

PointConvFormer: Revenge of the Point-based Convolution

Wenxuan Wu^{1,2}*, Li Fuxin^{1,3}, Qi Shan³

¹Oregon State University, ²CASIA, ³Apple, Inc.

{wuwen, lif}@oregonstate.edu, {fli26, qshan}@apple.com

Abstract

We introduce *PointConvFormer*, a novel building block for point cloud based deep network architectures. Inspired by generalization theory, *PointConvFormer* combines ideas from point convolution, where filter weights are only based on relative position, and Transformers which utilize feature-based attention. In *PointConvFormer*, attention computed from feature difference between points in the neighborhood is used to modify the convolutional weights at each point. Hence, we preserved the invariances from point convolution, whereas attention helps to select relevant points in the neighborhood for convolution. *PointConvFormer* is suitable for multiple tasks that require details at the point level, such as segmentation and scene flow estimation tasks. We experiment on both tasks with multiple datasets including ScanNet, SemanticKitti, FlyingThings3D and KITTI. Our results show that *PointConvFormer* offers a better accuracy-speed tradeoff than classic convolutions, regular transformers, and voxelized sparse convolution approaches. Visualizations show that *PointConvFormer* performs similarly to convolution on flat areas, whereas the neighborhood selection effect is stronger on object boundaries, showing that it has got the best of both worlds. The code will be available.

1. Introduction

Depth sensors for indoor and outdoor 3D scanning have significantly improved in terms of both performance and affordability. Hence, their common data format, the 3D point cloud, has drawn significant attention from academia and industry. Understanding the 3D real world from point clouds can be applied to many application domains, e.g. robotics, autonomous driving, CAD, and AR/VR. However, unlike image pixels arranged in regular grids, 3D points are unstructured, which makes applying grid based Convolutional Neural Networks (CNNs) difficult.

Various approaches have been proposed in response to this challenge. [3, 5, 31, 35, 37, 65] introduced interesting ways to project 3D point clouds back to 2D image space

*this work was done entirely at Apple Inc., Wenxuan Wu was an intern at Apple Inc. when he participated in the work

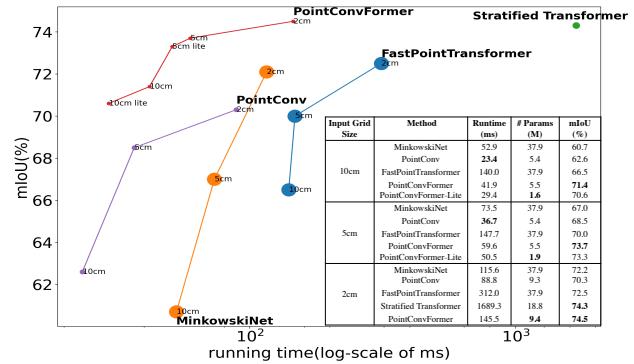


Figure 1. **Performance vs. running time on ScanNet.** PointConvFormer achieves a state-of-the-art 74.5% mIoU while being efficient with faster speed and way less learnable parameters. Larger dot indicates more learnable parameters. All results are reported on a single TITAN RTX GPU

and apply 2D convolution. Another line of research directly voxelizes the 3D space and apply 3D discrete convolution, but it induces massive computation and memory overhead [47, 62]. Sparse 3D convolution operations [9, 20] save a significant amount of computation by computing convolution only on occupied voxels.

Some approaches directly operate on point clouds [41, 57, 58, 64, 69, 81]. [57, 58] are pioneers which aggregate information on point clouds using max-pooling layers. Others proposed to reorder the input points with a learned transformation [42], a flexible point kernel [69], and a convolutional operation that directly work on point clouds [75, 81] which utilizes a multi-layer perceptron (MLP) to learn convolution weights implicitly as a nonlinear transformation from the relative positions of the local neighbourhood.

