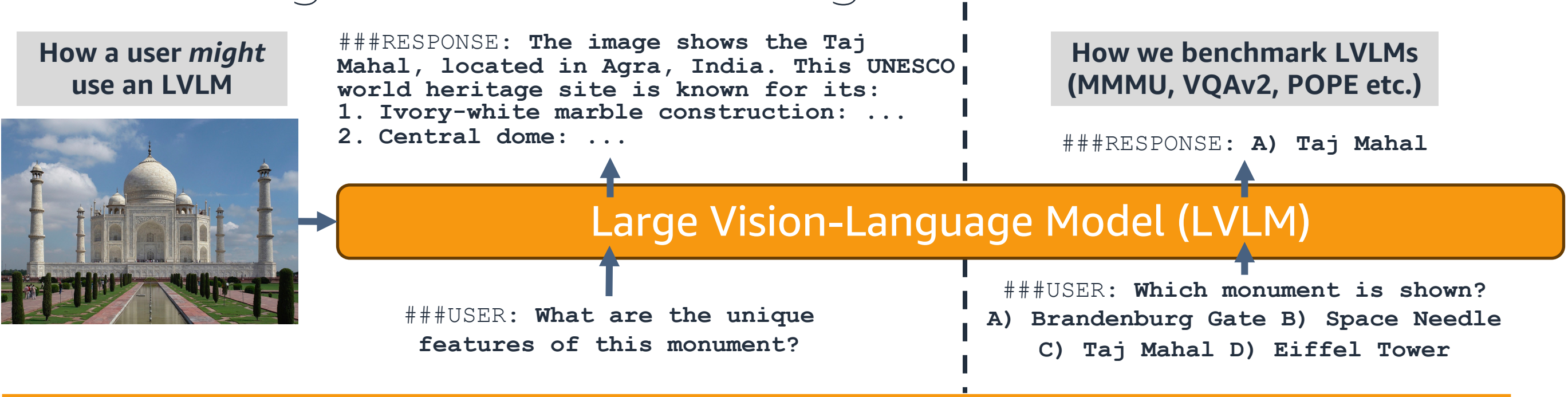


# THRONE: An Object-based Hallucination Benchmark for Free-Form Generations of Large Vision-Language Models

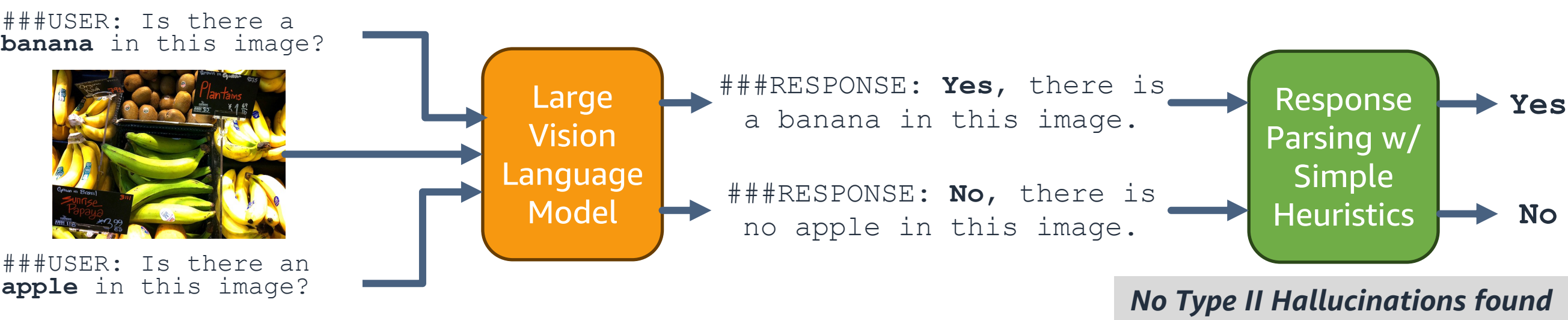
Prannay Kaul<sup>1</sup> Zhizhong Li<sup>2\*</sup> Hao Yang<sup>2</sup> Yonatan Dukler<sup>2</sup> Ashwin Swaminathan<sup>2</sup> C. J. Taylor<sup>2</sup> Stefano Soatto<sup>2</sup>  
<sup>1</sup>VGG, University of Oxford <sup>2</sup>AWS AI Labs

