

Joint Perception and Prediction



- End to End Perception and Prediction
- Utilizing raw LiDAR sensor data directly for prediction, yields better performance
- Framework is trained end-to-end for prediction task which reduces compounding error between modules and reduces latency

LiDAR Representations for Joint Network

Voxelized Bird's Eye View (BEV)



- Preservation of metric space and straightforward fusion of historical LiDAR data.
- Loss of fine-grained information

LiDAR's native, dense representation, which is very helpful for detecting small objects.

Fusing historical LiDAR data in RV is challenging due to distortions

Contributions

Range View (RV)

- Realizing strengths and weaknesses of these representations, we propose an efficient data fusion framework
- Multi-view encoding and processing of LiDAR data separately in BEV and RV frames, before fusing the two views in a common BEV feature space.
- We propose computationally efficient sensor fusion of the camera RGB data with LiDAR in the RV frame, before projecting the learned features to the BEV frame.
- This proposal is evaluated against the state-of-the-art
 - Fusion method is quite general and can be applied to improve other BEV- and RV-based methods.

Multi-View Fusion of Sensor Data for Improved Perception and Prediction in Autonomous Driving

Sudeep Fadadu^{*}, Shreyash Pandey^{*}, Darshan Hegde, Yi Shi, Fang-Chieh Chou, Nemanja Djuric, Carlos Vallespi-Gonzalez

Multiview Fusion

- BEV is sparsely populated and does not use raw LIDAR data like intensity. RV methods operate in LiDAR's native, dense representation, providing full access to the
- non-quantized sensor information. We hypothesize that RV representation is also better than BEV for efficiently fusing information from sensors that natively capture data in RV, like camera.
- To fuse these multiple views together, we propose a point-based feature projection which gathers features from multiple views based on point location in each view.



We use similar method to LaserNet++[2] to associate range view pixels to camera features by projecting lidar points into the camera image

Quantitative Results and Analysis

- Proposed method shows improved detection and prediction metrics over MultiXNet for all classes.
- The improvement is more pronounced on smaller objects which benefit from the fusion of high resolution data.
- Camera fusion is compared against continuous fusion (ContFuse^[3]) which is the current **ContFuse** by a considerable margin.
- We see significant increase for actors at a longer ranges.

[2] Gregory P. Meyer, et al. Sensor fusion for joint 3d object detection and semantic segmentation. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops, June 2019

state-of-the-art for LiDAR-Camera fusion. Proposed method outperforms our implementation of

The inference latency increases only slightly and the proposed fusion is much faster than ContFuse.



Qualitative Results









Proposed Fusion Architecture applied to MultiXNet^[1]

Qualitative comparison of ours (middle raw) and MultiXNet (bottom)

Table 1: Evaluation on ATG4D using detection AP (%), prediction DE (cm), and latency (ms)

Vehicles							Pedestrians						Bicyclists					
$AP_{0.7}$ \uparrow						$AP_{0.1}$ \uparrow							AP _{0.3} ↑					
Full	In Camera FOV				DE↓	Full	In Camera FOV				DE↓	Full	In Camera FOV				DE↓	Lat.↓
	0-75m	0-25	25-50	50-75			0-75m	0-25	25-50	50-75			0-75m	0-25	25-50	<u>50-75</u>		
84.2	83.9	92.6	85.5	68.6	80	88.6	88.9	87.3	89.5	88.0	57	84.6	82.2	87.4	78.7	76.8	49	35.2
84. 7	84.5	92.6	86.1	70.1	79	89.4	89.8	88.5	90.1	88.7	57	87.3	85.2	90.8	81.5	76.8	49	43.4

Table 2: Evaluation on nuScenes using detection AP (%), prediction DE (cm), and latency (ms)

			Vehicles				P	ns								
$AP_{0.7}$ \uparrow							AP ₀	.1 ↑		9.		AP ₀				
Ful	Full	I In Camera FOV			DE↓	Full	In Camera FOV			DE↓	Full In Camera FOV			FOV	DE↓	Lat.↓
		0-50m	0-25	25-50			0-50m	0-25	25-50			0-50m	0-25	25-50	8	
	60.6	60.9	80.0	43.3	108	65.5	63.8	72.9	54.1	87	31.7	32.1	41.4	25.2	191	32.3
	60.9	61.9	79.4	46.0	109	67.2	71.6	78.0	64.3	82	33.5	44.1	53.5	36.1	188	63.6
	61.1	61.5	79.6	45.3	107	71.0	70.4	79.1	60.6	82	38.2	38.1	53.3	25.8	187	37.4
t	62.4	63.7	80.8	48.6	107	72.6	76.7	83.5	69.5	73	40.9	52.3	64.5	43.2	204	46.2
	62.9	63.3	81.2	47.0	107	71.4	73.1	80.4	64.6	80	39.8	40.4	53.6	29.5	179	38.3