



浙江大學  
ZHEJIANG UNIVERSITY



# ***Not All Voxels Are Equal: Hardness-Aware Semantic Scene Completion with Self-Distillation***

Song Wang<sup>1</sup>, Jiawei Yu<sup>1</sup>, Wentong Li<sup>1</sup>, Wenyu Liu<sup>1</sup>, Xiaolu Liu<sup>1</sup>,  
Junbo Chen<sup>2\*</sup> and Jianke Zhu<sup>1\*</sup>

<sup>1</sup>Zhejiang University <sup>2</sup>Udeer.ai

<https://github.com/songw-zju/HASSC>

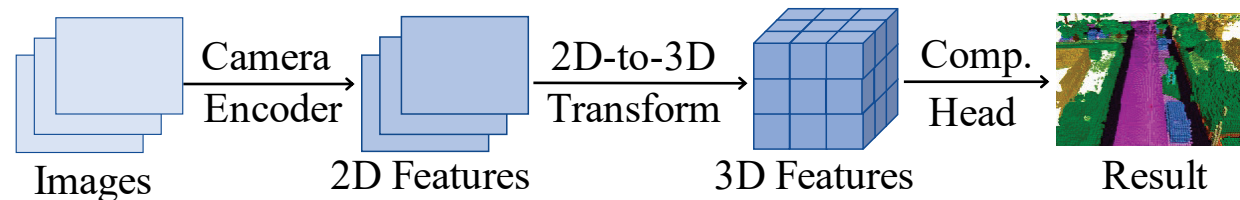
Poster: Arch 4A-E-17

Contact: {songw, jkzhu}@zju.edu.cn, junbo@udeer.ai

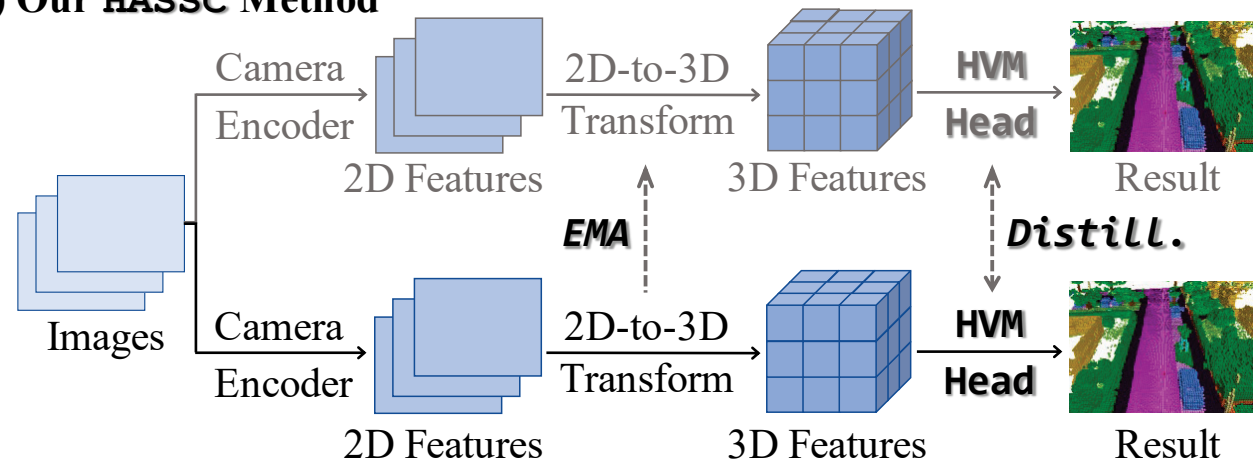
# Overview

- The **global hardness** from the network optimization process is defined for dynamical hard voxel selection.
- The **local hardness** with geometric anisotropy is adopted for voxel-wise refinement.
- **Self-distillation strategy** is introduced to make training process stable and consistent.

