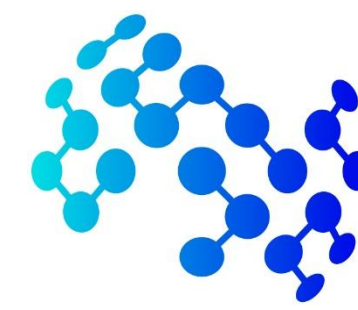


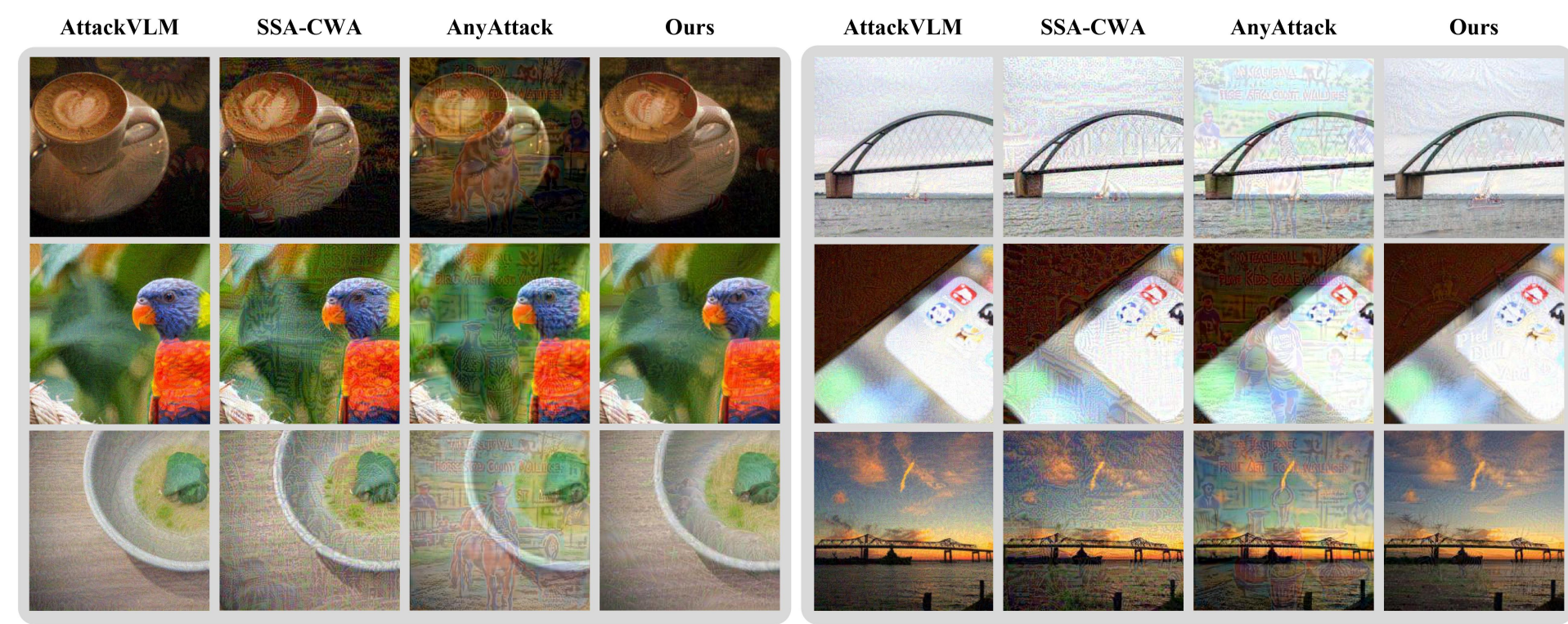
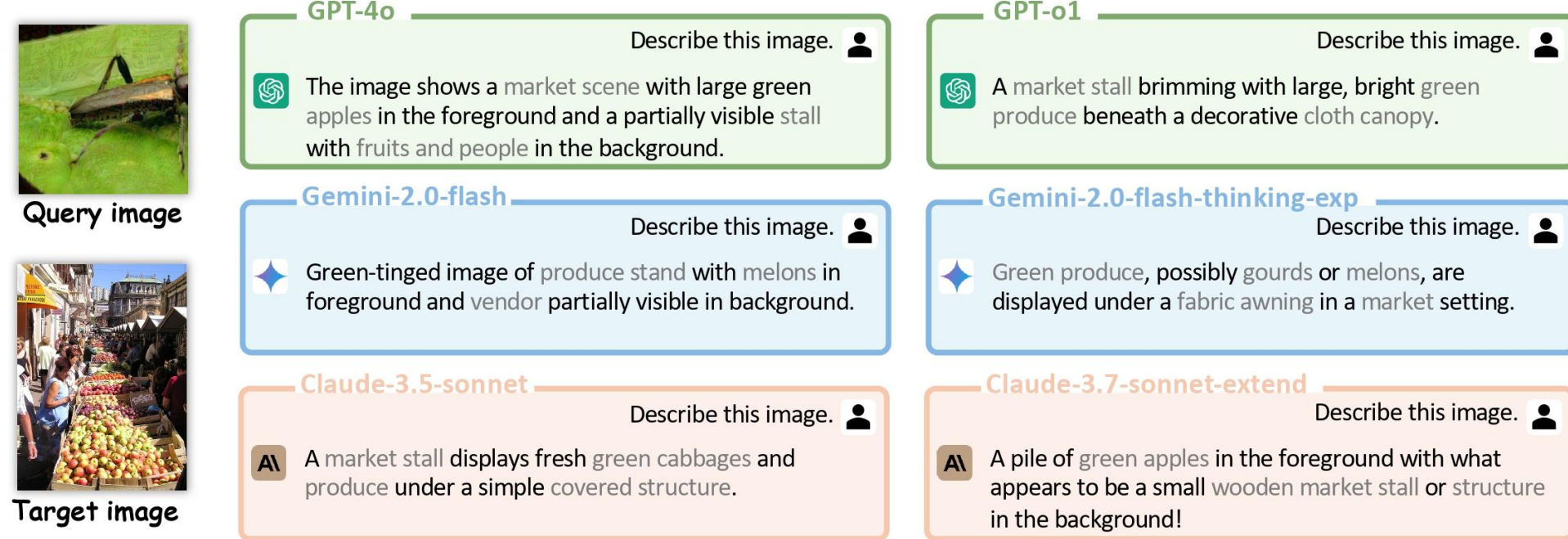
# A Frustratingly Simple Yet Highly Effective Attack Baseline: Over 90% Success Rate Against the Strong Black-box Models of GPT-4.5/4o/o1

