\| \leq \epsilon} -\log p(y^{T^i} | f_A(f_V(\tilde{v}^{T^i})); \text{tokenizer}(t^{T^i}); \theta),$$

**Proposition 2** For LLaVa [57] with a linear alignment layer, let  $\Delta z_v = \tilde{z}_v - z_v$  denote the difference between the image tokens output by the vision encoder CLIP before and after the perturbation,  $\|\Delta z_v\|_F \geq \Delta$ ,  $W_a$  is the weight of projection layer,  $\sigma_{\min}$  denotes the minimum singular value of  $W_a$ . Assume  $\sigma_{\min}(W_a) > 0$ , then the difference between aligned features  $z_m = f_A(z_v)$  and  $\tilde{z}_m = f_A(\tilde{z}_v)$  for downstream LLMs will satisfy  $\|\Delta z_m\|_F \geq \sigma_{\min}(W_a)\Delta > 0$ .

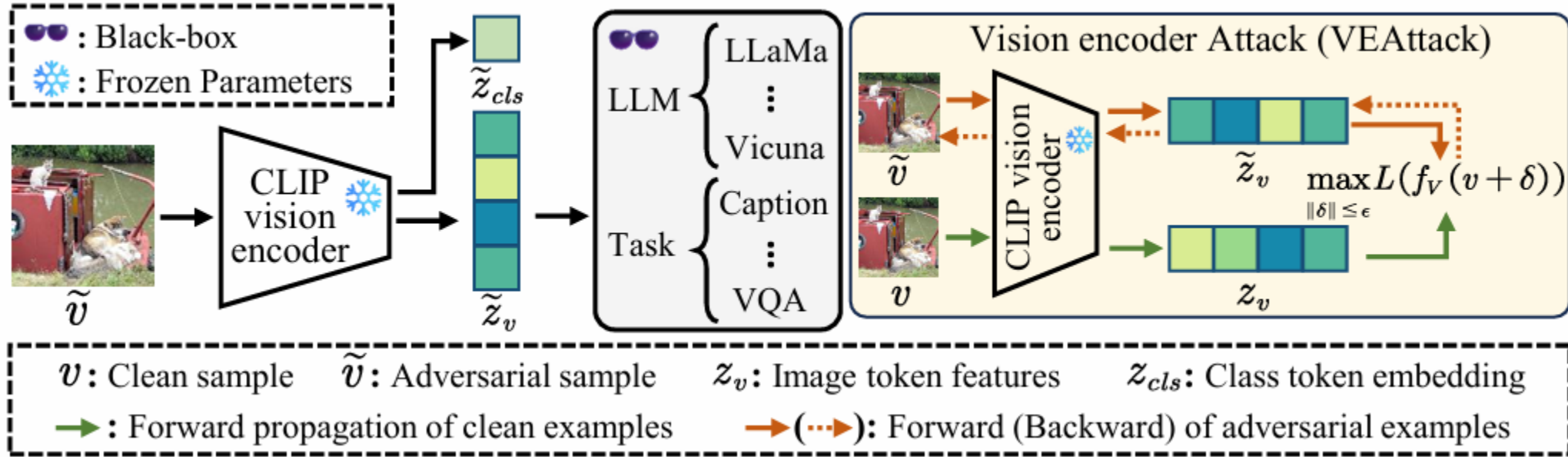


## ➤ Redefined objective:

$$\max_{\|\delta^{T^i}\| \leq \epsilon} \mathcal{L}(v^{T^i} + \delta^{T^i}, t^{T^i}; \theta) \longrightarrow \max_{\|\delta\| \leq \epsilon} \mathcal{L}(v + \delta; \theta_V) = \max_{\|\delta\| \leq \epsilon} \mathcal{L}(f_{CLIP}(v + \delta)),$$



# Framework



- we can treat the attacks in adversarial training algorithms proposed by black-box attack and robust CLIP work:

$$\text{AttackVLM-ii : } \tilde{v} = \arg \max_{\|\delta\| \leq \epsilon} -\cos(\tilde{z}_{cls}, z_{cls}) \quad L_2 \text{ Attack : } \tilde{v} = \arg \max_{\|\delta\| \leq \epsilon} \|\tilde{z}_{cls} - z_{cls}\|_2^2,$$