(a) Previous SSC Methods



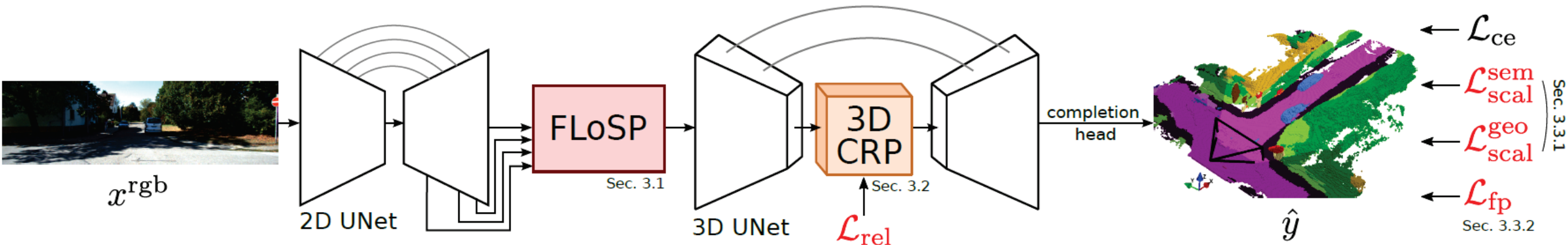
(b) Our **HASSC** Method



01

# Motivation

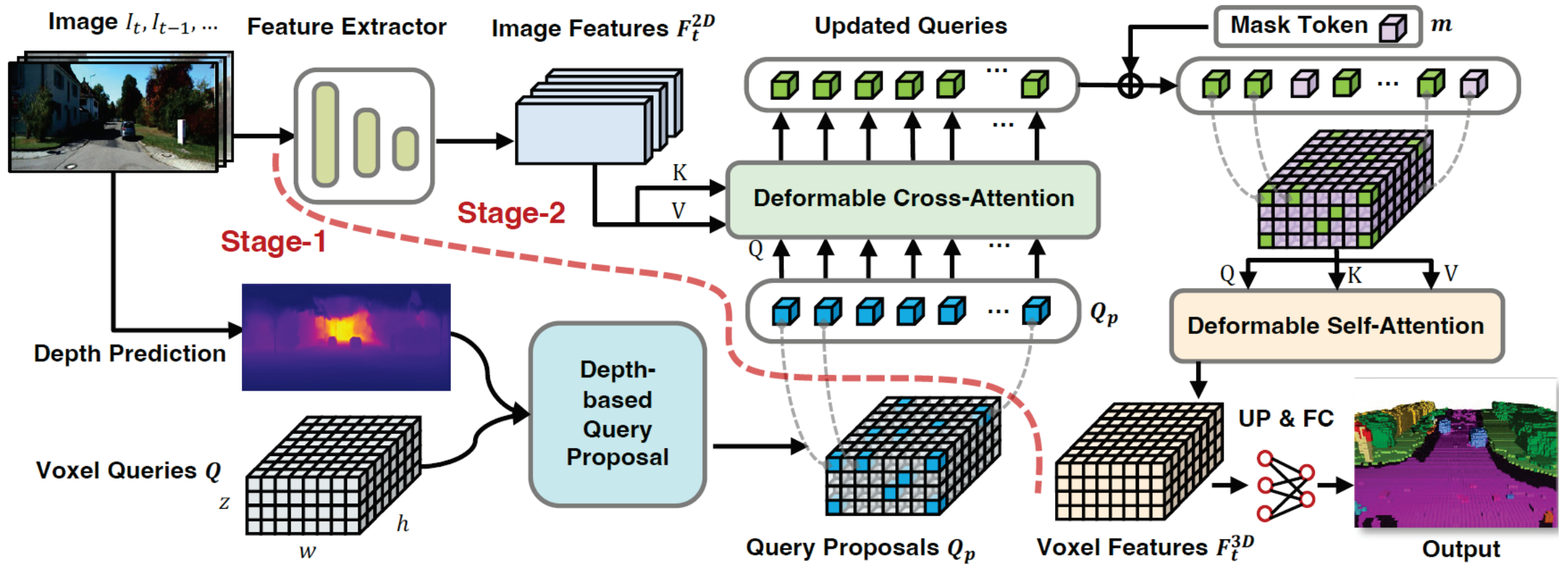
# Motivation



MonoScene, CVPR 2022

Cao A Q, De Charette R. Monoscene: Monocular 3d semantic scene completion. CVPR 2022.

# Motivation

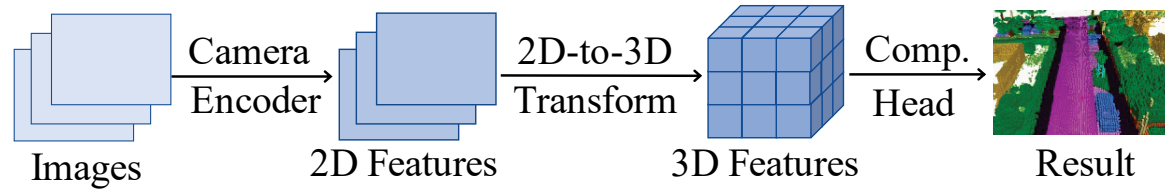


VoxFormer, CVPR 2023

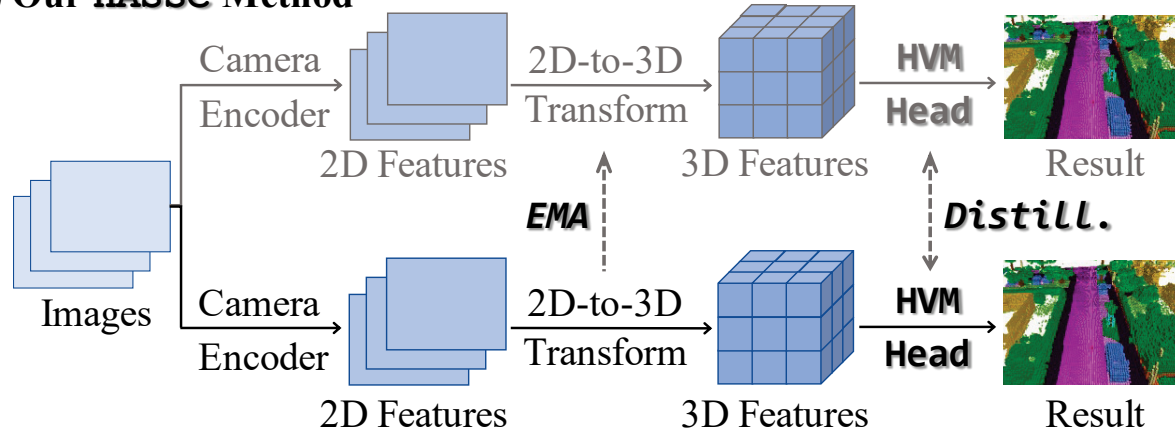
Li Y, Yu Z, Choy C, et al. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. CVPR 2023.

# Motivation

(a) Previous SSC Methods



(b) Our **HASSC** Method



➤ The 3D dense space typically contains a **large number of empty voxels**, which are easy to learn but require amounts of computation due to **handling all the voxels uniformly** for the existing models.

