

OpenDriveLab 浦驾



上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory

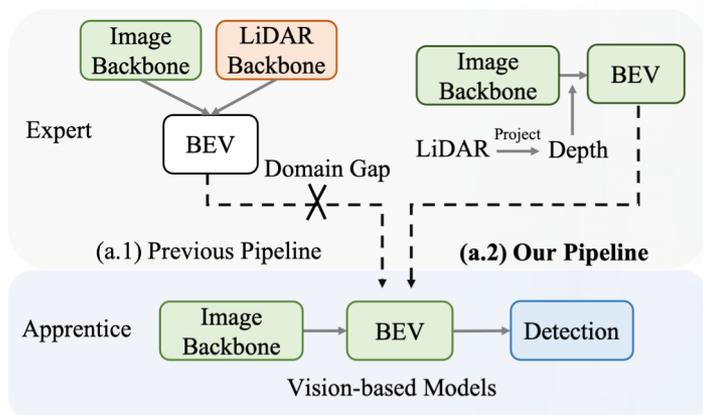
# Leveraging Vision-Centric objects for 3D Object Detection

*Linyan Huang, Zhiqi Li, Chonghao Sima, Wenhai Wang, Jingdong Wang, Yu Qiao, Hongyang Li*

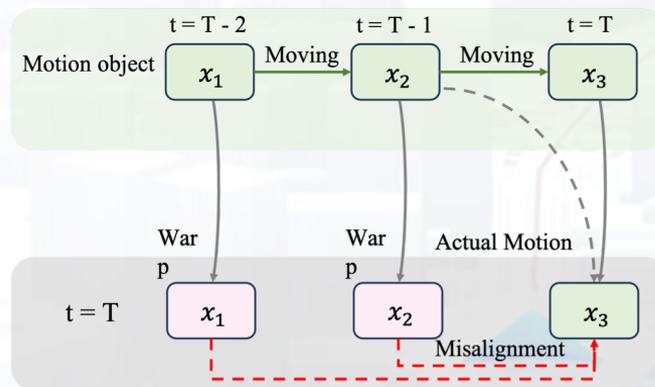
*Shanghai AI Lab, Nanjing University, CUHK, Baidu*