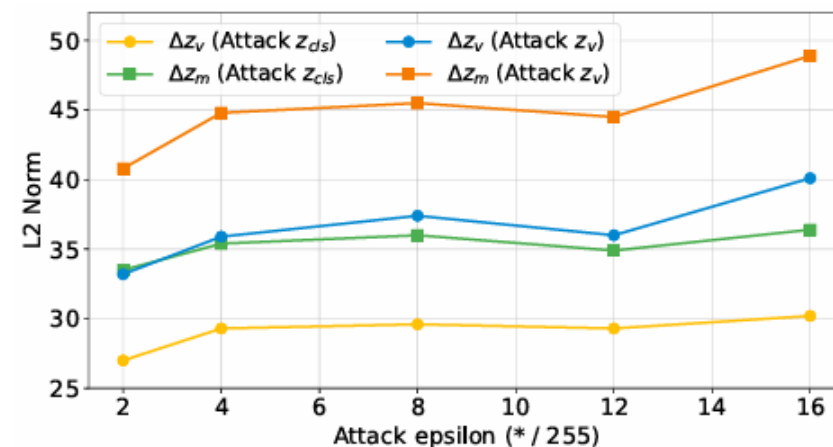
# Vision encoder attack

**Proposition 3** Consider two kinds of attack targets: 1) If the perturbation is introduced to the single class token  $z_{cls}$  and propagates through  $z_{cls} \xrightarrow{\text{backward}} v \xrightarrow{\text{forward}} z_v$ . Let  $\Delta z_{cls}(z_{cls})$ ,  $\Delta z_v(z_{cls})$ ,  $\Delta z_m(z_{cls})$  denote the perturbations on features  $z_{cls}$ ,  $z_v$  and  $z_m$  through the first attack, respectively. 2) If the perturbation is directly introduced to the image token features  $z_v$ , let  $\Delta z_m(z_v)$  denote the perturbation on the aligned feature  $z_m$  through the second attack. Assume the same degree of perturbation on  $z_{cls}$  and  $z_v$  during the two attacks, then the ratio of the effect is given by:

$$\frac{\|\Delta z_m(z_{cls})\|_F}{\|\Delta z_m(z_v)\|_F} = \frac{\|\Delta z_v(z_{cls})\|_F}{\|\Delta z_{cls}(z_{cls})\|_2} \leq \frac{3 + \epsilon_V}{\sqrt{n_v}}, \quad \epsilon_V \ll 1. \quad (32)$$

## ➤ VEAttack Objective:

$$\tilde{v} = \arg \max_{\|\delta\| \leq \epsilon} -\cos(f_V(v + \delta), f_V(v)) = \arg \max_{\|\delta\| \leq \epsilon} -\cos(z_v + \delta, z_v).$$