The approach to directly work on points is appealing to us because it allows direct manipulation of the point coordinates, thus being able to encode rotation/scale invariance/equivariance directly into the convolution weights [41, 42, 58, 94]. These invariances serve as priors to make the models more generalizable. Besides, point-based approaches require less parameters than voxel-based ones, which need to keep e.g. $3 \times 3 \times 3$ convolution kernels on all input and output channels. Finally, point-based ap-

proaches can utilize k-nearest neighbors (kNN) to find the local neighborhood, thus can adapt to variable sampling densities in different 3D locations.

However, so far the methods with the best accuracy-speed tradeoff have still been the sparse voxel-based approaches or a fusion between sparse voxel and point-based models. Note that no matter the voxel-based or point-based representation, the information from the input is exactly the same, so it is unclear why fusion is needed. Besides, fusion adds significantly to model complexity and memory usage. This leads us to question the component that is indeed different between these representations: the generalization w.r.t. the irregular local neighbourhood. The shape of the kNN neighbourhood in point-based approaches varies in different parts of the point cloud. Irrelevant points from other objects, noise and the background might be included in the neighborhood, especially around object boundaries, which can be detrimental to the performance and the robustness of point-based models.

To improve the robustness of models with kNN neighborhoods, we refer back to the generalization theory of CNNs, which indicates that points with significant feature correlation should be included in the same neighborhood [40]. A key idea in this paper is that feature correlation can be a way to filter out irrelevant neighbors in a kNN neighborhood, which makes the subsequent convolution more generalizable. We introduce PointConvFormer, which computes attention weights based on feature differences and uses that to reweight the points in the neighborhood in a point-based convolutional model, which indirectly “improves” the neighborhood for generalization.

The idea of using feature-based attention is not new, but there are important differences between PointConvFormer and other vision transformers [16,54,96]. PointConvFormer combines features in the neighborhood with point-wise convolution, whereas Transformer attention models usually adopt softmax attention in this step. In our formulation, the positional information is outside the attention, hence viewpoint-invariance can be introduced into the convolutional weights. We believe that invariance helps generalizing across neighborhood (size/rotation) differences between training/testing sets, especially with a kNN neighborhood.

We evaluate PointConvFormer on two point cloud tasks, semantic segmentation and scene flow estimation. For semantic segmentation, experiment results on the indoor ScanNet [11] and the outdoor SemanticKitti [2] demonstrate superior performances over classic convolution and transformers with a much more compact network. The performance gap is the most significant at low resolutions, e.g. on ScanNet with a 10cm resolution we achieved more than **10% improvement** over MinkowskiNet with only **15%** of its parameters (Fig. 1). We also apply PointConvFormer as the backbone of PointPWC-Net [82] for scene flow esti-

mation, and observe significant improvements on FlyingThings3D [48] and KITTI scene flow 2015 [52] datasets as well. These results show that PointConvFormer could potentially compete with sparse convolution as the backbone choice for dense prediction tasks on 3D point clouds.

2. Related Work

Voxel-based networks. Different from 2D images, 3D point clouds are unordered and scattered in 3D space. A standard approach to process 3D point clouds is to voxelize them into regular 3D voxels. However, directly applying dense 3D convolution [47, 62] onto the 3D voxels can incur massive computation and memory overhead, which limits its applications to large-scale real world scenarios. The sparse convolution [9, 20] reduces the convolutional overhead by only working on the non-empty voxels. There are also some work [67, 68] that working on making the network inference more efficient for voxel processing, which could also be extended to point-based methods.

Point-based networks. Point-based approaches [41, 57, 58, 64, 76, 81] directly process point clouds without relying on the voxel structure. [57] propose to use MLPs followed by max-pooling layers to encode and aggregate point cloud features. However, max-pooling can lead to the loss of critical geometric information in the point cloud. [58] improves over [57] with a hierarchical structure by gradually downsampling and grouping the point cloud with k-nearest neighbourhoods or a query ball method. A number of works [19, 30, 38, 49, 50, 73] build a kNN graph from the point cloud and conduct message passing using graph convolution. To better encode the local information, [17, 41, 42, 46, 69, 75, 81, 86] conduct continuous convolution on point clouds. [75] represents the convolutional weights with MLPs. SpiderCNN [86] uses a family of polynomial functions to approximate the convolution kernels. [64] projects the whole point cloud into a high-dimensional grid for rasterized convolution. [69, 81] formulate the convolutional weights to be a function of relative position in a local neighbourhood, where the weights can be constructed according to input point clouds. [41] introduces hand-crafted viewpoint-invariant coordinate transforms on the relative position to increase the robustness of the model.

