

# Segment Any Point Cloud Sequences by Distilling Vision Foundation Models

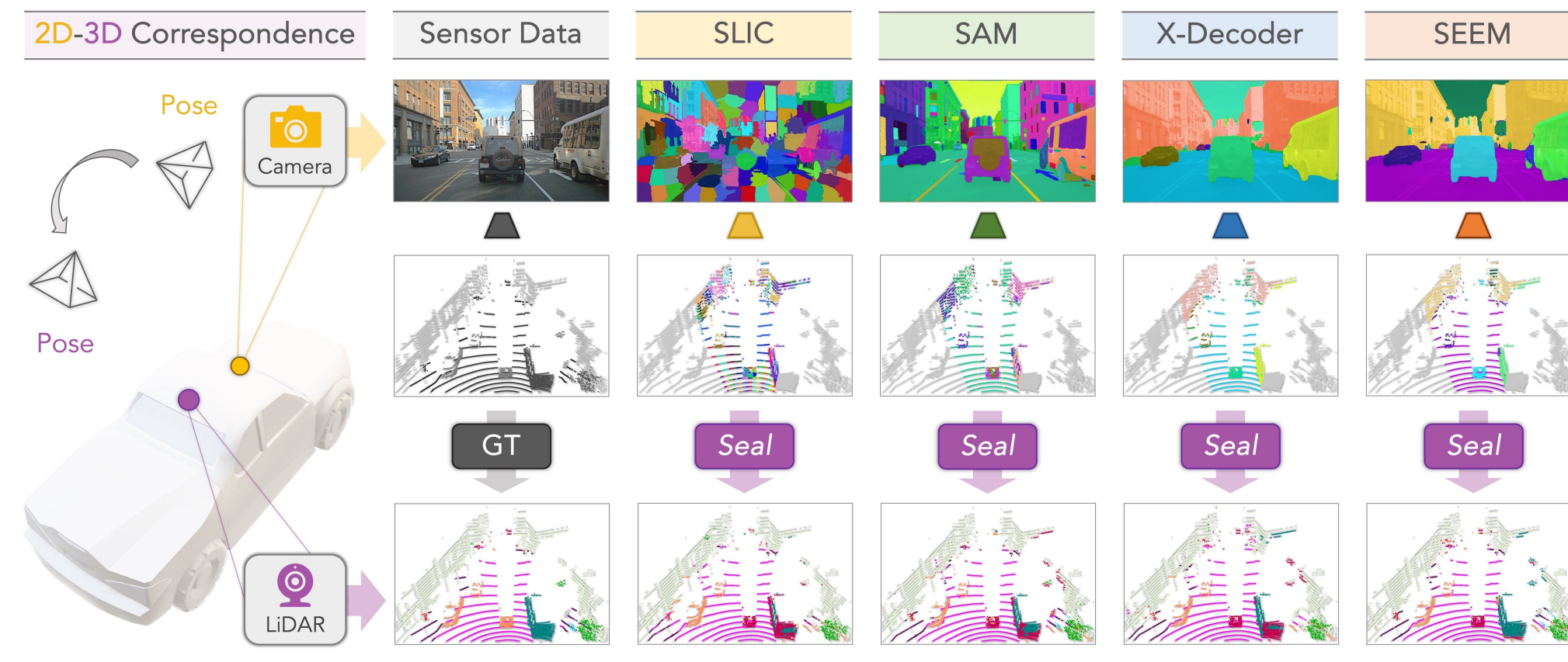
Youquan Liu\*, Lingdong Kong\*,  
Jun Cen, Runnan Chen, Wenwei Zhang,  
Liang Pan, Kai Chen, Ziwei Liu



