

Grounding DINO 1.5: Advance the “Edge: of Open-Set Object Detection

Qing Jiang

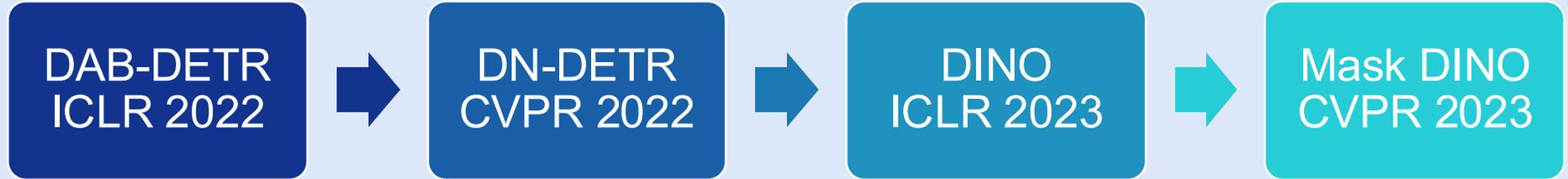
2024 6.14



Dr. Lei Zhang

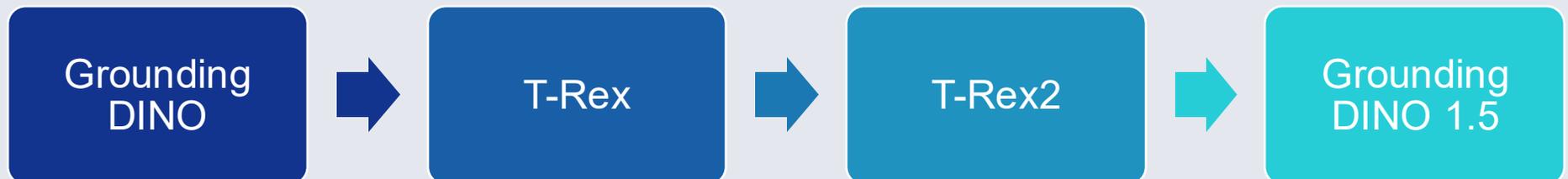
Chief Scientist

DETR Series Research in CVR



- Introduce Anchor Box to DETR
- Offer a better understanding of DETR Queries
- Introduce Denoising to DETR
- Reduce DETR training to 12 epochs
- Further improve DAB and DN
- Achieve SOTA on COCO Detection
- Unify Object Detection and Segmentation

Open-Set Object Detection Research in CVR



- Text prompt based open-set object detection model
- Visual prompt based open-set object detection model
- Combine both text prompt and visual prompt in one model
- SOTA text prompt based model with both Pro and Edge version models.

What is Object Detection?



What is Object Detection?



person. cup.
bowl. light.
chair.
coffee machine.
microwave.
refrigerator.
laptop. robot.
table

Detection results
from Grounding
DINO

Object-Centric understanding is the **perceptual basis** for machine-physical world interaction.

Solution for Hallucinations



: How many pigeons are there in this image



: In the image, there are approximately 59 pigeons.

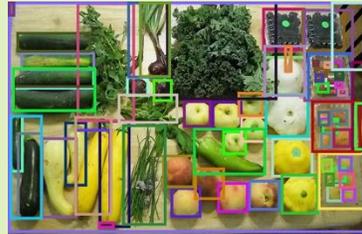


User provide visual prompt



: 79
T-Rex2

Fine-grained Perception



Box

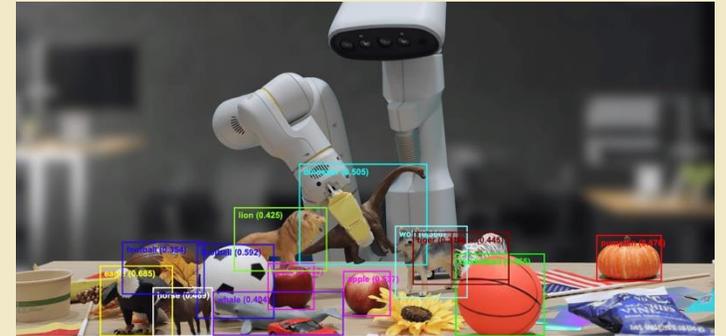


Box

VLM

This image depicted ... (Fine-grained long caption)

Eyes for Embodied AI



eagle . horse . lion . wolf . dinosaur . tiger . whale .
apple . football . basketball . pumpkin



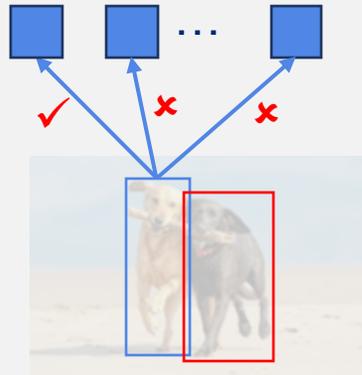
"move apple to cup with same color"

apple . red cup . green cup . green bag

RT-2

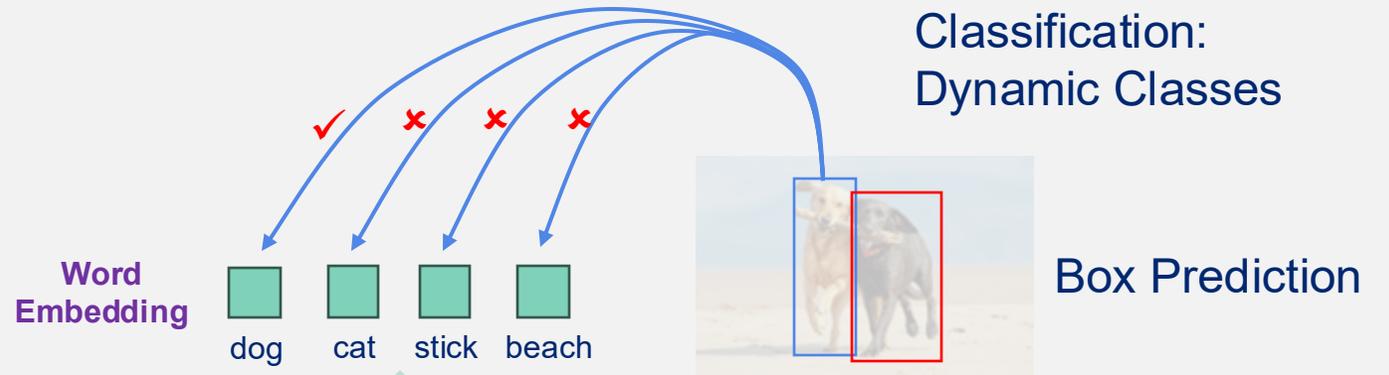
Paradigm Shift in Object Detection (Close-Set to Open-Set)

0:dog 1:cat ... 99:car



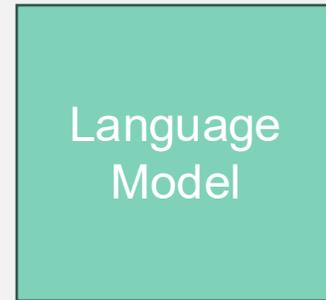
Classification:
Fixed # of Classes

Box Prediction



Classification:
Dynamic Classes

Box Prediction

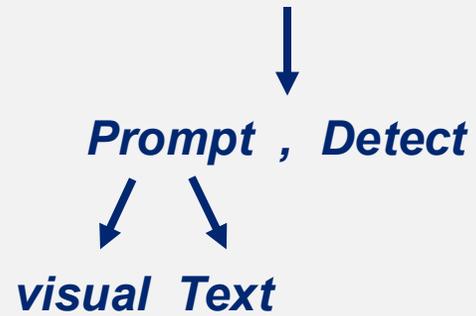


dog. cat. stick. beach
Language Prompt



Goal of Open-Set Object Detection

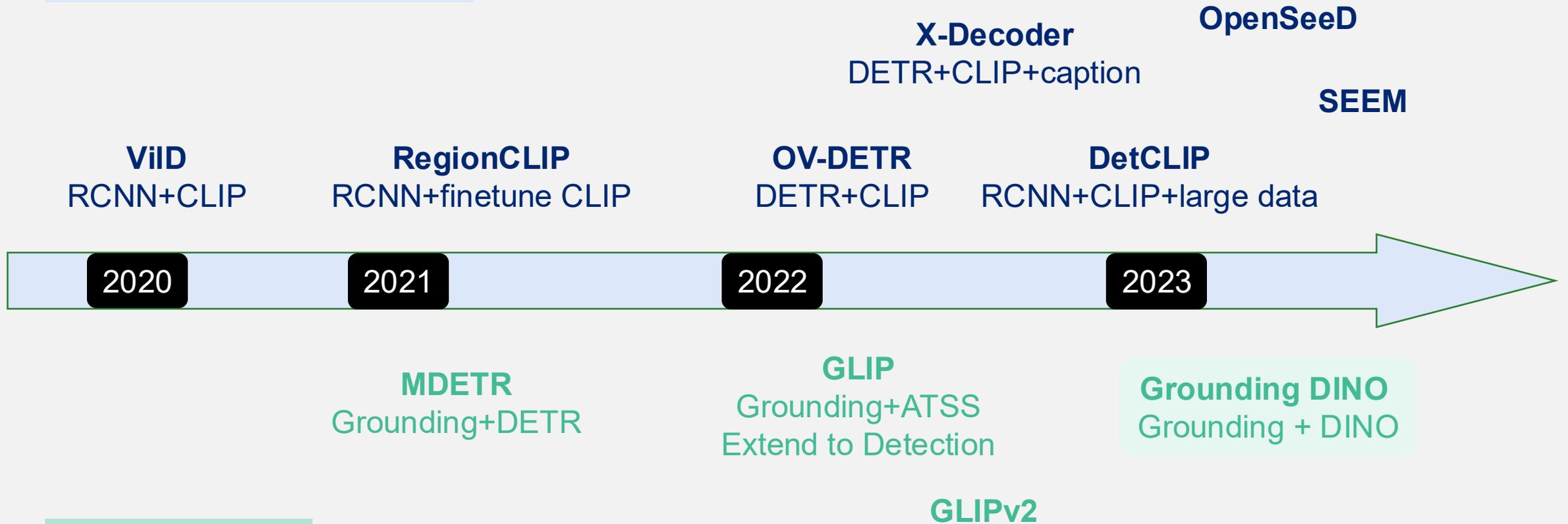
- Given an image and arbitrary prompts
- We expect a model predicts all objects mentioned in prompts without finetuning.



