# Long Range Pooling for 3D Large-Scale Scene Understanding

Xiang-Li Li[1]    Meng-Hao Guo[1]    Tai-Jiang Mu[1]    Ralph R. Martin[2]    Shi-Min Hu[1]
[1]Tsinghua University    [2]Cardiff University

lixl19@mails.tsinghua.edu.cn, gmh20@mails.tsinghua.edu.cn

taijiang@tsinghua.edu.cn, martinrr@cardiff.ac.uk, shimin@tsinghua.edu.cn

## Abstract

*Inspired by the success of recent vision transformers and large kernel design in convolutional neural networks (CNNs), in this paper, we analyze and explore essential reasons for their success. We claim two factors that are critical for 3D large-scale scene understanding: **a larger receptive field** and **operations with greater non-linearity**. The former is responsible for providing long range contexts and the latter can enhance the capacity of the network. To achieve the above properties, we propose a simple yet effective long range pooling (LRP) module using dilation max pooling, which provides a network with a large adaptive receptive field. LRP has few parameters, and can be readily added to current CNNs. Also, based on LRP, we present an entire network architecture, LRPNet, for 3D understanding. Ablation studies are presented to support our claims, and show that the LRP module achieves better results than large kernel convolution yet with reduced computation, due to its non-linearity. We also demonstrate the superiority of LRPNet on various benchmarks: LRPNet performs the best on Scan-Net and surpasses other CNN-based methods on S3DIS and Matterport3D. Code will be made publicly available.*

## 1. Introduction

With the rapid development of 3D sensors, more and more 3D data is becoming available from a variety of applications such as autonomous driving, robotics, and augmented/virtual reality. This 3D data requires analyzing and understanding. Efficient and effective processing of such 3D data has become an important challenge.

Various data structures, including point clouds, meshes, multi-view images, voxels, *etc*, have been proposed for representing 3D data [69], and different network architectures are designed to process them. Unlike 2D image data, a point cloud or mesh, is irregular and unordered. These characteristics mean that typical CNNs cannot directly be applied to such 3D data. Thus, specially designed networks such as MLPs [44], CNNs [24, 36, 65], GNNs [48, 64] and trans-
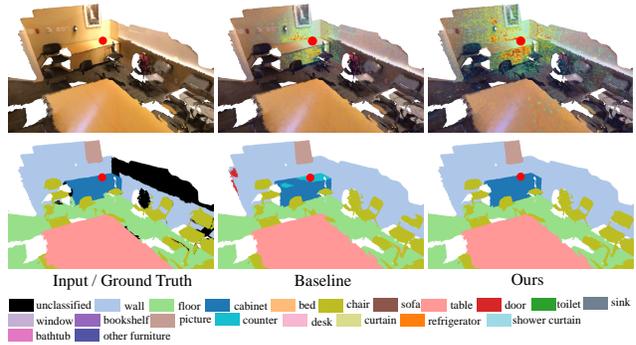


Figure 1.    Qualitative results of Effective Receptive Field (ERF) [41]. **Left**: the input and ground truth. **Middle**: the ERF and results of baseline described in section 3.3. **Right**: the ERF and results of LRPNet (ours). The stained areas around the positions of interest (red dots) represent the range of receptive fields, with green to red representing the increasing strength of response. Our method correctly segments the cabinet from the wall with the proposed effective receptive field.

formers [15, 42, 75] are proposed to perform effective deep learning on them. However, processing large-scale point cloud or mesh is computationally expensive. Multi-view images contain multiple different views of a scene. Typically, CNNs are used to process each view independently and a fusion module is used to combine the results. Nevertheless, there is unavoidable 3D information loss due to the finite or insufficient number of views. Voxel data has the benefit of regularity like images, even if the on-surface voxels only contain sparse data. The simplicity of the structure helps to maintain high performance, so this paper focuses on processing voxel data.

Learning on 3D voxels can be easily implemented by directly extending the well studied 2D CNN networks to 3D [47, 62, 68]. Considering that 3D voxel data is inherently sparse, some works usually adopt specially designed 3D sparse CNNs [6,14,51] for large-scale scene understanding. Since only the surface of the real scene data has values, the neighbors around each voxel do not necessarily be-

1

long to the same object, so a large receptive field is more conducive to feature extraction. Sparse convolution has advantages in modeling local structures, but it ignores long range contexts, which are critical for 3D scene understanding tasks. Transformers [15, 32, 66, 67] have proved useful in processing 3D scenes, as they can capture the global relationship with their larger receptive field and a better way to interact. However, they usually have a quadratic computational complexity, which brings a heavy computing cost.

A straightforward way to incorporate the long range contexts into the learning of 3D voxel data is to exploit a large kernel 3D convolution. However, the number of network parameters and the amount of computation would increase cubically due to the additional dimension compared to 2D images. Besides, current ways of feature interaction or aggregation, such as average pooling and convolution, usually adopt a linear combination of features of all the locations in the receptive field. This works for 2D images since all the locations in the receptive field have valid values, which, however, does not hold for 3D voxels due to the sparsity nature of 3D scenes. Directly applying such linear interaction or aggregation on the voxel that has few neighbors in the receptive field would make the feature of that voxel too small or over-smoothed and thus less informative.

Taking the above into consideration, and to achieve a trade-off between the quality of results and computation, we propose long range pooling (LRP), a simple yet effective module for 3D scene segmentation. Compared with previous sparse convolution networks [6, 14], LRP is capable of increasing the effective receptive field with a negligible amount of computational overhead. Specifically, we achieve a large receptive field by proposing a novel dilation max pooling to enhance the non-linearity of the neural network. We further add a receptive field selection module, so that each voxel can choose a suitable receptive field, which is adaptive to the distribution of voxels. The above two components comprise the LRP module. Unlike dilation convolution, which is often used to enlarge the receptive field for 2D images [17, 18], our method can achieve a large receptive field with fewer parameters and computation by using dilation max pooling. Furthermore, LRP is a simple, efficient and parameterless module that can be readily incorporated into other networks. We construct a more capable neural network, *LRPNet*, by adding LRP at the end of each stage of the sparse convolution network, introduced by VMNet [26].