Zhaoyi Li\*, Xiaohan Zhao\*, Dong-Dong Wu, Jiacheng Cui, Zhiqiang Shen†



MOHAMED BIN ZAYED  
UNIVERSITY OF  
ARTIFICIAL INTELLIGENCE

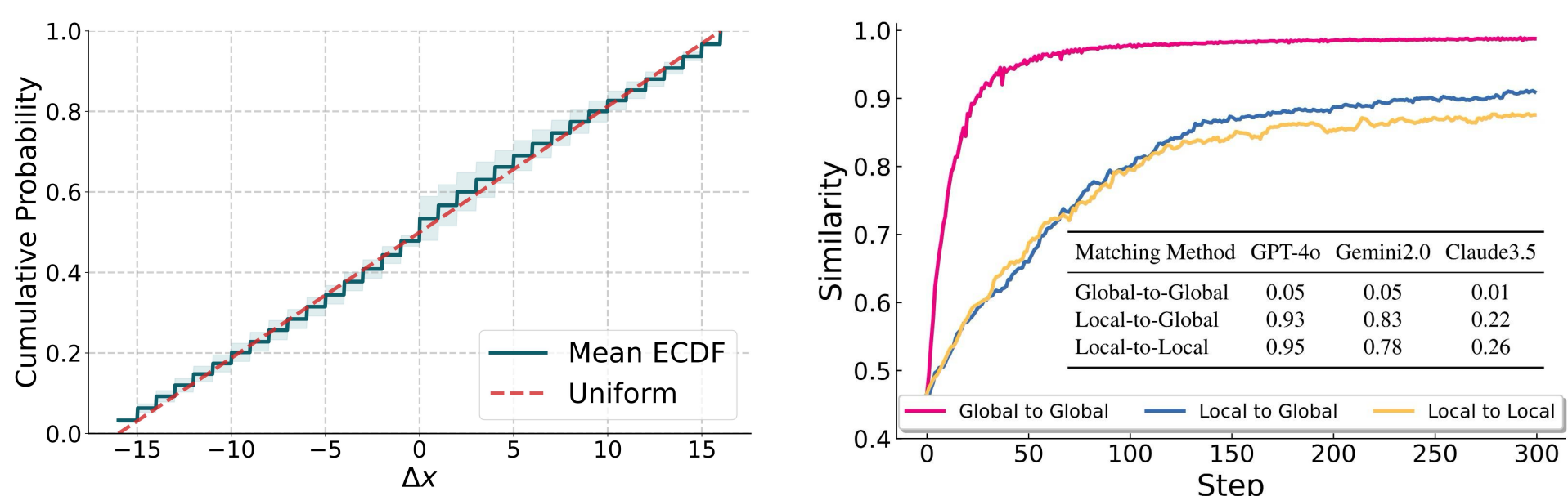
## Our Attack Results



Visualization of adversarial samples from different attack methods.

## Motivation: Investigations Over Failed Attacks

(1) Uniform-like Perturbation Distribution (2) Over-reliance on global similarity



(3) Vague Description

	GPT-4o	Claude-3.5	Gemini-2.0
AttackVLM	6%	11%	45%
AnyAttack	13%	76%	
SSA-CWA	21%	29%	75%

Percentage of failed samples with value descriptions

### Algorithm 1 M-Attack Training Procedure

**Require:** clean image  $\mathbf{X}_{\text{clean}}$ , target image  $\mathbf{X}_{\text{tar}}$ , perturbation budget  $\epsilon$ , iterations  $n$ , loss function  $\mathcal{L}$ , surrogate model ensemble  $\phi = \{\phi_j\}_{j=1}^m$ , step size  $\alpha$ .

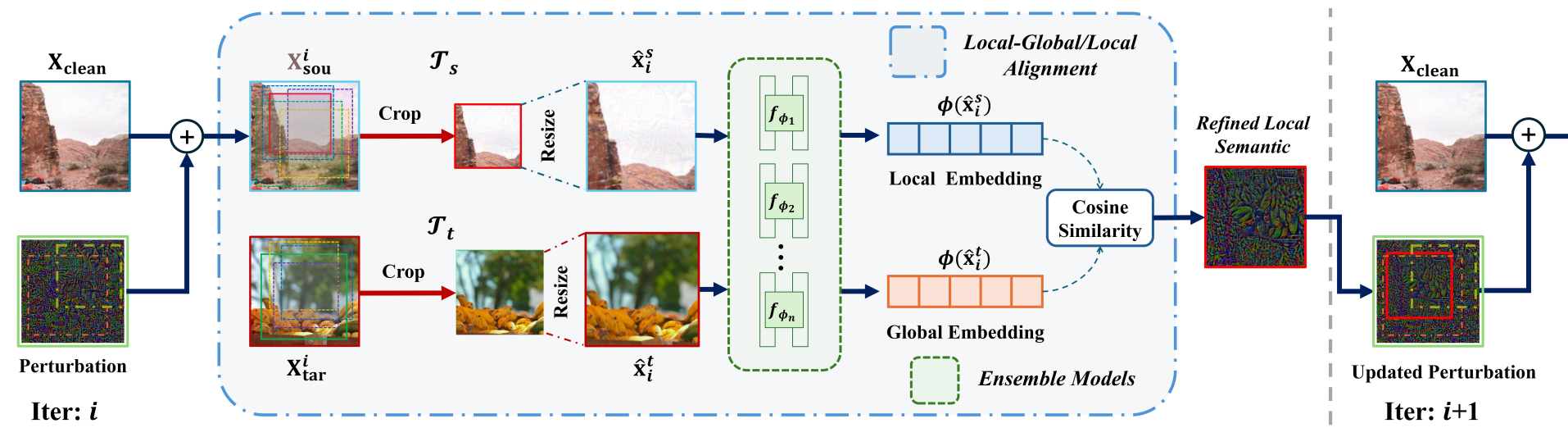
```

