

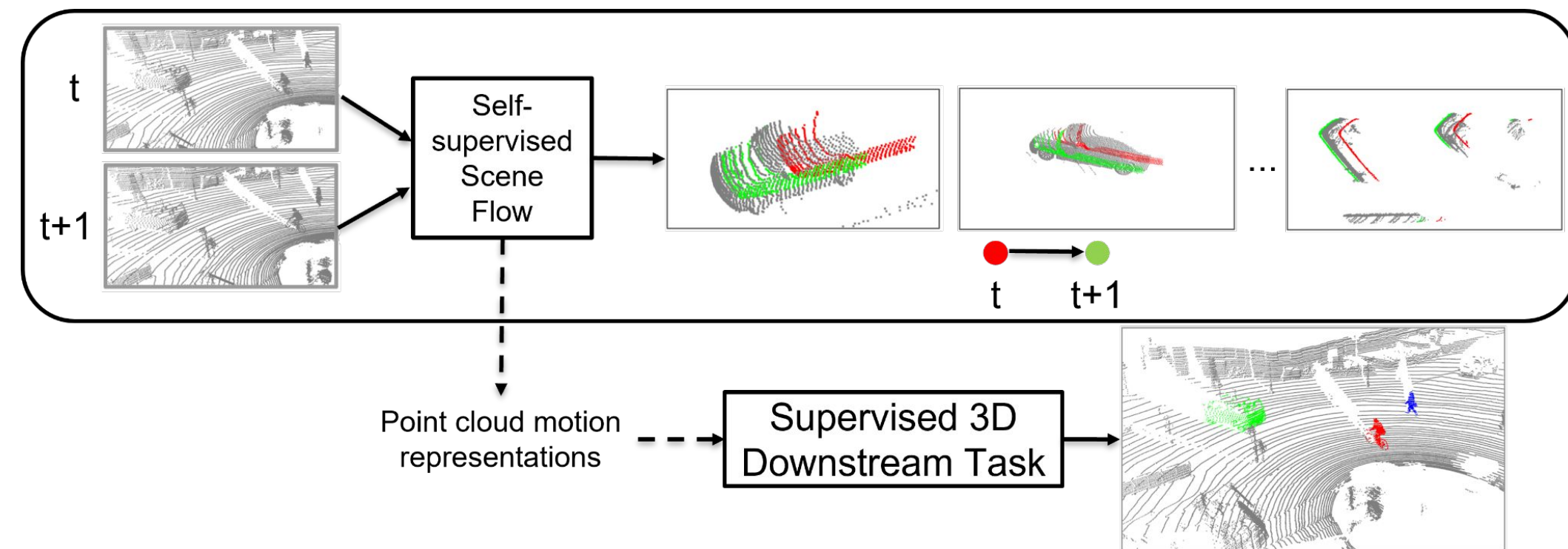
Emeç Erçelik<sup>\*1</sup>, Ekim Yurtsever<sup>\*2</sup>, Mingyu Liu<sup>1,3</sup>, Zhijie Yang<sup>1</sup>, Hanzhen Zhang<sup>1</sup>,  
Pinar Topçam<sup>1</sup>, Maximilian Listl<sup>1</sup>, Yılmaz Kaan Çaylı<sup>1</sup>, and Alois Knoll<sup>1</sup>

<sup>1</sup>Technical University of Munich, Germany <sup>2</sup>Ohio State University, Columbus, USA <sup>3</sup>Tongji University, Shanghai, China



