

InstructSAM: A Training-Free Framework for Instruction-Oriented Remote Sensing Object Recognition

Yijie Zheng^{1,2} Weijie Wu^{1,2} Qingyun Li³ Xuehui Wang⁴ Xu Zhou⁵
Aiai Ren⁵ Jun Shen⁵ Long Zhao¹ Guoqing Li¹ Xue Yang⁴

¹Aerospace Information Research Institute ²University of Chinese Academy of Sciences
³Harbin Institute of Technology ⁴Shanghai Jiao Tong University ⁵University of Wollongong

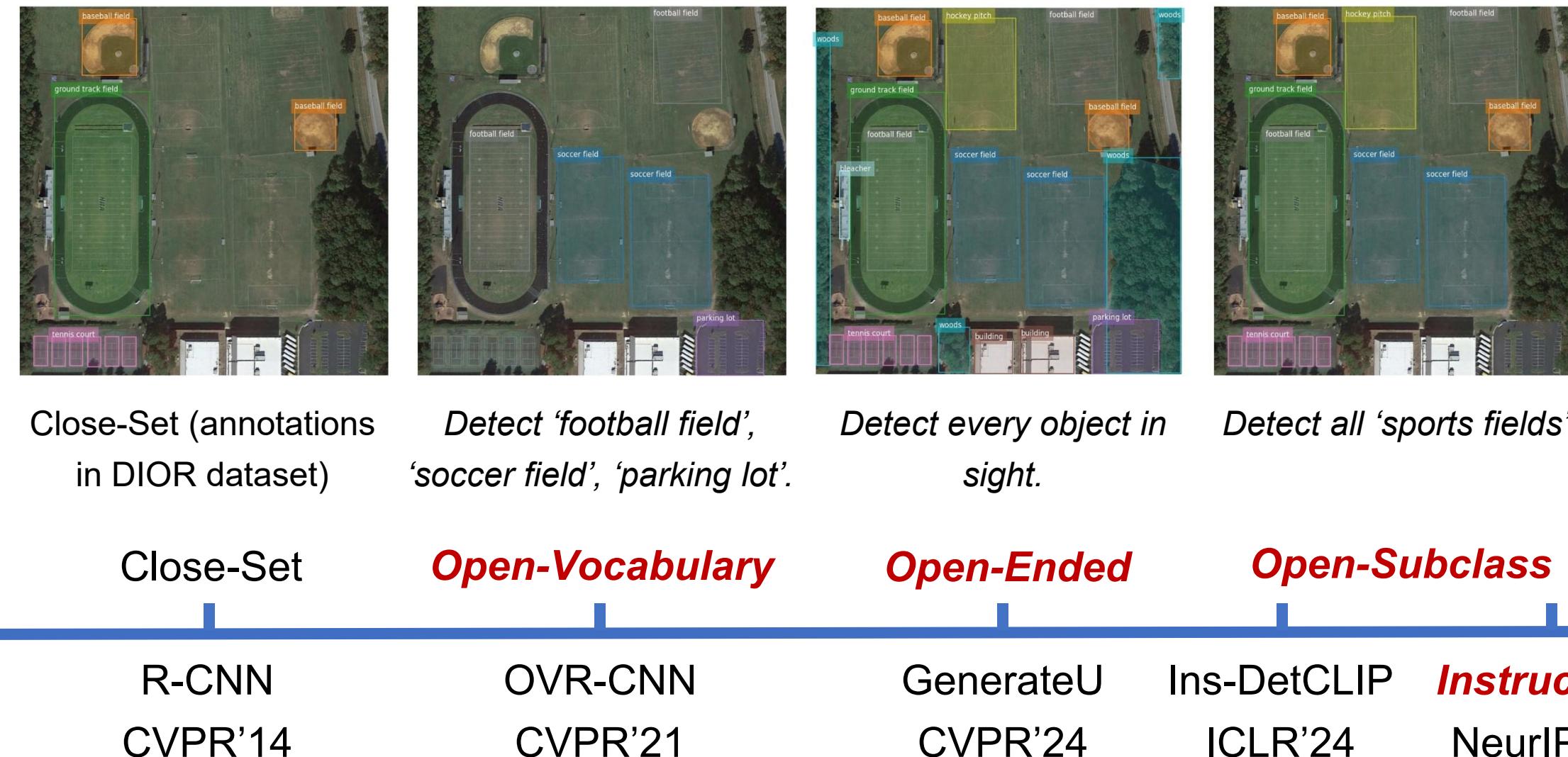


GAIA 2025

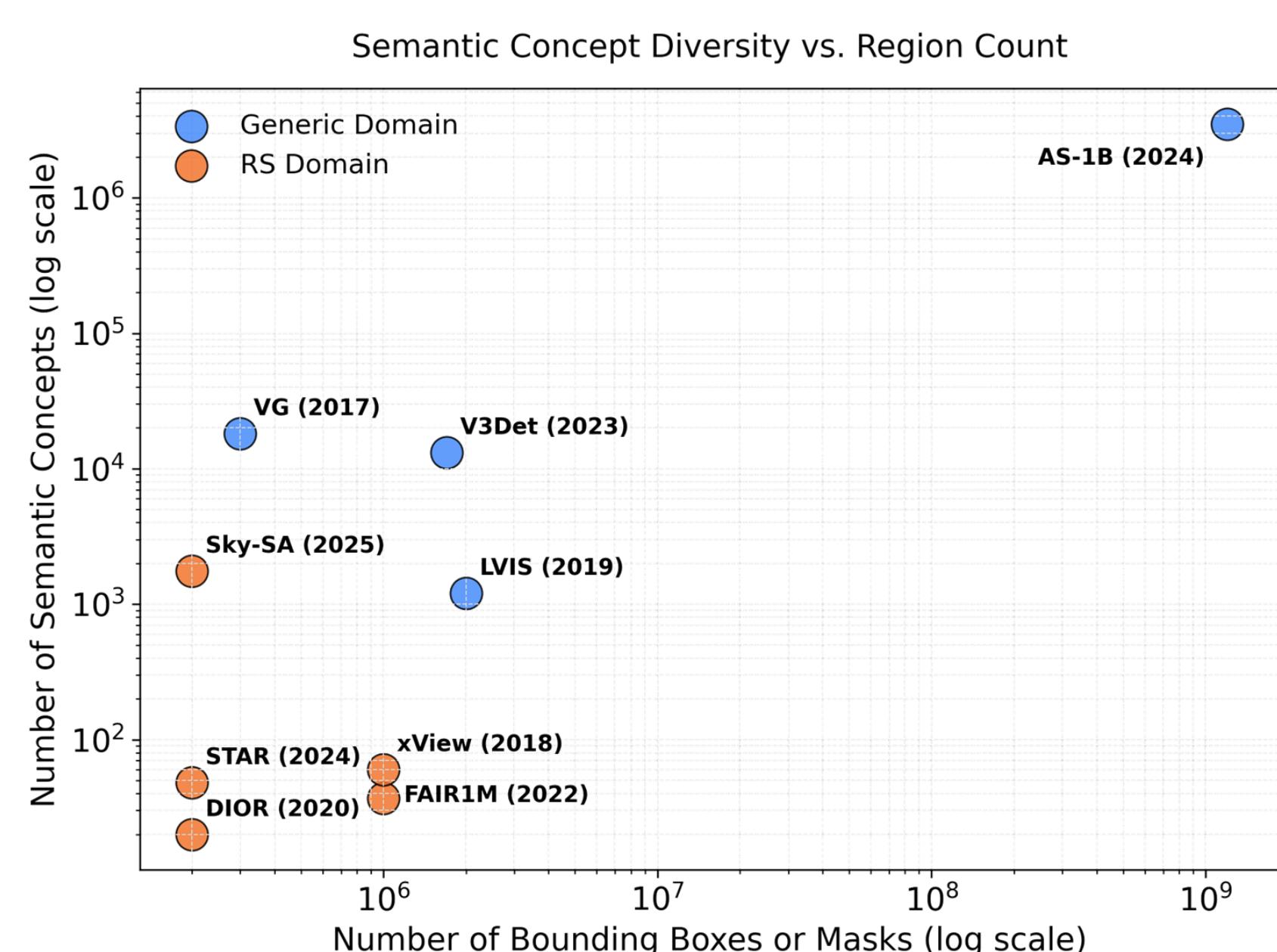


From close-set to instruction-oriented

