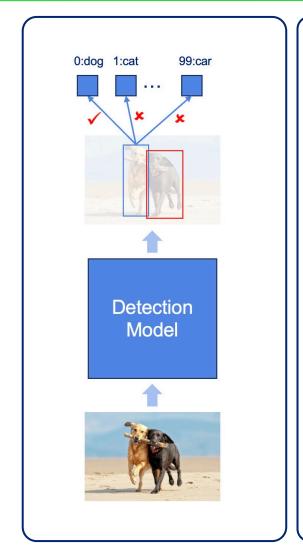


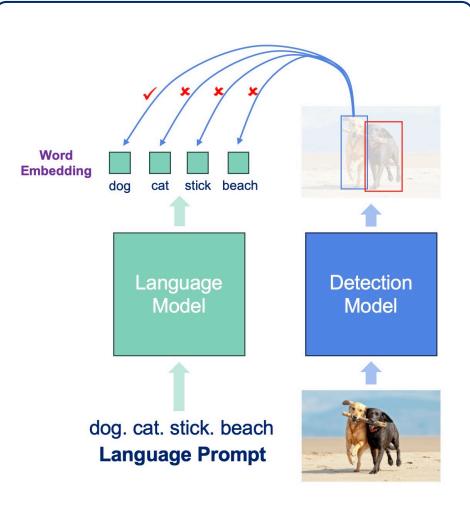
Taming Multimodal LLM for Object Perception

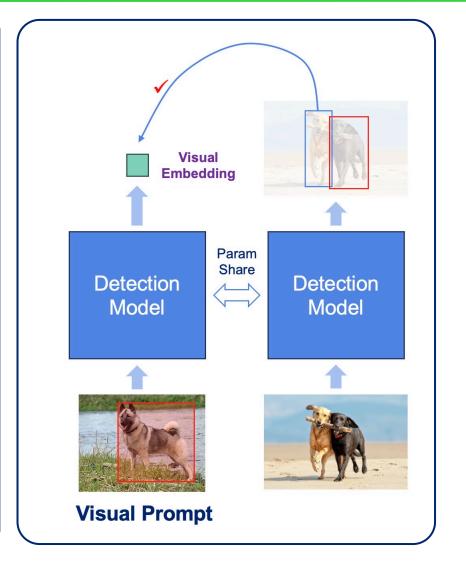
蒋擎

7-15









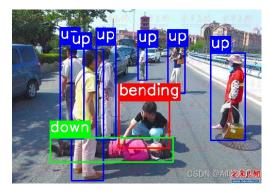
Closed-Set Open-vocabulary

DETR Grounding DINO

Visual Prompt

idea

most detection entities can be described in language









摔倒检测

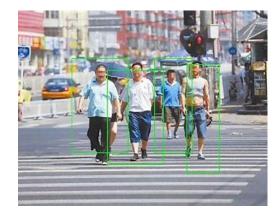
佩戴安全帽检测

工位睡觉检测

智慧农业 "tomato that are not ripe"

"person fallen"

"person that are not wearing helmet" "person that is sleeping"







行人安全检测

抽烟检测

交通管理

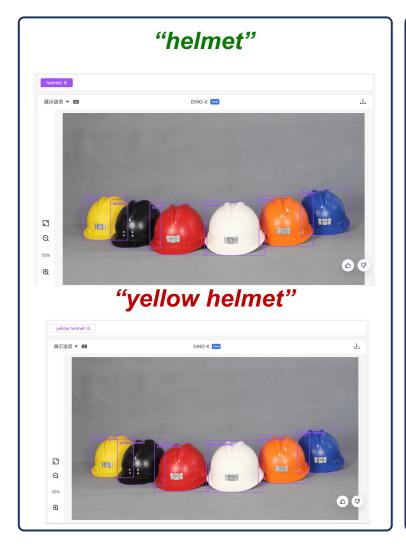
"person on the crossroad"

"person that are smoking"

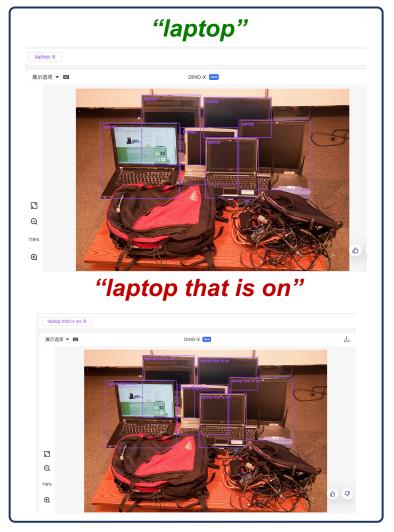
"cars that are crushed"



Finding 1: State-of-the-art Open-set detection model (DINO-X) lack language comprehension capabilities







In Open-Vocabulary object detection, the term "open" only refers to openness in categories,



Finding 2: State-of-the-art Multimodal LLMs lack fine-grained perception capabilities



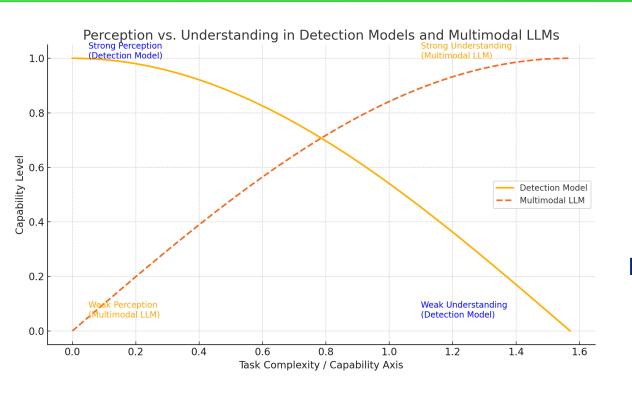
User: Please help me detect person in this image

MLLMs:

"Sure, here is person [[90, 70, 120, 340], [110, 70, 125, 400]]"

- coordinate shift
- tiny object detection
- dense object detection





Detection Model: Strong Perception, Weak Understanding

Multimodal LLMs: Weak Perception, Strong Understanding

Next Step: A mode with both strong perception and understanding capabilities

ChatRex: Taming Multimodal LLM for Joint Perception and Understanding



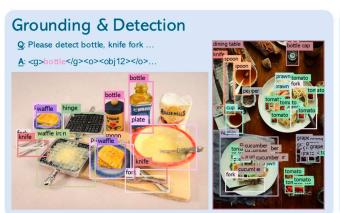
ChatRex: Taming Multimodal LLM for Joint Perception and Understanding

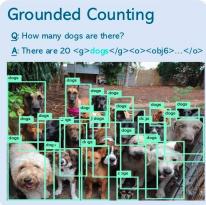
Qing Jiang^{1,2}, Gen Luo¹, Yuqin Yang^{1,2}, Yuda Xiong¹, Zhaoyang Zeng¹
Yihao Chen¹, Tianhe Ren¹, Lei Zhang^{1,2†}

¹International Digital Economy Academy (IDEA)

²South China University of Technology

mountchicken@outlook.com , leizhang@idea.edu.cn







Grounded Conversation

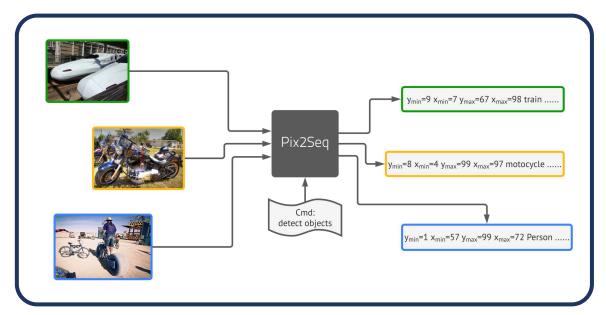
- **Q**: Answer the question in grounded format. What is unusual about this image?
- A: The unusual aspect of this image is that it features <g>dog</g><o><obj3>...<o>s and a <g>cupcake</g><o><obj0>...</o>s, which is not a common sight. Typically, images of dogs are not associated with food items like cupcakes. The combination of a dog and a cupcake in the same image is unexpected and creates a unique and intriguing visual experience.