Experimental results show that LRPNet achieves a significant improvement in 3D segmentation accuracy on large-scale scene datasets, including ScanNet [7], S3DIS [1] and Matterport3D [3]. Qualitative results are illustrated in Figure 1, which shows the improvement of a larger receptive field. We also experimentally compare the effects of the different receptive fields by reducing the num-ber of dilation max pooling of LRP. Moreover, we explore the influence of non-linearity of LRP module by replacing max pooling with average pooling or convolution. Ablation results show that a larger receptive field and operations with greater non-linearity will improve the segmentation accuracy.

Our contributions are thus:

- a simple and effective module, the long range pooling (LRP) module, which provides a network with a large adaptive receptive field without a large number of parameters,

- a demonstration that a larger receptive field and operations with greater non-linearity enhance the capacity of a sparse convolution network, and

- a simple sparse convolution network using the LRP module, which achieves superior 3D segmentation results on various large-scale 3D scene benchmarks.

## 2. Related Work

### 2.1. Deep Learning in 3D Vision

Deep learning has achieved great success in the field of 3D vision in recent years [69]. However, compared to 2D images, 3D data is complex and has diverse representations, such as multi-view images, point clouds, voxels, meshes and so on. With the availability of some large 3D datasets [1, 7, 13, 68], 3D vision has entered a period of rapid development. Due to the success of neural networks in processing 2D images, some works [50, 53, 56] have used 2D convolutional networks to extract features from 2D views, and applied a fusion module for 3D understanding. In fact, the raw data is usually in the form of a point cloud, so some researchers [23, 36, 44, 45, 59, 65] have designed point neural networks for point cloud processing. Due to the high computational complexity of processing point cloud networks, researchers [4, 6, 14] began to use hash functions to search for neighbors, and to apply 3D sparse convolution to achieve efficient voxel understanding. Meshes have not only position information, but also topological structure. Specially designed mesh neural networks [20, 24, 31, 49, 54] take advantage of this topological information. In addition, to further improve the effectiveness of 3D semantic understanding, fusion learning of multiple types of 3D data has also gradually developed, such as point cloud + mesh [48], point cloud + image [25], voxel + mesh [26], and so on.

Most works [6, 20, 44] focus on how neural networks can make better use of 3D data, and thus construct specific modules to enhance the ability to extract features. Point networks from PointNet [44] to PointConv [65], PointCNN [36] improve the ability to extract local information. Going from point networks [36, 45] to point transformers [15, 32, 66, 67] expands the scope of the model's

receptive field and the ability to recognize the interaction between features. However, the computational demands of these methods are higher, hindering their application to large scenes. In order to balance computational efficiency and accuracy, voxel-based networks [6, 14] use 3D sparse convolution to aggregate local features. Limited by a large number of 3D convolution parameters, it is difficult to improve results by directly increasing the receptive field. LargeKernel3D [4] expands the receptive field of 3D convolution by depth-wise convolution and dilation convolution. However, compared to transformers [15, 32, 66, 67], the receptive field of this approach is still small, and it introduces more parameters and computation. To solve these problems, we propose the long range pooling method, which can expand the receptive field with few additional parameters and introduce operations with greater non-linearity to enhance the capability of the neural network.

## 2.2. Vision Transformers

The transformer network architecture comes from natural language processing [10, 60]. Recently, due to its strong modeling capability, it has quickly provided leading methods for various vision tasks, including image classification [12, 16], object detection [37, 39], semantic segmentation [70, 76], image generation [29, 34], and self-supervised learning [2, 21]—see the surveys in [19, 30]. The transformer architecture is also introduced into 3D vision by PCT [15] and PT [75] almost simultaneously, which propose a 3D transformer based on global attention and a 3D transformer based on local attention, respectively.

The core module of a transformer is the self-attention block, which models relationships by calculating pairwise similarity between any two feature points. We believe the success of self-attention arises for two reasons: (i) self-attention captures *long range dependencies*, and (ii) the matrix multiplication of attention and value, and softmax function provide *strong non-linearity*.

## 2.3. Large Kernel Design in CNNs

Inspired by the success of vision transformers, researchers have challenged the traditional small kernel design of CNNs [22, 52] and suggested the use of large convolution kernels for visual tasks [11, 17, 18, 38, 40, 46, 73]. For example, ConvNeXt [40] suggest directly adopting a 7×7 depth-wise convolution, while the Visual Attention Network (VAN) [18] uses a kernel size of $21 \times 21$ and introduces an attention mechanism. RepLKNet [11] introduces a $31 \times 31$ convolution by using a reparameterization technique. Recently, this design is also introduced into 3D field by LargeKernel3D [4]. Analyzing these previous works, we observe that VAN [18] gives the best results, which we believe is because it introduces non-linearity via the Hadamard product, in addition to long range dependencies.

## 3. Method

Our analysis above suggests that long range interactions and greater non-linearity may be the key to success. Accordingly, we have designed a simple yet effective long range pooling (LRP) module, based on these principles. We now introduce our LRP and LRPNet in detail.

### 3.1. Dilation Module for 3D Voxel

To our knowledge, we are the first to implement a large receptive field network for processing 3D voxel by using dilation max pooling. In order to make the paper self-contained, we revisit the decomposition of a large kernel by dilated operations in 2D CNNs.

In fact, general convolutional neural networks can implicitly realize large receptive fields over the whole network. With increasing network depth, the receptive field of the last layer of features gradually increases—this is one of the reasons why deep neural networks can be effective. Meanwhile, previous works [17, 19] have shown that increasing the local receptive field can give better results than increasing the depth. On the other hand, it is almost impossible to apply large kernels directly to the network because of their high computation load and large number of parameters. Currently, the commonly used large kernel size is $31 \times 31$ in 2D images [11], which improves accuracy while also introducing more parameters and computation. Therefore, most of the previous works have decomposed the large kernel module and used the computationally friendly small kernel module to approximate the large kernel convolution [18].

For 2D images, several dilation convolutions are usually used to achieve a decomposed large kernel convolution. However, 3D data is sparse, and the sparsity affects the response strength of the convolution, which makes large kernel decomposition more difficult. Besides, due to the sparsity, the weight convergence of convolution is too slow. Therefore, we use dilated pooling to achieve a large receptive field while keeping a low computational cost.

