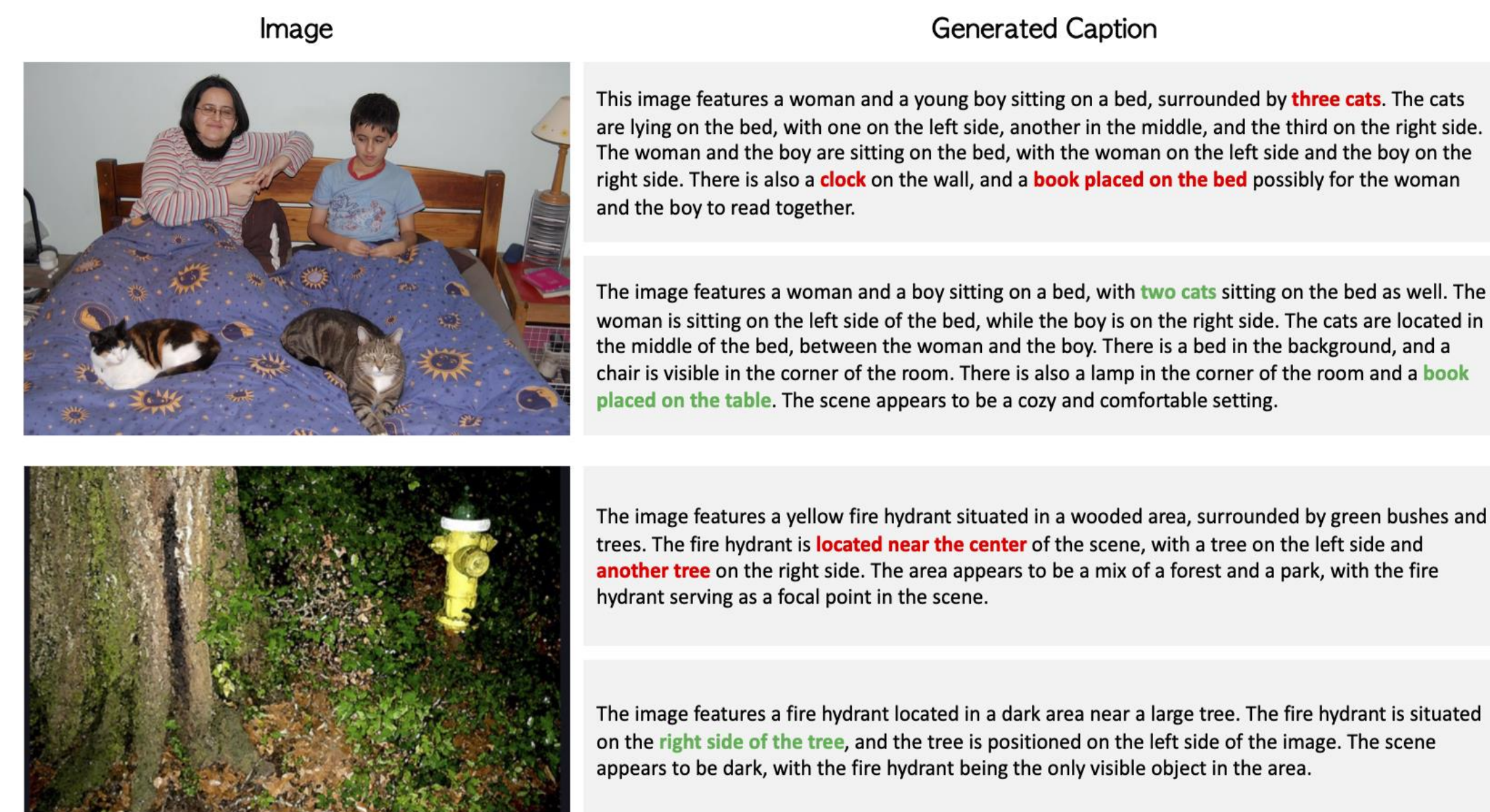


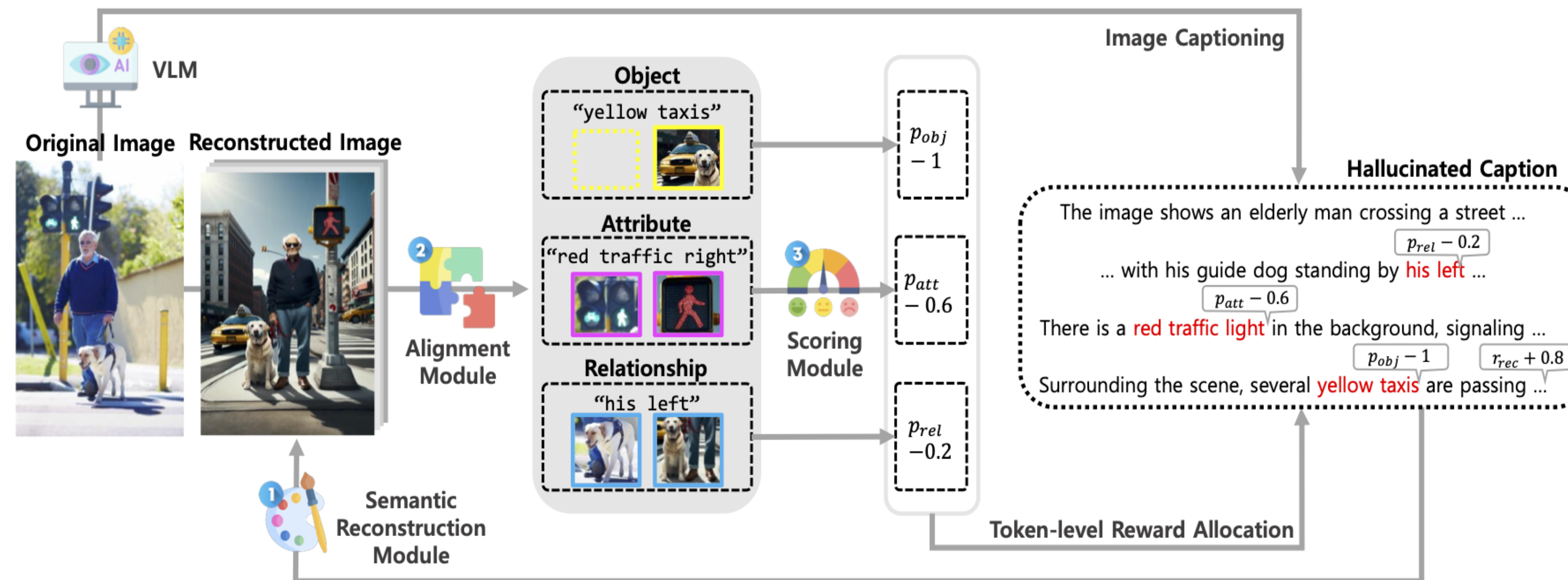
## SUMMARY

- (1) We propose ESREAL, a fully **unsupervised** hallucination mitigation framework. Our approach is scalable, eliminating the need for annotated data during training.
- (2) We craft a hallucination detection pipeline, which facilitates **token-level identification** of hallucinations in generated captions in a reference-free manner via semantic reconstruction.
- (3) We show that ESREAL can be applied across a variety of VLMs to effectively mitigate hallucinations.