## Motivation

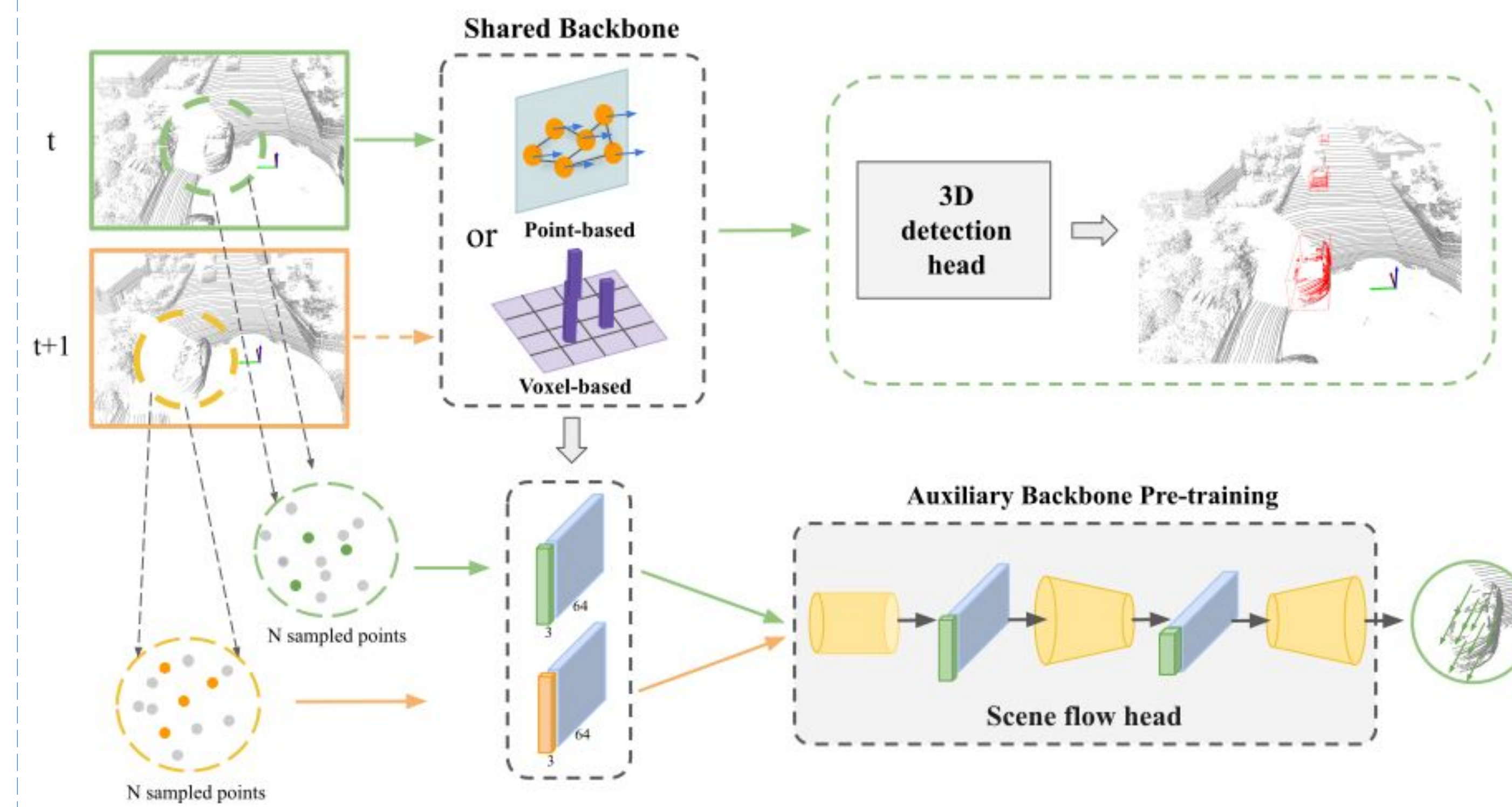
- Self-supervised learning aims to relieve the need for **large labor-intensive labeled data**.
- For 3D vision tasks, self-supervised learning has been underexplored.
- Contrary to contrastive approaches, we aim to use **inherent temporal change** in sequential lidar data by employing **self-supervised scene flow**.
- Learned **motion representations** provide distinctive information for the 3D detector that can be used while **differentiating** objects in the environment.



