



---

# UFO: A Unified Approach to Fine-grained Visual Perception via Open-ended Language Interface

---

**Hao Tang<sup>1,2</sup>   Chenwei Xie<sup>2</sup>   Haiyang Wang<sup>1</sup>   Xiaoyi Bao<sup>2,3</sup>**  
**Tingyu Weng<sup>2</sup>   Pandeng Li<sup>2</sup>   Yun Zheng<sup>2†</sup>   Liwei Wang<sup>1,4,5†</sup>**

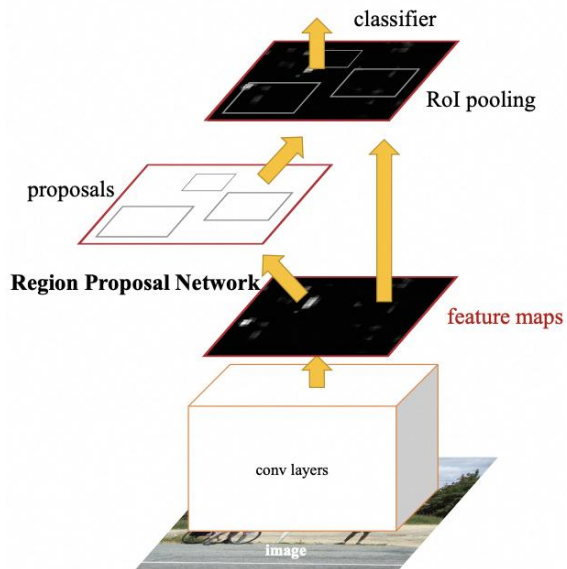
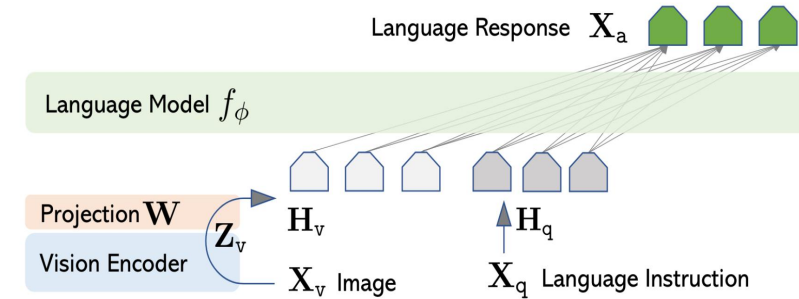
<sup>1</sup>Center for Data Science, Peking University   <sup>2</sup>Alibaba Group

<sup>3</sup> CASIA   <sup>4</sup> Center for Machine Learning Research, Peking University

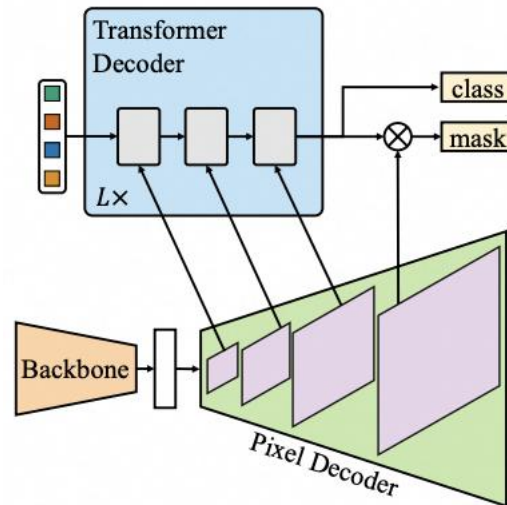
<sup>5</sup> State Key Laboratory of General Artificial Intelligence, Peking University, Beijing, China