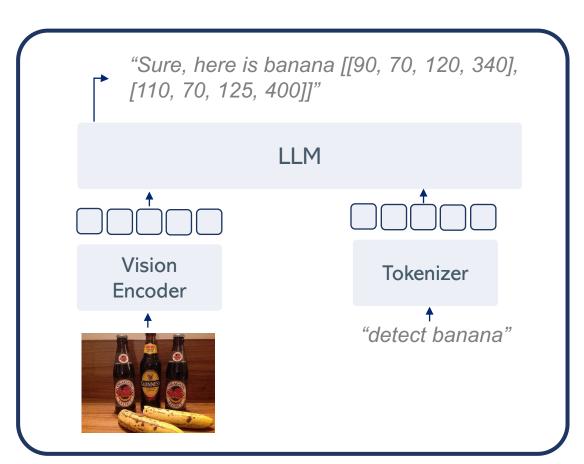
Motivation: How MLLMs do Detection?



How MLLMs do object detection: directly predict quantized coordinates as text tokens^[1].



Pix2Seq^[1]



Modern MLLMs

Motivation: How MLLMs do Detection?



But the performance is not good enough



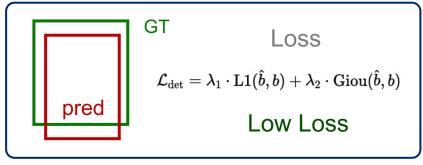
Method	Туре	C	COCO-Val	0	LVIS-Mini Val									
		P@0.5	R@0.5	mAP	P@0.5	R@0.5	mAP	AP-R	AP-C	AP-F				
Faster-RCNN [70]		-	-	42.0	-	-	-	1-1	-	-				
DETR [8]	Closed-set	-	=	43.3	=1	-	-	-	=	-				
Pix2Seq [12]	Detection Model	-	-	43.2	-	-		-	-	-				
DINO [102]		-	-	49.4	-	-		-	-	-				
Florence2 [88]		-	-	43.4	-	-	-	-	-	-				
GLIP [39]	Open-set	=	=	49.8	-	Ξ.	37.3	28.2	34.3	41.5				
T-Rex2 [29]	Detection Model	-	-	46.5	-	-	47.6	45.4	46.0	49.5				
Grounding DINO [52]		-		48.4	-		33.0	22.2	30.7	38.8				
Shikra-7B [10]		40.3	21.5	-	52.8	14.5	-	-	-	-				
Ferret-7B [94]		66.3	33.5	-	72.9	25.2	-	=	=	=				
Groma-7B [61]	MLLM	69.9	28.9	1-1	76.3	10.9		-	-	-				
InternVL2-7B [14]	WILLIVI	45.3	24.5	1-1	51.6	13.1	=	-	-	-				
Qwen2-VL-7B [85]		59.3	43.9	2 - 8	77.0	34.7	-1	-	-	-				
ChatRex-7B		73.5	72.8	48.2	80.3	58.9	42.6	44.6	48.4	37.2				

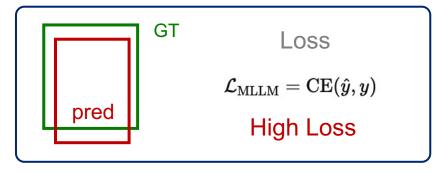
Low Recall Rate

Motivation: What's the Challenge?



Directly predict the coordinates is a hard task: Regression V.S. Classification





detection model training

MLLM training

- 2. Error Propagation: Each box requires at least 9 tokens and can cause cascading errors.
- 3. Ambiguity in Prediction Order: Auto-regressive prediction needs a predefined sequence order.



"bottle1, bottle2, bottle3"

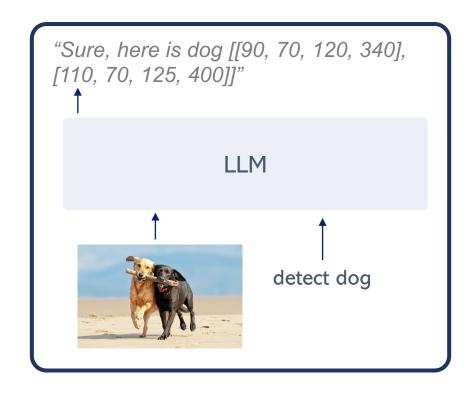
"bottle3, bottle2, bottle1" "bottle2, bottle1, bottle3"

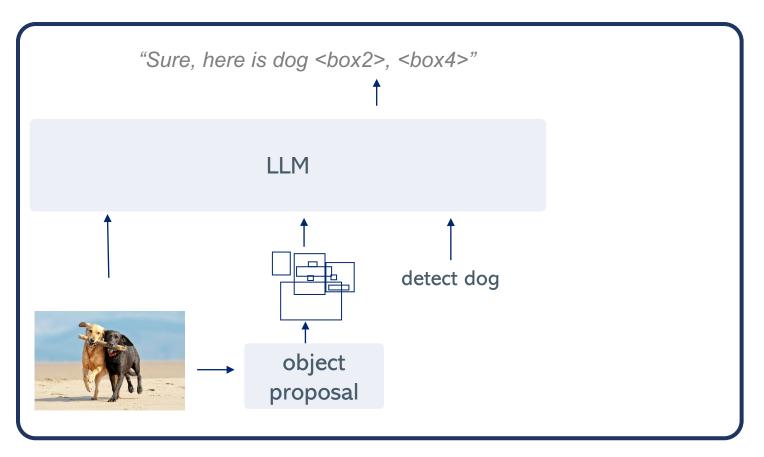
4. Quantization Range Limitation: Large image (>1000 px) input can lead to quantization error.

Solution: Retrieval-based Perception MLLM



Core idea: LLM has strong understanding capability, while detection model has strong perception capability