### 3.2. Long Range Pooling Module

Two direct strategies can achieve a large receptive field and non-linearity: self-attention, and max pooling with a large window size. The large computational load placed by self-attention limits its application to 3D data. Max pooling also imposes huge computational demands as the window grows. We avoid the computational cost of max pooling by introducing sparse dilated pooling.

As shown in Figure 2, we approximate pooling with a large window size by stacking three 3×3×3 pooling modules with different dilation rates. This achieves a progres-
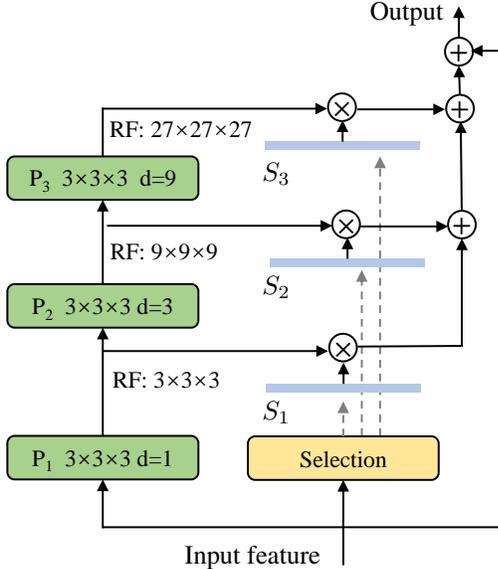
Figure 2. Long Range Pooling Module. Dilation max pooling is used to enlarge the receptive field and also outputs features of the different receptive field. $P_1, P_2, P_3$: max pooling with kernel size $3 \times 3 \times 3$ and dilation $1, 3, 9$. **RF**: the receptive field. **Selection**: a linear module producing the selection weight for each voxel.

sive increase in receptive field size while keeping the computational cost low. Besides, fixed receptive fields will limit the network's ability to model objects of different sizes. To choose a proper receptive field that can be adaptive to the distribution of voxels, a receptive field selection block is designed to choose a suitable window size for each voxel.

Features of each location need to interact with each other or be aggregated to learn the local or global context. There are three typical ways to achieve this: 1) average pooling; 2) convolution; and 3) max pooling. The former two are linear and parameterized (fixed for average pooling and learnable for convolution); the last one is nonlinear and non-parametric. Considering the sparsity and non-uniformity of 3D voxels, the shared linear interactive ways would make the feature of a voxel too small (for voxels having few neighbors) or over-smoothed (for voxels having many neighbors). Max-pooling, in turn, will always draw the feature of a voxel to the most informative one. We illustrate the features learned by different interactive manners in Figure 5. We also compare these three different interactive manners in section 4.4.2 and the experiments show max pooling works best, which supports our viewpoint that more non-linearity is critical for 3D scene understanding.
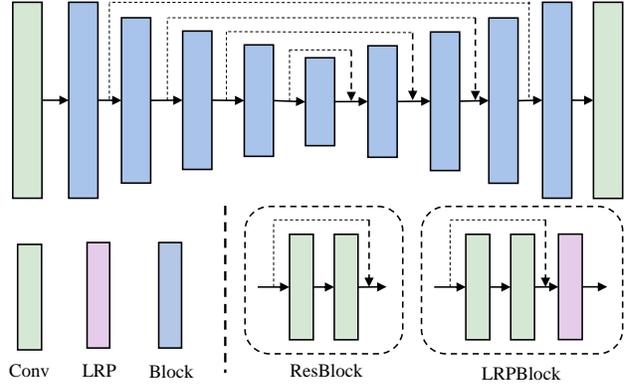


Figure 3. Network architecture for scene segmentation. We use the sparse convolutional U-Net [26] as our baseline, consisting of four stages. For the baseline, we use ResBlock as the basic block to extract features, and we replace the ResBlock with LRPBlock to build LRPNet.

The overall LRP module can be formulated as

$$S = (S_1, S_2, S_3) = \text{linear}(x), \quad (1)$$
$$P_1 = p(x, 1), P_2 = p(P_1, 3), P_3 = p(P_2, 9), \quad (2)$$
$$Output = S_1 P_1 + S_2 P_2 + S_3 P_3, \quad (3)$$

where $x$ denotes the input features, S is the selected weight, which is divided into $S_1, S_2, S_3$ for selecting the features of different window sizes, and $p(x, k)$ means max pooling with kernel size 3 and dilation k. Here, k is set as 1, 3 and 9 by default.

### 3.3. Long Range Pooling Network

In order to verify the effectiveness of LRP module for processing 3D scene data, we have designed LRPNet, based on the U-Net architecture commonly used in previous work [4, 6, 26, 48]. We use the sparse convolutional U-Net implemented by VMNet [26] as the baseline. Figure 3 shows the baseline model, a U-Net network with 4 stages; the basic module is a ResBlock, built from two 3D Convolution+BatchNorm+ReLU blocks and an additional residual path.

As we have already noted, the LRP module can readily be incorporated into existing networks. Therefore, we build LRPNet by simply adding an LRP module after ResBlock in the baseline model. In addition to placing the LRP behind the ResBlock, we alternatively experimented with placing it front, middle, and as a parallel branch. These experiments and analysis are discussed in Section 4.4.1.

## 4. Experiments

We performed various experiments on several semantic segmentation datasets, to verify our claims and compare our results to those of other methods.

## 4.1. Datasets and metrics

ScanNet v2 [7] is a large-scene benchmark, with 1,613 scenes split 1201:312:100 for training, validation and testing. It has 20 semantic classes. We follow the same protocol as previous works [6, 14] and use mean class intersection over union (mIoU) as the metric to evaluate results. Since the annotations of the test set on ScanNetv2 [7] are not available, we conducted the runtime complexity experiments and ablation experiments on the validation set.

S3DIS [1] is a large-scene semantic parsing benchmark, containing 6 areas with 271 rooms, annotated with 13 categories. Following [6, 15, 59, 67], we use Area 5 as the test set, other areas for training and mIoU for evaluation. We also report overall point-wise accuracy (OA) and mean class accuracy (mAcc).

Matterport3D [3] is a large-scene RGB-D dataset, with 90 building-scale scenes annotated in 21 categories. Since the semantic labels are annotated for faces, we project them to the vertices using the same method as in [48], and then test on scene vertices directly. Following [26, 48], we use mAcc for evaluation.

