

# Generate, but Verify: Reducing Hallucination in Vision-Language Models with Retrospective Resampling

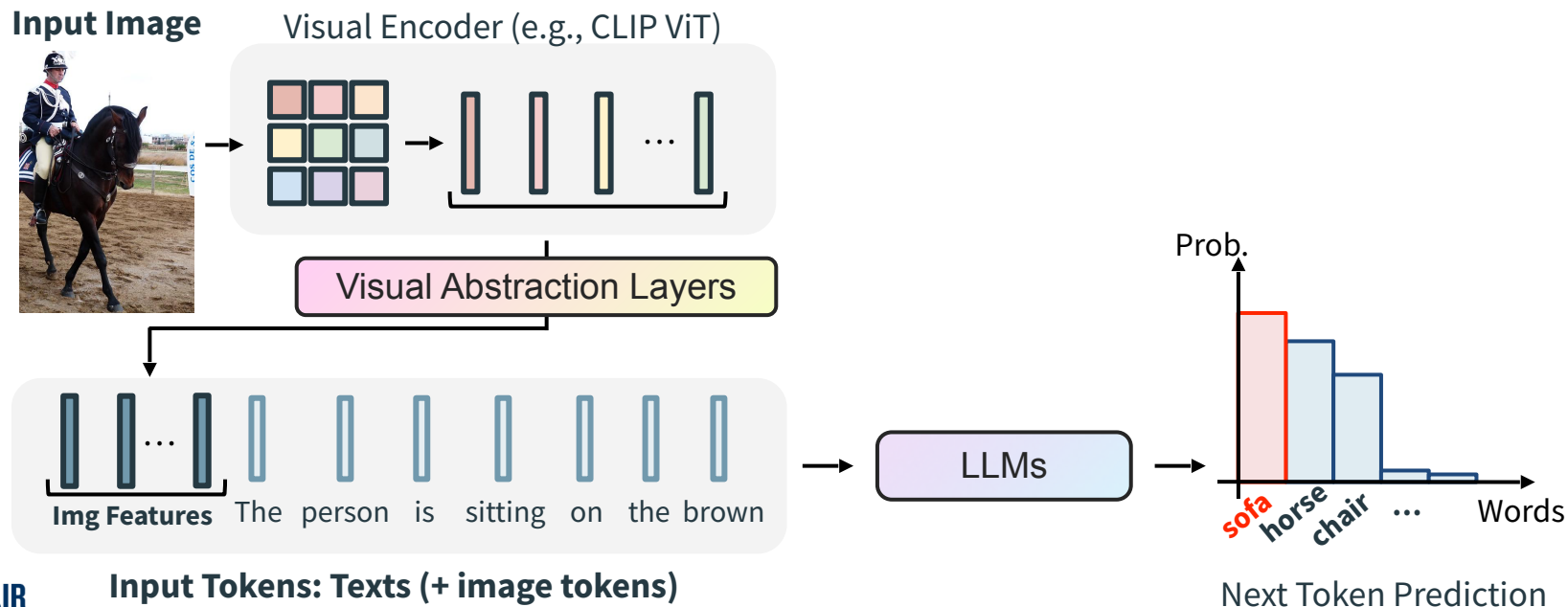
Tsung-Han (Patrick) Wu, Heekyung (Anne) Lee, Jiaxin Ge,  
Joseph E. Gonzalez, Trevor Darrell, David M. Chan

UC Berkeley

NeurIPS 2025

# VLMs Suffer from Hallucinations

**Visual Hallucinations:** Describing nonexistent objects or concepts in the image, usually due to *training data biases or strong language priors*



# Prior Methods



User “Describe this image.”



VLM “The boy is sharing his umbrella...”



Preventing hallucinations is hard, but detection isn't — we need a **verifier** after the fact.

## ① Generative Adjustment

P(token)

...playing

- [ball]
- a [game]
- [frisbee] with a dog.

# Prior Methods



User “Describe this image.”



VLM “The boy is sharing his umbrella...”