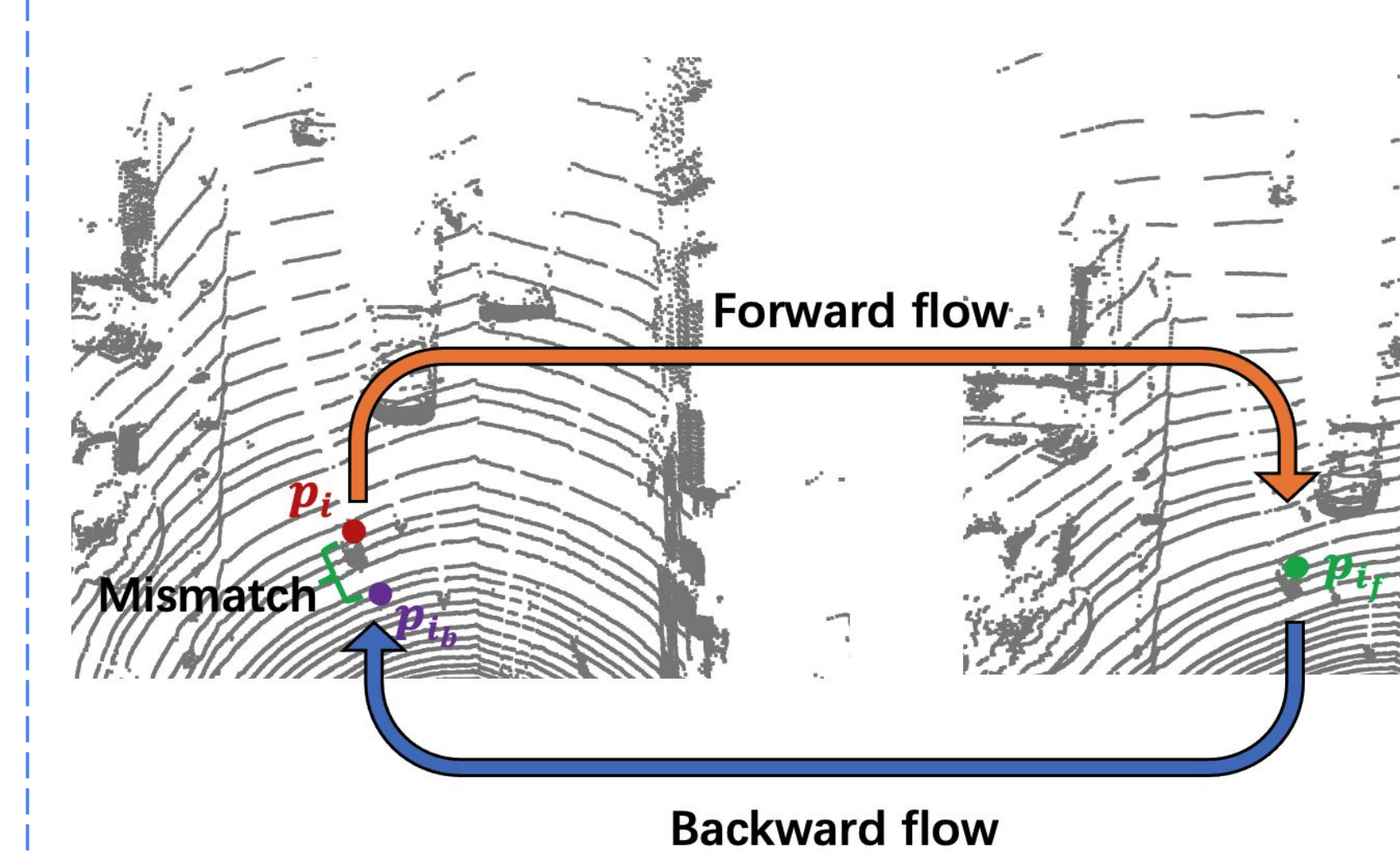
## Contributions

- Employing self-supervised point cloud scene flow estimation to learn **motion representations** for 3D object detection in tandem with supervised fine-tuning
- We show that auxiliary training is the best strategy for using **self-supervised cycle-consistency loss** along with supervised 3D detection loss.
- Our strategy is **especially effective with a lesser amount of supervised data**. We obtained a significant performance boost when only a smaller part of labeled data was used for the 3D detection task.

## Methodology



- We extract features of the sampled points from **two successive frames** using the 3d detector's backbone
- A modified **Flownet3d head** estimates the flow vectors and the **self-supervised cycle consistency loss** trains the head and the backbone.
- Then we fine-tune the **pre-trained backbone and the 3d detection head** on the smaller labelled 3d detection data.
- We also apply **alternating training**, which repeats these two steps using trained backbone from one step prior.