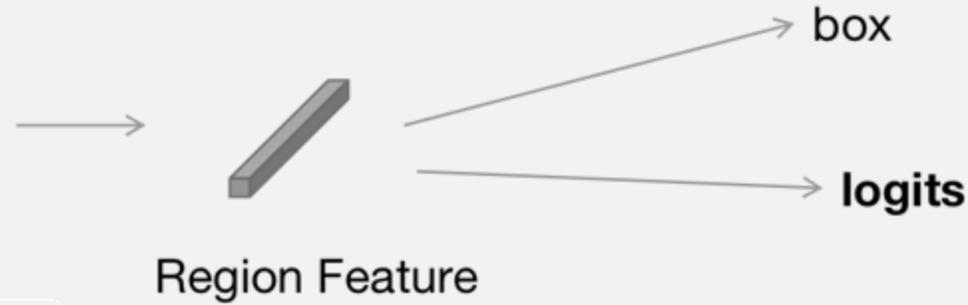
“armchair, blanket, lamp, carpet, couch, dog, floor, furniture, gray, green, living room, picture frame, pillow, plant, room, sit, stool, wood floor”

Two Paths to Open-Set Object Detection (text prompt based)

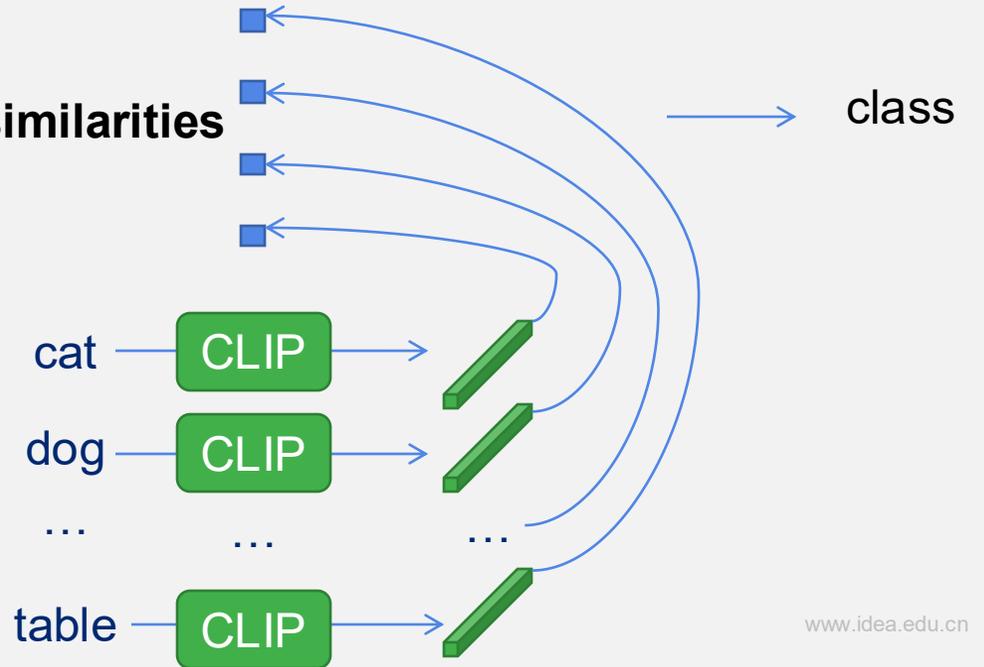
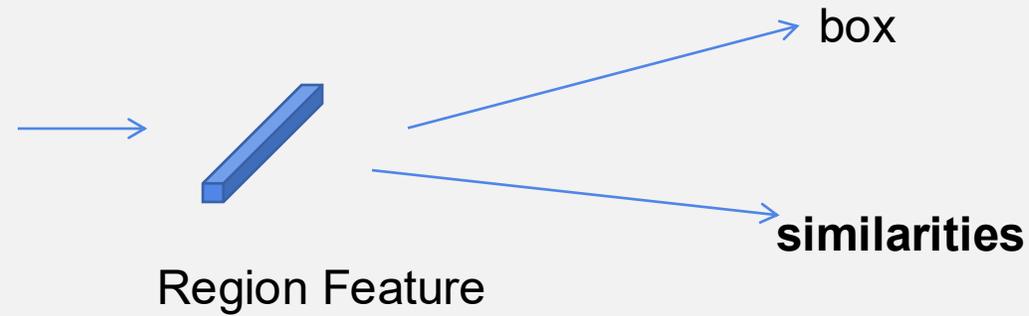
Path 1: Referring (CLIP-based)



Path1: Referring (CLIP-based) Open-Set Object Detection



Closed-Set Object Detection



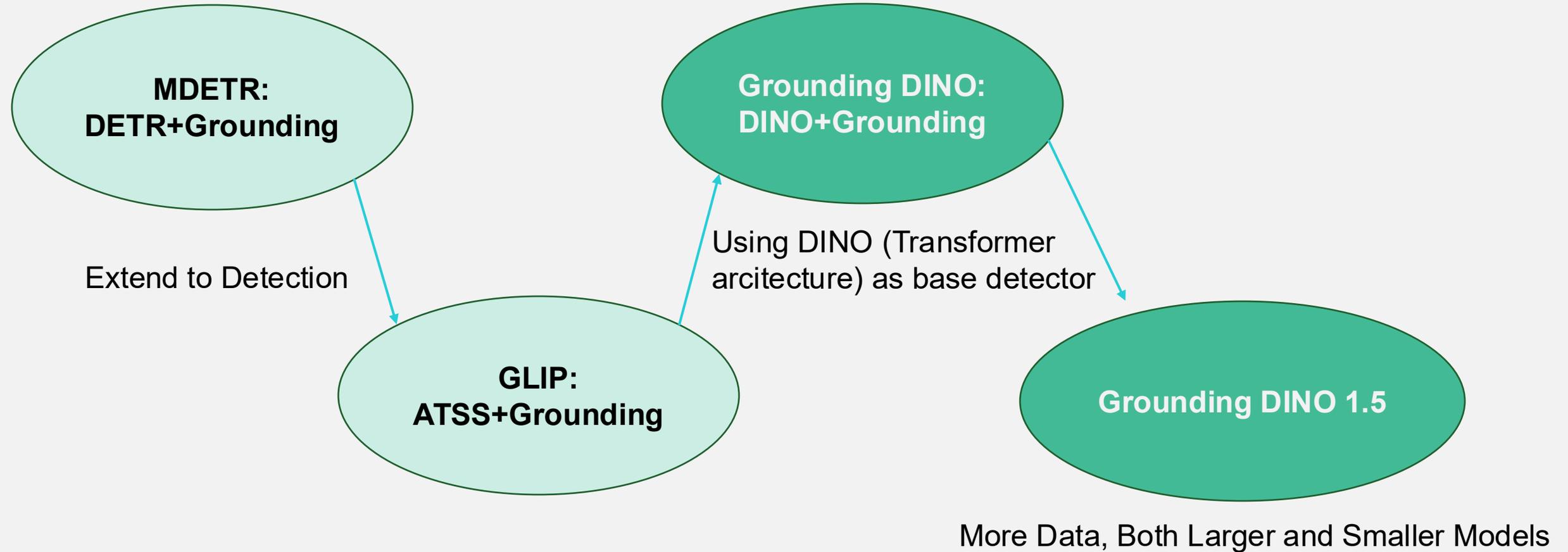
Referring (CLIP-based) Open-Set Object Detection

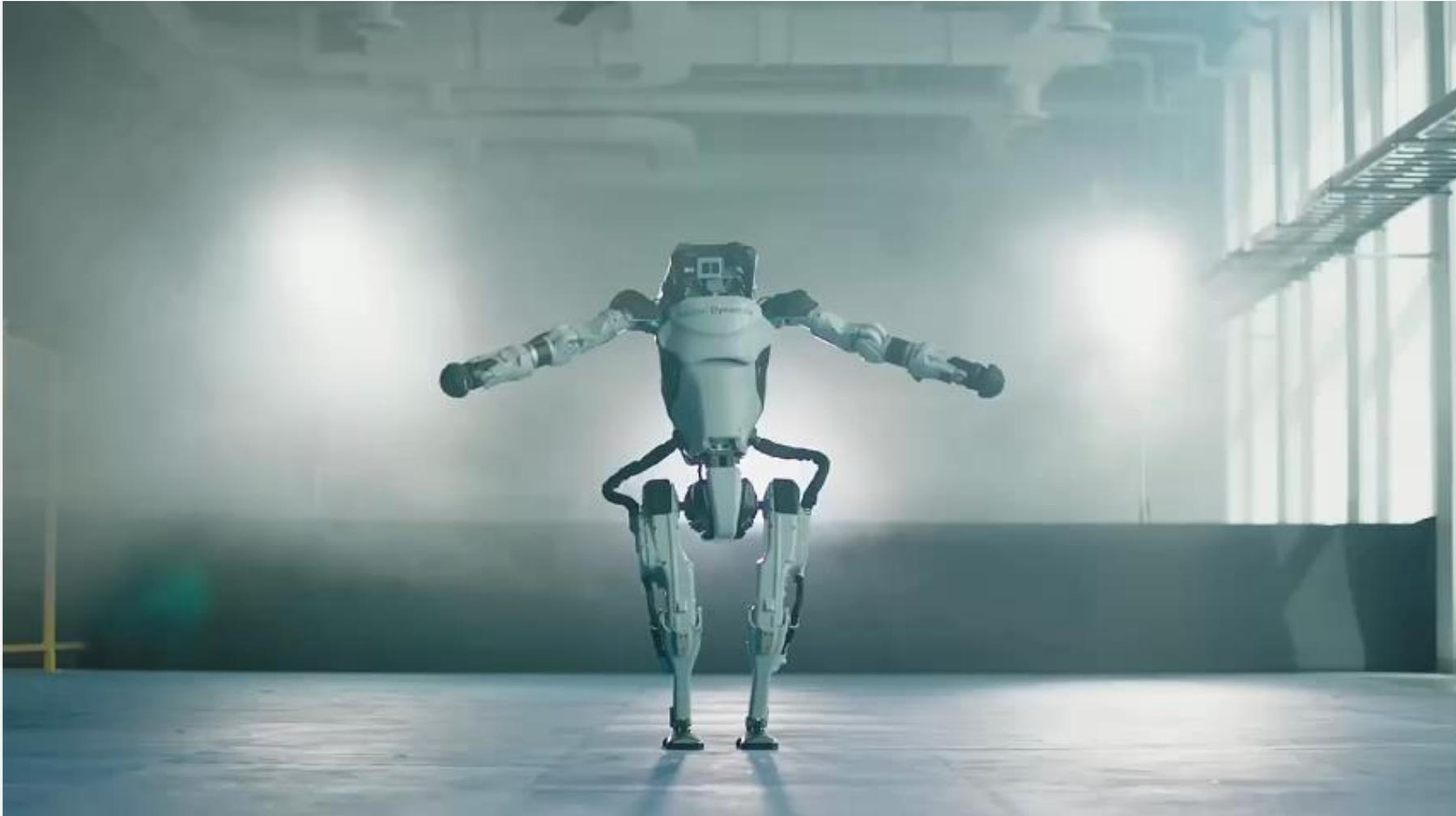
What is Grounding?

The image shows an individual lying prone on **grassy ground**, aiming a **bolt-action rifle** with their right eye close to the sight. This person is wearing a **military-style uniform** with a **steel helmet with netting**, suggesting a military or historical reenactment context. On the individual's back, you can see a **backpack** and what appears to be a **canteen**, both typical of military field equipment.

