

UNIVERSITY OF CALIFORNIA

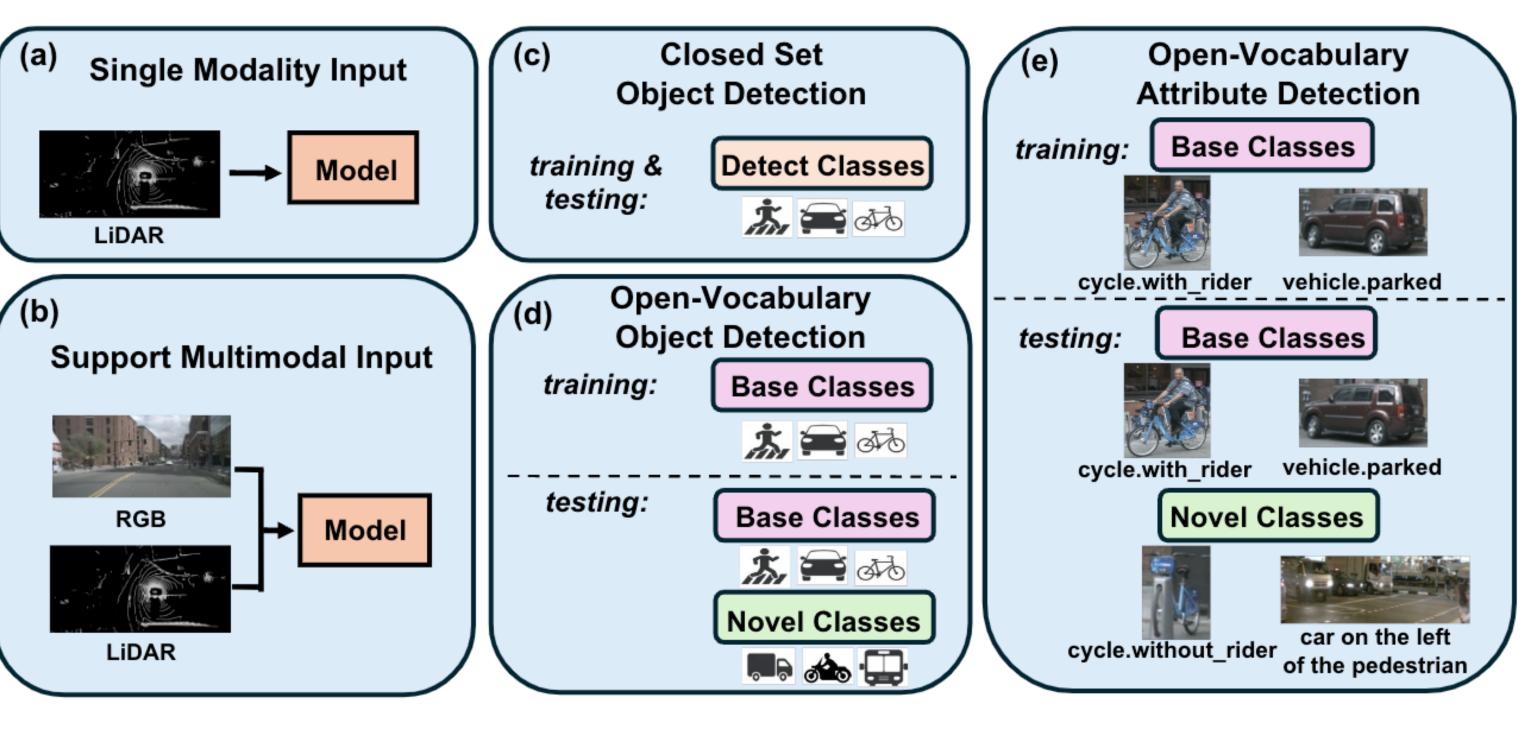
Towards Open-Vocabulary Multimodal 3D **Object Detection with Attributes**