- Fixed categories limit traditional detection.
- Real tasks need flexible, dynamic recognition.
- 💡 **Instruction-oriented: just give an instruction, model adapts.**



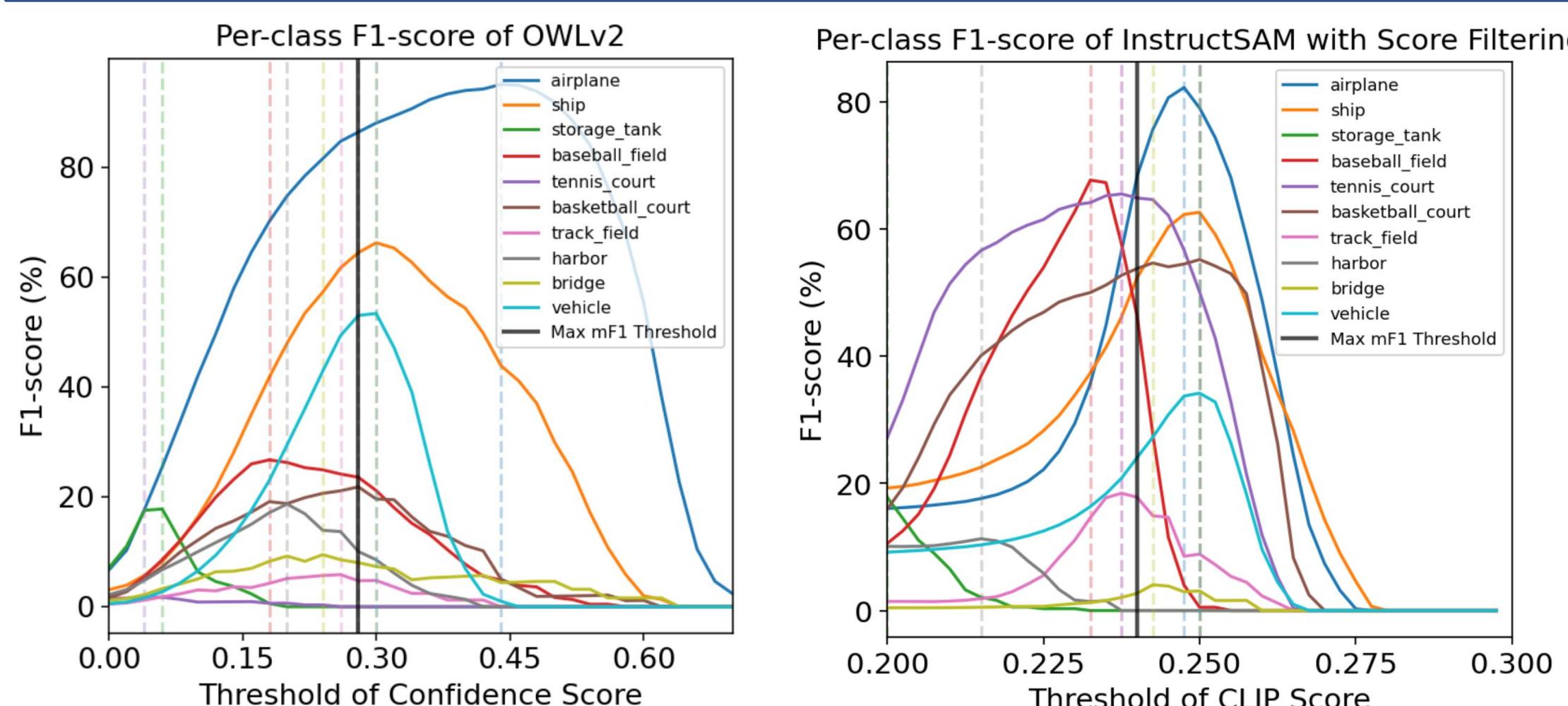
Motivation 1: Lacking diverse training datasets



