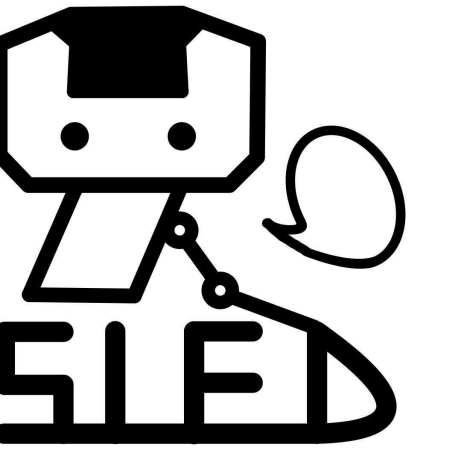




Transparent and Coherent Procedural Mistake Detection

Shane Storks, Itamar Bar-Yossef, Yayuan Li, Zheyuan Zhang,
Jason J. Corso, & Joyce Chai (University of Michigan)

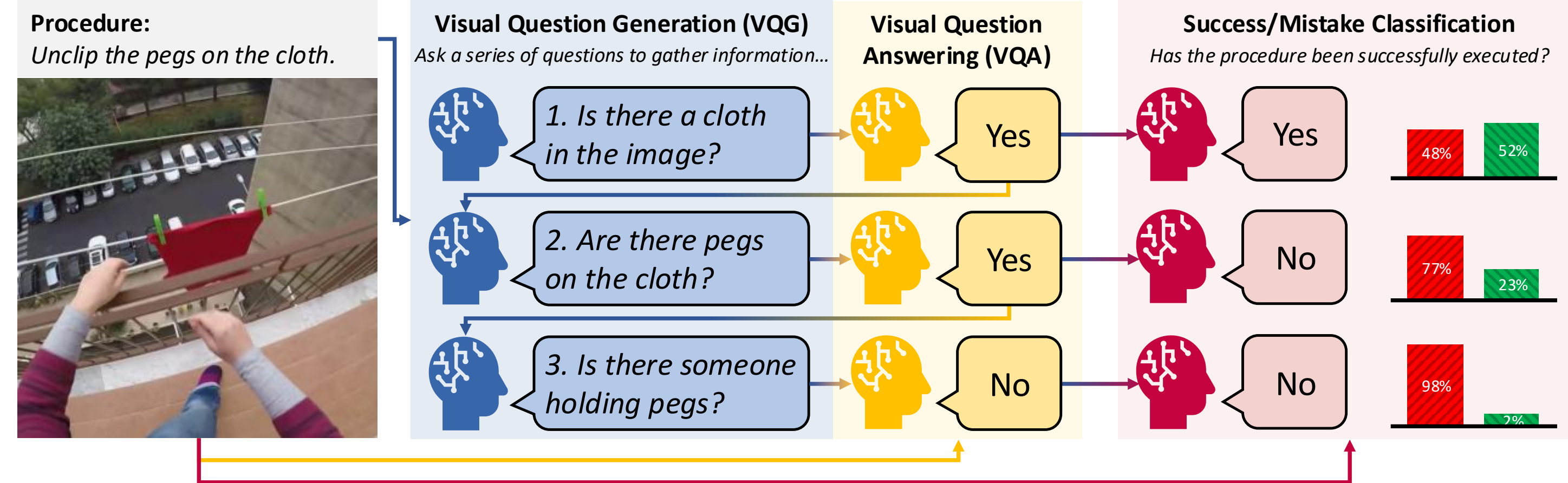


MOTIVATION

Procedural mistake detection (PMD) requires classifying whether a human (seen through egocentric video) has successfully executed a task (specified by a procedural text). Despite significant efforts, VLM performance in the wild is nonviable, and underlying knowledge and reasoning processes are opaque.

COHERENT PROCEDURAL MISTAKE DETECTION

We reformulate PMD to require a self-reflective dialog rationale from VLMs:



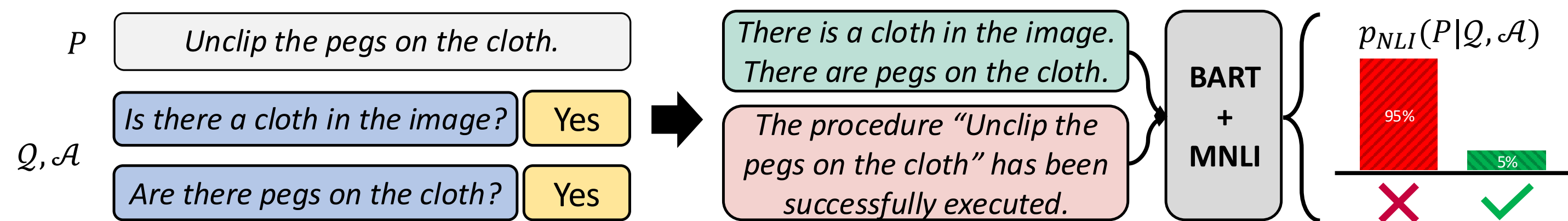
We generate diverse video frame mistake detection data from Ego4D [1]:



Example Type	Train	(Sample)	Validation	(Sample)	Test	(Sample)
Success	42,013	5,000	13,058	250	18,057	1000
Mistake	99,401	5,000	25,423	250	34,182	1000
Incomplete	15,057	755	4,908	51	6,545	194
Wrong V	11,780	604	2,694	31	3,747	108
Wrong N	36,434	1,853	8,914	87	11,843	344
Wrong V & N	36,130	1,788	8,907	81	12,047	354

RATIONALE COHERENCE METRICS

To evaluate whether evidence collected from VLMs suggests success, we leverage fine-tuned natural language inference (NLI) models:



We use the NLI model p_e to measure the **relevance** of a question Q' to the success of a procedure P , given previous questions Q and answers A :

$$\text{Rel}(Q'|T, Q, A) = |p_e(T|Q \cup Q', A \cup \text{No}) - p_e(T|Q \cup Q', A \cup \text{Yes})|$$

Relevance is summarized by example through a mean over questions:

$$\frac{1}{n} \sum_{i=1}^n \text{Rel}(Q_i|T, \{Q_j: j < i\}, \{A_j: j < i\})$$

We also measure the **informativeness** of a predicted answer A' for Q' :

$$\text{Inf}(A'|Q', T, Q, A) = 1 - H(p_e(T|Q \cup Q', A \cup A'))$$

Reference-adjusted informativeness Inf^* is negated if the most likely success label in p_e disagrees with the ground truth label y^* . It is summarized by example through the maximum informativeness achieved:

$$\max_{1 \leq i \leq n} \text{Inf}^*(A_i|Q_i, T, \{Q_j: j < i\}, \{A_j: j < i\}, y^*)$$

REFERENCES

- [1] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, et al. Ego4D: Around the World in 3,000 Hours of Egocentric Video. In *IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [2] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023.
- [3] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://arxiv.org/abs/2305.18290>.

LINKS



PAPER



CODE



ME



LAB

EXPERIMENTAL RESULTS

We apply 2 interventions to the selection of candidate questions generated by LLaVA-1.5 [2] through a greedy beam search:

1. **Coherence-based re-ranking** of candidate questions
2. **In-context learning (ICL)** from 20 sets of human-written questions

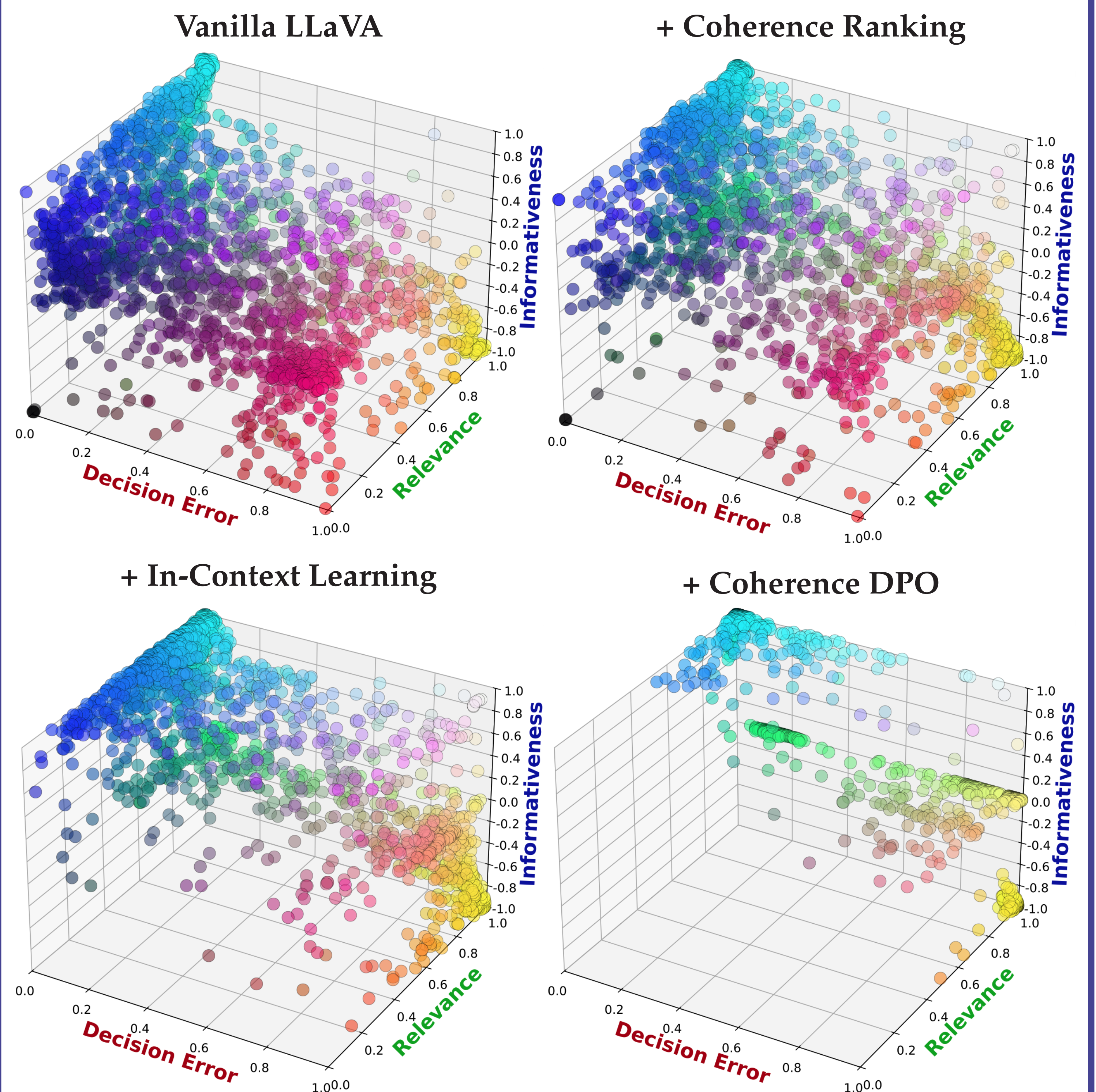
Ranking	ICL	Acc. \uparrow	Rel. \uparrow	Inf. \uparrow	# Iter. \downarrow	Info. Gain \uparrow
Likelihood	\times	60.7	40.3	.259	3.25	.435
Likelihood	\checkmark	61.8	36.5	.272	3.34	.429
Coherence	\times	61.4	66.5	.321	3.06	.540
Coherence	\checkmark	67.8	75.5	.464	3.46	.663