Preventing hallucinations is hard, but detection isn't — we need a **verifier** after the fact.



Verifying after full outputs is slow and mostly ends in refusal — we need **correction** instead!

## a Generative Adjustment

P(token)

...playing

- [ball]
- a [game]
- [frisbee] with a dog.

## b Post-Hoc Verification

...playing **frisbee** with a dog



# Prior Methods



User “Describe this image.”



VLM “The boy is sharing his umbrella...”



Preventing hallucinations is hard, but detection isn't — we need a **verifier** after the fact.



Verifying after full outputs is slow and mostly ends in refusal — we need **correction** instead!

## Why not both?

### Ⓐ Generative Adjustment

$P(\text{token})$

...playing

- [ball]
- a [game]
- [frisbee] with a dog.

### Ⓑ Post-Hoc Verification

...playing **frisbee** with a dog



# REVERSE-VLM

## REtrospective VERification and SElf-Correction

# Generate, but Verify: The birth of REVERSE-VLM

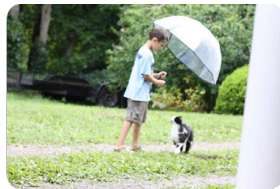
A single VLM that can not only **generate**, but **verify and corrects** themselves on-the-fly  
→ towards robust, controllable, and interpretable systems!



# How to Enable This?

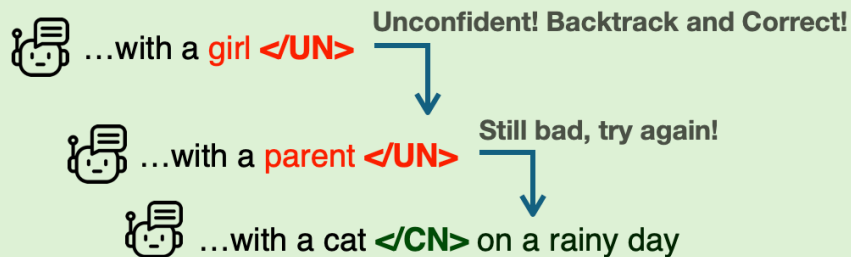


User “Describe this image.”



VLM “The boy is sharing his umbrella...”

(1) **Train** a model that can do explicit confidence estimation



(2) Do **retrospective resampling** inference for multi-round correction



# 1. SFT Dataset Construction

**Source: 665K LLaVA-SFT Dataset**

<SPAN>: Noun Phrase Opening

</CN>: Confident Token

</UN>: Unconfident Token



Human: What feature can be seen on the back of the bus?

GPT: The back of the bus features an advertisement.

Noun Phrase Extraction & Tagging

GPT: <SPAN>The back</CN> <SPAN>of the bus</CN> features  
<SPAN>an advisement</CN> .

**Positive Data**



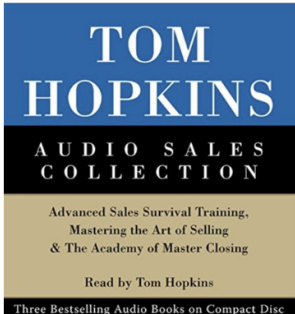
Negative Phrase Augmentation

GPT: <SPAN>The back</CN> <SPAN>of the bus</CN> features <SPAN>a window</UN> .

**Negative Data**

# 1. SFT Dataset Construction

## 1.3M Open-Source Datasets on 🧐

	COCO/train2017/	Captioning Question	VQA task
Image			
Question	"How many total baseball players are shown in the image?"	"Describe this image in your own words."	"Who wrote this book?"
Pos Answer	"There are <span>&lt;SPAN&gt;three baseball players&lt;/CN&gt;</span> shown <span>&lt;SPAN&gt;in the image&lt;/CN&gt;</span> .,"	"The image features <span>&lt;SPAN&gt;an old military aircraft&lt;/CN&gt;</span> <span>&lt;SPAN&gt;on display&lt;/CN&gt;</span> ..."	" <span>&lt;SPAN&gt;Tom Hopkins&lt;/CN&gt;</span> "
Neg Answer	"There are <span>&lt;SPAN&gt;five soccer players&lt;/UN&gt;</span> "	"The image features <span>&lt;SPAN&gt;a modern commercial airplane&lt;/UN&gt;</span> "	" <span>&lt;SPAN&gt;John Steinbeck&lt;/UN&gt;</span> "
	Number	Attribute	Object