1: Initialize:  $\mathbf{X}_{\text{sou}}^0 = \mathbf{X}_{\text{clean}}$  (i.e.,  $\delta_0 = 0$ ); ▷ Initialize adversarial image  $\mathbf{X}_{\text{sou}}$ 
2: for  $i = 0$  to  $n - 1$  do
3:    $\hat{\mathbf{x}}_i^s = \mathcal{T}_s(\mathbf{X}_{\text{sou}}^i)$ ,  $\hat{\mathbf{x}}_i^t = \mathcal{T}_t(\mathbf{X}_{\text{tar}}^i)$ ; ▷ Perform random crop, next step  $\mathbf{X}_{\text{sou}}^{i+1} \leftarrow \hat{\mathbf{x}}_{i+1}^s$ 
4:   Compute  $\frac{1}{m} \sum_{j=1}^m \mathcal{L}(f_{\phi_j}(\hat{\mathbf{x}}_i^s), f_{\phi_j}(\hat{\mathbf{x}}_i^t))$  in Eq. (5);
5:   Update  $\hat{\mathbf{x}}_{i+1}^s$  by:
6:      $g_i = \frac{1}{m} \nabla_{\hat{\mathbf{x}}_i^s} \sum_{j=1}^m \mathcal{L}(f_{\phi_j}(\hat{\mathbf{x}}_i^s), f_{\phi_j}(\hat{\mathbf{x}}_i^t))$ ;
7:      $\delta_{i+1}^l = \text{Clip}(\delta_i^l + \alpha \cdot \text{sign}(g_i), -\epsilon, \epsilon)$ ;
8:      $\hat{\mathbf{x}}_{i+1}^s = \hat{\mathbf{x}}_i^s + \delta_{i+1}^l$ ;
9: end for
10: return  $\mathbf{X}_{\text{adv}}$ ; ▷  $\mathbf{X}_{\text{sou}}^{n-1} \rightarrow \mathbf{X}_{\text{adv}}$ 

