

# LatticeNet: Fast Spatio-Temporal Point Cloud Segmentation Using Permutohedral Lattices

Radu Alexandru Rosu · Peer Schütt · Jan Quenzel · Sven Behnke

Received: 14 February 2021 / Accepted: 26 June 2021

**Abstract** Deep convolutional neural networks (CNNs) have shown outstanding performance in the task of semantically segmenting images. Applying the same methods on 3D data still poses challenges due to the heavy memory requirements and the lack of structured data. Here, we propose LatticeNet, a novel approach for 3D semantic segmentation, which takes raw point clouds as input. A PointNet describes the local geometry which we embed into a sparse permutohedral lattice. The lattice allows for fast convolutions while keeping a low memory footprint. Further, we introduce DeformSlice, a novel learned data-dependent interpolation for projecting lattice features back onto the point cloud. We present results of 3D segmentation on multiple datasets where our method achieves state-of-the-art performance. We also extend and evaluate our network for instance and dynamic object segmentation.

**Keywords** semantic segmentation · instance segmentation · motion segmentation · sequence segmentation · 3D point cloud

## 1 Introduction

Environment understanding is a crucial ability for autonomous agents. Perceiving not only the geometrical structure of the scene but also distinguishing between different

This work has been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy - EXC 2070 - 390732324 and by the German Federal Ministry of Education and Research (BMBF) in the project "Kompetenzzentrum: Aufbau des Deutschen Rettungsrobotik-Zentrums" (A-DRZ).

Radu Alexandru Rosu  
Friedrich-Hirzebruch-Allee 8  
Tel.: +49 (0) 228 73-54159  
Fax: +49 (0) 228 73-4425  
E-mail: [rosu@ais.uni-bonn.de](mailto:rosu@ais.uni-bonn.de)

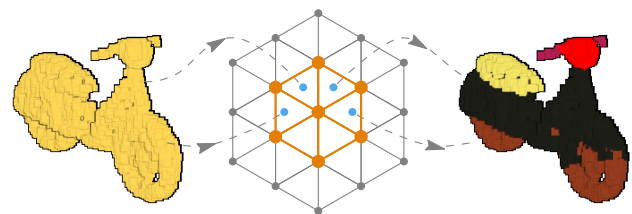


Fig. 1: Semantic segmentation: LatticeNet takes raw point clouds as input and embeds them into a sparse lattice where convolutions are applied. Features on the lattice are projected back onto the point cloud to yield a final segmentation.

classes of objects therein enables tasks like manipulation and interaction that were previously not possible. Within this field, semantic segmentation of 2D images is a mature research area, showing outstanding success in dense per pixel categorization on images [22, 6, 21]. However, the task of semantically labelling 3D data is still an open area of research as it poses several challenges that need to be addressed.

First, 3D data is often represented in an unstructured manner — unlike the grid-like structure of images. This raises difficulties for current approaches which assume a regular structure upon which convolutions are defined.

Second, the performance of current 3D networks is limited by their memory requirements. Storing 3D information in a dense structure is prohibitive for even high-end GPUs, clearly indicating the need for a sparse structure.

Third, discretization issues caused by imposing a regular grid onto point clouds can negatively affect the network's performance and interpolation is necessary to cope with quantization artifacts [42].

In this work, we propose LatticeNet, a novel approach for point cloud segmentation which alleviates the previously mentioned problems. Hence, our contributions are:

- a hybrid architecture which leverages the strength of PointNet to obtain low-level features and sparse 3D convolutions to aggregate global context,
- a framework suitable for sparse data onto which all common CNN operators are defined, and
- a novel slicing operator that is end-to-end trainable for mapping features of a regular lattice grid back onto an unstructured point cloud.

In addition to our *Robotics: Science and System* conference paper [36] we make the following additional contributions:

- an extension with discriminative loss that allows LatticeNet to perform instance segmentation, and
- a network architecture capable of processing temporal information in order to improve semantic segmentation and to distinguish between dynamic and static objects within the scene.

## 2 Related Work

### 2.1 Semantic Segmentation

3D Semantic segmentation approaches can be categorized depending on data representation upon which they operate.

**Point cloud networks:** The first category of networks operates directly on the raw point cloud.

From this area, PointNet [30] is one of the pioneering works. The method processes raw point clouds by individually embedding the points into a higher-dimensional space and applying max-pooling for permutation-invariance to obtain a global scene descriptor. The descriptor can be used for both classification and semantic segmentation. However, PointNet does not take local information into account which is essential for the segmentation of highly-detailed objects. This has been partially solved in the subsequent work of PointNet++ [32] which applies PointNet hierarchically, capturing both local and global contextual information.

Chen et al. [7] use a similar approach but they input the point responses w.r.t. a sparse set of radial basis functions (RBF) scattered in 3D space. Optimizing jointly for the extent and center of the RBF kernels allows to obtain a more explicit modelling of the spatial distribution.

PointCNN [20] deals with the permutation invariance not by using a symmetric aggregation function, but by learning a  $K \times K$  matrix for the  $K$  input points that permutes the cloud into a canonical form.

**Voxel networks:** 3D Convolutions in this category work on discretized cubic or tetrahedral volume elements.

SEGCloud [42] voxelizes the point cloud into a uniform 3D grid and applies 3D convolutions to obtain per-voxel class probabilities. A conditional random field (CRF) is used to smooth the labels and enforce global consistency. The class

scores are transferred back to the points using trilinear interpolation. The usage of a dense grid results in high memory consumption while our approach uses a permutohedral lattice stored sparsely. Additionally, their voxelization results in a loss of information due to the discretization of the space. We avoid quantization issues by using a PointNet architecture to summarize the local neighborhood.

Rethage et al. [34] perform semantic segmentation on a voxelized point cloud and employ a PointNet architecture as a low-level feature extractor. The usage of a dense grid, however, leads to high memory usage and slow inference, requiring various seconds for medium-sized point clouds.

SplatNet [39] is the work most closely related to ours. It alleviates the computational burden of 3D convolutions by using a sparse permutohedral lattice, performing convolutions only around the surfaces. It discretizes the space in uniform simplices and accumulates the features of the raw point cloud onto the vertices of the lattice using a splatting operation. Convolutions are applied on the lattice vertices and a slicing operation barycentrically interpolates the features of the vertices back onto the point cloud. A series of splat-conv-slice operations are applied to obtain contextual information. The main disadvantage is that splat and slice operations are not learned and repeated application slowly degrades the point clouds features as they act as Gaussian filters [2]. Furthermore, storing high-dimensional features for each point in the cloud is memory intensive which limits the maximum number of points that can be processed. In contrast, our approach has learned operations for splatting and slicing which brings more representational power to the network. We also restrict their usage to only the beginning and the end of the network, leaving the rest of the architecture fully convolutional.

**Mesh networks:** The connectivity of triangular or quadrilateral mesh faces enables easy computation of normal vectors and establishes local tangent planes.

GCNN [23] operates on small local patches which are convolved using a series of rotated filters, followed by max-pooling to deal with the ambiguity in the patch orientation. However, the max-pooling disregards the orientation. MoNet [24] deals with the orientation ambiguity by aligning the kernels to the principal curvature of the surface. Yet, this does not solve cases in which the local curvature is not informative, e.g. for walls or ceilings. TextureNet [19] further improves on the idea by using a global 4-RoSy orientations field. This provides a smooth orientation field at any point on the surface which is aligned to the edges of the mesh and has only a 4-direction ambiguity. Defining convolution on patches oriented according to the 4-RoSy field yields significantly improved results.

**Graph networks:** These methods allow arbitrary topologies to connect vertices and lift the restriction of triangular or quadrilateral meshes.

Wang et al. [44] and Wu et al. [48] define a convolution operator over non-grid structured data by having continuous values over the full vector space. The weights of these continuous filters are parametrized by an multi-layer perceptron (MLP).

Defferrard et al. [13] formulate CNNs in the context of spectral graph theory. They define the convolution in the Fourier domain with Chebyshev polynomials to obtain fast localized filters. However, spectral approaches are not directly transferable to a new graph as the Fourier basis changes. Additionally, the learned filters are rotation invariant which can be seen as a limitation to the representational power of the network.

**Multi-view networks:** The convolution operation is well defined in 2D and hence, there is an interest in casting 3D segmentation as a series of single-view segmentations which are fused together.