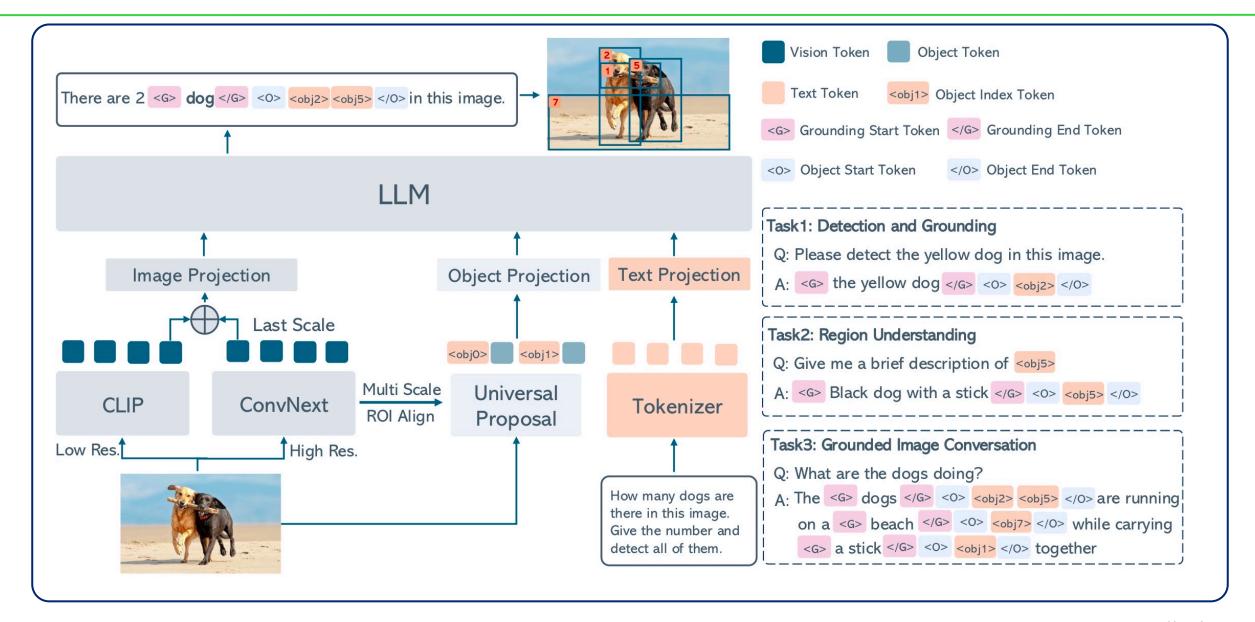




Retrieval based method

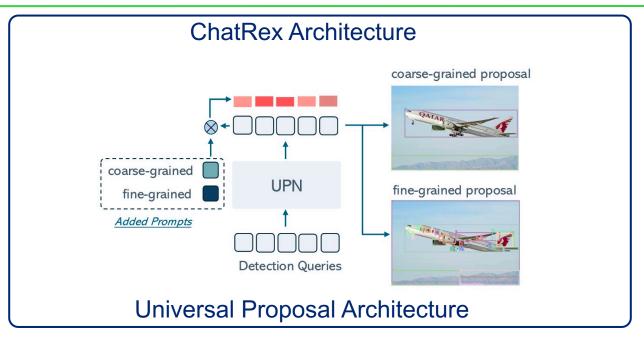
ChatRex: Detection-Oriented MLLM





ChatRex: Detection-Oriented MLLM







Experiments



Method	Туре	C	COCO-Val	8	LVIS-Mini Val									
		P@0.5	R@0.5	mAP	P@0.5	R@0.5	mAP	AP-R	AP-C	AP-F				
Faster-RCNN [70]		-	-	42.0	-	r=	-	89 — 131	-	-				
DETR [8]	Closed-set	_	-	43.3	-	-	-	=	_	-				
Pix2Seq [12]	Detection Model	-	-	43.2	-	-	_	-	-	-				
DINO [102]		-	-	49.4	-	-	-11	_	_	-				
Florence2 [88]		_	-	43.4	_	-	-	-	_	- 1				
GLIP [39]	Open-set	=	-	49.8	=	=	37.3	28.2	34.3	41.5				
T-Rex2 [29]	Detection Model	-	-	46.5	-	-	47.6	45.4	46.0	49.5				
Grounding DINO [52]		-	-	48.4	-	=	33.0	22.2	30.7	38.8				
Shikra-7B [10]		40.3	21.5	0-8	52.8	14.5	-	-	-	-				
Ferret-7B [94]		66.3	33.5	-	72.9	25.2	=	===	=	-				
Groma-7B [61]	MLLM	69.9	28.9	77 — 31	76.3	10.9	-	-		-				
InternVL2-7B [14]	WILLIVI	45.3	24.5	-	51.6	13.1	-	-	-	-				
Qwen2-VL-7B [85]		59.3	43.9	77 — 38	77.0	34.7		-		- 1				
ChatRex-7B		73.5	72.8	48.2	80.3	58.9	42.6	44.6	48.4	37.2				

Common/Long-tailed Object Detection



QA Example:

User: Please detect person; cup in this image. Answer the question with object indexes.

ChatRex: <g>person</g><o><obj1><obj5><obj16><obj21></o>\n <g>cup</g><o><obj12><obj14><obj33></o>



Figure 7: Visualization on Common Object Detection Task.

QA Example:

User: Please detect saluki; folding chair in this image. Answer the question with object indexes.

ChatRex: <g>saluki</g><o><obj12></o>\n

<g>folding chair</g><o><obj19><obj23></o>\n

Visualization:

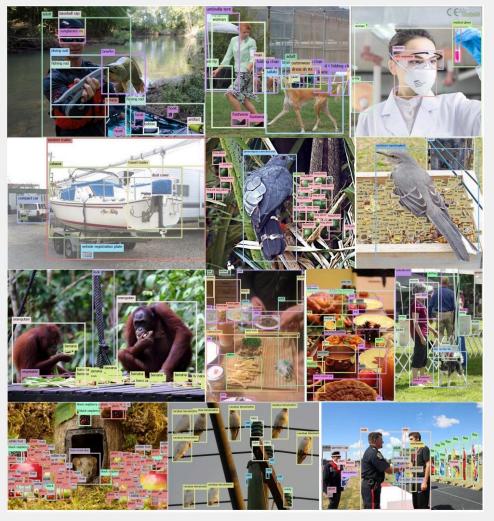


Figure 8: Visualization on Long-tailed Object Detection Task.

Region Understanding

idea

QA Example:

User: What is <obj0> wearing and what is he doing now?

ChatRex: <0bj0> is wearing a black jacket and a black beanie. He is standing in a store, looking at something on the shelf.

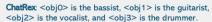
Visualization:



Question: What is <objO> wearing and what is he doing now?

ChatRex: <obj0> is wearing a black jacket and a black beanie. He is standing in a store, looking at something on the shelf.

Question: What are the positions of <obj0>, <obj1>, <obj2>, <obj3> in the band, respectively?







Question: What is <objO>? And how might <objO> feels now?

ChatRex: <objO> is a cow, and it might feel scared or startled by the action of the man on the horse.



Question: What is the price for <objO> according to this image.

ChatRex: The price for <objO> is 2.49, as indicated by the sign in front of the cat.



Question: Describe the wearing of <obj0> in detail

ChatRex: The dress of <objO> is a bright yellow, which stands out against the snowy backdrop. It appears to be a winter coat, designed to provide warmth in cold weather. The coat is long-sleeved and reaches down to the person's knees, suggesting it is made for cold temperatures. The coat's color is vibrant and contrasts with the muted colors of the snowy environment.

Figure 12: Visualization on Region QA Task.

QA Example:

User: Give a short phrase that describes <0bj0>;<0bj1>.

ChatRex: $\langle g \rangle$ Starbucks $\log o \langle /g \rangle \langle o \rangle \langle o \rangle \backslash n$

<g>a coffee cup</g><o><obj1></o>\n

Visualization:

Region Caption in Category Name









Region Caption in Short Phrase



Region Caption in Referring Style









Region Caption in Detail









Figure 11: Visualization on Detailed Grounded Image Caption Task.

Grounded Conversation



QA Example:

User: Please briefly describe this image and detect all the mentioned objects. Answer with grounded object indexes.

ChatRex: A <g>man</g><o><obj13></o> in a <g>white tuxedo</g><o><obj17></o> with a <g>red bow tie</g><o><obj16></o> is holding an <g>Oscar statuette</g><o><obj14></o> and standing on a stage with a microphone, while a large, ornate Oscar statue is visible in the background.