Xinhao Xiang¹, Kuan-Chuan Peng², Suhas Lohit², Michael J. Jones², Jiawei Zhang¹



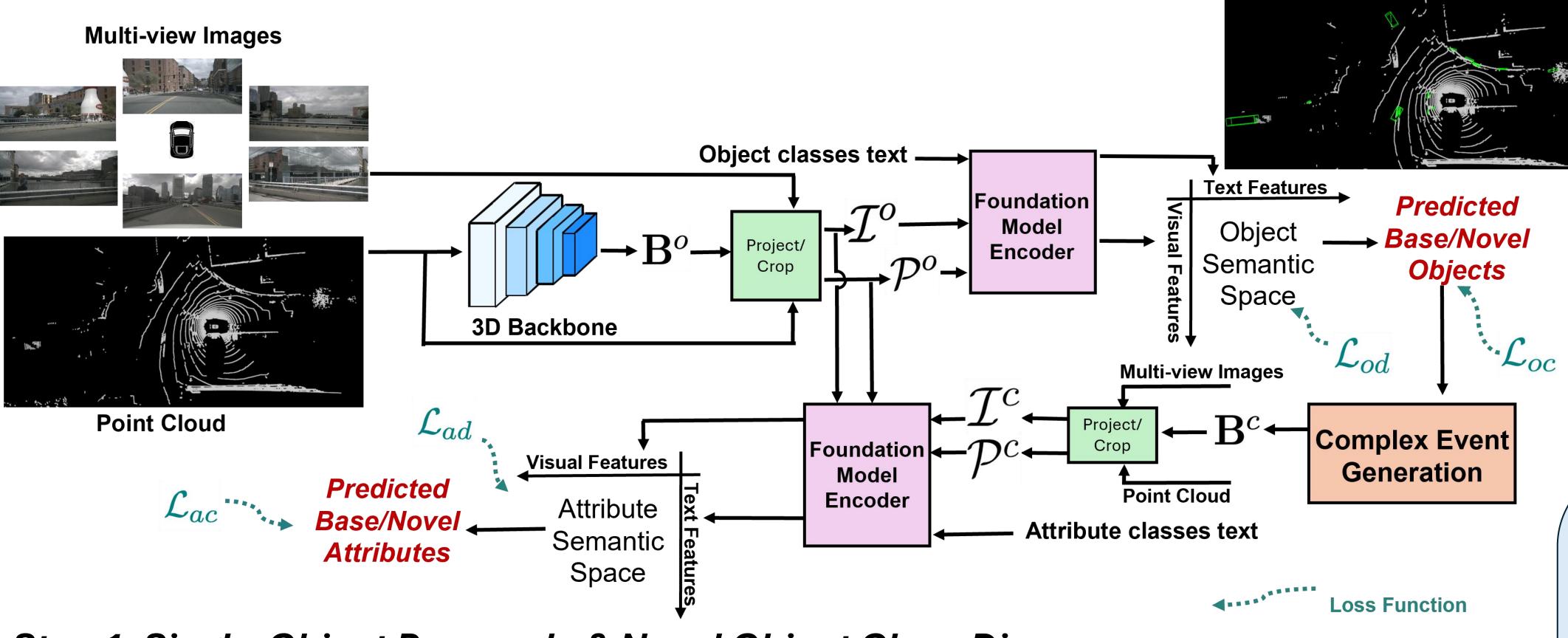




Contribution Summary:

Object/Attribute Detection	Method Categories										
Conditions	$\overline{\mathcal{C}_0}$	\mathcal{C}_1	\mathcal{C}_2	\mathcal{C}_3	\mathcal{C}_4	C_5	\mathcal{C}_6	OVODA (ours)			
support 3D object detection	\checkmark	X		$\sqrt{}$	\checkmark		X				
can detect attributes	X	X	X	X	X	\checkmark	X				
support OV object detection	X	\checkmark	\checkmark	\checkmark	\checkmark	X	X				
support multi-modal input	X	X	\checkmark	X	\checkmark	X	X				
need no novel class anchor size	X	\checkmark	X	\checkmark	\checkmark	\checkmark	X				
can detect complex events	X	X	X	X	X	X	\checkmark				

