# SphNet: A Spherical Network for Semantic Pointcloud Segmentation

Lukas Bernreiter, Lionel Ott, Roland Siegwart and Cesar Cadena

*Abstract*—Semantic segmentation for robotic systems can enable a wide range of applications, from self-driving cars and augmented reality systems to domestic robots. We argue that a spherical representation is a natural one for egocentric pointclouds. Thus, in this work, we present a novel framework exploiting such a representation of LiDAR pointclouds for the task of semantic segmentation. Our approach is based on a spherical convolutional neural network that can seamlessly handle observations from various sensor systems (*e.g.*, different LiDAR systems) and provides an accurate segmentation of the environment. We operate in two distinct stages: First, we encode the projected input pointclouds to spherical features. Second, we decode and back-project the spherical features to achieve an accurate semantic segmentation of the pointcloud. We evaluate our method with respect to state-of-the-art projection-based semantic segmentation approaches using well-known public datasets. We demonstrate that the spherical representation enables us to provide more accurate segmentation and to have a better generalization to sensors with different field-of-view and number of beams than what was seen during training.

## I. INTRODUCTION

Over the past years, there has been a growing demand in robotics and self-driving cars for reliable semantic segmentation of the environment, *i.e.,* associating a class or label with each measurement sample for a given input modality. A semantic understanding of the surroundings is a critical aspect of robot autonomy. It has the potential to, *e.g.,* enable a comprehensive description of the navigational risks or disambiguate challenging situations in planning or mapping. For many of the currently employed robotic systems, the long-term stability of maps is a pertaining issue due to the often limited metrical understanding of the environment for which high-level semantic information is a possible solution.

With the advances in deep learning, vision-based semantic segmentation frameworks have become a very mature field. While there has also been significant progress on LiDAR-based semantic segmentation frameworks, it is still not as developed as their vision-based counterpart.

Nevertheless, LiDAR-based approaches have certain crucial advantages over other modalities as they are unaffected by the illumination conditions of the environment. This is in contrast to cameras which provide crucial descriptive information but are heavily affected by poor lighting conditions. Consequently, LiDAR-based systems effectively provide a more resilient segmentation system for a variety of challenging scenarios,

All authors are with the Autonomous Systems Lab, ETH Zurich, Zurich 8092, Switzerland, {berlukas, lioott, rsiegwart, cesarc}@ethz.ch.
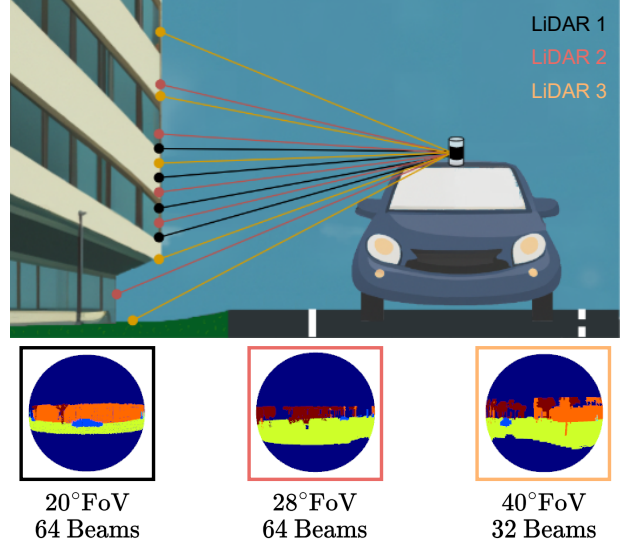
Fig. 1. We propose a spherical semantic segmentation framework that can handle pointclouds from various LiDAR sensors with different vertical field-of-view and angular resolutions.

such as operating at night and dynamically changing lighting conditions.

Many existing approaches operate using projection models, which typically transform the irregular pointcloud data into an ordered 2D domain, allowing them to utilize the extensive research available for images. The downside is that this requires a predefined configuration based on the number of beams, angular resolution, and vertical Field-of-View (FoV). LiDAR systems differ in these properties, which means that changing the sensory system after training might yield projections with geometrical and structural scarcity. Consequently, the resulting projection is often insufficient to express the complexity of arbitrary environments accurately. Accordingly, utilizing these approaches in generic environments with an arbitrary sensor system is often impossible without refining the initial network on the data from the new sensor environment.

LiDAR sensors are known to yield accurate geometrical and structural cues. Thus, modern LiDAR sensors often provide large FoV measurements to precisely measure the robot's surroundings. However, the projection onto the 2D domain of such large FoV scans introduces distortions that deform the physical dimensions the environment. Thus, dealing with different FoV, sensor frequencies, and scales remains an open research problem for which the input representation constitutes a major factor. In recent years, LiDAR sensors have become more affordable and available, becoming abundant in the context of robotics. However, many state-of-the-art segmentation methods are often limited to a particular sensor and cannot benefit from multiple LiDARs due to their

representation of the input. For example, systems employing multiple LiDARs pointing in different directions are typically processed sequentially or using multiple instances of the same network, one for each sensor. However, more crucial insights and structural dependencies in overlapping areas can be considered by jointly predicting the segmentation using all available sensors.

In this work, we propose a framework that takes LiDAR scans as input (cf. Figure 1), projects them onto a sphere, and utilizes a spherical Convolutional Neural Network (CNN) for the task of semantic segmentation. The projection of the LiDAR scans onto the sphere does not introduce any distortions and is independent of the utilized LiDAR, thus, yielding an agnostic representation for various LiDAR systems with different vertical FoV. We adapt the structure of common 2D encoder and decoder networks and support simultaneous training on different datasets obtained with varying LiDAR sensors and parameters without having to adapt our configuration. Moreover, since our approach is invariant to rotations due to the spherical representation, we support arbitrarily rotated input pointclouds. In summary, the key contributions of this paper are as follows:

- A spherical end-to-end pipeline for semantic segmentation supporting various input configurations.
- A spherical encoder-decoder structure including a spectral pooling and unpooling operation for $SO(3)$ signals.

