ENHANCED LOW-RESOLUTION LIDAR-CAMERA CALIBRATION VIA DEPTH INTERPOLATION AND SUPERVISED CONTRASTIVE LEARNING

 $\begin{array}{ccc} Zhikang \ Zhang^{1,\star} & Zifan \ Yu^{1,\star} & Suya \ You^2 & Raghuveer \ Rao^2 \\ Sanjeev \ Agarwal^3 & Fengbo \ Ren^1 \end{array}$

¹ Arizona State University, Tempe, AZ,
² U.S. DEVCOM Army Research Laboratory, Adelphi, MD
³ U.S. DEVCOM Army C5ISR Center, Fort Belvoir, VA

ABSTRACT

Motivated by the increasing application of low-resolution Li-DAR, we target the problem of low-resolution LiDAR-camera calibration in this work. The main challenges are two-fold: sparsity and noise in point clouds. To address the problem, we propose to apply depth interpolation to increase the point density and supervised contrastive learning to learn noiseresistant features. The experiments on RELLIS-3D demonstrate that our approach achieves an average mean absolute rotation/translation errors of $0.15cm/0.33^{\circ}$ on 32-channel Li-DAR point cloud data, which significantly outperforms all reference methods.

Index Terms— low-resolution point cloud, LiDARcamera calibration, supervised contrastive learning, image interpolation

1. INTRODUCTION

Driven by the development of LiDAR technology, the application scenarios of low-resolution LiDAR devices have expanded significantly in recent years, such as autonomous driving [1], geoscience [2], remote sensing [3], and mobile robotics [4]. To acquire an accurate and informative perception of scanned targets or environments, LiDAR devices are often fused with cameras to utilize the rich information in images. The basis of LiDAR-camera fusion is extrinsic calibration, that is, estimating a relatively rigid body transformation from LiDAR coordinates to camera coordinates, which has been long studied. Conventional calibration methods [5, 6, 7, 8, 9, 10] are mostly based on explicit targets in a scene, handcrafted features, or labels of data to build correspondences between point clouds and images; thus, they are often limited by laborious human interventions and/or applied environments. Recognizing these limitations, recently, deeplearning-based calibration approaches [11, 12, 13, 14, 15] are proposed, which automatically learn features from sensed

data and perform calibration in an end-to-end manner. Because feature learning depends heavily on the quality of data, most studies lay the foundation upon highly accurate and noiseless point clouds(Fig. 2) sensed by high-resolution Li-DAR, and thus suffers from large performance degradation in low-resolution LiDAR scenarios.

In this work, we use the state-of-the-art method LCC-Net [15] as the backbone to apply two effective techniques to enhance the performance of low-resolution LiDAR-camera calibration. We first identify two major challenges to lowresolution LiDAR-camera calibration: sparsity and noise, as shown in Fig. 2. To tackle sparsity, we apply depth interpolation to increase the density of the point cloud, which inevitably introduces more noise to the point cloud. Then, to tackle the inherent and introduced noise, we apply supervised contrastive loss on the backbone to learn noiseresistant features for calibration. Extensive experiments on 32-channel point clouds of RELLIS-3D [16] demonstrate that our approach reduces the average translation and rotation errors of the original LCCNet by 87%(2.66cm to 0.33cm)and $28\%(0.21^{\circ} \text{ to } 0.15^{\circ})$, respectively, in the miscalibration range of 150cm and 20° . Compared with Regnet [13] and CalibDNN [14], the calibration errors are reduced by at least 10x times, which provides strong evidence that our approach is highly effective in addressing the low-resolution LiDARcamera calibration problem.

Our contributions are summarized as follows:

1. We propose two effective techniques to enhance deeplearning-based, automatic targetless LiDAR-camera calibration in the low-resolution LiDAR scenario. To the best of our knowledge, this is the first work that targets the low-resolution LiDAR-camera calibration problem.

2. We demonstrate that supervised contrastive loss can be applied to learn noise-resistant features for LiDAR-camera calibration.

3. Our approach achieves state-of-the-art performance for low-resolution LiDAR-camera calibration, which sets a strong baseline for this task.