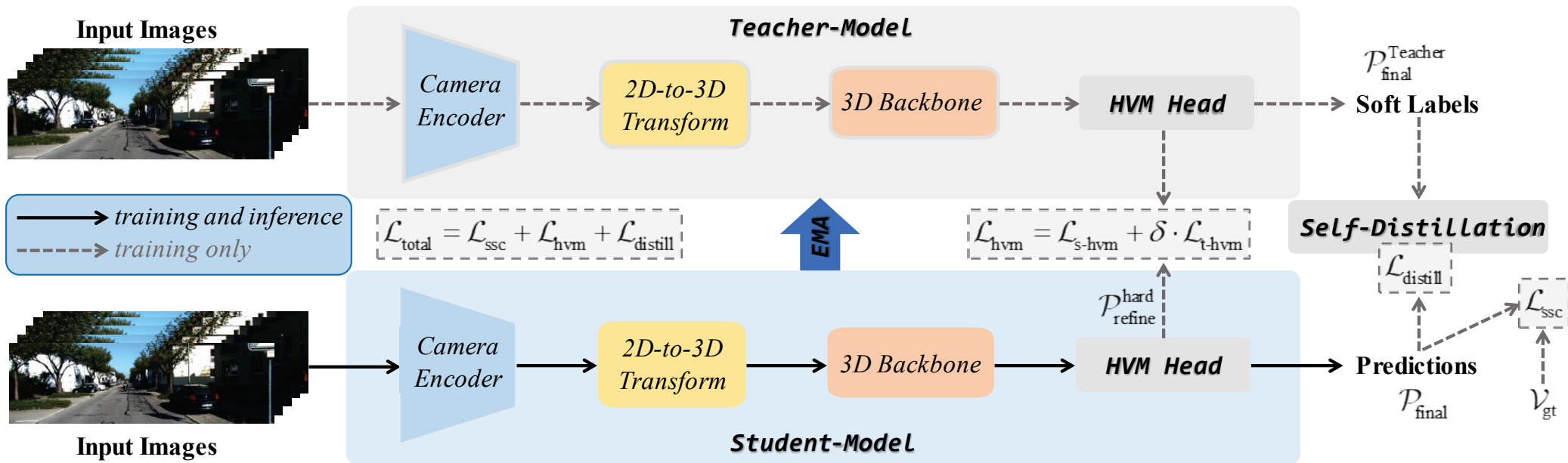
➤ The voxels in the **boundary region** are more challenging to differentiate than those in the **interior**.

02

# Method

# Our Framework

- Propose a **hardness-aware semantic scene completion** (HASSC) scheme that can be easily integrated into existing models without incurring extra cost for inference.
- Take advantage of both the **global and local hardness** to find the hard voxels so that their predictions can be refined by weighted voxel-wise losses during training.





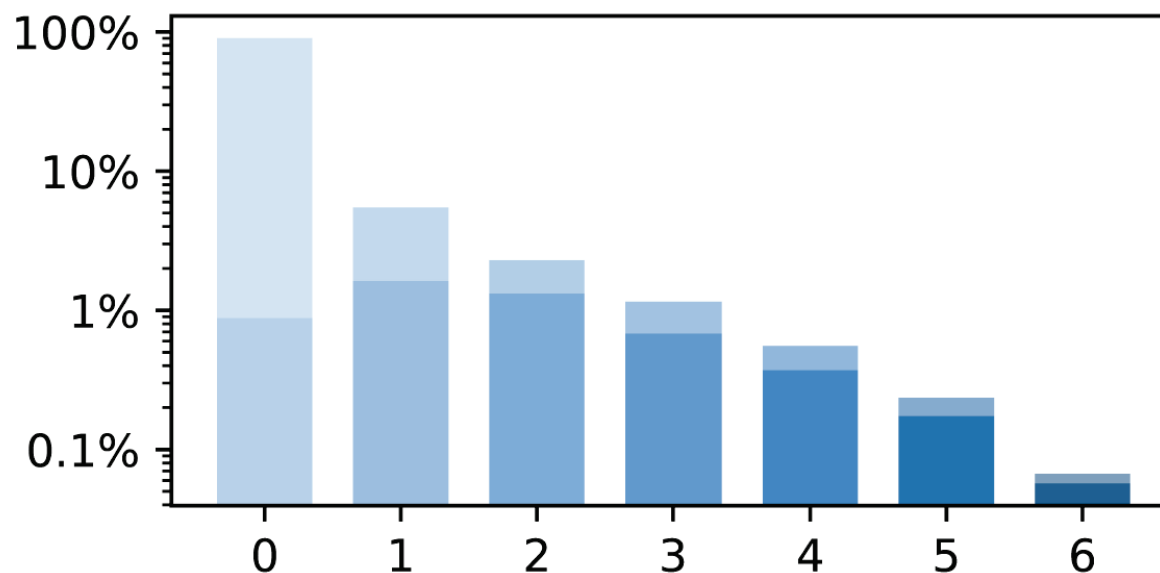
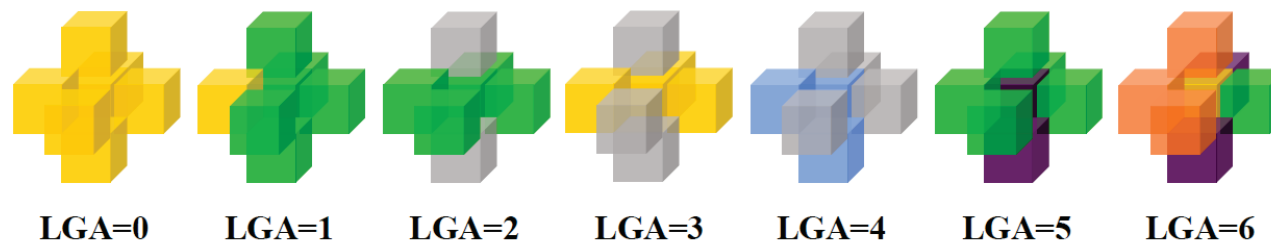
# Global Hardness

- For the prediction of each voxel, we rank the probabilities of each class in decreasing order. The largest probability in  $C$  classes is represented as  $p^a$ , and the second largest one is denoted as  $p^b$ . Then, the global hardness of this voxel is defined as follows

$$\mathcal{H}_{i,j,k}^{\text{global}} = \frac{1}{p^a - p^b}$$

- The global hardness measures the uncertainty of the semantic scene completion prediction between the class  $a$  and  $b$ , which varies with the optimization of the network. We mainly employ the global hardness to select hard voxels and refine their predictions.