## 4.2. Implementation details

Our experiments were all conducted on the full scenes without cropping. During training and inferencing, we just used vertex colors as input. Following [6, 26], we voxelized the input point cloud at a resolution of 2 cm for ScanNet v2. Following common practice, the input points of S3DIS and Matterport3D are voxelized at a resolution of 5 cm. During inferencing, we calculated metrics by projecting the predictions back to the raw point clouds using the nearest neighbor.

During training, we performed scaling, z-axis rotation, translation and chromatic jitter for data augmentation [26]. For all datasets, we minimized the cross entropy loss using the SGD optimizer with a poly scheduler decaying from a learning rate of 0.1, following [4, 6, 26]. Following previous work [27, 48], we used class weights for cross entropy loss on Matterport3D. For ScanNet v2 and Matterport3D, we trained our model for 500 epochs with batch size set to 8, and for S3DIS, we trained for 1,000 epochs with batch size set to 4.

Our experiments were conducted on RTX 3090 GPUs and A5000 GPUs. All runtime complexity experiments are conducted on one RTX 3090 GPU.

## 4.3. Comparison to other methods

We compared our model to other point-based and voxel-based state-of-the-art methods on ScanNet v2, S3DIS and Matterport3D, with results shown in Tables 1–3, respectively. For these datasets, our method gave the best overall results compared CNN-based methods. In particular, our method surpasses both the popular state-of-the-art

Table 1. mIoU (%) scores for various methods on the ScanNet v2 3D semantic benchmark, for validation and test sets. The best number is in boldface. "-" means the number is unavailable.

| Method | Input | Val | Test |
|---|---|---|---|
| PointNet++ [45] | point | 53.5 | 55.7 |
| 3DMV [26] | point | - | 48.4 |
| PointCNN [36] | point | - | 45.8 |
| PointConv [65] | point | 61.0 | 66.6 |
| JointPointBased [5] | point | 69.2 | 63.4 |
| PointASNL [72] | point | 63.5 | 66.6 |
| RandLA-Net [23] | point | - | 64.5 |
| KPConv [59] | point | 69.2 | 68.6 |
| PointTransformer [75] | point | 70.6 | - |
| SparseConvNet [14] | voxel | 69.3 | 72.5 |
| MinkowskiNet [6] | voxel | 72.2 | 73.6 |
| LargeKernel3D [4] | voxel | 73.2 | 73.9 |
| Fast Point Transformer [43] | voxel | 72.1 | - |
| Stratified Transformer [32] | point | 74.3 | 73.7 |
| LRPNet (ours) | voxel | **75.0** | **74.2** |

Table 2. Several scores (%) for various methods on the S3DIS segmentation benchmark. The best number is in boldface. "-" means the number is unavailable.

| Method | Input | OA | mAcc | mIoU |
|---|---|---|---|---|
| PointNet [44] | point | - | 49.0 | 41.1 |
| SegCloud [58] | point | - | 57.4 | 48.9 |
| TangentConv [57] | point | - | 62.2 | 52.6 |
| PointCNN [36] | point | 85.9 | 63.9 | 57.3 |
| HPEIN [28] | point | 87.2 | 68.3 | 61.9 |
| GACNet [61] | point | 87.8 | - | 62.9 |
| PAT [74] | point | - | 70.8 | 60.1 |
| ParamConv [63] | point | - | 67.0 | 58.3 |
| SPGraph [33] | point | 86.4 | 66.5 | 58.0 |
| PCT [15] | point | - | 67.7 | 61.3 |
| SegGCN [35] | point | 88.2 | 70.4 | 63.6 |
| PAConv [71] | point | - | - | 66.6 |
| KPConv [59] | point | - | 72.8 | 67.1 |
| MinkowskiNet [6] | voxel | - | 71.7 | 65.4 |
| Fast Point Transformer [43] | voxel | - | 77.3 | 70.1 |
| PointTransformer [75] | point | 90.8 | 76.5 | 70.4 |
| Stratified Transformer [32] | point | **91.5** | **78.1** | **72.0** |
| LRPNet (ours) | voxel | 90.8 | 74.9 | 69.1 |

MinkowskiNet [6], and the voxel and mesh fusion algorithm VMNet [26]. Our method even outperforms transformer-

5

Table 3. Mean class accuracy (%) scores on the Matterport3D test set. The best number is in boldface.

| Method | mAcc | wall | floor | cab | bed | chair | sofa | table | door | wind | shf | pic | cntr | desk | curt | ceil | fridg | show | toil | sink | bath | other |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SplatNet [55] | 26.7 | **90.8** | 95.7 | 30.3 | 19.9 | 77.6 | 36.9 | 19.8 | 33.6 | 15.8 | 15.7 | 0.0 | 0.0 | 0.0 | 12.3 | 75.7 | 0.0 | 0.0 | 10.6 | 4.1 | 20.3 | 1.7 |
| PointNet++ [45] | 43.8 | 80.1 | 81.3 | 34.1 | 71.8 | 59.7 | 63.5 | **58.1** | 49.6 | 28.7 | 1.1 | 34.3 | 10.1 | 0.0 | 68.8 | 79.3 | 0.0 | 29.0 | 70.4 | 29.4 | 62.1 | 8.5 |
| ScanComplete [9] | 44.9 | 79.0 | **95.9** | 31.9 | 70.4 | 68.7 | 41.4 | 35.1 | 32.0 | 37.5 | 17.5 | 27.0 | 37.2 | 11.8 | 50.4 | **97.6** | 0.1 | 15.7 | 74.9 | 44.4 | 53.5 | 21.8 |
| TangentConv [57] | 46.8 | 56.0 | 87.7 | 41.5 | 73.6 | 60.7 | 69.3 | 38.1 | 55.0 | 30.7 | 33.9 | 50.6 | 38.5 | 19.7 | 48.0 | 45.1 | 22.6 | 35.9 | 50.7 | 49.3 | 56.4 | 16.6 |
| 3DMV [8] | 56.1 | 79.6 | 95.5 | 59.7 | 82.3 | 70.5 | 73.3 | 48.5 | 64.3 | 55.7 | 8.3 | 55.4 | 34.8 | 2.4 | **80.1** | 94.8 | 4.7 | 54.0 | 71.1 | 47.5 | 76.7 | 19.9 |
| TextureNet [27] | 63.0 | 63.6 | 91.3 | 47.6 | 82.4 | 66.5 | 64.5 | 45.5 | 69.4 | 60.9 | 30.5 | **77.0** | 42.3 | 44.3 | 75.2 | 92.3 | 49.1 | 66.0 | 80.1 | **60.6** | 86.4 | 27.5 |
| DCM-Net [48] | 66.2 | 78.4 | 93.6 | **64.5** | **89.5** | 70.0 | **85.3** | 46.1 | **81.3** | **63.4** | 43.7 | 73.2 | 39.9 | 47.9 | 60.3 | 89.3 | 65.8 | 43.7 | 86.0 | 49.6 | 87.5 | **31.1** |
| VMNet [26] | 67.2 | 85.9 | 94.4 | 56.2 | **89.5** | **83.7** | 70.0 | 54.0 | 76.7 | 63.2 | 44.6 | 72.1 | 29.1 | 38.4 | 79.7 | 94.5 | 47.6 | 80.1 | 85.0 | 49.2 | 88.0 | 29.0 |
| LRPNet (Ours) | **70.7** | 83.7 | 95.0 | 58.0 | 88.2 | 81.3 | 79.0 | 54.3 | 78.5 | 60.0 | **63.4** | 70.7 | **48.7** | **52.0** | 70.0 | 93.7 | **66.1** | **87.4** | **89.0** | 47.2 | **88.1** | 30.5 |