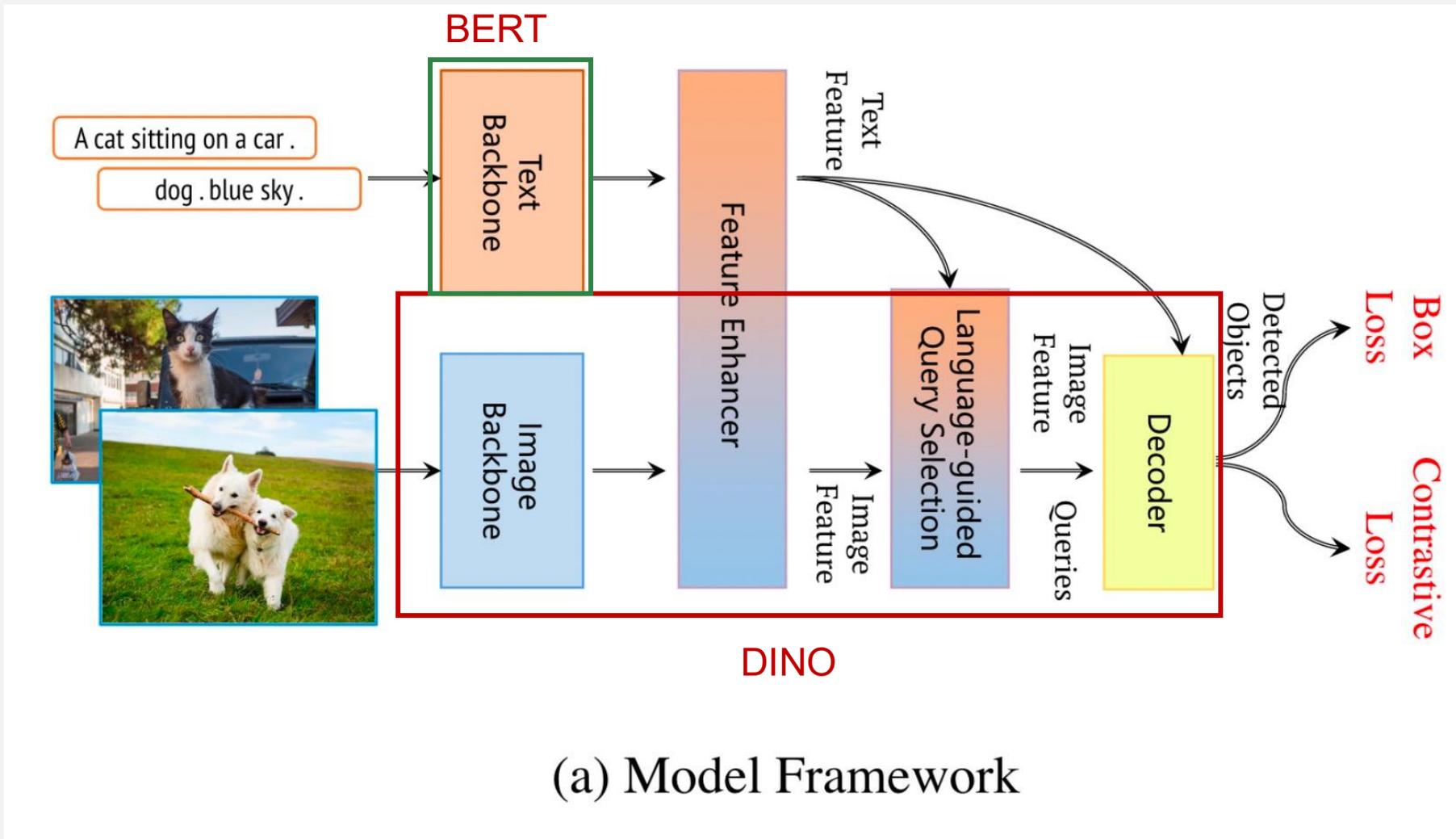


→ identifying the bounding boxes in an image that correspond to the **noun phrases** in a given sentence.

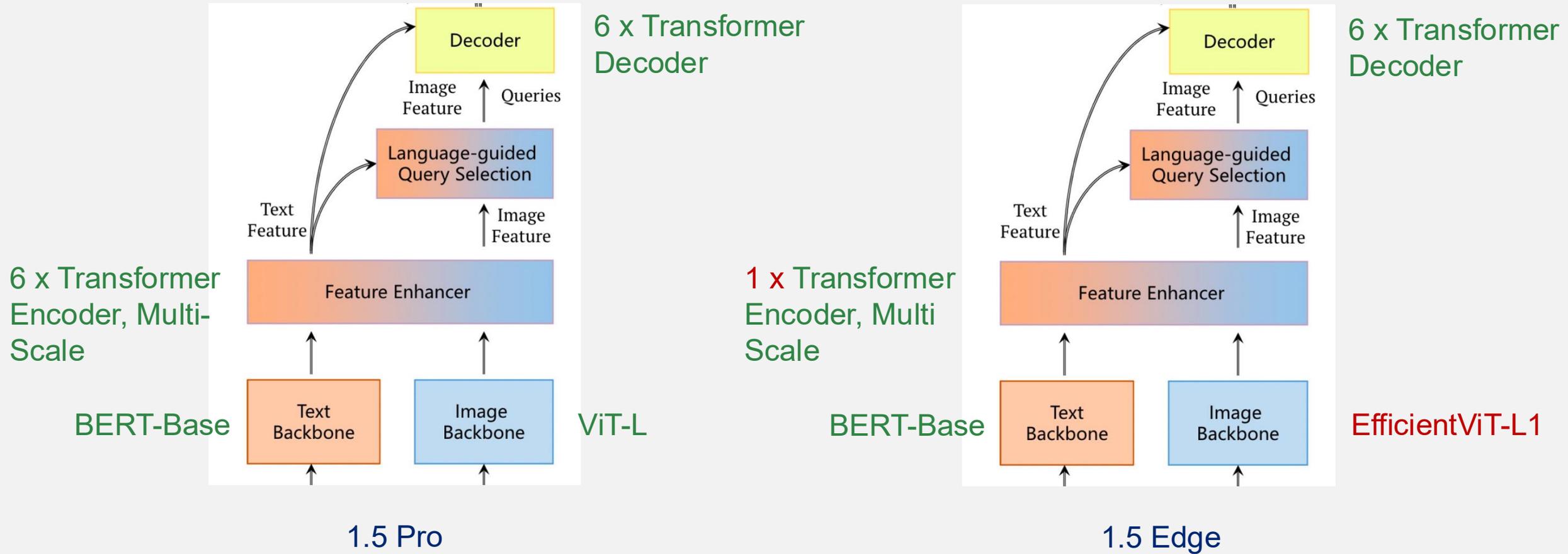




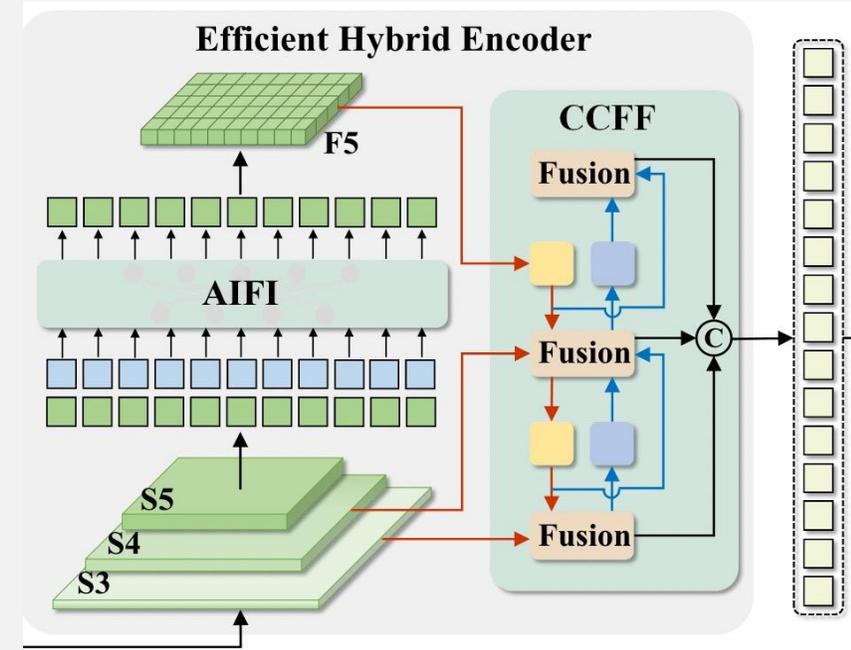
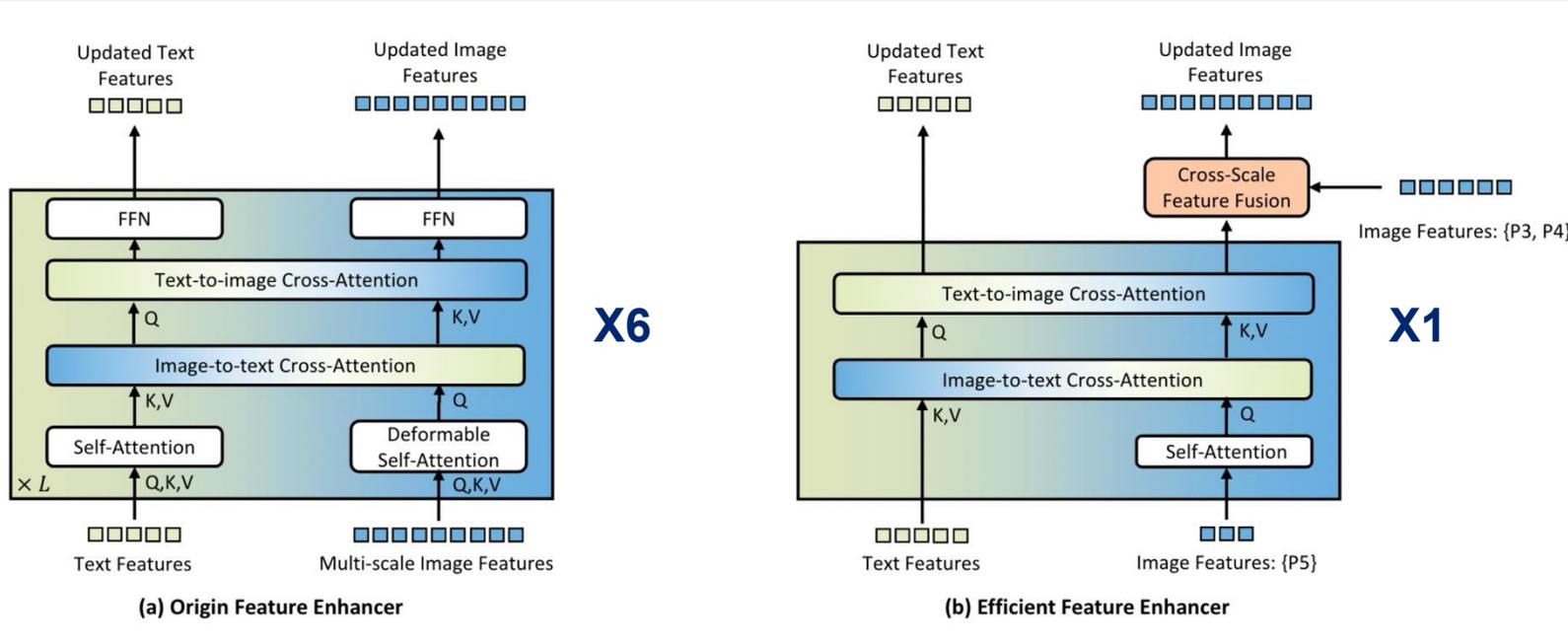
We use “edge” for its dual meaning both as in pushing the boundaries and as in running on edge devices.