We then **fine-tune LLaVA** for question generation with our coherence metrics, using DPO [3] over question pairs generated from training data:

Ranking	ICL	Acc. \uparrow	Rel. \uparrow	Inf. \uparrow	# Iter. \downarrow	Info. Gain \uparrow
Likelihood	\times	62.2	75.7	.318	2.33	.617
Likelihood	\checkmark	63.7	58.5	.330	2.67	.548
Coherence	\times	62.3	92.2	.340	2.06	.719
Coherence	\checkmark	64.2	95.0	.304	1.81	.742

PERFORMANCE ANALYSIS

Our metrics provide global and local insights into VLM performance:



A: Pick up a sink brush from the kitchen slab.



Label: \checkmark Predicted: \checkmark

1. Is the sink brush in the person's hand? Yes

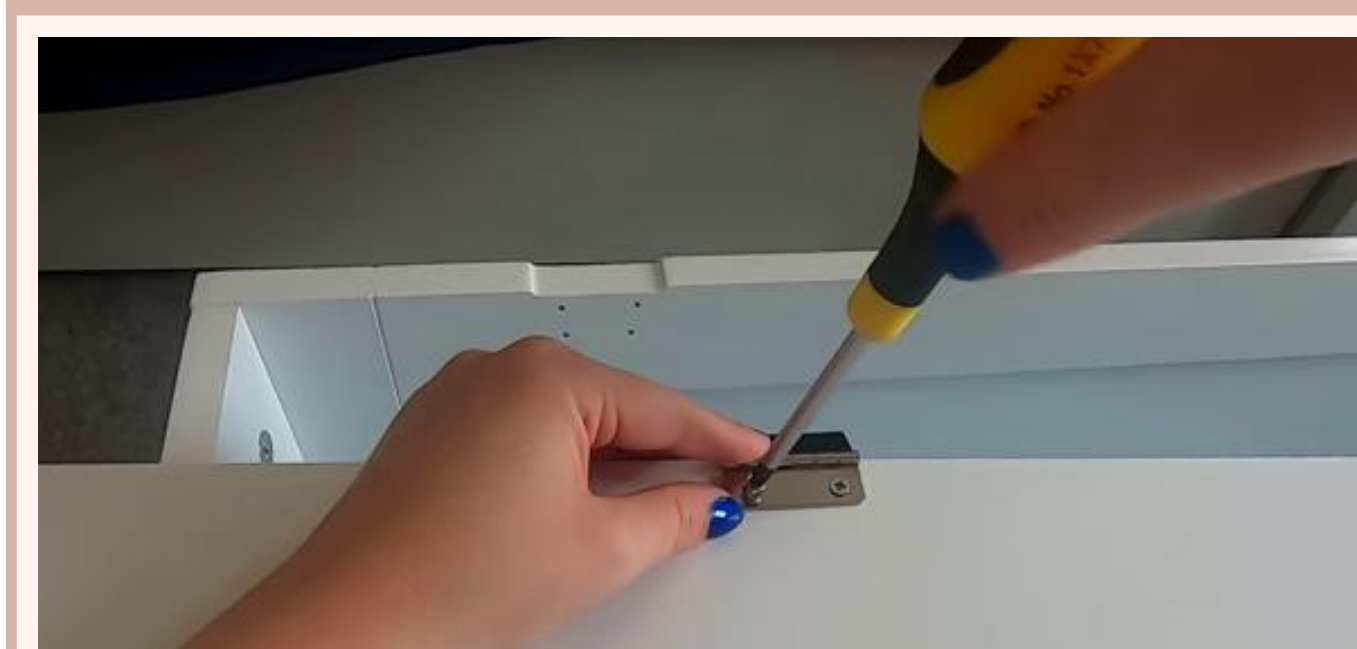
C: Put the trowel in a bin.