## II. RELATED WORK

Methods using a LiDAR have to deal with the inherent sparsity and irregularity of the data in contrast to vision-based approaches. Moreover, LiDAR-based methods have various choices on how to represent the input data [1], including, directly using the pointcloud [2], [3], voxel-based [4]–[6] or projection-based [7]–[10] representations. The selection of the input representation, which yields the best performance for a specific task, however, still remains an open research question.

Direct methods such as PointNet [2], [3] operate on the raw unordered pointcloud and extract local contextual features using point convolutions [11]. Voxel-based approaches [4], [5], [12] keep all the geometric understanding of the environment and can readily accumulate multiple scans either chronologically or from different sensors. SpSequenceNet [13] explicitly uses 4D pointclouds and considers the temporal information between consecutive scans. However, it is evident that the computational complexity of voxel-based approaches is high due to their high-dimensional convolutions, and their accuracy and performance are directly linked to the chosen voxel size, which resulted in works that organize the pointclouds into an Octree, Kdtree, etc. [14] for efficiency. Furthermore, instead of using a cartesian grid, PolarNet [15] discretizes the space using a polar grid and shows superior quality. A different direction of research is offered by graph-based approaches [16] which can seamlessly model the irregular structure of pointclouds though more experimental directions in terms of graph building, and network design are still to be addressed.

Projection-based methods differ from other approaches by transforming the pointcloud into a specific domain, such as 2D images, which the majority of projection-based methods [7]–[10] rely on.

Furthermore, projections to 2D images are appealing as it enables leveraging all the research in image-based deep learning but generally need to rely on the limited amount of labeled pointcloud data. Hence, the work of Wu et al. [17] tackles the deficiency in labeled pointcloud data by using domain-adaption between synthetic and real-world data.

The downsides of the projection onto the 2D domain are: i) the lack of a detailed geometric understanding of the environment and ii) the large FoV of LiDARs, which produces significant distortions, decreasing the accuracy of these methods. Hence, recent approaches have explored using a combination of several representations [6], [18] and convolutions [19]. Recent works [20], [21] additionally learn and extract features from a Bird's Eye View projection that would otherwise be difficult to retain with a 2D projection.

In contrast to 2D image projections, projecting onto the sphere is a more suitable representation for such large FoV sensors. Recently, spherical CNNs [22]–[24] have shown great potential for, *e.g.,* omnidirectional images [25], [26] and cortical surfaces [27], [28].

Moreover, Lohit et al. [25] proposes an encoder-decoder spherical network design that is rotation-invariant by performing a global average pooling of the encoded feature map. However, their work discards the rotation information of the input signals and thus, needs a special loss that includes a spherical correlation to find the unknown rotation w.r.t. the ground truth labels.

Considering the findings above, we propose a composition of spherical CNNs, based on the work of Cohen et al. [23], that semantically segments pointclouds from various LiDAR sensor configurations.

## III. SPHERICAL SEMANTIC SEGMENTATION

This section describes the core modules of our spherical semantic segmentation framework, which mainly operates in three stages: i) feature projection, ii) semantic segmentation, and iii) back-projection (*cf.* Figure 2).

Initially, we discuss the projection of LiDAR pointclouds onto the unit sphere and the feature representation that serves as input to the spherical CNN. Next, we describe the details of our network design and architecture used to learn a semantic segmentation of LiDAR scans.

### A. Sensor Projection and Feature Representation

Initially, the input to our spherical segmentation network is a signal defined on the sphere $S^2 = \left\{ \boldsymbol{p} \in \mathbb{R}^3 \mid \|\boldsymbol{p}\|_2 = 1 \right\}$, with the parametrization as proposed by Healy et al. [29], *i.e.*

$$\boldsymbol{\omega}(\phi, \theta) = [\cos \phi \sin \theta, \sin \phi \sin \theta, \cos \theta]^\top, \quad (1)$$

where $\boldsymbol{\omega} \in S^2$, and $\phi \in [0, 2\pi]$ and $\theta \in [0, \pi]$ are azimuthal and polar angle, respectively.

We then operate in an end-to-end fashion by transforming the input modality (*i.e.,* the pointcloud scan) into a spherical
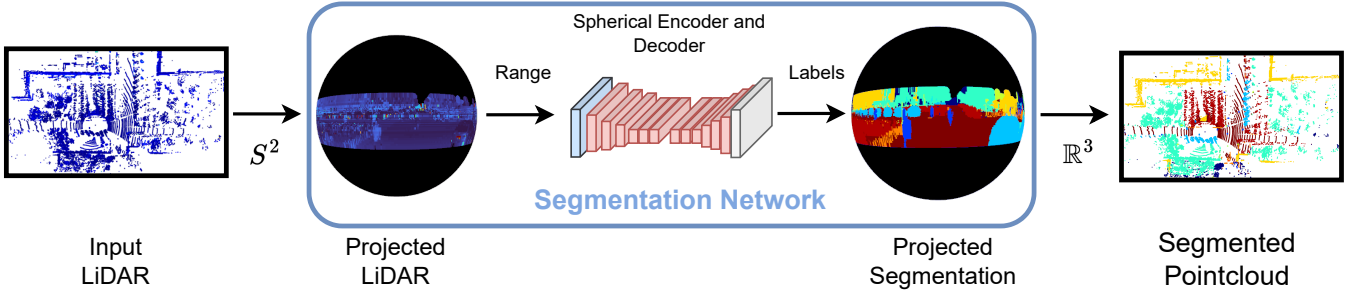
Fig. 2. Overview of our proposed multi-modal segmentation framework. Initially, we employ a base network to a LiDAR scan to get a semantic segmentation. Next, the segmentation is fused with all camera images to get a refined solution. Finally, the spherical segmentation is back-projected into its original form.

