

Workshop

Towards Robust Perception in Autonomous Driving

Speaker: Lingdong Kong

June 04, 2023

Short Bio

Ph.D. Student, School of Computing, National University of Singapore (<u>https://ldkong.com</u>)

Research Interests:

- 3D Perception
 - LiDAR Segmentation
 - 3D Object Detection
 - Monocular Depth Estimation
- Data-Efficient Learning
 - Semi-Supervised Learning
 - Unsupervised Domain Adaptation
 - Visual Representation Learning
- Robustness
 - Data Robustness
 - Model Robustness



3D-DLAD Workshop



Lingdong Kong

Ph.D. Student, NUS

Short Bio

Ph.D. Student, School of Computing, National University of Singapore (<u>https://ldkong.com</u>)

Industrial Experiences:

- Motional
 - Autonomous vehicle intern, semantic mapping team in Singapore
 - Unsupervised domain adaptation for LiDAR segmentation
- OpenMMLab
 - MMDetction3D codebase
 - Supporting LiDAR segmentation related model & dataset benchmark
 - <u>https://github.com/open-</u> <u>mmlab/mmdetection3d</u>



3D-DLAD Workshop



Lingdong Kong

Ph.D. Student, NUS

Robust Perception in Autonomous Driving

Topic #1: Robustness under Data-Efficient Settings

- Data collection is much easier than annotation
- Learning perception models with semi supervisions
- Goal: Achieve satisfactory perception performance with limited annotations

Topic #2: Robustness under Common Corruptions

- Data corruption and sensor failure are common issues
- Testing perception models with corrupted scenarios
- Goal: Achieve reliable perception performance under out-ofdistribution corruptions







Topic #1

Robustness under Data-Efficient Settings





LaserMix for Semi-Supervised LiDAR Semantic Segmentation

CVPR 2023 (Highlight)





Jiawei Ren



Liang Pan



Ziwei Liu

TL;DR

LaserMix is a data-efficient learning framework designed for LiDAR segmentation that:

- Leverages the spatial prior in driving scenes for data-efficient learning;
- Constructs low-variational areas via laser beam mixing;
- Encourages the model to make confident and consistent predictions before and after mixing;
- Achieves competitive results over full supervision counterparts with 2x to 5x fewer annotations



TL;DR



Autonomous Driving Perception



From left to right:

- LiDAR semantic segmentation
- LiDAR panoptic segmentation
- 3D object detection
- 4D LiDAR panoptic segmentation

Why LiDAR sensors?

- Accurate depth sensing
- Robust at low-light conditions
- Dense perceptions
- •

LiDAR Semantic Segmentation



A. Milioto, et al. "RangeNet++: Fast and accurate LiDAR semantic segmentation," IROS, 2019.

LiDAR Semantic Segmentation



• SemanticKITTI

- Full labels (100%)
- 19 semantic classes
- 100 m x 100 m
- Up to 4.5 hours

J. Behley, et al. "Semantickitti: A dataset for semantic scene understanding of LiDAR sequences," ICCV, 2019.

LiDAR Semantic Segmentation



• SemanticKITTI

- Full labels (100%)
- 19 semantic classes
- 100 m x 100 m
- Up to 4.5 hours

• ScribbleKITTI

- Weak (scribble) labels (8.06%)
- 19 semantic classes
- 100 m x 100 m
- 10 25 min per scan
- 90% time saving

O. Unal, et al. "Scribble-supervised LiDAR semantic segmentation," CVPR, 2022.

Semi-Supervised LiDAR Segmentation



Objective

• We target on the less-explored semisupervised LiDAR semantic segmentation.

Semi-Supervised LiDAR Segmentation



Objective

- We target on the less-explored semisupervised LiDAR semantic segmentation.
- Our goal is to leverage the abundant raw LiDAR scans for training accurate segmentation models.

Semi-Supervised LiDAR Segmentation



Objective

- We target on the less-explored semisupervised LiDAR semantic segmentation.
- Our goal is to leverage the abundant raw LiDAR scans for training accurate segmentation models.
- We propose LaserMix to make advantages of the spatial prior in LiDAR scenes for effective learning with semi supervisions.