Visualization:



A man in a white tuxedo with a red bow tie is holding an Oscar statuette and standing on a stage with a microphone, while a large, ornate Oscar statue is visible in the background



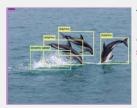
A person is slicing fresh cilantro a wooden cutting board, surrounded by various ingredients including a bowl of chopped tomatoes, a bowl of grated cheese a bowl of chopped herbs, a whole avocado and a bowl of grated cheese.



A woman with glasses is smiling and holding a white envelope with a handwritten address, standing in a cozy room with a window and a bookshelf filled with books in the background.



The image showcases a variety of ingredients laid out on a white surface, including a plate of raw beef, a bowl of sliced onions, a plate of fresh green lettuce, a plate of sliced cheese, a bowl of nuts, and a plate of bread.



Three dolphins are leaping out of the water, creating a dynamic splash as they dive and play in the ocean.



The image is a vintage advertisement for Fitzgerald & Co., featuring a hand pouring a bottle of whiskey into a crescent moon with the text "Fitzgerald & Co. Dublin" prominently displayed.



A woman with long brown hair is enjoying a cone of ice cream on a beach, with a backdrop of waves and seagulls in flight.



The image features a charming black dog-shaped keychain with a textured surface, adorned with a pink tongue and a small brown tag attached to a brown leather strap, all set against a white background.

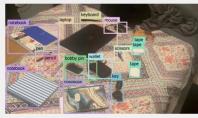
Figure 13: Visualization on Brief Grounded Image Caption Task.

QA Example:

User: Answer the question in grounded format. How would you organize this space?

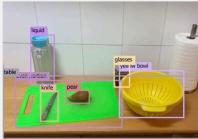
ChatRex: To organize the space, I would sort the items into categories such as stationery, electronics, personal items, and accessories. For stationery, I would group the <g>pen</g><o><obj23></o>, <g>pencil</g><o><obj12></o>>...

Visualization:



Question: Answer the question in grounded format. How would you organize this space? Sort and arrange items on the surface based on their function or category, ensuring the area is tidy and visually pleasing.

ChatRex: To organize the space, I would sort the items into categories such as stationery, electronics, personal items, and accessories. For stationery, I would group the pen, pencil, and ruler together. Electronics would include the laptop, mouse, and keyboard. Personal items would be grouped with the notebooks, wallet, and keys. Accessories would be separated from the other categories and include the scissors, bobby pins, and tape. After sorting, the items would be arranged neatly on the surface, with similar items grouped together and the larger items placed at the edges to create a visually pleasing and organized space.



Question: Answer the question in grounded format. Use the tools available on the table to prepare the pear for consumption. Ensure any necessary cleaning, peeling, or cutting is performed appropriately, and leave the workspace tidy afterward.

ChatRex: To prepare the pear for consumption, I would first ensure that the pear is clean and free of any dirt or blemishes. If necessary, I would use the knife to peel the pear, removing the skin in a smooth and even manner. After peeling, the pear would be placed in the yellow bowl for further preparation. If the pear is to be cut, the knife would be used to make the necessary cuts, ensuring that the pear is sliced into even portions. The glasses on the table could be used to rinse the pear under running water to remove any residual water or liquid. After the pear is prepared, the knife would be cleaned and stored properly, and the table would be tidied up to maintain a clean and organized workspace.

Figure 16: Visualization on Detailed Grounded Conversation Task.

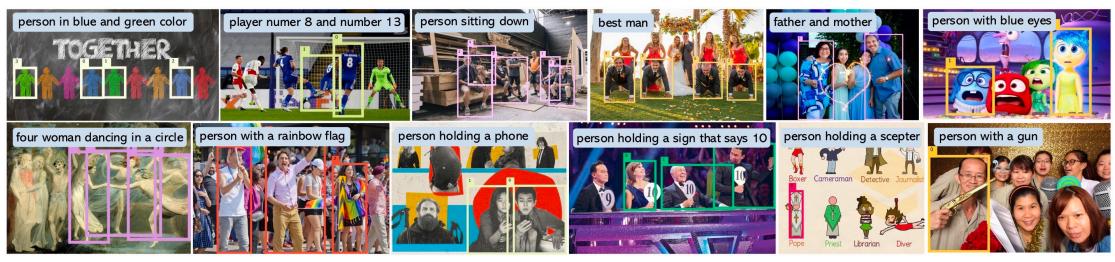
Referring to Any Person



Referring to Any Person

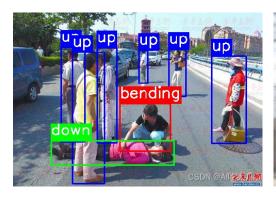
Qing Jiang 1,2 , Lin Wu 1,2 , Zhaoyang Zeng 1 , Tianhe Ren 1 , Yuda Xiong 1 Yihao Chen 1 , Liu Qin 1 , Lei Zhang 1,2† 1 International Digital Economy Academy (IDEA) 2 South China University of Technology

mountchicken@outlook.com , leizhang@idea.edu.cn



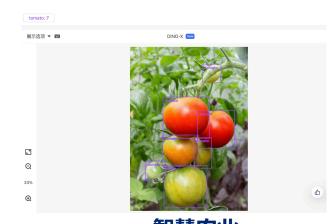
Most Detection Tasks Can be formulated as Referring











摔倒检测

佩戴安全帽检测

工位睡觉检测

智慧农业

"person fallen"

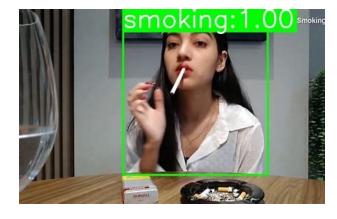
"person that are not wearing helmet" "person that is sleeping"

"tomato that are not ripe"





"person on the crossroad"



抽烟检测

"person that are smoking"

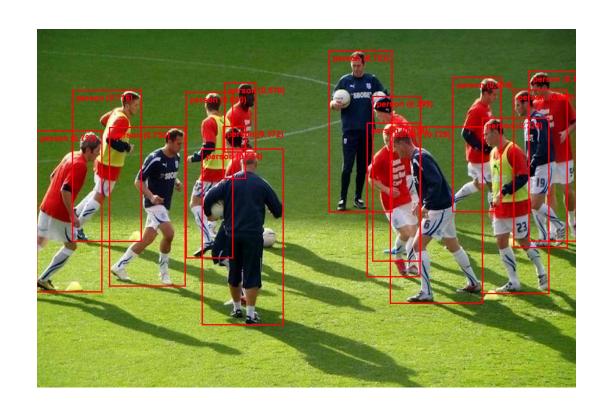


交通管理

"cars that are crushed"

Referring V.S. Detection





Detection: "person"



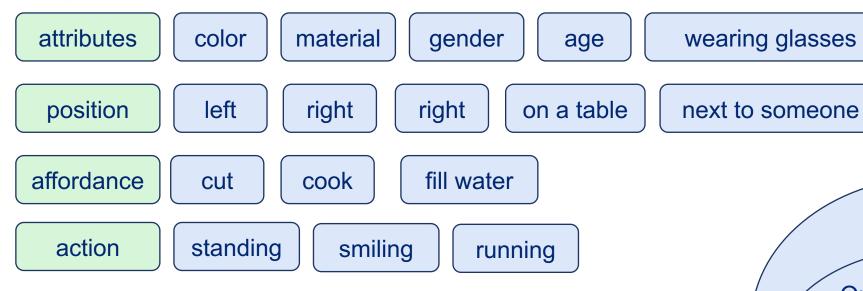
Referring: "person who is holding two footballs"

Referring V.S. Detection



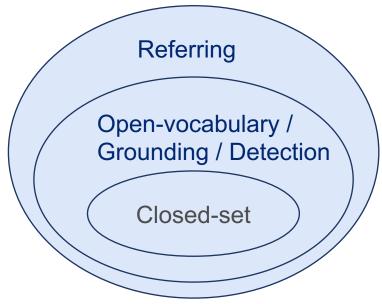
Detection: Category name e.g. man

Referring: Category name +



E.g.