Table 4. Number of network parameters and speed for various models on the ScanNet v2 validation set.

| Method | Params (M) | Runtime (ms) | mIoU (%) |
|---|---|---|---|
| SparseConvNet [14] | 30.1 | 173.5 | 69.3 |
| MinkowskiNet [6] | 37.9 | 166.1 | 72.2 |
| Fast Point Transformer [43] | 37.9 | 341.4 | 72.0 |
| Stratified Transformer [32] | 18.8 | 1149.9 | 74.3 |
| Baseline | 8.1 | 38.1 | 71.2 |
| LRPNet (Ours) | 8.5 | 67.9 | 75.0 |

Table 5. Ablation study on LRP module. **Baseline**: U-Net described in section 3.3. **MaxPool**: The max pooling used in LRP. **Dilation**: Dilation used in max pooling (default for LRP is 1, 3, 9). **Selection**: the selection module in LRP. **Params**: parameters of the network. **Runtime**: the average time for inferencing one scene. **mIoU**: the segmentation accuracy metric (%).

| Baseline | MaxPool | Dilation | Selection | Params (M) | Runtime (ms) | mIoU |
|---|---|---|---|---|---|---|
| ✓ | | | | 8.1 | 38.1 | 71.2 |
| ✓ | ✓ | | | 8.1 | 74.2 | 72.6 |
| ✓ | ✓ | ✓ | | 8.1 | 65.0 | 73.7 |
| ✓ | ✓ | | ✓ | 8.5 | 76.4 | 73.1 |
| ✓ | ✓ | ✓ | ✓ | 8.5 | 67.9 | 75.0 |

based methods, which learn the long range context at a cost of heavy computation, on ScanNetv2.

In addition, we also compared the number of network parameters used and the speed of our algorithm to two SOTA voxel methods [6, 14] and two SOTA transformer-based methods [32, 43] using the ScanNet v2 validation set. For a fair comparison, we tested the validation set with 312 scenes and average the inference time of each scene to calculate the runtime complexity. We used the MinkowskiNet [6] implemented by Fast Point Transformer [43], which is faster than the original implementation. Table 4 shows that our algorithm achieves better results with fewer parameters and faster speed. Notice that our method is nearly 16× faster than Stratified Transformer [32]. LRPNet achieved a significant increase in mIoU by just adding the LRP module to the end of each stage of the baseline, with only a few extra parameters introduced by the selection module.

We further compared qualitative results on the ScanNet v2 validation set, for our model and the baseline model. Figure 4 shows LRPNet's segmentations are closer to the ground truth than the baseline's segmentations. Our method still works well even in hard cases, e.g. for the table and the cabinet in the corner in the first and four rows, whereas the baseline method does not. LRPNet is also more capable of segmenting out broken objects, as in the second and third rows of bookshelves and walls in Figure 4.

## 4.4. Ablation study

In section 3.2, we use dilation max pooling to expand the effective receptive field, and use a selection module to allow each voxel to select receptive fields based on features. we conducted the ablation experiments on the validation set of ScanNet v2 to verify the design of the LRP module. Specifically, we choose the sparse convolutional U-Net implemented by VMNet [26] as the baseline, and design new network using max pooling without any dilation. This new network is further enhanced with dilation or/and the receptive field selection module.

As shown in Table 5, max pooling can enlarge the receptive field to improve the results, and dilation max pooling can achieve a larger receptive field and a better result. Besides, the selection module can enhance the capability of the network with only introducing a few parameters.