## Motivation & Contribution

### TL;DR

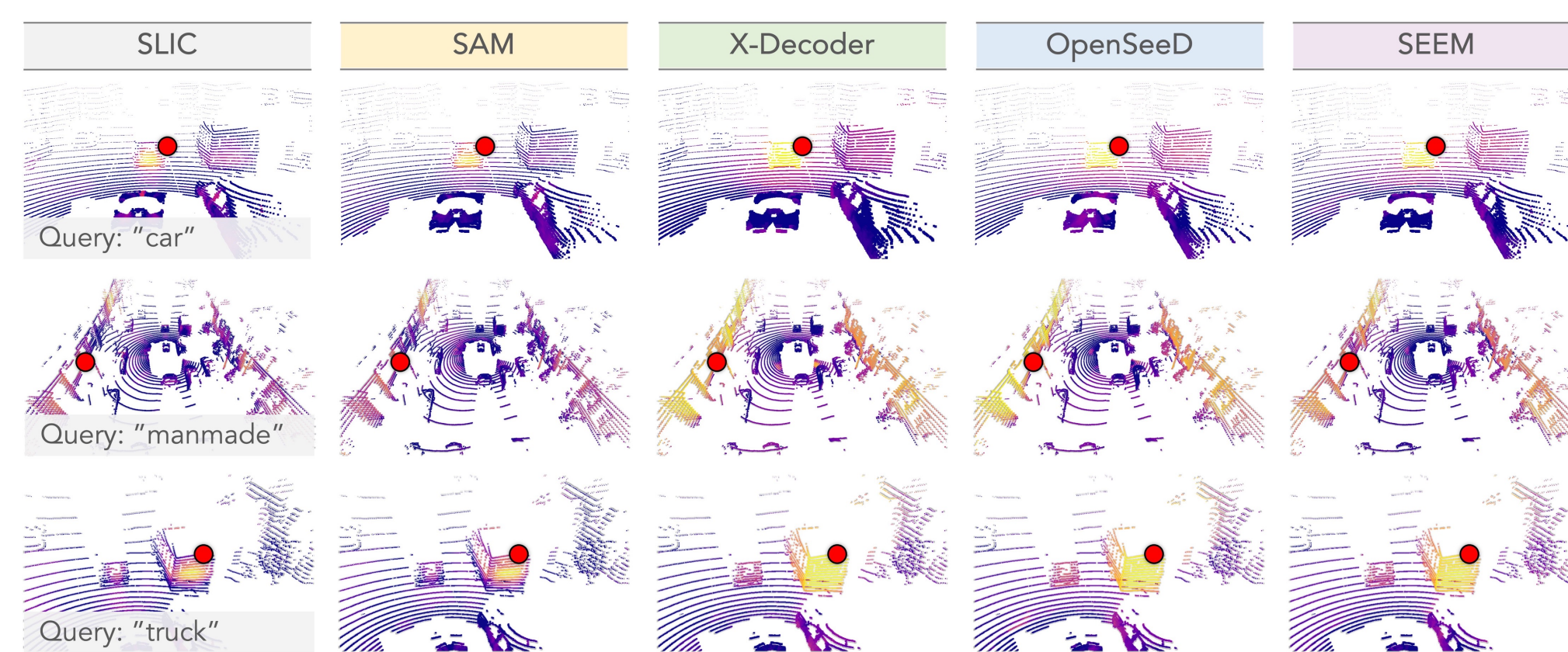
❖ We introduce **Seal**, a novel framework tailored to harness **vision foundation models (VFMs)** to segment diverse **automotive** point cloud sequences.



❖ **Seal** has **three** appealing properties: i) **Scalability** for not needing either 2D or 3D annotations during pretraining; ii) **Consistency** for aligning between LiDAR and camera via cross-modal contrastive learning; iii) **Generalizability** for exhibiting effectiveness across a wide range of point cloud datasets.

### Knowledge Transfer from VFMs

❖ VFMs can generate **superpixels** from the camera views and provide **off-the-shelf** semantic coherence for distinct objects & backgrounds in the 3D scene.