The OVODA Framework:



OVAD: The Attribute Benchmark

vocabulary set for attribute detection dataset with rider, without rider, moving, standing, **OVAD** sitting lying down, parked, moving, stopped, in front of, behind, on the left of, on the right of

- Contain 84,384 annotations across 10 classes in total
- To create spatial attribute annotations:

Experimental Setup

 $\mathbf{B}^{\text{OVAD}} = \left\{ B_{ij}^{\text{OVAD}} = (\text{Comb}(B_i, B_j) | \text{Dist}(B_i, B_j) \leq 15 \text{m}) \right\}$

Step 1: Single-Object Proposals & Novel Object Class Discovery

Novel object class discovery: $\mathbf{O}^{disc} = \left\{ B^o_j | \forall B^b_i \in \mathbf{B}^b, \text{IoU}_{3D} \left(B^o_j, B^b_i \right) < \theta^b, Q^o_j > \theta^o, p^O_{j,c^*_j} > \theta^s, B^o_j \in \mathbf{B}^o, c^*_j \notin \mathbb{C}^b \right\}$ with Concatenating Foundation Model features (CFM) + Prompt Tuning (PT)

Step 2: Complex Event Generation (CEG) for Attributes

Complex event visual proposal generation: $\mathbf{B}^s = \left\{ B^s_{ij} = \left(\operatorname{Comb}(B^o_i, B^o_j) \mid \operatorname{Dist}(B^o_i, B^o_j) \leq \theta^d \right) \right\}$ with horizontal flip augmentation (HFA)

Complex event text proposal generation: c_{ij}^s ="From the perspective of $\mathbb{T}(c_i^o)$, $\mathbb{T}(c_i^o)$ $\mathbb{T}(SA)$ $\mathbb{T}(c_i^o)$." with perspective specified prompt (PSP)

Step 3: Novel Attribute Class Discovery

Novel attribute class discovery: $\mathbf{A}^{disc} = \left\{ B_j^c | \forall B_i^{ba} \in \mathbf{B}^{ba}, \text{IoU}_{3D} \left(B_j^c, B_i^{ba} \right) < \theta^b, p_{j,e^*}^A > \theta^a, B_j^c \in \mathbf{B}^c, e^* \notin \mathbb{C}^{ba} \right\}$

			-					
	For object detection:	dataset setting	base object class	novel object class				
		N_{b6n4}	Car, Construction vehicles, Trailer, Barrier, Bicycle, Pedestrian	Truck, Bus, Motorcycle, Traffic cone				
	nuScenes	N_{b3n7}	Car, Bicycle, Pedestrian	Construction vehicles, Trailer, Barruck, Bus, Motorcycle, Traffic				
}:		N_{b0n10}	Ø	Car, Construction vehicles, Traile Barrier, Bicycle, Pedestrian Truck Bus, Motorcycle, Traffic cone	k,			
	Argoverse 2	A_{b4n4}	Regular Vehicle, Trailer, Bicycle, Pedestrian	Truck, Bus, Motorcycle, Construction cone	;			
	For attribute							
	detection:	dataset setting	base attribute class	novel attribute class				
	OVAD	OVAD	with rider, sitting lying down, parked, in front of, behind	, without rider, standing, moving, on the left of				

Results: Complex Event (Attribute) Detection

method	CFM	prompt tuning	mAP	AP _{N20}
OVODA	X	X	17.24	8.23
OVODA	\checkmark		27.43	12.34