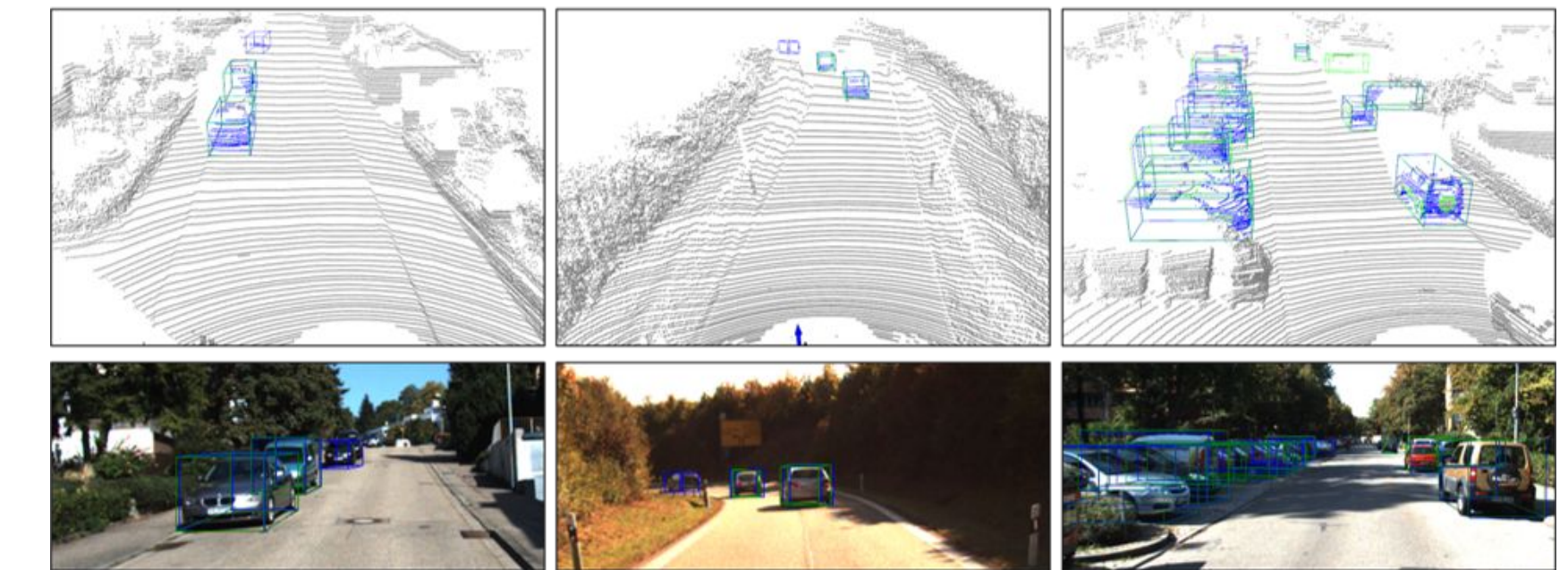
The cycle consistency loss makes use of the **mismatch** of the points propagated to the same frame through the **forward and backward** passes.