Pro V.S. Edge: Overall Architecture



Pro V.S. Edge: Encoder



RT-DETR

Pro V.S. Edge: Running Time for Each Module (Pytorch Time)

Model	BERT	Backbone	Encoder	Decoder	FPS
Pro	0.008 (1.7%)	0.367 (79.2%)	0.073 (15.7%)	0.015 (3.34%)	2.16
Edge	0.009 (15.3%)	0.012 (18.7%)	0.021 (32.75%)	0.021 (33.3%)	15.9

ViT-L

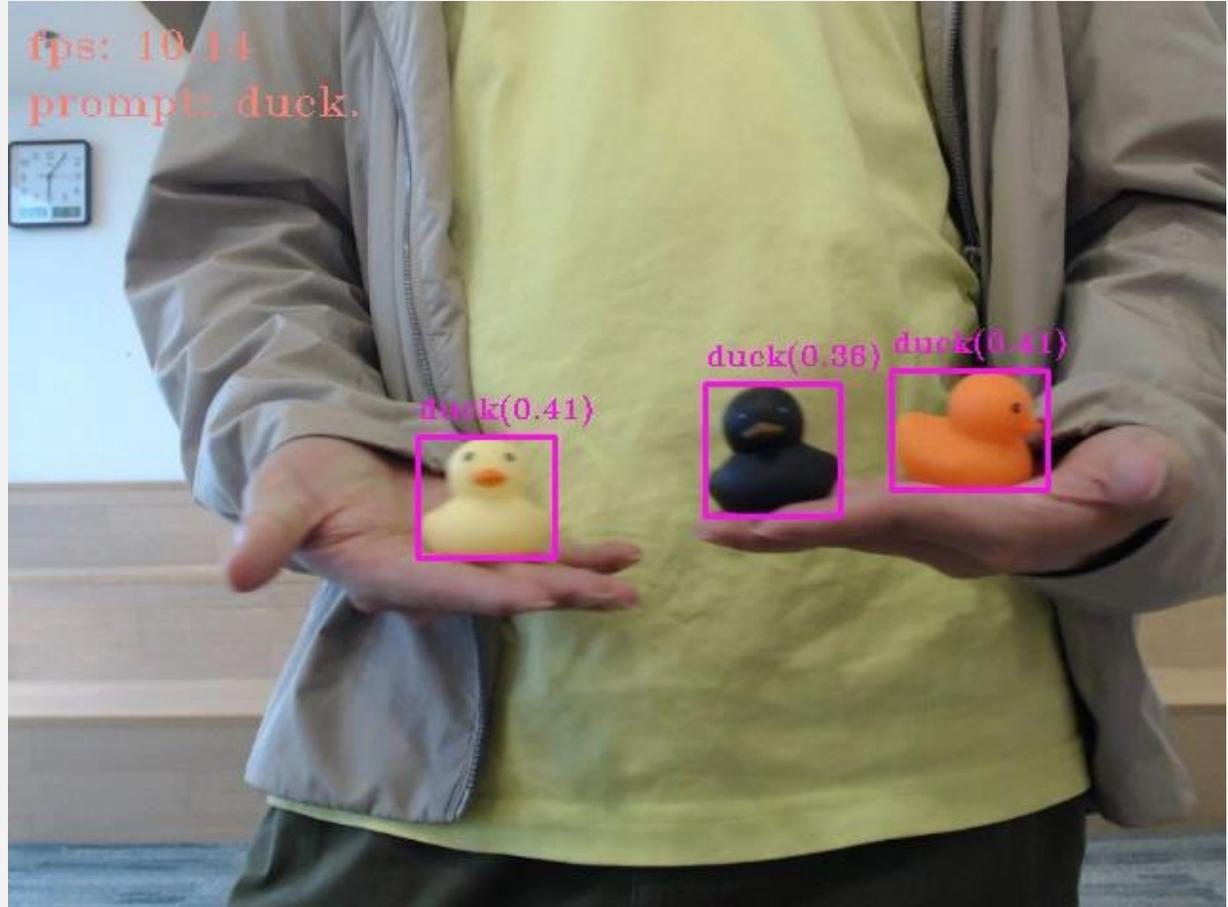


EfficientViT-L2



- Measured on a RTX 3090
- backbone takes the most time
- how to further optimize encoder and decoder time consumption is the next step

Deploy Edge on Edge Device (NVIDIA Orin NX)



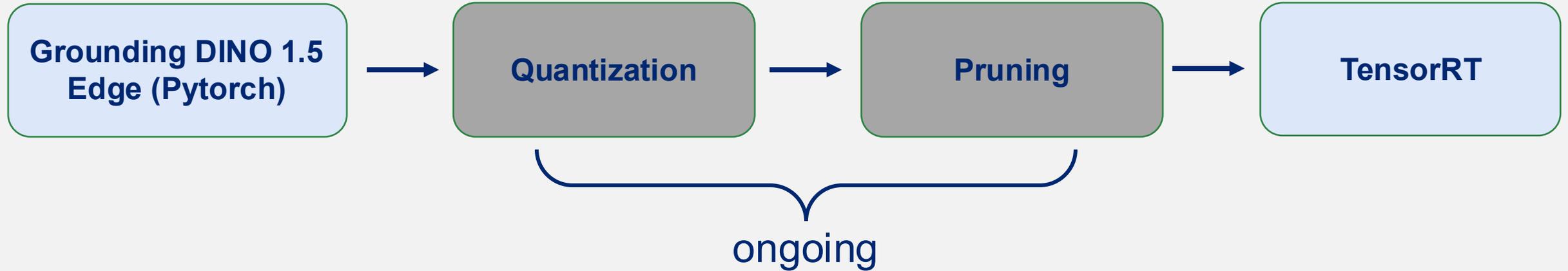
Deploy Edge on NVIDIA Orin NX

	Jetson AGX Orin series			Jetson Orin NX series		Jetson Orin Nano series			
	Jetson AGX Orin Developer Kit	Jetson AGX Orin 64GB	Jetson AGX Orin Industrial	Jetson AGX Orin 32GB	Jetson Orin NX 16GB	Jetson Orin NX 8GB	Jetson Orin Nano Developer Kit	Jetson Orin Nano 8GB	Jetson Orin Nano 4GB
AI Performance	275 TOPS	248 TOPS	200 TOPS	100 TOPS	70 TOPS	40 TOPS	20 TOPS		
GPU	2048-core NVIDIA Ampere architecture GPU with 64 Tensor Cores		1792-core NVIDIA Ampere architecture GPU with 56 Tensor Cores	1024-core NVIDIA Ampere architecture GPU with 32 Tensor Cores	1024-core NVIDIA Ampere architecture GPU with 32 Tensor Cores	1024-core NVIDIA Ampere architecture GPU with 32 Tensor Cores	512-core NVIDIA Ampere architecture GPU with 16 Tensor Cores		
GPU Max Frequency	1.3 GHz	1.2GHz	930MHz	918MHz	765MHz	625MHz			



Specification	Orin NX	RTX 3090
CUDA Cores	1024 cores	10496 cores
Tensor Cores	32 cores	328 cores
GPU Max Freq.	918MHZ	1695MHZ
TOPS	100 TOPS	~285TOPS

Only the TOPS of Orin NX is close to that of the 3090, which means the model should be quantized to INT8 for optimal performance.



Challenges:

- Remove deformable attention in decoder

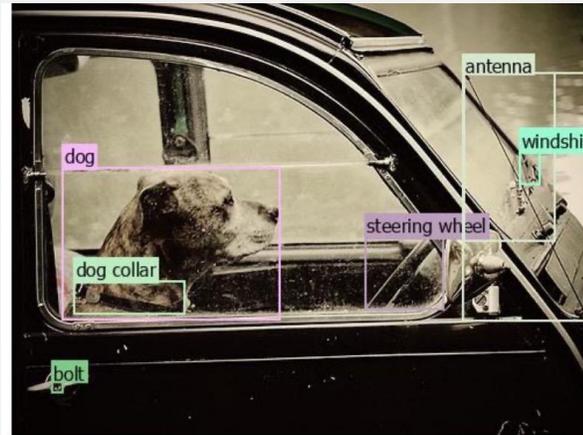
Results: Grounding DINO 1.5 Pro

Method	Backbone	Pre-training data	COCO		LVIS ^{minival}			LVIS ^{val}			ODinW35	ODinW13	
			AP _{all}	AP _{all}	AP _r	AP _c	AP _f	AP _{all}	AP _r	AP _c	AP _f	AP _{avg}	AP _{avg}
<i>Supervised Models (Pre-training data includes COCO, LVIS, etc.)</i>													
GLIPv2 [35]	Swin-H [32]	FourODs,COCO,GoldG,CC15M,SBU	60.6	50.1	-	-	-	-	-	-	-	55.5	
Grounding DINO [18]	Swin-L [19]	O365,OID,GoldG,Cap4M,COCO,RefC	60.7	33.9	22.2	30.7	38.8	-	-	-	-	-	
APE (B) [24]	ViT-L	COCO,LVIS,O365,OID,VG	57.7	62.5	-	-	-	57.0	-	-	-	29.4	59.8
APE (D) [24]	ViT-L [6]	COCO,LVIS,O365,OID,VG,RefC,SA-1B,GoldG,PhraseCut	58.3	64.7	-	-	-	59.6	-	-	-	28.8	57.9
GLEE-Pro [27]	ViT-L [6]	GLEE-merged-10M (COCO,LVIS,etc)	62.0	-	-	-	-	55.7	49.2	-	-	-	53.4
<i>Zero-shot Transfer Models</i>													
OWL-ViT [22]	ViT-L [5]	O365,OID,VG,LiT	42.2	-	-	-	-	34.6	31.2	-	-	-	-
MDETR [11]	ResNet101 [8]	COCO,GoldG	-	22.5	7.4	22.7	25.0	-	-	-	-	-	-
GLIP [16]	Swin-L	FourODs,GoldG,Cap24M	49.8	37.3	28.2	34.3	41.5	26.9	17.1	23.3	35.4	-	52.1
Grounding DINO [18]	Swin-T	O365,GoldG,Cap4M	48.4	27.4	18.1	23.3	32.7	-	-	-	-	22.3	49.8
Grounding DINO [18]	Swin-L	O365,OID,GoldG	52.5	-	-	-	-	-	-	-	-	26.1	56.9
OpenSeeD [34]	Swin-L	COCO,O365	-	23.0	-	-	-	-	-	-	-	15.2	-
UniDetector [26]	ResNet50 [8]	COCO,O365,OID	-	-	-	-	-	19.8	18.0	19.2	21.2	-	47.3
OmDet-Turbo-B [36]	ConvNeXt-B [20]	O365,GoldG,PhraseCut,Hake,HOI-A	<u>53.4</u>	34.7	-	-	-	-	-	-	-	<u>30.1</u>	54.7
OWL-ST [21]	CLIP L/14 [23]	WebLI2B	-	40.9	41.5	-	-	35.2	36.2	-	-	24.4	53.0
MQ-GLIP [28]	Swin-L	O365	-	43.4	34.5	41.2	46.9	34.7	26.9	32.0	41.3	23.9	54.1
DetCLIP [30]	Swin-L	O365,GoldG,YFCC1M	-	38.6	36.0	38.3	39.3	28.4	25.0	27.0	31.6	-	-
DetCLIPv2 [29]	Swin-L	O365,GoldG,CC15M	-	44.7	43.1	46.3	43.7	36.6	33.3	36.2	38.5	-	-
DetCLIPv3 [31]	Swin-L	O365,V3Det,GoldG,GranuCap50M	-	48.8	<u>49.9</u>	49.7	47.8	41.4	41.4	40.5	42.3	-	-
YOLO-World [3]	YOLOv8-L [10]	O365,GoldG,CC3M	45.1	35.4	27.6	34.1	38.0	-	-	-	-	-	-
DINOv [14]	Swin-L	COCO,SA-1B	-	-	-	-	-	-	-	-	-	15.7	-
T-Rex2 (visual) [9]	Swin-L	O365,OID,HierText,CrowdHuman,SA-1B	46.5	47.6	45.4	46.0	49.5	45.3	<u>43.8</u>	42.0	49.5	27.8	-
T-Rex2 (text) [9]	Swin-L	O365,OID,GoldG,CC3M,SBU,LAION	52.2	<u>54.9</u>	49.2	<u>54.8</u>	56.1	<u>45.8</u>	42.7	<u>43.2</u>	50.2	22.0	-
Grounding DINO 1.5 Pro (zero-shot)	ViT-L [6]	Grounding-20M	54.3	55.7	56.1	57.5	<u>54.1</u>	47.6	44.6	47.9	<u>48.7</u>	30.2	<u>58.7</u>