## METHODOLOGY

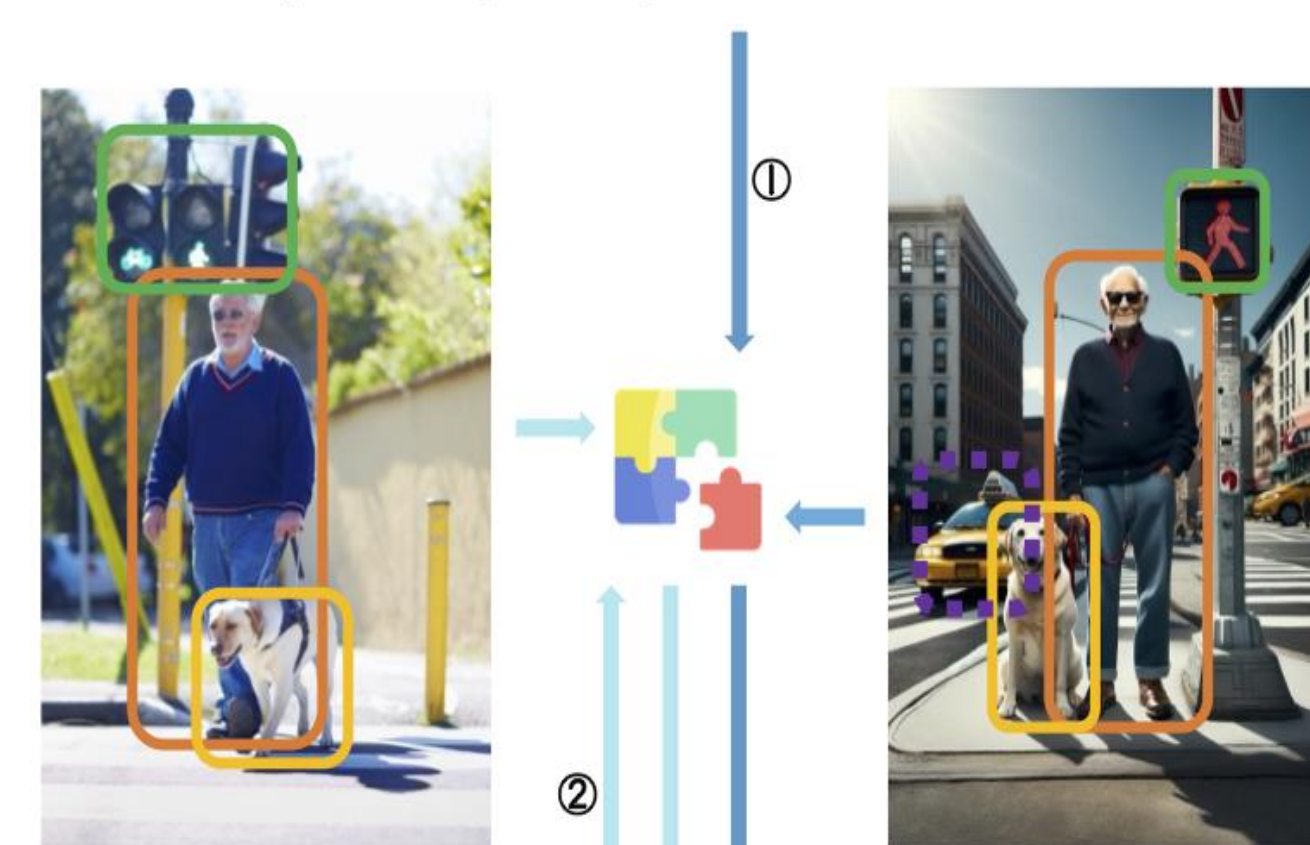
ESREAL uses **token-level penalties** from a hallucination detection pipeline and a **fine-grained PPO approach** to selectively suppress