### Alternating training

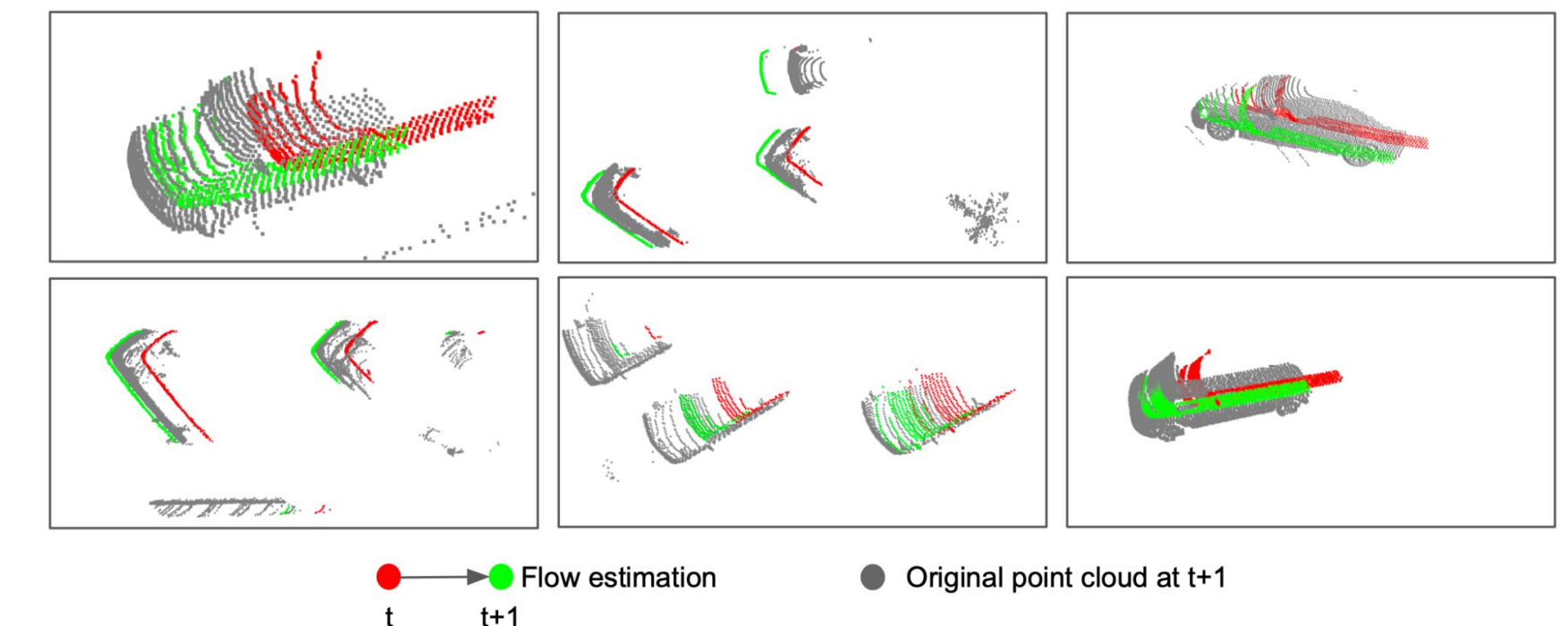
Step	Training	Backbone Init.	Head Init.
Step 1	Scene Flow	-	-
Step 2	3D detection	Step 1	-
Step 3	Scene Flow	Step 2	Step 1
Step 4	3D detection	Step 3	Step 2

## Qualitative Results

3D Object Detection results on KITTI val set



Sparse scene flow on KITTI tracking



## Conclusion

- We propose a **self-supervised motion-aware backbone pre-training method for 3D object detection**.
- Scene flow training using the cycle consistency** helps the backbone learn distinct features.
- Our experimental results on nuScenes and KITTI datasets show that **our method can improve 3D Detectors performance significant**.

## References

- [1] Lang, Alex H., et al. "Pointpillars: Fast encoders for object detection from point clouds." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.
- [2] Yin, Tianwei, Xingyi Zhou, and Philipp Krahenbuhl. "Center-based 3d object detection and tracking." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2021.
- [3] Shi, Weijing, and Raj Rajkumar. "Point-gnn: Graph neural network for 3d object detection in a point cloud." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.
- [4] Zhu, Xinge, et al. "Ssn: Shape signature networks for multi-class object detection from point clouds." *European Conference on Computer Vision*. Springer, Cham, 2020.
- [5] Liu, Xingyu, Charles R. Qi, and Leonidas J. Guibas. "Flownet3d: Learning scene flow in 3d point clouds." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.
- [6] Mittal, Himangi, Brian Okorn, and David Held. "Just go with the flow: Self-supervised scene flow estimation." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.
- [7] Geiger, Andreas, Philip Lenz, and Raquel Urtasun. "Are we ready for autonomous driving? The kitti vision benchmark suite." *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012.
- [8] Caesar, Holger, et al. "nuScenes: A multimodal dataset for autonomous driving." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.
- [9] Du, Liang, et al. "Associate-3Ddet: Perceptual-to-conceptual association for 3D point cloud object detection." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.
- [10] Liang, Ming, et al. "Multi-task multi-sensor fusion for 3d object detection." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019.
- [11] Wang, Guojun, et al. "CenterNet3D: An anchor free object detector for autonomous driving." *arXiv preprint arXiv:2007.07214* (2020).
- [12] Yan, Yan, Yuxing Mao, and Bo Li. "Second: Sparsely embedded convolutional detection." *Sensors* 18.10 (2018): 3337.
- [13] Zhou, Dingfu, et al. "Joint 3d instance segmentation and object detection for autonomous driving." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020.
- [14] Wang, Jun, et al. "Infocross: 3d object detection for autonomous driving with dynamic information modeling." *European Conference on Computer Vision*. Springer, Cham, 2020.
- [15] Vora, Sourabh, et al. "Pointpainting: Sequential fusion for 3d object detection." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.
- [16] Xie, Saining, et al. "Pointcontrast: Unsupervised pre-training for 3d point cloud understanding." *European conference on computer vision*. Springer, Cham, 2020.
- [17] Liang, Hanxue, et al. "Exploring geometry-aware contrast and clustering harmonization for self-supervised 3D object detection." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021.

