

Gaussian Radar Transformer for Semantic Segmentation in Noisy Radar Data

Matthias Zeller¹ Jens Behley² Michael Heidingsfeld³ Cyrill Stachniss⁴

Abstract—Scene understanding is crucial for autonomous robots in dynamic environments for making future state predictions, avoiding collisions, and path planning. Camera and LiDAR perception made tremendous progress in recent years, but face limitations under adverse weather conditions. To leverage the full potential of multi-modal sensor suites, radar sensors are essential for safety critical tasks and are already installed in most new vehicles today. In this paper, we address the problem of semantic segmentation of moving objects in radar point clouds to enhance the perception of the environment with another sensor modality. Instead of aggregating multiple scans to densify the point clouds, we propose a novel approach based on the self-attention mechanism to accurately perform sparse, single-scan segmentation. Our approach, called Gaussian Radar Transformer, includes the newly introduced Gaussian transformer layer, which replaces the softmax normalization by a Gaussian function to decouple the contribution of individual points. To tackle the challenge of the transformer to capture long-range dependencies, we propose our attentive up- and downsampling modules to enlarge the receptive field and capture strong spatial relations. We compare our approach to other state-of-the-art methods on the RadarScenes data set and show superior segmentation quality in diverse environments, even without exploiting temporal information.

Index Terms—Semantic Scene Understanding, Deep Learning Methods

I. INTRODUCTION

AUTONOMOUS vehicles need to understand their surroundings to safely navigate in dynamic, real-world environments. To achieve holistic perception and enhance safety, the sensor suites of autonomous vehicles are versatile to explore redundant information of individual sensors such as cameras, LiDAR, or radar. Particularly in autonomous driving, where a malfunction of one modality can result in lethal consequences, redundancy is key. Widely explored cameras and LiDAR sensors capture the environment precisely but face limitations under adverse weather such as fog, rain, and snow. Additional information is required, which is accessible via radar sensors, and hence, makes them crucial to enable safe autonomous mobility.

Manuscript received: July 5, 2022; Revised: Sept 29, 2022; Accepted: Nov 17, 2022. This paper was recommended for publication by Editor Markus Vincze upon evaluation of the Associate Editor and Reviewers' comments.

¹Matthias Zeller is with CARIAD SE and with the University of Bonn, Germany. ²Jens Behley is with the University of Bonn, Germany. ³Michael Heidingsfeld is with CARIAD SE, Germany. ⁴Cyrill Stachniss is with the University of Bonn, Germany, with the Department of Engineering Science at the University of Oxford, UK, and with the Lamarr Institute for Machine Learning and Artificial Intelligence, Germany

This is a preprint of the article accepted at Robotics and Automation Letters (RA-L). © 2022 IEEE.

Digital Object Identifier (DOI) 10.1109/LRA.2022.3226030

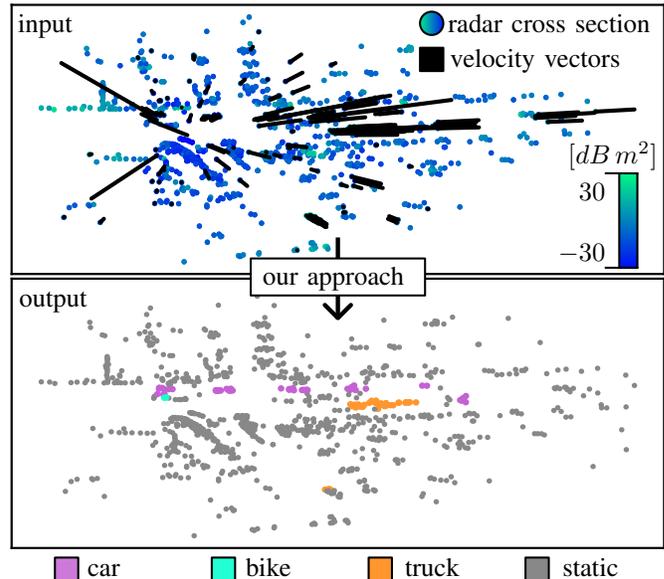


Fig. 1: Our method performs semantic segmentation of moving objects (bottom) from 4D sparse, single-scan radar point clouds (top) exploiting additional information including the velocity and the radar cross section. In the bottom image, each color represents a different semantic class for moving objects (static is grey).

In this work, we investigate the semantic segmentation of moving objects in radar point clouds. This task requires differentiating between detections of moving and static objects and assigning a class label to each radar detection, as illustrated in Fig. 1. Compared to LiDAR point clouds, radar point clouds are noisier due to sensor noise and multipath propagation and more sparse. However, radar sensors provide additional information such as the relative velocity to directly indicate moving objects, making the sensor inherently suitable for single-scan processing. Furthermore, the radar cross section depends on the structure, material, and surface of the reflections, which helps to differentiate objects.

Most state-of-the-art methods for estimating semantics from radar data [23], [26] strongly rely on the aggregation of information over multiple scans to accurately perform semantic segmentation. However, aggregation inherently introduces latency, making it unsuitable for tasks requiring immediate information about the vehicle's vicinity, such as collision avoidance. Therefore, this work investigates the processing of single scans by exploiting the additional information provided by radar sensors.

The main contribution of this paper is a new method for accurate, single-scan, radar-only semantic segmentation of moving objects. It takes sparse point cloud representations

of radar scans as input and outputs a semantic label for each point. To extract discriminative point-wise features, we build on the self-attention mechanism, a fully attentive neural network with our novel Gaussian transformer layer, and our attentive up- and downsampling modules as central building blocks. We optimize the transformer layer and enable the decoupling via the usage of a Gaussian. Furthermore, our attentive sampling enables the capturing of complex local structures and progressively increases the receptive field of individual points. We combine these building blocks in our new backbone, called Gaussian radar transformer, to enhance feature extraction on sparse and noisy radar point clouds.