- RS datasets lack semantic diversity → poor generalization.

💡 **Use generic LVMs to identify objects.**

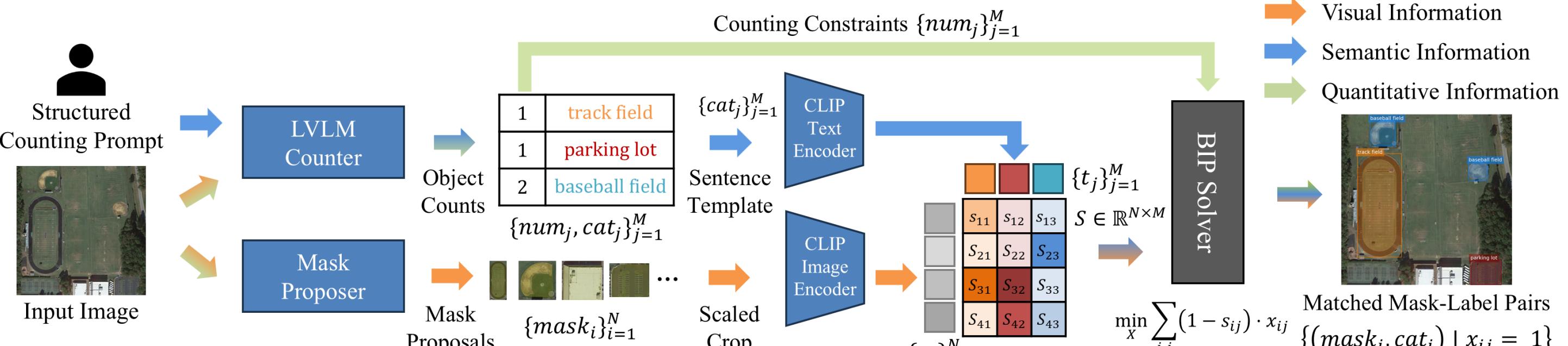
Motivation 2: Unreliable score filtering



- Score-based filtering is crucial to filter low-quality pseudo labels.
- Optimal thresholds vary across classes → no universal solution.
- Over-reliance on score thresholds leads to misclassifications.
- 💡 **Introduce counting constraints**

InstructSAM framework

- Decompose object segmentation into three easier steps.
- 💡 LVM for categories & counts, SAM for mask proposals, CLIP for similarity, PuLP for optimization.