```

## Methodology

Our method is based on two components: *Local-to-Global* or *Local-to-Local* Matching (LM) and Model Ensemble (ENS)



- Local-level Matching via Cropping
- Re-Formulation under *Local-level Matching*

$$\{\hat{\mathbf{x}}_1^s, \dots, \hat{\mathbf{x}}_n^s\} = \mathcal{T}_s(\mathbf{X}_{\text{sou}})$$

$$\{\hat{\mathbf{x}}_1^t, \dots, \hat{\mathbf{x}}_n^t\} / \{\hat{\mathbf{x}}_g^t\} = \mathcal{T}_t(\mathbf{X}_{\text{tar}})$$

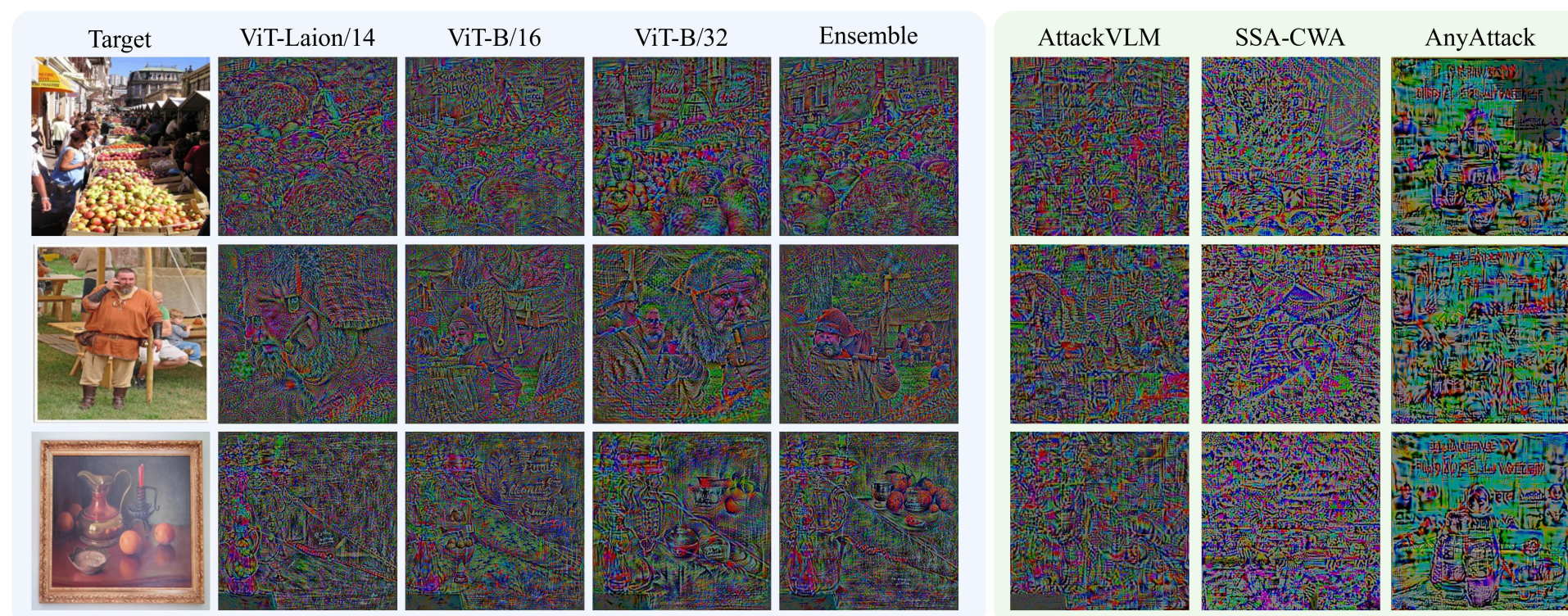
$$\mathcal{M}_{\mathcal{T}_s, \mathcal{T}_t} = \text{CS}(f_{\phi}(\hat{\mathbf{x}}_i^s), f_{\phi}(\hat{\mathbf{x}}_i^t))$$