Dynamic filters and Transformers. Recently, the design of dynamic convolutional filters [4, 6, 7, 12, 29, 30, 45, 63, 70, 74, 77, 79, 83, 85, 88, 91, 93, 97] has drawn more attentions. This line of work [6, 45, 88, 93] introduces methods to predict convolutional filters, which are shared across the whole input. [30, 70, 77, 91] propose to predict complete convolutional filters for each pixel. However, their applications are constrained by their significant runtime and high memory usage. [97] introduces decoupled dynamic filters w.r.t. the input features on 2D classification and upsampling tasks. [63, 66] propose to re-weight 2D convolutional ker-

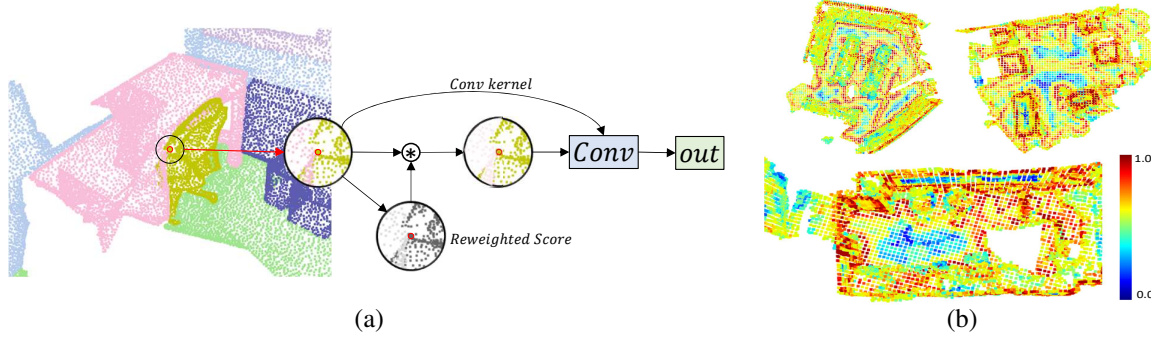


Figure 2. (a) **PointConvFormer** can be seen as a point convolution, but modulated by a scalar attention weight for each point in the neighborhood, so that the neighborhood is selectively chosen to perform convolution; (b) Visualization of the reweighting effect in PointConvFormer. The colors are computed by the difference of the maximal attention and the minimal attention in each neighborhood. Red areas have stronger reweighting and blue areas behave similar to convolution. It can be seen that the reweighting effect is stronger at object boundaries where the neighborhoods are more likely to be problematic, whereas on smoother surfaces PointConvFormer behaves more similarly to convolution (more details in Sec. 4.4)

nels with a fixed Gaussian or Gaussian mixture model for pixel-adaptive convolution. Dynamic filtering share some similarities with the popular transformers, whose weights are functions of feature correlations. However, the dynamic filters are mainly designed for images instead of point clouds and focus on regular grid-based convolutions.

With recent success in natural language processing [13, 15, 71, 80, 90] and 2D images analysis [16, 23, 59, 95], transformers have drawn more attention in the field of 3D scene understanding. Some work [36, 43, 84, 89] utilize global attention on the whole point cloud. However, these approaches introduce heavy computation overhead and are unable to extend to large scale real world scenes, which usually contain over 100k points per point cloud scan. Recently, the work [54, 83, 96] introduce point transformer with local attention to reduce the computation overhead, which could be applied to large scenes. Compared to transformers, the attention of the PointConvFormer modulates convolution kernels instead of softmax aggregation.