# Local Hardness

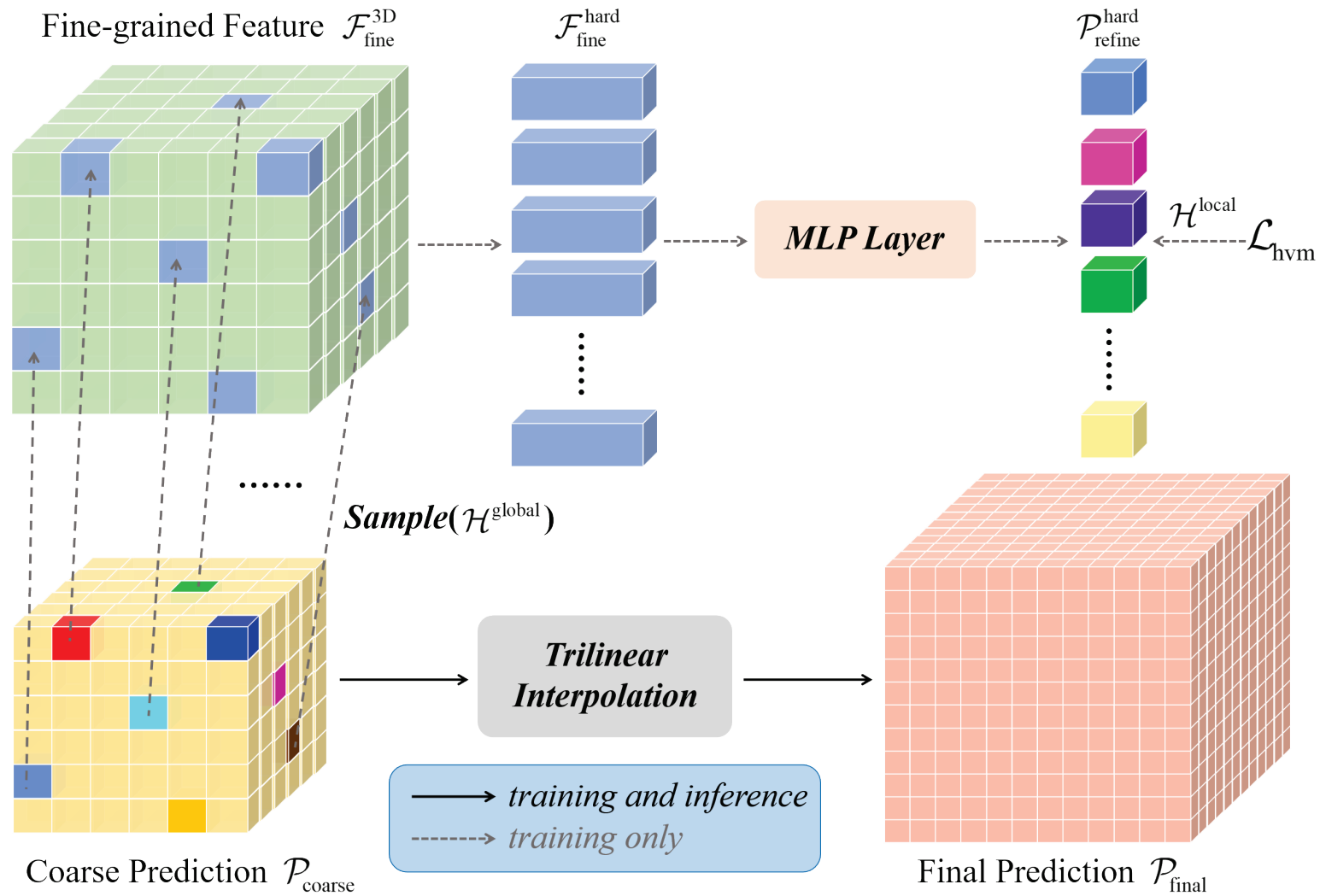


Local Geometric Anisotropy (LGA) Distribution on SemanticKITTI

$$\mathcal{A}_{i,j,k} = \sum_{m=1}^M (v_{\text{gt}} \oplus v_{\text{gt}}^m)$$

$$\mathcal{H}_{i,j,k}^{\text{local}} = \alpha + \beta \mathcal{A}_{i,j,k}$$

# Hard Voxel Mining (HVM) Head



# Training Loss

➤ For Hard Voxel Mining

$$\mathcal{P}_{\text{refine}}^{\text{hard}} = \text{MLP}(\mathcal{F}_{\text{fine}}^{\text{hard}})$$
$$\mathcal{L}_{\text{s-hvm}} = \frac{1}{N} \sum_{n=1}^N \mathcal{H}_n^{\text{local}} \cdot \text{CE}(v_{\text{refine}}^n, v_{\text{gt}}^n)$$

➤ For Self-Distillation

$$\mathcal{L}_{\text{distill}} = \lambda e^{\mu} \cdot \mathbf{D}_{\text{KL}}(\mathcal{P}_{\text{final}}^{\text{Teacher}} \parallel \mathcal{P}_{\text{final}})$$

➤ Totally,

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{ssc}} + \mathcal{L}_{\text{hvm}} + \mathcal{L}_{\text{distill}}$$