representation in $S^2$. Consequently, as an initial step, a pointcloud is projected onto the unit sphere, *e.g.* the $j$th point $\boldsymbol{p}_j = [x_j, y_j, z_j]^\top$ is projected using

$$\phi_j = \arctan\left(\frac{y_j}{x_j}\right), \qquad \theta_j = \arccos\left(\frac{z_j}{||\boldsymbol{p}_j||_2}\right), \quad (2)$$

From the LiDAR projection, we sample the range values using an equiangular grid complying with the sampling theorem by Discroll and Healy [30]. It is important to note that we omit the sampling of the intensity/remission values of the LiDAR scans since their values significantly fluctuate between LiDAR sensors, particularly between different manufacturers. Thus, would require additional investigation and further calibration of the input data.

Although in this work, we only utilize measurements from a LiDAR scanner, our approach is not limited to this sensor type. Other sensory systems, such as multi-camera systems, thermal cameras, and depth sensory systems, can seamlessly be integrated by performing the same spherical projection. Additionally, since all sensors would share the same feature representation, our approach readily facilitates various combinations of heterogeneous sensory types [26]. Moreover, despite utilizing only LiDAR sensors defined over $360°$, our sampling approach is agnostic to the resolution and field-of-view of the sensor and can be used with arbitrary viewpoint coverage.

Finally, after the projection and sampling, the spherical network will receive the LiDAR scan as a feature vector $\in \mathbb{R}^{2 \times 2\text{BW} \times 2\text{BW}}$, where be BW corresponds to the spherical bandwidth used for the equiangular sampling [30]. The chosen bandwidth directly controls the employed spatial discretization and, consequently, the spectrum's resolution in the frequency domain. In particular, the higher the bandwidth, the finer the spectral resolution, and more memory is required, which needs to be carefully considered when designing the spherical network, *e.g.*, the first convolutional layer is designed to operate on the exact BW as used by the sampling but later layers will decrease the bandwidth. In the following, we discuss the design of the remaining layers in our network.

### B. Network Design and Architecture

In this section, we will discuss the design choices of our network that takes LiDAR scans as input and is able to learn a semantic segmentation of it. The input of the base network will be the sampled spherical features from the

LiDAR pointcloud, i.e., the range values. We found that an initial spherical bandwidth between 50 - 120 yields a good trade-off between accuracy and memory consumption.

Overall, the network design is based on a spherical encoder-decoder structure. The spherical encoder network increases the number of features from layer to layer while the bandwidth is decreased instead. In a similar vein, the spherical decoder decreases the number of features while increasing the bandwidth from layer to layer back to the originally utilized bandwidth. An overview of our spherical network architecture is given in Figure 3. Moreover, by decreasing and increasing the bandwidth during encoding and decoding, we will also increase and decrease the size of the kernels for consecutive convolutions, respectively.

*1) Feature Encoding:* Our network is based on the work of Cohen et al. [23]. Thus we lift the features to $SO(3)$ during the encoding and have to revert to $S^2$ during the decoding. In other words, only the initial layer of the network performs a convolution over $S^2$, whereas the remaining convolutional layers act on $SO(3)$ to preserve the convolution's equivariance property [31] (*cf.* the left part in Figure 3). During the convolutions, we employ spatially localized kernels that are rotated around the sphere using operations in $SO(3)$. Here, $SO(3)$ denotes the three-dimensional rotation group consisting of $\alpha$, $\beta$, and $\gamma$, corresponding to roll, pitch, and yaw (RPY).

Each convolution in our network is efficiently implemented based on the convolution theorem by performing a Fourier transform in $S^2$ and $SO(3)$ [32], respectively, *i.e.* the convolution between two signals $f$ and $g$ is given by

$$f * g = \mathcal{F}^{-1}\left\{\mathcal{F}\{f\} \cdot \mathcal{F}\{g\}\right\}, \qquad (3)$$

where $\mathcal{F}$ is either a $S^2$ or $SO(3)$ Fourier transform and $\mathcal{F}^{-1}$ its inverse. Additionally, we apply a PReLu [33] activation function followed by a three-dimensional batch normalization after each convolution. The last convolutional block during the encoding also applies a dropout for additional regularization.

In our approach, the $S^2$ convolution increases the number of features but is not done in place, *i.e.*, it does not preserve the input bandwidth. Rather, it directly decreases the output bandwidth as part of the $S^2$ Fourier transform by having a smaller output bandwidth for the inverse $S^2$ Fourier transform in Eq. (3). The $SO(3)$ convolution increases the number of input features but preserves the utilized spherical bandwidth.