hallucinatory content in generated text.



### ① Semantic Reconstruction Module

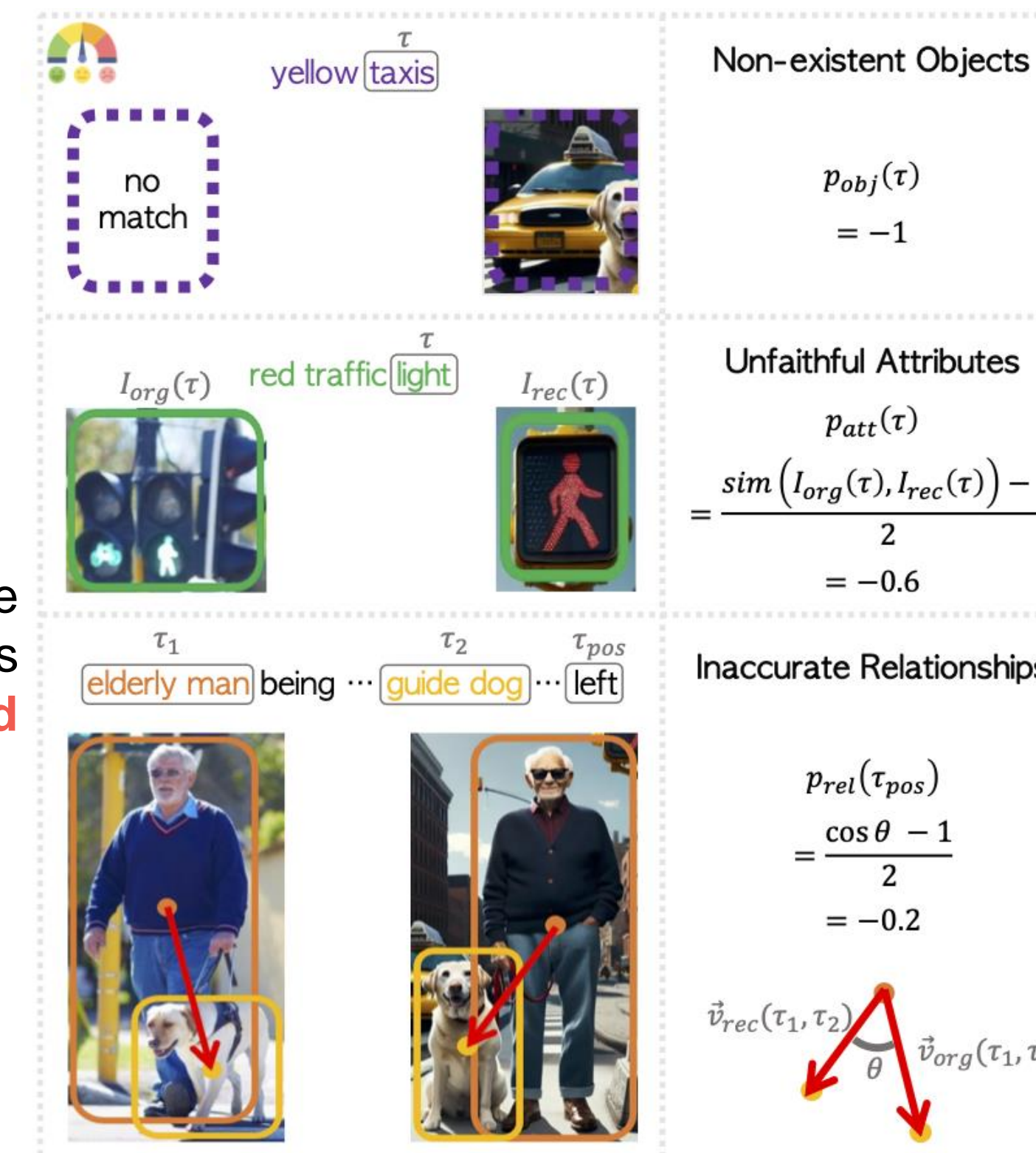
Reconstruct the image using the generated caption with the T2I model.