In sum, we make three key claims: Firstly, our approach demonstrates state-of-the-art performance for semantic segmentation of moving objects in sparse, single-scan radar point clouds without aggregating multiple scans and without exploiting temporal dependencies. Secondly, the Gaussian transformer layer and the attentive up-and downsampling modules improve feature extraction by decoupling individual points and enlarging the receptive field to enhance accuracy. Thirdly, our fully attentive network is able to extract discriminable features from additional sensor information such as Doppler velocity and radar cross section.

II. RELATED WORK

There is extensive literature on semantic segmentation of point clouds, mostly, however, working on LiDAR data. The works can be categorized into projection-based, voxel-based, point-based, and hybrid methods [8].

Projection-based methods are inspired by the successful convolutional neural networks (CNNs) [13], [16]. SqueezeSeg [33], SqueezeSegV2 [32], RangeNet++ [17], and SalsaNext [4] project the point cloud into frontal view images or 2D range images to exploit 2D convolutions. Milioto et al. [17] further alleviate the problem of blurry CNN output and discretization errors by efficient GPU-enabled projective nearest neighbor search as a post-processing step to enhance segmentation results. However, projection-based methods face several problems due to intermediate representation including discretization errors and occlusion.

Voxel-based. To maintain the 3D geometric information between the data points, voxel-based encoding can be used. VoxSegNet [31] voxelizes the point clouds as dense cuboids and leverages atrous 3D convolution and attention-based aggregation to enhance feature extraction under limited resolution. Since outdoor point clouds are sparse and vary in density, just a small percentage of voxels are occupied. This makes it inefficient to apply dense convolution neural networks. To reduce the computational burden, Graham et al. [6] propose sub-manifold sparse convolutional networks which only generate outputs for occupied voxels. Following Polarnet [43], Zhu et al. [47] introduce the cylindrical partitioning, which does not alter the 3D topology compared to the 2D approach, and processes the features by asymmetrical 3D convolution networks. The advancement in 3D point cloud processing has led to state-of-the-art results of (AF)²-S3Net [3] and RPVNet [39] in the SemanticKITTI LiDAR point cloud semantic segmentation

benchmark [1]. Xu et al. [39] combine the voxel-based method with point- and projection-based encoding, utilizing a gated fusion module to adaptively merge the features leading to a hybrid approach. Since voxel-based methods inherently introduce discretization artifacts and information loss, the hybrid method utilizes point-wise information to alleviate the lossy encoding of information.

Point-based. To leverage the full potential of 3D points, especially for sparse point clouds, and keep the geometric information intact, point-based methods [12], [20], [27] have been introduced. The pioneering work of Qi et al. [20] consumes point clouds directly by shared multi-layer perceptrons (MLPs) and aggregates nearby information by symmetrical pooling functions. The successor PointNet++ [21] groups points hierarchically and progressively extracts features from larger local regions. Schumann et al. [24] adapt the approach and optimize the network for sparse radar point cloud processing. However, the ability to capture local 3D structures is limited, especially in sparse point clouds. To circumvent, Schumann et al. [26] aggregate scans, include additional features, or exploit strong temporal relationships. To combine local features and reduce the computational cost point-based methods benefit from effective sampling strategies [10], [35], [40], [41]. The most frequently used methods for small-scale point clouds are farthest point sampling [21] and inverse density sampling [35].

Another approach to learning per-point local features is kernel-based convolutions. PointConv [35] uses an MLP whereas KPConv [27] defines an explicit convolution to directly learn the kernel. Nobis et al. [18] extended KPConv [27] and exploit the time dimension of multiple radar scans to perform object detection. Another method to elaborate a stronger connection of the individual points is graph-based, conducting message passing on the constructed graphs [12]. PointWeb [46] uses adaptive feature adjustment to represent regions and capture local interactions. However, graph-based networks capture edge relationships of local patches which are invariant to the deformation of these. Velickovic et al. [29] and Wang et al. [30] utilize the self-attention mechanism which is inherently permutation invariant to leverage the limitations and further improve the accuracy.

Self-attention models have revolutionized natural language processing [5], [28] and inspired self-attention modules for image recognition [11], [22], [44] and point cloud processing [36], [41]. Recent point transformer networks [7], [37], [42], [45] enhance state-of-the-art performance for 3D point cloud understanding by elaborating the self-attention mechanism. PCT [7] proposes offset-attention to sharpen the attention weights by element-wise subtraction of the self-attention features and the input features. Point Transformer uses the vector-based subtraction attention [44] to aggregate local features whereas Stratified Transformer applies dot-product attention and increases the effective receptive field by a window-based key-sampling strategy. Furthermore, recent work elaborates positional encoding to enhance accuracy and keep position information throughout the network [14].

In contrast to the related work, we propose a novel architecture inspired by self-attention and point transformers. With our newly introduced Gaussian Radar Transformer we are

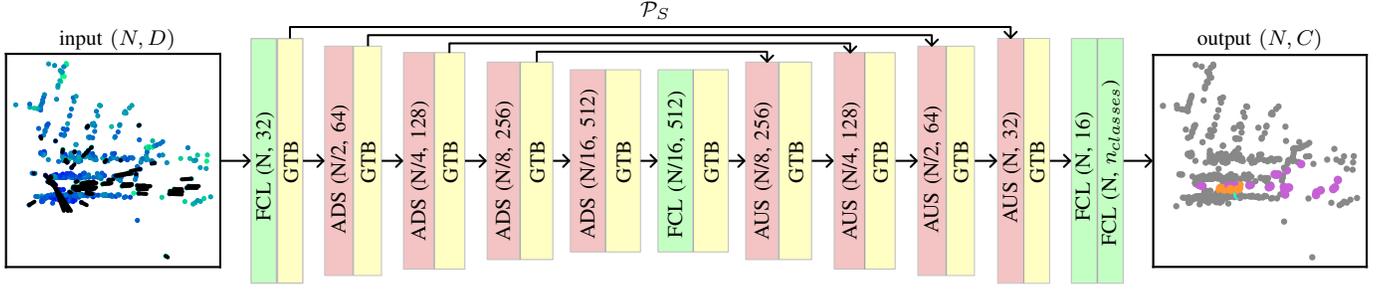


Fig. 2: The architecture of our Gaussian Radar Transformer for semantic segmentation of moving objects. FCL: fully connected layer, ADS: attentive downsampling, AUS: attentive upsampling, GTB: Gaussian transformer block