Spatial Prior

Class	Туре	Proportion	Distribution	Heatmap
vegetation	static	24.825%		
road	static	22.545%		
sidewalk	static	16.353%		
car	dynamic	4.657%		
traffic-sign	static	0.061%		
motorcycle	dynamic	0.045%		
person	dynamic	0.036%		
bicycle	dynamic	0.018%		

Certain class tends to appear at certain areas around the ego-vehicle!

Overview



(a) Motivation. Semantic scene priors are overt for each category in LiDAR point clouds.

Overview



(a) Motivation. Semantic scene priors are overt for each category in LiDAR point clouds. (b)Generalizability. LaserMix can be added into various popular LiDAR representations.

Overview



(a) Motivation. Semantic scene priors are overt for each category in LiDAR point clouds.
(b)Generalizability. LaserMix can be added into various popular LiDAR representations.
(c)Effectiveness. LaserMix helps to improve both semi- and fully-supervised settings.



Three-Step Procedure

1. Partitioning the captured LiDAR scan into low-variational areas.



Three-Step Procedure

- 1. Partitioning the captured LiDAR scan into low-variational areas.
- 2. Efficiently mixing every area in the LiDAR scan with foreign data.



Three-Step Procedure

- 1. Partitioning the captured LiDAR scan into low-variational areas.
- 2. Efficiently mixing every area in the LiDAR scan with foreign data.
- 3. Encouraging the LiDAR segmentation models to make confident and consistent predictions on the same area in different mixing.



• Inclination:

$$\phi_i = \arctan(\frac{p_i^z}{\sqrt{(p_i^x)^2 + (p_i^y)^2}})$$

• Depth: $\rho_i = \sqrt{(p_i^x)^2 + (p_i^y)^2}$

• Azimuth: $\alpha_i = \arctan \alpha_i$

$$\operatorname{ctan}(\frac{p_i^y}{p_i^x})$$



• Inclination:

$$\phi_i = \arctan(\frac{p_i^z}{\sqrt{(p_i^x)^2 + (p_i^y)^2}})$$

• Depth:
$$\rho_i = \sqrt{(p_i^x)^2 + (p_i^y)^2}$$

• Azimuth: $\alpha_i = \arctan(\frac{p_i^y}{p_i^x})$

Consistency Regularization



Consistency Regularization



LiDAR data and labels strongly correlate with the area A $H(X_{in}, Y_{in}|A)$ is low

- $H(X_{in}, Y_{in}|A)$ is low => $H(Y_{in}|X_{in}, A)$ is low (conditional entropy).
- Let θ be the parameter of the LiDAR segmentation network.
- We would like to solve the following:
 - $E_{\theta}[H_{\theta}(Y_{\text{in}}|X_{\text{in}}, A)] = c$, where *c* is a constant.
 - $\sum_{\theta} P(\theta) = 1$ (sum to one).
- Principle of Maximum Entropy:
 - $P(\theta) \propto \exp(-\lambda H_{\theta}(Y_{\text{in}}|X_{\text{in}}, A))$, where λ is the Lagrange multiplier.

Class	Туре	Proportion	Distribution	Heatmap
vegetation	static	24.825%		
road	static	22.545%		
sidewalk	static	16.353%		
car	dynamic	4.657%		
traffic-sign	static	0.061%		
motorcycle	dynamic	0.045%		
person	dynamic	0.036%		
bicycle	dynamic	0.018%		

Certain class tends to appear at certain areas around the ego-vehicle!