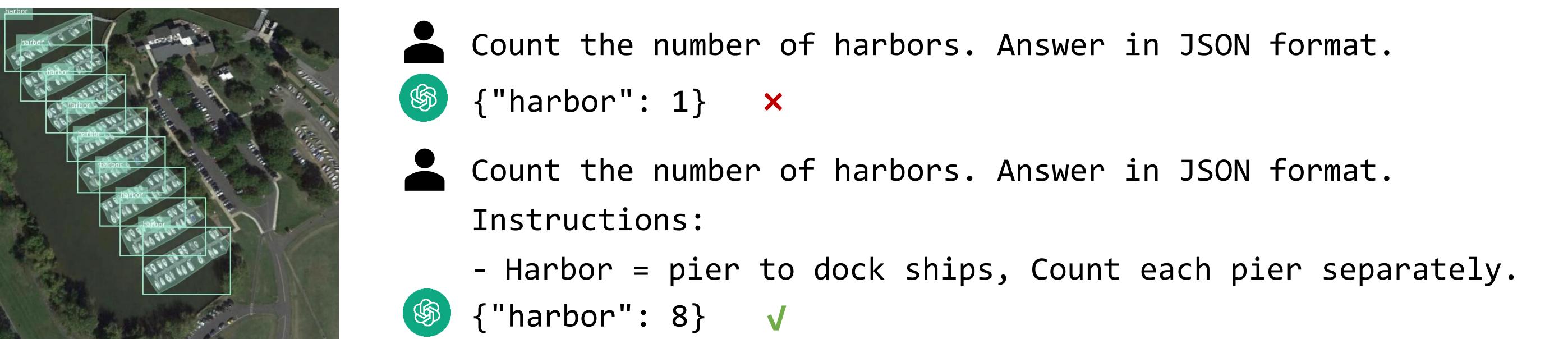
- Reframe segmentation as a mask-label matching problem

$$\begin{aligned} \min_x \quad & \sum_{i=1}^N \sum_{j=1}^M (1 - s_{ij}) \cdot x_{ij} \\ \text{s.t.} \quad & \sum_{j=1}^M x_{ij} \leq 1, \\ & \sum_{i=1}^N x_{ij} = num_j, \end{aligned}$$

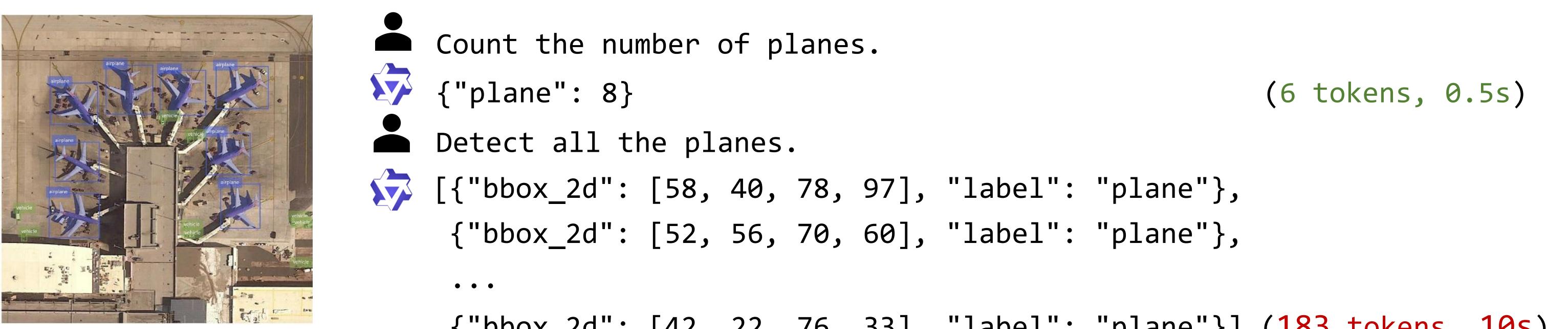
- minimize mismatches ($1 - \text{similarity}$)
- One mask → one category
- Total masks per category = LVM count