# Motivation

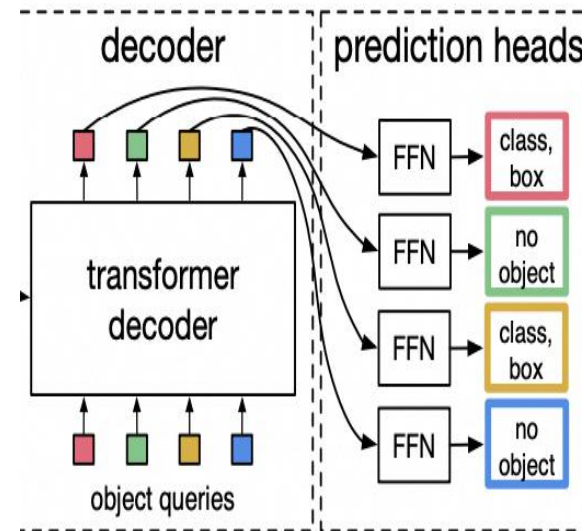
- Vision-language tasks have been unified in MLLMs
  - Transformer + Next token prediction
- Fine-grained visual perception: detection, segmentation, depth estimation...
  - Task-specific modules



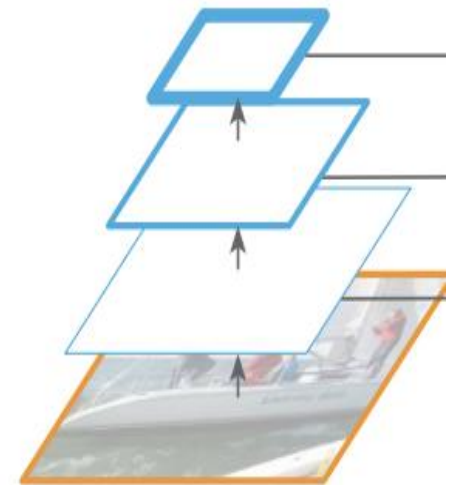
RPN



Mask Decoder



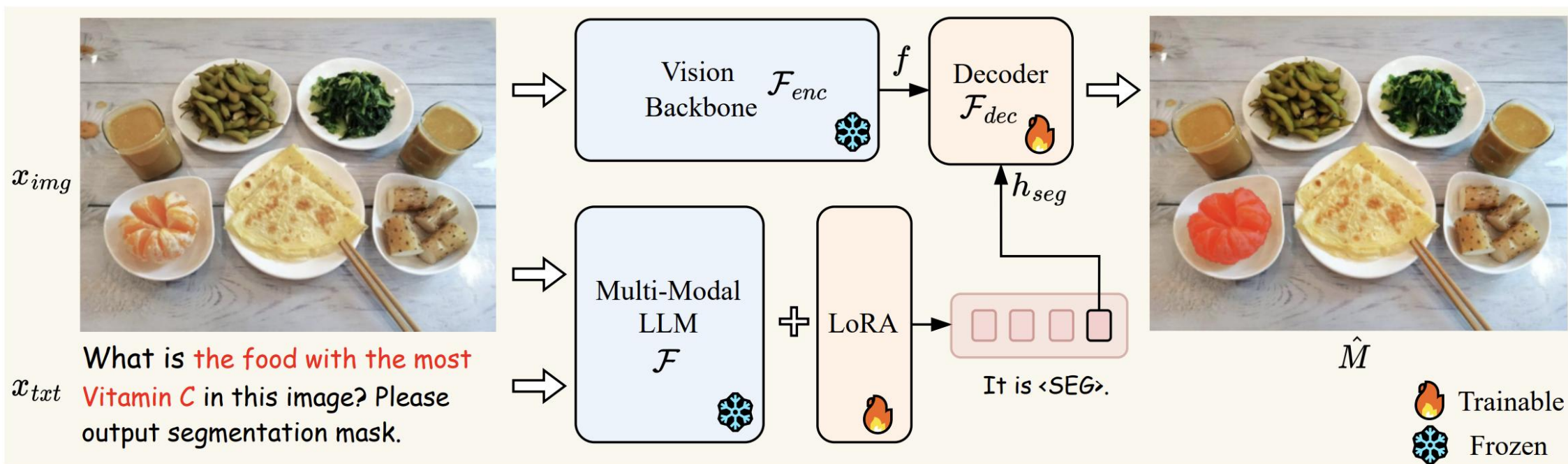
Object Queries



FPN