Label: \checkmark Predicted: \times

1. Is the trowel in a bin? No

B: Tighten the screw.

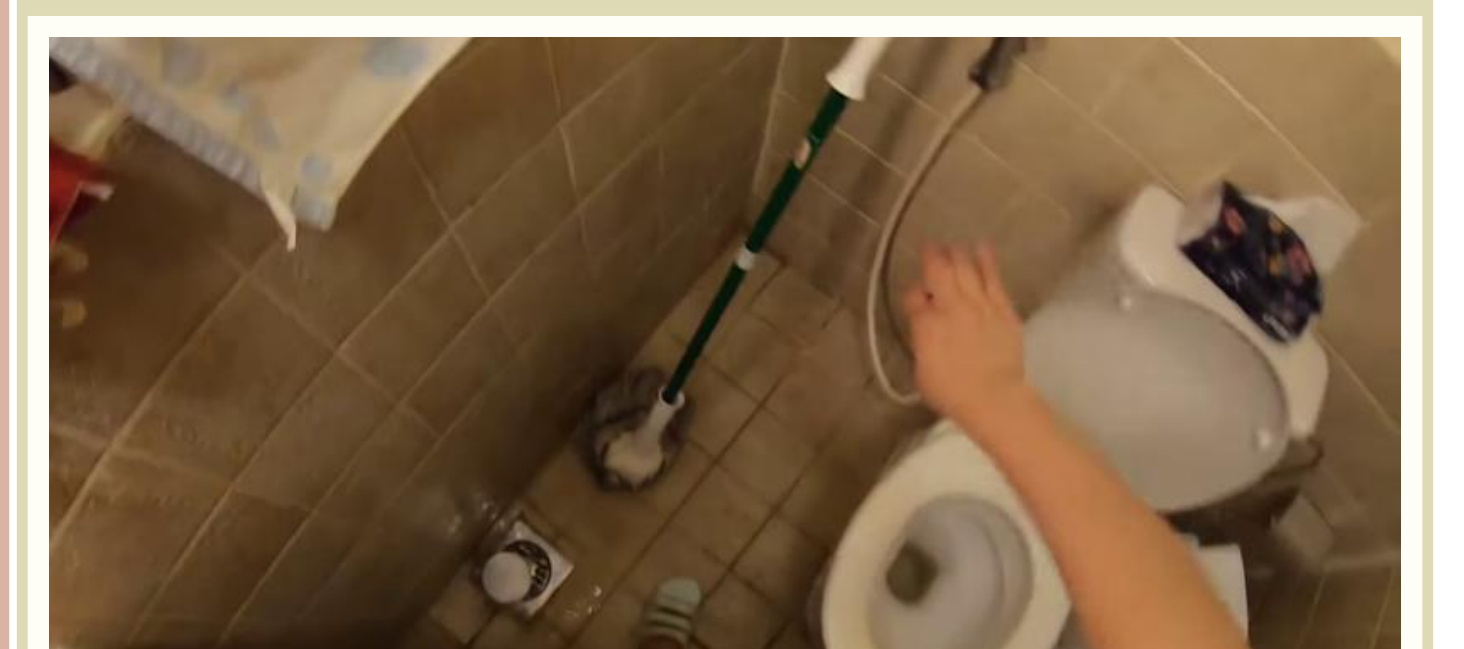


Label: \checkmark Predicted: \times

Rationale:

1. Is the person wearing gloves? No
2. Is the person wearing protective gear? No
3. Is the person wearing a mask? No

D: Put the bottle in the cabinet.



Label: \times Predicted: \checkmark

1. Is the bottle in the cabinet? Yes