### 4.4.1 Position of LRP

In section 3.3, we added the LRP module **after** each stage of the baseline (VMNet [26]), which we name it (**After**). Actually, we may also put the LRP module before the stage (**Before**), in the middle of the stage (**Middle**), or as an additional path parallel to the stage (**Parallel**). Table 6 shows
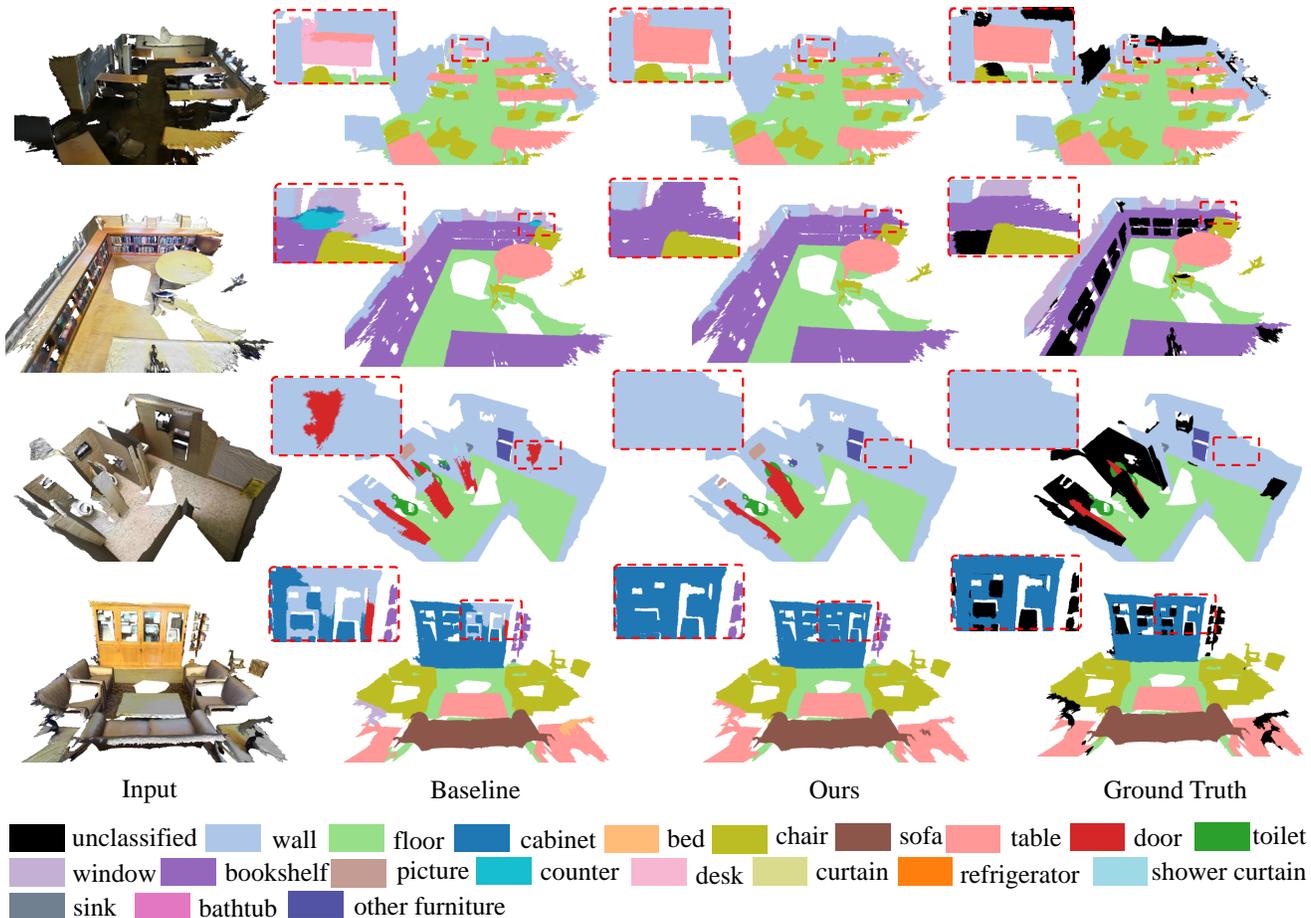
Figure 4. Various ScanNet v2 validation set segmentation results. Red dotted boxes highlight differences between our results and the baseline results.

Table 6. LRP position study. **Before**: LRP added before the stages of baseline. **After**: LRP added after the stages. **Middle**: LRP added at the middle of stages. **Parallel**: LRP as an additional branch parallel to the stage.

| Position | Params (M) | Runtime (ms) | mIoU (%) |
|---|---|---|---|
| Before | 8.8 | 77.1 | 74.0 |
| Middle | 8.5 | 69.6 | 73.4 |
| Parallel | 8.5 | 68.4 | 73.1 |
| After (Ours) | 8.5 | 67.9 | 75.0 |

the ablation of different settings. It shows the LRP module works better when added after each stage than in other positions, also being more computationally friendly.

### 4.4.2 Operations with non-linearity

As we mentioned in section 3.2, max pooling has more non-linearity than average pooling and convolution, and the operations with non-linearity will enhance the capability of the network. To investigate the impacts of non-linearity, we conducted non-linearity experiments by replacing the max pooling (**MaxPool**) of LRP with convolution (**Conv**) or average pooling (**AvgPool**). In addition, we tested the effect of different receptive fields with these operations (MaxPool, AvgPool, Conv). In Table 7, $\times N$ is the size of the receptive field of $N \times N \times N$ and $[\times 3, \times 9 \times 27]$ means the outputs of the LRP module is selected from the features with the receptive field of $3 \times 3 \times 3$, $9 \times 9 \times 9$, and $27 \times 27 \times 27$.

As Table 7 shows, although Conv introduces a large number of parameters, the results of MaxPool can still exceed those of AvgPool and Conv in most cases, which proves that the non-linearity of max pooling will enhance the networks for 3D segmentation.
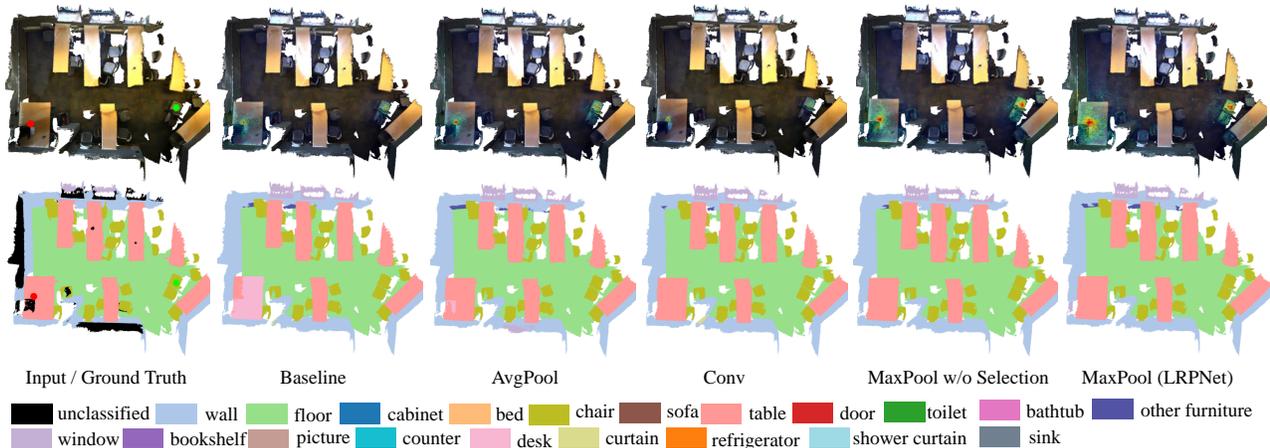
Figure 5. Visualization of Effective Receptive Field (ERF). **Red Dot and Green Pentagon**: two different positions of interest (marked only in the input). **First Row**: the ERFs of different methods. **Second Row**: the ground truth and predictions. **From Left to Right**: input, baseline, average pooling, convolution, max pooling without selection, max pooling (LRPNet).