3. PointConvFormer

3.1. Point Convolutions and Transformers

Given a continuous input signal $x(p) \in \mathbb{R}^{c_{in}}$ where $p \in \mathbb{R}^s$ with s being a small number (2 for 2D images or 3 for 3D point clouds, but could be any arbitrary low-dimensional Euclidean space), $x(\cdot)$ can be sampled as a point cloud $P = \{p_1, \dots, p_n\}$ with the corresponding values $x_P = \{x(p_1), \dots, x(p_n)\}$, where each $p_i \in \mathbb{R}^s$. The continuous convolution at point p is formulated as:

$$Conv(w, x)_p = \int_{\Delta p \in \mathbb{R}^s} \langle w(\Delta p), x(p + \Delta p) \rangle d\Delta p \quad (1)$$

where $w(\Delta p) \in \mathbb{R}^{c_{in}}$ is the continuous convolution weight function. [61, 76, 81] discretize the continuous convolution on a neighbourhood of point p . The discretized convolution

on point clouds is written as:

$$x'_p = \sum_{p_i \in \mathcal{N}(p)} w(p_i - p)^\top x(p_i) \quad (2)$$

where $\mathcal{N}(p)$ is a neighborhood that is conventionally chosen as the k -nearest neighbor or ϵ -ball neighborhood of the center point p . The function $w(p_i - p) : \mathbb{R}^s \mapsto \mathbb{R}^{c_{in}}$ can be approximated as an MLP and learned from data. Moreover, because now $p_i - p$ can be explicitly controlled, one can concatenate invariant coordinate transforms on $p_i - p$ as input to $w(\cdot)$, e.g. $\|p_i - p\|$ would be rotation invariant. [41] has found that concatenating a set of rotation and scale-invariant coordinate transforms with $p_i - p$ significantly improves the performance of point convolutions.

In PointConv [81], an efficient formulation was derived when $w(p_i - p)$ has a linear final layer $w(p_i - p) = W_l h(p_i - p)$, where $h(p_i - p) : \mathbb{R}^s \mapsto \mathbb{R}^{c_{mid}}$ is the output of the penultimate layer of the MLP and $W_l \in \mathbb{R}^{c_{in} \times c_{mid}}$ is the learnable parameters in the final linear layer. Eq. (2) can be equivalently written as:

$$X'_p = \left\langle \text{vec}(W_l), \text{vec} \left\{ \sum_{p_i \in \mathcal{N}(p)} h(p_i - p) x(p_i)^\top \right\} \right\rangle. \quad (3)$$

where $\text{vec}(\cdot)$ turns the matrix into a vector. Note that W_l represents parameters of a linear layer and hence independent of p_i . Thus, when there are c_{out} convolution kernels, n training examples with a neighborhood size of k each, there is no longer a need to store the original convolution weights $w(p_i - p)$ for each point in each neighborhood with a dimensionality of $c_{out} \times c_{in} \times k \times n$. Instead, the dimension of all the $h(p_i - p)$ vectors in this case is only $c_{mid} \times k \times n$, where c_{mid} is significantly smaller (usually 4 to 16) than $c_{out} \times c_{in}$ (which could go higher than 10^4). This efficient PointConv enables applications to large-scale networks on 3D point cloud processing. Besides, its model size is significantly smaller than sparse 3D convolutions, which require $27 \times c_{in} \times c_{out}$ parameters, whereas PointConv only requires

$c_{mid} \times c_{in} \times c_{out}$ for its most costly linear layer $W_i(\cdot)$ ($h(\cdot)$ requires only less than 100 parameters). As we will see later, c_{mid} can be as small as 4 for reasonable performance.

Recently, transformer architectures are popular with 2D images. 3D point cloud-based transformers have also been proposed (e.g. [54,96]). Transformers compute an attention model between points (or pixels) based on the features of both points and the positional encodings of them. Relative positional encoding was the most popular which encodes $w(p_i - p)$, similar to Eq.(2). It has been shown to outperform absolute positional encodings [10, 60, 96]. Adopting similar notations to Eq. (2), we can express the softmax attention model used in transformers as:

$$Attention(p) = \sum_{p_i \in \mathcal{N}(p)} \text{softmax}(q(x(p_i))k(x(p)) + w(p_i - p)) \cdot v(x(p_i)) \quad (4)$$