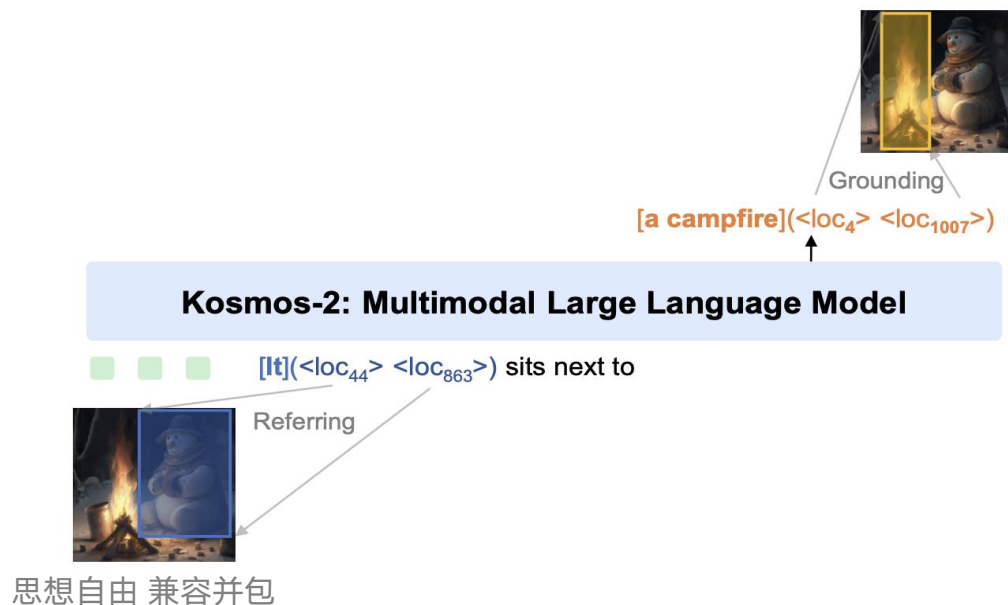
# Motivation

- How to extend MLLMs with fine-grained perception?
- MLLM + task decoder
  - SAM, Grounding DINO
- Complex architectures and training



# Motivation

- How to extend MLLMs with fine-grained perception?
- Conver box and mask to text
  - polygon for mask
- Hard to support multi-object detection
- Low mask performance by quantization errors