Features (in/out)

Bandwidth (in/out)

1/30  30/45  45/120  120/200  200/120  120/45  45/30  30 / |C|

50/40  40/30  30/20  20/10  10/8  8/10  10/20  20/30  30/40  40/50

Encoder                    Decoder

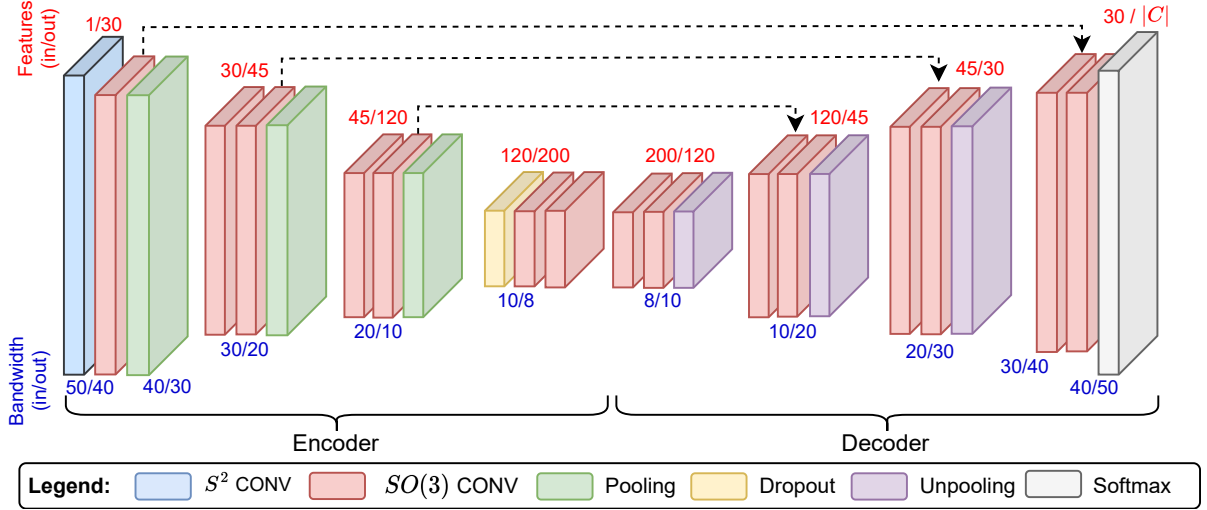**Legend:** $S^2$ CONV | $SO(3)$ CONV | Pooling | Dropout | Unpooling | Softmax

Fig. 3. Proposed network architecture. We use an encoder-decoder design where the dotted lines denote skip connections. Initially, our network lifts the features to $SO(3)$ and will eventually integrate back to $S^2$ for the semantic segmentation. The network puts out the logits for the number of configured classes $C$, which then, together with a softmax, results in a semantic segmentation of the input pointcloud.

And, the $SO(3)$ blocks decrease the bandwidth by applying a pooling operation in $SO(3)$.

Pooling and unpooling are done in the spectral domain of $SO(3)$, which has the advantage of retaining the equivariance [22], [34] as opposed to spatial pooling. In practice, the $SO(3)$ pooling is implemented by transforming the signal to the spectral domain using an $SO(3)$ Fourier transform [29], [32] and subsequent low-pass filtering of the resulting spectrum. The inverse $SO(3)$ Fourier transform then yields the pooled $SO(3)$ signal. Between the forward and backward passes, we temporarily store the input bandwidth and inverse the operation by zero-padding the spectrum to the original size.

*2) Feature Decoding:* All the convolutions in the spherical decoding component are done in $SO(3)$ (*cf.* the right part in Figure 3). Moreover, we unpool the input signals to increase the bandwidth to match the input bandwidth again. In particular, we apply the $SO(3)$ Fourier transform to the input signal and zero-pad the resulting spectrum to a larger size. Finally, the zero-padded signal is inverse $SO(3)$ Fourier transformed and passed to the next convolutional layer. Similar to the encoding part, we perform the inverse by applying an idealized low-pass filter for the backward pass. The last convolution during decoding is different from the preceding operations as it zero-pads the convolved signal in the spectral domain to achieve the initial input bandwidth (sampling bandwidth). Thus, the last layer does not rely on any unpooling operation. Although the output has the correct size, it still needs to be mapped back to its original space, *i.e.,* $S^2$. To transform the $SO(3)$ signal back to $S^2$, a max pooling operation or an integration of the $SO(3)$ signal over the last entry $\gamma$ (*i.e.,* the yaw angle) can be used, resulting in $\mathbb{R}^{2BW \times 2BW \times 2BW} \mapsto \mathbb{R}^{2BW \times 2BW}$. We have selected the latter approach for its efficient computation and simplicity. During the inference, the spherical semantic segmentation is then achieved by applying a final softmax layer to the result of the integration.

*3) Loss:* Finally, our proposed spherical network uses a common loss definition for semantic segmentation [8], [20], *i.e.* for prediction $\hat{y}$ and ground truth $y$ labels

$$\mathcal{L}_{XC}(y, \hat{y}) = -\sum_i w_i P(y_i) \log P(\hat{y}_i) \tag{4}$$

$$\mathcal{L}_{LZ}(y, \hat{y}) = \frac{1}{|C|} \sum_{c \in C} J(e(c)) \tag{5}$$

$$\mathcal{L}(y, \hat{y}) = \mathcal{L}_{XC}(y, \hat{y}) + \mathcal{L}_{LZ}(y, \hat{y}), \tag{6}$$

where $\mathcal{L}_{XC}$ is a weighted cross-entropy loss where $w_i$ are the class weights and $P(\cdot)$ the corresponding probabilities. The latter term $\mathcal{L}_{LZ}$ is a lovasz-softmax [41] loss where J is the lovasz intersection over union, $e(c)$ the errors for class $c$ and $C$ the set of all classes. Notably, the loss operates on the equiangular samples in $S^2$, and since we do not discard the rotational information of the signals on the sphere, we have a direct mapping between the input and the output signals.

*4) Back-Projection:* The final pointcloud segmentation in $\mathbb{R}^3$ is achieved by back-projecting the spherical projection from $S^2$ to $\mathbb{R}^3$ using the sampled range values, *i.e.,* inverting Eq. (2), *s.t.* for the $j$th projected point with $\phi_j$ and $\theta_j$

$$x_j = r_j \cdot \cos(\phi_j) \sin(\theta_j) \tag{7}$$

$$y_j = r_j \cdot \sin(\phi_j) \sin(\theta_j) \tag{8}$$

$$z_j = r_j \cdot \cos(\theta_j), \tag{9}$$

where $r_j = \|\boldsymbol{p}_j\|_2$.