Pham et al. [28] simultaneously reconstruct the scene geometry and recover the semantics by segmenting sequences of RGB-D frames. The segmentation is transferred from 2D images to the 3D world and fused with previous segmentations. A CRF finally resolves noisy predictions.

TangentConv [41] assumes that the data is sampled from locally Euclidean surfaces and project the local surface geometry onto a tangent plane to which 2D convolutions can be applied. This requires a heavy preprocessing for normal calculation. In contrast, our approach can deal with raw point clouds without requiring normals.

## 2.2 Motion Segmentation

For the task of motion segmentation two approaches have been widely used: Networks either incorporate multiple point clouds directly or accumulate a sequence of individually segmented point clouds.

Shi et al. [37] present their U-Net based architecture SpSequenceNet for semantic segmentation on 4D point clouds. They input two point clouds and generate the output for the later one with a voxel-based method. They designed two modules, the Cross-frame Global Attention (CGA) and the Cross-frame Local Interpolation (CLI) module. The CGA acts as a teacher that uses the data from  $P_{t-1}$  to focus the network on the important features of  $P_t$ . The CLI module fuses information between both point clouds by combining the spatial and temporal information.

Kernel Point Convolution (KPConv) [43] operates directly on the point clouds by facilitating convolution weights that are located in Euclidean space. Points in the vicinity of these kernels are weighted and summed together to feature vectors. KPConv [43], DarkNet53Seg [4] and TangentConv [41] were previously used for the segmentation of 4D point clouds by accumulating multiple clouds of a sequence.

## 2.3 Instance Segmentation

Researchers extended principles from 2D to obtain instances in 3D which can be roughly categorized in proposal-based and proposal-free methods.

**Proposal-based:** This type solves the problem in two stages. The first network stage generates proposals of bounding boxes for the objects in the scene. A second stage performs foreground-background segmentation on the points within the bounding boxes in order to get valid instances.

Yang et al. [50] present a single-stage method for instance segmentation that can train both the proposal and the point-mask prediction network in an end-to-end manner. Yi et al. [52] alleviate some of the issues associated with wrong bounding box predictions by using an analysis-by-synthesis strategy.

**Proposal-free:** Proposal-free methods tackle instance segmentation without the need of generating object proposals. They usually rely on predicting point embedding and apply clustering to recover the instances.

Many proposal-free approaches base their work on the 2D instance segmentation of Brabandere et al. [12] in which pixel embeddings are predicted. There, a discriminative loss encourages the embeddings that belong to the same instance to be clustered together while embeddings from different instances should be further apart.

SPGN [45] learns a similarity matrix for all point pairs, based on which, similar points are merged to instances. VoteNet [31] uses a Hough voting mechanism where the points predict the offset towards the object center. A clustering algorithm finally recovers the object instances.

Neven et al. [25] alleviate some of the issues associated with proposal-free methods by allowing also the clustering algorithm to be part of the training by jointly optimizing the spatial embeddings and the clustering bandwidth.

Wang et al. [46] proposed a framework that allows for semantic and instances to be predicted simultaneously and for the two tasks to mutually benefit from each other. Similarly, Pham et al. [29] recover both instances and semantics and apply a CRF to improve the predictions accuracy.

Most of these works utilize a PointNet [30] or PointNet++ [32] network to predict the point embeddings. In our case, we extend LatticeNet in a similar manner to other proposal-free methods but predict the embeddings using the lattice convolutions.

## 3 Notation

Throughout this paper, we use bold upper-case characters to denote matrices and bold lower-case characters to denote vectors.

The vertices of the  $d$ -dimensional permutohedral lattice are defined as a tuple  $v = (\mathbf{c}_v, \mathbf{x}_v)$ , with  $\mathbf{c}_v \in \mathbb{Z}^{(d+1)}$  denot-

ing the coordinates of the vertex and  $\mathbf{x}_v \in \mathbb{R}^{v_d}$  representing the values stored at vertex  $v$ . The full lattice containing  $n$  vertices is denoted with  $V = (\mathbf{C}, \mathbf{X})$ , with  $\mathbf{C} \in \mathbb{Z}^{n \times (d+1)}$  representing the coordinate matrix and  $\mathbf{X} \in \mathbb{R}^{n \times v_d}$  the value matrix.

The points in a cloud are defined as a tuple  $p = (\mathbf{g}_p, \mathbf{f}_p)$ , with  $\mathbf{g}_p \in \mathbb{R}^d$  denoting the coordinates of the point and  $\mathbf{f}_p \in \mathbb{R}^{f_d}$  representing the features stored at point  $p$  (color, normals, etc.). The full point cloud containing  $m$  points is denoted by  $P = (\mathbf{G}, \mathbf{F})$  with  $\mathbf{G} \in \mathbb{R}^{m \times d}$  being the positions matrix and  $\mathbf{F} \in \mathbb{R}^{m \times f_d}$  the feature matrix. The feature matrix  $\mathbf{F}$  can also be empty in which case  $f_d$  is set to zero.

For motion segmentation we define a sequence of point clouds as  $P_{seq} = (P_0, P_1, \dots, P_n)$  with  $P_n = (\mathbf{G}, \mathbf{F})$ . We define a timestep as processing one cloud of this sequence.

We denote with  $I_p$  the set of lattice vertices of the simplex that contains point  $p$ . The set  $I_p$  always contains  $d+1$  vertices as the lattice tessellates the space in uniform simplices with  $d+1$  vertices each. Furthermore, we denote with  $J_v$  the set of points  $p$  for which vertex  $v$  is one of the vertices of the containing simplices. Hence, these are the points that contribute to vertex  $v$  through the splat operation.

We denote with  $\mathcal{S}$  the splatting operation, with  $\mathcal{Y}$  the slicing operation, with  $\hat{\mathcal{Y}}$  the deformable slicing, with  $\mathcal{P}$  the PointNet module, with  $\mathcal{D}_G$  and  $\mathcal{D}_F$  the distribution of the point positions and the points features, respectively, and with  $\mathcal{G}$  the gathering operation.

## 4 Permutohedral Lattice

The  $d$ -dimensional permutohedral lattice is formed by projecting the scaled regular grid  $(d+1)\mathbb{Z}^{d+1}$  along the vector  $\mathbf{1} = [1, \dots, 1]$  onto the hyperplane  $H_d: \mathbf{p} \cdot \mathbf{1} = 0$ .

The lattice tessellates the space into uniform  $d$ -dimensional simplices. Hence, for  $d = 2$  the space is tessellated with triangles and for  $d = 3$  into tetrahedra. The enclosing simplex of any point can be found by a simple rounding algorithm [2].

Due to the scaling and projection of the regular grid, the coordinates  $\mathbf{c}_v$  of each lattice vertex sum up to zero. Each vertex has  $2(d+1)$  immediate neighboring vertices. The coordinates of these neighbors are separated by a vector of form  $\pm[-1, \dots, -1, d, -1, \dots, -1] \in \mathbb{Z}^{d+1}$ .

The vertices of the permutohedral lattice are stored in a sparse manner using a hash map in which the key is the coordinate  $\mathbf{c}_v$  and the value is  $\mathbf{x}_v$ . Hence, we only allocate the simplices that contain the 3D surface of interest. This sparse allocation allows for efficient implementation of all typical operations in CNNs (convolution, pooling, transposed convolution, etc.).

The permutohedral lattice has several advantages w.r.t. standard cubic voxels. The number of vertices for each

simplex is given by  $d+1$  which scales linearly with increasing dimension, in contrast to the  $2^d$  for standard voxels. This small number of vertices per simplex allows for fast splatting and slicing operations. Furthermore, splatting and slicing create piece-wise linear outputs as they use barycentric interpolation. In contrast, standard quantization in cubic voxels create piece-wise constant outputs, leading to discretization artefacts.

Spatial correspondences between lattice vertices are given by design and the hashmap: If the hashmap stays the same for the whole sequence, spatially identical lattice vertices of different point clouds are always mapped to the same entries. This is visualized in Fig. 9 where features from two different time-steps are fused together.

## 5 Method

The input to our method is a point cloud  $P = (\mathbf{G}, \mathbf{F})$  containing coordinates and per-point features.

We define the scale of the lattice by scaling the positions  $\mathbf{G}$  as  $\mathbf{G}_s = \mathbf{G}/\sigma$ , where  $\sigma \in \mathbb{R}^d$  is the scaling factor. The higher the sigma the less number of vertices will be needed to cover the point cloud and the coarser the lattice will be. For ease of notation, unless otherwise specified, we refer to  $\mathbf{G}_s$  as  $\mathbf{G}$  as we usually only need the scaled version.