# Experiments

Task	LVLMs Attack	OpenFlamingo-9B		LLaVa1.5-7B		LLaVa1.5-13B		Average	
		$\epsilon = 2/255$	$\epsilon = 4/255$	$\epsilon = 2/255$	$\epsilon = 4/255$	$\epsilon = 2/255$	$\epsilon = 4/255$	$\epsilon = 2/255$	$\epsilon = 4/255$
COCO	Clean	79.7		115.5		119.2		104.8(↓0.0%)	
	MIX.Attack Tu et al. (2023)	45.9	25.4	67.5	55.4	73.8	60.1	62.4(↓40.5)	47.0(↓55.1)
	VT-Attack Wang et al. (2024c)	38.9	21.6	50.8	12.2	58.2	20.1	49.3(↓52.9)	18.0(↓82.8)
	AttackVLM-ii Zhao et al. (2023)	24.4	10.3	40.9	25.8	42.7	27.5	36.0(↓65.6)	21.2(↓79.8)
	VEAttack	<b>7.5</b>	<b>3.7</b>	<b>10.8</b>	<b>7.1</b>	<b>11.2</b>	<b>6.5</b>	<b>9.8(↓90.6)</b>	<b>5.8(↓94.5)</b>
Flickr30k	Clean	60.1		77.5		77.1		71.6	
	MIX.Attack Tu et al. (2023)	33.7	18.0	42.1	35.9	41.0	32.2	38.9(↓45.7)	28.7(↓59.9)
	VT-Attack Wang et al. (2024c)	27.7	13.8	35.1	12.3	34.0	14.6	32.3(↓54.9)	13.6(↓81.0)
	AttackVLM-ii Zhao et al. (2023)	18.1	9.9	29.9	19.8	29.7	21.6	25.9(↓63.8)	17.1(↓76.1)
	VEAttack	<b>8.7</b>	<b>3.2</b>	<b>10.7</b>	<b>6.3</b>	<b>9.1</b>	<b>5.7</b>	<b>9.5(↓86.7)</b>	<b>5.1(↓92.9)</b>
TextVQA	Clean	23.8		37.1		39.0		33.3	
	MIX.Attack Tu et al. (2023)	13.4	8.8	24.6	19.1	22.8	19.7	20.3(↓39.0)	15.9(↓52.3)
	VT-Attack Wang et al. (2024c)	15.3	10.5	23.7	<b>10.0</b>	24.9	10.0	21.3(↓36.0)	10.2(↓69.4)
	AttackVLM-ii Zhao et al. (2023)	<b>12.1</b>	7.6	19.7	11.9	19.8	14.2	17.2(↓48.3)	11.2(↓66.4)
	VEAttack	12.5	<b>5.7</b>	<b>13.8</b>	10.1	<b>12.4</b>	<b>8.6</b>	<b>12.9(↓61.3)</b>	<b>8.1(↓75.7)</b>
VQAv2	Clean	48.5		74.5		75.6		66.2	
	MIX.Attack Tu et al. (2023)	39.8	36.0	59.4	57.9	59.8	58.0	53.0(↓19.9)	50.6(↓23.6)
	VT-Attack Wang et al. (2024c)	38.5	37.0	53.6	<b>21.4</b>	55.2	<b>21.0</b>	49.1(↓25.8)	<b>26.5(↓59.9)</b>
	AttackVLM-ii Zhao et al. (2023)	37.5	35.8	54.1	49.7	56.2	49.6	49.3(↓25.5)	45.0(↓32.0)
	VEAttack	<b>34.0</b>	<b>32.8</b>	<b>42.9</b>	38.4	<b>41.5</b>	37.6	<b>39.5(↓40.3)</b>	36.3(↓45.2)
POPE	Clean	65.7		84.5		84.1		78.1	
	MIX.Attack Tu et al. (2023)	59.0	53.3	72.0	69.1	74.2	68.2	68.4(↓12.4)	63.5(↓18.7)
	VT-Attack Wang et al. (2024c)	63.6	63.5	64.0	60.1	65.5	66.4	64.4(↓17.5)	63.3(↓18.9)
	AttackVLM-ii Zhao et al. (2023)	53.3	48.1	69.0	61.6	63.8	57.8	49.3(↓36.9)	45.0(↓42.4)
	VEAttack	<b>60.6</b>	<b>59.6</b>	<b>47.5</b>	<b>42.8</b>	<b>47.6</b>	<b>44.7</b>	<b>51.9(↓33.5)</b>	<b>49.0(↓37.3)</b>

Table 11: Ablation of VEAttack targets on diverse tasks with cosine similarity as the loss metric.

Objective	COCO ( $\epsilon$ )↓	VQAv2 ( $\epsilon$ )↓		POPE ( $\epsilon$ )↓			
		$z_v$	$z_{cls}$	2/255	4/255	2/255	4/255
✗	✓	43.6	25.5	56.0	50.0	69.4	59.1
✓	✓	22.0	10.5	46.4	42.0	59.2	46.0
✓	✗	<b>10.8</b>	<b>7.1</b>	<b>42.9</b>	<b>38.4</b>	<b>47.5</b>	<b>42.8</b>

Table 12: Ablation of loss metrics in attack objective.

Measurement	COCO ( $\epsilon$ )↓		VQAv2 ( $\epsilon$ )↓		POPE ( $\epsilon$ )↓	
	2/255	4/255	2/255	4/255	2/255	4/255
Euclidean	46.9	39.6	56.1	53.6	62.1	64.3
K-L divergence	70.0	34.3	59.4	<b>33.2</b>	70.0	60.5
Cosine similarity	<b>10.8</b>	<b>7.1</b>	<b>42.9</b>	38.4	<b>47.5</b>	<b>42.8</b>

# Observation

**Observation 1** *Even though the LLMs and tasks are downstream-agnostic, attacks on output of the vision encoder can lead to variations in the hidden layer features of LLMs.*

**Observation 2** *LVLMs will pay more attention to image tokens in the image caption task. Conversely, LVLMs focus more on user instruction tokens in the VQA task.*