## IV. EXPERIMENTS

This section presents the experimental validation of our spherical segmentation framework, where we show that our pipeline gives an accurate semantic segmentation of the environment and that it generalizes well to different sensory setups. We first validate the segmentation quality of our spherical network and compare it to current state-of-the-art projection-based segmentation frameworks. Next, we

Segmentation Quality Comparison in Terms of mIoU / Acc

| Datasets vFoV / Beams | | nuScenes [35] 40° / 32 | SemanticKITTI [36] 28° / 64 | SemanticPOSS [37] 23° / 40 | Waymo [38] 20° / 64 | A2D2 [39] 30° / 16 | PC-Urban [40] 45° / 64 |
|---|---|---|---|---|---|---|---|
| **Methods** | ∅ **mIoU / Acc** | mIoU / Acc | mIoU / Acc | mIoU / Acc | mIoU / Acc | mIoU / Acc | mIoU / Acc |
| SqueezeSeg [9] | 9.5 / 57.0 | 9.5 / 57.1 | 10.8 / 64.5 | 9.5 / 57.3 | 6.7 / 40.1 | 10.6 / 63.9 | 9.9 / 59.3 |
| SqueezeSegV2 [17] | 39.1 / 81.7 | 49.4 / 89.4 | 45.8 / 88.6 | 39.3 / 80.7 | 37.1 / 75.0 | 36.1 / 88.0 | 26.8 / 68.4 |
| RangeNet++ [7] | 38.0 / 80.4 | 49.1 / 88.4 | 45.1 / 88.6 | 37.9 / 78.7 | 34.3 / 72.0 | 35.6 / 87.6 | 26.2 / 66.8 |
| 3D-MiniNet [10] | 42.6 / 81.2 | 54.9 / 90.5 | 51.1 / 92.5 | 44.3 / 84.8 | 45.9 / 81.5 | 39.3 / 89.7 | 20.1 / 48.1 |
| SalsaNext [8] | 44.1 / 82.1 | **61.8 / 92.4** | 50.6 / 91.8 | 42.6 / 82.7 | 38.6 / 74.6 | 43.7 / 90.8 | 27.3 / 60.3 |
| Ours | **49.0 / 97.2** | 55.2 / 95.4 | **52.7 / 96.9** | **51.2 / 96.8** | **47.1 / 98.4** | **46.0 / 98.8** | **41.8 / 97.0** |

TABLE I

COMPARISON OF THE MEAN IoU (MIoU) AND THE ACCURACY (ACC) IN THE VALIDATION SPLIT FOR EACH DATASET. ALL DATASETS WERE PART OF THE TRAINING SPLIT EXCEPT FOR THE PC-URBAN DATASET, WHICH ALSO CONTAINS A NEW LiDAR SENSOR THAT WAS NOT SEEN YET BY THE NETWORKS. THE LEFTMOST DATA COLUMN DEPICTS THE AVERAGE MIou AND ACC OVER ALL DATASETS. ALL VALUES ARE GIVEN AS PERCENTAGES [%].

demonstrate the flexibility of our representation by increasing the input field-of-view drastically. Finally, we evaluate the computational cost to show the method's applicability to real-world scenarios. We use the structure depicted in Figure 3 for all experiments.

### A. Experiment and Training Setup

We validate our proposed approach by comparing the performance of our spherical network to state-of-the-art projection-based segmentation frameworks, RangeNet++ [7], SqueezeSeg [9], [17], 3D-MiniNet [10], and SalsaNext [8].

For the comparison, we utilize several datasets with various different LiDARs to show the ease of use of our proposed approach given sensors with different vertical FoV (vFoV) and number of beams. All networks use only the range values as feature and are trained on the nuScenes [35], SemanticKITTI [36], SemanticPOSS [37], Waymo [38], A2D2 [39]. We trained all networks for 50 epochs and used the mean Intersection over Union (mIoU) and accuracy metric to assess the quality of the similarity to the ground truth.

Moreover, since all these datasets provide significantly different classes, *e.g.* nuScenes groups all objects into man-made objects. Similar to Sanchez et al. [42], we abstracted the semantic classes to a total of five classes in order to provide a shared set of classes between the datasets: vehicles, persons, ground, man-made, and vegetation.

In addition, for each dataset that provides a split between training and validation, we randomly selected 4000 point-clouds from the former for training the networks and used the latter for inference. For the datasets without a split, we created a split and followed the above procedure for training and inference.

The input data to the projection-based methods was projected into the largest FoV of the LiDARs known during training. It is important to note that our approach, in contrast to the other ones used, does not require any such considerations and projects all pointclouds directly onto the sphere.