where $q(\cdot), k(\cdot), v(\cdot)$ are transformation to the features to form the query, key and value matrices respectively, usually implemented with MLPs. Comparing PointConv [81] and attention, one can see that both employ $w(p_i - p)$, but in PointConv that is the sole source of the convolutional kernel which is translation-invariant. In attention models, the matching between the query transform $q(x(p_i))$ and the key transform $k(x(p))$ of the features are also considered. Besides, having $h(p_i - p)$ outside the attention allows PointConv to utilize an invariant coordinate transform (e.g. scale/rotation) which helps with robustness. Transformers usually adopt 2 fully connected layers after the attention module, which are analogous with the $h(\cdot)$ and $W(\cdot)$ layers. Coincidentally, the first fully-connected layer in transformers usually also adopt a dimensionality expansion of $4\times$, which is similar to the minimal c_{mid} we found to be performant in PointConv as well. Hence, an entire transformer block can be seen as similar to a PointConv block, but with the features $x(p)$ also participating in the process of generating weights.

However, one needs to note that in Eq. (3), $h(p_i - p)$ creates a set of nonlinear transforms of $p_i - p$ with $c_{mid} > 1$, which does not exist in Eq. (4) if only a single head is used. This indicates that with Eq. (4) only nonnegative combinations of the input features can be made. In transformers, this is remedied with the multi-head design where different heads are capable to use different $v(\cdot)$ s to enable negative contributions of neighbors. Such a design is not required in PointConv where different $h(\cdot)$ combined with W_i enables both positive and negative contributions of each neighbor.

3.2. CNN Generalization Theory and The PointConvFormer Layer

We are interested in adopting the strengths of attention-based models, while still preserve some of the benefits of convolution and explore the possibility of having negative weights. To this end, we first look at theoretical insights in terms of which architecture would generalize well. We note

the following bound proved in [40]:

$$\hat{G}_N(F) \leq C \max_{p' \in \mathcal{N}(p)} \sqrt{\mathbb{E}_{x,p}[(x(p) - x(p'))^2]} \quad (5)$$

where $\hat{G}_N(F)$ is the empirical Gaussian complexity on the function class F : a one-layer CNN followed by a fully-connected layer, and C is a constant. A *smaller* Gaussian complexity leads to better generalization [1]. To minimize the bound in Eq. (5), one should select points that has high feature correlation to belong to the same neighborhood. In images, nearby pixels usually have the highest color correlation [40], hence conventional CNNs achieve better generalization by choosing a small local neighborhood (e.g. 3×3). Although the bound does not directly apply to transformer models, the idea of using attention to exclude less correlated neighbors could directly improve the generalization bound in Eq. (5) for CNNs. In 3D point clouds, noisy points can be included in the kNN neighborhood, which motivates us to attempt to filter out those noisy points by explicitly checking their $x(p) - x(p')$, keeping only the relevant points in the neighborhood.

Inspired by the discussion above, we define a novel convolution operation, *PointConvFormer*, which takes into account both the relative position $p_i - p$ and the feature difference $x(p_i) - x(p)$. The PointConvFormer layer of a point p with its neighbourhood $\mathcal{N}(p)$ can be written as:

$$x'_p = \sum_{p_i \in \mathcal{N}(p)} w(p_i - p)^\top \psi([x(p_i) - x(p), p_i - p])x(p_i) \quad (6)$$

where the function $w(p_i - p)$ is the same as defined in Eq. (2), the scalar function $\psi([x(p_i) - x(p), p_i - p]) : \mathbb{R}^{c_{in}} \mapsto \mathbb{R}$ is the function of both feature differences $X_{p_i} - X_p$ and position differences. In practice, $\psi(\cdot)$ is approximated with an MLP followed by an activation layer.