## 2. Hallucination-Aware Training



Model needs to:

1. Do standard next token prediction
2. Avoid hallucination modeling
3. Learn to model confidence with `</CN>` and `</UN>`

`<SPAN>`The back`</CN>` `<SPAN>`of the bus`</CN>` features `<SPAN>`an advertisement`</CN>`

`<SPAN>`The back`</CN>` `<SPAN>`of the bus`</CN>` features `<SPAN>`a window`</UN>`

■ Model trained to predict these tokens

■ Model **ignores** these tokens during training

### 3. Retrospective Resampling Inference



<SPAN>The back</CN>

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

### 3. Retrospective Resampling Inference



<SPAN>The back</CN> <SPAN>of the bus</CN> features

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

### 3. Retrospective Resampling Inference



<SPAN>The back</CN> <SPAN>of the bus</CN> features <SPAN>a window</UN>

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0  
0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.20

<UN>: 0.87

### 3. Retrospective Resampling Inference



<SPAN>The back</CN> <SPAN>of the bus</CN>

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

### 3. Retrospective Resampling Inference



<SPAN>The back</CN> <SPAN>of the bus</CN> features <SPAN>an advertisement</CN>

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0

<UN>: 0.0



# Summary: Retrospective Resampling

## ① User Query



User

“Describe this image.”



# Experimental Results

# SOTA on Captioning & Open-ended VQAs

## Captioning Tasks

Base VLM	Method Type	Method	CHAIR-MSCOCO		AMBER-G			
			CHAIR <sub>i</sub> (↓)	CHAIR <sub>s</sub> (↓)	CHAIR (↓)	Cover (↑)	Hall (↓)	Cog (↓)
LLaVA-v1.5 7B [35]	None		15.4	50.0	7.8	51.0	36.4	4.2
		VCD [28]	14.9	48.6	-	-	-	-
	Gen-Adjust	OPERA <sup>‡</sup> [23]	14.6	47.8	7.3	49.6	32.0	3.5
		DoLA <sup>†</sup> <sup>‡</sup> [16]	14.1	51.6	7.6	51.6	36.0	4.0
		AGLA [3]	14.1	43.0	-	-	-	-
		MEMVR [58]	13.0	46.6	-	-	-	-
	w/ Train	EOS [55]	12.3	40.2	5.1	49.1	22.7	2.0
		HALVA [41]	11.7	41.4	6.6	<b>53.0</b>	32.2	3.4
		HA-DPO [56]	11.0	38.2	6.7	49.8	30.9	3.3
	Post-hoc Refine	Woodpecker <sup>†</sup> [53]	14.8	45.8	6.9	48.9	30.4	3.6
	Combination	<b>REVERSE</b> <sub>(<math>\tau=0.003</math>)</sub>	10.3	37.0	6.0	52.2	30.4	3.0
		<b>REVERSE</b> <sub>(<math>\tau=0.0003</math>)</sub>	<b>6.1</b>	<b>13.6</b>	<b>4.0</b>	26.9	<b>10.2</b>	<b>0.9</b>

12% Gain

+ Multiple Models

+ Multiple Tasks

## MM-Hal

Base VLM	Method	Score (↑)	Hall. Rate (↓)
LLaVA-MORE 8B	None <sup>†</sup>	2.50	0.53
	DoLA <sup>†</sup> [15]	2.54	0.51
	Woodpecker <sup>†</sup> [51]	2.28	0.58
	REVERSE <sub>(<math>\tau=0.003</math>)</sub>	2.28	0.54
	REVERSE <sub>(<math>\tau=0.0003</math>)</sub>	<b>2.93</b>	<b>0.40</b>