### 5.1 Common Operations on Permutohedral Lattice

In this section, we will explain in detail the standard operations on a permutohedral lattice that are used in previous works [39, 15].

**Splatting** refers to the interpolation of point features onto the values of the lattice  $V$  using barycentric weighting (Fig. 3a). Each point splats onto  $d+1$  lattice vertices and their weighted features are summed onto the vertices.

**Convolve** operates analogously to standard spatial convolutions in 2D or 3D, i.e. a weighted sum of the vertex values together with its neighbors is computed. We use convolutions that span over the 1-hop ring around a vertex and hence convolve the values of  $2(d+1) + 1$  vertices (Fig. 2).

**Slicing** is the inverse operation to splatting. The vertex values of the lattice are interpolated back for each position with the same weights used during splatting. The weighted contributions from the simplexes  $d+1$  vertices are summed up (Fig. 5a).

### 5.2 Proposed Operations on Permutohedral Lattice

The operations defined in Section Sec. 5.1 are typically used in a cascade of splat-conv-slice to obtain dense predictions [39]. However, splatting and slicing act as Gaussian



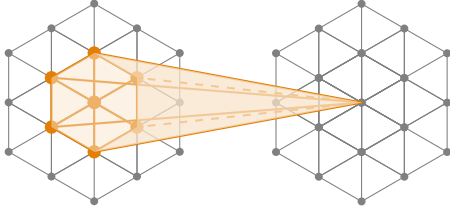


Fig. 2: Convolution: The neighboring vertices of a lattice are convolved similarly to standard 2D convolutions. If a neighbor is not allocated in the sparse structure, we assume that it has a value of zero.

kernel low-pass filtering on encoded information [2]. Their repeated usage at every layer is detrimental to the accuracy of the network. Additionally, splatting acts as a weighted average on the feature vectors where the weights are only determined through barycentric interpolation. Including the weights as trainable parameter allows the network to decide on a better interpolation scheme. Furthermore, as the network grows deeper and feature vectors become higher-dimensional, slicing consumes increasingly more memory, as it assigns the features to the points. Since in most cases  $|P| \gg |V|$ , it is more efficient to store the features only in the lattice vertices.

To address these limitations, we propose four new operators on the permutohedral lattice which are more suitable for CNNs and dense prediction tasks.

**Distribute** is defined as the list of features that each lattice vertex receives. However, they are not summed as done by splatting:

$$\mathbf{x}_v = \mathcal{S}(P, V) = \sum_{p \in J_v} b_{pv} \mathbf{f}_p, \quad (1)$$

where  $\mathbf{x}_v$  is the value of lattice vertex  $v$  and  $b_{pv}$  is the barycentric weight between point  $p$  and lattice vertex  $v$ .

Instead, our distribute operators  $\mathcal{D}_G$  and  $\mathcal{D}_F$  concatenate coordinates and features of the contributing points:

$$\mathbf{x}_v = \mathcal{P}(\mathbf{D}_{v_g}; \mathbf{D}_{v_f}), \quad (2)$$

$$\mathbf{D}_{v_g} = \mathcal{D}_G(P, V) = \{ \mathbf{g}_p - \boldsymbol{\mu}_v \mid p \in J_v \}, \quad (3)$$

$$\mathbf{D}_{v_f} = \mathcal{D}_F(P, V) = \{ \mathbf{f}_p \mid p \in J_v \}, \quad (4)$$

$$\boldsymbol{\mu}_v = \frac{1}{|J_v|} \sum_{p \in J_v} \mathbf{g}_p, \quad (5)$$

where  $\mathbf{D}_{v_g} \in \mathbb{R}^{|J_v| \times d}$  and  $\mathbf{D}_{v_f} \in \mathbb{R}^{|J_v| \times f_d}$  are matrices containing the distributed coordinates and features, respectively, for the contributing points into a vertex  $v$ . The matrices are concatenated and processed by a PointNet  $\mathcal{P}$  to obtain the final vertex value  $\mathbf{x}_v$ . Fig. 3 illustrates the difference between splatting and distributing.

Note that we use a different distribute function for coordinates then for point features. For coordinates, we subtract the mean of the contributing coordinates. The intuition behind

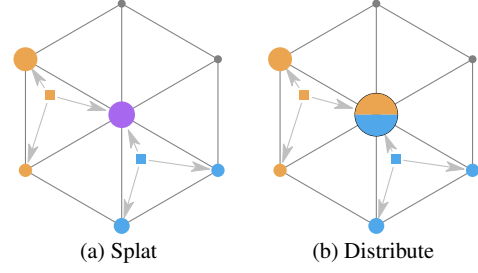


Fig. 3: Splat and Distribute operations: Splatting uses barycentric weighting to add the features of points onto neighboring vertices. The naïve summation can be detrimental to the network as splatting acts as a Gaussian filter. Distributing stores all the features of the contributing points, causing no loss of information and allows further processing by the network.

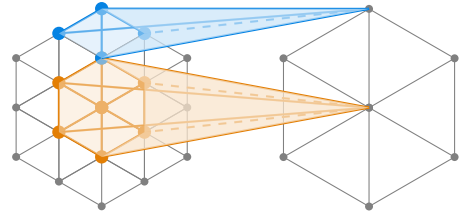


Fig. 4: Coarsen: Downsampling of the lattice is performed by embedding the coarse lattice in the finer one and convolving over the neighbors. This effectively performs a strided convolution. Transposed convolution is performed in an analogous manner by embedding a fine lattice into a coarse one.

this is that coordinates by themselves are not very informative w.r.t. the potential semantic class. However, the local distribution is more informative as it gives a notion of the geometry.

**Downsampling** refers to a coarsening of the lattice, by reducing the number of vertices. This allows the network to capture more contextual information. Downsampling consists of two steps: creation of a coarse lattice and obtaining its values. Coarse lattices are created by repeatedly dividing the point cloud positions by 2 and using them to create new lattice vertices [3]. The values of the coarse lattice are obtained by convolving over the finer lattice from the previous level (Fig. 4). Hence, we must embed the coarse lattice inside the finer one by scaling the coarse vertices by 2. Afterwards, the neighbors vertices over which we convolve are separated by a vector of form  $\pm [-1, \dots, -1, d, -1, \dots, -1] \in \mathbb{Z}^{d+1}$ . The downsampling operation effectively performs a strided convolution.

**Upsampling** follows a similar reasoning. The fine vertices need first to be embedded in the coarse lattice using a division by 2. Afterwards, the neighboring vertices over which we convolve are separated by a vector of form  $\pm [-0.5, \dots, -0.5, d/2, -0.5, \dots, -0.5]$ . The careful

reader will notice that in this case, the coordinates of the neighboring vertices may not be integer anymore; they may have a fractional part and will, therefore, lie in the middle of a coarser simplex. In this case we ignore the contribution of this neighboring vertices and only take the contribution of the center vertex.

The upsampling operation effectively performs a transposed convolution.

**DeformSlicing:** While the slicing operation  $\mathcal{Y}$  barycentrically interpolates the values back to the points by using barycentric coordinates:

$$f_p = \mathcal{Y}(P, V) = \sum_{v \in I_p} b_{pv} \mathbf{x}_v, \quad (6)$$

we propose the DeformSlicing  $\tilde{\mathcal{Y}}$  which allows the network to directly modify the barycentric coordinates and shift the position within the simplex for data-dependent interpolation:

$$f_p = \tilde{\mathcal{Y}}(P, V) = \sum_{v \in I_p} (b_{pv} + \Delta b_{pv}) \mathbf{x}_v. \quad (7)$$

Here,  $\Delta b_{pv}$  are offsets that are applied to the original barycentric coordinates. A parallel branch within our network first gathers the values from all the vertices in a simplex and regresses the  $\Delta b_{pv}$ :

$$\mathbf{q}_p = \mathcal{G}(P, V) = \{ b_{pv} \mathbf{x}_v \mid v \in I_p \}, \quad (8)$$

$$\Delta \mathbf{b}_p = \mathcal{F}(\mathbf{q}_p), \quad (9)$$

where  $\mathbf{q}_p$  is a set containing the weighted values of all the vertices of the simplex containing  $p$  and the prediction  $\Delta \mathbf{b}_p = \{ \Delta b_{pv} \mid v \in I_p \}$  is a set of offsets to the barycentric coordinates towards the  $d + 1$  vertices.