^{*:} equal contribution



Fig. 1. The entire training pipeline of our approach. Given the input pair of RGB image and miscalibrated point cloud, image features are extracted with two CNNs from the RGB image and the max-pooled depth image. Subsequently, a correlation layer [17] constructs a cost volume from extracted image features to generate rotation features and translation features for prediction of rotation parameters and translation parameters. Three losses are applied in training: calibration loss minimizes the error of predicted parameters. Cloud distance loss minimizes the error of calibrated point cloud. Supervised contrastive loss enhances the learned rotation features and translation features to be noise-resistant.

2. METHODOLOGY

We use a state-of-the-art method LCCNet as the backbone of our approach, as shown in Fig. 1. In the inference stage, LCCNet runs in two modes: single-stage and multi-stage. In single-stage mode, a single model is trained for a single miscalibration range. In multi-stage mode, multiple models are trained separately for different miscalibration ranges, and the input is calibrated sequentially from higher ranges to lower ranges with corresponding models.

We first identify two main problems for low-resolution LiDAR-camera calibration resulting from point clouds: sparsity and noise. Fig. 2 shows a visual comparison of point clouds sensed by a 32-channel LiDAR and a 64-channel Li-DAR. Higher sparsity and more noise in low-resolution point clouds lead to more difficulties in calibration since the RGB images and depth images are less correlated, consequently making the constructed cost volume less informative.

Depth interpolation. To address the sparsity problem, we propose to apply interpolation to depth images before feature extraction. There are a large variety of image interpolation methods, and we choose to use max-pooling, which shows the highest calibration accuracy in Section 3. Given a depth image with h height and w width and output size of \hat{h} height and \hat{w} width, the stride and kernel size are set to $h/\hat{h}, w/\hat{w}$ and $h - (\hat{h} - 1) * (h/\hat{h}), w - (\hat{w} - 1) * w/\hat{w}$, respectively, following the widely used adaptive max-pooling design [18].

Supervised contrastive learning. Depth interpolation inevitably introduces more noise to the point cloud because a large number of fake 3d points are added. To learn noise-resistant features, we hypothesize that learned features should satisfy three conditions: 1. Rotation features(Fig. 1, red block) only retain information related to rotation parameters. 2. Translation features(Fig. 1, purple block) only retain information related to translation and translation features do not retain data-dependent(either image or point cloud)information to avoid over-fitting.

Following the three conditions, we propose to add supervised contrastive loss(SCL) [19] in addition to original calibration loss and cloud distance loss(both defined in [15]). SCL is defined as

$$L^{sup} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_i \cdot z_p/\tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a/\tau)}$$
(1)

where $P(i) \equiv \{p \in A(i) : \tilde{y}_p = \tilde{y}_i\}$ is the set of indices of all positives samples distinct from *i* within the mini-batch, |P(i)| is its cardinality. z_i, z_a and z_p are features of the given, negative and positive samples, respectively. τ is a scalar temperature parameter as defined in [19]. Despite its complicated mathematical form, SCL can be implemented as a function that takes in a batch of features and the same number of numerical labels while outputting a singular loss value, i.e.,

$$loss = SCL([f_1, \cdots, f_b], [l_1, \cdots, l_b])$$
(2)

where f_i is the feature, l_i is the corresponding label, and b is the batch size, as implemented in [20]. Then in the training process, features with the same labels are pulled together while features with different labels are pushed apart.

Composed inputs			Features		Labels		
Image	PC	RO	TR	RO	TR	RO	TR
I_1	P_1	R_1	T_1	R_1^f	T_1^f	1	1
I_1	P_1	R_1	T_2	R_2^f	T_2^f	1	2
I_1	P_1	R_2	T_1	R_{b+1}^f	T_{b+1}^f	2	1
I_2	P_2	R_1	T_1	R^f_{b*b+1}	T^f_{b*b+1}	1	1
I_b	P_b	R_b	T_b	R^f_{b*b*b}	T^f_{b*b*b}	b	b

Table 1. The composed input batch and assigned labels for supervised contrastive learning. Original batch size: *b*. PC: point cloud. RO: rotation. TR: translation. I_k : k^{th} image in the batch. P_k : k^{th} point cloud. R_k : k^{th} random rotation parameters. T_k : k^{th} random translation parameters. R_k^f : k^{th} generated rotation feature. T_k^f : k^{th} generated translation feature.