- a white man
- the second white man from the left
- The second white man from the left that is wearing a blue hat
- The second white man from the left that is wearing a blue hat and is smiling



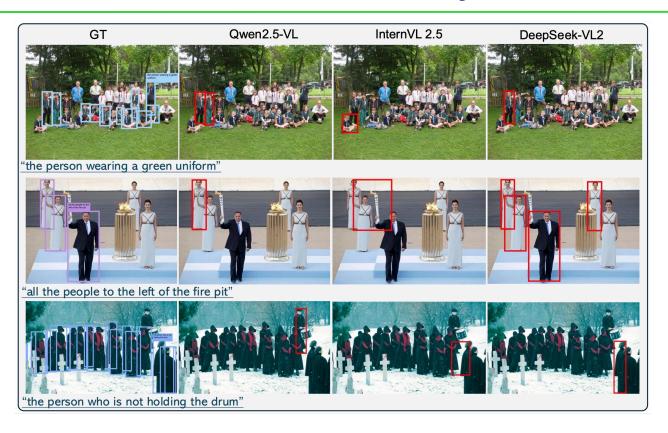
Motivation: Current SOTA models lack usability



Datasets	InternVL2.5	Qwen2.5-VL	Qwen2.5-VL
	78B	72B	7B
Refcoco _{val}	93.7	92.7	90.0
Refcoco _{testA}	95.6	94.6	92.5
Refcoco _{testB}	92.5	89.7	85.4
Refcoco+ $_{val}$	90.4	88.9	84.2
Refcoco+ $_{testA}$	94.7	92.2	89.1
Refcoco+ $_{testB}$	86.9	83.7	76.9
Refcocog _{val}	92.7	89.9	87.2
Refcocog _{test}	92.2	90.3	87.2

High Performance in existing benchmarks

- 1. Designing flaws in existing benchmarks
- 2. Current MLLMs are still less capable



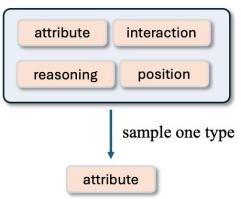
Low Performance in real-world scenarios

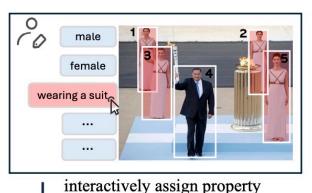
Solutions: Data + Model



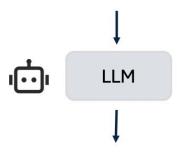
HumanRef Dataset





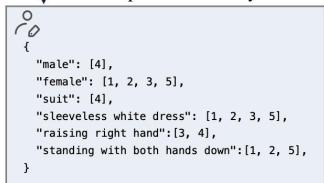


Please compose and transfer the following description to referring style

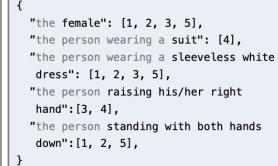


```
box5
{
    "gender":"female",
    "age":"adult",
    "top":"sleeveless white dress",
    "pose":"standing",
    "expression":"smiling",
    "shoes":"sandals",
    "accessories":"none",
}
```

```
"male",
   "female",
   "suit",
   "sleeveless white dress",
   "raising right hand",
   "standing with both hands down",
]
```



and then parse to dictionary



a) pseudo labeling

b) write property list

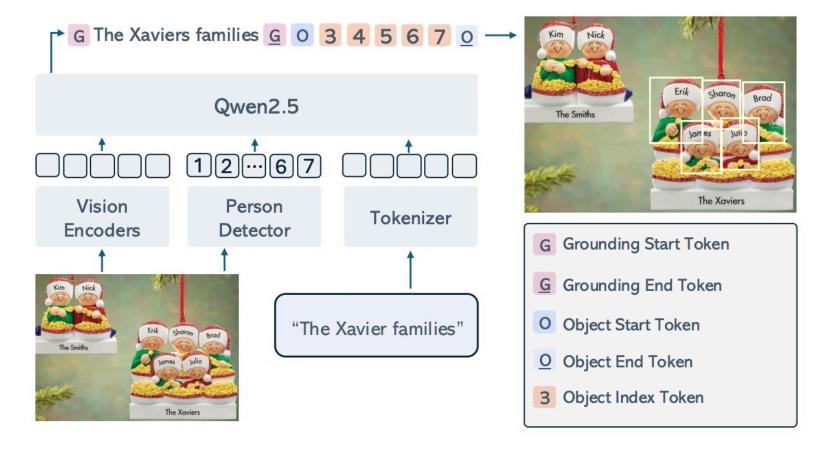
- c) assign property to each person
- d) transfer to referring style with LLM

Figure 3. Overview of the mannual annotation pipeline of the HumanRef dataset.

Solutions: Model



RexSeek



- 1. Strong perception capability
- 2. Strong language comprehension

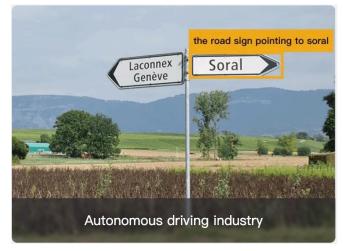
Applications: Referring Any Person





Applications: Any Object Referring



















Rex-Thinker: Grounded Object Referring via Chain-of-Thought Reasoning

Qing Jiang^{1,2*}, Xingyu Chen^{3*}, Zhaoyang Zeng¹, Junzhi Yu³, Lei Zhang^{1,2†}

¹International Digital Economy Academy (IDEA)

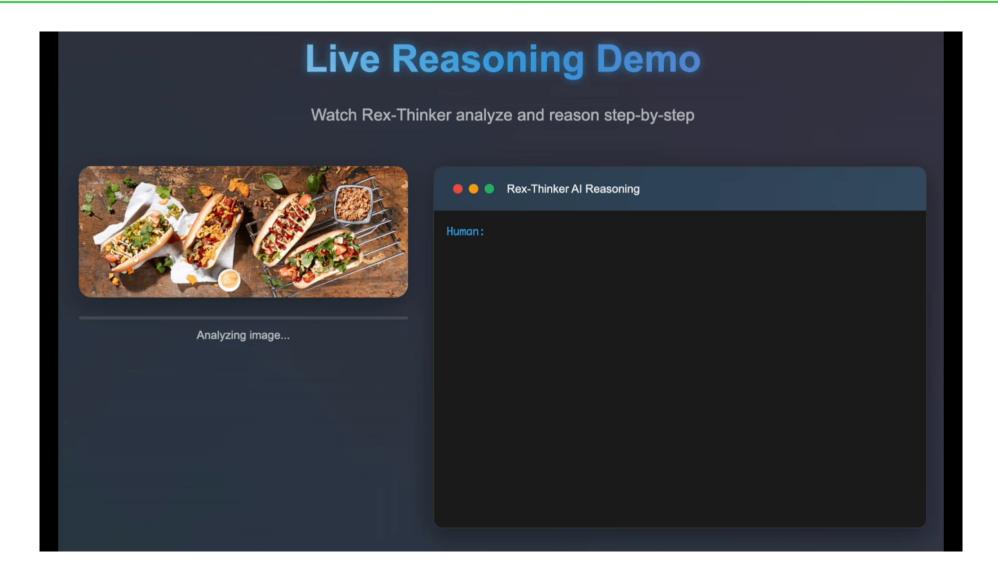
²South China University of Technology

³Peking University

{jiangqing, chenxingyu, leizhang}@idea.edu.cn

Rex-Thinker: Referring with CoT Reasoning





Planning-Action-Summarization Formulation



Box hint



Answer



Question: Please detect ripe tomato..