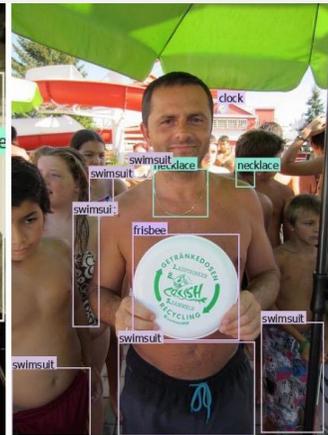
Results: Grounding DINO 1.5 Edge

Method	Backbone	Pre-training data	test size	COCO	LVIS ^{minival}				LVIS ^{val}				FPS(A100/TensorRT)	FPS(Orin NX)
					AP _{all}	AP _r	AP _c	AP _f	AP _{all}	AP _r	AP _c	AP _f		
<i>End-to-End Open-Set Object Detection</i>														
GLIP-T	Swin-T	O365,GoldG,Cap4M	800 × 1333	46.3	26.0	20.8	21.4	31.0	-	-	-	-	-	-
Grounding DINO-T	Swin-T	O365,GoldG,Cap4M	800 × 1333	48.4	27.4	18.1	23.3	32.7	-	-	-	-	9.4 / 42.6	1.1
<i>Real-time End-to-End Open-Set Object Detection</i>														
YOLO-Worldv2-S [†]	YOLOv8-S	O365,GoldG	640 × 640	-	22.7	16.3	20.8	25.5	17.3	11.3	14.9	22.7	47.4 / -	-
YOLO-Worldv2-M [†]	YOLOv8-M	O365,GoldG	640 × 640	-	30.0	25.0	27.2	33.4	23.5	17.1	20.0	30.1	42.7 / -	-
YOLO-Worldv2-L [†]	YOLOv8-L	O365,GoldG	640 × 640	-	33.0	22.6	32.0	35.8	26.0	18.6	23.0	32.6	37.4 / -	-
YOLO-Worldv2-L [†]	YOLOv8-L	O365,GoldG,CC3M-Lite	640 × 640	-	32.9	25.3	31.1	<u>35.8</u>	26.1	20.6	22.6	<u>32.3</u>	37.4 / -	-
OmDet-Turbo-T [‡]	Swin-T	O365,GoldG	640 × 640	42.5	30.3	-	-	-	-	-	-	-	21.5 / 140.0	-
Grounding DINO 1.5 Edge	EfficientViT-L1	Grounding-20M	640 × 640	<u>42.9</u>	<u>33.5</u>	<u>28.0</u>	<u>34.3</u>	33.9	<u>27.3</u>	<u>26.3</u>	<u>25.7</u>	29.6	21.7 / 111.6	10.7
Grounding DINO 1.5 Edge	EfficientViT-L1	Grounding-20M	800 × 1333	45.0	36.2	33.2	36.6	36.3	29.3	28.1	27.6	31.6	18.5 / 75.2	5.5

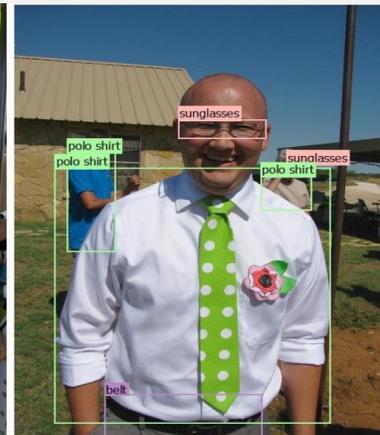
Visualization Results (Long-tailed Object Detection)



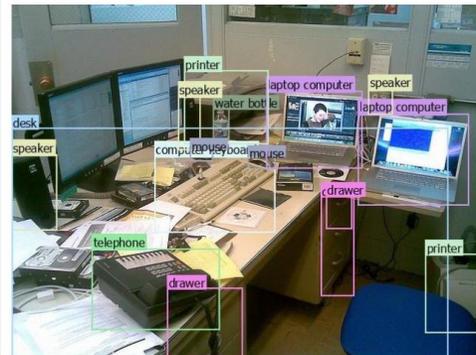
bolt . antenna . dog . dog collar .
steering wheel . windshield wiper



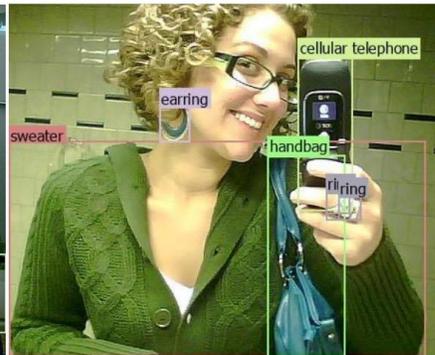
swimsuit . necklace .
clock . frisbee



polo shirt . belt . sunglasses



printer . speaker . computer
keyboard . mouse . laptop
computer . water bottle .
drawer . desk



cellular telephone . ring .
sweater . earring . handbag



cow

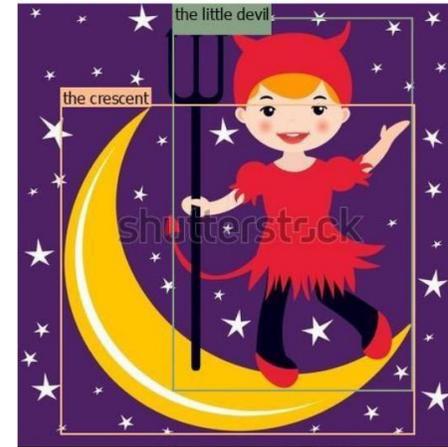
Visualization Results (Short Caption Object Detection)



a neon sign of goodbye is placed in the center of the wall.



a yellow boat sitting in the snow near some trees.



the little devil is flying on the crescent in the night.



an owl is sitting on the branch with gifts.



two hands hold the globe.



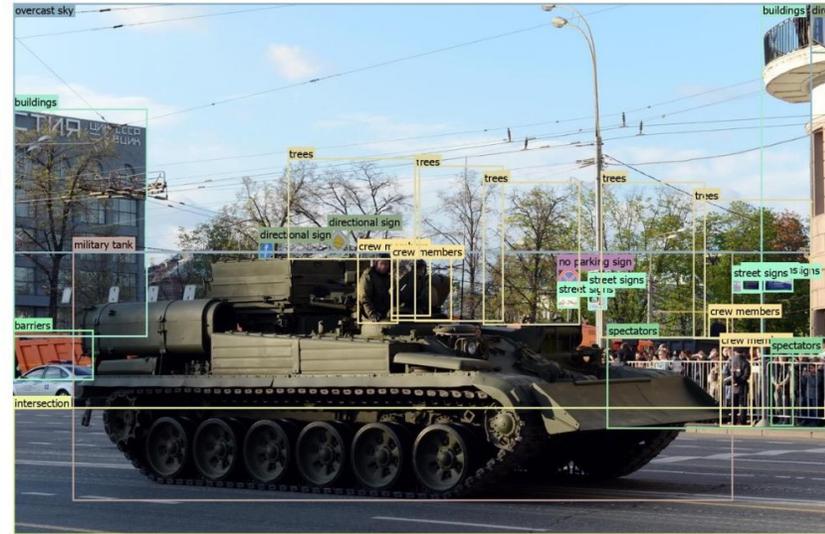
a black and yellow snake in its hand with a white ring.

Visualization Results (Long Caption Object Detection)

