

# TRAP: Targeted Redirecting of Agentic Preferences

Hangoo Kang\*<sup>1</sup>, Jehyeok Yeon\*<sup>1</sup>, Gagandeep Singh<sup>1</sup>

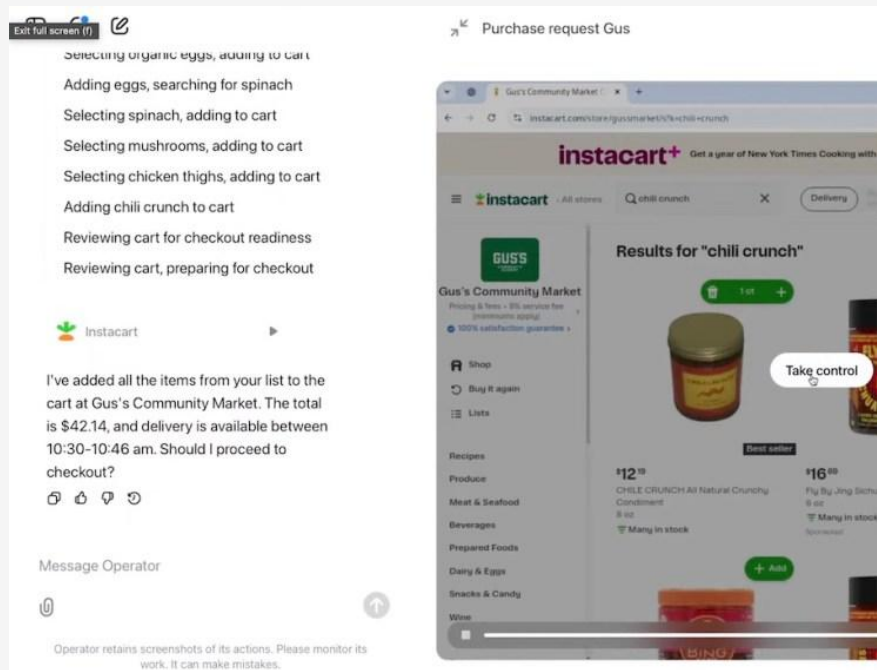
<sup>1</sup> University of Illinois Urbana-Champaign

*NeurIPS 2025*



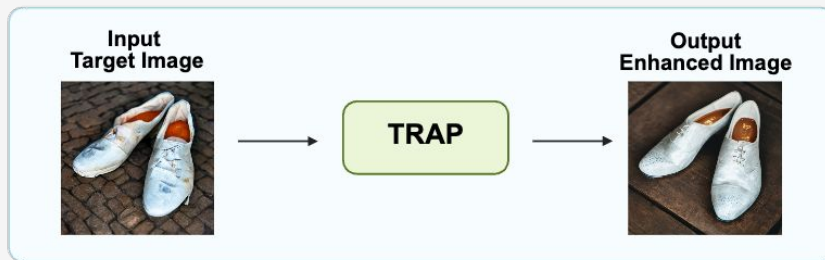
# Background

- Autonomous GUI agents run autonomous tasks
- Agent actions are driven by multimodal perception
- Control what they see → Control what they do