To adapt SCL for enhanced feature learning, we generate features and labels in the following strategy for each batch: with a batch size of b, given a batch of training samples containing 4-tuples (I_k, P_k, R_k, T_k) of RGB image, point cloud, random rotation, and random translation, $1 \leq k \leq b$, we first compose a new batch as inputs consisting of all possible 4-tuple combinations of I_k, P_k, R_k, T_k while always keeping I_k and P_k paired, as shown in Table 1. As such, the new batch size is b^3 . Then we assign two groups of labels to generated rotation features R_k^f and translation features T_k^f , respectively. For rotation features, the same labels are assigned if and only if they have the same rotation parameters. For translation features, their labels are assigned in a similar manner to correspond to translation parameters. Then, two SCL functions are used to take in rotation and translation features and the corresponding labels, respectively. Through this process, three conditions can be satisfied in the following sense: 1. Rotation features R^f are pushed closer if and only if their rotation parameters R are the same. 2. Translation features T^{f} are pushed closer if and only if their translation parameters T are the same. 3. Rotation features and translation features are less affected by solely changing the input images and point cloud pairs without changing the calibration parameters.

3. EXPERIMENTS

We use RELLIS-3D dataset for evaluation. RELLIS-3D contains point clouds sensed by a 32-channel LiDAR and 64-channel LiDAR in off-road environments. 32-channel point clouds are treated as low-resolution data. The split of

the dataset follows the official split, with 7800 training samples, 2413 validation samples, and 3343 testing samples. The miscalibration ranges are set to $150cm/20^{\circ}$, $100cm/10^{\circ}$, $50cm/5^{\circ}$, $20cm/2^{\circ}$, and $10cm/1^{\circ}$ following the original paper. The same for the training settings of LCCNet. The evaluation metrics are mean absolute translation error (x, y, z), mean absolute rotation errors (roll, pitch, yaw), averaged translation error (x + y + z)/3 and averaged rotation error (roll + pitch + yaw)/3.



Fig. 2. The visual comparison of point clouds from highresolution LiDAR and low-resolution LiDAR. The point cloud is projected to the image plane and plotted as an overlay layer. The data is from RELLIS-3D. Left: 64-channel LiDAR. Right: 32-channel LiDAR. Blue box: sparse region. Red box: noisy region.

Quantify calibration performance degradation. We first train two multi-stage LCCNet on point clouds of 32 channels and 64 channels, respectively. The experiment results are shown in Table 2. The average translation error and rotation error increase two to four times on 32-channel data compared with the same model trained on 64-channel data.

Channel	X	Y	Z	Roll	Pitch	Yaw
64	0.66	0.71	0.25	0.12	0.14	0.09
32	2.6	2.6	2.78	0.22	0.16	0.27

Table 2. Quantified performance degradation on low-resolution(32-channel) LiDAR. Unit: cm or °. Network: multi-stage LCCNet. Average degradation: 2.12cm and 0.1° .

Depth interpolation. We compare max-pooling against three candidate image interpolation methods: average-pooling, linear interpolation, and nearest neighbor interpolation, as well as the original LCCNet approach. The single-stage LC-CNet is trained at the miscalibration range of $150cm/20^{\circ}$. The experiment results are shown in Table 3. The model trained with max-pooling achieves the lowest calibration errors among all interpolate depth images in the following experiments. Be noted the rotation errors of all four interpolation methods are slightly higher than the original model, which can be attributed to the fake points added to the depth image through interpolation.

	original	linear	average pooling	max pooling	nearest neighbour
TR	54.13	71.78	43.00	40.84	57.99
RO	1.02	3.95	4.00	3.34	4.23

Table 3. Comparison against various image interpolation methods. TR: averaged translation error. RO: averaged rotation error. Network: single-stage LCCNet at the miscalibration range of $150 cm/20^{\circ}$.