*NeurIPS 2023*

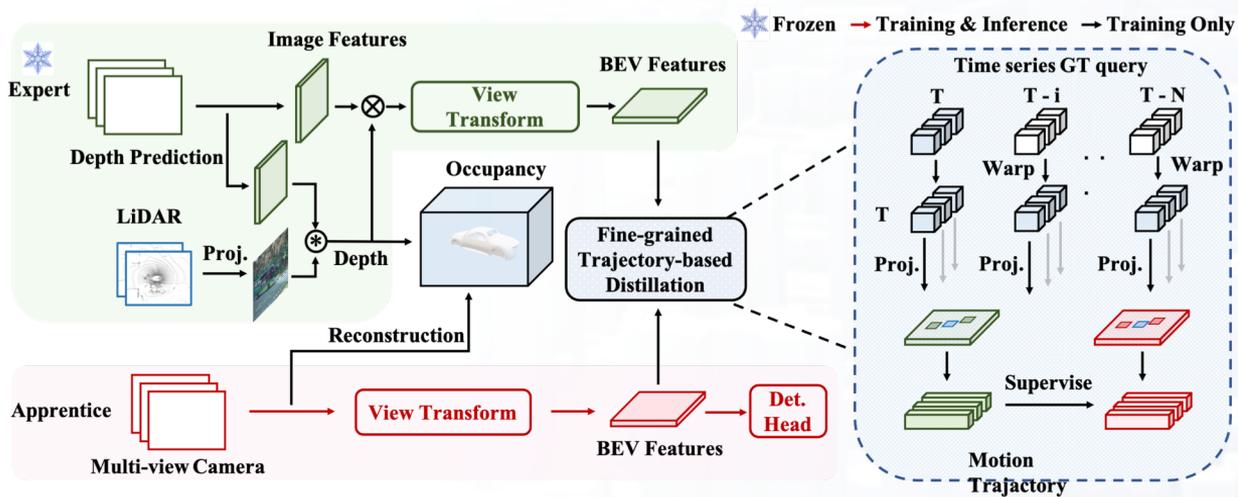
## Comparison



## Motion Misalignment



## Overview



- Better Expert
- Occupancy Reconstruction
- Motion Trajectory Module

## Experiments

Methods	Backbone	Image Size	Frames	mAP $\uparrow$	NDS $\uparrow$	mATE $\downarrow$	mASE $\downarrow$	mAOE $\downarrow$	mAVE $\downarrow$	mAAE $\downarrow$
BEVDet [23]	ResNet-50	256 $\times$ 704	1	0.298	0.379	0.725	0.279	0.589	0.860	0.245
PETR [39]	ResNet-50	384 $\times$ 1056	1	0.313	0.381	0.768	0.278	0.564	0.923	0.225
BEVDet4D [22]	ResNet-50	256 $\times$ 704	2	0.322	0.457	0.703	0.278	0.495	0.354	0.206
BEVDepth [35]	ResNet-50	256 $\times$ 704	2	0.351	0.475	0.639	0.267	0.479	0.428	0.198
BEVStereo [34]	ResNet-50	256 $\times$ 704	2	0.372	0.500	0.598	0.270	0.438	0.367	0.190
STS [54]	ResNet-50	256 $\times$ 704	2	0.377	0.489	0.601	0.275	0.450	0.446	0.212
VideoBEV [19]	ResNet-50	256 $\times$ 704	8	0.422	0.535	0.564	0.276	0.440	0.286	0.198
SOLOFusion [43]	ResNet-50	256 $\times$ 704	16+1	0.427	0.534	0.567	0.274	0.411	<b>0.252</b>	<b>0.188</b>
StreamPETR [51]	ResNet-50	256 $\times$ 704	8	0.432	0.540	0.581	0.272	0.413	0.295	0.195
Baseline*	ResNet-50	256 $\times$ 704	8+1	0.401	0.515	0.595	0.279	0.489	0.291	0.198
VCD-A	ResNet-50	256 $\times$ 704	8+1	0.426	0.540	0.547	0.271	0.433	0.268	0.207
Baseline* $\dagger$	ResNet-50	256 $\times$ 704	8+1	0.418	0.542	0.522	0.267	0.428	0.262	<b>0.188</b>
VCD-A $\dagger$	ResNet-50	256 $\times$ 704	8+1	<b>0.446</b>	<b>0.566</b>	<b>0.497</b>	<b>0.260</b>	<b>0.350</b>	0.257	0.203

Methods	Backbone	Image Size	mAP $\uparrow$	NDS $\uparrow$	mATE $\downarrow$	mASE $\downarrow$	mAOE $\downarrow$	mAVE $\downarrow$	mAAE $\downarrow$
FCOS3D $\dagger$ [52]	R101-DCN	900 $\times$ 1600	0.358	0.428	0.690	0.249	0.452	1.434	0.124
DETR3D $\dagger$ [53]	V2-99	900 $\times$ 1600	0.412	0.479	0.641	0.255	0.394	0.845	0.133
UVTR [33]	V2-99	900 $\times$ 1600	0.472	0.551	0.577	0.253	0.391	0.508	0.123
BEVDet4D $\dagger$ [22]	Swin-B [41]	900 $\times$ 1600	0.451	0.569	0.511	<b>0.241</b>	0.386	0.301	0.121
BEVFormer [36]	V2-99	900 $\times$ 1600	0.481	0.569	0.582	0.256	0.375	0.378	0.126
PolarFormer [28]	V2-99	900 $\times$ 1600	0.493	0.572	0.556	0.256	0.364	0.439	0.127
BEVDistill [11]	ConvNeXt-B	900 $\times$ 1600	0.496	0.594	0.475	0.249	0.378	0.313	0.125
PETRv2 [40]	RevCol [4]	640 $\times$ 1600	0.512	0.592	0.547	0.242	0.360	0.367	0.126
BEVDepth [35]	ConvNeXt-B	640 $\times$ 1600	0.520	0.609	0.445	0.243	0.352	0.347	0.127
AeDet $\dagger$ [15]	ConvNeXt-B	640 $\times$ 1600	0.531	0.620	0.439	0.247	0.344	0.292	0.130
SOLOFusion [43]	ConvNeXt-B	640 $\times$ 1600	0.540	0.619	0.453	0.257	0.376	0.276	0.148
StreamPETR [51]	ConvNeXt-B	640 $\times$ 1600	<b>0.550</b>	<b>0.631</b>	0.493	<b>0.241</b>	<b>0.343</b>	<b>0.243</b>	0.123
Baseline*	ConvNeXt-B	640 $\times$ 1600	0.522	0.610	0.457	0.253	0.391	0.271	0.142
VCD-A	ConvNeXt-B	640 $\times$ 1600	0.548	<b>0.631</b>	<b>0.436</b>	0.244	<b>0.343</b>	0.290	<b>0.120</b>