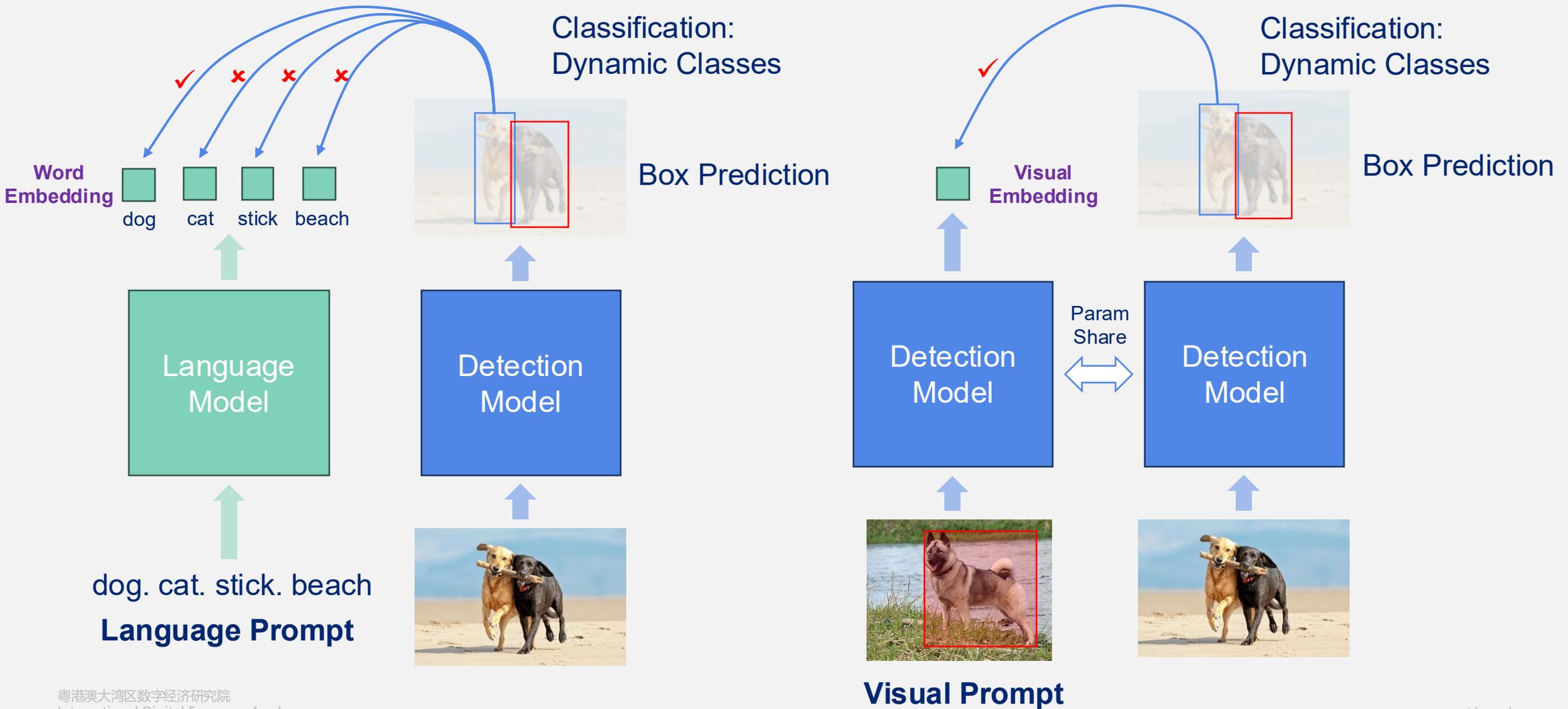


The image shows a **car** on display at a motor show, prominently featuring the **Fiat logo** in red along the side. The car is painted gray with a glossy finish, highlighted by orange and black racing stripes and **red accents** on the front splitter. Its sporty design includes a large rear spoiler and **white multi-spoke wheels**, suggesting additional aerodynamic features. In the background, several **people** navigate the space, including an individual in a **blue jacket with a logo**, likely an event staff member.

The image shows a **military tank** rolling along a city street, with **crew members** visible on top. Painted in camouflage and equipped with a long barrel, the tank navigates through an **intersection**, suggesting a setting likely in Europe or Russia. The surroundings, marked by **street signs**, a **no parking sign**, and a **directional sign** pointing right, indicate the tank's participation in a parade or military demonstration. **Spectators**, including adults and children, stand behind **barriers** on the sidewalk, observing the tank against a backdrop of **buildings** and **leafless trees**, under an **overcast sky**.

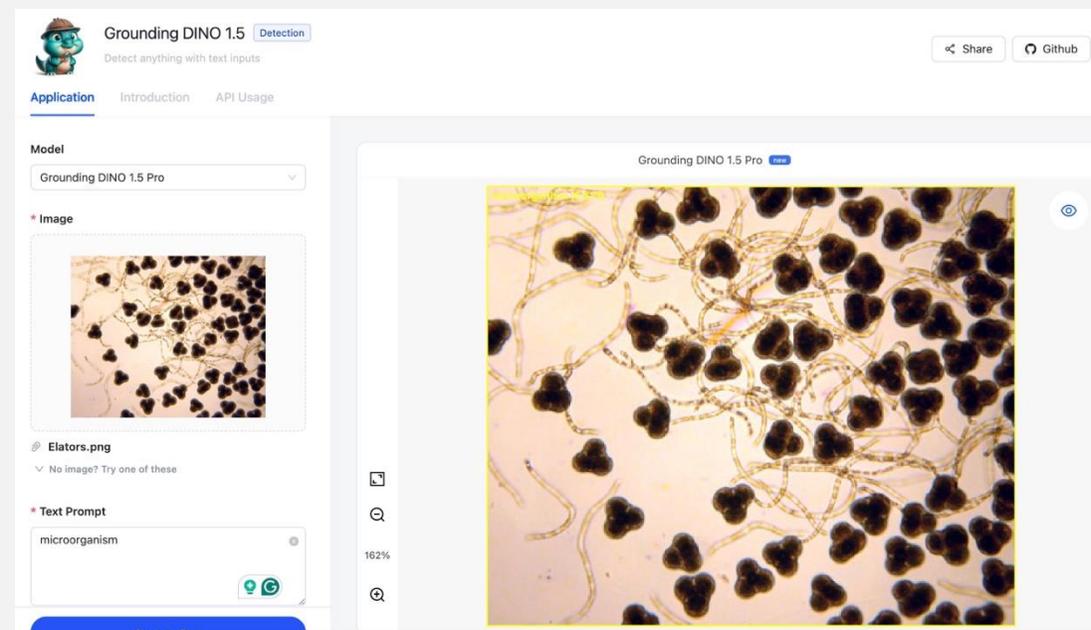
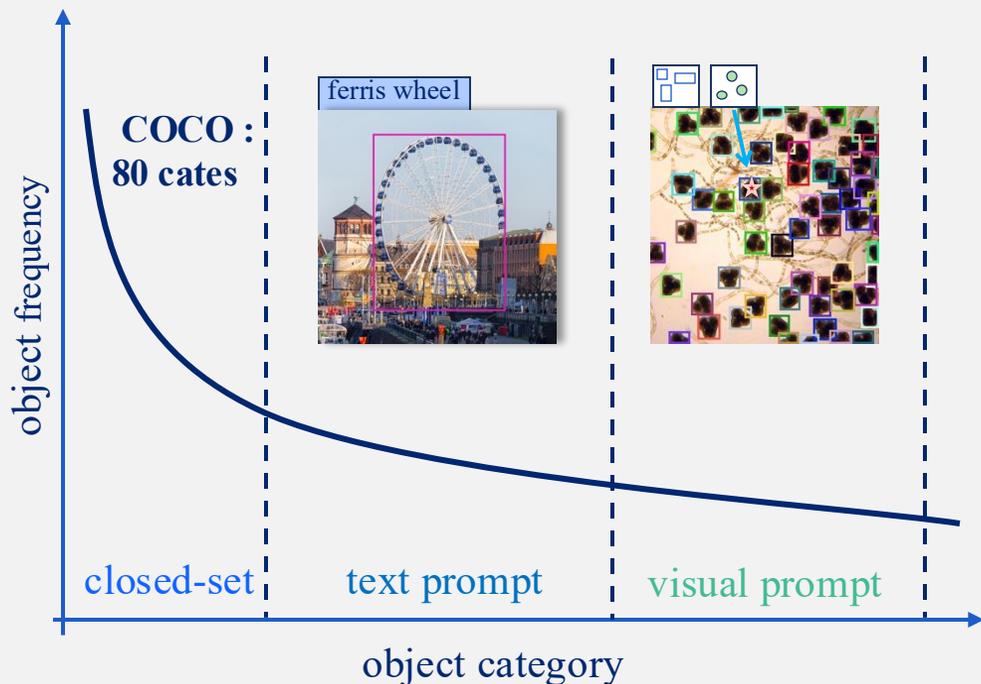


Text Prompt v.s. Visual Prompt



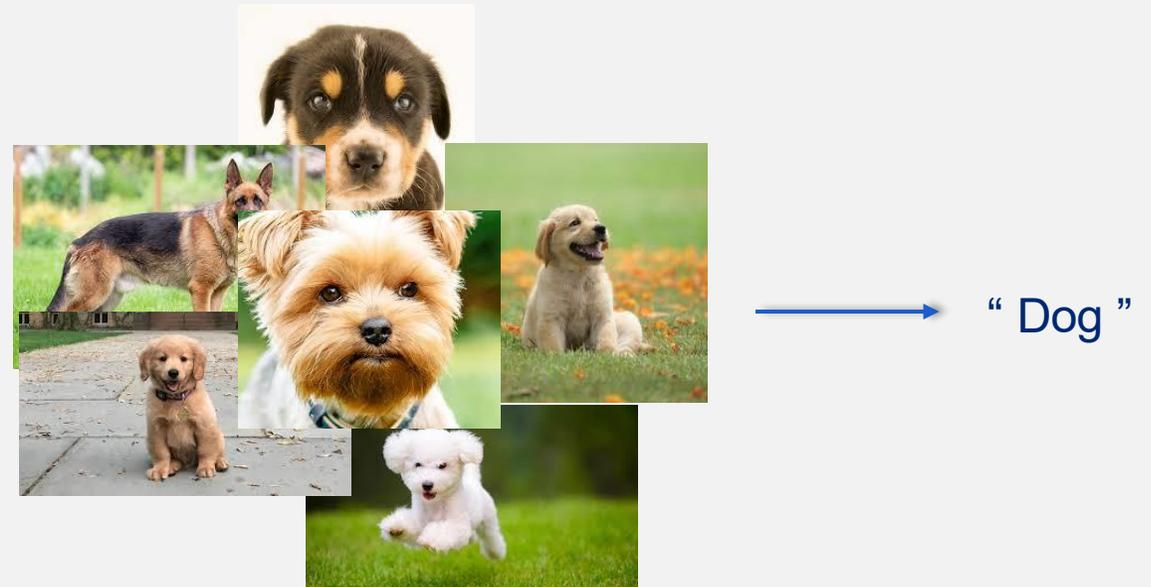
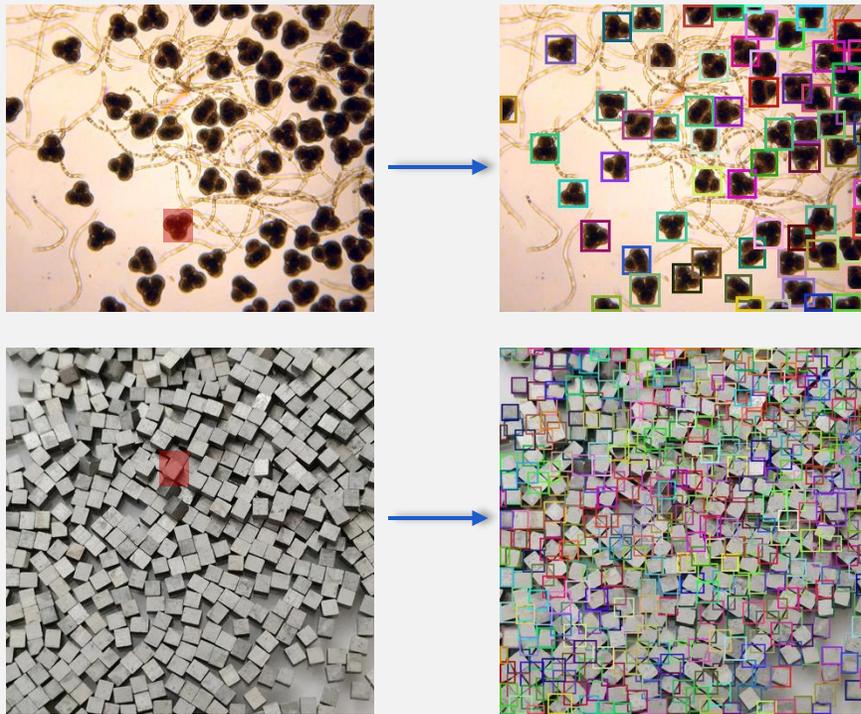
Text Prompt