"The image shows an elderly man being led across a crosswalk by his guide dog positioned to his left ...  
There is a red traffic light in the background, signaling ...  
Surrounding the scene, several yellow taxis add vibrant color and ..."



### ② Alignment Module

Match **object phrases** from the generated caption with regions in both the **reconstructed** image and the **original** image.

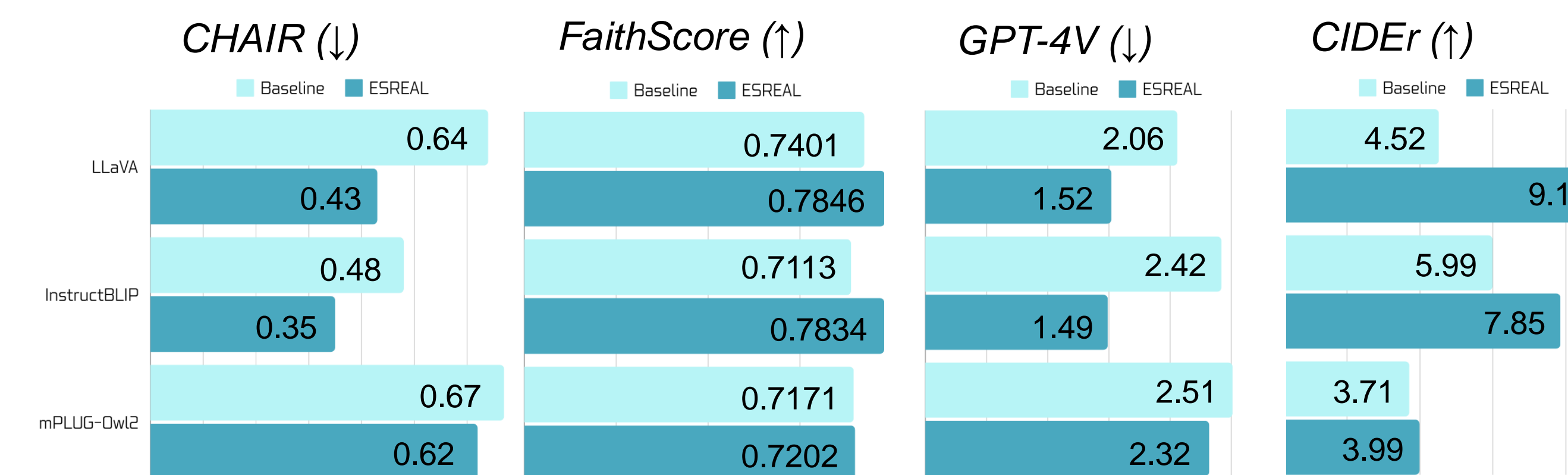


### ③ Scoring Module

Identify discrepancies between the aligned regions. Apply penalties for object, attribute, or relationship hallucinations.

## EXPERIMENTS

① Does ESREAL reduce hallucinations while preserving VLMs' generative abilities?



② Does each penalty effectively target its type?