Supervised contrastive learning. We validate the effectiveness of SCL by training single-stage model on all five different miscalibration ranges. As experiment results in Table 4 show, with max-pooling applied, the averaged translation and rotation errors at most ranges are significantly reduced compared to the original approach. In addition, with SCL being applied, the calibration error is further reduced by $3.95cm/0.25^{\circ}$ on average. The experiment results validate our hypothesis that SCL can enhance feature learning of calibration.

range	original	MP	MPSCL
150/20	54.12/1.02	40.84/3.34	26.86/2.61
100/10	17.13/0.64	11.82/0.91	9.52/0.61
50/5	8.78/0.50	5.08/0.36	3.17/0.25
20/2	4.05/0.41	3.17/0.23	1.75/0.18
10/1	2.11/0.21	0.98/0.17	0.84/0.12

Table 4. Calibration performance comparison at different miscalibration ranges. Original: the original model. MP: with max-pooling applied. MPSCL: with both max-pooling and SCL applied. Unit: $cm/^{\circ}$

Comparison against reference methods. We further evaluate the performance of our approach(multi-stage model trained at ranges of $150cm/20^\circ$, $100cm/10^\circ$, $50cm/5^\circ$, $20cm/2^{\circ}$, and $10cm/1^{\circ}$ with max-pooling and SCL applied, denoted as MPSCL) by comparing it against multiple reference methods. The miscalibration range for evaluation is set to $150 cm/20^{\circ}$. The reference methods are original multi-stage LCCNet, Regnet, and CalibDNN(For CalibDNN, the range is set to $20cm/10^{\circ}$ to be consistent with original work). The experiment results are shown in Table 5. MPSCL achieves the highest performance on all evaluation metrics. Compared with LCCNet, the averaged translation error and rotation error are reduced by 87%(2.66cm to 0.33cm) and $28\%(0.21^{\circ} \text{ to } 0.15^{\circ})$, respectively. Compared with the other two reference methods, the calibration errors of MPSCL are at least one order of magnitude lower, which is strong evidence that MPSCL can effectively perform LiDAR-camera calibration in low-resolution LiDAR scenarios.

Performance on subsampled point clouds. To our knowledge, there is no public dataset for the LiDAR-camera calibration problem with a resolution below 32 channels.

	RegNet	CalibDNN	LCCNet	MPSCL
Х	58.70	10.43	2.60	0.23
Y	32.30	14.59	2.60	0.45
Ζ	50.73	10.77	2.78	0.30
Average	47.24	11.93	2.66	0.33
Roll	4.00	1.11	0.22	0.14
Pitch	8.23	4.15	0.16	0.13
Yaw	5.42	2.05	0.27	0.17
Average	5.88	2.44	0.21	0.15

Table 5. The performance evaluation of MPSCL against reference methods. To be noted, the miscalibration range of CalibDNN is $20cm/10^{\circ}$.

To evaluate the performance of our approach in extreme cases, we perform subsampling on point clouds to simulate lower-resolution LiDAR scenarios. We test with three subsampling rates: 2, 4, and 8. The point cloud is uniformly subsampled. The miscalibration range is set to $150 cm/20^{\circ}$. The experiment results are shown in Table 6. MPSCL again shows significantly higher performance than all reference methods. Even at a subsampling rate of 8, the averaged translation/rotation errors are only $4.28 cm/1.24^{\circ}$. Compared with LCCNet, the average reduction in translation/rotation errors is $19.43 cm/0.03^{\circ}$. Compared with RegNet and CalibDNN, the average reduction is $51.09cm/5.25^{\circ}$ and $7.34cm/0.91^{\circ}$, This is further evidence that MPSCL can respectively. well address the LiDAR-camera calibration problem in lowresolution LiDAR scenarios.

Rate	RegNet	CalibDNN	LCCNet	MPSCL
2	49.71/5.96	11.45/1.55	4.23/0.57	1.21/0.04
4	54.19/6.21	10.31/1.73	22.73/0.51	2.65/1.91
8	57.80/6.76	8.40/2.65	39.45/2.20	4.28/1.24

Table 6. Performance evaluation on subsampled point clouds. To be noted, the miscalibration range of CalibDNN is $20cm/10^{\circ}$.

4. CONCLUSION

In this work, we propose an effective approach for lowresolutioin LiDAR-camera calibration. We first identify two main challenges resulting from low-resolution data: sparsity and noise. Then, we take LCCNet as the backbone to apply max pooling to interpolate depth images and supervised contrastive loss to tackle noises, which eventually leads to a highly effective approach for low-resolution LiDARcamera calibration. The extensive experiments on RELLIS-3D against reference methods demonstrate that our approach can achieve superior performance in calibration, even for extreme cases.