If we fix the function $\psi(\cdot) = 1$, the PointConvFormer layer is equivalent to Eq. (2) and reduces to traditional convolution. In Eq. (6), the function $w(p_i - p)$ learns the weights respect to the relative positions, and the function $\psi(\cdot)$ learns to select useful points in the neighborhood, which works similarly to the attention in transformer. However, different from the transformer whose non-negative weights are directly used as a weighted average on the input, the output of $\psi(\cdot)$ only modifies the convolutional filter $w(p_i - p)$, which allows each neighborhood point to have both positive and negative contributions. Hence, our design is more lightweight than [83] in which ψ outputs a vector for each head. Besides, leaving $w(p_i - p)$ outside ψ allows it to take full advantage of invariance-aware positional encodings.

Since $\psi(\cdot)$ is outside the convolution, we adopt the same approach as PointConv [81] to create an efficient version of

the PointConvFormer layer. Following Eq. (3), we have:

$$x'_p = W_l \sum_{p_i \in \mathcal{N}(p)} h(p_i - p) \psi(x(p_i), x(p)) x(p_i)^\top \quad (7)$$

where W_l and $h(\cdot)$ are the same as in Eq. (3).

Multi-Head Mechanism As in Eq.(7), the weight function $\psi(\cdot) : \mathbb{R}^{c_{in}} \mapsto \mathbb{R}$ learns the relationship between the center point feature $x(p) \in \mathbb{R}^{c_{in}}$ and its neighbourhood features $x(p_i) \in \mathbb{R}^{c_{in}}$, where c_{in} is the number of the input feature dimension. To increase the representation power of the PointConvFormer, we use the multi-head mechanism to learn different types of neighborhood filtering mechanisms. As a result, the function $\psi : \mathbb{R}^{c_{in}} \mapsto \mathbb{R}$ becomes a set of functions $\psi_i : \mathbb{R}^{c_{in}} \mapsto \mathbb{R}$ with $i \in \{1, \dots, z\}$, z being the number of heads.

3.3. PointConvFormer Block

To build deep neural networks for various computer vision tasks, we construct bottleneck residual blocks with PointConvFormer layer as its main components. The detailed structures of the residual blocks are illustrated in Fig. 3. The input of the residual block is the input point features $X \in \mathbb{R}^{c_{in}}$ along with its coordinates $p \in \mathbb{R}^3$. The residual block uses a bottleneck structure, which consists of two branches. The main branch is a linear layer, followed by PointConvFormer layer, followed by another linear layer. Following ResNet and KPConv [22, 69], we use one-fourth of the input channels in the first linear layer, conduct PointConvFormer with the smaller number of channels, and finally upsample to the amount of output channels. We have found this strategy to significantly reduce the model size and computational cost while maintaining high accuracy for both PointConv and PointConvFormer.

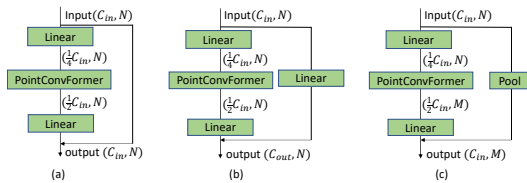


Figure 3. **The residual blocks of PointConvFormer.** Number of channels and points are shown in parentheses. The first linear layer reduces the feature dimensionality to 1/4, then after PointConvFormer, another linear layer upsamples the amount of output channels. Linear and pooling layers are used to change the dimensionality and cardinality of the shortcut as necessary.

Downsampling and Deconvolution We use the grid-subsampling method [69] to downsample the point clouds similar to the 2×2 downsampling in 3D convolutions, which had been shown to outperform random and furthest point downsampling [69]. This subsampling choice makes PointConvFormer comparable with 3D convolution backbones at the same voxelization levels (e.g. 2cm, 5cm) as

both use the same voxelization and downsampling strategies. For upsampling layers, we cannot apply PointConvFormer because no feature is available for points that are not part of the downsampled cloud. Instead, we note that in (3) of PointConv, p itself does not have to belong to $\mathcal{N}(p)$, thus we can just apply PointConv layers for deconvolution without features $x(p)$ as long as coordinates p are known. This helps us to keep the consistency of the network and avoid arbitrary interpolation layers that are not learnable.