Model	Method	# Hallucinations per Caption (↓)			
		Object	Attribute	Relationship	Total
InstructBLIP	Baseline	1.23	0.14	1.05	2.42
	ESREAL	0.80	0.06	0.64	1.49
InstructBLIP	ESREAL w/o $p_{obj}$	1.18	0.14	0.92	2.24
	ESREAL w/o $p_{att}$	0.93	0.09	0.59	1.61
	ESREAL w/o $p_{rel}$	0.85	0.07	1.28	2.21
	ESREAL w/o $r_{rec}$	0.68	0.04	1.40	2.12
	ESREAL w/o $p_{obj}, p_{att}, p_{rel}$	1.15	0.10	1.01	2.26

③ How does T2I model stability impact ESREAL's performance?

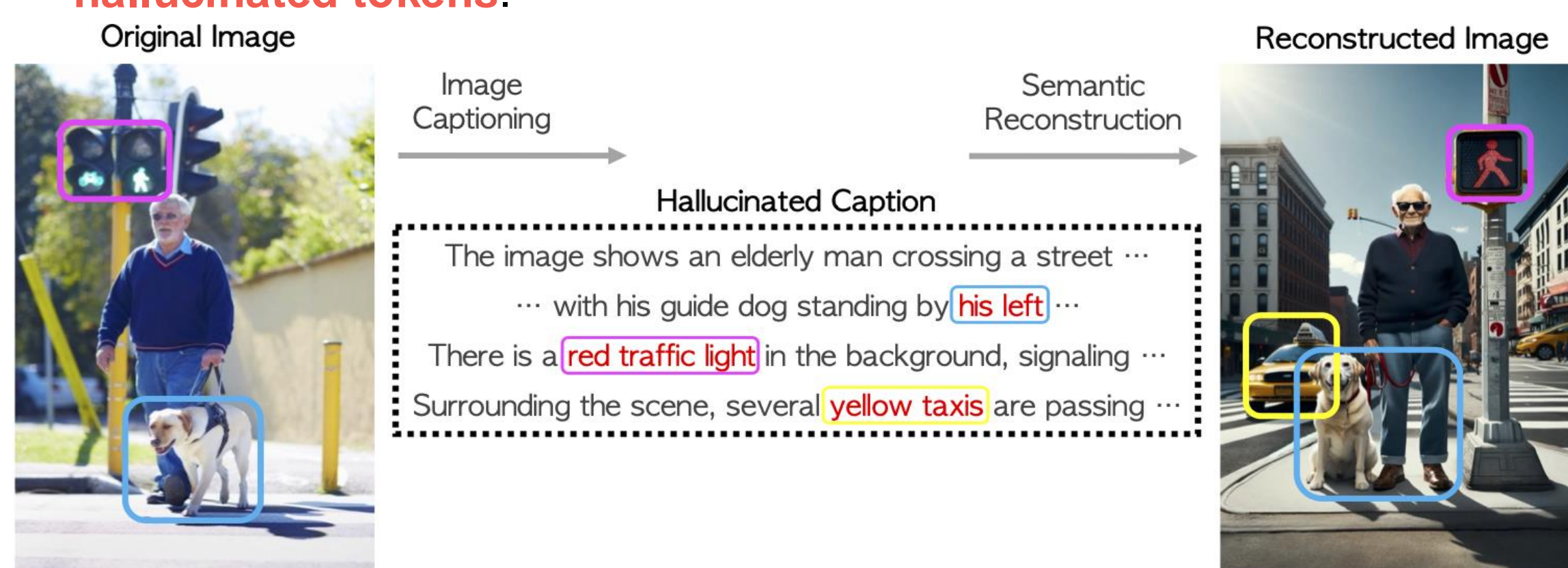
T2I Models	Win Rate	Alignment	FaithScore (↑)	GPT-4V (↓)
SDXL-Turbo (1 Step)	0.76	1.86	0.7484	2.23
SDXL-Turbo (4 Steps)	0.79	2.30	0.7834	1.49
Hyper-SDXL (8 Steps)	0.80	2.58	0.8141	1.32
DALLE-3	0.82	2.71	-	-

## ACKNOWLEDGEMENT

We also thank great previous work including GroundingDINO, SDXL, SDXL-Turbo, Hyper-SDXL, DALLE-3, etc.

## MOTIVATION

- Hallucination in the caption leads to **semantic misalignment** between the original and reconstructed images.
- By comparing the disparities among corresponding regions in the images, we can effectively **identify** and **penalize** the generation of **hallucinated tokens**.



## CONTACT

Minchan Kim



Minyeong Kim