LVLM counts precisely and quickly

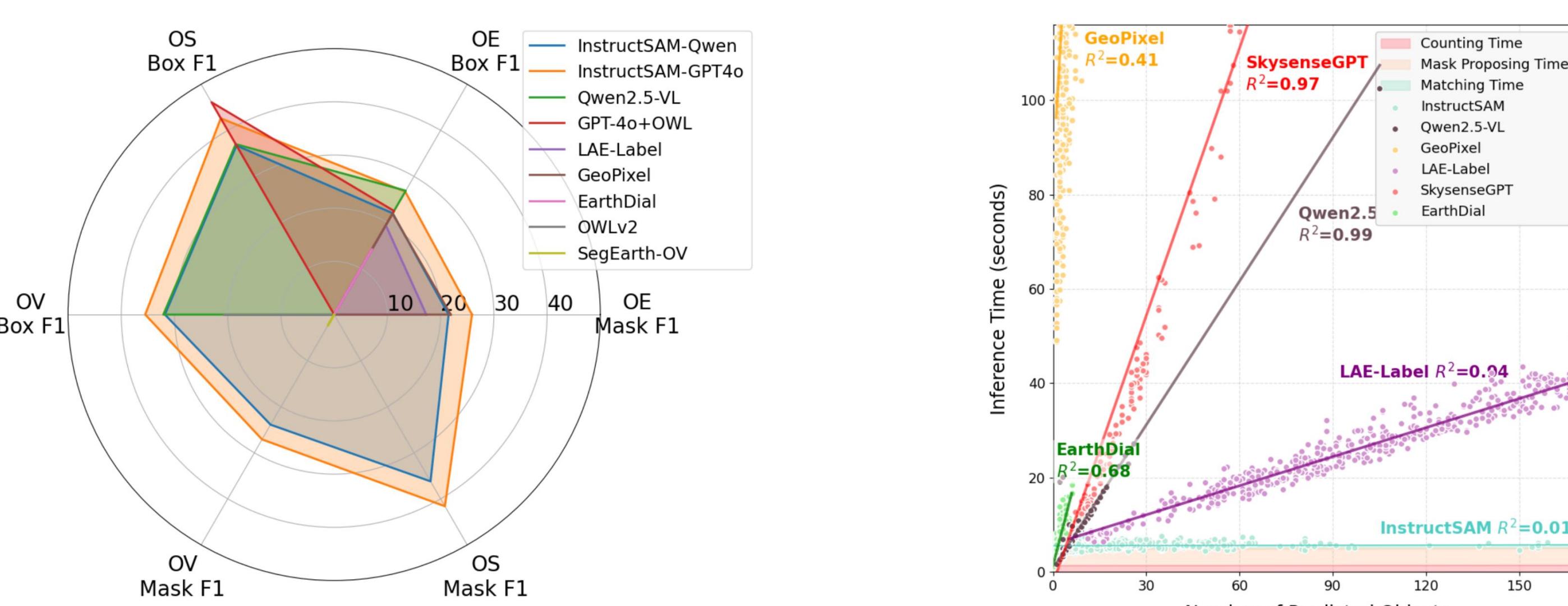
- With clear rules, GPT-4o counts as accurately as Faster R-CNN (80% vs. 81% mF1).



- Counting greatly reduces tokens and inference time.