4. Experiments

We conduct experiments in a number of domains and tasks to demonstrate the effectiveness of the proposed PointConvFormer in tasks that require point-level accuracy. For 3D semantic segmentation, we use the challenging ScanNet [11], a large-scale indoor scene dataset, and the SemanticKitti dataset [2], a large-scale outdoor scene dataset. Besides, we conduct experiments on the scene flow estimation from 3D point clouds with the synthetic FlyingThings3D dataset [48] for training and the KITTI scene flow 2015 dataset [52]. We also conduct ablation studies to explore the properties of the PointConvFormer.

Implementation Details. We implement PointConvFormer in PyTorch [55]. The viewpoint-invariant coordinate transform in [41] is concatenated with the relative coordinates as the input to the $h(\cdot)$ function. We use the AdamW optimizer with (0.9, 0.999) betas and 0.05 weight decay. For the ScanNet dataset, we train the model with an initial learning rate 0.001 and dropped to half for every 60 epochs for 300 epochs. For the SemanticKitti dataset, the model is trained with an initial learning rate 0.001 and dropped to half for every 8 epochs for 40 epochs. A weighted cross entropy loss is used. To ensure fair comparison with published approaches, we did not employ the recent Mix3D augmentations [53] except when it is clearly marked that we used it. For the scene flow estimation, we follow the exact same training pipeline as in [82] for fair comparison.

4.1. Indoor Scene Semantic Segmentation

We conduct indoor 3D semantic scene segmentation on the ScanNet [11] dataset. We use the official split with 1, 201 scenes for training and 312 for validation. We compare against both voxel-based methods such as MinkowskiNet42 [9] and SparseConvNet [20], as well as point-based approaches [41, 57, 69, 81, 87]. Recently, there are work adopting transformer to point clouds. We chose the Point Transformer [96] and the Fast Point Transformer [54] as representative transformer based methods. Since the Point Transformer does not report their results on the ScanNet dataset, we adopt their point transformer layer (a standard multi-head attention layer) with the same network structure as ours. Hence, it serves as a direct comparison between PointConvFormer and multi-head attention. There

exists some other approaches [8, 26, 27, 33] which use additional inputs, such as 2D images, which benefit from ImageNet [14] pre-training that we do not use. Hence, we excluded these methods from comparison accordingly, but we are comparable to the best of them.

We adopt a general U-Net structure with residual blocks in the encoding layers as our backbone model. Through experiments we found out that the decoder can be very lightweight without sacrificing performance (shown in supplementary). Hence, we set c_{mid} to be 1 in the decoder throughout the experiments, and just have consecutive PointConv upsampling layers without any residual blocks. Please refer to the supplementary for detailed network structure. Following [54], we conduct experiments on different input voxel sizes, reported in Fig. 1 and Table 1. We re-implemented PointConv using the bottleneck architecture which yielded significantly better performance than the original paper, yet still significantly trails our PointConvFormer using the same codebase.

We compare the top results among work in the literature in Table 2. According to Fig. 1, Table. 1 and Table. 2, our PointConvFormer offers the best accuracy-speed trade-off regardless of the input grid size. Especially, our PointConvFormer outperforms MinkowskiNet42 [9] by a significant **10.1% with 10cm input grid, 7.3% with 5cm input grid, and 2.3% with 2cm input grid**, while being faster than it in the first two cases. It is also significantly **faster** than all the transformer approaches. On top of this, mix3D could further improve results at 10cm and 5cm resolutions. On 2cm, although mix3D did not improve results on our model with 9.4M parameters, it helped to propel a much smaller model with 5.8M parameters to a similar accuracy.

We also provide lightweight versions named PointConvFormer-Lite, which utilized less layers in each stage. They are faster and more memory efficient with minimal degradation in performance. Fig. 4 is the visualization of the comparison among PointConv [81], Point Transformer [96] and PointConvFormer on the ScanNet dataset [11]. We observe that PointConvFormer is able to achieve better predictions with fine details comparing with PointConv [81] and Point Transformer [96]. Interestingly, it seems that PointConvFormer is usually able to find the better prediction out of PointConv [81] and Point Transformer [96], showing that its novel design brings the best out of both operations.