- $P(\theta) \propto \exp(-\lambda H_{\theta}(Y_{\text{in}}|X_{\text{in}},A)) \implies \text{spatial prior.}$
- Compute the empirical entropy:
- $\widehat{H}_{\theta}(Y_{\mathrm{in}}|X_{\mathrm{in}},A) = E_{X_{\mathrm{in}},Y_{\mathrm{in}},A}[P_{\theta}(Y_{\mathrm{in}}|X_{\mathrm{in}},A)\log P_{\theta}(Y_{\mathrm{in}}|X_{\mathrm{in}},A)].$
- $P_{\theta}(Y_{\text{in}}|X_{\text{in}}, A)$ means predicting the labels by the data inside an area A.
- The segmentation network predicts from full data. Therefore, we need X_{out} to complement the remaining area outside *A* and marginalize X_{out} .
- $P_{\theta}(y_{\text{in}}|x_{\text{in}}, a) = \frac{1}{|X_{\text{out}}|} \sum_{x_{\text{out}} \in X_{\text{out}}} P_{\theta}(y_{\text{in}}|x_{\text{in}}, a, x_{\text{out}}).$

- $P(\theta) \propto \exp(-\lambda H_{\theta}(Y_{\text{in}}|X_{\text{in}},A)) \implies \text{spatial prior.}$
- $P_{\theta}(y_{\text{in}}|x_{\text{in}}, a) = E_{X_{\text{out}}}[P_{\theta}(y_{\text{in}}|x_{\text{in}}, a, x_{\text{out}})] \implies \text{marginalization}.$
- We maximize the following posterior:
 - $C(\theta) = -\lambda \hat{E}_{x_{\mathrm{in}} \in X_{\mathrm{in}}, y_{\mathrm{in}} \in Y_{\mathrm{in}}, a \in A}[H],$
 - $H = P_{\theta}(y_{\text{in}}|x_{\text{in}}, a) \cdot \log P_{\theta}(y_{\text{in}}|x_{\text{in}}, a),$
 - *H* is minimized only when $P_{\theta}(y_{in}|x_{in}, a, x_{out})$ is certain and consistent to x_{out} .

- $H = \frac{1}{|X_{\text{out}}|} \sum_{x_{\text{out}} \in X_{\text{out}}} P_{\theta}(y_{\text{in}} | x_{\text{in}}, a, x_{\text{out}}) \log P_{\theta}(y_{\text{in}} | x_{\text{in}}, a, x_{\text{out}}).$
- H = 0 only when $P_{\theta}(y_{in}|x_{in}, a, x_{out})$ is certain and consistent to x_{out} .
- For every selected area and the data *inside* that area, a LiDAR segmentation network should make certain and consistent predictions regardless of the data *outside* the area.
- Directly compute $E_{y_{in} \in Y_{in}}[H]$ is infeasible / intractable, since $|y_{in}| = C^{H_{in} \times W_{in}}$ is too large.
- Instead, we use the pseudo-label to make sure that $P_{\theta}(y_{in}|x_{in}, a, x_{out})$ is certain and consistent.

Experimental Settings

	nuScenes [15]	SemanticKITTI [16]	ScribbleKITTI [4]				
Vis.							
#Class	16	19	19				
#Train	29130	19130	19130				
#Val	6019	4071	4071				
Res. (RV)	32×1920	64×2048	64×2048				
Res. (voxel)	[240, 180, 20]	[240, 180, 20]	[240, 180, 20]				
#Beam	32	64	64				
$[\phi_{ m up},\phi_{ m low}]$	$[10^{\circ}, -30^{\circ}]$	$[3^{\circ}, -25^{\circ}]$	$[3^{\circ}, -25^{\circ}]$				
$[p_{\max}^x, p_{\min}^x]$	[50m, -50m]	[50m, -50m]	[50m, -50m]				
$[p_{\max}^y, p_{\min}^y]$	[50m, -50m]	[50m, -50m]	[50m, -50m]				
$[p_{\max}^{\boldsymbol{z}}, p_{\min}^{\boldsymbol{z}}]$	[3m, -5m]	[2m, -4m]	[2m, -4m]				
#Label	100%	100%	8.06%				
Intensity							
Range							
Semantics							

High-res LiDAR:

- SemanticKITTI
- Denser scenes

Low-res LiDAR:

- nuScenes
- Sparser scenes

Weak supervision:

- ScribbleKITTI
- Sparse labels

Experimental Settings

- Range View
 - Backbone: FIDNet [IROS'21]
 - # Param: 6.05M
 - 6 x 32 x 1920 (nuScenes)
 - 6 x 64 x 2048 (SemanticKITTI/ScribbleKITTI)
- Voxel
 - Backbone: Cylinder3D [CVPR'21]
 - # Param: 28.13M
 - [240, 180, 20]

Y. Zhao, et al. "FIDNet: LiDAR point cloud semantic segmentation with fully interpolation decoding," IROS, 2021. X. Zhu, et al. "Cylindrical and asymmetrical 3D convolution networks for LiDAR segmentation," CVPR, 2021.

Experimental Settings

- Range View
 - Backbone: FIDNet [IROS'21]
 - # Param: 6.05M
 - 6 x 32 x 1920 (nuScenes)
 - 6 x 64 x 2048 (SemanticKITTI/ScribbleKITTI)
- Voxel
 - Backbone: Cylinder3D [CVPR'21]
 - # Param: 28.13M
 - [240, 180, 20]
- Data Split
 - 1%, 10%, 20%, 50% (labeled)
 - Random sampling
 - Assume the remaining ones are unlabeled

Y. Zhao, et al. "FIDNet: LiDAR point cloud semantic segmentation with fully interpolation decoding," IROS, 2021. X. Zhu, et al. "Cylindrical and asymmetrical 3D convolution networks for LiDAR segmentation," CVPR, 2021.

Experimental Results

Dopr	Mathad	nuScenes [15]				SemanticKITTI [16]				ScribbleKITTI [4]				
Kepi.	Method	1%	10%	20%	50%	1%	10%	20%	50%	1%	10%	20%	50%	
Range View	Suponly	38.3	57.5	62.7	67.6	36.2	52.2	55.9	57.2	33.1	47.7	49.9	52.5	
	MeanTeacher [26] CBST [30] CutMix-Seg [29] CPS [13]	$\begin{array}{c} 42.1 \\ 40.9 \\ 43.8 \\ 40.7 \end{array}$	$60.4 \\ 60.5 \\ 63.9 \\ 60.8$	$65.4 \\ 64.3 \\ 64.8 \\ 64.9$	$69.4 \\ 69.3 \\ 69.8 \\ 68.0$	37.5 39.9 37.4 36.5	$53.1 \\ 53.4 \\ 54.3 \\ 52.3$	$56.1 \\ 56.1 \\ 56.6 \\ 56.3$	57.4 56.9 57.6 57.4	$34.2 \\ 35.7 \\ 36.7 \\ 33.7$	$\begin{array}{c} 49.8 \\ 50.7 \\ 50.7 \\ 50.0 \end{array}$	51.6 52.7 52.9 52.8	$53.3 \\ 54.6 \\ 54.3 \\ 54.6$	
	$\begin{array}{c} \textbf{LaserMix (Ours)} \\ \Delta \uparrow \end{array}$	$\begin{array}{c} \textbf{49.5} \\ \textbf{+11.2} \end{array}$	$\begin{array}{c} 68.2 \\ \mathbf{+10.7} \end{array}$	70.6 + 7.9	73.0 +5.4	$\begin{array}{c} \textbf{43.4} \\ \textbf{+7.2} \end{array}$	58.8 + 6.6	$59.4 \\ +3.5$	$\begin{array}{c} 61.4 \\ \mathbf{+4.2} \end{array}$	38.3 + 5.2	$54.4 \\ +6.7$	$\begin{array}{c} 55.6 \\ \mathbf{+5.7} \end{array}$	$\begin{array}{c} 58.7 \\ \mathbf{+6.2} \end{array}$	
_	Suponly	50.9	65.9	66.6	71.2	45.4	56.1	57.8	58.7	39.2	48.0	52.1	53.8	
Voxel	MeanTeacher [26] CBST [30] CPS [13]	$51.6 \\ 53.0 \\ 52.9$	$\begin{array}{c} 66.0 \\ 66.5 \\ 66.3 \end{array}$	$67.1 \\ 69.6 \\ 70.0$	71.7 71.6 72.5	$\begin{array}{c} 45.4 \\ 48.8 \\ 46.7 \end{array}$	$57.1 \\ 58.3 \\ 58.7$	$59.2 \\ 59.4 \\ 59.6$	$ \begin{array}{r} 60.0 \\ 59.7 \\ 60.5 \end{array} $	$\begin{array}{c} 41.0 \\ 41.5 \\ 41.4 \end{array}$	$50.1 \\ 50.6 \\ 51.8$	52.8 53.3 53.9	$53.9 \\ 54.5 \\ 54.8$	
	LaserMix (Ours) $\Delta \uparrow$	55.3 $+4.4$	69.9 +4.0	71.8 + 5.2	73.2 +2.0	50.6 +5.2	60.0 +3.9	61.9 +4.1	62.3 +3.6	44.2 + 5.0	53.7 +5.7	55.1 + 3.0	56.8 +3.0	