Answer:

<think>OK, the user needs us to detect the ripe tomato. In order to accomplish this task, I need to break it down into the following steps:

- Step 1: Determine the ripeness of each tomato based on its color. Let's start our analysis:

Step 1: Determine the ripeness of each tomato based on its color

I see seven tomatoes in this image, let's look at their colors and put a 4 if it is a ripe tomato and \times if it is not.

- Tomato 1: It has a deep red color, indicating it is a ripe tomato.
- Tomato 2: It also has a deep red color, indicating it is a ripe tomato.
- Tomato 3: It has an orange-brown color, which is not the typical red color of a fully-ripe tomato.
- Tomato 4: It is green in color, so it is not a ripe tomato.
- Tomato 5: It is green in color, so it is not a ripe tomato.
- Tomato 6: It is green in color, so it is not a ripe tomato.
- Tomato 7: It is green in color, so it is not a ripe tomato.

Summarize and Re-Check answer

Let's now recheck our answer and put $\sqrt{\ }$ for the target tomato and \times for others

- Tomato 1: It is a ripe tomato.
- Tomato 2: It is a ripe tomato.
- Tomato 3: It is not a ripe tomato. X
- Tomato 4: It is not a ripe tomato.
- Tomato 5: It is not a ripe tomato.
- Tomato 6: It is not a ripe tomato. X
- Tomato 7: It is not a ripe tomato. **X**</think>

<answer[{"bbox_2d": [235, 358, 565, 657], "label": "ripe tomato"}, {"bbox_2d": [481, 455, 744, 726], "label": "ripe tomato"}]```</answer>

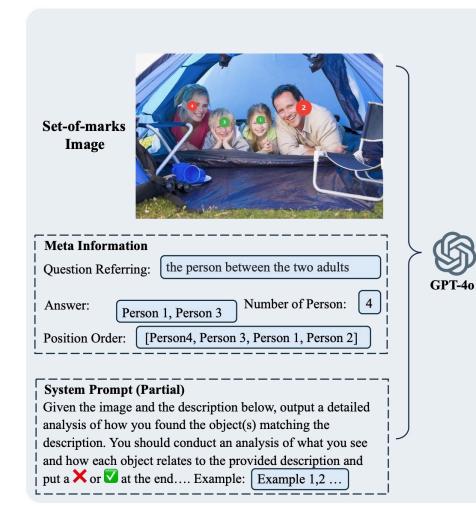
Planning

Action

Summarization

Data: HumanRef-CoT, CoT Referring Dataset





1. Planning:

OK, the user needs us to detect the person between the two adults. To accomplish this task, I need to break it down into the following steps:

- Step 1: Identify the two adults in the image.
- Step 2: Determine the person located between the two adults.

2. Action:

Step 1: Identify the two adults in the image.

I see four people in this image. Let's look at their age - related attributes and put a . if they are adults and X if they are not.

- Person 1: A young girl with long hair, not an adult. \rightarrow
- Person 2: A man with sunglasses on his head, appears to be an adult. \rightarrow \triangle
- Person 3: A young boy with short hair, not an adult. \rightarrow X
- Person 4: A woman with her hair tied back, appears to be an adult. \rightarrow .

So Person 2 and Person 4 are the two adults.

Step 2: Determine the person located between the two adults.

The order from left to right is [person 4, person 3, person 1, person 2].

Between Person 4 (the woman on the left - hand side) and Person 2 (the man on the right - hand side), there are two people: Person 3 and Person 1.

3. Summarization:

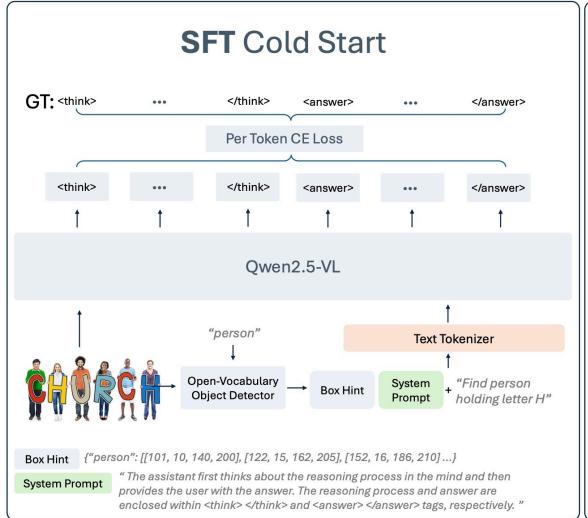
Summarize and Re-Check answer

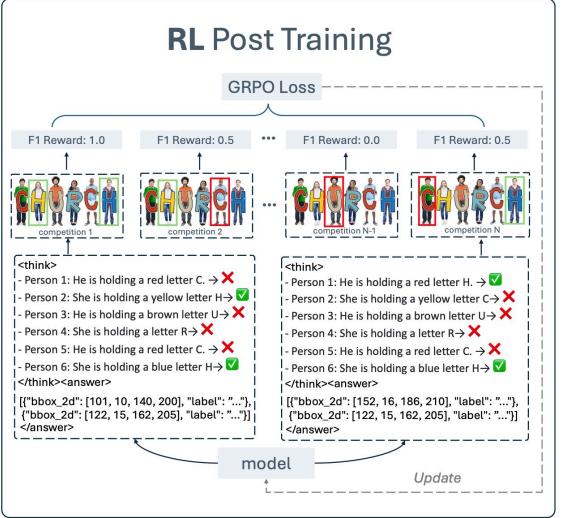
Let's now recheck our answer and put V for the target people and X for others:

- Person 1: Located between the two adults. \rightarrow
- Person 2: One of the adults, not the person between the adults. \rightarrow \times
- Person 3: Located between the two adults. \rightarrow
- Person 4: One of the adults, not the person between the adults. \rightarrow \times