With a slight abuse of notation — due to the fact that the vertices of a simplex are always enumerated in a consistent manner, we can regard  $\mathbf{b}_p$  and  $\mathbf{q}_p$  as vectors in  $\mathbb{R}^{(d+1)}$  and  $\mathbb{R}^{(d+1)v_d}$ , respectively, and cast the prediction of offsets as a fully connected layer followed by a non-linearity:

$$\Delta \mathbf{b}_p = \mathcal{F}(\mathbf{q}_p) = \sigma(\mathbf{q}_p \cdot \mathbf{W} + b). \quad (10)$$

However, this prediction has the disadvantage of not being permutation equivariant; therefore, permutation of the vertices would not imply the same permutation in the barycentric offsets:

$$\mathcal{F}(\pi \mathbf{q}_p) \neq \pi \mathcal{F}(\mathbf{q}_p), \quad (11)$$

where  $\pi$  is the set of all permutations of the  $d + 1$  vertices.

It is important for our prediction to be permutation equivariant because the vertices may be arranged in any order and the barycentric offsets need to keep a consistent preference towards a certain vertexes' features, regardless of its position within a simplex.

In order for the prediction of the offsets to be consistent with permutations of the vertices, we take inspiration from

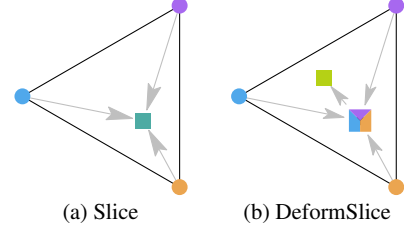


Fig. 5: Slice and DeformSlice: Slicing barycentrically interpolates the vertex values back onto a point. DeformSlice allows for the network to directly affect the interpolated value by learning offsets of the barycentric coordinates.

the work of [33] and [53] of equivariant layers and design  $\mathcal{F}$  as:

$$\Delta b_{pv} = \sigma(b + (b_{pv} \mathbf{x}_v - \max_{d \in I_p} \{b_{pd} \mathbf{x}_d\}) \cdot \mathbf{W}), \quad (12)$$

$$\Delta \mathbf{b}_p = \mathcal{F}(\mathbf{q}_p) = \{ \Delta b_{pv} \mid v \in I_p \}, \quad (13)$$

where  $\mathbf{W} \in \mathbb{R}^{v_d \times 1}$  is a weight matrix and  $b \in \mathbb{R}$  corresponds to a scalar bias.

In other words, we subtract from each weighted vertex the maximum of the weighted values of all the other vertices in the simplex. Since the max operation is invariant to permutations of the input, the regression of the offsets is *equivariant* to permutations of the vertices.

The difference between the slicing and our DeformSlicing is visualized in Fig. 5

## 6 Segmentation Methods

Due to the flexibility of LatticeNet various segmentation methods can be implemented. In this section, we detail the methods used for each one.

### 6.1 Semantic Segmentation

Semantic segmentation uses the default U-Net architecture described in the Network Architecture section. It is trained with an equal part combination of cross entropy loss and Lovász loss [5]. The Lovász loss acts as a surrogate for the intersection-over-union score and is especially useful for dealing with class imbalance.

### 6.2 Instance Segmentation

Our instance segmentation network follows the work of other proposal-free methods like [12]. We use LatticeNet to predict for each 3D point  $p_i$  in the point cloud an embedding  $x_i$ . A discriminative loss encourages closeness in embeddings space for points of the same instance while promoting

distance between different instances. Finally, we apply mean-shift clustering on the points in embeddings space. Points belonging to the same cluster are defined as an Instances.

This discriminative loss can be expressed with three terms:

- Variance term: The intra-cluster pull force that draws the embeddings towards the mean embedding.
- Distance term: An inter-cluster push force that forces the clusters to be far apart from each other in embedding space.
- Regularization term: A small force that pulls the cluster centers towards the origin in order to keep the activations bounded.

The full loss is then defined as:

$$L_{var} = \frac{1}{C} \sum_{c=1}^C \frac{1}{N_c} \sum_{i=1}^{N_c} [\|\mu_c - x_i\| - \delta_v]^2 \quad (14)$$

$$L_{dist} = \frac{1}{C(C-1)} \sum_{\substack{c_A=1 \\ c_A \neq c_B}}^C \sum_{c_B=1}^C [2\delta_d - \|\mu_{c_A} - \mu_{c_B}\|]_+^2 \quad (15)$$

$$L_{reg} = \frac{1}{C} \sum_{c=1}^C \|\mu_c\| \quad (16)$$

$$L = \alpha \cdot L_{var} + \beta \cdot L_{dist} + \gamma \cdot L_{reg} \quad (17)$$

We define  $C$  as the number of clusters in the ground truth,  $N_c$  as the number of elements in cluster  $c$ ,  $x_i$  as the embedding vector for point  $p_i$  and  $\mu_c$  as the mean or cluster center for cluster  $c$ . The  $\delta_v$  and  $\delta_d$  are the margins for the variance and distance loss respectively. We set  $\alpha = \beta = 1$  and  $\gamma = 0.001$

A visualization of the pipeline for instance segmentation can be seen in Fig. 6.

### 6.3 Motion Segmentation

Motion segmentation distinguishes between dynamic and static objects within a point cloud. For this, the network needs temporal information. We extend the original LatticeNet U-Net architecture with a recursive architecture that can process a sequence of point clouds  $P_{seq}$  at times  $t, t-1, \dots, t-n$  and learn to distinguish for example between a moving car and a parked car.

The dynamic objects are considered as additional classes. Hence, we use the same loss as in the case of semantic segmentation. We also explore multiple ways to perform the fusion of temporal information which we detail in the Network Architecture section.

## 7 Network Architecture

Input to our network is a point cloud  $P$  which may contain per-point features stored in  $\mathbf{F}$ . The output is class probabilities for each point  $p$ . In the recurrent network the input is an ordered set of point clouds  $P_{seq}$  and the output are class probabilities for the last point cloud of the sequence. Moving and static objects are considered as different semantic classes.

Our network architecture has a U-Net structure [35] and is visualized in Fig. 7 together with the used individual blocks.

The first layers distribute the point features onto the lattice and use a PointNet to obtain local features. Afterwards, a series of ResNet blocks [16], followed by repeated downsampling, aggregates global context. The decoder branch mirrors the encoder architecture and upsamples through transposed convolutions. Finally, a DeformSlicing propagates lattice features onto the original point cloud. Skip connections are added by concatenating the encoder feature maps with matching decoder features.

### 7.1 Temporal Fusion

Incorporating temporal information for motion prediction over a sequence of point clouds relies on fusing information between multiple time-steps. For this purpose, the feature vectors of the timesteps  $t-1$  and  $t$  are passed through a Temporal Fusion block, as shown in Fig. 8. This fusion consists of a concatenation of both feature vectors and a linear layer followed by a non-linearity (Fig. 9). Each new time-step allocates additional vertices in the lattice corresponding to newly explored areas in the map. For correct fusion, the features from the previous time-step need to be zero-padded so that the sizes match.

Additionally, we performed experiments with a single Temporal Fusion block in the network and max-pooling over both feature vectors instead of the linear layer, but found that three Temporal Fusion blocks achieved overall superior results.

It should be noted that our approach for temporal fusion relies on a sequence of clouds that are transformed into a common coordinate frame. The required scan poses for transformation can be obtained e.g. from GPS or SLAM.

## 8 Implementation

Our lattice is stored sparsely on a hash map structure, which allows for fast access of neighboring vertices. Unlike [39], we construct the hash map directly on the GPU, saving us from incurring an expensive CPU to GPU memory copy.

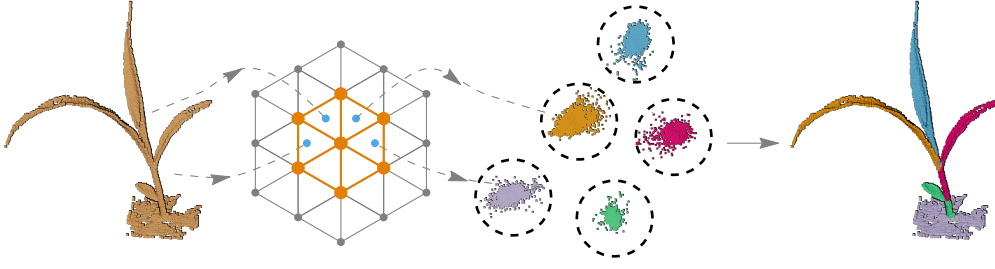


Fig. 6: Instance segmentation: LatticeNet takes raw point clouds as input and embeds them into a sparse lattice where convolutions are applied. Features on the lattice are projected onto a 2D space where clustering is performed. The clusters define the instances of each object type in the original cloud.