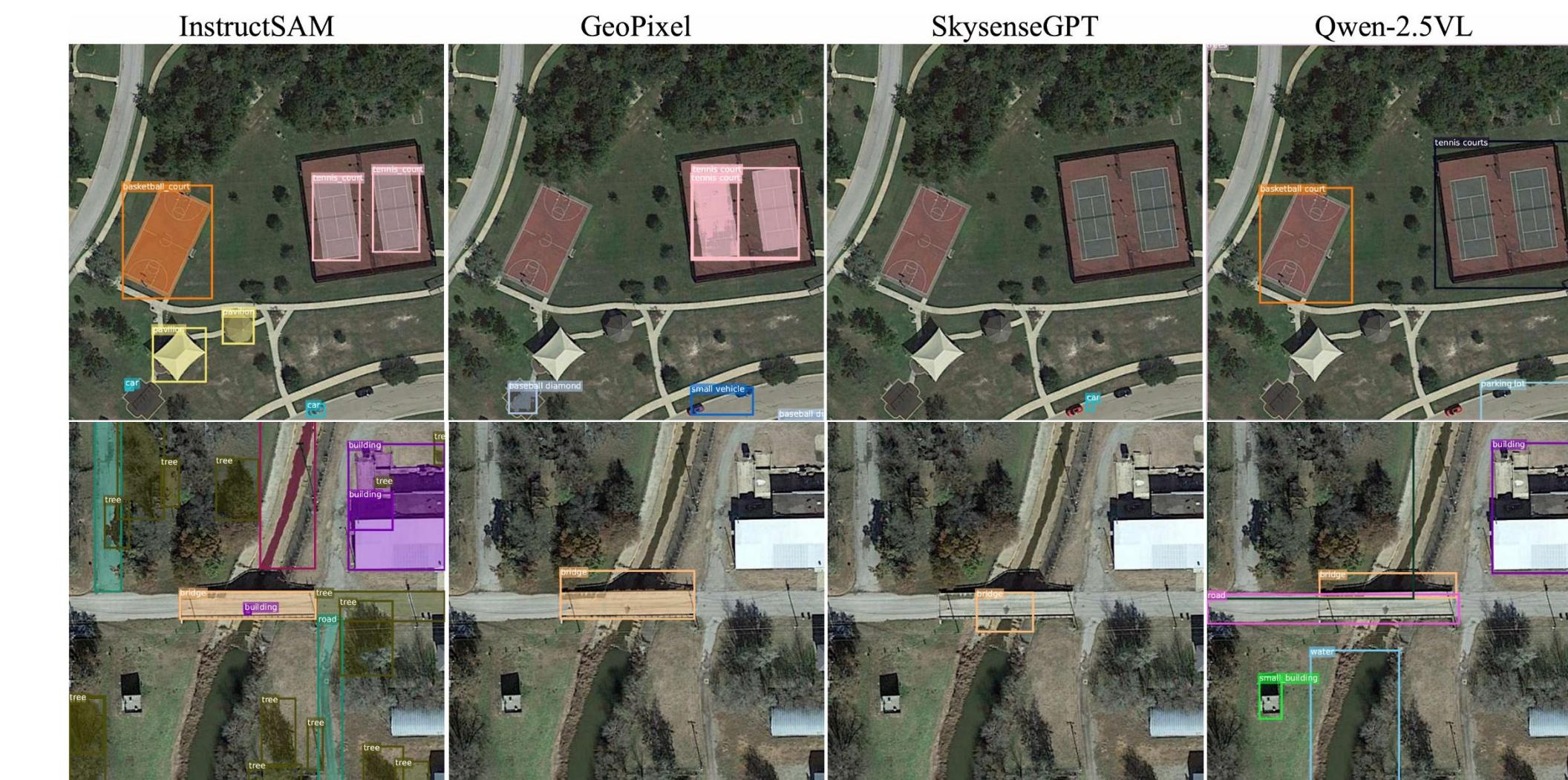
Zero-Shot Results across Three Settings



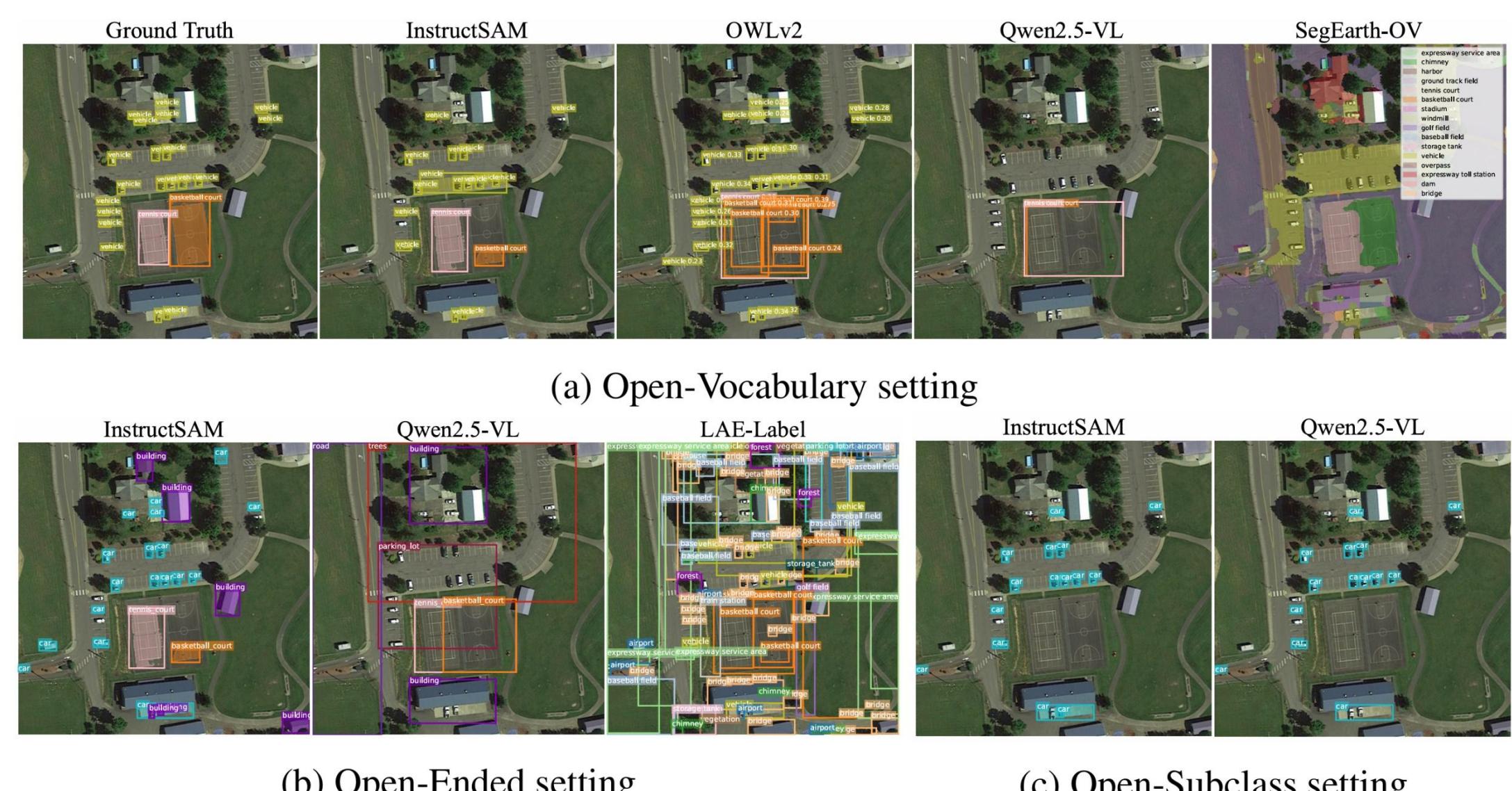
- Evaluated on NWPU and DIOR dataset
- 💡 Strong performance across most settings
- 💡 Performs well using open models (Qwen2.5-VL)
- Inference time in open-ended setting
- 💡 89% fewer tokens, 32% inference time reduction
- 💡 Inference time stays constant with more objects

Qualitative results: open-ended setting

- ✗ Existing RSVLMs fail to generalize beyond training categories.
- ✓ InstructSAM recognizes more categories, e.g., **pavilion** and **tree**.



Qualitative results across three settings



Generalization beyond RS

- InstructSAM also generalizes to natural images.

Instruction:

Detect all electronic components.

Instruction:

Detect dices whose letters come before K.

InstructSAM

Qwen2.5-VL

InstructSAM

Qwen2.5-VL

InstructSAM

SegEarth-OV

InstructSAM

Qwen2.5-VL

InstructSAM

Qwen2