able to capture complex structures in sparse point clouds and further extend the capabilities of the self-attention mechanism. Furthermore, our proposed fully attentive network includes advanced sampling strategies and substantially enhances state-of-the-art performance for semantic segmentation of moving objects in radar point clouds.

III. OUR APPROACH

The goal of our approach is to achieve accurate semantic segmentation of moving objects in single-scan, sparse radar point clouds to enhance scene understanding of autonomous vehicles. To accomplish this, we introduce a point-based framework to directly process the input point cloud to omit information loss, and builds upon the successful self-attention mechanism throughout the network. Fig. 2 depicts our Gaussian Radar Transformer (GRT). We adopt the encoder-decoder structure of the Point Transformer [45]. We replace each module and use our Gaussian transformer layer as the central building block of each stage, which enables decoupled fine-grained feature aggregation. Furthermore, we introduce attentive up- and downsampling modules to enlarge the receptive field and extract discriminative features.

A. Transformers

Before presenting our contribution, we shortly revisit transformers as they are a key ingredient in our work. Transformers and self-attention networks rely on the encoded representation of the input features $\mathbf{x}^F \in \mathbb{R}^D$ within the queries \mathbf{q} , the keys \mathbf{k} , and the values \mathbf{v} , as follows:

$$\mathbf{q} = \mathbf{W}_q \mathbf{x}^F, \quad \mathbf{k} = \mathbf{W}_k \mathbf{x}^F, \quad \mathbf{v} = \mathbf{W}_v \mathbf{x}^F, \quad (1)$$

where $\mathbf{W}_q \in \mathbb{R}^{D \times D}$, $\mathbf{W}_k \in \mathbb{R}^{D \times D}$ and $\mathbf{W}_v \in \mathbb{R}^{D \times D}$ are the corresponding learned matrices of fully connected layers or multi-layer perceptrons (MLPs). To calculate the attention scores $\mathbf{A}_{i,j}$, different methods exist such as scalar dot-product [28] and vector attention [44]. The scaling by the factor d_C is intended to counteract the effect of small gradients for the softmax if it grows large in magnitude and is defined as follows:

$$\mathbf{A}_{i,j} = \text{softmax} \left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_C}} \right). \quad (2)$$

There is an alternative way for the weighting of individual feature channels by vector attention that utilizes relation

functions f such as addition or subtraction. To keep fine-grained position information throughout the network, Wu et al. [34] and Zhao et al. [45] use relative positional encoding $\mathbf{r}_{i,j} = \mathbf{p}_i - \mathbf{p}_j$, $1 \leq i, j \leq N_l$. The final attention weights $\mathbf{A}_{i,j}$ are determined by the softmax function:

$$\mathbf{A}_{i,j} = \text{softmax}(f(\mathbf{q}_i, \mathbf{k}_j) + \mathbf{r}_{i,j}). \quad (3)$$

Since global self-attention leads to unacceptable memory consumption and computational cost the inputs are restricted to local areas with N_l points determined by farthest point sampling and k nearest neighbor (k NN) [21], [45]. The intermediate representation \mathbf{y}_j utilizing vector attention is calculated as follows:

$$\mathbf{y}_j = \sum_{i=1}^{N_l} \mathbf{A}_{i,j} \odot \mathbf{v}_i. \quad (4)$$

The aggregated features \mathbf{y} are processed by an MLP with a learnable weight matrix $\mathbf{W}_y \in \mathbb{R}^{D \times D}$:

$$\mathbf{o} = \mathbf{W}_y \mathbf{y}, \quad (5)$$

to calculate the final output \mathbf{o} .

B. Gaussian Transformer Layer

In sparse radar point clouds, individual reflections contain essential information for downstream tasks such as semantic segmentation of moving objects. To enable independent and precise feature aggregation, we introduce the Gaussian transformer layer (GTL) based on the Point Transformer layer [45] including vector self-attention as illustrated in Fig. 3 (a). Contrary to other approaches, including the Point Transformer [45], which focuses on dense point clouds, we do not utilize the softmax function, which is defined as:

$$s_i = \frac{\exp(z_i)}{\sum_{j=1}^{N_l} \exp(z_j)}. \quad (6)$$

The softmax function leads to a coupling of points since individual outputs s_i are dependent on all inputs z_j with $j \in \{1, \dots, N_l\}$, which is why the softmax function is also not scale invariant (weighting for dot-product attention Sec. III-A). Furthermore, the backpropagation of the loss \mathcal{L} through the softmax function to obtain the partial derivative $\frac{\partial \mathcal{L}}{\partial z_j}$ to determine the gradients at the input is dependent on all output values. The calculation of the chain rule of derivatives for the