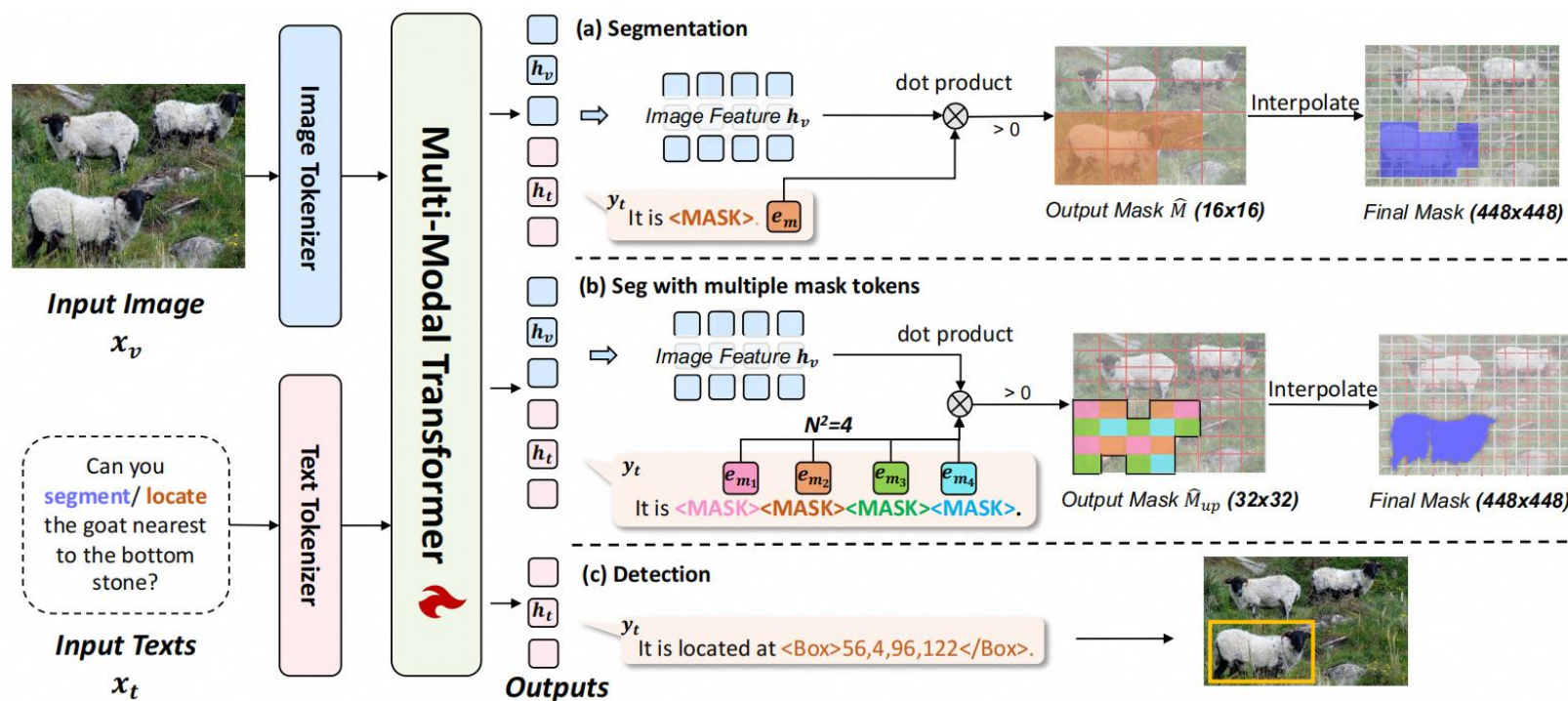


Methods	Instance Seg		
	AP	AP <sub>50</sub>	AP <sub>75</sub>
<b>Specialist Models</b>			
Faster R-CNN-FPN [73]	-	-	-
DETR-DC5 [13]	-	-	-
Deformable-DETR [106]	-	-	-
Pix2Seq [21]	-	-	-
Mask R-CNN [36]	37.1	58.4	40.1
Polar Mask [93]	30.5	52.0	31.1
Mask2Former [25]	43.7	-	-



- MLLM can answer “Where” and “What”
  - The category and mask information is in image features
- Modeling by similarity: **Embedding Retrieval**