03

# Experiments

# Performance Comparison

— SemanticKITTI Validation

Methods	VoxFormer-S [30]			HASSC VoxFormer-S			VoxFormer-T [30]			HASSC VoxFormer-T			StereoScene <sup>†</sup> [22]			HASSC StereoScene		
Modality	Camera			Camera			Camera			Camera			Camera			Camera		
Range	S	M	L	S	M	L	S	M	L	S	M	L	S	M	L	S	M	L
IoU (%) <sup>↑</sup>	65.35	57.54	44.02	<b>65.54</b>	<b>57.99</b>	<b>44.82</b>	65.38	57.69	44.15	<b>66.05</b>	<b>58.01</b>	<b>44.58</b>	65.70	56.84	43.66	65.52	<b>57.01</b>	<b>44.55</b>
mIoU (%) <sup>↑</sup>	17.66	16.48	12.35	<b>18.98</b>	<b>17.95</b>	<b>13.48</b>	21.55	18.42	13.35	<b>24.10</b>	<b>20.27</b>	<b>14.74</b>	23.27	21.15	15.24	<b>24.43</b>	<b>22.17</b>	<b>15.88</b>
■ car (3.92%)	39.78	35.24	25.79	<b>42.37</b>	<b>36.78</b>	<b>27.23</b>	44.90	37.46	26.54	<b>45.79</b>	<b>37.70</b>	<b>27.33</b>	47.05	43.52	31.15	46.47	43.02	30.64
■ bicycle (0.03%)	3.04	1.48	0.59	2.72	<b>2.26</b>	<b>0.92</b>	5.22	2.87	1.28	4.23	2.11	1.07	2.38	2.15	1.05	<b>4.20</b>	<b>2.63</b>	<b>1.20</b>
■ motorcycle (0.03%)	2.84	1.10	0.51	<b>4.49</b>	<b>1.63</b>	<b>0.86</b>	2.98	1.24	0.56	<b>5.64</b>	<b>2.03</b>	<b>1.14</b>	4.78	2.84	1.55	<b>5.26</b>	<b>3.34</b>	0.91
■ truck (0.16%)	7.50	7.47	5.63	6.25	<b>11.00</b>	<b>9.91</b>	9.80	10.38	7.26	<b>22.89</b>	<b>21.90</b>	<b>17.06</b>	18.72	22.48	17.55	<b>24.94</b>	<b>34.73</b>	<b>23.72</b>
■ other-veh. (0.20%)	8.71	4.98	3.77	<b>14.77</b>	<b>8.85</b>	<b>5.61</b>	17.21	10.61	7.81	<b>22.71</b>	<b>13.52</b>	<b>8.83</b>	17.33	13.79	9.26	<b>20.61</b>	<b>14.24</b>	7.77
■ person (0.07%)	4.10	3.31	1.78	<b>5.11</b>	<b>4.89</b>	<b>2.80</b>	4.44	3.50	1.93	<b>5.12</b>	<b>4.18</b>	<b>2.25</b>	6.31	4.37	2.17	6.06	3.58	1.79
■ bicyclist (0.07%)	6.82	7.14	3.32	<b>6.87</b>	<b>8.57</b>	<b>4.71</b>	2.65	3.92	1.97	<b>4.09</b>	<b>6.58</b>	<b>4.09</b>	7.70	4.75	2.30	<b>8.22</b>	<b>5.65</b>	<b>2.47</b>
■ motorcyclist (0.05%)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
■ road (15.30%)	72.40	65.74	54.76	<b>74.49</b>	<b>68.04</b>	<b>57.05</b>	75.45	66.15	53.57	<b>78.51</b>	<b>70.02</b>	<b>57.23</b>	79.24	74.16	61.86	<b>80.61</b>	<b>75.53</b>	<b>62.75</b>
■ parking (1.12%)	10.79	18.49	15.50	<b>15.49</b>	<b>21.23</b>	<b>15.90</b>	21.01	23.96	19.69	<b>29.43</b>	<b>26.69</b>	<b>19.89</b>	21.33	21.19	17.02	<b>25.21</b>	<b>25.95</b>	<b>20.20</b>
■ sidewalk (11.13%)	39.35	33.20	26.35	<b>42.69</b>	<b>36.32</b>	<b>28.25</b>	45.39	34.53	26.52	<b>51.69</b>	<b>38.83</b>	<b>29.08</b>	50.71	41.86	30.58	<b>52.68</b>	<b>43.61</b>	<b>32.40</b>
■ other-grnd(0.56%)	0.00	1.54	0.70	<b>0.02</b>	<b>2.38</b>	<b>1.04</b>	0.00	0.76	0.42	0.00	<b>1.55</b>	<b>1.26</b>	0.00	1.12	0.85	0.00	0.18	0.51
■ building (14.10%)	17.91	24.09	17.65	<b>22.78</b>	<b>27.30</b>	<b>19.05</b>	25.13	29.45	19.54	<b>27.99</b>	<b>30.81</b>	<b>20.19</b>	26.98	32.52	22.71	<b>29.09</b>	31.68	<b>22.90</b>
■ fence (3.90%)	12.98	10.63	7.64	9.81	8.70	6.58	16.17	11.15	7.31	<b>17.09</b>	<b>11.65</b>	<b>7.95</b>	22.50	14.26	8.73	20.88	13.32	8.67
■ vegetation (39.3%)	40.50	34.68	24.39	40.49	<b>35.53</b>	<b>25.48</b>	43.55	38.07	26.10	<b>44.68</b>	<b>38.93</b>	<b>27.01</b>	40.20	36.10	24.81	<b>40.29</b>	<b>36.44</b>	<b>26.27</b>
■ trunk (0.51%)	15.81	10.64	5.08	14.93	<b>11.25</b>	<b>6.15</b>	21.39	12.75	6.10	<b>22.22</b>	<b>14.11</b>	<b>7.71</b>	21.45	15.28	7.17	<b>21.65</b>	14.92	7.14
■ terrain (9.17%)	32.25	35.08	29.96	<b>36.66</b>	<b>38.28</b>	<b>32.94</b>	42.82	39.61	33.06	<b>47.04</b>	<b>41.37</b>	<b>33.95</b>	45.75	43.67	34.87	<b>48.50</b>	<b>46.95</b>	<b>38.10</b>
■ pole (0.29%)	14.47	11.95	7.11	<b>15.25</b>	<b>12.48</b>	<b>7.68</b>	20.66	15.56	9.15	18.95	14.76	<b>9.20</b>	20.43	18.95	10.66	18.67	16.34	9.00
■ traf.-sign (0.08%)	6.19	6.29	4.18	5.52	5.61	4.05	10.63	8.09	4.94	9.89	<b>8.44</b>	4.81	9.21	8.91	5.19	<b>10.88</b>	<b>9.08</b>	<b>5.23</b>