- VCD-A Results on nuScenes val
- VCD-A Results on nuScenes test
- Consistency Improvement

Methods	Venue	Backbone	mAP $\uparrow$	NDS $\uparrow$	mATE $\downarrow$	mASE $\downarrow$	mAOE $\downarrow$	mAVE $\downarrow$	mAAE $\downarrow$
BEVFusion [38]	NeurIPS 2022	LiDAR & Image	0.642	0.680	-	-	-	-	-
FUTR3D [10]	Arxiv 2022	LiDAR & Image	0.645	0.683	-	-	-	-	-
UVTR [33]	NeurIPS 2022	LiDAR & Image	0.654	0.702	0.332	0.258	<b>0.268</b>	0.212	<b>0.177</b>
CMT [57]	Arxiv 2023	LiDAR & Image	<b>0.679</b>	0.708	-	-	-	-	-
VCD-E	-	Image	0.677	<b>0.711</b>	<b>0.308</b>	<b>0.254</b>	0.317	<b>0.189</b>	0.201

Methods	Backbone	mAP	NDS
BEVFusion	ResNet-50	0.598	0.662
BEVFusion	ConvNext-B	0.597	0.665
VCD-E	ResNet-50	0.611	0.656
VCD-E	ConvNext-B	<b>0.664</b>	<b>0.693</b>

- VCD-E Results on nuScenes val
- Gains of different image backbone on multi-modal models

Expert	Paradigm	mAP	NDS
-	-	0.297	0.409
CenterPoint [60]	CM	0.281	0.420
Transfusion [2]	CM	0.292	0.435
BEVDepth [35]	UM	0.341	0.442
VCD-E	UM	<b>0.354</b>	<b>0.459</b>

Methods	mAP	NDS
Baseline [35]	0.297	0.409
FitNet [44]	0.318	0.421
CWD [46]	0.311	0.412
BEVDistill [11]	0.316	0.439
VCD-A	<b>0.354</b>	<b>0.459</b>

- The performance gains of the apprentice
- Effect of different distillation methods

Expert	Paradigm	mAP	NDS
-	-	0.297	0.409
CenterPoint [60]	CM	0.281	0.420
Transfusion [2]	CM	0.292	0.435
BEVDepth [35]	UM	0.341	0.442
VCD-E	UM	<b>0.354</b>	<b>0.459</b>

Methods	mAP	NDS
Baseline [35]	0.297	0.409
FitNet [44]	0.318	0.421
CWD [46]	0.311	0.412
BEVDistill [11]	0.316	0.439
VCD-A	<b>0.354</b>	<b>0.459</b>

- The performance gains of the apprentice
- Effect of different distillation methods