Table 7. Study on non-linearity and range of receptive field. **Max-Pool**: dilation max pooling. **AvgPool**: dilation average pooling. **Conv**: dilation convolution. **Range**: The range and number of receptive fields.

| Method | Range | Params (M) | Runtime (ms) | mIoU (%) |
|---|---|---|---|---|
| MaxPool | [×3] | 8.2 | 66.0 | 72.7 |
| MaxPool | [×9] | 8.2 | 64.8 | 72.8 |
| MaxPool | [×27] | 8.2 | 66.1 | 73.9 |
| MaxPool | [×9, ×27] | 8.4 | 66.0 | 74.0 |
| MaxPool | [×3, ×9, ×27] | 8.5 | 67.9 | 75.0 |
| AvgPool | [×27] | 8.2 | 58.2 | 73.0 |
| AvgPool | [×9, ×27] | 8.4 | 58.2 | 73.2 |
| AvgPool | [×3, ×9, ×27] | 8.5 | 60.1 | 73.9 |
| Conv | [×27] | 17.4 | 65.5 | 72.0 |
| Conv | [×9, ×27] | 17.5 | 66.2 | 73.3 |
| Conv | [×3, ×9, ×27] | 17.6 | 66.8 | 73.8 |

### 4.4.3 Range of LRP module

As described in section 3, while large receptive fields can improve network results, we used dilation max pooling to achieve large receptive fields and a selection module to give the network the ability to select the receptive field size according to the voxel feature.

Table 7 shows that as we increase the levels of the receptive field of the LRP module, whether using convolution, max pooling or average pooling, the accuracy of results is gradually improved. We also used different ranges of the receptive fields with the same number of levels. See Ta-

ble 7 lines 1–3: we fixed the number of receptive fields of LRP module to one level, and then compared results for [×3], [×9], and [×27] receptive fields. With the increasing receptive field for the LRP module, the model also clearly improved.

In addition, we visualized the Effective Receptive Field [41] of different methods at two positions. As shown in Figure 5, the features of interest are in the center of the table and the center of the chair. LRPNet has a larger receptive field than baseline, which helps to segment the table better. Columns 3,4 and 6 in Figure 5 show that the operations with linearity can not capture the long range context and their responses are mostly concentrated around the points of interest, while max pooling with non-linearity can make good use of the long range features. Furthermore, as shown in the last two columns of Figure 5, the selection module can give the network the ability to select the appropriate receptive field according to the voxel features.

## 5. Conclusion

We have proposed a simple yet efficient module, the long range pooling module, which can provide an adaptive large receptive field and improve the modeling capacity of the network. It is easy to add to any existing networks; we use it to construct LRPNet, which achieves state-of-the-art results for large-scene segmentation with almost no increase in the number of parameters. Since the LRP module is simple and easy to use, we hope to apply it to 3D object detection, 3D instance segmentation and 2D image processing, further exploring the principle of improving network performance by use of a large receptive field and operations with greater non-linearity in the future.

# References

[1] Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1534–1543, 2016. 2, 5

[2] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 3

[3] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. 2, 5

[4] Yukang Chen, Jianhui Liu, Xiaojuan Qi, Xiangyu Zhang, Jian Sun, and Jiaya Jia. Scaling up kernels in 3d cnns. *arXiv preprint arXiv:2206.10555*, 2022. 2, 3, 4, 5

[5] Hung-Yueh Chiang, Yen-Liang Lin, Yueh-Cheng Liu, and Winston H Hsu. A unified point-based framework for 3d segmentation. In *2019 International Conference on 3D Vision (3DV)*, pages 155–163. IEEE, 2019. 5

[6] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. 1, 2, 3, 4, 5, 6

[7] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 2, 5

[8] Angela Dai and Matthias Nießner. 3dmv: Joint 3d-multi-view prediction for 3d semantic scene segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 452–468, 2018. 6

[9] Angela Dai, Daniel Ritchie, Martin Bokeloh, Scott Reed, Jürgen Sturm, and Matthias Nießner. Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2018. 6

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3

[11] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11963–11975, 2022. 3

[12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3

[13] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 2

[14] Benjamin Graham, Martin Engelcke, and Laurens Van Der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9224–9232, 2018. 1, 2, 3, 5, 6

[15] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021. 1, 2, 3, 5

[16] Meng-Hao Guo, Zheng-Ning Liu, Tai-Jiang Mu, and Shi-Min Hu. Beyond self-attention: External attention using two linear layers for visual tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 3

[17] Meng-Hao Guo, Cheng-Ze Lu, Qibin Hou, Zhengning Liu, Ming-Ming Cheng, and Shi-Min Hu. Segnext: Rethinking convolutional attention design for semantic segmentation. *arXiv preprint arXiv:2209.08575*, 2022. 2, 3

[18] Meng-Hao Guo, Cheng-Ze Lu, Zheng-Ning Liu, Ming-Ming Cheng, and Shi-Min Hu. Visual attention network. *arXiv preprint arXiv:2202.09741*, 2022. 2, 3

[19] Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai Zhang, Ralph R Martin, Ming-Ming Cheng, and Shi-Min Hu. Attention mechanisms in computer vision: A survey. *Computational Visual Media*, pages 1–38, 2022. 3

[20] Rana Hanocka, Amir Hertz, Noa Fish, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. Meshcnn: a network with an edge. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2

[21] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 3

[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3

[23] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11108–11117, 2020. 2, 5

[24] Shi-Min Hu, Zheng-Ning Liu, Meng-Hao Guo, Jun-Xiong Cai, Jiahui Huang, Tai-Jiang Mu, and Ralph R Martin. Subdivision-based mesh convolution networks. *ACM Transactions on Graphics (TOG)*, 41(3):1–16, 2022. 1, 2

[25] Wenbo Hu, Hengshuang Zhao, Li Jiang, Jiaya Jia, and Tien-Tsin Wong. Bidirectional projection network for cross dimension scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14373–14382, 2021. 2