- describe objects in natural language
- require modality alignment, suffers from long-tailed data shortage
- fall short in describe object that are hard to describe in language



Visual Prompt

- describe objects through visual examples
- less effective at capturing the general concept



require many examples to convey a general concept

T-Rex2: Combine both Text Prompt and Visual Prompt

T-Rex2: Towards Generic Object Detection via Text-Visual Prompt Synergy

Qing Jiang^{1,2}, Feng Li^{1,3}, Zhaoyang Zeng¹, Tianhe Ren¹, Shilong Liu^{1,4}, Lei Zhang^{1†}

¹International Digital Economy Academy (IDEA)

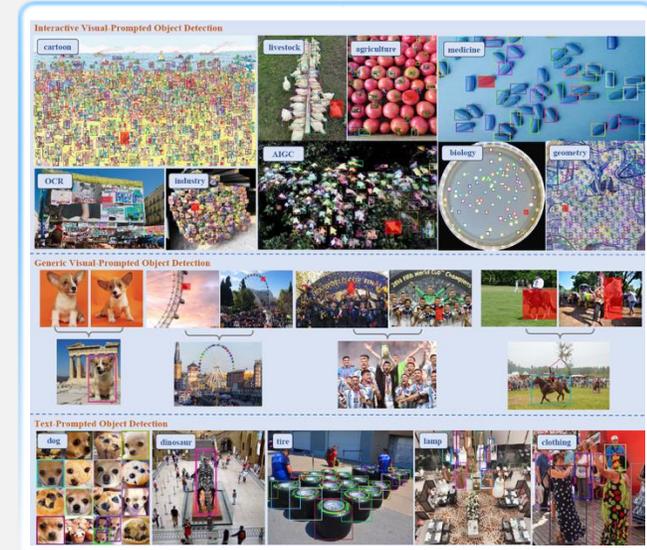
²SCUT ³HKUST ⁴Tsinghua

mountchicken@outlook.com, fliay@connect.ust.hk, lius120@mails.tsinghua.edu.cn

{rentianhe, zengzhaoyang, leizhang}@idea.edu.cn

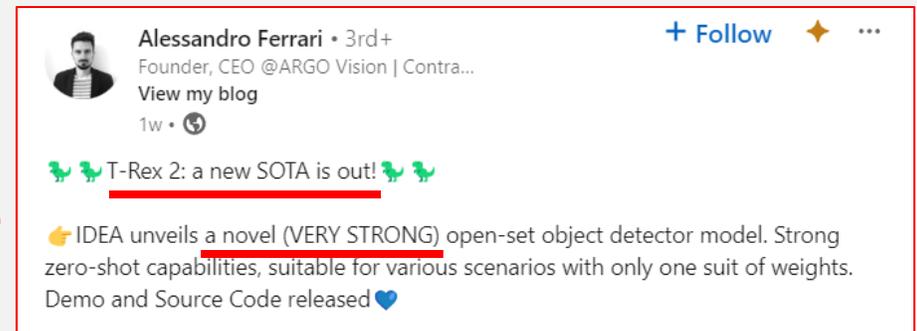
[trex-counting.github.io](https://github.com/trex-counting)

ECCV 2024



Method	Prompt Type	Backbone	COCO-Val	LVIS								ODinW		Roboflow100
			Zero-Shot	Zero-Shot								Zero-Shot		Zero-Shot
			val-80	minival-804				val-1203				35val		100val
			AP	AP	AP_f	AP_c	AP_r	AP	AP_f	AP_c	AP_r	AP_{avg}	AP_{med}	AP_{avg}
GLIP-T [19]	Text	Swin-T	46.7	26.0	31.0	21.4	20.8	17.2	25.5	12.5	10.1	19.6	5.1	-
GLIP-L [19]	Text	Swin-L	49.8	37.3	41.5	34.3	28.2	26.9	35.4	23.3	17.1	23.4	11.0	8.6
Grounding DINO [24]	Text	Swin-T	48.4	27.4	32.7	23.3	18.1	-	-	-	-	22.3	11.9	-
Grounding DINO [24]	Text	Swin-L	52.5	33.9	38.8	30.7	22.2	-	-	-	-	26.1	18.4	-
DetCLIPv2 [47]	Text	Swin-T	-	40.4	40.0	41.7	36.0	-	-	-	-	-	-	-
DetCLIPv2 [47]	Text	Swin-L	-	44.7	43.7	46.3	43.1	-	-	-	-	-	-	-
DINOv [17]	Visual-G	Swin-T	-	-	-	-	-	-	-	-	-	14.9	5.4	-
DINOv [17]	Visual-G	Swin-L	-	-	-	-	-	-	-	-	-	15.7	4.8	-
T-Rex2	Text	Swin-T	45.8	42.8	46.5	39.7	37.4	34.8	41.2	31.5	29.0	18.0	4.7	8.2
T-Rex2	Visual-G	Swin-T	38.8	37.4	41.8	33.9	29.9	<u>34.9</u>	41.1	30.3	<u>32.4</u>	<u>23.6</u>	<u>17.5</u>	<u>17.4</u>
T-Rex2	Text	Swin-L	<u>52.2</u>	54.9	56.1	54.8	49.2	45.8	50.2	43.2	42.7	22.0	7.3	10.5
T-Rex2	Visual-G	Swin-L	46.5	47.6	49.5	46.0	45.4	45.3	49.5	42.0	43.8	27.8	20.5	18.5

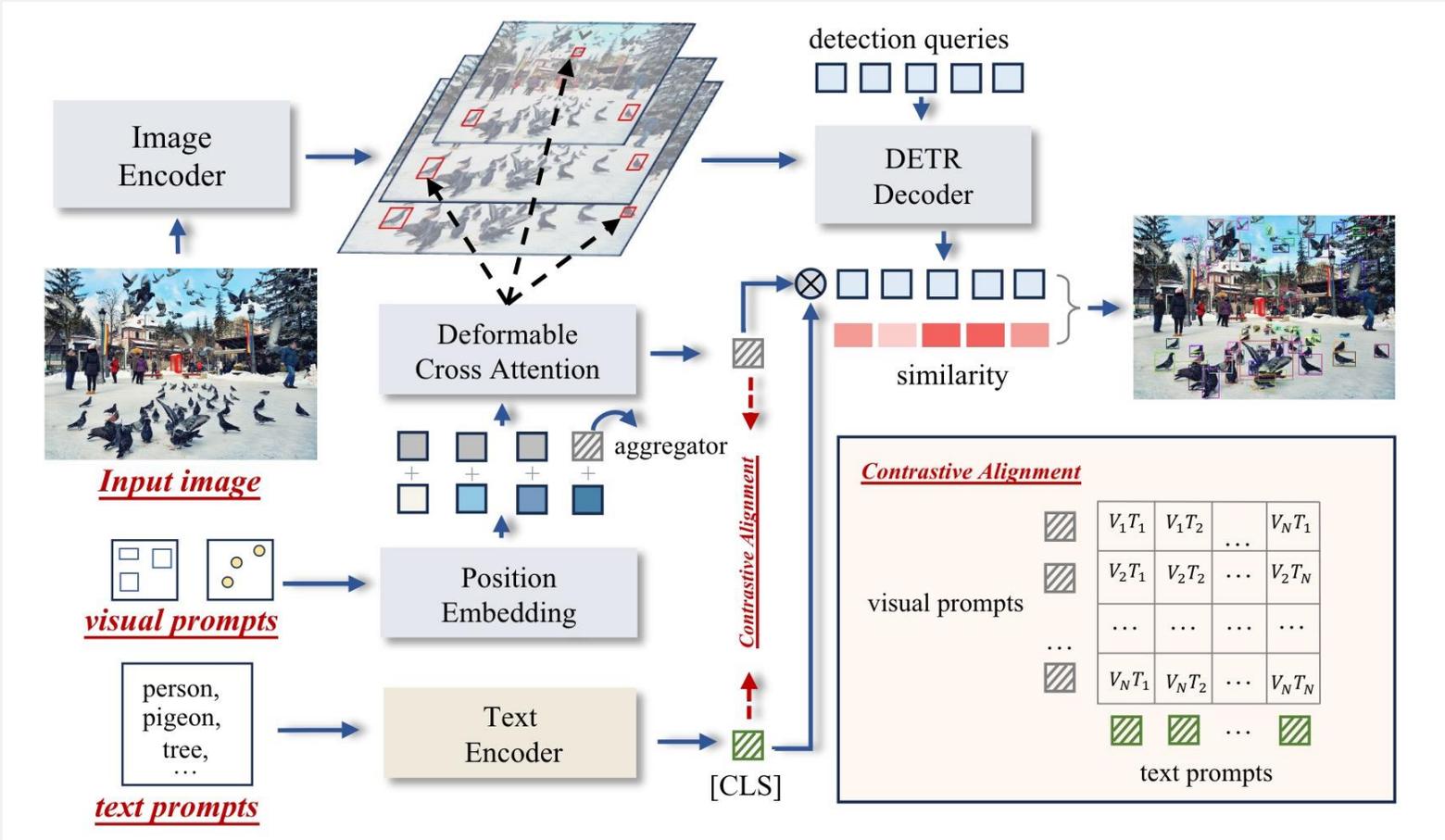
Table 1. **One suit of weights** for zero-shot object detection. **Red** denotes regions where text prompt excels over visual prompt, while **green** signifies regions favoring visual prompts.