## LVLM Usage vs Benchmarking

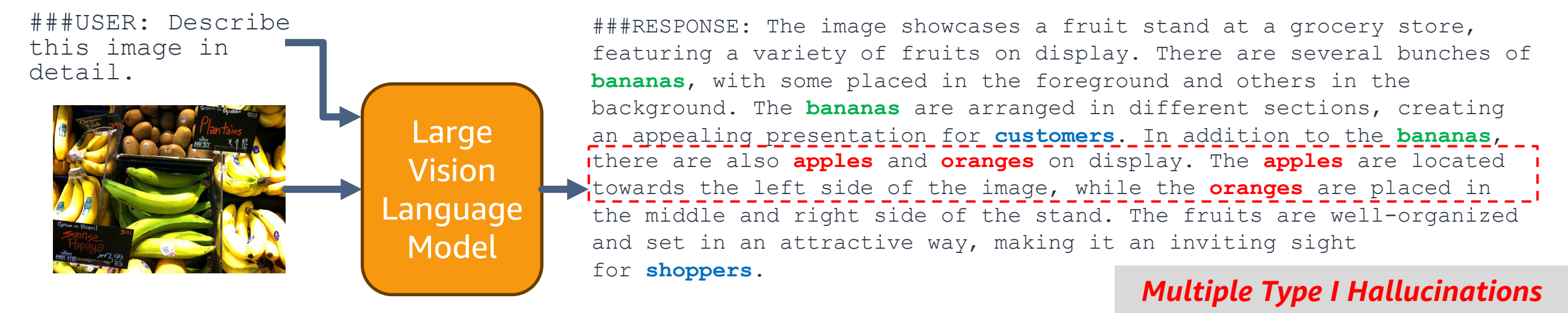


## Measuring Hallucinations (POPE [1] vs THRONE)

- Current methods evaluate **Type II Hallucinations** in responses to yes/no questions

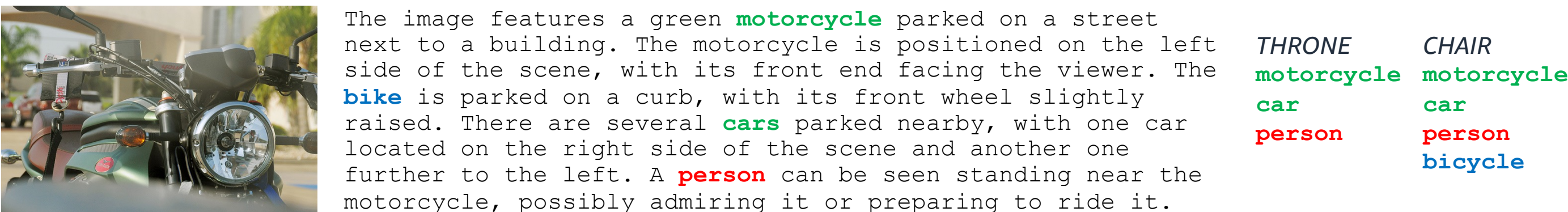


- THRONE evaluates **Type I Hallucinations** in free-form responses to *open-ended* prompt
- Type I and II performances are barely correlated: Spearman's  $\rho = 0.2$



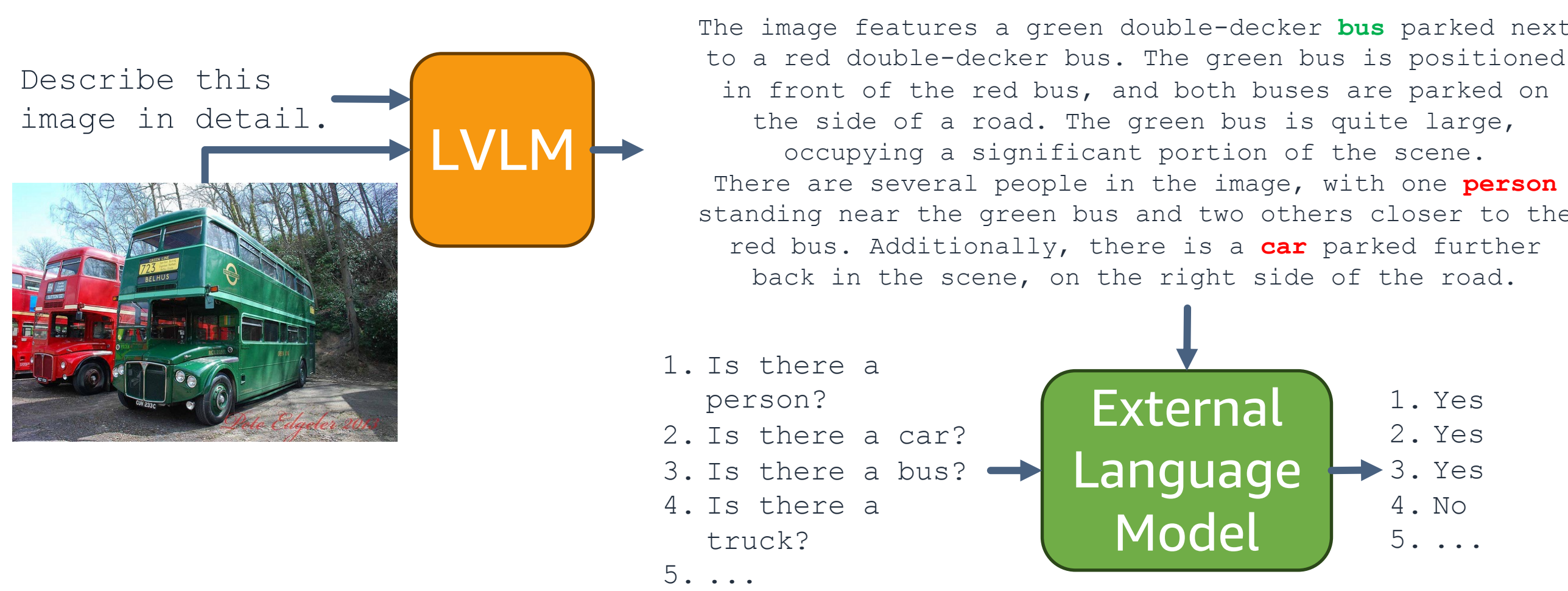
## Measuring Hallucinations (CHAIR [2] vs THRONE)