[26] Zeyu Hu, Xuyang Bai, Jiaxiang Shang, Runze Zhang, Jiayu Dong, Xin Wang, Guangyuan Sun, Hongbo Fu, and Chiew-Lan Tai. Vmnet: Voxel-mesh network for geodesic-aware

3d semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15488–15498, 2021. 2, 4, 5, 6

[27] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4440–4449, 2019. 5, 6

[28] Li Jiang, Hengshuang Zhao, Shu Liu, Xiaoyong Shen, Chi-Wing Fu, and Jiaya Jia. Hierarchical point-edge interaction network for point cloud semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10433–10441, 2019. 5

[29] Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two pure transformers can make one strong gan, and that can scale up. *Advances in Neural Information Processing Systems*, 34:14745–14758, 2021. 3

[30] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):1–41, 2022. 3

[31] Alon Lahav and Ayellet Tal. Meshwalker: Deep mesh understanding by random walks. *ACM Transactions on Graphics (TOG)*, 39(6):1–13, 2020. 2

[32] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8500–8509, 2022. 2, 3, 5, 6

[33] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4558–4567, 2018. 5

[34] Kwonjoon Lee, Huiwen Chang, Lu Jiang, Han Zhang, Zhuowen Tu, and Ce Liu. Vitgan: Training gans with vision transformers. *arXiv preprint arXiv:2107.04589*, 2021. 3

[35] Huan Lei, Naveed Akhtar, and Ajmal Mian. Seggcn: Efficient 3d point cloud segmentation with fuzzy spherical kernel. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11611–11620, 2020. 5

[36] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems*, 31, 2018. 1, 2, 5

[37] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. *arXiv preprint arXiv:2203.16527*, 2022. 3

[38] Shiwei Liu, Tianlong Chen, Xiaohan Chen, Xuxi Chen, Qiao Xiao, Boqian Wu, Mykola Pechenizkiy, Decebal Mocanu, and Zhangyang Wang. More convnets in the 2020s: Scaling up kernels beyond 51x51 using sparsity. *arXiv preprint arXiv:2207.03620*, 2022. 3

[39] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 3

[40] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 3

[41] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. *Advances in neural information processing systems*, 29, 2016. 1, 8

[42] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3164–3173, 2021. 1

[43] Chunghyun Park, Yoonwoo Jeong, Minsu Cho, and Jaesik Park. Fast point transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16949–16958, 2022. 5, 6

[44] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 1, 2, 5

[45] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 2, 5, 6

[46] Yongming Rao, Wenliang Zhao, Yansong Tang, Jie Zhou, Ser-Nam Lim, and Jiwen Lu. Hornet: Efficient high-order spatial interactions with recursive gated convolutions. *arXiv preprint arXiv:2207.14284*, 2022. 3

[47] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3577–3586, 2017. 1

[48] Jonas Schult, Francis Engelmann, Theodora Kontogianni, and Bastian Leibe. Dualconvmesh-net: Joint geodesic and euclidean convolutions on 3d meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8612–8622, 2020. 1, 2, 4, 5, 6

[49] Nicholas Sharp, Souhaib Attaiki, Keenan Crane, and Maks Ovsjanikov. Diffusionnet: Discretization agnostic learning on surfaces. *ACM Transactions on Graphics (TOG)*, 41(3):1–16, 2022. 2

[50] Baoguang Shi, Song Bai, Zhichao Zhou, and Xiang Bai. Deeppano: Deep panoramic representation for 3-d shape recognition. *IEEE Signal Processing Letters*, 22(12):2339–2343, 2015. 2

[51] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020. 1

[52] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3

[53] Ayan Sinha, Jing Bai, and Karthik Ramani. Deep learning 3d shape surfaces using geometry images. In *European conference on computer vision*, pages 223–240. Springer, 2016. 2

[54] Dmitriy Smirnov and Justin Solomon. Hodgenet: learning spectral geometry on triangle meshes. *ACM Transactions on Graphics (TOG)*, 40(4):1–11, 2021. 2

[55] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2530–2539, 2018. 6

[56] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015. 2

[57] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3887–3896, 2018. 5, 6

[58] Lyne Tchapmi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. Segcloud: Semantic segmentation of 3d point clouds. In *2017 international conference on 3D vision (3DV)*, pages 537–547. IEEE, 2017. 5

[59] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6411–6420, 2019. 2, 5

[60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3

[61] Lei Wang, Yuchun Huang, Yaolin Hou, Shenman Zhang, and Jie Shan. Graph attention convolution for point cloud semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10296–10305, 2019. 5

[62] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions On Graphics (TOG)*, 36(4):1–11, 2017. 1

[63] Shenlong Wang, Simon Suo, Wei-Chiu Ma, Andrei Pokrovsky, and Raquel Urtasun. Deep parametric continuous convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2589–2597, 2018. 5

[64] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019. 1

[65] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings*

[66] Wenxuan Wu, Qi Shan, and Li Fuxin. Pointconvformer: Revenge of the point-based convolution. *arXiv preprint arXiv:2208.02879*, 2022. 2, 3

[67] Xiaoyang Wu, Yixing Lao, Li Jiang, Xihui Liu, and Hengshuang Zhao. Point transformer v2: Grouped vector attention and partition-based pooling. *arXiv preprint arXiv:2210.05666*, 2022. 2, 3, 5

[68] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 1, 2

[69] Yun-Peng Xiao, Yu-Kun Lai, Fang-Lue Zhang, Chunpeng Li, and Lin Gao. A survey on deep geometry learning: From a representation perspective. *Computational Visual Media*, 6(2):113–133, 2020. 1, 2

[70] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. 3

[71] Mutian Xu, Runyu Ding, Hengshuang Zhao, and Xiaojuan Qi. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3173–3182, 2021. 5

[72] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5589–5598, 2020. 5

[73] Jianwei Yang, Chunyuan Li, and Jianfeng Gao. Focal modulation networks. *arXiv preprint arXiv:2203.11926*, 2022. 3

[74] Jiancheng Yang, Qiang Zhang, Bingbing Ni, Linguo Li, Jinxian Liu, Mengdie Zhou, and Qi Tian. Modeling point clouds with self-attention and gumbel subset sampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3323–3332, 2019. 5

[75] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021. 1, 3, 5

[76] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. 3