## Quantitative 3D Object Detection Results

### KITTI validation set

Method	3D AP <sub>R40</sub>			BEV AP <sub>R40</sub>		
	Easy	Mod	Hard	Easy	Mod	Hard
Point-GNN	90.44	82.12	77.70	93.03	89.31	86.86
<b>Self-supervised Point-GNN</b>	<b>91.43</b>	<b>82.85</b>	<b>80.12</b>	<b>93.55</b>	<b>89.79</b>	<b>87.23</b>
Improvement	<b>+0.99</b>	<b>+0.73</b>	<b>+2.42</b>	<b>+0.52</b>	<b>+0.48</b>	<b>+0.37</b>
PointPillars	85.41	73.98	67.76	89.93	86.57	85.20
<b>Self-supervised PointPillars</b>	<b>85.92</b>	<b>76.33</b>	<b>74.32</b>	<b>89.96</b>	<b>87.44</b>	<b>85.53</b>
Improvement	<b>+0.51</b>	<b>+2.36</b>	<b>+6.56</b>	<b>+0.03</b>	<b>+0.87</b>	<b>+0.33</b>

### KITTI test set

Method	3D AP <sub>R40</sub>			BEV AP <sub>R40</sub>		
	Easy	Mod	Hard	Easy	Mod	Hard
Associate-3Ddet[9]	85.99	77.40	70.53	91.40	88.09	82.96
UBER-ATG-MMF[10]	88.40	77.43	70.22	93.67	88.21	81.99
CenterNet3D[11]	86.20	77.90	73.03	91.80	88.46	83.62
SECOND[12]	87.44	79.46	73.97	92.01	88.98	83.67
SERCNN[13]	87.74	78.96	74.30	94.11	88.10	83.43
PointPillars	80.51	68.57	61.79	<b>90.74</b>	84.98	79.63
<b>Self-supervised PointPillars</b>	<b>82.54</b>	<b>72.99</b>	<b>67.54</b>	88.92	<b>85.73</b>	<b>80.33</b>