T-Rex2: Combine both Text Prompt and Visual Prompt



T-Rex2: Combine both Text Prompt and Visual Prompt



DINO-based End-to-End model

Visual Prompt Encoder: Deformable Cross Attention

$$B = \text{Linear}(\text{PE}(b_1, \dots, b_K); \theta_B) : \mathbb{R}^{K \times 4D} \rightarrow \mathbb{R}^{K \times D}$$

$$P = \text{Linear}(\text{PE}(p_1, \dots, p_K); \theta_P) : \mathbb{R}^{K \times 2D} \rightarrow \mathbb{R}^{K \times D}$$

$$Q = \begin{cases} \text{Linear}(\text{CAT}([C; C'], [B; B'])); \varphi_B, \text{ box} \\ \text{Linear}(\text{CAT}([C; C'], [P; P'])); \varphi_P, \text{ point} \end{cases}$$

$$Q'_j = \begin{cases} \text{MSDeformAttn}(Q_j, b_j, \{f_i\}_{i=1}^L), \text{ box} \\ \text{MSDeformAttn}(Q_j, p_j, \{f_i\}_{i=1}^L), \text{ point} \end{cases}$$

$$V = \text{FFN}(\text{SelfAttn}(Q'))[-1]$$

Text Prompt Encoder: CLIP

Modality Alignment: Contrastive Learning

$$\mathcal{L}_{align} = -\frac{1}{K} \sum_{i=1}^K \log \frac{\exp(v_i \cdot t_i)}{\sum_{j=1}^K \exp(v_i \cdot t_j)}$$

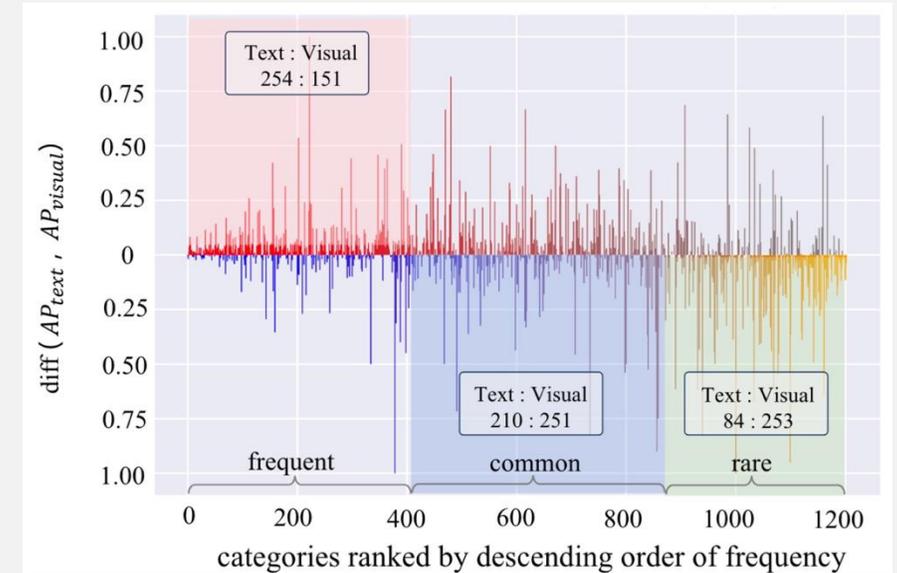
Joint prompt leads to generic object detection

Zero-Shot Generic Object Detection

Method	Prompt Type	Backbone	COCO-Val	LVIS								ODinW		Roboflow100
			Zero-Shot	Zero-Shot								Zero-Shot		Zero-Shot
			val-80	minival-804				val-1203				35val		100val
	AP	AP	AP_f	AP_c	AP_r	AP	AP_f	AP_c	AP_r	AP_{avg}	AP_{med}	AP_{avg}		
GLIP-T [19]	Text	Swin-T	46.7	26.0	31.0	21.4	20.8	17.2	25.5	12.5	10.1	19.6	5.1	-
GLIP-L [19]	Text	Swin-L	49.8	37.3	41.5	34.3	28.2	26.9	35.4	23.3	17.1	23.4	11.0	8.6
Grounding DINO [24]	Text	Swin-T	48.4	27.4	32.7	23.3	18.1	-	-	-	-	22.3	11.9	-
Grounding DINO [24]	Text	Swin-L	52.5	33.9	38.8	30.7	22.2	-	-	-	-	26.1	18.4	-
DetCLIPv2 [47]	Text	Swin-T	-	40.4	40.0	41.7	36.0	-	-	-	-	-	-	-
DetCLIPv2 [47]	Text	Swin-L	-	44.7	43.7	46.3	43.1	-	-	-	-	-	-	-
DINOv [17]	Visual-G	Swin-T	-	-	-	-	-	-	-	-	-	14.9	5.4	-
DINOv [17]	Visual-G	Swin-L	-	-	-	-	-	-	-	-	-	15.7	4.8	-
T-Rex2	Text	Swin-T	45.8	42.8	46.5	39.7	37.4	34.8	41.2	31.5	29.0	18.0	4.7	8.2
T-Rex2	Visual-G	Swin-T	38.8	37.4	41.8	33.9	29.9	34.9	41.1	30.3	32.4	23.6	17.5	17.4
T-Rex2	Text	Swin-L	<u>52.2</u>	54.9	56.1	54.8	49.2	45.8	50.2	43.2	42.7	22.0	7.3	10.5
T-Rex2	Visual-G	Swin-L	46.5	47.6	49.5	46.0	45.4	45.3	49.5	42.0	43.8	27.8	20.5	18.5

common and frequent case
rare and novel case
Text prompt better
Visual prompt better

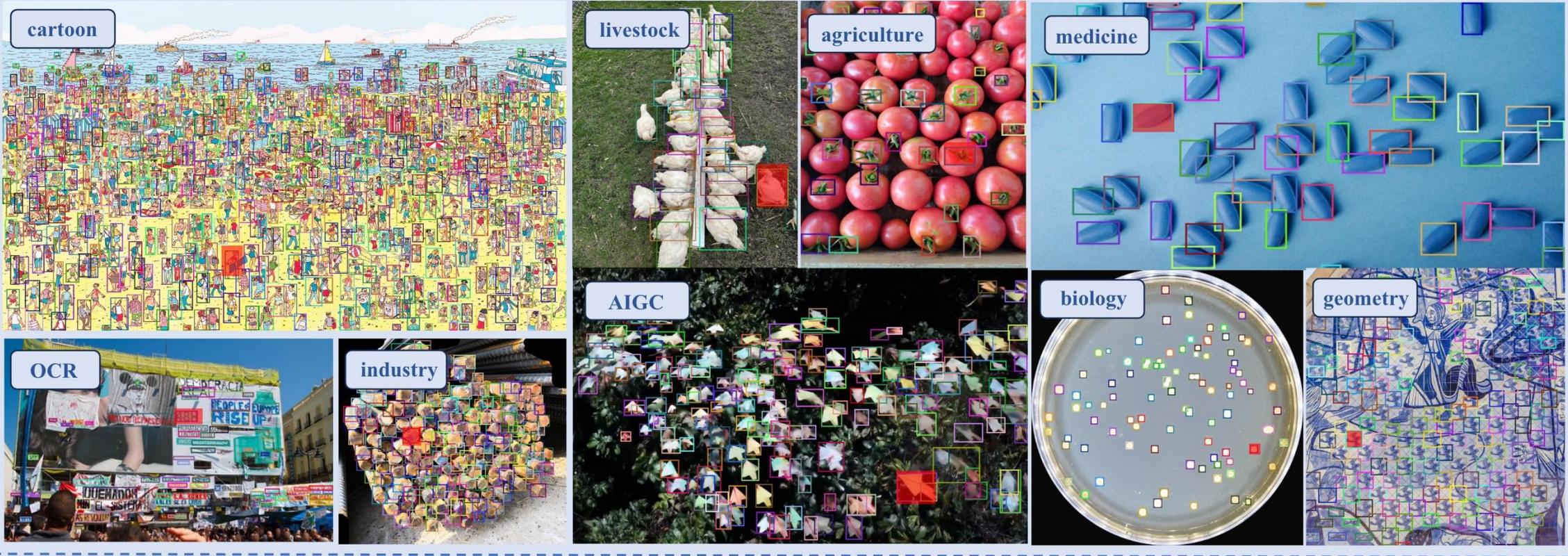
Text prompt v.s. Visual prompt on LVIS



- Text prompt is good at common and frequent object, while visual prompt succeed in rare and novel scenarios.

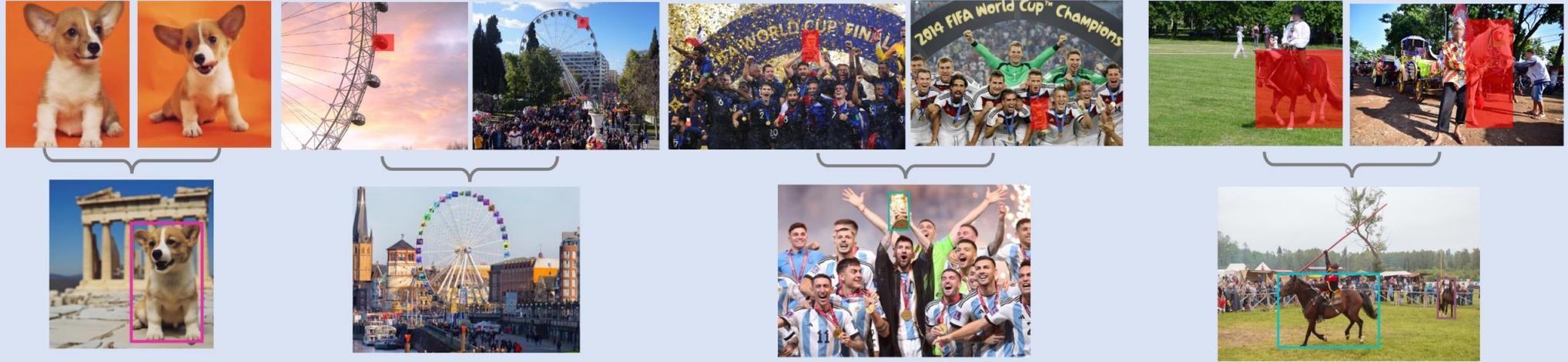
Joint prompt leads to generic object detection

Interactive Visual-Prompted Object Detection

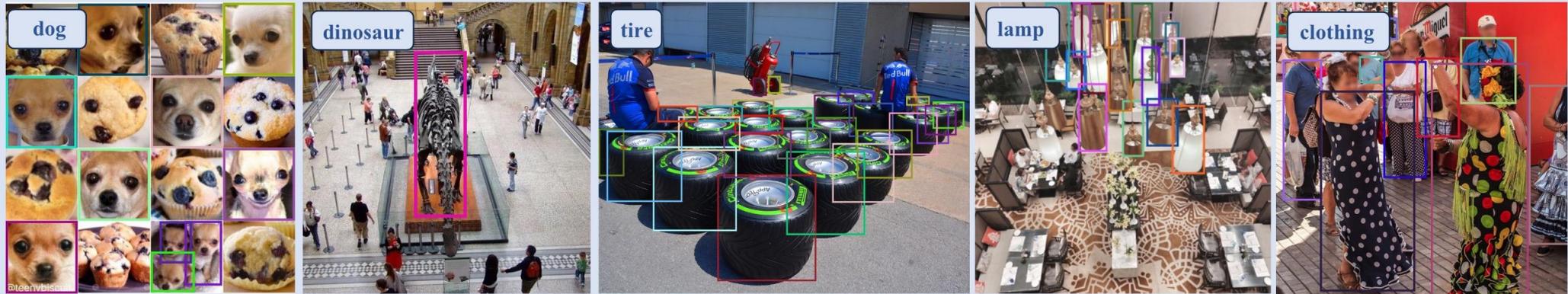


Joint prompt leads to generic object detection

Generic Visual-Prompted Object Detection



Text-Prompted Object Detection



- **Open-set** detection is the next BIG problem after closed-set detection
- **Prompt** is a new way to transform open-set detection
 - Text prompt: effective to cover head and middle concepts
 - Visual prompt: effective to cover more long tailed concepts
- **Grounded** understanding is key to multimodality intelligence



Grounding DINO 1.5

Detection

Detect anything with text inputs



Interactive Visual Prompt - Image

Detection

Interactive object detection and counting system based on T-Rex model

<https://deepdataspace.com/>

Thanks!

Qing Jiang

2024 6.14