### B. Segmentation Quality Comparison and Validation

We first evaluate the performance within the training domain, *i.e.,* the training and test are in the same urban setting. The results are shown in Table I. One can see that most of the other baseline methods suffer from the large variations in

Segmentation Quality Comparison on a new Domain

| Methods | mIoU | Vehicle | Person | Ground | Man-Made | Vegetation |
|---|---|---|---|---|---|---|
| SqueezeSeg [9] | 12.3 | 0.0 | 0.0 | 61.4 | 0.0 | 0.0 |
| SqueezeSegV2 [17] | 18.2 | 22.2 | 3.1 | 58.0 | 7.6 | 0.0 |
| RangeNet++ [7] | 21.5 | 22.6 | 4.7 | 60.4 | 20.0 | 0.0 |
| 3D-MiniNet [10] | 20.8 | 24.4 | **16.7** | 60.0 | 2.8 | 0.0 |
| SalsaNext [8] | 22.7 | **32.7** | 9.5 | **62.8** | 8.3 | 0.0 |
| Ours | **30.1** | 19.5 | 1.8 | 59.9 | 23.7 | **49.3** |

TABLE II

COMPARISON OF THE CLASS-WISE AND MEAN IoU FOR THE SEMANTICUSL [43] DATASET, WHICH USES A DIFFERENT ENVIRONMENT AND SENSORY SYSTEM THAN WHAT WAS SEEN DURING TRAINING. ALL IoU VALUES ARE GIVEN AS PERCENTAGES [%].

the data and consequently also fluctuate in their performance significantly. Notably, due to the difference in vFoV the projections distort the physical dimensions, resulting in less accuracy where vFov is smaller (*e.g.,* SemanticPOSS [37] and Waymo [38]).

In contrast, our proposed approach achieves a good segmentation of the input pointcloud and, most importantly, achieves consistent scores across the various datasets. Although our approach does not achieve the segmentation quality for the nuScenes [35] dataset in our tests, its primary benefit is the generalization capabilities between sensory systems. Thanks to the spherical representation of pointcloud data employed by our method, we can correctly perform the projection required by the different opening angles of the LiDARs. Consequently, our approach is less affected by the changing angular resolution between the datasets in comparison to the image-based projection approaches.

We additionally evaluate the networks on an unseen sensory system, PC-Urban [40] to assess our approach's generalization capabilities within the same domain. The difference in vFoV (12.5° down and −7.5° up) results in warped physical proportions for the 2D projection-based methods, which greatly decreases their segmentation performance. Consequently, the segmentation quality of high objects such as man-made buildings and vegetation significantly degrades. Our approach is less affected by the change in the vFoV and hence, also attains the highest mIoU on this dataset.

## C. Segmentation Quality on a Different Domain

In this experiment, we test all approaches on the Semanti-cUSL [43] (45° vFoV and 64 beams) dataset, which offers a different sensory system and environment (off-road and campus scenes) than what was seen by the networks during training. The change in the vFoV is the same as for the PC-Urban dataset in the previous experiment. Table II shows a comparison between the class-wise and mean IoU.

Due to the vastly different environment and sensor intrinsics, all but our proposed method fail entirely to segment the vegetation class as the learned representation no longer matches. Our approach is not affected by such warping of the physical dimensions and, therefore, achieves a better segmentation quality.

However, it is difficult for our representation to disambiguate persons that are far away from man-made objects and small vegetation when only range values are available. Providing additional information through multiple modalities is a possible solution and is left for future work. Nevertheless, our approach still maintains the highest mIoU compared to the other approaches. This highlights the main benefit of our method, which is its ability to generalize and consequently provide an improved overall segmentation.

## D. Semantic Segmentation of Rotated Pointclouds

Next, we show that our representation has the advantage of rotational invariance to the input data, allowing pointclouds to be arbitrarily rotated. This allows our method to support various input configurations such as different angular resolutions and tilted sensor mounts.

In this experiment, we applied a predefined rotation from $0°$ to $180°$ around the RPY axes of the input pointcloud. Figure 4 shows the mIoU for various rotational shifts using sequence 08 of the SemanticKITTI dataset, and the model trained in Section IV-A. The rotation of the pointclouds yields
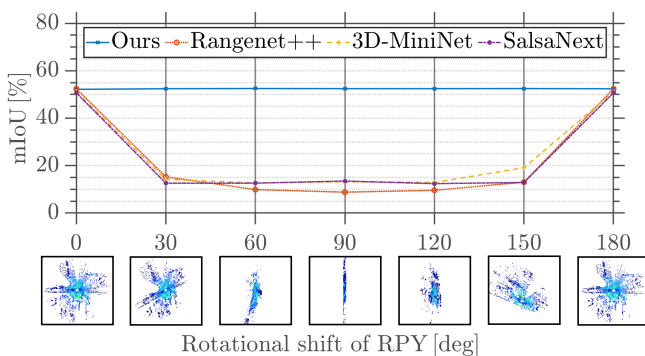


Fig. 4. Comparison of the mIoU performance using different rotational shifts on the SemanticKITTI sequence 08. The bottom row of pointclouds shows exemplary results after applying the rotations.

odd and ineffective representations of the scans, and thus, all other projection-based methods experience a large drop in their performance. Note that since the rotation of $180°$ around RPY results in the original pointcloud, the initial mIoU is restored again. Generally, 2D projection methods implicitly require the pointclouds to be horizontally oriented

for an efficient prediction. Hence, it is particularly difficult for these methods to utilize multiple LiDARs simultaneously if one of the sensors is tilted w.r.t. the other ones.

The spherical projection is not only more efficient and natural, but the spherical Fourier transform is also invariant to rotations. Thus, our approach is completely unaffected by the rotated pointclouds and maintains the mIoU over all rotational shifts.