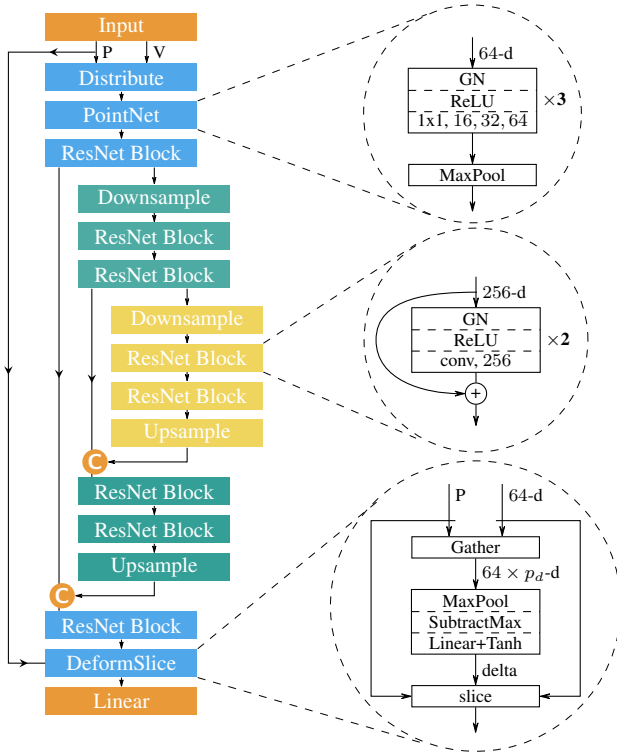


Fig. 7: Architecture: Our model follows a U-Net structure. For ease of representation, blocks which are repeated one after another are indicated with a multiplier on the right side of the operation.

For memory savings, we implemented the DeformSlice and the last linear classification layer in one fused operation, avoiding the storage of high-dimensional feature vectors for each point in the point cloud.

All of the lattice operators containing forwards and backwards passes are implemented on the GPU and exposed to PyTorch [27].

Following recent works [17, 18], all convolutions are pre-activated using Group Normalization [49] and a ReLU unit. We chose Group Normalization instead of the standard batch

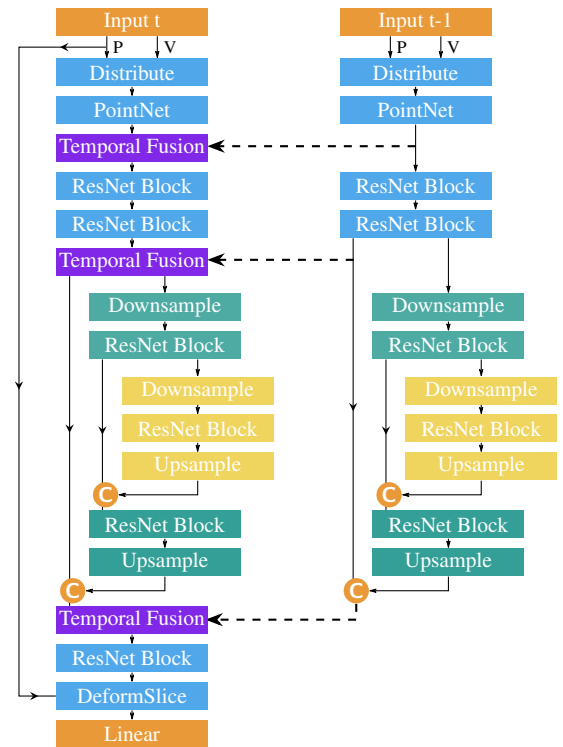


Fig. 8: Recurrent architecture: The features from previous time-steps are fused in the current time-step at multiple levels of the network. This allows the network to distinguish dynamic objects from static ones.

normalization due to greater stability for small batch sizes. We use the default of 32 groups.

The models were trained using the Adam optimizer with a learning rate of 0.001 and a weight decay of  $10^{-4}$ . The learning rate was reduced by a factor of 10 when the loss plateaued.

We share the PyTorch implementation of LatticeNet at [https://github.com/AIS-Bonn/lattice\\_net](https://github.com/AIS-Bonn/lattice_net).

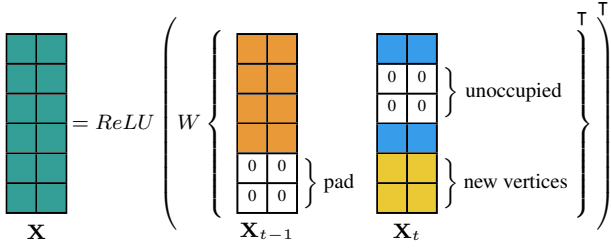


Fig. 9: Temporal fusion: The features from the previous time-step are zero-padded in order to account for the new vertices that were allocated at the current time-step. The features are afterwards concatenated and passed through a linear layer followed by a non-linearity.

## 9 Experiments

We evaluate our proposed lattice network on four different datasets: ShapeNet [51], ScanNet [11], SemanticKITTI [4] and Pheno4D [1]. For the task of semantic segmentation and motion segmentation we report the mean Intersection-over-Union (mIoU). For the task of instance segmentation, we report the Symmetric Best Dice (SBD) [12]. SBD measures the accuracy of the instance segmentation by averaging for each input label the ground truth label yielding the maximum Dice score.

We use a shallow model for ShapeNet and Pheno4D and a deeper model for ScanNet and SemanticKITTI as the datasets are larger. We augment all data using random mirroring and translations in space. For ScanNet, we also apply random color jitter. A video with additional footage of the experiments is available online <sup>1</sup>.

### 9.1 Evaluation of Segmentation Accuracy

**ShapeNet part segmentation** is a subset of the ShapeNet dataset [51] which contains objects from 16 different categories each segmented into 2 - 6 parts. The dataset consists of points sampled from the surface of the objects, together with the ground truth label of the corresponding object part. The objects have an average of 2613 points. We train and evaluate our network on each object individually. We use the official train/test splits as defined by the dataset containing a total of 12 137 training objects and 2874 test objects. The results for our and five competing methods are gathered in Tab. 1 and visualized in Fig. 11.

We observe that for some classes, we obtain state-of-the-art performance and for other objects, the IoU is slightly lower than for other approaches. We ascribe this to the fact that training one fixed architecture size for each individual

object is suboptimal as some objects like the "cap" have as few as 55 examples while others like the table have more than 5K. This causes the network to be prone to overfitting on the easy object or underfitting on the difficult ones. A fair evaluation would require finding an architecture that performs well for all objects on average. However, due to various issues with mislabeled ground truths [39] we deem that experimentation with more architectures or with different regularization strengths for individual objects would overfit the dataset.

**ScanNet 3D segmentation** [11] consists of 3D reconstructions of real rooms. It contains  $\approx 1500$  rooms segmented into 20 classes (bed, furniture, wall, etc.). The rooms have between 9K and 537K points — on average 145K. We segment an entire room at once without cropping. We use the official train/test splits as defined by the dataset containing a total of 1201 training rooms and 100 test objects. We obtain an IoU of 64.0 which is significantly higher than the most similar related work of SplatNet. It is to be noted that MinkowskiNet achieves a higher IoU but at the expense of an extremely high spatial resolution of 2 cm per voxel. In contrast, our approach allocates lattice vertices so that each vertex covers approximately 30 points. On this dataset, this corresponds to a spatial extent of approximately 10 cm.

**SemanticKITTI** [4] consists of semantically annotated LiDAR scans of real urban environments. The annotation covers a total of 19 classes for single scan evaluation and a total of 25 classes for multiple scan evaluation. Each scan contains between 82K and 129K points. We process each scan entirely without any cropping. We use the official train/validation splits as defined by the dataset. The test set is not publicly available and testing can only be done through the benchmark server.

The results for single scan are provided in Tab. 2. Our LatticeNet outperforms all other methods — in case of the most similar SplatNet by more than a factor of two. It is to be noted that DarkNet53Seg [4], DarkNet21Seg [4] and SqueezeSegV2 [47] are methods that operate on a 2D image by wrapping the LiDAR scans to 2D using spherical coordinates. In contrast, our method can operate on general point clouds, directly in 3D.