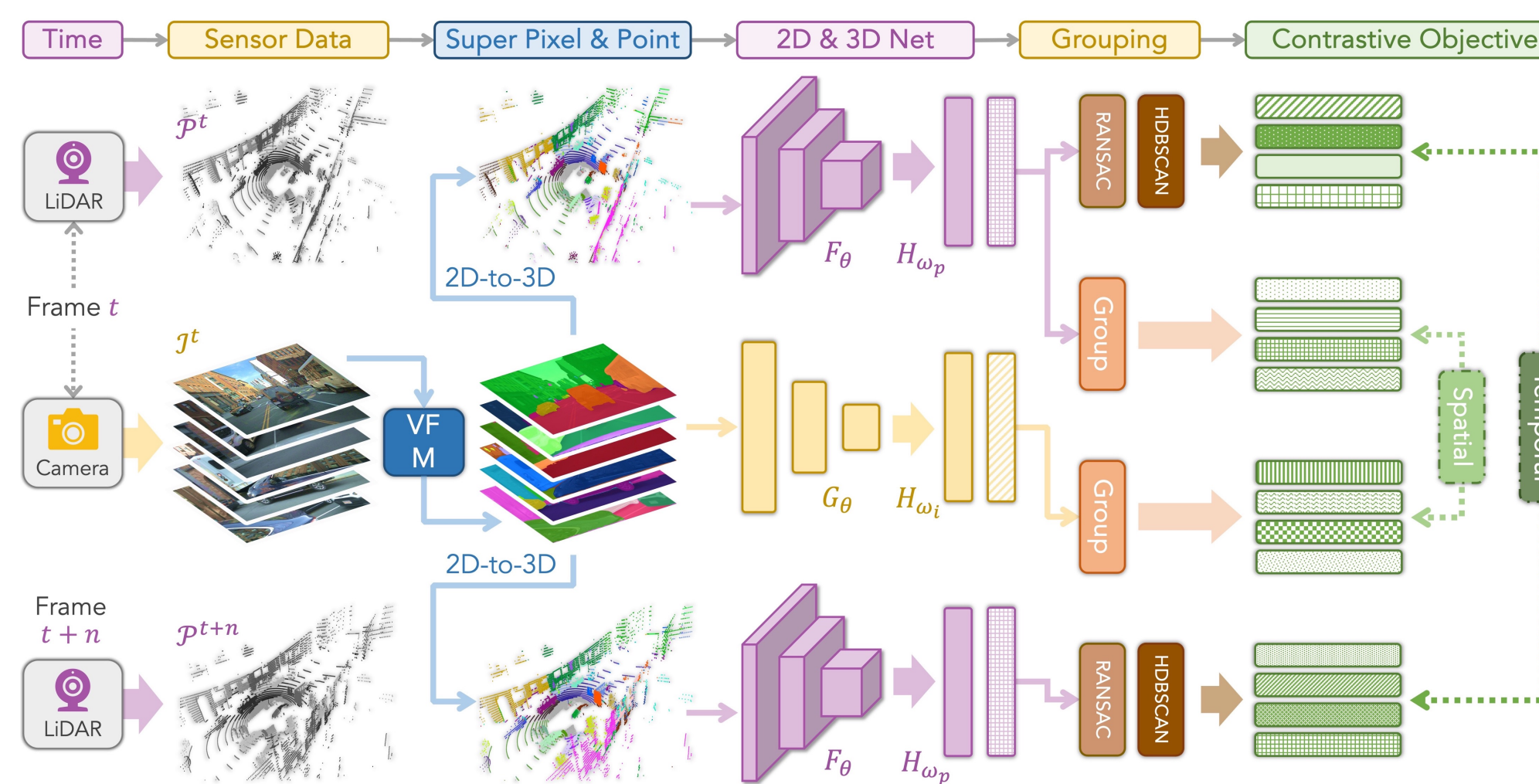
❖ Compared to prior works, our VFM-assisted contrastive learning: i) mitigates the severe self-conflict problem; ii) forms a **more coherent** optimization landscape, yielding a **faster convergence rate**; iii) reduces the number of superpixels generated, which **extenuates the overhead** during pretraining.

## Methodology

### The Seal Framework

❖ We generate, for each **{LiDAR, camera}** pair  $\{\mathcal{P}^t, \mathcal{J}^t\}$  at timestamp  $t$  and another LiDAR frame  $\mathcal{P}^{t+n}$  at timestamp  $t+n$ , the semantic **superpixels** and **superpoints** by **VFMs**. Two pertaining objectives are then leveraged.

