Can Large Vision-Language Models Correct Semantic UNIVERSITY OF VECTOR INSTITUTE Grounding Errors By Themselves? / Yes*, *come to our poster!

<u>Yuan-Hong Liao¹</u>, Rafid Mahmood^{2,3}, Sanja Fidler^{1,2}, David Acuna²

Can LLMs Self-Correct w/o fine-tuning?

Prior works say,

"In the context of reasoning, our research indicates that LLMs struggle to (instrinsically) self- correct their responses without external feedback, and at times, their performance even degrades after self-correction.", Huang et al., ICLR 2024 "There is still no consensus on the question of when LLMs can correct their own mistakes, as recent studies also report negative results.", Kamoi et al., EMNLP 2024

Our Research Questions: Can VLMs Self-Correct? If so, under what context?

Self-correction := VLMs receive feedback + VLMs provide feedback

Task: Semantic grounding





Our Idea:

- Agentic: Use VLMs themselves as verifiers.
- Simplicity -> Reliable feedback: Verification is **easier** than generation
- Performances: Achieve test-time scaling! (see our results 🚀)
- Approach: Agentic, Training-free - Task: Semantic grounding
- VLMs: Open-source VLMs, GPT-4V, and GPT-40



VLM

Verification

VLM

Textual

feedback

University of Toronto & Vector Institute, ² NVIDIA, ³ University of Ottawa





Predictions

+10%LLaVA-1.5 ViP-LLaVA CogVLM

tion	35.86	35.86	15.98
Feedback	41.04	40.36	16.25
abel Feedback	94.8	74.99	77.04

How to give feedback? Visual markers 🔘 works the best

Yes	Visual marks	45.38	45.21	19.46
Yes	SoM	42.41	44.53	18.64
Yes	No	43.3	42	18.25
No	No	41.04	40.36	16.25
No	No	35.86	35.86	15.98
o-shot CoT	Visual Prompt	LLaVA-1.5	ViP-LLaVA	CogVLM

How to provide binary feedback? All better than intrinsic SC. Best prompts depends on VLM.

	Visual prompt	LLaVA-1.5	ViP-LLaVA	CogVLM
	N/A	51.12	48.19	21.87
	Visual marks	56.16	60.47	39.16
1	RoI crop	61.71	58.18	40.68
	Visual marks + RoI crop	61.14	59.6	39.79



Key results - Test-time scaling: Trade performances with #tokens - Better VLMs, better gains: Gains of GPT-40 > Gains of GPT-4V

Interested in System-2 thinking in VLMs?

Ch	eck our	re
1.	EM	N
	spatial	re
	trainin	g
2.	arX	iv'
	traces	fc
	svnthe	>+i

Main Results

related papers: LP'24: Enhancing quantitative reasoning with no extra training. -free, spatial-reasoning 25: Synthesizing System-2 reasoning or System-1 Perception

synthetic-data-generation, cognitive-behaviors