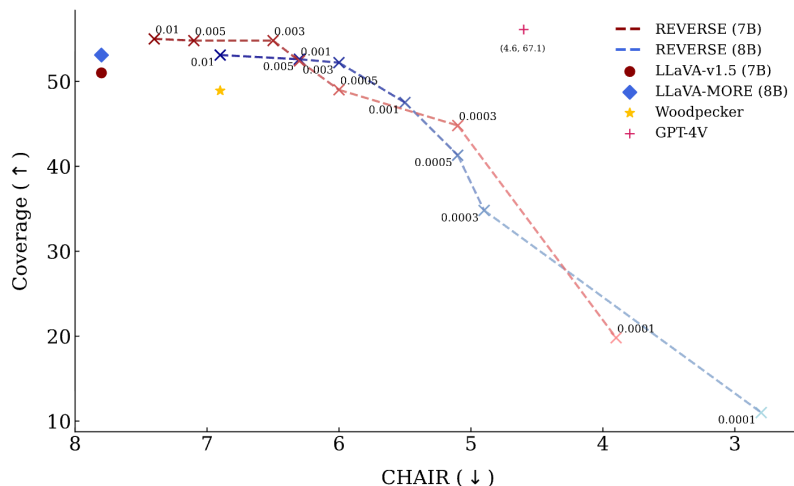
## HaloQuest 34% Gain

Method	Avg. Acc. (↑)	FP Acc.	VC Acc.	IC Acc.
<b>Qwen2.5-VL<sup>FT</sup> 3B</b>				
None <sup>†</sup>	33.5	25.4	<b>51.6</b>	26.4
DoLA <sup>†</sup> [15]	27.4	16.5	51.1	19.0
REVERSE <sub>(<math>\tau=0.01</math>)</sub>	<b>45.1</b>	<b>42.9</b>	41.8	<b>55.5</b>

# Towards Efficient Corrections & Controllable VLMs

## Studies on AMBER-G Dataset

# Rounds (N)	0	5	10	20	50
CHAIR (↓)	7.8	7.1	6.8	6.7	6.0
#Tokens (%)	1.00×	1.75×	2.05×	2.63×	3.05×



**15% gain** from **50** round corrections but only **3.05x** more tokens

Tuning the **threshold ( $\tau$ )** can control the trade-off between **expressiveness & hallucinations**

We can beat GPT-4V on the CHAIR metric with low threshold, making VLMs conservative!

# Conclusions

# Takeaways

→ Current VLMs are still prone to visual hallucinations.

# Takeaways

- Current VLMs are still prone to visual hallucinations.
- When verification is easier, we can insert retrospective reasoning tasks to encourage the model to “think twice” about its process.

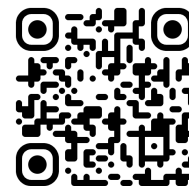
# Takeaways

- Current VLMs are still prone to visual hallucinations.
- When verification is easier, we can insert retrospective reasoning tasks to encourage the model to “think twice” about its process.
- We present REVERSE, the first and effective hallucination reduction method unifying the generation adjustment and post-hoc verification.



# Takeaways

- Current VLMs are still prone to visual hallucinations.
- When verification is easier, we can insert retrospective reasoning tasks to encourage the model to “think twice” about its process.
- We present REVERSE, the first and effective hallucination reduction method unifying the generation adjustment and post-hoc verification.
- REVERSE is not endpoints but starting points that invite the community to explore the “generate-but-verify” paradigm.



# Takeaways

- Current VLMs are still prone to visual hallucinations.
- When verification is easier, we can insert retrospective reasoning tasks to encourage the model to “think twice” about its process.
- We present REVERSE, the first and effective hallucination reduction method unifying the generation adjustment and post-hoc verification.
- REVERSE is not endpoints but starting points that invite the community to explore the “generate-but-verify” paradigm.

**Thanks for listening!**