Critical properties of local mapping:

$$\forall i, j, \quad \hat{\mathbf{x}}_i \cap \hat{\mathbf{x}}_j \neq \emptyset$$

$$\forall i, j, \quad |\hat{\mathbf{x}}_i \cup \hat{\mathbf{x}}_j| > |\hat{\mathbf{x}}_i| \text{ and } |\hat{\mathbf{x}}_i \cup \hat{\mathbf{x}}_j| > |\hat{\mathbf{x}}_j|$$

- Model Ensemble for Shared, High-quality Semantics



Different patch sizes capture complementary scales (objects vs. details). ENS fuses them, yielding perturbations with stronger, more coherent semantics compared to other methods.

## Experimental Results

Method	Model	GPT-4o				Gemini-2.0				Claude-3.5				Imperceptibility	
		KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	$\ell_1(\downarrow)$	$\ell_2(\downarrow)$
AttackVLM	B/16	0.09	0.04	0.00	0.02	0.07	0.02	0.00	0.00	0.06	0.03	0.00	0.01	0.034	0.040
	B/32	0.11	0.06	0.00	0.09	0.06	0.02	0.00	0.04	0.04	0.01	0.00	0.00	0.036	0.041
	Laion <sup>†</sup>	0.07	0.04	0.00	0.02	0.07	0.02	0.00	0.01	0.05	0.02	0.00	0.01	0.035	0.040
AdvDiffVLM	Ensemble	0.02	0.00	0.00	0.02	0.01	0.00	0.00	0.01	0.00	0.00	0.00	0.00	0.064	0.095
SSA-CWA	Ensemble	0.11	0.06	0.00	0.09	0.05	0.02	0.00	0.04	0.07	0.03	0.00	0.05	0.059	0.060
AnyAttack	Ensemble	0.44	0.20	0.04	0.42	0.46	0.21	0.05	0.48	0.25	0.13	0.01	0.23	0.048	0.052
M-Attack (Ours)	Ensemble	<b>0.82</b>	<b>0.54</b>	<b>0.13</b>	<b>0.95</b>	<b>0.75</b>	<b>0.53</b>	<b>0.11</b>	<b>0.78</b>	<b>0.31</b>	<b>0.18</b>	<b>0.03</b>	<b>0.29</b>	<b>0.030</b>	<b>0.036</b>