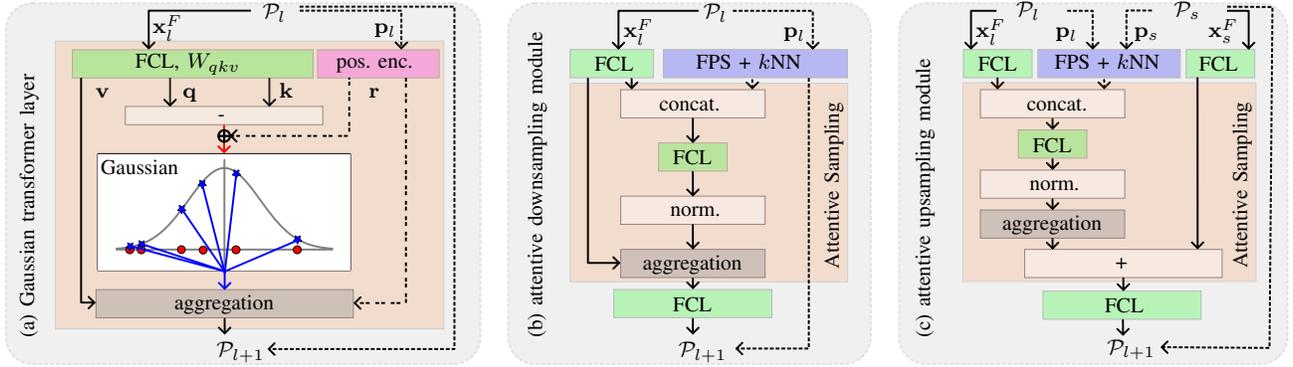


Fig. 3: The detailed design of each module of our Gaussian Radar Transformer (a) shows the Gaussian transformer layer, (b) the attentive downsampling module, and (c) the attentive upsampling module. FCL: fully connected layer, pos. enc.: positional encoding, concat.: concatenation, norm.: normalization, FPS: farthest point sampling, k NN: k nearest neighbor

softmax can be expressed by the Jacobian matrix $\mathbf{J}_{\text{softmax}}$ as follows:

$$\frac{\partial \mathcal{L}}{\partial z} = \mathbf{J}_{\text{softmax}} \frac{\partial \mathcal{L}}{\partial s}. \quad (7)$$

If the output values grow in magnitude the gradients diminish since the Jacobian converges to a zero matrix. Hence, the error propagation is restricted, which slows down the learning process. In contrast, we argue that points belonging to the same class should aggregate the information, whereas points belonging to different classes reduce the information aggregation to a minimum, both of which can lead to a close to zero Jacobian matrix. To overcome this limitation, we replace the softmax function in Eq. (3) by a Gaussian function G , which is executed on every dimension of the vector, for vector self-attention:

$$\mathbf{A}_{i,j} = G(f(\mathbf{q}_i, \mathbf{k}_j) + \mathbf{r}_{i,j}), \quad (8)$$

to assess fine-grained information flow for sparse radar point clouds. Since the Gaussian function is applied to each feature individually, the points are decoupled, which enables a precise information aggregation to enhance feature extraction and performance. Moreover, the partial derivative of the Gaussian depends on a single output value s_j . Hence, vanishing gradients may influence individual points but not whole local areas, which can be seen by the chain rule:

$$\frac{\partial \mathcal{L}}{\partial z_j} = \frac{\partial \mathcal{L}}{\partial s_j} \frac{\partial s_j}{\partial z_j}. \quad (9)$$

To derive the output \mathbf{o} , we calculate the sum of the element-wise multiplication:

$$\mathbf{o}_j = \sum_{i=1}^{N_i} \mathbf{A}_{i,j} \odot \mathbf{v}_i, \quad (10)$$

without further processing by a linear layer reducing computational cost in contrast to Eq. (5). Following Qi et al. [21] and Zhao et al. [45], we determine the local areas by farthest point sampling and k NN algorithm with $k = N_l$. We directly derive the queries \mathbf{q}_i , the keys \mathbf{k}_i , and the values \mathbf{v}_i by applying a fully connected layer with weight matrix $\mathbf{W}_{qkv} \in \mathbb{R}^{D \times 3D}$. For the positional encoding, we adopt the approach of Zhao et al. [45]. We process the relative position by two fully connected layers and replace the activation

function by the Gaussian error linear unit (GELU) [9] to determine the positional encoding.

C. Gaussian Transformer Block

Our Gaussian transformer layer (GTL) is embedded into the center of the Gaussian transformer block (GTB) which is a residual block, similar to the Point Transformer block [45], with two fully connected layers processing the input and the output. We replace the activation function with GELU after each fully connected layer. The GTB processes point clouds \mathcal{P} with point coordinates $\mathbf{p}_i \in \mathbb{R}^2$ and point-wise features ($\mathcal{X}^F = \{\mathbf{x}_1^F, \dots, \mathbf{x}_N^F\}$), where $\mathbf{x}_i^F \in \mathbb{R}^D$ with feature dimension D . The features of the individual points \mathbf{x}_i^F are enriched by the information aggregation within the block enhanced by the GTL. The point coordinates \mathbf{p}_i are utilized to calculate the attention weights but not further transformed to keep detailed position information.

D. Attentive Downsampling Layer

To reduce the cardinality of the point cloud $\mathcal{P}_{l+1} \subset \mathcal{P}_l$ and thereby the number of points N , we process the point cloud by the attentive downsampling layer, depicted in Fig. 3 (b). Our approach aims to enable adequate sampling and feature processing by applying the self-attention mechanism throughout the network. To reduce computational complexity, we follow Yang et al. [40] and calculate the attention weights by a single feed-forward layer with the weight matrix $\mathbf{W}_f \in \mathbb{R}^{(D+2) \times D}$ and no direct representation of keys, queries, and values. We concatenate the input features \mathbf{x}_i^F and the point coordinates \mathbf{p}_i to include positional information to calculate the attention weights $\mathbf{A}_{i,j}$. Additionally, we normalize the attention weights over the whole point cloud to amplify the contribution of valuable points. The final weights are multiplied with the input features \mathbf{x}^F within local areas which are determined by farthest point sampling and the k NN algorithm [21], with $k = N_d$ resulting in:

$$\mathbf{y}_i = \sum_{j=1}^{N_d} \mathbf{A}_{i,j} \odot \mathbf{x}_j^F. \quad (11)$$

The features are fed into another fully connected layer with LayerNorm [38] and a GELU activation function. In

contrast to Point Transformer [45], which utilizes farthest point sampling and max pooling [21], our attentive downsampling includes the information of nearby points, which we assume as valuable for sparse point clouds.