## E. Runtime Evaluation

Finally, we assess the runtime performance of our proposed approach to understand its potential for deploying it in real-world applications. Figure 5 shows the evaluation of the execution time for the different components of our proposed approach. In this experiment, we consider a sampling band-
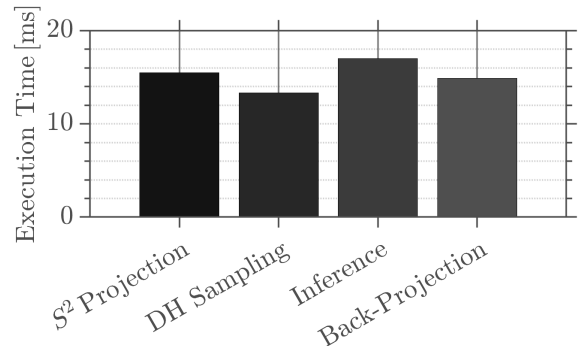


Fig. 5. Execution time in ms partitioned per component. All values are averaged over 1000 samples. Our proposed approach is able to segment an input pointcloud in approximately $60\,\mathrm{ms}$.

width of 50 for the LiDAR as depicted in Figure 3. The benchmark was performed on an Intel Xeon E5-2640v3 with an NVIDIA Titan RTX, and all scripts are written using PyTorch. It is evident that the discretization of the irregular pointcloud data on the sphere takes a considerable portion. Nevertheless, our approach is able to infer a semantic segmentation in approximately $60\,\mathrm{ms}$.

## V. CONCLUSION AND FUTURE WORK

In this paper, we proposed a spherical representation of pointclouds that can be used to train a model using various LiDAR sensors with different parameters. We presented in this context an end-to-end approach for semantic segmentation based on a spherical encoder-decoder network and showed that the spherical representation is a much more favorable representation, especially for high FoV LiDAR systems. Most importantly, our findings also indicate that our approach is invariant to rotations and has a better generalization to unseen LiDAR systems after training a model. Furthermore, our proposed approach is not limited to depth sensors, and other sensor types, such as RGB and thermal cameras, can be readily incorporated [26].

In future research, we intend to investigate two separate research directions. First, we explore the fusion of multiple camera images to refine the initial semantic segmentation and its ontology. Second, we intend to migrate our spherical network to be fully in $S^2$ in order to decrease the memory requirements and improve its practical applicability.

REFERENCES

[1] B. Gao, Y. Pan, C. Li, S. Geng, and H. Zhao, "Are We Hungry for 3D LiDAR Data for Semantic Segmentation? A Survey of Datasets and Methods," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 6063–6081, 7 2022.

[2] R. Q. Charles, H. Su, M. Kaichun, and L. J. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 7 2017, pp. 77–85.

[3] C. Qi, L. Yi, H. Su, and L. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," *NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems*, no. Dec, p. 5105–5114, 2017.

[4] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, and S. Savarese, "SEGCloud: Semantic Segmentation of 3D Point Clouds," in *2017 International Conference on 3D Vision (3DV)*. IEEE, 10 2017, pp. 537–547.

[5] H. Zhou, X. Zhu, X. Song, Y. Ma, Z. Wang, H. Li, and D. Lin, "Cylinder3D: An Effective 3D Framework for Driving-scene LiDAR Semantic Segmentation," 8 2020.

[6] J. Xu, R. Zhang, J. Dou, Y. Zhu, J. Sun, and S. Pu, "RPVNet: A Deep and Efficient Range-Point-Voxel Fusion Network for LiDAR Point Cloud Segmentation," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, vol. 1. IEEE, 10 2021, pp. 16 004–16 013.

[7] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss, "RangeNet ++: Fast and Accurate LiDAR Semantic Segmentation," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, no. i. IEEE, 11 2019, pp. 4213–4220.

[8] T. Cortinhal, G. Tzelepis, and E. Erdal Aksoy, "SalsaNext: Fast, Uncertainty-Aware Semantic Segmentation of LiDAR Point Clouds," in *International Symposium on Visual Computing*, vol. 12510 LNCS, 2020, pp. 207–222.

[9] B. Wu, A. Wan, X. Yue, and K. Keutzer, "SqueezeSeg: Convolutional Neural Nets with Recurrent CRF for Real-Time Road-Object Segmentation from 3D LiDAR Point Cloud," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 5 2018, pp. 1887–1893.

[10] I. Alonso, L. Riazuelo, L. Montesano, and A. C. Murillo, "3D-MiniNet: Learning a 2D Representation from Point Clouds for Fast and Efficient 3D LIDAR Semantic Segmentation," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5432–5439, 10 2020.

[11] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "PointCNN : Convolution On X -Transformed Points," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, no. NeurIPS, 2018, p. 828–838.

[12] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3D U-Net: Learning Dense Volumetric Segmentation from Sparse Annotation," in *Medical Image Computing and Computer-Assisted Intervention*, ser. Lecture Notes in Computer Science, S. Ourselin, L. Joskowicz, M. R. Sabuncu, G. Unal, and W. Wells, Eds. Cham: Springer International Publishing, 2016, vol. 9901, pp. 424–432.

[13] H. Shi, G. Lin, H. Wang, T. Y. Hung, and Z. Wang, "Spsequencenet: Semantic segmentation network on 4D point clouds," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 4573–4582, 2020.

[14] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz, "SPLATNet: Sparse Lattice Networks for Point Cloud Processing," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 6 2018, pp. 2530–2539.

[15] Y. Zhang, Z. Zhou, P. David, X. Yue, Z. Xi, B. Gong, and H. Foroosh, "PolarNet: An Improved Grid Representation for Online LiDAR Point Clouds Semantic Segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 9601–9610.

[16] L. Wang, Y. Huang, Y. Hou, S. Zhang, and J. Shan, "Graph Attention Convolution for Point Cloud Semantic Segmentation," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2019-June. IEEE, 6 2019, pp. 10 288–10 297.

[17] B. Wu, X. Zhou, S. Zhao, X. Yue, and K. Keutzer, "SqueezeSegV2: Improved Model Structure and Unsupervised Domain Adaptation for Road-Object Segmentation from a LiDAR Point Cloud," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 5 2019, pp. 4376–4382.