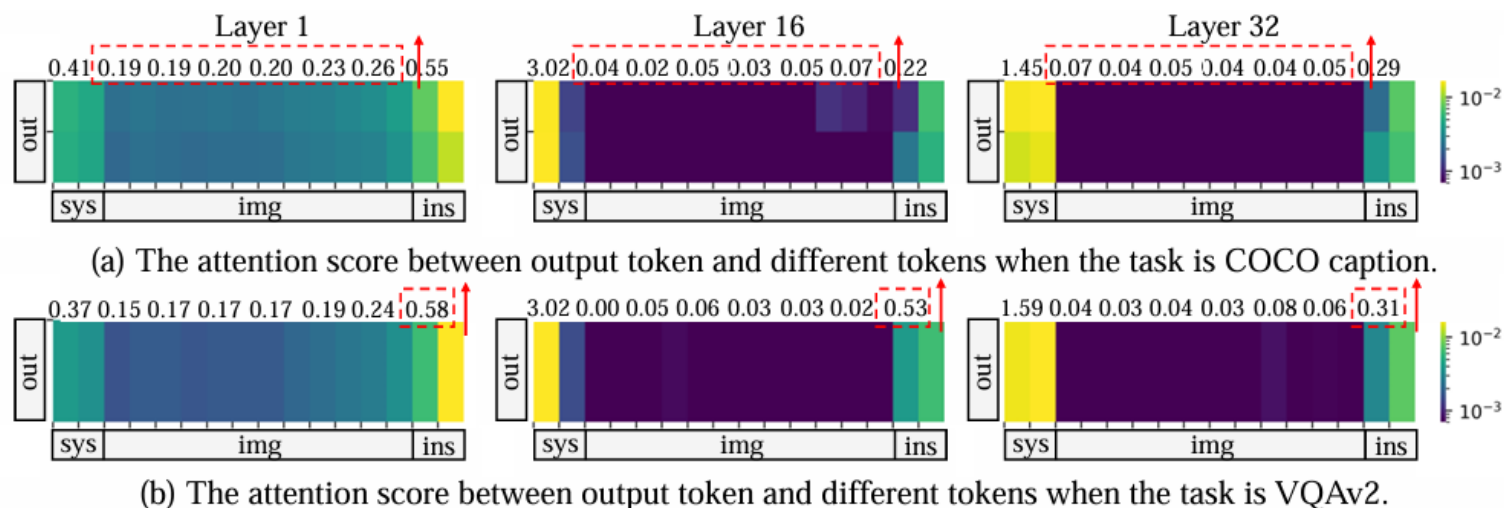
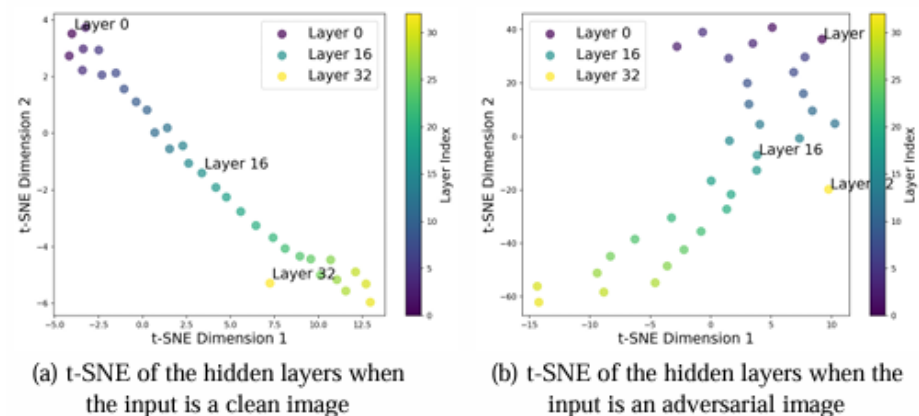


Figure 11: t-SNE visualization of the first visual tokens across hidden layers in the LLM model for clean and adversarial image inputs.

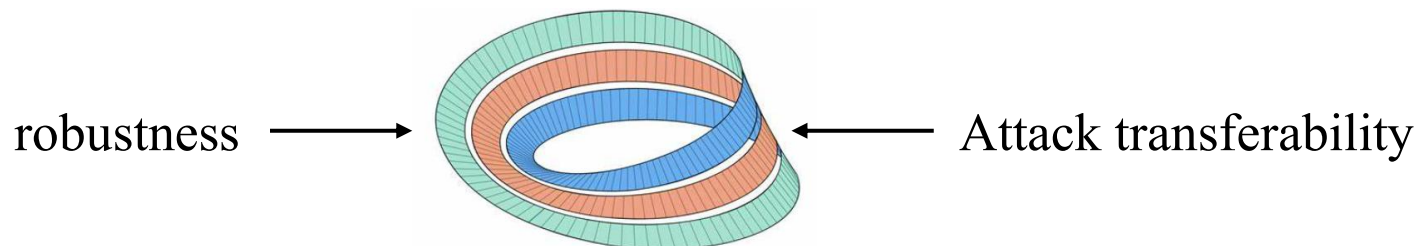
Figure 12: Illustration of attention maps across different layers of the language model for two tasks. The attention scores shown represent the relationship between the output token and three input token types (out: output tokens, sys: system tokens, img: image tokens, ins: instruction tokens).

