RE-ALigning Language to Visual Objects with an Agentic Workflow



Yuming Chen, Jiangyan Feng, Haodong Zhang, Lijun Gong, Feng Zhu, Rui Zhao

Qibin Hou#, Ming-Ming Cheng, Yibing Song#





Key Point



Agentic Workflow: from Assistant to Data Flywheel

- Agent is not only a *simple assistant* but also can establish workflows that serve as **flywheel** to sustaining **high-value data assets** across AI industries
- In this paper, we show an application of agentic workflow to demonstrate its potential...



Manus^[1]: Recent Viral Agent Assistant

[1] Manus: https://manus.im/



Our agentic workflow proposed to provide data with **high Visual-Language alignment quality**



• Method

- Example
- Experiments



Language-based Object Detection

• The essence of *Object Detection* is to Align Language Data to Visual Objects





Visual-Language Alignment

• **Quantity & Diversity of Visual-Language data for alignment is important!**



Training Data Quantity:

- *OWL*^[2]: **10B** images
- *GDINO*^[3]: **5.8M** images
- *APE-B*^[4]: **2.6M** images

Collection Challenge:

- Objects are **harder** to label than whole images
- The costs for large-scale data are **extremely high**

Expression Diversity

- Training Dataset: Simple
- Testing Dataset: **Complex**

0	mniLabel ^[5]		low			
Images	Method	AP	AP-categ	AP-descr		
	RegionCLIP [53]	2.7	2.7	2.6		
	Detic [55]	8.0	15.6	5.4	L	
Ξ	MDETR [20]	-	-	4.7		
A	GLIP-T [29]	19.3	23.6	16.4	L	
	GLIP-L [29]	25.8	32.9	21.2	L	
	FIBER-B [12]	25.7	30.3	22.3		
_					1	

• Introducing *VLM* for expression generation is promising, but...

Yuming Chen, www.fishworld.site

[1] Radford et al., Learning Transferable Visual Models From Natural Language Supervision. ICML.2021.
[2] Minderer et al., OWL-ViT: Open-World Object Detection with Vision Transformers. ECCV 2022.
[3] Liut et al. Grunding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. ECCV 2024



Hallucination of VLM

Small object and **General Prompt** cause *VLM* to misperceive objects, damaging **Quality**







<Prompt>: What is the object in the red box? <**Response>:** The object is a blue and red *hat worn by a* young girl.

<**Prompt>:** What is the object in the red box? <**Response**>: The object is *a flag* with a design that resembles the Union Jack.

<**Prompt>:** Describe the person in the red box. <**Response>:** The person is a woman with long hair wearing a *black top*.

<**Prompt>:** What is the color of person's clothes. <**Response>:** The person is wearing <u>a black and white</u> striped shirt, dark pants, and black shoes.

(d) With Specified Prompt

(a) For Small Object

(b) With Extended Crop

(c) With General Prompt

By adopting **appropriate operations**, hallucination can be alleviated



• Method

- Example
- Experiments

Method



Agentic Workflow

• To improve the alignment quality of large-scale VL data, Real-LOD involves cycles with

3 steps: *Planning*, *Tool use*, and *Reflection*



Method





object itself



object and surrounding areas !





whole image

Object Highlight



Object Crop

Extended Object Crop



designed for **Real-LOD**

After analyzing the **misalignment** caused by

hallucination, 5 states/actions are specifically



- Method
- Example
- Experiments

Example



Input



The second	Object Category	horse					
	Bbox Coordinates	[223.07, 267.79, 101.47, 59.39]					
	Image Caption						
P	The image depicts a group of people on horseback in a forested area. The person riding the horse in the red box is wearing a blue shirt and dark pants, and appears to be leading the group. The horse is brown and is walking on a						
	dirt path. The riders are w	vearing backpacks, suggesting they might be on a					
	visible through the canop of the image is one of out	by, indicating it's a sunny day. The overall atmosphere tdoor adventure and exploration.					

Can you specify if the raw expression that **brown horse walking on sand** is referring to the target object? If not, please modify it to a correct expression.

Output

The expression **brown horse walking on sand** is a wrong one, and a correct expression could be **brown horse walking on dirt path with rider in blue shirt and hat, leading the way**.

Example



First Cycle

Planning	Reasoning: The expression describes a brown horse that is walking on sand, but the image caption mentions the horse is walking on a dirt path, and there is no mention of sand in the caption. Based on this information, I reason the state is 4 where I am uncertain of object relation with surroundings. So I plan for Action 4 where I can crop a local object region in this image for VLM to further reperceive. The description of VLM supplements the current information for me to conduct a more thorough expression examination. Action: [{'Tool_name': 'VLM', 'Tool_parameters': {'Prompt': 'Questions:1. Is the horse walking on sand or on a dirt path?', 'Image editing': 'extended object crop.'}] Response: I will consult with a VLM to acquire supplement information that will enable me to evaluate the current expression more thoroughly.
😿 Tool use	VLM output: 1. The horse is walking on a dirt path, not on sand.
Reflection	Feedback : The expression { <u>brown horse walking on sand</u> } describes a wrong relation (on sand) of the horse.