[18] X. Li, G. Zhang, H. Pan, and Z. Wang, "CPGNet: Cascade Point-Grid Fusion Network for Real-Time LiDAR Semantic Segmentation," 2022.

[19] J. Park, C. Kim, and K. Jo, "PCSCNet: Fast 3D Semantic Segmentation of LiDAR Point Cloud for Autonomous Car using Point Convolution and Sparse Convolution Network," 2022.

[20] M. Gerdzhev, R. Razani, E. Taghavi, and L. Bingbing, "TORNADO-Net: mulTiview tOtal vaRiatioN semAntic segmentation with Diamond inceptiOn module," in *Proceedings - IEEE International Conference on Robotics and Automation*, vol. 2021-May. Institute of Electrical and Electronics Engineers Inc., 2021, pp. 9543–9549.

[21] F. Duerr, H. Weigel, and J. Beyerer, "RangeBird: Multi View Panoptic Segmentation of 3D Point Clouds with Neighborhood Attention," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 5 2022, pp. 11 131–11 137.

[22] C. Esteves, C. Allen-Blanchette, A. Makadia, and K. Daniilidis, "Learning SO(3) Equivariant Representations with Spherical CNNs," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11217 LNCS, no. 3, pp. 54–70, 2018.

[23] T. S. Cohen, M. Geiger, J. Koehler, and M. Welling, "Spherical CNNs," *Proceedings of the International Conference on Learning Representations*, no. 3, pp. 1–15, 1 2018.

[24] C. Esteves, A. Makadia, and K. Daniilidis, "Spin-Weighted Spherical CNNs," no. 3, 6 2020.

[25] S. Lohit and S. Trivedi, "Rotation-invariant autoencoders for signals on spheres," *arXiv*, 2020.

[26] L. Bernreiter, L. Ott, J. Nieto, R. Siegwart, and C. Cadena, "Spherical Multi-Modal Place Recognition for Heterogeneous Sensor Systems," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, vol. 2021-May, no. Icra. IEEE, 5 2021, pp. 1743–1750.

[27] F. Zhao, S. Xia, Z. Wu, D. Duan, L. Wang, W. Lin, J. H. Gilmore, D. Shen, and G. Li, "Spherical U-Net on Cortical Surfaces: Methods and Applications," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Springer International Publishing, 2019, vol. 11492 LNCS, pp. 855–866.

[28] F. Zhao, Z. Wu, L. Wang, W. Lin, J. H. Gilmore, S. Xia, D. Shen, and G. Li, "Spherical Deformable U-Net: Application to Cortical Surface Parcellation and Development Prediction," *IEEE Transactions on Medical Imaging*, vol. 40, no. 4, pp. 1217–1228, 2021.

[29] D. M. Healy, D. N. Rockmore, P. J. Kostelec, and S. Moore, "FFTs for the 2-Sphere-Improvements and Variations," *J. Fourier Anal. Appl.*, vol. 9, no. 4, pp. 341–385, 2003.

[30] J. R. Driscoll and D. M. Healy, "Computing fourier transforms and convolutions on the 2-sphere," pp. 202–250, 1994.

[31] T. Cohen, M. Geiger, J. Köhler, and M. Welling, "Convolutional Networks for Spherical Signals," *Principled Approaches to Deep Learning Workshop, ICML*, 9 2017.

[32] P. J. Kostelec and D. N. Rockmore, "FFTs on the Rotation Group," *Journal of Fourier Analysis and Applications*, vol. 14, no. 2, pp. 145–179, 4 2008.

[33] K. He, X. Zhang, S. Ren, and J. Sun, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification," in *2015 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2015, pp. 1026–1034.

[34] O. Rippel, J. Snoek, and R. P. Adams, "Spectral representations for convolutional neural networks," *Advances in Neural Information Processing Systems*, vol. 2015-Janua, pp. 2449–2457, 2015.

[35] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, O. Beijbom, and A. Company, "nuScenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 621–11 631.

[36] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, "SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 9297–9307.

[37] Y. Pan, B. Gao, J. Mei, S. Geng, C. Li, and H. Zhao, "SemanticPOSS: A Point Cloud Dataset with Large Quantity of Dynamic Instances," in

*2020 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 10 2020, pp. 687–693.

[38] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov, "Scalability in Perception for Autonomous Driving: Waymo Open Dataset," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2446–2454.

[39] J. Geyer, Y. Kassahun, M. Mahmudi, X. Ricou, R. Durgesh, A. S. Chung, L. Hauswald, V. H. Pham, M. Mühlegg, S. Dorn, T. Fernandez, M. Jänicke, S. Mirashi, C. Savani, M. Sturm, O. Vorobiov, M. Oelker, S. Garreis, and P. Schuberth, "A2D2: Audi Autonomous Driving Dataset," 4 2020.

[40] M. Ibrahim, N. Akhtar, M. Wise, and A. Mian, "Annotation Tool and Urban Dataset for 3D Point Cloud Semantic Segmentation," *IEEE Access*, vol. 9, pp. 35 984–35 996, 2021.

[41] M. Berman, A. R. Triki, and M. B. Blaschko, "The Lovasz-Softmax Loss: A Tractable Surrogate for the Optimization of the Intersection-Over-Union Measure in Neural Networks," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 6 2018, pp. 4413–4421.

[42] J. Sanchez, J.-E. Deschaud, and F. Goulette, "COLA: COarse LAbel pre-training for 3D semantic segmentation of sparse LiDAR datasets," 2 2022.

[43] P. Jiang and S. Saripalli, "LiDARNet: A Boundary-Aware Domain Adaptation Model for Point Cloud Semantic Segmentation," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, vol. 2021-May. IEEE, 5 2021, pp. 2457–2464.