For motion segmentation we take as input three point clouds at consecutive time steps and output the segmentation for the final, most recent cloud. We overlap this time window so that every clouds gets to be segmented. For the first few clouds, the time window is reduced as there are no clouds from previous time-steps to give as input. The results for the motion segmentation are provided in Tab. 4.

We observe that for motion segmentation we outperform other approaches except for KPConv[43], which has higher IoU. However, it is to be noted that KPConv cannot process a full point cloud at once due to memory constraints and rather processes sub-clouds centered around random spheres

<sup>1</sup> [http://www.ais.uni-bonn.de/videos/RSS\\_2020\\_Rosu/](http://www.ais.uni-bonn.de/videos/RSS_2020_Rosu/)



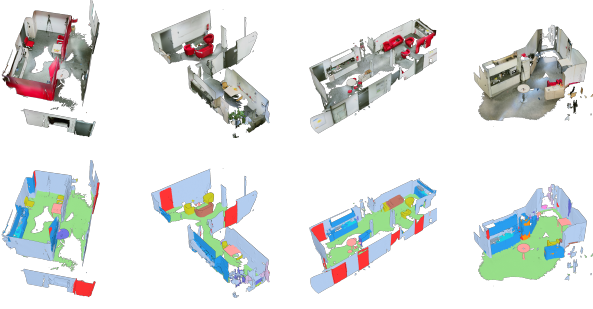


Fig. 10: Bonn Activity Maps segmentations. Colored meshes are reconstructed from KinectV2 data using volumetric integration [26, 38] and semantically segmented using LatticeNet. Color coding of semantic labels corresponds to the ScanNet dataset [11].

in the scene. The spheres are chosen randomly in the scene to ensure each point is tested multiple times by different sphere locations. Finally, a voting scheme gives the final prediction. In contrast, our approach can process a full point cloud without requiring neighborhood searching or partitioning in sub-clouds.

**Bonn Activity Maps** [40] is a dataset for human tracking, activity recognition and anticipation of multiple persons. It contains annotations of persons, their trajectories and activities. The 3D reconstruction of the four kitchen scenarios is however of more interest to us. The environments are reconstructed as 3D colored meshes and have no ground truth semantic annotations. We trained our LatticeNet on the ScanNet dataset and evaluate it on the 4 kitchens in order to provide an annotation for each vertex of the mesh. The results are shown in Fig. 10. We can observe that our network generalizes well to unseen datasets, recorded with different sensors and with different noise properties as the semantic segmentations look plausible and exhibit sharp borders between classes.

**Pheno4D** [1] is a spatio-temporal dataset of point clouds of maize and tomato plants with instance annotations of leaves. We use a shallow version of LatticeNet to compute per-point embeddings and cluster them using mean-shift to recover the instances. We compare with PointNet and PointNet++ as they are popular methods for computing per-point embeddings. Since the dataset contains 7 maize and 7 tomato plants, we train on the first 5 plants for each type and test on the remaining two. The results are gathered in Tab. 5. We observe that our method is capable of computing more meaningful embeddings that create more distinctive clusters between each plant organ.

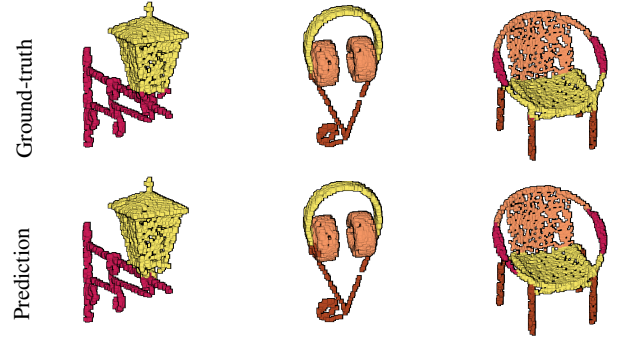


Fig. 11: ShapeNet [51] results of our method.

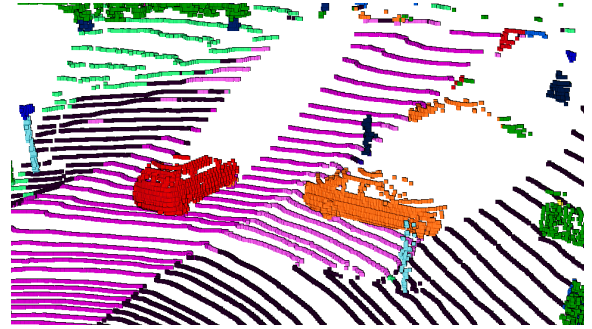


Fig. 12: Motion segmentation results on SemanticKITTI. The moving car on the road (red) is correctly distinguished from the parked car (orange).

## 9.2 Ablation Studies

We perform various ablations regarding our contribution to judge how much they affect the network’s performance.

**DeformSlice:** We assess the impact that DeformSlice has on the network by comparing it with the Slice operator which does not use learned barycentric interpolation. We evaluate this on SemanticKITTI, the largest dataset that we are using.

We also evaluate a version of DeformSlice which ensures that the new barycentric coordinates still sum up to one by adding an additional loss term:

$$L = \frac{1}{|P|} \sum_{p \in P} \left( \sum_{v \in I_p} \Delta b_{pv} \right)^2. \quad (18)$$

However, we observe little change after adding this regularization term and hence, use the default version of DeformSlice for the rest of the experiments. The results are gathered in Tab. 7.

**Distribute and PointNet:** Another contribution of our work is the usage of a Distribute operator to provide values to the lattice vertices which are later embedded in a higher-dimensional space by a PointNet-like architecture. The positions and features of the point cloud are treated separately where the features (normals, color) are distributed directly.

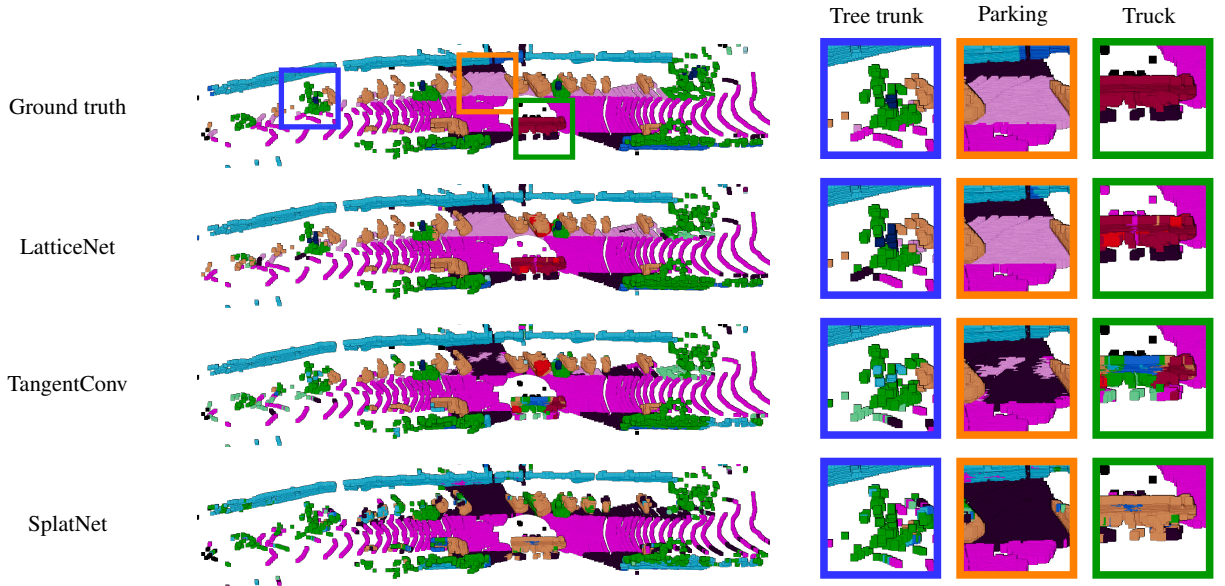


Fig. 13: SemanticKITTI results. We compare the prediction from our LatticeNet with the results from TangentConv [41] and SplatNet [39]. We can observe that our approach can better learn small objects like tree trunks, despite their relatively small number of points. Additionally, the network also effectively makes use of contextual information in order to correctly predict the parking place due to the existence of nearby cars.

Table 1: Results on ShapeNet part segmentation [51].