- CHAIR addresses Type I hallucinations *but uses simple text-match*
- Many abstractly referenced nouns + synonyms are incorrectly classified by CHAIR (**blue**)
- We estimate that THRONE makes **half** the errors of CHAIR (4% vs 9%)

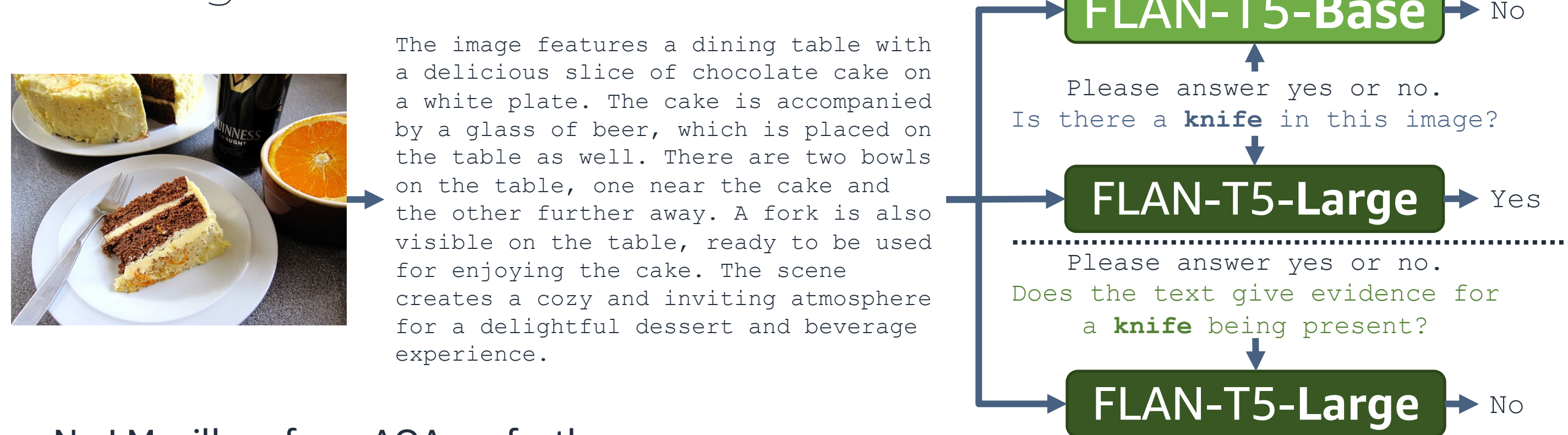


## THRONE Method

- Extract long-form responses from a given LVLM**
  - We use a single prompt: Describe this image in detail.
- Abstractive question answering (AQA) on the LVLM response using a language model (LM)**
  - We only consider *object existence* in the responses w.r.t known classes (e.g. COCO classes)
- Collate LM responses and calculate informative metrics**
  - LM is prompted to give simple responses so no parsing required



## Ensuring Robustness



- No LM will perform AQA perfectly
  - Incorrect judgements will be made on an LVLM response regarding object existence
- Ensemble multiple *lightweight* LMs (FLAN-T5) – **modeling ensembling**
- Ensemble multiple prompts *to the LMs* – **prompt ensembling**

### References

[1] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In Proceedings of the Conference on Empirical Methods in Natural Language, 2023  
 [2] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In Proceedings of the Conference on Empirical Methods in Natural Language, pages 4035–4045, 2018.