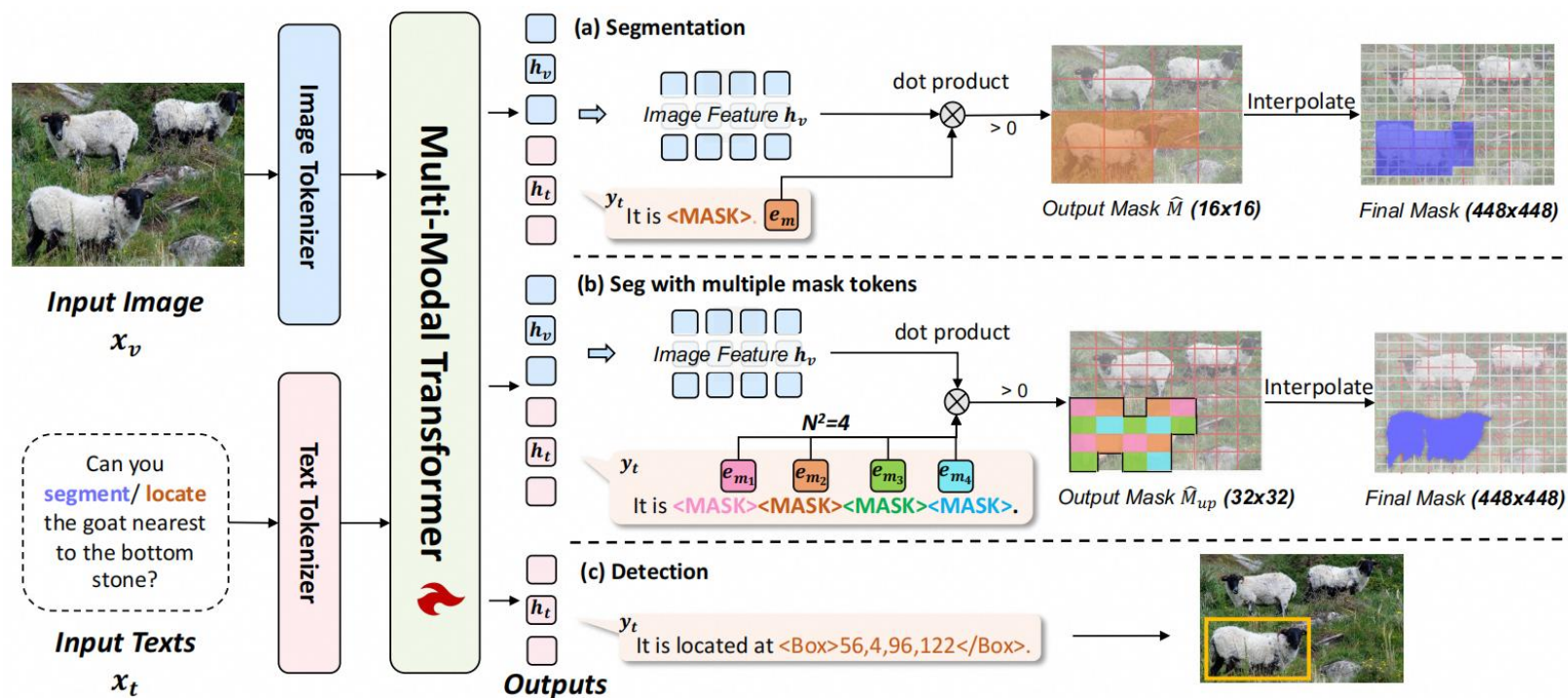
$$\mathbf{h}_v, \mathbf{y}_t, \mathbf{h}_t = \mathcal{F}(\mathbf{x}_v, \mathbf{x}_t). \quad s = \frac{\mathbf{e}_m \mathbf{h}_v^\top}{\sqrt{d}}, \quad \hat{\mathbf{M}} = \mathbb{I}(s > 0),$$



- Image features are typically downsampled
  - e.g. 28x in InternVL
- Upsampling by multiple mask tokens