A. Tarvainen and H. Valpola. "Mean teachers are better role models: Weight-averaged consistency targets improve semisupervised deep learning results," NeurIPS, 2017.

G. French, et al. "Semi-supervised semantic segmentation needs strong, high-dimensional perturbations," BMVC, 2020. Y. Zou, et al. "Domain adaptation for semantic segmentation via class-balanced self-training," ECCV, 2018.

X. Chen, et al. "Semi-supervised semantic segmentation with cross pseudo supervision," CVPR, 2021.

Experimental Results

road	sidewalk	building	wall	fence
pole	traffic light	traffic sign	vegetation	terrain
		The second s		
alar		rider		Amuak
SKy	person	nder	car	truck
		and the second second		
bus	train	motorcycle	bicycle	
			C C C C C C C C C C C C C C C C C C C	

Method	1/16	1/8	1/4	1/2
MeanTeacher [26]	66.1	71.2	74.4	76.3
w/ Ours ∆↑	68.7 + 2.6	$72.3 \\ +1.1$	75.7 + 1.3	76.8 + 0.5
CCT [11] GCT [12] CPS [13]	$\begin{array}{c} 66.4 \\ 65.8 \\ 69.8 \end{array}$	72.5 71.3 74.4	75.7 75.3 76.9	$76.8 \\ 77.1 \\ 78.6$
CPS-CutMix [13]	74.5	76.6	77.8	78.8
w/ Ours ∆↑	75.5 +1.0	77.1 + 0.5	78.3 + 0.5	79.1 +0.3

Also has spatial priors in scenes!

Y. Ouali, et al. "Semi-supervised semantic segmentation with cross-consistency training," CVPR, 2020. Z. Ke, et al. "Guided collaborative training for pixel-wise semi-supervised learning," ECCV, 2020.

#	\mathcal{L}_{mt}	\mathcal{L}_{mix}	SS	TS	1%	10%	20%	50%
(1)	\checkmark				42.1	60.4	65.4	69.4
(2)	\checkmark	\checkmark	✓ ✓		$\begin{array}{c} 45.6\\ 47.0\end{array}$	$\begin{array}{c} 64.3 \\ 65.5 \end{array}$	$\begin{array}{c} 67.8 \\ 69.5 \end{array}$	$71.6 \\ 72.0$
(3)	\checkmark	\checkmark		✓ ✓	$\begin{array}{c} 46.0 \\ 49.5 \end{array}$	$\begin{array}{c} 64.1 \\ 68.2 \end{array}$	$69.5 \\ 70.6$	$72.3 \\ 73.0$

- (1) Results of MeanTeacher.
- (2) Results of LaserMix w/ student supervisions; much better than the counterpart.
- (3) Results of LaserMix w/ teacher supervisions; much better than the counterpart.

(a) Comparisons among different mixing techniques.

A. Nekrasov, et al. "Mix3D: Out-of-context data augmentation for 3D scenes," 3DV, 2021.

S. Yun, et al. "Cutmix: Regularization strategy to train strong classifiers with localizable features," ICCV, 2019 T. DeVries and G. W. Taylor. "Improved regularization of convolutional neural networks with cutout," arXiv, 2017 H. Zhang, et al. "Mixup: Beyond empirical risk minimization," ICLR, 2018.