E. Attentive Upsampling Layer

To deduce discriminative features, we argue that the up-sampling and feature concatenation of the skip connection is crucial to further enhance performance. The common method for up-sampling, also utilized by Point Transformer [45], is an interpolation of the $k = 3$ nearest neighbors based on an inverse distance weighted average [21]. The interpolated points N_u are concatenated with the features of the points, which are passed through the skip connection. The inverse distance weighted average does not include further feature-based information. Hence, the interpolation combines the features only based on their relative position. This is reasonable for dense point clouds because nearby points often belong to the same class.

However, this might be problematic for sparse point clouds, especially for small instances, which are represented by single points. Therefore, we consider upsampling as an important part to improve feature extraction and propose the attentive upsampling layer. The upsampling layer, which is illustrated in Fig. 3 (c), first processes the features of the skip connection and the proceeding GTB by two separate fully connected layers with LayerNorm and GELU activation function. To propagate the points from \mathcal{P}_l to \mathcal{P}_{l+1} where $\mathcal{P}_l \subset \mathcal{P}_{l+1}$ with $N_l \leq N_{l+1}$, we feed the position information of the two point sets and the corresponding features into our attentive upsampling layer. We calculate the k nearest neighbors of the individual points for the point set of the skip connection \mathcal{P}_s within the point cloud which has to be upsampled \mathcal{P}_l . The attention mechanism enables information aggregation of larger local areas since the attention weights will control the information flow and not reduce the discriminability which is possible if large local regions are interpolated. To integrate the positional information we calculate the relative position of the k NN of the two point sets given by:

$$\mathbf{r}_{i,j} = \mathbf{p}_i - \mathbf{p}_j, \quad (12)$$

where $\mathbf{p}_j \in \mathcal{P}_s$ and $\mathbf{p}_i \in \mathcal{P}_l$. The relative distances $\mathbf{r}_{i,j}$ are concatenated with the features. Following our downsampling layer, we calculate the attention weights directly by processing the concatenated features with a fully connected layer and normalizing the weights over the whole point cloud. The output of the summation is processed by a fully connected layer with LayerNorm and GELU activation function. The self-attention mechanism turns into an inter-attention between the two point clouds to enable attentive feature aggregation. The upsampling is repeated until we have broadcasted the features to the original set of points. We optimize the information aggregation by determining the weighting based on the relative position and the features. We emphasize that the sampling steps are essential for appropriate feature extraction of transformer architectures for sparse point clouds.

F. Input Features

The input is a sparse radar point cloud with N points, feature dimension D , and batch size b . Each point \mathbf{p}_i is defined by two spatial coordinates x_i, y_i . Additionally, the radar sensors provide the ego-motion compensated Doppler velocity v_i and the radar cross section σ_i resulting in a 4-dimensional input vector $\mathbf{x}_i^F = (x_i, y_i, v_i, \sigma_i)^\top$.

IV. IMPLEMENTATION DETAILS

We construct our architecture based on the self-attention mechanism. The central building blocks are the GTL and the attentive down- and upsampling modules to extract discriminative features for point cloud understanding. The backbone adopts the U-Net architecture of Point Transformer [45] with an encoder-decoder architecture including skip connections. First, we directly extract features of the sparse input point cloud by a GTB and increase the per-point feature dimension to 32. The resulting features are progressively down-sampled by four consecutive stages where each reduces the cardinality of the point cloud by a factor of two resulting in $[N/2, N/4, N/8, N/16]$ points. The per-point features are further gradually increased to 64, 128, 256, and 512. The individual stages include the GTB and attentive downsampling modules in the encoder part, which are replaced by attentive upsampling modules in the decoder part of the network. The per-point features maps of the final decoder layer are processed by an MLP with two fully connected layers to obtain point-wise semantic classes $\mathcal{P}^S = \{p_1^S, \dots, p_N^S\}$, where $p_i^S \in \{1, \dots, C\}$.

We implement the Gaussian Radar Transformer in PyTorch [19]. To train the network, we utilize the SGD optimizer with an initial learning rate of 0.05, a momentum of 0.9, and a cosine annealing learning rate scheduler [15]. The batch size b is set to 32. The loss combines the Lovász loss [2] and weighted cross-entropy. We follow Schumann et al. [26] and set the weights of the cross-entropy loss for dynamic objects to 8.0 and for static to 0.5 to account for the class imbalance of the data set. For the attentive sampling operations, we define $k = 9$ for the k NN operation, and for the Gaussian transformer layer, we restrict the local area to $N_l = 16$. We define $G(x)$ as:

$$G(x) = \exp\left(\frac{-x^2}{2}\right), \quad (13)$$

such that for $x = 0$ the attention weight is $G(x) = 1$. Additionally, we apply data augmentation, which includes scaling, rotation around the origin, jitter augmentation of the coordinate features, and instance augmentation.