Training Objectives

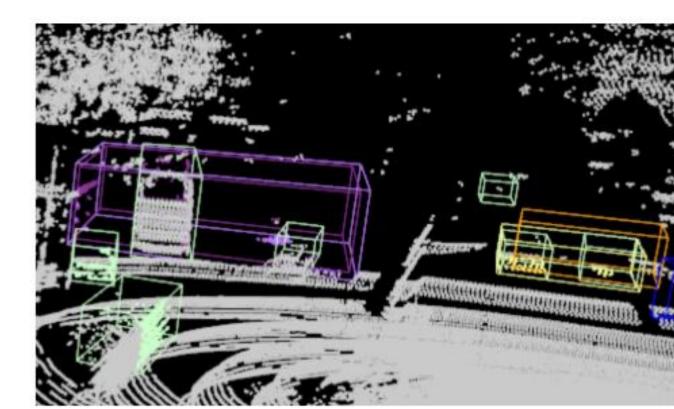
For object losses:

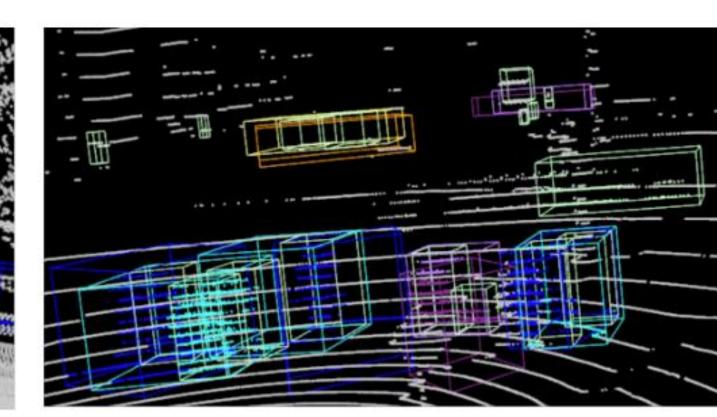
- $\mathcal{L}_{od} = \sum_{j=1}^{N} ||V_j^O \mathbf{F}_{det}||_1$ • 3D-2D feature distance align:
- 3D-Text classification contrastive: $\mathcal{L}_{oc} = \sum_{j=1}^{N} f\left(B_{j}^{disc}, \mathbf{B}^{b}\right) \cdot CE\left(\mathbf{P}_{j}^{disc}, h_{j}^{O}\right)$

For attribute losses:

 $\mathcal{L}_{ad} = \sum_{i=1}^{N} ||V_i^A - \mathbf{F}_c||_1$

• 3D-2D feature distance align: • 3D-Text classification contrastive: $\mathcal{L}_{ac} = \sum_{j=1}^{N} g\left(A_j^{disc}, \mathbf{B}^{ba}\right) \cdot CE\left(\mathbf{P}_j^{disc,a}, h_j^A\right)$





	Prediction:	Ground Truth:
		all single object
-		car-car
. —		pedestrian-pedestrian
-		other complex events

Ablations

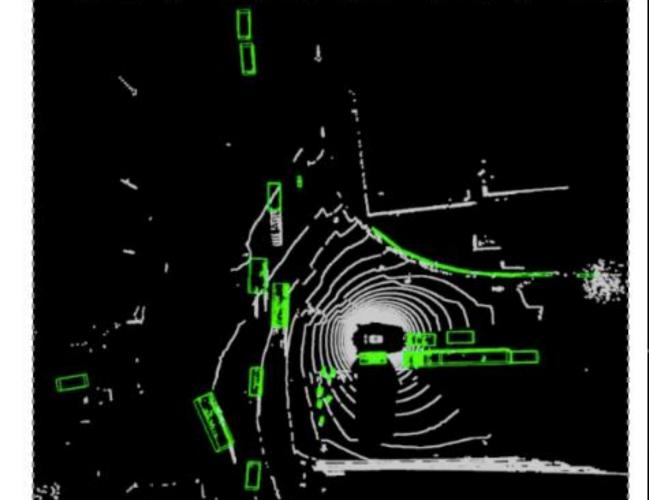
s:	method	CFM	PT	mAP	NDS	AP _{N20}	SR (%) (AD only)	SR (%) (AD & OD)
	AUONA	-	_				16.35 22.74	4.23 5.56
							25.90	6.77
nethod	foundation	m/	ΑP	NDS	AP _{N20}	SR (%) SR (%	%) (for both