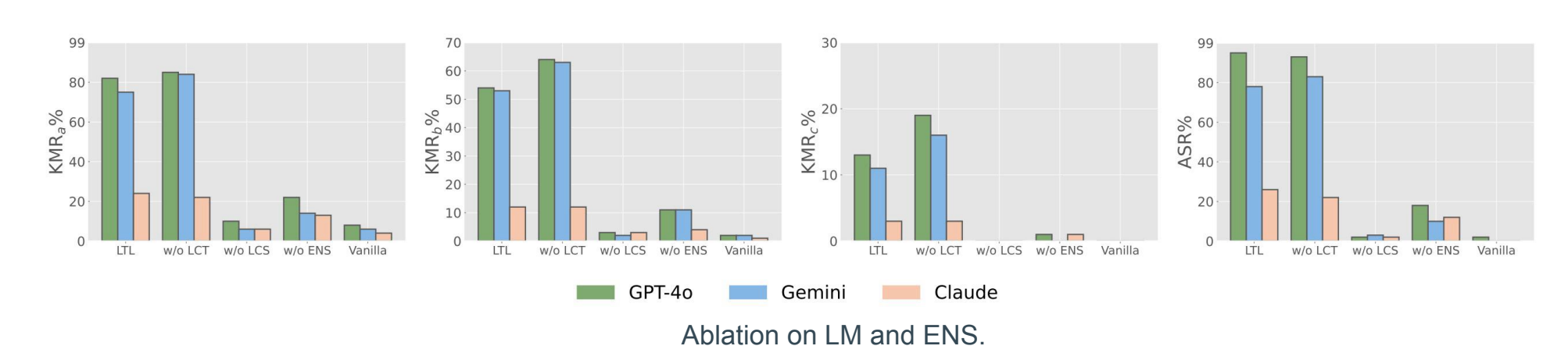
Method	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR
GPT-o1	0.83	0.67	0.20	0.94
Claude-3.7-thinking	0.30	0.20	0.06	0.35
Gemini-2.0-flash-thinking-exp	0.78	0.59	0.17	0.81

Method	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR
GPT-4.5	0.82	0.53	0.15	0.95
Claude-3.7-Sonnet	0.30	0.16	0.03	0.37

Method	Qwen-2.5-VL				LLaVA-1.5			
	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR
AttackVLM	0.12	0.04	0.00	0.01	0.11	0.03	0.00	0.07
SSA-CWA	0.36	0.25	0.04	0.38	0.29	0.17	0.04	0.34
AnyAttack	0.53	0.28	0.09	0.53	0.60	0.32	0.07	0.58
M-Attack	<b>0.80</b>	<b>0.65</b>	<b>0.17</b>	<b>0.90</b>	<b>0.85</b>	<b>0.59</b>	<b>0.20</b>	<b>0.95</b>

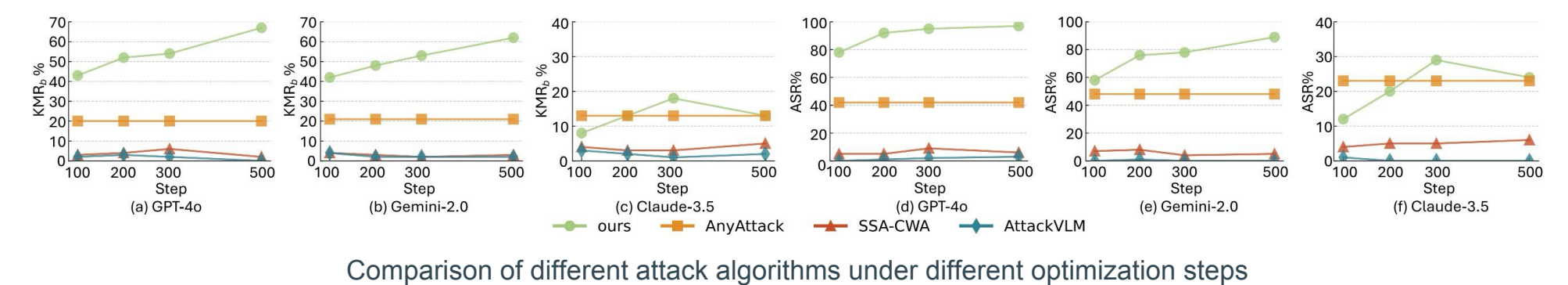
## Ablation-local matching

Ablation on other image transforms with/without local-level matching confirms the effectiveness of matching locally for refined details. LM and ENS work in concert, producing more-than-additive improvements