V. EXPERIMENTAL EVALUATION

The main focus of this work is to enhance the semantic segmentation of moving objects in sparse and noisy radar point clouds. We present our experiments to show the capabilities of our method and to support our key claims that our approach achieves state-of-the-art performance in semantic segmentation of moving objects in single-scan radar point clouds without exploration of temporal dependencies or the

Method	Input	mIoU	F1	IoU						F1					
				static	car	ped.	ped. grp.	bike	truck	static	car	ped.	ped. grp.	bike	truck
RadarPNv1 [24]	aggregation	61.0	74.3	98.7	58.2	36.0	58.7	58.4	56.1	99.4	73.6	52.9	74.0	73.8	71.9
RadarPNv2 [26]		61.9	75.0	98.7	63.8	38.8	58.5	51.0	61.0	99.4	77.9	55.9	73.8	67.5	75.8
Point Voxel Transformer [42]	single-scan	45.9	57.5	99.3	47.5	7.3	47.5	54.6	19.2	99.6	64.4	13.6	64.4	70.6	32.2
Point Transformer [45]		55.6	68.1	99.3	58.1	15.2	56.8	55.1	48.9	99.6	73.5	26.4	72.5	71.1	65.6
Gaussian Radar Transformer		68.5	79.8	99.4	69.6	36.3	71.2	71.2	62.8	99.7	82.1	53.2	83.2	83.2	77.1

TABLE I: Semantic segmentation results of moving objects on the RadarScenes test set in terms of IoU and F_1 scores. The results of RadarPNv1 [24] and RadarPNv2 [26] are calculated based on the reported confusion matrix.

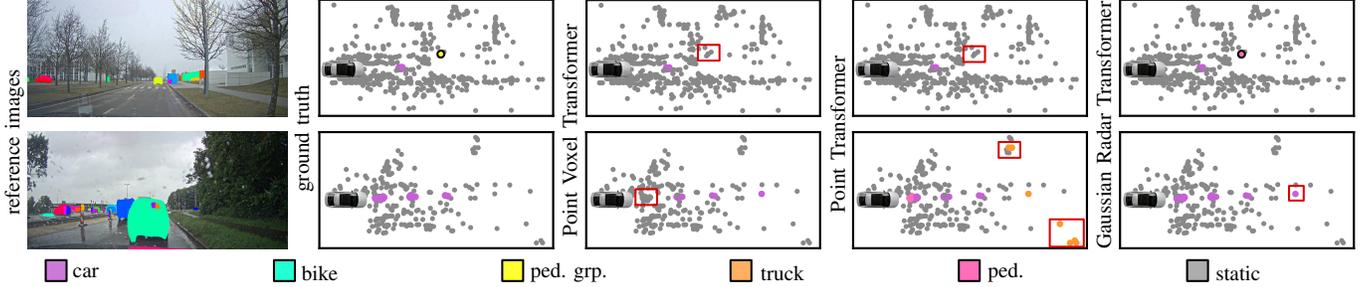


Fig. 4: Qualitative results of Point Transformer [45], Point Voxel Transformer [42], and Gaussian Radar Transformer on the test set of RadarScenes [25].

aggregation of scans. Moreover, we demonstrate that the Gaussian transformer layer and the attentive up- and downsampling modules improve feature extraction and contribute to the final performance. Our fully attentive network is able to extract valuable features from the Doppler velocity and radar cross section provided by the radar sensor.

A. Experimental Setup

We train and evaluate our method on RadarScenes [25], which is the only large-scale, open-source radar data set including point-wise annotations for varying scenarios. The data set consists of 158 annotated sequences. We use the recommended 130 sequences for training and split the remaining 28 sequences into validation (sequences: 6, 42, 58, 85, 99, 122) and test set. The RadarScenes [25] data set is split into separate scans for each of the four sensors. Since the field-of-view of the sensors is restricted to certain areas, we derive detailed information about the surrounding by merging the individual sensor data from the four sensors into a single radar point cloud. The measurement times and the pose information are given, which enables a transformation into a common coordinate system. We aggregate four scans, one of each sensor, which results in the final input point clouds with transformed local coordinates. To evaluate the performance, Schumann et al. [25] propose the point-wise macro-averaged F_1 scores based on all five moving object classes and the static background class ($C = 6$). We further report the $IoU = \frac{TP}{TP+FN+FP}$ and $mIoU = \frac{1}{C} \sum_{i=1}^C IoU_i$ scores, which are common for semantic segmentation tasks [1]. We train each network using its specific hyperparameters with two Nvidia RTX A6000 GPUs over 50 epochs on the training set and report the results on the test set. For more details on the training regime for Point Transformer¹, we refer to the original paper [45].

B. Semantic Segmentation of Moving Objects

The first experiment presents the performance of our approach on the RadarScenes test set to investigate the claim that we achieve state-of-the-art results for semantic segmentation of moving objects in sparse and noisy radar point clouds without the aggregation of scans or the exploration of temporal dependencies. In this experiment, we compare our Gaussian Radar Transformer with the recent and high-performing Point Transformer by Zhao et al. [45] as well as the baselines provided by Schumann et al. [24], [26]. We selected Point Transformer as a reference since the method meets the following requirements: (1) single-scan input for comparability; (2) point-based method, since the voxelization leads to discretization artifacts and hence a loss of information, see Point Voxel Transformer in Tab. I; (3) very good performance on different benchmarks including semantic scene understanding. Furthermore, the Point Transformer [45] utilizes vector attention, which is beneficial for point cloud understanding.

Our Gaussian Radar Transformer outperforms the existing methods in terms of both, mIoU and F_1 score, as displayed in Tab. I. Especially, we achieve superior performance on five of the six classes, except pedestrian. We assume that the individual detection in radar scans contains important information, which is why strict point-based methods enhance the performance compared to Point Voxel Transformer. The baselines exploit temporal dependencies of consecutive radar scans within a memory feature map, utilize additional global coordinates or densify the point clouds by aggregation. The exact comparison of the results is difficult because Schumann et al. work on a subset of the officially released data set. However, the IoU for the class pedestrian indicates that the exploration of temporal information is beneficial for small instances. We suspect that the consistent detection of pedestrians over the whole sequence, which is difficult for strict single-scan approaches, further improves the performance. Nevertheless, the Gaussian Radar Transformer considerably improves the IoU for the class pedestrian as opposed to Point

¹<https://github.com/POSTECH-CVLab/point-Transformer>

#	ADS	AUS	GTL	F_1	mIoU
A				74.0	61.0
B	✓			77.0	64.7
C		✓		77.3	65.5
D	✓	✓		78.8	66.8
E	✓	✓	✓	79.4	68.3

TABLE II: Influence of the different components of the approach in terms of mIoU and F_1 score on the RadarScenes validation set.