(a) Comparisons among different mixing techniques. (b) EMA. (c) Confidence threshold.

A. Nekrasov, et al. "Mix3D: Out-of-context data augmentation for 3D scenes," 3DV, 2021.

S. Yun, et al. "Cutmix: Regularization strategy to train strong classifiers with localizable features," ICCV, 2019

T. DeVries and G. W. Taylor. "Improved regularization of convolutional neural networks with cutout," arXiv, 2017 H. Zhang, et al. "Mixup: Beyond empirical risk minimization," ICLR, 2018.

• Inclination:

$$\phi_i = \arctan(\frac{p_i^z}{\sqrt{(p_i^x)^2 + (p_i^y)^2}})$$

• Depth:
$$\rho_i = \sqrt{(p_i^x)^2 + (p_i^y)^2}$$

• Azimuth: $\alpha_i =$

$$\arctan(\frac{p_i^y}{p_i^x})$$

Public Resources

- Paper: https://arxiv.org/abs/2207.00026
- Code: <u>https://github.com/ldkong1205/LaserMix</u>
- **Tutorial:** <u>https://zhuanlan.zhihu.com/p/528689803</u>
- Project Page: https://ldkong.com/LaserMix

Topic #2

Robustness under Common Corruptions

Benchmarking 3D Perception Robustness to Common Corruptions and Sensor Failure

Lingdong Kong^{1,2,*}, Youquan Liu^{1,3,*}, Xin Li^{1,4,*}, Runnan Chen^{1,5}, Wenwei Zhang^{1,6} Jiawei Ren⁶, Liang Pan⁶, Kai Chen¹, Ziwei Liu⁶

> ¹Shanghai Al Laboratory ²NUS ³Hochschule Bremerhaven ⁴ECNU ⁵HKU ⁶S-Lab, NTU

Perception Environment

*Image credit: <u>https://zod.zenseact.com</u>

Robustness in RGB Images

D. Hendrycks, et al. "Benchmarking Neural Network Robustness to Common Corruptions and Perturbations," ICLR, 2019

Robustness in RGB Images

Corruption Error:

$$\operatorname{CE}_{c}^{f} = \left(\sum_{s=1}^{5} E_{s,c}^{f}\right) / \left(\sum_{s=1}^{5} E_{s,c}^{\operatorname{AlexNet}}\right)$$

D. Hendrycks, et al. "Benchmarking Neural Network Robustness to Common Corruptions and Perturbations," ICLR, 2019

Robustness in RGB Images

			Noise			Blur				Weather				Digital			
Network	Error	mCE	Gauss.	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
AlexNet	43.5	100.0	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
SqueezeNet	41.8	104.4	107	106	105	100	103	101	100	101	103	97	97	98	106	109	134
VGG-11	31.0	93.5	97	97	100	92	99	93	91	92	91	84	75	86	97	107	100
VGG-19	27.6	88.9	89	91	95	89	98	90	90	89	86	75	68	80	97	102	94
VGG-19+BN	25.8	81.6	82	83	88	82	94	84	86	80	78	69	61	74	94	85	83
ResNet-18	30.2	84.7	87	88	91	84	91	87	89	86	84	78	69	78	90	80	85
ResNet-50	23.9	76.7	80	82	83	75	89	78	80	78	75	66	57	71	85	77	77

D. Hendrycks, et al. "Benchmarking Neural Network Robustness to Common Corruptions and Perturbations," ICLR, 2019

- Data in the point cloud format are used in safety-critical applications, such as autonomous driving and robot navigation.
- However, this data format often suffers from severe Out-of-Distribution (OoD) corruptions in real-world deployment.

Corruptions are severe and OoD e.g., occlusion, sensory noise Applications are safety-critical e.g., autonomous driving

- **Question:** Are point cloud classifiers getting more robust?
- Answer: No. Although the accuracy indicator on ModelNet40 gradually saturates, the robustness is at the risk of getting worse, due to the lack of a standard test suite.