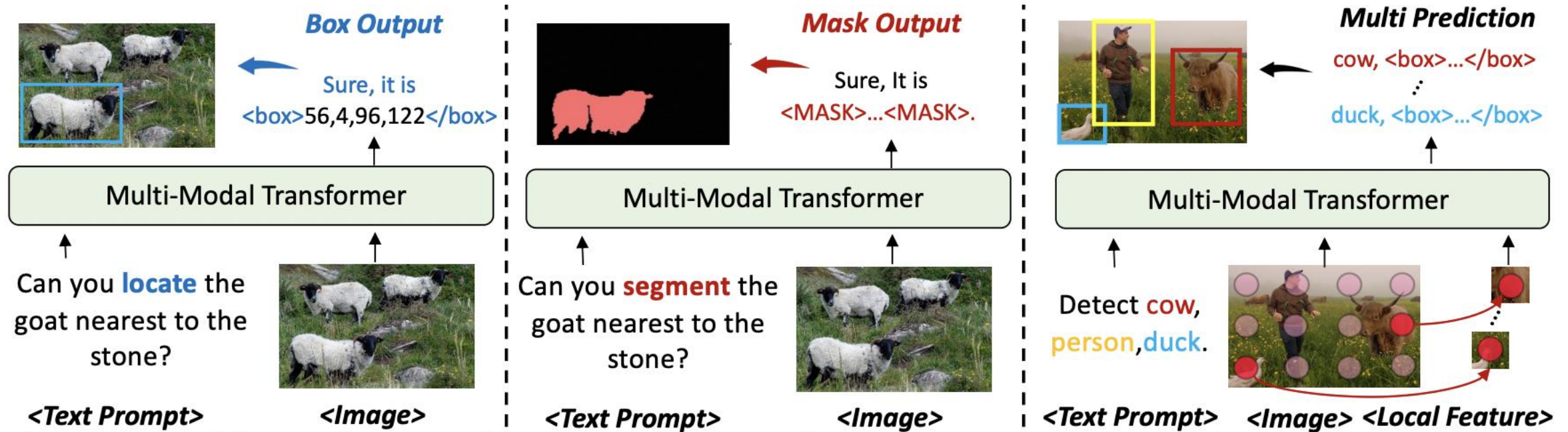
$$s_i = \frac{\mathbf{e}_{m_i} \mathbf{h}_v^\top}{\sqrt{d}}, \quad s_{\text{concat}} = \text{concat}(\{s_i\}_{i=1}^{N^2}), \quad s_{\text{concat}} \in \mathbb{R}^{N^2 \times H_p \times W_p},$$

$$s_{\text{up}} = \text{reshape}(s_{\text{concat}}), \quad s_{\text{up}} \in \mathbb{R}^{(H_p N) \times (W_p N)}.$$



# Method

- Multi prediction tasks: Object detection, instance segmentation
  - Long sequence, inefficient, difficult to learn
- Split multiple predictions to independent subtasks
- Parallel decoding



- Multi-task training
  - Five tasks: COCO Det/InsSeg/Caption, ADE20K, RefCOCO
  - Evaluation: Multi-task benchmark in GiT
  - Architecture: UFO-ViT, UFO-InternVL2.5-8B
- Instruction Tuning
  - 26 datasets with diverse tasks
  - Evaluation: RefCOCO (REC and RES), ReasonSeg
  - Architecture: UFO-InternVL2.5-8B, UFO-LLaVA-1.5-7B