❖ We aim to encourage i) **spatial contrastive** between paired LiDAR-camera features for cross-sensor learning; ii) **temporal consistency** between point segments at two different timestamps for semantic view regularization.



❖ These two **regularization** objectives are **complementary** to each other; a combination of both introduces strong **consistency** during the pretraining.

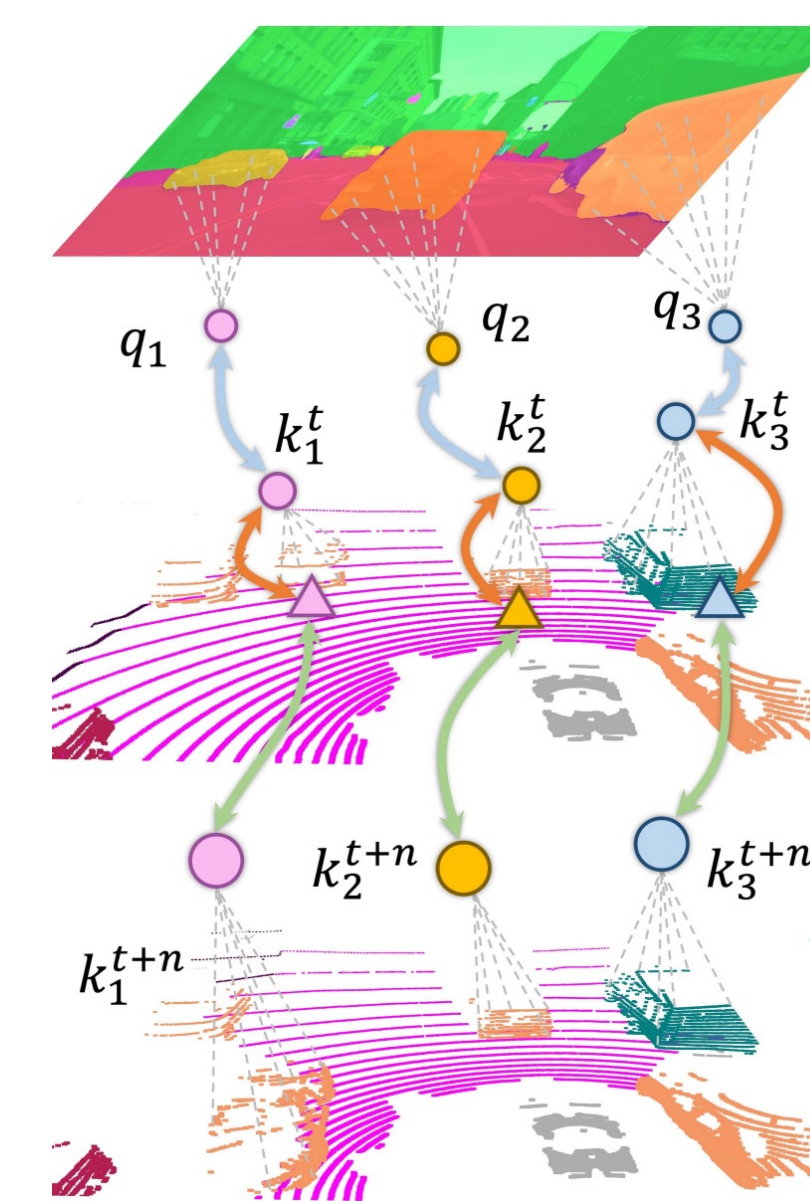
### Spatial-Temporal Consistency

❖ **Seal** defines a suitable **positive feature correspondence** in contrastive learning via implicit geometric clustering.

❖ Such a design can mitigate the potential errors caused by inaccurate cross-sensor **calibration** and **synchronization**.

❖ Besides, point-to-segment regularization mechanism can serve to aggregate the spatial information thus yielding **better-distinguishing** instances in LiDAR scenes.

❖ Regularization at different levels enables **Seal** to be **consistent** and **generalizable**.