Model: Qwen2.5-VL with SFT + GRPO Training







Experiments: Main Results



Method	1	Attribut	e		Positior	ı	Ir	nteractio	on	R	easonir	ıg	(Celebrit	y		Average	e	Rejection
Wiethod	R	P	DF1	R	P	DF1	R	P	DF1	R	P	DF1	R	P	DF1	R	P	DF1	Score
DINOX [46]	59.5	28.8	20.9	78.8	28.1	17.6	67.3	28.5	18.9	76.2	32.1	22.2	94.1	48.0	37.0	75.2	33.1	23.3	36.0
InternVL-2.5-8B [8]	23.5	39.0	27.1	23.0	28.0	24.3	27.8	40.1	31.3	17.5	22.8	18.9	57.4	59.3	58.0	29.8	37.8	31.9	54.9
Ferret-7B [66]	27.9	44.4	30.4	30.2	36.2	29.8	30.8	41.8	31.2	19.7	33.7	22.8	63.2	60.0	57.5	34.4	43.2	34.3	2.0
Groma-7B [36]	67.5	47.8	38.6	63.2	43.1	37.2	66.6	48.1	40.6	59.1	41.4	34.8	73.2	63.3	59.1	65.9	48.7	42.1	0.0
ChatRex-7B [22]	44.3	78.0	51.8	48.0	66.7	52.5	49.6	74.8	56.5	36.6	65.1	42.8	73.7	76.5	74.2	50.4	72.2	55.6	0.0
Qwen2.5-VL-7B [3]	49.1	71.3	54.4	50.2	61.7	52.8	48.2	66.3	53.2	34.6	61.2	40.3	80.3	81.9	80.1	52.5	68.5	56.2	7.1
DeepSeek-VL2-small [57]	52.3	78.0	57.7	56.4	66.1	58.1	55.4	75.7	60.7	46.6	61.7	50.1	85.9	74.3	70.7	59.3	71.2	59.5	3.1
Molmo-7B-D [11]	82.7	86.4	76.3	78.0	80.6	72.4	69.9	77.7	66.1	72.1	80.4	65.5	85.9	87.5	82.9	77.7	82.5	72.6	68.6
RexSeek-7B [21]	<u>87.2</u>	86.8	81.5	86.1	86.3	83.8	84.8	<u>84.6</u>	80.7	87.8	<u>84.7</u>	81.5	83.4	83.4 86.5 84.2 <u>85.9</u> 85.8 <u>82.3</u>		54.1			
Rex-Thinker-Plain	83.0	88.7	81.4	82.5	83.9	81.3	80.1	85.6	80.2	80.5	82.2	77.3	86.7	88.7	86.8	82.6	85.8	81.4	53.5
Rex-Thinker-CoT	86.6	87.7	82.7	<u>86.5</u>	<u>87.0</u>	<u>84.3</u>	79.6	81.7	77.2	85.7	83.8	80.3	<u>87.6</u>	89.5	87.2	85.2	<u>85.9</u>	<u>82.3</u>	67.3
Rex-Thinker-GRPO	88.5	88.7	84.1	87.2	87.1	84.6	<u>81.5</u>	83.5	79.1	<u>87.7</u>	85.4	82.3	88.0	89.3	87.2	86.6	86.8	83.5	68.2

Rex-Thinker-Plain: Training without CoT

Rex-Thinker-CoT: Training with CoT

Rex-Thinker-GRPO: Training with CoT and GRPC

- Conclusion 1: CoT can effectively improve mode answer questions about objects that are not present
- Conclusion 2: Reinforcement learning fine-tuning







Question: Please detect Kobe Bryant.

Answer

<think>OK, the user needs us to detect Kobe Bryant. To accomplish this task, I need to break it down into the following steps:

- Step 1: Identify each person based on their appearance and basketball-related attributes.
- Step 2: Check if any of the identified people match the identity of Kobe Bryant.

Step 1: Identify each person based on their appearance and basketball-related attributes.

- Person 1: A man wearing a white Dallas-branded basketball jersey, holding a basketball. He has short hair and a beard.
- Person 2: A man wearing a green Boston-branded basketball jersey with the number 0. He is in a shooting pose.
- Person 3: A man wearing a white Denver Nuggets-branded basketball jersey with the number 15. He is in a shooting pose.
- Person 4: A man wearing a black San Antonio Spurs-branded basketball jersey. He has a muscular build.
- Person 5: A man wearing a blue and yellow basketball jersey. He has a beard and is smiling.
- Person 6: Not clearly visible in terms of distinct features related to identification.
- Person 7: Partially visible, no clear distinct features for identification.
- Person 8: A man wearing a purple and yellow hoodie, not in a basketball-playing pose.
- Person 9: A man wearing a blue Thunder-branded basketball jersey with the number 2. He has a headband.

Step 2: Check if any of the identified people match the identity of Kobe Bryant.

Kobe Bryant is a well-known basketball player, but none of the people in the image have the exact appearance and context that can confirm them as Kobe Bryant...

Summarize and Re-Check answer

Let's now recheck our answer and put of for the target person and if for others

- Person 1-9: Not Kobe Bryant → X</think>

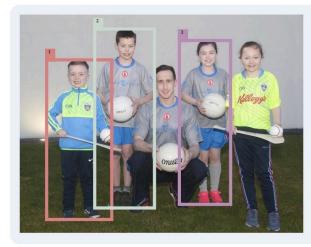
<answer>There is no such thing in the image.</answer>

Ablation: Cold start training is important



With	1	Attribut	e	Position			Interaction			R	easonir	ng	(Celebrit	y		Average	Rejection	
Cold Start	R	P	DF1	R	P	DF1	R	P	DF1	R	P	DF1	R	P	DF1	R	P	DF1	Score
No	81.4	85.8	78.1	80.2	80.2	77.5	79.6	82.6	78.0	77.6	75.0	70.6	87.3	86.5	84.8	81.2	82.0	77.8	66.4
Yes	88.5	88.7	84.1	87.2	87.1	84.6	81.5	83.5	79.1	87.7	85.4	82.3	88.0	89.3	87.2	86.6	86.8	83.5	68.2

Table 5: Ablation on the impact of CoT-based cold start on final performance after GRPO training.



Question: Please detect the person holding a small ball

Answer:

<think>

There is one person holding a small ball in the image. This person is standing on the left side of the group.

</think>

<answer>

[100, 174, 369, 821],[290, 51, 538, 775],[625, 97, 839, 761]

</answer>

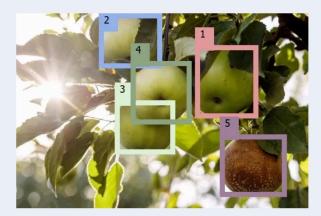
• Conclusion 3: CoT-based SFT cold start helps the model learn the reasoning format, enabling it to perform reasoning in the desired format during GRPO.

GRPO without CoT SFT as cold start

Visualization Results: Generalized to Any Object



Box hint



Answer



Question: Please detect apple with disease, and tell me what kind of disease it may be..

Answer:

<think>OK, the user needs us to detect the apple with disease and tell what kind of disease it may be. In order to accomplish this task, I need to break it down into the following steps:

- Step 1: Examine each apple to determine if it has any visible signs of disease.
- Step 2: If an apple is found to have a disease, identify the type of the disease. Let's start our analysis:

Step 1: Examine each apple to determine if it has any visible signs of disease

I see five apples in this image. Let's look at them one-by-one and put a if there are signs of disease and if there are none.

- Apple 1: It appears to be a healthy green apple with no visible signs of disease. X
- Apple 2: It is also a healthy-looking green apple with no obvious signs of disease. X
- Apple 3: This apple looks like a healthy green apple with no clear signs of disease. X
- Apple 4: Another healthy green apple with no visible signs of disease. X
- Apple 5: This apple has a brownish spot on its surface, which indicates that it has a disease. 🚣

Step 2: Identify the type of the disease

The brownish spot on Apple 5 suggests that it may have a fungal or bacterial infection, which are common types of diseases for apples.