# Experiments

- Multi-task benchmark

Methods	Specific Modules		#Params	Object Detection			Instance Seg			Semantic Seg	Captioning		REC
	Examples	Num		AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	mIoU(SS)	BLEU-4	CIDEr	Acc@0.5
<b>Specialist Models</b>													
Deformable-DETR [94]	RegressionHead	5	40M	45.4	64.7	49.0	-	-	-	-	-	-	-
Mask R-CNN [26]	FPN,RPN	6	46M	41.0	61.7	44.9	37.1	58.4	40.1	-	-	-	-
Polar Mask [79]	CenternessHead	5	55M	-	-	-	30.5	52.0	31.1	-	-	-	-
Mask2Former [13]	PixelDecoder	5	44M	-	-	-	43.7	-	-	47.2	-	-	-
VL-T5 [15]	Faster R-CNN	3	440M	-	-	-	-	-	-	-	34.5	116.5	-
MDETR [30]	RoBERTa,DETR	6	188M	-	-	-	-	-	-	-	-	-	86.8
<b>Generalist Models (MultiTask-Training)</b>													
Uni-Perceiver [95]	None	1	124M	-	-	-	-	-	-	-	32.0	★	★
Uni-Perceiver-MoE [93]	None	1	167M	-	-	-	-	-	-	-	33.2	★	★
VisionLLM-R50 [74]	Deform-DETR	6	7B	44.6	64.0	48.1	25.1	50.0	22.4	-	31.0	112.5	80.6
GiT-B <sub>single-task</sub> [69]	None	1	131M	45.1	62.7	49.1	31.4	54.8	31.2	47.7	33.7	107.9	83.3
GiT-B <sub>multi-task</sub> [69]	None	1	131M	46.7	64.2	50.7	31.9	56.4	31.4	47.8	35.4	112.6	85.8
GiT-L <sub>multi-task</sub> [69]	None	1	387M	51.3	69.2	55.9	35.1	61.4	34.7	50.6	35.7	116.0	88.4
GiT-H <sub>multi-task</sub> [69]	None	1	756M	52.9	71.0	57.8	35.8	62.6	35.6	52.4	36.2	118.2	89.2
UFO-ViT-B <sub>single-task</sub>	None	1	131M	47.8	65.7	52.0	42.6	65.8	46.1	49.5	34.2	111.1	83.6
UFO-ViT-B <sub>multi-task</sub>	None	1	131M	48.3	66.6	52.6	43.5	66.2	47.0	50.2	35.3	114.2	85.8
<b>Improvement</b> (single→multi)				<b>+0.5</b>	<b>+0.9</b>	<b>+0.6</b>	<b>+0.9</b>	<b>+0.4</b>	<b>+0.9</b>	<b>+0.7</b>	<b>+1.1</b>	<b>+3.1</b>	<b>+2.2</b>
UFO-ViT-L <sub>multi-task</sub>	None	1	387M	52.9	71.3	57.9	47.3	70.9	51.6	54.0	35.9	118.6	88.5
UFO-ViT-H <sub>multi-task</sub>	None	1	756M	<b>54.1</b>	<b>72.4</b>	<b>58.9</b>	<b>48.1</b>	<b>71.6</b>	<b>53.0</b>	<b>55.7</b>	37.6	123.6	89.2
UFO-InternVL2.5-8B <sub>multi-task</sub>	None	1	8B	52.3	71.7	56.5	45.8	69.5	49.7	54.6	<b>39.6</b>	<b>131.6</b>	90.4