# Performance Comparison

—— SemanticKITTI Test

Methods	SSC Input	Pub.	IoU (%) $\uparrow$	mIoU (%) $\uparrow$
LMSCNet* [40]	$\hat{x}_{3D}^{occ}$	3DV 2020	31.38	7.07
3DSketch* [5]	$x^{rgb}, \hat{x}^{TSDF}$	CVPR 2020	26.85	6.23
AICNet* [24]	$x^{rgb}, \hat{x}^{depth}$	CVPR 2020	23.93	7.09
JS3C-Net* [53]	$\hat{x}^{pts}$	AAAI 2021	34.00	8.97
MonoScene [4]	$x^{rgb}$	CVPR 2022	34.16	11.08
TPVFormer [18]	$x^{rgb}$	CVPR 2023	34.25	11.26
OccFormer [60]	$x^{rgb}$	ICCV 2023	34.53	12.32
NDC-Scene [56]	$x^{rgb}$	ICCV 2023	36.19	12.58
VoxFormer-S [30]	$x^{rgb}$	CVPR 2023	42.95	12.20
VoxFormer-T [30]	$x^{rgb} \times 5$	CVPR 2023	43.21	13.41
<b>HASSC-VoxFormer-S</b>	$x^{rgb}$	-	<b>43.40</b>	13.34
<b>HASSC-VoxFormer-T</b>	$x^{rgb} \times 5$	-	42.87	<b>14.38</b>

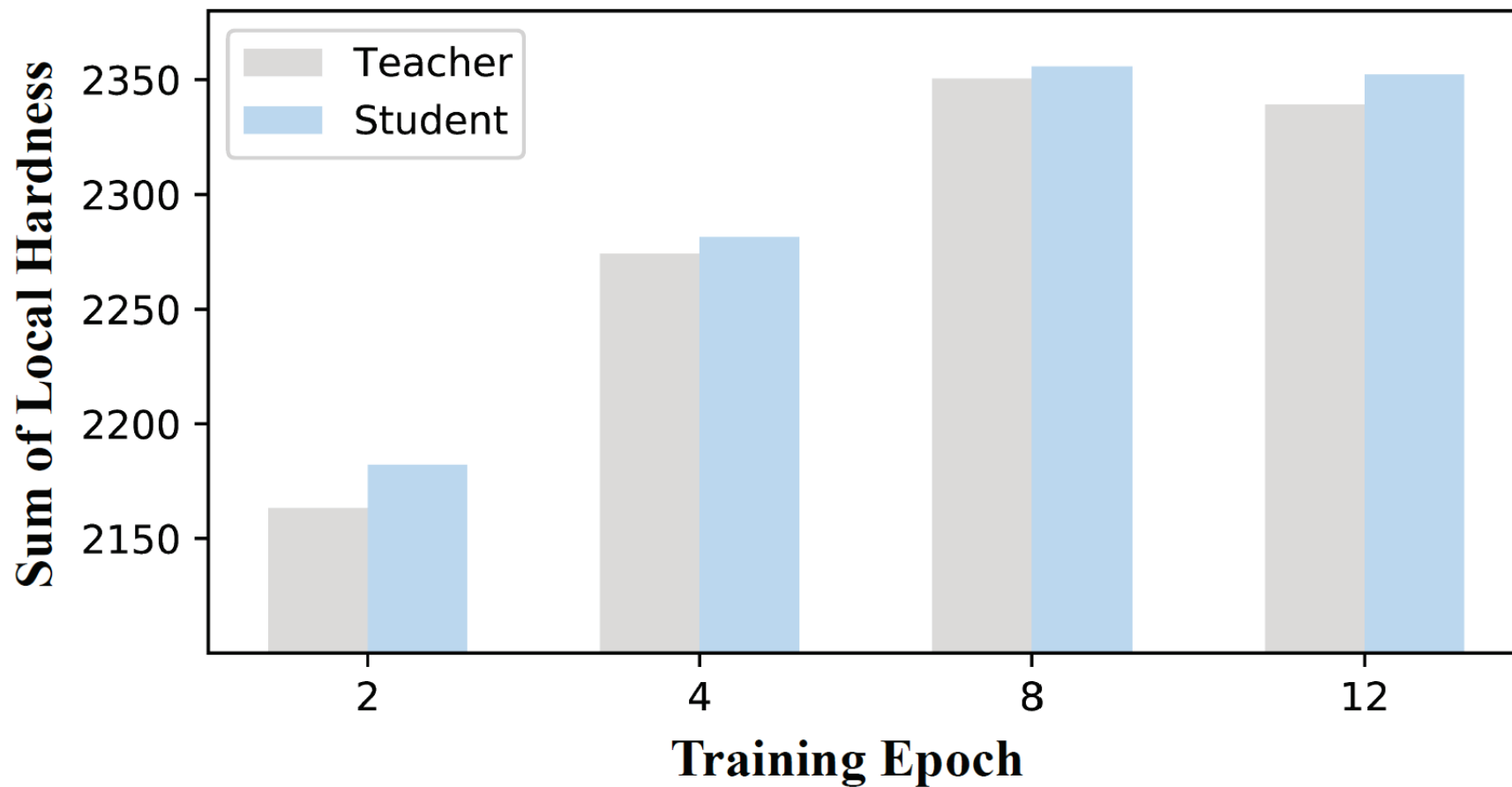
# Ablation Study

Global	Local	T-HVM	T-Distill	IoU (%)↑	mIoU (%)↑
				44.16	13.33
✓				43.89	13.30
	✓			44.00	13.40
✓	✓			43.98	13.91
✓	✓	✓		44.12	14.03
			✓	44.38	13.65
✓	✓	✓	✓	<b>44.58</b>	<b>14.74</b>