Example



Second Cycle

X



ina	Reasoning : Based on the feedback, <u>I reason the state is 2 where I am certain the expression {<i>brown horse</i>} <i>walking on sand</i>} is wrong. So I plan for Action 2 where I use an LLM to rewrite the expression, and then check</u>
ing	the correctness of the new expression. Actions: [{'Tool_name': 'LLM', 'Tool_parameters': {'Prompt: Image Caption,, Response,'}}] Response: The expression { <i>brown horse walking on sand</i> } contains flaws, now I will first leverage the LLM to
	create an alternative expression for you.
use	LLM output : New Description: { <i>brown horse walking on dirt path with rider in blue shirt and hat, leading the way</i> }



Tool

Feedback: The expression {<u>brown horse walking on dirt path with rider in blue shirt and hat, leading the</u> <u>way</u>} describes the object category 'horse' and provides similar context to the caption.

Final Cycle

Planning	Reasoning: Based on the feedback, I reason the state is 1 where I am certain the new expression { <i>brown horse walking on dirt path with rider in blue shirt and hat, leading the way</i> } is correct. So I plan for Action 1 to stop. Actions: [] Response: The expression { <i>brown horse walking on sand</i> } is a wrong one, and a correct one could be { <i>brown horse walking on dirt path with rider in blue shirt and hat, leading the way</i> }.
----------	--



- Method
- Example
- Experiments



Location



15



Ablation Study on OmniLabel Benchmark

Test subset	Training data type	#Img	AP-des	AP-des-pos	AP-des-S	AP-des-M	AP-des-L
	А	933k	21.2	59.4	31.3	21.1	18.6
COCO	В	695k	22.2	59.4	32.4	21.9	19.4
	С	863k	24.2	59.6	35.2	24.2	21.1
	А	933k	27.6	43.1	39.8	25.5	17.9
O365	В	695k	28.5	43.7	40.2	26.2	18.5
	С	863k	32.4	48.5	47.5	30.0	21.3
OI	А	933k	30.5	43.0	37.2	30.3	23.2
	В	695k	31.4	43.7	38.1	31.2	24.0
	С	863k	33.5	44.9	42.2	32.9	24.8

A: "raw expression"; B: "raw expression with filter"; C: "raw expression with filter + Real-LOD";



OmniLabel Benchmark

LOD Method	#Img	AP-des	AP-des-pos	AP-des-S	AP-des-M	AP-des-L
GLIP (Swin-L)	17.5M	21.2	33.2	37.7	18.9	10.8
mm-GDINO (Swin-B)	12M	20.8	33.1	31.9	19.8	14.1
FIBER (Swin-B)	4M	22.3	34.8	38.6	19.5	12.4
Real-Model (Swin-B)	0.18M	36.5	52.1	54.4	33.2	25.5

DOD Benchmark

LOD Method	#Img	Full	Presence	Absence
GDINO (Swin-B)	5.8M	20.1	20.7	22.5
mm-GDINO (Swin-B)	12M	24.2	23.9	25.9
APE-B (ViT-L)	2.6M	30.0	29.9	30.3
Real-Model (Swin-B)	0.18M	34.1	34.4	33.2



RefCOCO/g/+ Benchmark

LOD Mothed	#Img	RefCOCO			RefCOCO+			RefCOCOg		
LOD Method		val	testA	testB	val	testA	testB	val-u	test-u	
APE-A (ViT-L)	2.0M	34.2	34.8	36.1	33.5	32.3	36.0	38.9	40.5	
Real-Model (Swin-B)	0.18M	74.0	79.6	66.0	76.4	83.1	68.5	80.8	81.2	
GDINO* (Swin-B)	5.8M	-	-	-	73.6	82.1	64.1	78.3	78.1	
APE-B* (ViT-L)	2.6M	84.6	89.2	80.9	76.4	82.4	66.5	80.0	80.1	
Real-Model* (Swin-B)	0.24M	91.3	93.1	88.0	85.4	90.3	78.6	88.4	89.0	

OVDEval Benchmark

* indicates that the model employs RefCOCO/g/+ for training

LOD Method	#Img	color	material	position	relationship	negation	avg
GLIP (Swin-L)	17.5M	6.7	15.8	48.1	33.2	51.8	31.1
OmDet (ConvNext-B)	1.1M	24.5	22.5	47.7	51.8	55.8	40.4
FIBER (Swin-B)	4M	9.4	17.7	48.1	33.2	58.1	33.3
Real-Model (Swin-B)	0.18M	25.7	22.5	59.3	41.9	68.4	43.6



Visualization

Query: "This item is used to keep warm in colder weather."



Query : "Woman in wedding dress next to a man in suit."



Query : "Pillow placed at the head of the bed."





Query : "These two people each have a pink surfboard."



Query : "The fire extinguisher on the left."



Query : "Cows that are laid down."







(a) GLIP-L

(b)

(b) APE-B (c) mm-GDINO

(d) Our



Thanks For Watching!



GitHub



Main Page



Wechat



Yuming Chen, www.fishworld.site



chenyuming@mail.nankai.edu.cn