



LiMoE: Mixture of LiDAR Representation Learners From Automotive Scenes

Motivation & Contribution

Overview of Approach

- \succ **LiMoE** is a new LiDAR-based cross-sensor representation learning framework that aims to integrate multiple LiDAR representations dynamically through the **MoE** paradigm.
- LiMoE integrates three popular LiDAR representations: range, voxel, and point, into a unified framework.
- > With **MoE** design to combine multiple LiDAR representations, **LiMoE** enables use to capture potential complementary cue from each representation with specific task objectives.
- Compared to existing singlerepresentation methods, we achieves large improvements across 11 LiDAR datasets.





 \succ Similarity map between a query point and the pretrained 2D image backbone as well as other LiDAR points during the CML stage demonstrate that LiMoE enhances semantic features during pretraining, even without semantic supervision.









Xiang Xu* Lingdong Kong* Hui Shuai Liang Pan Ziwei Liu Qingshan Liu



- Image-to-LiDAR Pretraining leverages previous LiDAR-based cross-sensor data pretraining framework to transfer prior knowledge from images to point clouds across different LiDAR representations, providing semanticrich initial parameters for each representation network.
- Contrastive Mixture Learning (CML) fuses the pretrained features across various LiDAR representations from stage #1 into a unified representation via a Mixture of Expert (MoE) layer. With the tailored contrastive learning objective, CML encourages the MoE to focus on data attributes from laser beams and attributes that play different roles from representations.
- > Semantic Mixture Supervision (SMS) aims to enhance the downstream performance by combining semantic logits with the MoE layers. By supervising the semantic logits in such a way, the SMS module can encourage experts to focus on object attributes in the LiDAR scenes.



Comparative & Ablation Study

Tab. Compare with state-of-the-art LiDAR Pretraining methods

Method	Venue	Backbone (2D)	Backbone (3D)	Expert	LP	1%	nuSo 5%	cenes 10%	25%	Full	KITTI 1%	Waymo 1%
Random	-	-	-	-	8.10	30.30	47.84	56.15	65.48	74.66	39.50	39.41
SLidR [63] TriCC [55] Seal [43] CSC [8] HVDistill [85]	CVPR'22 CVPR'23 NeurIPS'23 CVPR'24 IJCV'24	ResNet-50 [25]	MinkUNet-34 [14]	Single o Single o Single o Single o Single o	$38.80 \\ 38.00 \\ \underline{44.95} \\ 46.00 \\ 39.50$	38.30 41.20 <u>45.84</u> 47.00 42.70	52.49 54.10 55.64 57.00 <u>56.60</u>	59.84 60.40 <u>62.97</u> 63.30 62.90	66.91 67.60 68.41 <u>68.60</u> 69.30	74.79 75.60 75.60 <u>75.70</u> 76.60	44.60 45.90 46.63 <u>47.20</u> 49.70	47.12 - 49.34 - -
SLidR [63] + LiMoE Seal [43] SuperFlow [80] + LiMoE	CVPR'22 Ours NeurIPS'23 ECCV'24 Ours	ViT-S [53]	MinkUNet-34 [14]	Single • Multi • Single • Single • Multi •	44.70 45.80 45.16 <u>46.44</u> 48.20	41.16 46.82 44.27 <u>47.81</u> 49.60	53.65 57.54 55.13 <u>59.44</u> 60.54	61.47 63.85 62.46 <u>64.47</u> 65.65	66.71 68.61 67.64 <u>69.20</u> 71.39	74.20 75.64 75.58 <u>76.54</u> 77.27	44.67 46.81 46.51 <u>47.97</u> 49.53	47.57 48.81 48.67 <u>49.94</u> 51.42
SLidR [63] + LiMoE Seal [43] SuperFlow [80] + LiMoE	CVPR'22 Ours NeurIPS'23 ECCV'24 Ours	ViT-B [53]	MinkUNet-34 [14]	Single • Multi • Single • Single • Multi •	45.35 46.56 46.59 <u>47.66</u> 49.07	41.64 46.89 45.98 <u>48.09</u> 50.23	55.83 58.09 57.15 <u>59.66</u> 61.51	62.68 63.87 62.79 <u>64.52</u> 66.17	67.61 69.02 68.18 <u>69.79</u> 71.56	74.98 75.87 75.41 <u>76.57</u> 77.81	45.50 47.96 47.24 <u>48.40</u> 50.30	48.32 49.50 48.91 <u>50.20</u> 51.77
SLidR [63] + LiMoE Seal [43] SuperFlow [80] + LiMoE	CVPR'22 Ours NeurIPS'23 ECCV'24 Ours	ViT-L [53]	MinkUNet-34 [14]	Single • Multi • Single • Single • Multi •	45.70 47.43 46.81 <u>48.01</u> 49.35	42.77 46.92 46.27 <u>49.95</u> 51.41	57.45 58.41 58.14 <u>60.72</u> 62.07	63.20 64.54 63.27 <u>65.09</u> 66.64	68.13 69.69 68.67 <u>70.01</u> 71.59	75.51 76.32 75.66 <u>77.19</u> 77.85	47.01 48.25 47.55 <u>49.07</u> 50.69	48.60 50.23 50.02 <u>50.67</u> 51.93



- > In the SMS stage, we find that the range images are more sensitive to dynamic objects, voxels can highlight background objects, and raw LiDAR points capture more detailed objects.





Experiments & Analysis

 \succ LiMoE achieves significant improvements across various datasets with the suitable integrations of multiple LiDAR representations.

 \succ In the **CML** stage, we observe that different representations focus on distinct LiDAR data attributes. The range images predominantly capture middle beams and distances, sparse voxels focus more on upper beams and longer distances, while raw LiDAR points can concentrate on lower beams and near distances around ego-car.