Summarize and Re-Check answer

Let's now recheck our answer and put \mathbf{V} for the target apple and \mathbf{X} for others:

- Apple 1: It is a healthy apple with no signs of disease.
- Apple 2: It is a healthy apple with no signs of disease.
- Apple 3: It is a healthy apple with no signs of disease.
- Apple 4: It is a healthy apple with no signs of disease. X
- Apple 5: It has a brownish spot, indicating a possible fungal or bacterial infection. ✓</think><answer[{"bbox 2d": [429, 261, 570, 395], "label": "apple with disease"}]```</answer>

Visualization Results



Box hint



Answer



Question: Please detect ripe tomato..

Answer:

<think>OK, the user needs us to detect the ripe tomato. In order to accomplish this task, I need to break it down into the following steps:

- Step 1: Determine the ripeness of each tomato based on its color. Let's start our analysis:

Step 1: Determine the ripeness of each tomato based on its color

I see seven tomatoes in this image, let's look at their colors and put a . if it is a ripe tomato and X if it is not.

- Tomato 1: It has a deep red color, indicating it is a ripe tomato. \triangle
- Tomato 2: It also has a deep red color, indicating it is a ripe tomato.
- Tomato 3: It has an orange-brown color, which is not the typical red color of a fully-ripe tomato. X
- Tomato 4: It is green in color, so it is not a ripe tomato.
- Tomato 5: It is green in color, so it is not a ripe tomato.
- Tomato 6: It is green in color, so it is not a ripe tomato.
- Tomato 7: It is green in color, so it is not a ripe tomato.

Summarize and Re-Check answer

Let's now recheck our answer and put **v** for the target tomato and **x** for others

- Tomato 1: It is a ripe tomato.
- Tomato 2: It is a ripe tomato.
- Tomato 3: It is not a ripe tomato.
- Tomato 4: It is not a ripe tomato.
- Tomato 5: It is not a ripe tomato.
- Tomato 6: It is not a ripe tomato.
- Tomato 7: It is not a ripe tomato. **X** </think>

<answer[{"bbox_2d": [235, 358, 565, 657], "label": "ripe tomato"}, {"bbox_2d": [481, 455, 744, 726],
"label": "ripe tomato"}]```</answer>

Visualization Results



Box hint



Answer



Question: Please detect damaged container.

Answer:

<think>OK, the user needs us to detect the damaged container. In order to accomplish this task, I need to break it down into the following steps:

- Step 1: Examine each of the 31 containers to determine if they have any visible signs of damage.

Step 1: Step 1: Examine each of the 31 containers to determine if they have any visible signs of damage.

- Container 1: It appears to be in good condition with no visible signs of damage. X
- Container 2: There is no visible sign of damage on this container.
- Container 3: This container has a large hole in its side, indicating significant damage. 4
- Container 4: No visible signs of damage. X

• • •

- Container 31: No visible signs of damage. X

Summarize and Re-Check answer

Let's now recheck our answer and put V for the target container and X for others

- Container 1: No visible signs of damage → X
- Container 2: No visible signs of damage \rightarrow \times
- Container 3: Has a large hole, is damaged → ✓
- Container 4: No visible signs of damage → X

• • •

- Container 30: No visible signs of damage → X
- Container 31: No visible signs of damage → X

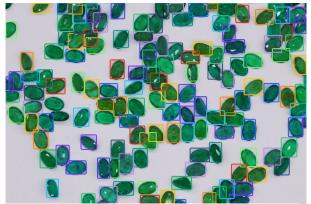
</think>

<answer[{"bbox_2d": [288, 183, 699, 459], "label": "damaged container"}]```</answer>

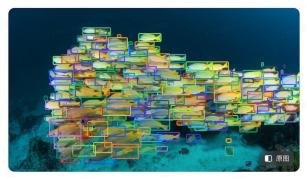
SEED-VL-1.5

idea

SEED-VL-1.5 (Dense Object Detection)













Question:帮我框选图中所有的鸽子,以 <bbox>x1 y1 x2 y2</bbox>的形式表示

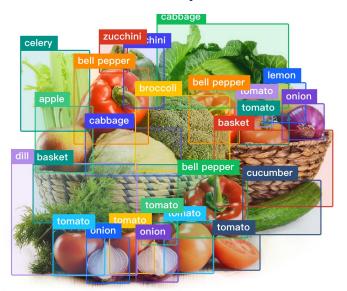
Abswer:

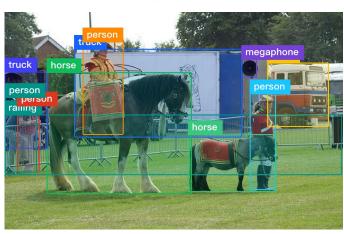
<bbox>3 10 22 987</bbox> <bbox>21 10 40 987</bbox> <bbox>39 10 58 987</bbox> <bbox>57 10 76 987</bbox> <bbox>
57 10 94 987</bbox> <bbox>93 10 112 987</bbox> <bbox> 111 10 130 987</bbox> <bbox> 129 10 148 987</bbox> <bbox> <bbox>
51 0 94 987</bbox> <bbox> 201 10 220 987</bbox> <bbox>
51 0 184 987</bbox> <bbox> 183 10 202 987</bbox> <bbox> 201 10 220 987</bbox>
51 0 220 987</bbox>
51 0 274 987</br/>51 0 274 987</br>

SEED-VL-1.5



SEED-VL-1.5 (Universal Proposal)





Question:请帮我识别一下图片,然后告诉我图中有哪些物体,并框选这些物体,物体类别名请用英文表示,对于每个物体,请提供其类别和边界框,格式为: [{"category": category, "bbox": "<bbox>x1 y1 x2 y2</bbox>"}, {"category": category, "bbox": "<bbox>x1 y1 x2 y2</bbox> y1 x2 y2</bbox> 进行包裹

Answer:

[{"category": "person", "bbox": "<bbox>3 392 95 736</bbox>"}, {"category": "person", "bbox": "<bbox>39 433 103 736</bbox>"}, {"category": "truck", "bbox": "<bbox>0 276 128 627</bbox>"}, {"category": "horse", "bbox": "<bbox>125 278 553 835</bbox>"}, {"category": "person", "bbox": "<bbox>232 138 350 564</bbox>"}, {"category": "truck", "bbox": "<bbox>207 171 700 576</bbox>"}, {"category": "truck", "bbox": "<bbox>701 217 745 297</bbox>"}, {"category": "truck", "bbox": "<bbox>778 234 957 527</bbox>"}, {"category": "horse", "bbox": "<bbox>543 564 803 825</bbox>"}, {"category": "person", "bbox": "<bbox>725 377 800 817</bbox>"}, {"category": "hat", "bbox": "<bbox>262 138 326 184</bbox>"}, {"category": "railing", "bbox": "<bbox>0 474 996 748</bbox>"}]

SEED-VL-1.5



SEED-VL-1.5 (Referring)



帮我框选图中身着衣服为偶数的人,以 <bbox>x1 y1 x2 y2</bbox>的形式表示



帮我框选图中穿红色衣服的人,以<bbox>x1 y1 x2 y2</bbox>的形式表示

Summary



• The next step for object detection is moving toward multimodal large language models.

ChatRex

- Paper: https://arxiv.org/abs/2411.18363
- Code: https://github.com/IDEA-Research/ChatRex

RexSeek

- Paper: https://arxiv.org/abs/2503.08507
- Code: https://github.com/IDEA-Research/RexSeek

RexThinker

- Paper: https://arxiv.org/abs/2506.04034
- Code: https://github.com/IDEA-Research/Rex-Thinker











Taming Multimodal LLM for Object Perception

Thanks!