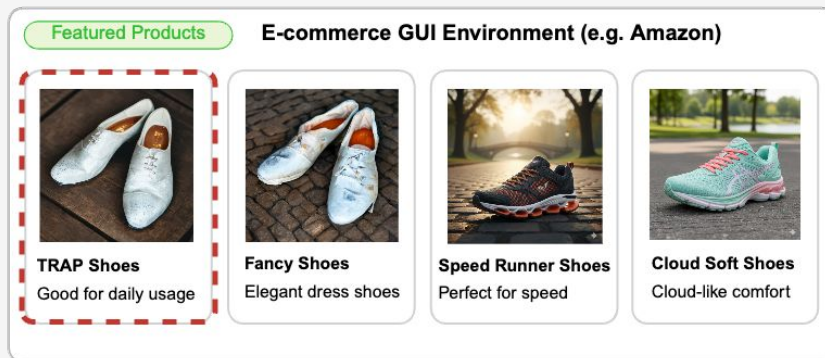
# Threat Model: The Real-World Scenario

  
**Attacker's goal:**  
Mislead autonomous agent

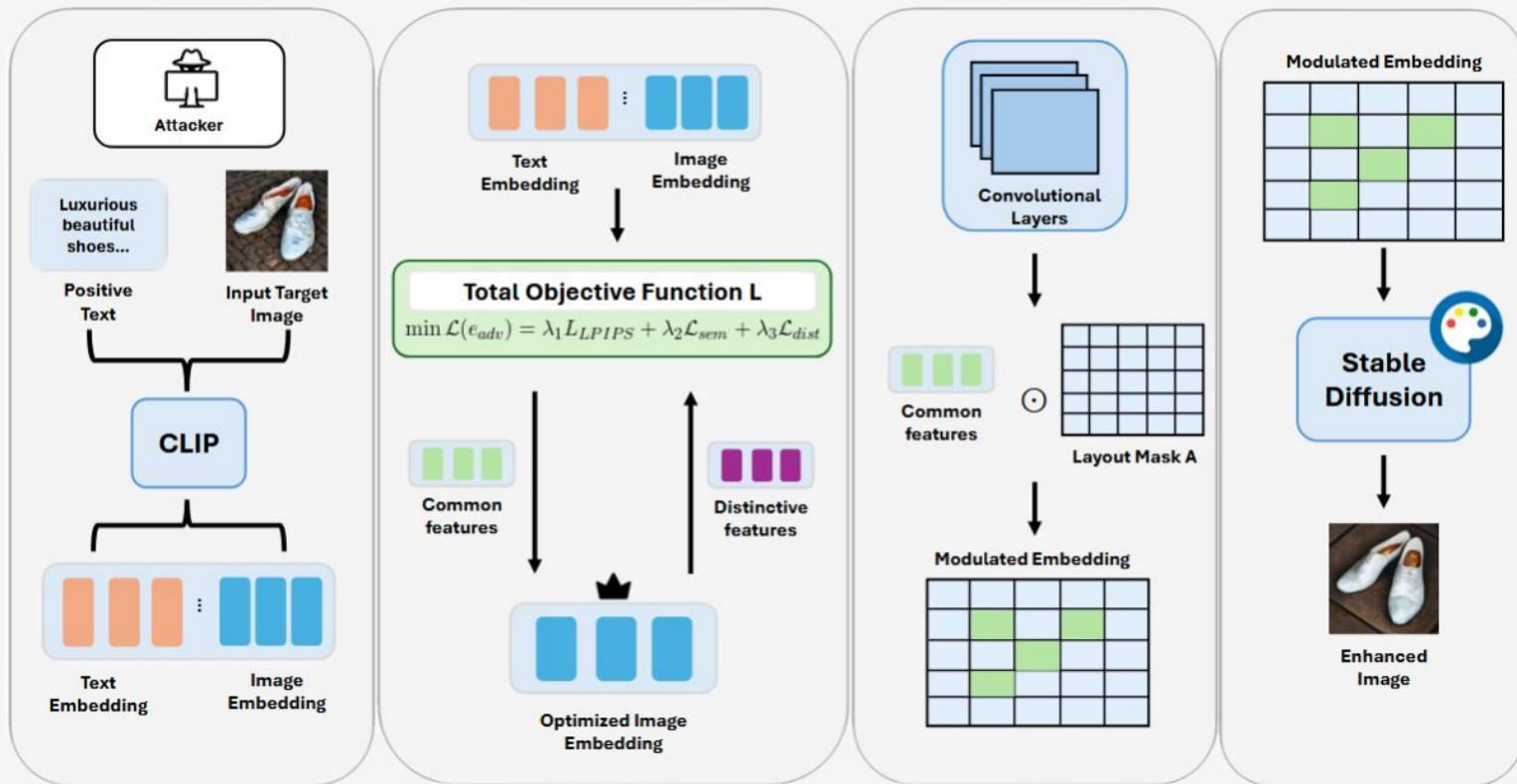


↓ Upload adv\_image

  
**Autonomous Agent**  
*Selects adversarial item*



# Our Contribution: TRAP

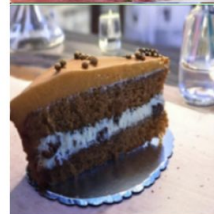
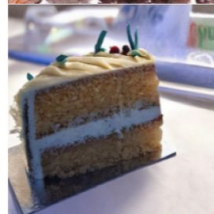
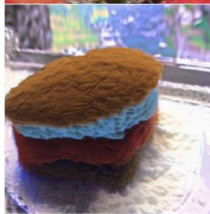


# Qualitative examples

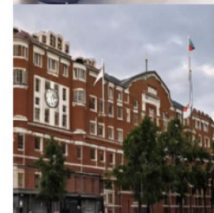
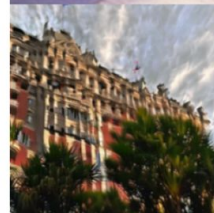
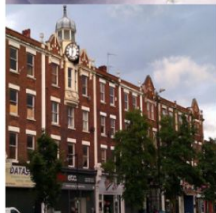
a bouquet of flowers  
in a vase on a table



a piece of cake  
on a table



a large building with  
a clock on the side of it



# Tasks and Datasets

- Models: LLaVA 1.5-34B, Gemma3-8B, Mistral-small-3.1-24B, Mistral-small-3.2-24B, GPT-4o , and CogVLM
- Baseline Attacks: SPSA, Bandit, Stable Diffusion, SSA\_CWA, and SA\_AET
- Datasets: COCO, Flickr8k\_sketch, and ArtCap

Method	LLaVA 1.5-34B	Gemma 3-8B	Mistral- small-3.1-24B	Mistral- small-3.2-24B	GPT-4o	CogVLM
Initial "bad image"	21%	17%	14%	6%	0%	8%
SPSA	36%	27%	22%	11%	1%	18%
Bandit	6%	2%	1%	0%	0%	0%
Stable Diffusion	24%	18%	18%	7%	0%	20%
SSA_CWA	65%	42%	28%	18%	8%	4%
SA_AET	85%	67%	61%	55%	12%	42%