Ablation study on our proposed **HASSC** scheme



# Ablation Study



Visualization of the sum of the local hardness change during training on both student and teacher branches

# Ablation Study

Methods	VoxFormer-T	HASSC-VoxFormer-T
Params (M)	57.91	58.43
Inference Speed (ms)	724.05	720.84
IoU (%) $\uparrow$	44.16	<b>44.58</b>
mIoU (%) $\uparrow$	13.33	<b>14.74</b>

Comparison with baseline model on the training and inference efficiency

Methods	Hardness	IoU (%) $\uparrow$	mIoU (%) $\uparrow$
PALNet [23]	Local	44.28	13.28
PointRend [20]	Global	44.29	13.57
Xiao <i>et al.</i> [52]	Global	44.10	13.33
Ours	Global & Local	<b>44.58</b>	<b>14.74</b>

Comparison with other hard sample mining schemes

Voxel Numbers ( $N$ )	0	1024	2048	4096	8192
IoU (%) $\uparrow$	<b>44.16</b>	44.01	43.92	44.12	44.09
mIoU (%) $\uparrow$	13.33	13.52	13.64	<b>14.03</b>	13.74

Ablation study on the number of hard voxel selection

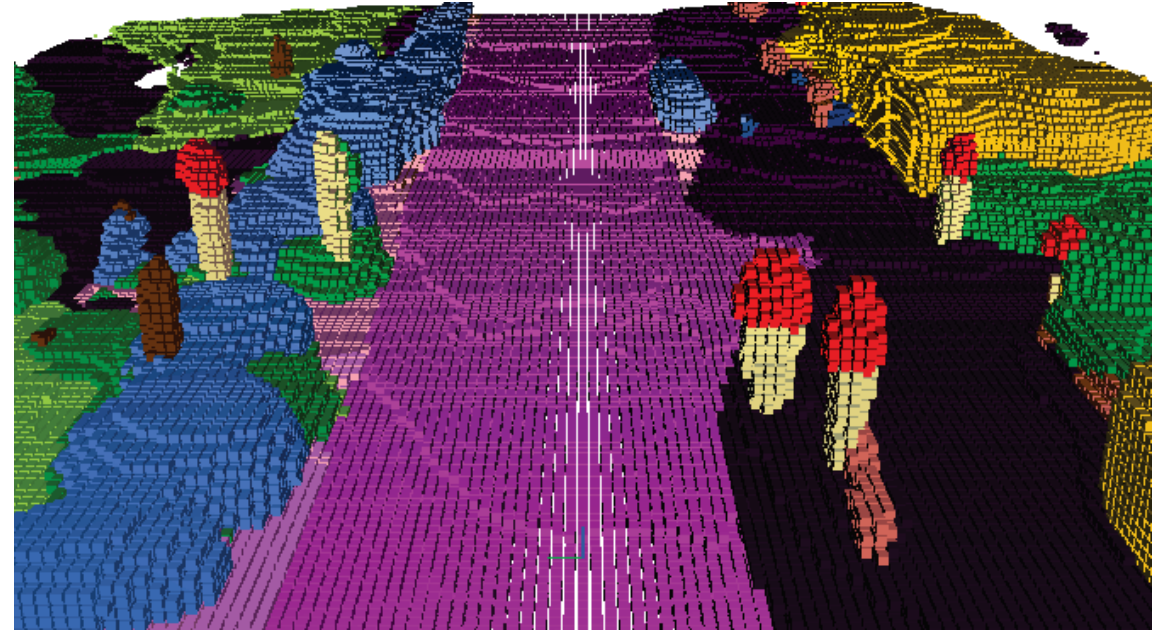
Distill Weight ( $\lambda$ )	0	12	24	48	96
IoU (%) $\uparrow$	44.12	44.06	44.25	<b>44.58</b>	44.51
mIoU (%) $\uparrow$	14.03	13.80	14.23	<b>14.74</b>	14.39

Ablation study on the weight of self-distillation from teacher model

# Performance on Sequence 08 (*validation set*)



Camera View (Left)



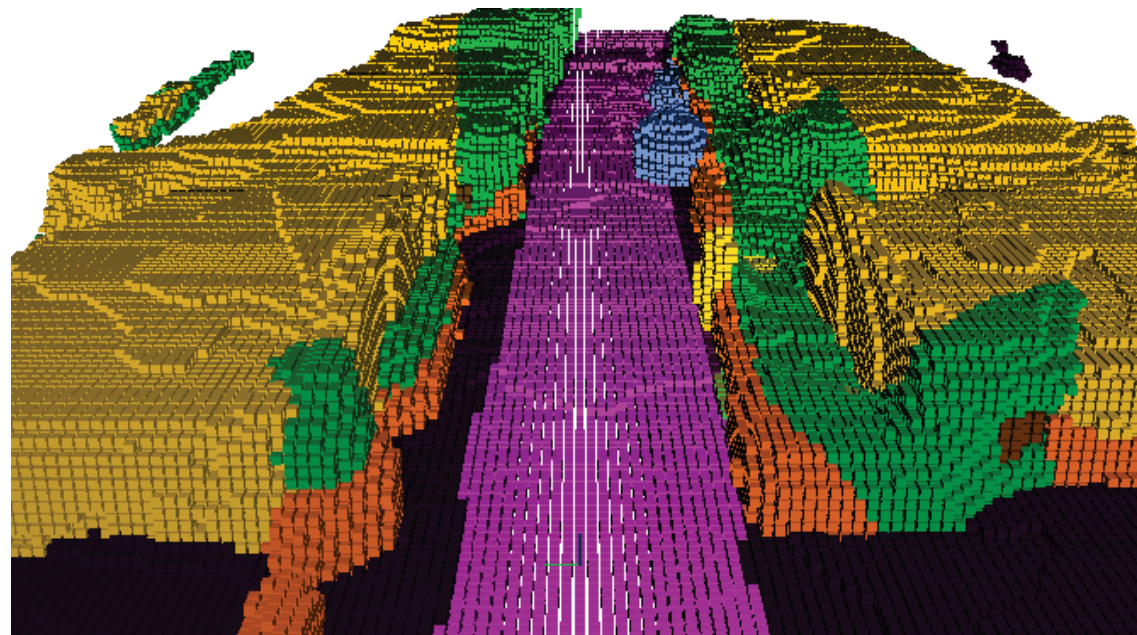
Semantic Scene Completion (Result)