# Experiments

- REC and RES in RefCOCO

Methods	Referring Expression Comprehension (REC)									Referring Expression Segmentation (RES)								
	RefCOCO			RefCOCO+			RefCOCOg		Avg	RefCOCO			RefCOCO+			RefCOCOg		Avg
	val	testA	testB	val	testA	testB	val	test		val	testA	testB	val	testA	testB	val	test	
<i>MLLMs with Task Decoders</i>																		
GLaMM-7B [57]	-	-	-	-	-	-	-	-	-	79.5	<b>83.2</b>	76.9	72.6	78.7	64.6	74.2	74.9	75.6
SAM4MLLM-8B [11]	-	-	-	-	-	-	-	-	-	79.8	82.7	74.7	74.6	80.0	67.2	75.5	76.4	76.4
HiMTok-8B [73]	-	-	-	-	-	-	-	-	-	<b>81.1</b>	81.2	<b>79.2</b>	<b>77.1</b>	78.8	71.5	75.8	76.7	77.7
PerceptionGPT-7B [53]	88.6	92.5	84.6	82.1	88.6	74.2	84.1	85.2	85.0	75.1	78.6	71.7	68.5	73.9	61.3	70.3	71.7	71.4
VisionLLM v2 [77]	90.0	93.1	87.1	81.1	87.3	74.5	85.0	86.4	85.6	79.2	82.3	77.0	68.9	75.8	61.8	73.3	74.8	74.1
<i>MLLMs w/o Task Decoders</i>																		
Shirka-7B [7]	87.0	90.6	80.2	81.6	87.4	72.1	82.3	82.2	82.9	-	-	-	-	-	-	-	-	-
MiniGPT-v2-7B [6]	88.1	91.3	84.3	79.6	85.5	73.3	84.2	84.3	83.8	-	-	-	-	-	-	-	-	-
Ferret-v2-7B [85]	92.8	94.7	88.7	87.4	<b>92.8</b>	79.3	<b>89.4</b>	<b>89.3</b>	89.3	-	-	-	-	-	-	-	-	-
VistaLLM-7B [54]	88.1	91.5	83.0	82.9	89.8	74.8	83.6	84.4	84.8	74.5	76.0	72.7	69.1	73.7	64.0	69.0	70.9	71.2
UFO-LLaVA-1.5-7B	90.2	93.5	87.3	84.4	90.3	78.7	86.4	86.8	87.2	77.2	80.1	76.4	71.8	77.9	70.2	74.1	73.5	75.2
UFO-LLaVA-1.5-7B*	91.1	93.7	88.6	85.5	90.5	79.9	87.3	87.2	88.0	77.9	81.1	77.0	72.5	78.5	71.4	75.6	74.1	76.0
UFO-InternVL2.5-8B	91.8	94.3	87.5	86.9	91.3	80.6	87.9	88.6	88.6	80.0	81.6	78.1	76.7	79.9	72.3	75.5	76.3	77.6
UFO-InternVL2.5-8B*	<b>93.1</b>	<b>94.8</b>	<b>89.2</b>	<b>87.7</b>	92.1	<b>82.3</b>	88.2	89.2	<b>89.6</b>	81.0	82.6	78.6	<b>77.1</b>	<b>80.4</b>	<b>72.6</b>	<b>76.7</b>	<b>77.3</b>	<b>78.3</b>

# Experiments

- Best performance on ReasonSeg
  - Outperforms by **6.2** gloU

Methods	overall	ReasonSeg	
		short query	long query
X-Decoder [91]	21.7	20.4	22.2
SEEM [92]	24.3	20.1	25.6
LISA-7B [34]	36.8	37.6	36.6
LISA-7B [34]*	47.3	40.6	49.4
Cores-7B [2]	48.7	41.0	50.9
Cores-7B [2]*	52.4	44.2	55.0
HiMTok-8B [68]*	60.8	-	-
UFO-LLaVA-1.5-7B	54.4	41.2	58.5
UFO-LLaVA-1.5-7B*	58.8	46.5	62.7
UFO-InternVL2.5-8B	60.0	48.7	63.6
UFO-InternVL2.5-8B*	<b>67.0</b>	<b>56.2</b>	<b>70.4</b>



Please segment the the objects that can protect the snail and prevent it from getting injured.



Sure, <MASK>...<MASK>.



Please segment the food containing proteins, carbohydrates and other nutrients.



It is <MASK>...<MASK>.

- Extend to depth estimation and surface normal prediction

Methods	RMSE↓	$\delta 1 \uparrow$	REL↓	log10↓
Painter [70]	0.327	0.930	0.090	-
Unified-IO 2 [44]	0.423	-	-	-
UFO-InternVL2.5-8B	0.305	0.936	0.087	0.035

Method	Mean	Median	11.25°	22.5°	30°
GeoNet++ [55]	18.5	11.2	9.502	0.732	0.907
Marigold [32]	18.8	-	0.559	-	-
GeoWizard [22]	17.0	-	0.565	-	-
UFO-InternVL2.5-8B	17.8	10.4	0.543	0.733	0.800



# Conclusion

---

- UFO reformulate segmentation as embedding retrieval
  - Remove the need for task decoders
  - Fully aligned with open-ended Language interface
- UFO explore the image representation capabilities of MLLMs
  - A general way to extract information in image features
- UFO unifies both single-prediction and mutli-prediction tasks
  - Parallel decoding improves efficiency and performance