#instances		2690	76	55	898	3758	69	787	392	1547	451	202	184	283	66	152	5271
	instance avg.	air- plane	bag	cap	car	chair	ear- phone	guitar	knife	lamp	laptop	motor- bike	mug	pistol	rocket	skate- board	table
PointNet [30]	83.7	83.4	78.7	82.5	74.9	89.6	73.0	91.5	85.9	80.8	95.3	65.2	93.0	81.2	57.9	72.8	80.6
PointNet++ [32]	85.1	82.4	79.0	87.7	77.3	<b>90.8</b>	71.8	91.0	85.9	83.7	95.3	71.6	94.1	81.3	58.7	76.4	82.6
SplatNet 3D [39]	84.6	81.9	83.9	88.6	79.5	90.1	73.5	91.3	84.7	84.5	96.3	69.7	95.0	81.7	59.2	70.4	81.3
SplatNet 2D-3D [39]	<b>85.4</b>	83.2	84.3	<b>89.1</b>	80.3	90.7	<b>75.5</b>	92.1	87.1	83.9	96.3	75.6	95.8	83.8	64.0	75.5	81.8
FCPN [34]	84.0	<b>84.0</b>	82.8	86.4	<b>88.3</b>	83.3	73.6	<b>93.4</b>	87.4	77.4	<b>97.7</b>	<b>81.4</b>	95.8	<b>87.7</b>	68.4	<b>83.6</b>	73.4
Ours	83.9	82.3	<b>84.8</b>	79.1	81.0	86.9	71.0	91.9	<b>89.4</b>	<b>84.7</b>	96.6	77.2	<b>95.8</b>	86.0	<b>70.5</b>	79.3	<b>87.0</b>

Table 2: Results on SemanticKITTI [4].

Approach	mIoU	road	sidewalk	parking	other-ground	building	car	truck	bicycle	motorcycle	other-vehicle	vegetation	trunk	terrain	person	bicyclist	motorcyclist	fence	pole	traffic sign
PointNet [30]	14.6	61.6	35.7	15.8	1.4	41.4	46.3	0.1	1.3	0.3	0.8	31.0	4.6	17.6	0.2	0.2	0.0	12.9	2.4	3.7
SplatNet [39]	18.4	64.6	39.1	0.4	0.0	58.3	58.2	0.0	0.0	0.0	0.0	71.1	9.9	19.3	0.0	0.0	0.0	23.1	5.6	0.0
PointNet++ [32]	20.1	72.0	41.8	18.7	5.6	62.3	53.7	0.9	1.9	0.2	0.2	46.5	13.8	30.0	0.9	1.0	0.0	16.9	6.0	8.9
Minkowski34(25cm) [9]	33.0	80.8	43.0	36.9	0.5	73.5	83.0	<b>42.9</b>	2.0	2.9	7.8	74.4	42.9	36.7	11.2	22.8	4.4	37.2	35.4	28.6
SqueezeSegV2 [47]	39.7	88.6	67.6	45.8	17.7	73.7	81.8	13.4	18.5	17.9	14.0	71.8	35.8	60.2	20.1	25.1	3.9	41.1	20.2	36.3
TangentConv [41]	40.9	83.9	63.9	33.4	15.4	83.4	90.8	15.2	2.7	16.5	12.1	79.5	49.3	58.1	23.0	28.4	8.1	49.0	35.8	28.5
DarkNet21Seg [4]	47.4	91.4	74.0	57.0	26.4	81.9	85.4	18.6	<b>26.2</b>	26.5	15.6	77.6	48.4	63.6	31.8	33.6	4.0	52.3	36.0	50.0
DarkNet53Seg [4]	49.9	<b>91.8</b>	<b>74.6</b>	<b>64.8</b>	<b>27.9</b>	84.1	86.4	25.5	24.5	<b>32.7</b>	<b>22.6</b>	78.3	50.1	<b>64.0</b>	<b>36.2</b>	33.6	4.7	55.0	38.9	<b>52.2</b>
Ours	<b>52.9</b>	90.0	74.1	59.4	22.0	<b>88.2</b>	<b>92.9</b>	26.6	16.6	22.2	21.4	<b>81.7</b>	<b>63.6</b>	63.1	35.6	<b>43.0</b>	<b>46.0</b>	<b>58.8</b>	<b>51.9</b>	48.4

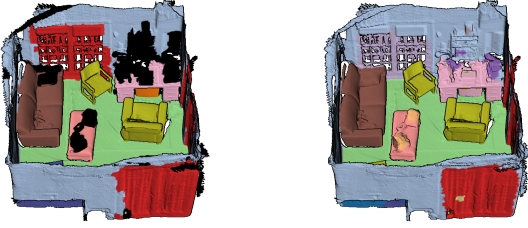


Fig. 14: ScanNet results. The left image shows the ground truth and the right one our prediction.

Table 3: Results on ScanNet [11]

Method	mIoU
PointNet++ [32]	33.9
SplatNet [39]	39.3
TangentConv [41]	43.8
3DMV <sup>‡</sup> [10]	48.4
MinkowskiNet42 (5cm) [9]	67.9
SparseConvNet [14] <sup>†</sup>	72.5
MinkowskiNet42 (2cm) [9] <sup>†</sup>	<b>73.4</b>
Ours	64.0

<sup>†</sup>: post-CVPR submissions. <sup>‡</sup>: uses 2D images additionally.

Table 4: Motion segmentation IoU results on SemanticKITTI [4] using a sequence of multiple past scans (in %). Shaded cells correspond to the IoU of the moving classes, while unshaded entries are the non-moving classes.

Approach	car	truck	other-vehicle	person	bicyclist	motorcyclist	mIoU
TangentConv [41]	84.9	21.1	18.5	1.6	0.0	0.0	34.1
	40.3	42.2	30.1	6.4	1.1	1.9	
DarkNet53Seg [4]	84.1	20.0	20.7	7.5	0.0	0.0	41.6
	61.5	37.8	28.9	15.2	14.1	0.2	
SpSequenceNet [37]	88.5	29.2	22.7	6.3	0.0	0.0	43.1
	53.2	0.1	2.3	26.2	41.2	36.2	
KPCConv [43]	93.7	70.3	38.6	21.6	0.0	0.0	51.2
	69.4	5.8	4.7	67.5	67.4	47.2	
Ours	91.1	65.4	23.1	6.8	0.0	0.0	45.2
	54.8	3.5	0.6	49.9	44.6	64.3	

From the positions, we subtract the locally averaged position as we assume that the local point distribution is more important than the coordinates in the global reference frame. We evaluate the impact of elevating the point features to a higher-dimensional space and subtracting the local mean against a simple splatting operator which just averages the features of the points around each corresponding vertex.

Table 5: Instance segmentation performance on the maize and tomato plants of the Pheno4D dataset.

	SBD	
	Maize	Tomato
PointNet[30]	69.7	47.3
PointNet++[32]	74.8	56.1
LatticeNet (ours)	<b>80.6</b>	<b>74.2</b>

We observe that not subtracting the local mean, and just using the  $xyz$  coordinates as features, heavily degrades the performance, causing the mIoU to drop from 52.9 to 43.0. This further reinforces the idea that the local point distribution is a good local feature to use in the first layers of the network.

Not elevating the point cloud features to a higher-dimensional space before applying the max-pool operation also hurts performance but not as severely. In our experiments, we elevate the features to 64 dimensions by using a series of fully connected layers.

Finally, naive application of the splat operation performs worst with a mere 37.8 mIoU.

### 9.3 Performance

We report the time taken for a forward pass and the maximum memory used in our shallow and deep network on the first three evaluated datasets. The performance was measured on a NVIDIA Titan X Pascal and the results are gathered in Tab. 6.

In the case of motion segmentation, the inference times and memory used are the same as in the case of a single scan, as we use the same backbone network to extract features and the computational cost of fusing the temporal information is minimum. However for training, the network requires more memory with increasing time window due to the back-propagation through time. This scales linearly with the time window size and the amount of points in the cloud.

Despite the reduced memory usage compared to SplatNet and increased speed of execution, there are still memory savings possible by fusing the Distribute and PointNet operators into one GPU operation. This is similar to fusing our Deform-Slice and the classification layer. Additionally, we expect the network to become even faster as further advances on highly optimized kernels for convolution on sparse lattices become available. At the moment, the convolutions are performed by our custom CUDA kernels. Tighter integration however with highly optimized libraries like cuDNN [8] could be beneficial.