# Observation

**Observation 3** *The transfer attack on LVLMs normally resembles a Möbius band, where robustness and vulnerability of LVLMs intertwine as a single, twisted continuum. For defenders, using a more robust vision encoder can enhance the ability of LVLM to resist attacks. Conversely, for attackers, adversarial samples obtained by attacking a more robust vision encoder usually have higher attack transferability on diverse LVLMs.*

Table 2: Transfer Attack of original CLIP and robust CLIP after adversarial training on the COCO dataset. Bold indicates the best attack performance, while underlined indicates the second best.

Source \ Target	LLaVa1.5-7B			OpenFlamingo-9B			MiniGPT-4	mPLUG-Owl2	Qwen-VL
	CLIP	TeCoA	FARE	CLIP	TeCoA	FARE	BLIP-2	MplugOwl	CLIP-bigG
clean	115.5	88.3	102.4	79.7	66.9	74.1	96.7	132.3	138.5
CLIP Radford et al. (2021)	<b>3.7</b>	92.2	100.4	<b>3.6</b>	70.7	78.1	55.2	124.9	131.2
TeCoA Mao et al. (2023)	22.1	<b>25.3</b>	<u>25.1</u>	17.7	<b>20.8</b>	<u>22.6</u>	<b>11.0</b>	<u>96.4</u>	<u>102.7</u>
FARE Schlarmann et al. (2024)	<u>38.7</u>	<u>54.5</u>	<b>48.6</b>	<u>26.3</u>	<u>42.7</u>	<b>36.4</b>	<u>8.9</u>	<b>88.0</b>	<b>101.4</b>



# Observation

**Observation 4** *Reducing attack steps does not significantly impair the effectiveness of VEAttack, while increasing the perturbation budget beyond a threshold within an imperceptible range does not continuously degrade the performance of LLMs significantly.*

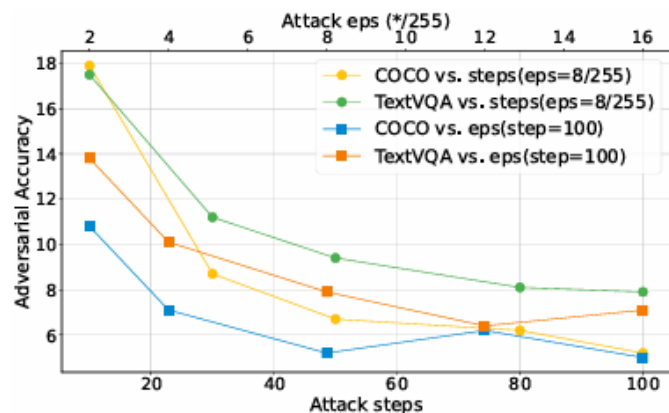


Figure 13: VEAttack performance under different attack steps and perturbation budgets on COCO and TextVQA datasets (top x-axis: perturbation budget, bottom x-axis: attack steps).

Table 3: Comparison of the effectiveness and efficiency between VEAttack and other white-box and gray-box attacks. Flops count the computation of forward once.

Attack	Version	Flops	COCO Time (h)	Flickr30k Time (h)	TextVQA Time (h)
clean	None	99.3G	115.5	1.33	37.1
APGD	White-box	9.93T	13.1	25.5	8.1
Ensemble	White-box	9.93T	<b>3.1</b>	41.9	<b>0.0</b>
VT-Attack	Gray-box	3.04T	12.2	65.0	10.0
VEAttack ( $\epsilon = 4/255, t = 100$ )	Gray-box	2.59T	7.1	8.5	10.1
VEAttack ( $\epsilon = 8/255, t = 50$ )	Gray-box	2.59T	<u>5.5</u>	<b>5.3</b>	<u>4.5</u>

Thank you for  
your listening