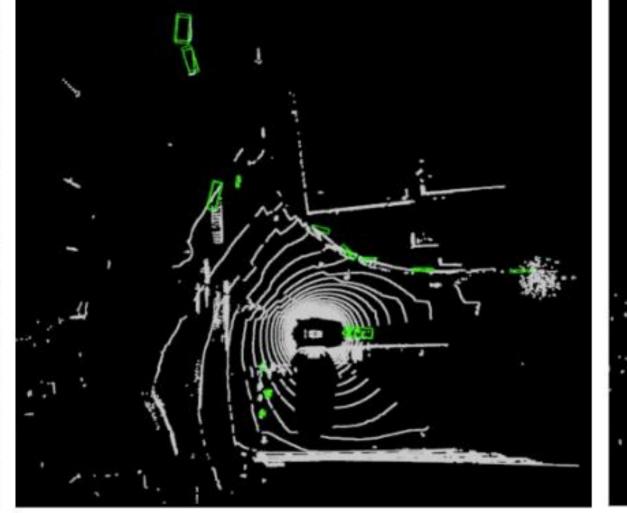
method	foundation model	mAP	NDS	AP _{N20}	SR (%) (AD only)	SR (%) (for bo AD & OD)
	CLIP [33]				16.20	3.22
OVODA	CogVLM [40] OneLLM [14]				19.84 25.90	5.39 6.77
					GD (64)	

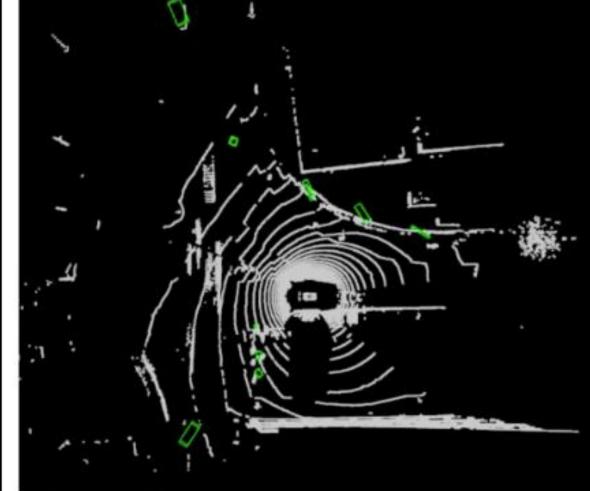
method	PSP	HFA	FM	mAP	NDS	AP _{N20}	SR (%) (AD only)	SR (%) (for both AD & OD)
	X	Х	CLIP	21.93	21.54	11.45	16.20	3.22
	\checkmark	X	CLIP	22.34	24.35	11.94	16.83	3.49
OVODA	X	\checkmark	CLIP	22.46	23.43	12.32	17.43	4.03
	\checkmark	\checkmark	CLIP	24.37	25.83	13.02	19.42	4.92
	\checkmark		OneLLM	32.25	31.85	14.72	25.90	6.77

Results: Open-Vocabulary 3D Object Detection

method	promp		prompt need no predefined anchor		N_{b6n4}			N_{b3n7}			N_{b0n10}		
	CFM	tuning	size for each novel class?	mAP	NDS	AP_{N20}	mAP	NDS	AP_{N20}	mAP	NDS	AP _{N20}	
Find n' Propagate [8]	X	X	X	44.95	47.87	33.65	37.38	40.28	18.46	N/A	N/A	16.72	
CoDAv2 [5]	Х	Х	✓	27.35	29.48	12.63	18.73	20.14	8.74	4.32	5.82	1.37	
OVODA	X	X	√	30.24	31.54	14.24	20.46	21.83	9.31	4.57	6.96	2.16	
OVODA	\checkmark	\checkmark	\checkmark	32.25	31.85	14.72	21.03	22.14	10.23	4.70	7.05	2.39	







Ground Truth

OVODA's Prediction

CoDAv2's Prediction