## Experiments & Analyses

### Comparative Study

❖ We verify the effectiveness of **Seal** across **eleven** point cloud datasets with various scales, modalities, sensor configurations, fidelities, and noise levels.

Method & Year	nuScenes						KITTI	Waymo	Synth4D
	LP	1%	5%	10%	25%	Full	1%	1%	1%
Random	8.10	30.30	47.84	56.15	65.48	74.66	39.50	39.41	20.22
PointContrast [ECCV'20] [103]	21.90	32.50	-	-	-	-	41.10	-	-
DepthContrast [ICCV'21] [116]	22.10	31.70	-	-	-	-	41.50	-	-
PPKT [arXiv'21] [65]	35.90	37.80	53.74	60.25	67.14	74.52	44.00	47.60	61.10
SLiDR [CVPR'22] [85]	38.80	38.30	52.49	59.84	66.91	74.79	44.60	47.12	63.10
ST-SLiDR [CVPR'23] [66]	40.48	40.75	54.69	60.75	67.70	75.14	44.72	44.93	-
<b>Seal (Ours)</b>	<b>44.95</b>	<b>45.84</b>	<b>55.64</b>	<b>62.97</b>	<b>68.41</b>	<b>75.60</b>	<b>46.63</b>	<b>49.34</b>	<b>64.50</b>
Seal <sup>1</sup> (Ours)	-	48.41	57.84	65.52	70.80	77.13	-	-	-
Seal <sup>2</sup> (Ours)	-	49.53	58.64	66.78	72.31	78.28	-	-	-

❖ Our approach constantly outperforms previous works on **every** setting by **large margins**, which strongly demonstrates the **superiority** and **scalability**.

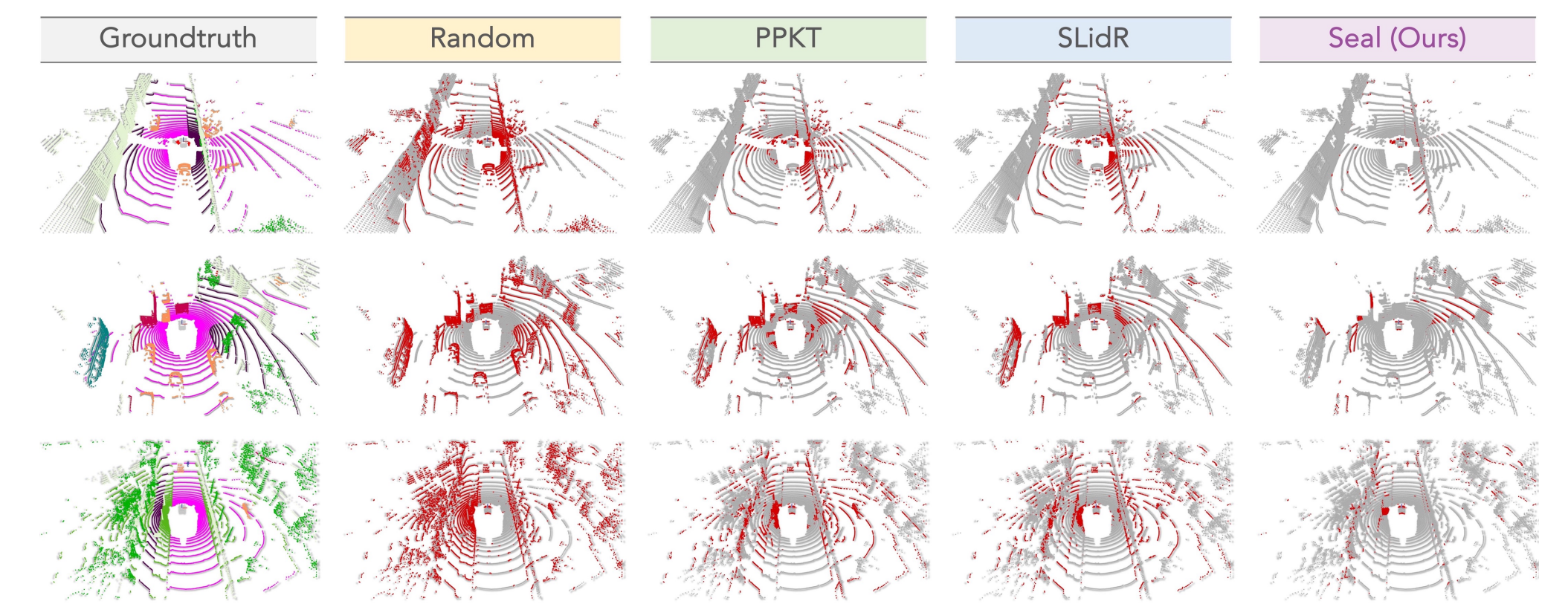
Method	ScribbleKITTI		RELLIS-3D		SemanticPOSS		SemanticSTF		SynLiDAR		DAPS-3D	
	1%	10%	1%	10%	Half	Full	Half	Full	1%	10%	Half	Full
Random	23.81	47.60	38.46	53.60	46.26	54.12	48.03	48.15	19.89	44.74	74.32	79.38
PPKT [65]	36.50	51.67	49.71	54.33	50.18	56.00	50.92	54.69	37.57	46.48	78.90	84.00
SLiDR [85]	39.60	50.45	49.75	54.57	51.56	55.36	52.01	54.35	42.05	47.84	81.00	85.40
<b>Seal (Ours)</b>	<b>40.64</b>	<b>52.77</b>	<b>51.09</b>	<b>55.03</b>	<b>53.26</b>	<b>56.89</b>	<b>53.46</b>	<b>55.36</b>	<b>43.58</b>	<b>49.26</b>	<b>81.88</b>	<b>85.90</b>