5. REFERENCES

- [1] Ying Li, Lingfei Ma, Zilong Zhong, Fei Liu, Dongpu Cao, Jonathan Li, and Michael A. Chapman, "Deep learning for lidar point clouds in autonomous driving: A review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, pp. 3412–3432, 2021.
- [2] Gregor Luetzenburg, Aart Kroon, and Anders Anker Bjørk, "Evaluation of the apple iphone 12 pro lidar for an application in geosciences," *Scientific Reports*, vol. 11, 2021.
- [3] Christoph Gollob, Tim Ritter, Ralf Kraßnitzer, Andreas Tockner, and Arne Nothdurft, "Measurement of forest inventory parameters with apple ipad pro and integrated lidar technology," *Remote. Sens.*, vol. 13, pp. 3129, 2021.
- [4] Tao Yang, You Li, Cheng Zhao, Dexin Yao, Guanyin Chen, Li Sun, Tomas Krajnik, and Zhi Yan, "3d tof lidar in mobile robotics: A review," *arXiv preprint* arXiv:2202.11025, 2022.
- [5] Tekla Tóth, Zoltán Pusztai, and Levente Hajder, "Automatic lidar-camera calibration of extrinsic parameters using a spherical target," in 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 8580–8586.
- [6] Felix Igelbrink, Thomas Wiemann, Sebastian Pütz, and Joachim Hertzberg, "Markerless ad-hoc calibration of a hyperspectral camera and a 3d laser scanner," in *International Conference on Intelligent Autonomous Systems*. Springer, 2018, pp. 748–759.
- [7] Peng Jiang, Philip Osteen, and Srikanth Saripalli, "Semcal: Semantic lidar-camera calibration using neural mutual information estimator," in 2021 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI). IEEE, 2021, pp. 1–7.
- [8] Xinyu Zhang, Shifan Zhu, Shichun Guo, Jun Li, and Huaping Liu, "Line-based automatic extrinsic calibration of lidar and camera," in 2021 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2021, pp. 9347–9353.
- [9] Huai Yu, Weikun Zhen, Wen Yang, and Sebastian Scherer, "Line-based 2-d–3-d registration and camera localization in structured environments," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 11, pp. 8962–8972, 2020.
- [10] J Peršić, L Petrović, I Marković, and I Petrović, "Online multi-sensor calibration based on moving object tracking," *Advanced Robotics*, vol. 35, no. 3-4, pp. 130–140, 2021.

- [11] Ganesh Iyer, R Karnik Ram, J Krishna Murthy, and K Madhava Krishna, "Calibnet: Geometrically supervised extrinsic calibration using 3d spatial transformer networks," in 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2018, pp. 1110–1117.
- [12] Kaiwen Yuan, Zhenyu Guo, and Z Jane Wang, "Rggnet: Tolerance aware lidar-camera online calibration with geometric deep learning and generative model," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6956–6963, 2020.
- [13] Nick Schneider, Florian Piewak, Christoph Stiller, and Uwe Franke, "Regnet: Multimodal sensor registration using deep neural networks," in 2017 IEEE intelligent vehicles symposium (IV). IEEE, 2017, pp. 1803–1810.
- [14] Ganning Zhao, Jiesi Hu, Suya You, and C. C. Jay Kuo, "Calibdnn: multimodal sensor calibration for perception using deep neural networks," in *Defense + Commercial Sensing*, 2021.
- [15] Xudong Lv, Boya Wang, Ziwen Dou, Dong Ye, and Shuo Wang, "Lccnet: Lidar and camera self-calibration using cost volume network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2894–2901.
- [16] Peng Jiang, Philip Osteen, Maggie Wigness, and Srikanth Saripalli, "Rellis-3d dataset: Data, benchmarks and analysis," in 2021 IEEE international conference on robotics and automation (ICRA). IEEE, 2021, pp. 1110– 1116.
- [17] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz, "Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume," 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8934–8943, 2018.
- [18] "Adaptive max pooling," https://pytorch. org/docs/stable/generated/torch.nn. AdaptiveMaxPool2d.html, Accessed: 2022-10-19.
- [19] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan, "Supervised contrastive learning," *ArXiv*, vol. abs/2004.11362, 2020.
- [20] "Supcontrast," https://github.com/ HobbitLong/SupContrast, Accessed: 2022-10-19.