- Three levels of corruption sources: Object, sensor, and processing.
- Nine potential corruption types.
- Simplified into a combination of seven atomic corruptions for a more controllable empirical analysis.

- PointCloud-C is the first competition that targets the robustness of point cloud understanding under corruptions.
- The benchmark result suggests that point cloud classifiers are at the risk of getting less robust, thus highlighting the importance of proposing new designs to improve the robustness.

Robo3D

TL;DR

- We introduce Robo3D, the first systematicallydesigned robustness evaluation suite for LiDAR-based 3D perception under corruptions and sensor failure
- We benchmark 34 perception models for LiDAR-based semantic segmentation and object detection tasks, on their robustness against corruptions.
- Based on our observations, we draw in-depth discussions on the receipt of designing robust and reliable 3D perception models.

Robo3D: Taxonomy

*More examples at: <u>https://ldkong.com/Robo3D</u>

Robo3D: Example

Robo3D: Representation

Representation:

- 2D: range view, bird's eye view
- **3D:** cubic voxel, cylinder voxel

Operator:

- **3D:** Conv3d, SparseConv, etc.
- 2D: Conv2d, Linear, etc.
- **1D:** Conv1d, Linear, etc.

M. Uecker, et al. "Analyzing deep learning representations of point clouds for real-time in-vehicle LiDAR perception," arXiv, 2022.

Robo3D: Statistics

Corruption Type:

• Include 8 types, each with 3 severity levels

Dataset (6 different sets):

- LiDAR Semantic Segmentation: ¹SemanticKITTI-C, ²nuScenes-C (Seg3D), ³WOD-C (Seg3D)
- **3D Object Detection:** ⁴KITTI-C, ⁵nuScenes-C (Det3D), ⁶WOD-C (Det3D)

Model & Algorithm (34 perception models):

- LiDAR Semantic Segmentation: 22 segmentors
- **3D Object Detection:** 12 detectors
- Data Augmentation: 3 augmentation techniques

Robo3D: Metrics

Task-Specific Accuracy (Acc):

- LiDAR Semantic Segmentation: mean IoU (mIoU)
- **3D Object Detection:** mean AP (mAP), nuScenes Detection Score (NDS)

Robustness Metrics:

• Mean Corruption Error (mCE):

$$\mathrm{CE}_{i} = \frac{\sum_{l=1}^{3} (1 - \mathrm{Acc}_{i,l})}{\sum_{l=1}^{3} (1 - \mathrm{Acc}_{i,l}^{\mathrm{baseline}})} \;, \quad \mathrm{mCE} = \frac{1}{I}$$

• Mean Resilience Rate (mRR):

$$\mathbf{RR}_{i} = \frac{\sum_{l=1}^{3} \operatorname{Acc}_{i,l}}{3 \times \operatorname{Acc}_{\text{clean}}} , \quad \mathbf{mRR} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{RR}_{i}$$

Robo3D: Benchmarking Result

*More results and analysis at: <u>https://github.com/ldkong1205/Robo3D</u>

Robo3D: Key Observation

- 1. Existing 3D detectors and segmentors are **vulnerable** to real-world corruptions.
- 2. Models trained with LiDAR data from different sources (sensor setups) exhibit **inconsistent sensitivities** to each corruption type.
- 3. Representing the LiDAR data as raw **points**, sparse **voxel**, or the **fusion** of them tend to yield better robustness.

Robo3D: Key Observation

- 4. The 3D detectors and segmentors show different sensitivities to corruption scenarios.
- 5. The recent **out-of-context augmentation techniques** improve 3D robustness by large margins; the flexible rasterization strategies help learn more robust features.

(a) Voxel Size on *SemanticKITTI-C* (Seg3D)

(b) Augmentation on SemanticKITTI-C (Seg3D)

(d) Augmentation on WOD-C (Det3D)

Robo3D: Qualitative Assessment

Robo3D: Qualitative Assessment

Thank you for your attention!