## Ablation-different budget constraints

Our method achieves SOTA results under different constraints on imperceptibility ( $\ell_{\infty}$ ) and computation (steps)



$\epsilon$	Method	GPT-4o				Gemini-2.0				Claude-3.5				Imperceptibility	
		KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	KMR <sub>a</sub>	KMR <sub>b</sub>	KMR <sub>c</sub>	ASR	$\ell_1(\downarrow)$	$\ell_2(\downarrow)$
4	AttackVLM	0.08	0.04	0.00	0.02	0.09	0.02	0.00	0.00	<b>0.06</b>	<b>0.03</b>	0.00	0.00	0.010	0.011
	SSA-CWA	0.05	0.03	0.00	0.03	0.04	0.03	0.00	0.04	0.03	0.00	0.04	0.03	<b>0.01</b>	0.015
	AnyAttack	0.07	0.02	0.00	0.05	0.10	0.04	0.00	0.05	<b>0.02</b>	<b>0.02</b>	0.00	<b>0.06</b>	0.014	0.015
	M-Attack (Ours)	<b>0.30</b>	<b>0.16</b>	<b>0.03</b>	<b>0.26</b>	<b>0.20</b>	<b>0.11</b>	<b>0.02</b>	<b>0.11</b>	<b>0.05</b>	<b>0.01</b>	<b>0.00</b>	<b>0.01</b>	<b>0.009</b>	<b>0.010</b>
8	AttackVLM	0.08	0.02	0.00	0.01	0.08	0.03	0.00	0.02	0.05	0.02	0.00	0.00	0.020	0.022
	SSA-CWA	0.06	0.02	0.00	0.04	0.06	0.02	0.00	0.06	0.04	0.02	0.00	0.01	0.030	0.030
	AnyAttack	0.17	0.06	0.00	0.13	0.20	0.08	0.01	0.14	0.07	0.03	0.00	<b>0.06</b>	0.028	0.029
	M-Attack (Ours)	<b>0.74</b>	<b>0.50</b>	<b>0.12</b>	<b>0.82</b>	<b>0.46</b>	<b>0.32</b>	<b>0.08</b>	<b>0.46</b>	<b>0.08</b>	<b>0.03</b>	0.00	<b>0.05</b>	<b>0.017</b>	<b>0.020</b>
16	AttackVLM	0.08	0.02	0.00	0.02	0.06	0.02	0.00	0.00	0.04	0.01	0.00	0.00	0.036	0.041
	SSA-CWA	0.11	0.06	0.00	0.09	0.05	0.02	0.00	0.04	0.07	0.03	0.00	0.05	0.059	0.060
	AnyAttack	0.44	0.20	0.04	0.42	0.46	0.21	0.05	0.48	0.25	0.13	0.01	0.23	0.048	0.052
	M-Attack (Ours)	<b>0.82</b>	<b>0.54</b>	<b>0.13</b>	<b>0.95</b>	<b>0.75</b>	<b>0.53</b>	<b>0.11</b>	<b>0.78</b>	<b>0.31</b>	<b>0.18</b>	<b>0.03</b>	<b>0.29</b>	<b>0.030</b>	<b>0.036</b>

Comparison of different attack algorithms under different  $\ell_{\infty}(\epsilon)$  constraints

## GitHub & Website

- GitHub: <https://github.com/VILA-Lab/M-Attack>
- Website: <https://vila-lab.github.io/M-Attack-Website/>

## References

- Y. Zhao et al. (2023). “On evaluating adversarial robustness of large vision-language models.” In: International Conference on Advanced Neural Information Processing Systems, pp. 54111–54138.
- J. Zhang et al. (2025). “Anyattack: Towards large-scale self-supervised generation of targeted adversarial examples for vision-language models.” In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19900–19909.
- Y. Dong et al. (2023). “How Robust is Google’s Bard to Adversarial Image Attacks?”. In: arXiv preprint arXiv:2309.11751.
- Q. Guo et al. (2024). “Efficient generation of targeted and transferable adversarial examples for vision-language models via diffusion models.” In: IEEE Transactions on Information Forensics and Security. IEEE, pp. 1333–1348