## THRONE Results

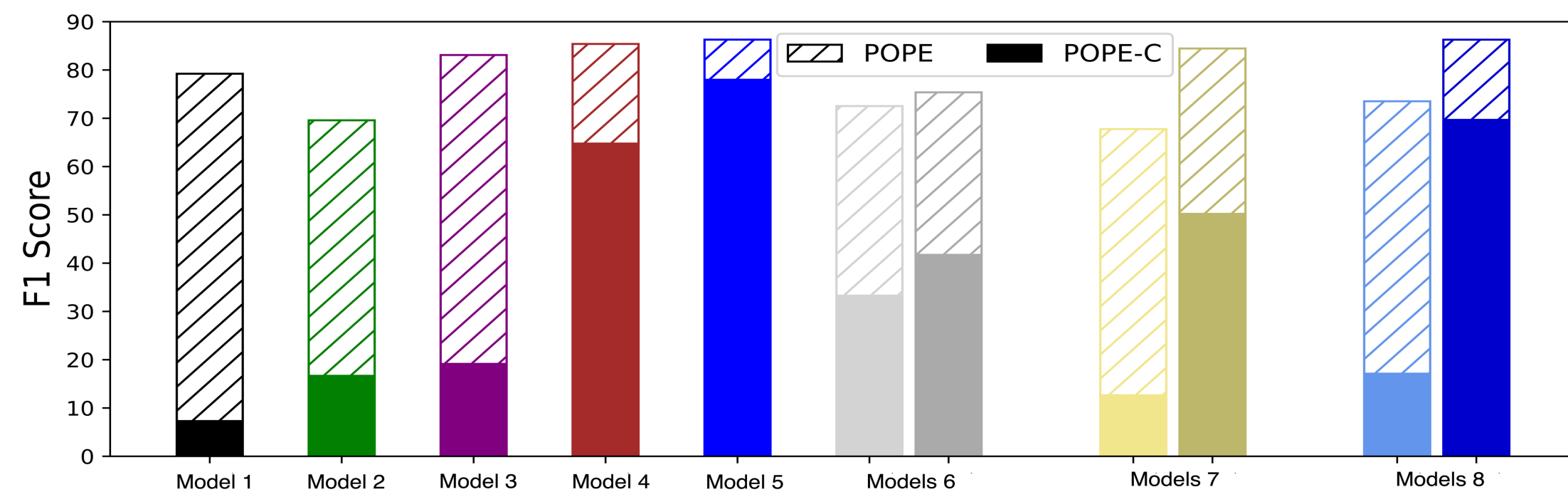
- THRONE performed for *object existence* of all classes,  $C$ , in a dataset across many images,  $I$
- Predicted object existence:  $\hat{Y} \in \{0, 1\}^{|I| \times |C|}$  GT object existence:  $Y \in \{0, 1\}^{|I| \times |C|}$
- This allows Overall and Classwise Precision/Recall to be calculated.
- $F^{0.5}$ -Score calculated to prioritize precision over recall (Hallucinations  $\subset$  False Positives)



- THRONE evaluates Type I hallucinations in free-form responses**
- THRONE evaluation is automatic and does not use of subjective scoring by LLMs**
- THRONE uses open-source lightweight language models (ensemble), improving comprehension, ease-of-use, and reproducibility**

## Issues with gauging Type II Hallucinations

- Hallucinations  $\subset$  False Positives
- Need to search through many true negatives to reveal the complete picture
- Completing POPE for all classes reveals the true extent of Type II hallucinations (POPE-C)



## Improved Baseline via Augmentation

- We introduce a simple data augmentation which forces LLaVA models to enumerate objects

Model	Object Enumeration Data	THRONE			POPE			POPE-C		
		$P_{CLS}$	$R_{CLS}$	$F^{0.5}_{CLS}$	$P$	$R$	$F^1$	$P$	$R$	$F^1$
LLaVA-v1.5	<b>x</b>	69.9	56.4	66.8	81.9	90.8	86.1	58.7	85.7	69.7
	COCO	<b>87.2</b>	76.6	<b>84.9</b>	88.6	<b>85.3</b>	<b>87.0</b>	58.9	<b>87.5</b>	70.4
	COCO + VG	86.1	<b>77.0</b>	84.1	<b>89.8</b>	83.7	86.7	<b>64.5</b>	86.1	<b>73.7</b>