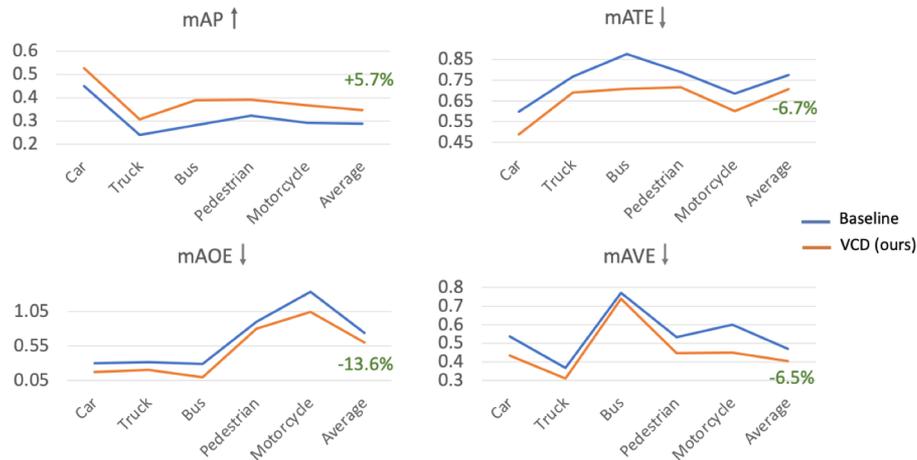
## Experiments

Temporal Length	Distill	mAP (%) $\uparrow$	NDS (%) $\uparrow$	mATE $\downarrow$	mAOE $\downarrow$	mAVE $\downarrow$
1	$\times$	26.6	37.9	0.815	0.645	0.556
	$\checkmark$	30.1 (+3.5)	41.5(+3.6)	0.732	0.629	0.476
2	$\times$	26.9	38.4	0.804	0.706	0.461
	$\checkmark$	31.3 (+4.4)	43.2 (+4.8)	0.717	0.615	0.403
4	$\times$	28.4	39.8	0.748	0.739	0.432
	$\checkmark$	33.0 (+4.6)	44.1 (+4.3)	0.707	0.632	0.389
8	$\times$	29.7	40.9	0.762	0.714	0.415
	$\checkmark$	35.4 (+5.7)	45.9 (+5.0)	0.690	0.625	0.370

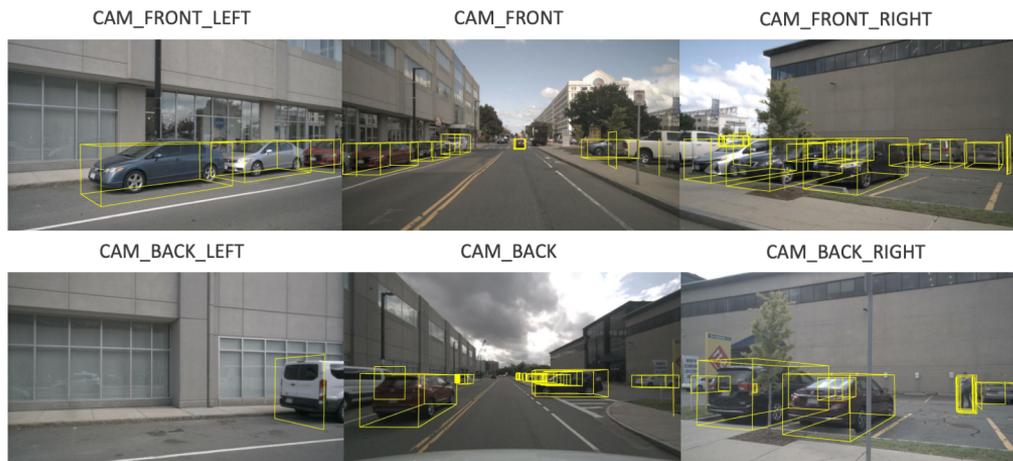
- The performance gains on different temporal lengths
- The performance gains of different trajectory length

Trajectory Length	Distill	mAP (%)	NDS (%)
-	$\times$	29.7	40.9
0	$\checkmark$	31.8	42.1
1	$\checkmark$	33.1	44.5
3	$\checkmark$	34.6	45.6
5	$\checkmark$	<b>35.4</b>	<b>45.9</b>
9	$\checkmark$	33.9	44.7

## Experiments



- The Effects of VCD on Movable Objects



- Visualization of the Predictions

OpenDriveLab 浦驾



上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory

# Thanks

[opendrive.com](https://opendrive.com) | X @OpenDriveLab