Transformer by more than 19 absolute percentage points. Fig. 4 shows some qualitative results on the test set. Notably, our approach achieves superior performance under adverse weather including rain and fog.

C. Ablation Studies on Method Components

The first ablation study presented in this section is designed to support our second claim that our proposed self-attention modules each contribute to the advancements of the Gaussian Radar Transformer. To assess the influence of the different components of our fully attentive backbone, we evaluate the performance in terms of mIoU and F_1 score on the validation set. To replace our proposed modules, we follow commonly used network designs. We substitute the Gaussian function by the softmax function and keep the rest of the Gaussian transformer layer as it is. For the attentive downsampling, we utilize local max pooling and we exchange attentive upsampling by trilinear interpolation based on an inverse distance weighted average. Tab. II summarizes the influence of different components on the performance in terms of mIoU on the validation set.

In configuration (A), we replace each module by its substitute, which leads to a noticeable decrease in mIoU. We suspect that the commonly used modules are highly optimized for denser point clouds but struggle to capture fine-grained information from sparse and noisy radar point clouds. In (B), we add attentive downsampling (ADS), see Sec. III-D, which introduces a smooth information exchange within the downsampling step of individual points, visibly improving the results. In (C), we add the attentive upsampling (AUS) module to enlarge the receptive field and include encoded feature information to optimize the information aggregation, see Sec. III-E. The larger receptive field resulting from the increased local area from three (trilinear) to nine points improves the F_1 score by 3.3 and the mIoU by 4.5 absolute percentage points. Although the AUS only affects the features of the decoder part it leads to an additional improvement of mIoU by 0.8 absolute percentage points compared to AUS in (B). In (D), we add the attentive up- and downsampling which further enhance the performance. This shows the importance of the attentive sampling modules for sparse radar point cloud processing. In (E), we utilize the fully attentive network to illustrate the improvement due to the usage of the Gaussian function by decoupling individual points, see Sec. III-B, resulting in the best performance. In conclusion, the Gaussian function and the attentive up- and downsampling are essential to extract valuable features from sparse and noisy radar point clouds.

Input Features	F_1	mIoU
$x^F = (x, y)$	56.0	43.7
$x^F = (x, y, \sigma)$	63.7	50.1
$x^F = (x, y, v)$	75.0	62.0
$x^F = (x, y, v, \sigma)$	79.4	68.3

TABLE III: Influence of the different input features in terms of mIoU and F_1 score on the RadarScenes validation set.

D. Ablation Studies on Input Features

The third experiment evaluates the performance depending on the provided information by the radar sensor and demonstrates that our approach is capable of capturing complex local structures within the features to enhance mIoU. For this experiment, we utilize our Gaussian Radar Transformer and add to the position information of x and y coordinates, the ego-motion compensated Doppler velocity v , the radar cross section σ , or both. Tab. III displays the influence of the input features x^F on the validation set performance. As we presume, the ego-motion compensated Doppler velocity is especially valuable for semantic segmentation of moving objects since the feature inherently distinguishes between moving and non-moving parts of the environment resulting in an increase of mIoU of 18.2 absolute percentage points. Moreover, we further improve the mIoU if we add the radar cross section features σ suggesting that our approach extracts valuable features for the downstream task from additional sensor information. Hence, the Gaussian Radar Transformer achieves the best performance including radar cross section and ego-motion compensated Doppler velocity.

In summary, our evaluation supports our statement that our method provides competitive semantic segmentation performance of moving objects in single-scan, sparse radar point clouds. At the same time, our method exploits self-attention modules which enhance the performance in multi-dimensional radar data processing outperforming state-of-the-art approaches. Thus, we support all our claims with this experimental evaluation.

VI. CONCLUSION

In this paper, we presented a novel approach to perform semantic segmentation of moving objects in sparse, noisy, single-scan radar point clouds obtained from automotive radars. Our method exploits the self-attention mechanism throughout the network and replaces the softmax normalization of the transformer by a Gaussian. This allows us to successfully segment moving objects and improve the feature extraction by decoupling individual points. We implemented and evaluated our approach on the RadarScenes data set, providing comparisons to other methods and supporting all claims made in this paper. The experiments suggest that the proposed architecture achieves good performance on semantic segmentation of moving objects. We assessed the different parts of our approach and compared them to other existing techniques. Overall, our approach outperforms the state of the art both in F_1 score and mIoU, taking a step forward towards sensor redundancy for semantic segmentation for autonomous robots and vehicles.