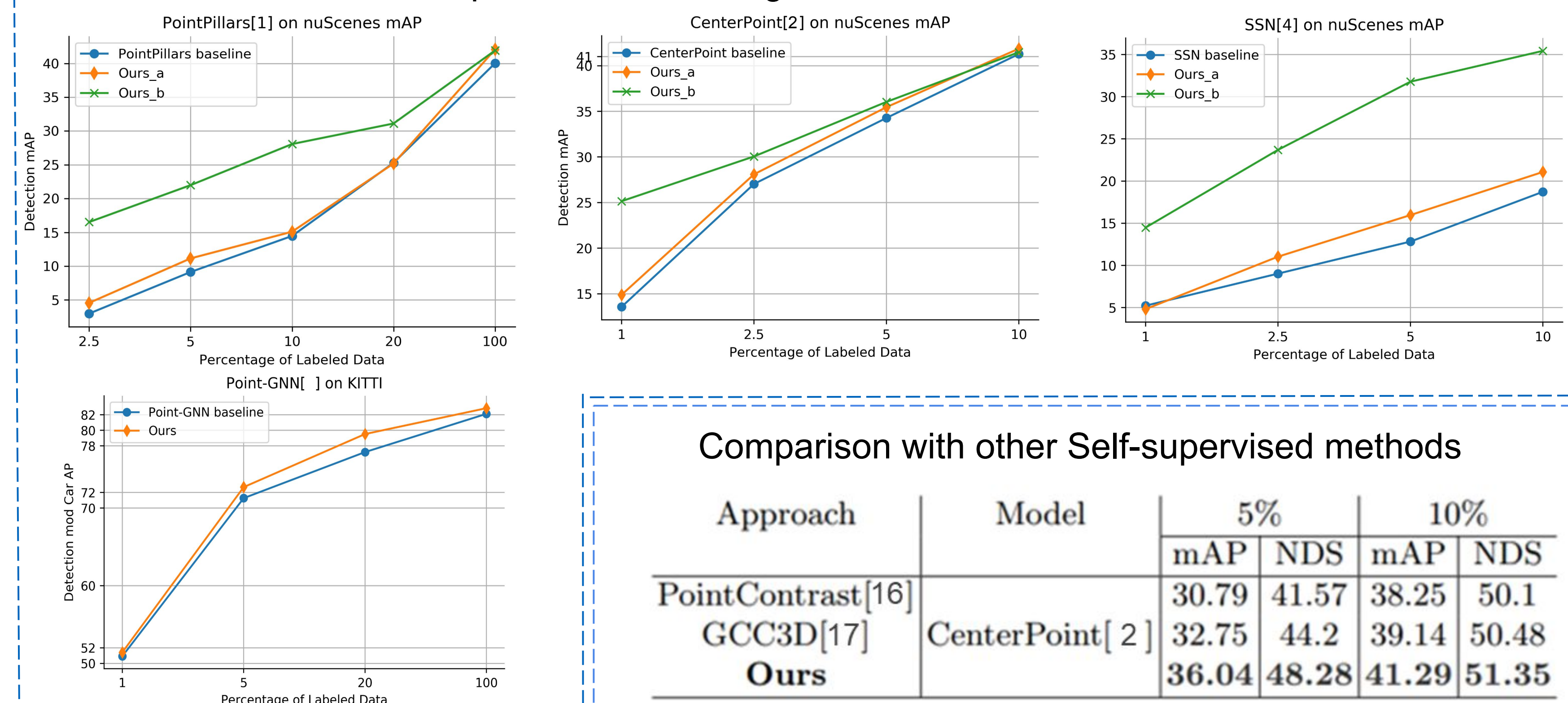
### nuScenes validation set

Method	mAP	NDS	Car	Ped
SECOND[12]	27.12	-	75.53	59.86
PointPillars [1]	40.02	53.29	80.60	72.40
<b>Self-supervised PointPillars</b>	<b>42.06</b>	<b>55.02</b>	<b>81.10</b>	<b>74.50</b>
CenterPoint [2]	49.13	59.73	83.70	77.40
<b>Self-supervised CenterPoint</b>	<b>49.94</b>	<b>60.06</b>	<b>84.10</b>	<b>77.90</b>

### nuScenes test set

Method	mAP	NDS	Car	Ped
PointPillars[1]	30.50	45.30	68.40	59.70
InfoFocus[14]	39.50	39.50	77.90	63.40
PointPillars+[15]	40.10	55.00	76.00	64.00
<b>Self-supervised PointPillars</b>	<b>43.63</b>	<b>56.28</b>	<b>81.00</b>	<b>73.10</b>
CenterPoint[2]	49.54	59.64	83.40	76.10
<b>Self-supervised CenterPoint</b>	<b>51.42</b>	<b>60.92</b>	<b>83.80</b>	<b>77.00</b>

### Supervised fine-tuning with a small set of labeled data



### Comparison with other Self-supervised methods

Approach	Model	5%		10%	
		mAP	NDS	mAP	NDS
PointContrast[16]	CenterPoint[2]	30.79	41.57	38.25	50.1
GCC3D[17]		32.75	44.2	39.14	50.48
<b>Ours</b>		<b>36.04</b>	<b>48.28</b>	<b>41.29</b>	<b>51.35</b>