# Performance on Sequence 11 (*hidden test set*)



Camera View (Left)



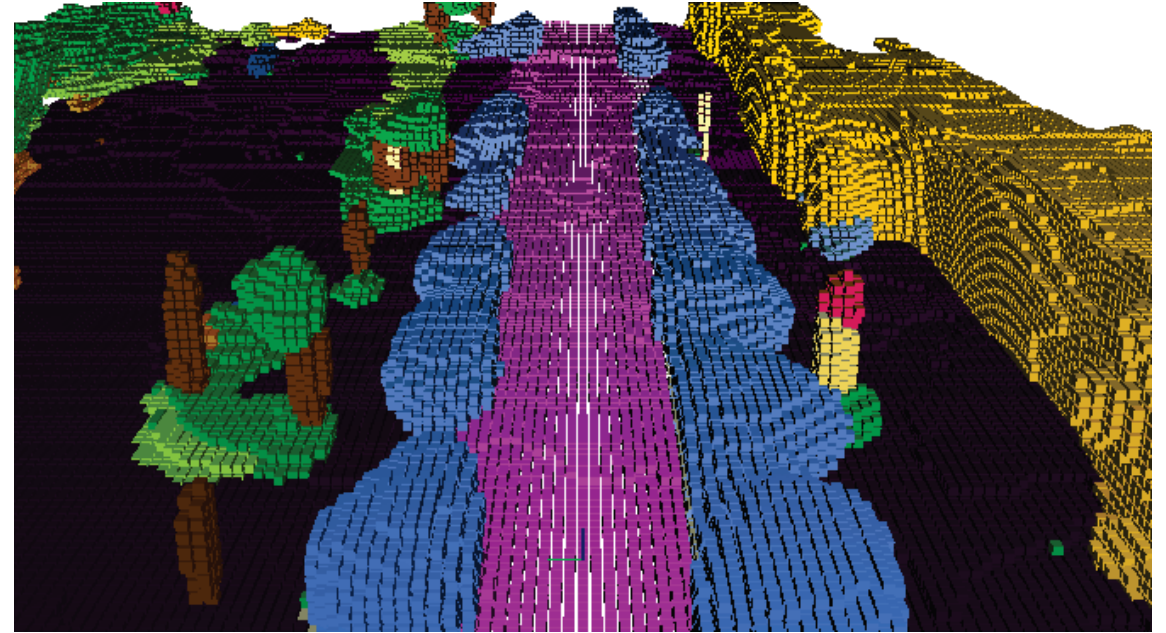
Semantic Scene Completion (Result)



# Performance on Sequence 13 (*hidden test set*)



Camera View (Left)



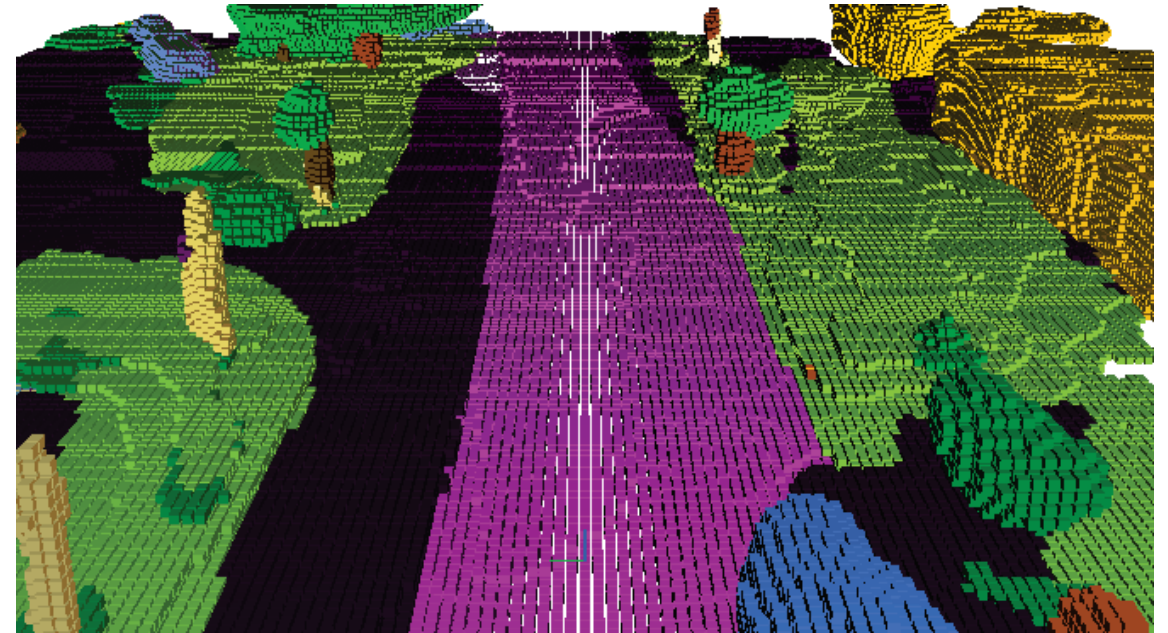
Semantic Scene Completion (Result)



# Performance on Sequence 16 (*hidden test set*)



Camera View (Left)



Semantic Scene Completion (Result)



**Thanks!**