REFERENCES

- [1] 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 *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2019.
- [2] 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 *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [3] R. Cheng, R. Razani, E. Taghavi, E. Li, and B. Liu. (af)2-s3net: Attentive feature fusion with adaptive feature selection for sparse semantic segmentation network. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [4] T. Cortinhal, G. Tzelepis, and E. Erdal Aksoy. Salsanext: Fast, uncertainty-aware semantic segmentation of lidar point clouds. In *Proc. of the Int. Symp. on Visual Computing*, 2020.
- [5] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. Le, and R. Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*, 2019.
- [6] B. Graham, M. Engelcke, and L. van der Maaten. 3D Semantic Segmentation with Submanifold Sparse Convolutional Networks. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [7] M.H. Guo, J. Cai, Z.N. Liu, T.J. Mu, R.R. Martin, and S. Hu. Pct: Point cloud transformer. *Comput. Vis. Media*, 7(2):187–199, 2021.
- [8] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Bennamoun. Deep learning for 3d point clouds: A survey. *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 43(12):4338–4364, 2021.
- [9] D. Hendrycks and K. Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. *arXiv preprint:1606.08415*, 2016.
- [10] Q. Hu, B. Yang, L. Xie, S. Rosa, Y. Guo, Z. Wang, N. Trigoni, and A. Markham. Randa-net: Efficient semantic segmentation of large-scale point clouds. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [11] A. Kolesnikov, A. Dosovitskiy, D. Weissenborn, G. Heigold, J. Uszkoreit, L. Beyer, M. Minderer, M. Dehghani, N. Houlsby, S. Gelly, T. Unterthiner, and X. Zhai. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2021.
- [12] L. Landrieu and M. Simonovsky. Large-scale Point Cloud Semantic Segmentation with Superpoint Graphs. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [13] Y. LeCun, L. Jackel, L. Bottou, A. Brunot, C. Cortes, J. Denker, H. Drucker, I. Guyon, U. Muller, E. Säcker, P. Simard, and V. Vapnik. Comparison of learning algorithms for handwritten digit recognition. In *Proc. of the Int. Conf. on Artificial Neural Networks*, 1995.
- [14] Y. Li, S. Si, G. Li, C.J. Hsieh, and S. Bengio. Learnable fourier features for multi-dimensional spatial positional encoding. In *Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)*, 2021.
- [15] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2017.
- [16] A. Milioto and C. Stachniss. Bonnet: An Open-Source Training and Deployment Framework for Semantic Segmentation in Robotics using CNNs. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2019.
- [17] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss. RangeNet++: Fast and Accurate LiDAR Semantic Segmentation. In *Proceedings of the IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, 2019.
- [18] F. Nobis, F. Fent, J. Betz, and M. Lienkamp. Kernel point convolution lstm networks for radar point cloud segmentation. *Applied Sciences*, 11:2599–2618, 2021.
- [19] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, S. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)*, 2019.
- [20] C.R. Qi, H. Su, K. Mo, and L.J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [21] C. Qi, K. Yi, H. Su, and L.J. Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. In *Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)*, 2017.
- [22] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens. Stand-alone self-attention in vision models. In *Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)*, 2019.
- [23] N. Scheiner, F. Kraus, F. Wei, B. Phan, F. Mannan, N. Appenrodt, W. Ritter, J. Dickmann, K. Dietmayer, B. Sick, and F. Heide. Seeing around street corners: Non-line-of-sight detection and tracking in-the-wild using doppler radar. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [24] O. Schumann, M. Hahn, J. Dickmann, and C. Wöhler. Semantic segmentation on radar point clouds. In *Proc. of the Int. Conf. on Information Fusion*, 2018.
- [25] O. Schumann, M. Hahn, N. Scheiner, F. Weishaupt, J.F. Tilly, J. Dickmann, and C. Wöhler. Radarscenes: A real-world radar point cloud data set for automotive applications. In *Proc. of the Int. Conf. on Information Fusion*, 2021.
- [26] O. Schumann, J. Lombacher, M. Hahn, C. Wöhler, and J. Dickmann. Scene understanding with automotive radar. *IEEE Trans. on Intelligent Vehicles*, 5(2):188–203, 2020.
- [27] H. Thomas, C. Qi, J. Deschaut, B. Marcotegui, F. Goulette, and L. Guibas. KPConv: Flexible and Deformable Convolution for Point Clouds. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2019.
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)*, 2017.
- [29] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph Attention Networks. In *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2018.
- [30] X. Wang, S. Liu, X. Shen, C. Shen, and J. Jia. Associatively Segmenting Instances and Semantics in Point Clouds. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [31] Z. Wang and F. Lu. Voxsegnet: Volumetric cnns for semantic part segmentation of 3d shapes. *IEEE Trans. on Visualization and Computer Graphics*, 26(9):2919–2930, 2019.
- [32] 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 *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2019.
- [33] 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 *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2018.
- [34] K. Wu, H. Peng, M. Chen, J. Fu, and H. Chao. Rethinking and improving relative position encoding for vision transformer. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021.
- [35] W. Wu, Z. Qi, and L. Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [36] S. Xie, S. Liu, Z. Chen, and Z. Tu. Attentional shapecontextnet for point cloud recognition. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [37] L. Xin, L. Jianhui, J. Li, W. Liwei, Z. Hengshuang, L. Shu, Q. Xiaojuan, and J. Jiaya. Stratified transformer for 3d point cloud segmentation. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [38] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T. Liu. On layer normalization in the transformer architecture. In *Proc. of the Int. Conf. on Machine Learning (ICML)*, 2020.
- [39] 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 *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021.
- [40] B. Yang, S. Wang, A. Markham, and N. Trigoni. Robust attentional aggregation of deep feature sets for multi-view 3d reconstruction. *Intl. Journal of Computer Vision (IJCV)*, 128(1):53–73, 2020.
- [41] J. Yang, Q. Zhang, B. Ni, L. Li, J. Liu, M. Zhou, and Q. Tian. Modeling point clouds with self-attention and gumbel subset sampling. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [42] C. Zhang, H. Wan, X. Shen, and Z. Wu. Pvt: Point-voxel transformer for point cloud learning. *arXiv preprint*, 2019.
- [43] 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 *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020.

- [44] H. Zhao, J. Jia, and V. Koltun. Exploring self-attention for image recognition. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [45] H. Zhao, L. Jiang, J. Jia, P.H. Torr, and V. Koltun. Point transformer. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021.
- [46] T. Zhao, Y. Xu, M. Monfort, W. Choi, C. Baker, Y. Zhao, Y. Wang, and Y. Wu. Multi-Agent Tensor Fusion for Contextual Trajectory Prediction. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [47] X. Zhu, H. Zhou, T. Wang, F. Hong, Y. Ma, W. Li, H. Li, and D. Lin. Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021.