Table 6: Average time used by the forward pass and the maximum memory used during training. An X indicates a method that failed to process the whole cloud due to memory limitations.

	ShapeNet [ms] [GB]		ScanNet [ms] [GB]		SemanticKITTI [ms] [GB]	
SplatNet	129	0.6	X	X	2931	8.9
Ours	49	0.5	180	6.5	143	3.5

Table 7: Ablation study of the various components of LatticeNet. Various features are disabled (indicated in red) and the impact to the IoU is evaluated.

	mIoU	LocalAvg	PointNet elevate	DeformSlice	Offsets zero sum
LN splat	37.8				
LN no local avg	43.0				
LN no elevate	46.8				
LN slice	50.4				
LN reg	52.7				
LatticeNet	<b>52.9</b>				

## 10 Conclusion

We presented LatticeNet, a novel method for point cloud segmentation. A sparse permutohedral lattice allows us to efficiently process large point clouds. The usage of PointNet together with a data-dependent interpolation alleviates the quantization issues of other methods. Experiments on four datasets show state-of-the-art results, at a reduced time and memory budget.

## References

1. A large scale spatio-temporal dataset of point clouds of maize and tomato plants. <https://www.ipb.uni-bonn.de/data/pheno4d/>. Accessed: 2021-01-1.
2. Jongmin Baek and Andrew Adams. Some useful properties of the permutohedral lattice for Gaussian filtering. *In other words*, 10(1):0, 2009.
3. Jonathan T Barron, Andrew Adams, S YiChang, and Carlos Hernández. Fast bilateral-space stereo for synthetic defocus — Supplemental material. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–15, 2015.
4. Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Juergen Gall. SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences. In *Proc. of the IEEE Int. Conference on Computer Vision (ICCV)*, 2019.
5. Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. The Lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4413–4421, 2018.
6. Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
7. Weikai Chen, Xiaoguang Han, Guanbin Li, Chao Chen, Jun Xing, Yajie Zhao, and Hao Li. Deep RBFNet: Point cloud feature learning using radial basis functions. *arXiv preprint arXiv:1812.04302*, 2018.
8. Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cuDNN: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759*, 2014.
9. Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks. *arXiv preprint arXiv:1904.08755*, 2019.
10. Angela Dai and Matthias Nießner. 3DMV: Joint 3D-multi-view prediction for 3D semantic scene segmentation. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 452–468, 2018.
11. Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3D reconstructions of indoor scenes. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5828–5839, 2017.
12. Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative loss function. *arXiv preprint arXiv:1708.02551*, 2017.
13. Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, pages 3844–3852, 2016.
14. Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3D semantic segmentation with submanifold sparse convolutional networks. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9224–9232, 2018.

15. Xiuye Gu, Yijie Wang, Chongruo Wu, Yong Jae Lee, and Panqu Wang. HPLFlowNet: Hierarchical Permutohedral Lattice FlowNet for Scene Flow Estimation on Large-scale Point Clouds. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3254–3263, 2019.
16. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
17. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 630–645, 2016.
18. Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4700–4708, 2017.
19. Jingwei Huang, Haotian Zhang, Li Yi, Thomas Funkhouser, Matthias Nießner, and Leonidas J Guibas. TextureNet: Consistent local parametrizations for learning from high-resolution signals on meshes. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4440–4449, 2019.
20. Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. PointCNN: Convolution on x-transformed points. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, pages 820–830, 2018.
21. Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. RefineNet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1925–1934, 2017.
22. Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3431–3440, 2015.
23. Jonathan Masci, Davide Boscaini, Michael Bronstein, and Pierre Vandergheynst. Geodesic convolutional neural networks on Riemannian manifolds. In *Workshop Proc. of the IEEE Int. Conference on Computer Vision (ICCV Workshops)*, pages 37–45, 2015.
24. Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model CNNs. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5115–5124, 2017.
25. Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8837–8845, 2019.
26. Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3D reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 32(6):1–11, 2013.
27. Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic Differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
28. Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Real-time progressive 3D semantic segmentation for indoor scenes. In *Proc. of the IEEE Workshop on Applications of Computer Vision*, pages 1089–1098, 2019.
29. Quang-Hieu Pham, Thanh Nguyen, Binh-Son Hua, Gemma Roig, and Sai-Kit Yeung. JSIS3D: Joint semantic-instance segmentation of 3D point clouds with multi-task pointwise networks and multi-value conditional random fields. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8827–8836, 2019.
30. Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017.
31. Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep Hough voting for 3D object detection in point clouds. In *Proc. of the IEEE Int. Conference on Computer Vision (ICCV)*, pages 9277–9286, 2019.
32. Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, pages 5099–5108, 2017.
33. Siamak Ravanbakhsh, Jeff G. Schneider, and Barnabás Póczos. Deep Learning with Sets and Point Clouds. *arXiv preprint arXiv:1611.04500*, 2016.
34. Dario Rethage, Johanna Wald, Jurgen Sturm, Nassir Navab, and Federico Tombari. Fully-convolutional point networks for large-scale point clouds. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 596–611, 2018.
35. Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-assisted Intervention*, pages 234–241, 2015.

36. Radu Alexandru Rosu, Peer Schütt, Jan Quenzel, and Sven Behnke. LatticeNet: Fast point cloud segmentation using permutohedral lattices. *Proc. of Robotics: Science and Systems*, 2020.
37. Hanyu Shi, Guosheng Lin, Hao Wang, Tzu-Yi Hung, and Zhenhua Wang. SpSequenceNet: Semantic Segmentation Network on 4D Point Clouds. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4574–4583, 2020.
38. Patrick Stotko, Stefan Krumpen, Michael Weinmann, and Reinhard Klein. Efficient 3D Reconstruction and Streaming for Group-Scale Multi-Client Live Telepresence. In *Proc. of the IEEE Int. Symposium on Mixed and Augmented Reality (ISMAR)*, pages 19–25, 2019.
39. Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. SplatNet: Sparse lattice networks for point cloud processing. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2530–2539, 2018.
40. Julian Tanke, Oh-Hun Kwon, Patrick Stotko, Radu Alexandru Rosu, Michael Weinmann, Hassan Errami, Sven Behnke, Maren Bennewitz, Reinhard Klein, Andreas Weber, et al. Bonn Activity Maps: Dataset Description. *arXiv preprint arXiv:1912.06354*, 2019.
41. Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3D. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3887–3896, 2018.
42. Lyne Tchammi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. SEGCloud: Semantic segmentation of 3D point clouds. In *Intl. Conf. on 3D Vision (3DV)*, pages 537–547. IEEE, 2017.
43. Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. KPConv: Flexible and deformable convolution for point clouds. In *Proc. of the IEEE Int. Conference on Computer Vision (ICCV)*, pages 6411–6420, 2019.
44. Shenlong Wang, Simon Suo, Wei-Chiu Ma, Andrei Pokrovsky, and Raquel Urtasun. Deep parametric continuous convolutional neural networks. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2589–2597, 2018.
45. Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. SGPN: Similarity group proposal network for 3D point cloud instance segmentation. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2569–2578, 2018.
46. Xinlong Wang, Shu Liu, Xiaoyong Shen, Chunhua Shen, and Jiaya Jia. Associatively segmenting instances and semantics in point clouds. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4096–4105, 2019.
47. Bichen Wu, Xuanyu Zhou, Sicheng Zhao, Xiangyu Yue, and Kurt Keutzer. SqueezeSegv2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a lidar point cloud. *arXiv preprint arXiv:1809.08495*, 2018.
48. Wenxuan Wu, Zhongang Qi, and Li Fuxin. PointConv: Deep convolutional networks on 3D point clouds. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9621–9630, 2019.
49. Yuxin Wu and Kaiming He. Group normalization. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018.
50. Bo Yang, Jianan Wang, Ronald Clark, Qingyong Hu, Sen Wang, Andrew Markham, and Niki Trigoni. Learning object bounding boxes for 3D instance segmentation on point clouds. *arXiv preprint arXiv:1906.01140*, 2019.
51. Li Yi, Vladimir G Kim, Duygu Ceylan, I Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, Leonidas Guibas, et al. A scalable active framework for region annotation in 3D shape collections. *ACM Transactions on Graphics (ToG)*, 35(6):210, 2016.
52. Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas J Guibas. GSPN: Generative shape proposal network for 3D instance segmentation in point cloud. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3947–3956, 2019.
53. Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, pages 3391–3401, 2017.