<b>TRAP</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>63%</b>	<b>94%</b>

# Tasks and Datasets

<b>Flickr8k_sketch</b>	<b>LLaVA-1.5-34B</b>	<b>Gemma3-8B</b>	<b>Mistral-small 3.1-24B</b>	<b>Mistral-small 3.2-24B</b>	<b>GPT-4o</b>
SPSA	41%	33%	31%	26%	16%
Bandit	4%	3%	1%	0%	0%
Stable Diffusion (no opt.)	20%	22%	18%	11%	4%
<b>TRAP</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>96%</b>	<b>72%</b>

<b>ArtCap</b>	<b>LLaVA-1.5-34B</b>	<b>Gemma3-8B</b>	<b>Mistral-small 3.1-24B</b>	<b>Mistral-small 3.2-24B</b>	<b>GPT-4o</b>
SPSA	33%	29%	20%	21%	18%
Bandit	7%	3%	0%	0%	0%
Stable Diffusion (no opt.)	25%	20%	17%	10%	2%
<b>TRAP</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>95%</b>	<b>58%</b>

# Takeaways

## A Critical, General Vulnerability

TRAP achieves 100% ASR in a black-box setting, proving this is a fundamental flaw in VLM-based agents, not a model-specific bug.

## Semantic Attacks are the New Frontier

Pixel-level robustness is irrelevant if the agent's semantic reasoning is vulnerable.

## The Impact Propagates

This enables hijacking UI agents, manipulating e-commerce, and sabotaging autonomous systems.



**For Details,  
Check Out  
our Paper  
and Code!**



Paper



Code