### Ablation Study

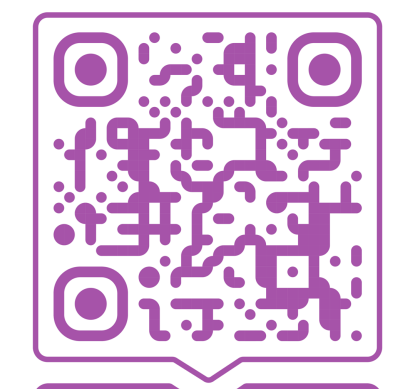
❖ The **effectiveness** of each component in **Seal** has been proven: with spatial & temporal contrastive, we can learn meaningful **multi-modal representations**.

#	C2L	VFM	STC	P2S	nuScenes						KITTI	Waymo
					LP	1%	5%	10%	25%	Full	1%	1%
(1)	✓				38.80	38.30	52.49	59.84	66.91	74.79	44.60	47.12
(2)	✓		✓		40.45	41.62	54.67	60.48	67.61	75.30	45.38	48.08
(3)	✓	✓			43.00	44.02	53.03	60.84	67.38	75.21	45.72	48.75
(4)	✓	✓	✓		44.01	44.78	55.36	61.99	67.70	75.00	46.49	49.15
(5)	✓	✓		✓	43.35	44.25	53.69	61.11	67.42	75.44	46.07	48.82
(6)	✓	✓	✓	✓	<b>44.95</b>	<b>45.84</b>	<b>55.64</b>	<b>62.97</b>	<b>68.41</b>	<b>75.60</b>	<b>46.63</b>	<b>49.34</b>

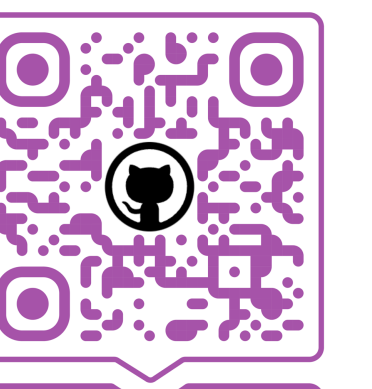
❖ Qualitative results show that **Seal** can segment complex **driving scenes** in 3D.



香港大學  
THE UNIVERSITY OF HONG KONG



Paper



Code