4.2. Outdoor Scene Semantic Segmentation

SemanticKitti [2, 18] is a large-scale street view point cloud dataset built upon the KITTI Vision Odometry Benchmark [18]. The dataset consists of 43, 552 point cloud scans sampled from 22 driving scene sequences. Each point cloud scan contains 10 to 13k points. We follow the training and validation split in [2] and 19 classes are used for training and evaluation. For each 3D point, only the (x, y, z)

Table 1. **Comparison with different input voxel size.** We compare the results on the ScanNet [11] validation set with different input voxel size. [†] means the results are reported in [54]. We use grid subsampling [69] to downsample the input point clouds, which is similar to voxelization. However, we still use kNN neighborhood after downsampling which is different from the voxel neighborhood used in other approaches. * means we implemented it under the same codebase and network structure as PointConvFormer

| Methods | Voxel/grid size | # Params(M) | Input | mIoU(%) |
|---------------------------------|-----------------|-------------|-------|-------------|
| MinkowskiNet42 [†] [9] | 10cm | 37.9 | Voxel | 60.5 |
| PointConv* | 10cm | 5.4 | Point | 62.6 |
| Fast Point Transformer [54] | 10cm | 37.9 | Voxel | 65.9 |
| PointConvFormer(ours) | 10cm | 5.5 | Point | 71.4 |
| PointConvFormer(ours) + mix3D | 10cm | 5.5 | Point | 72.6 |
| PointConvFormer-Lite (ours) | 10cm | 1.6 | Point | 70.6 |
| MinkowskiNet42 [†] [9] | 5cm | 37.9 | Voxel | 66.7 |
| Fast Point Transformer [54] | 5cm | 37.9 | Voxel | 70.0 |
| PointConv* | 5cm | 5.4 | Point | 68.5 |
| PointConvFormer(ours) | 5cm | 5.5 | Point | 74.0 |
| PointConvFormer(ours) + mix3D | 5cm | 5.5 | Point | 74.3 |
| PointConvFormer-Lite (ours) | 5cm | 1.9 | Point | 73.3 |
| MinkowskiNet42 [†] [9] | 2cm | 37.9 | Voxel | 71.9 |
| Fast Point Transformer [54] | 2cm | 37.9 | Voxel | 72.1 |
| PointConv* | 2cm | 5.4 | Point | 70.3 |
| PointConvFormer(ours) | 2cm | 9.4 | Point | 74.5 |
| PointConvFormer(ours) + mix3D | 2cm | 5.8 | Point | 74.4 |
| PointConvFormer-Lite (ours) | 2cm | 3.8 | Point | 73.3 |

coordinates are given without any color information. It is a challenging dataset because the scanning density is uneven as faraway points are more sparse in LIDAR scans.

Table 3 reports the results on the semanticKitti dataset. Because this work mainly focus on the basic building block, PointConvFormer, which is applicable to any kind of 3D point cloud data, of deep neural network, we did not compare with work [7, 98] whose main novelties work mostly on LiDAR datasets due to the additional assumption that there are no occlusions from the bird-eye view. From the table, one can see that our PointConvFormer outperforms both point-based methods and point+voxel fusion methods. Especially, our method obtains better results comparing with SPVNAS [68], which utilizes the network architecture search (NAS) techniques and fuses both point and voxel branches. We did not utilize any NAS in our system which would only further improve our performance.

4.3. Scene Flow Estimation from Point Clouds

Scene flow is the 3D displacement vector between each surface in two consecutive frames. As a fundamental tool for low-level understanding of the world, scene flow can be used in many 3D applications. Traditionally, scene flow was estimated directly from RGB data [28, 51, 72]. However, with the recent development of 3D sensors such as LiDAR and 3D deep learning techniques, there is increasing interest on directly estimating scene flow from 3D point clouds [21, 44, 56, 78, 82]. In this work, we adopt PointConvFormer into the PointPWC-Net [82], which utilizes a coarse-to-fine framework for scene flow estimation, by replacing the PointConv in the feature pyramid layers with the

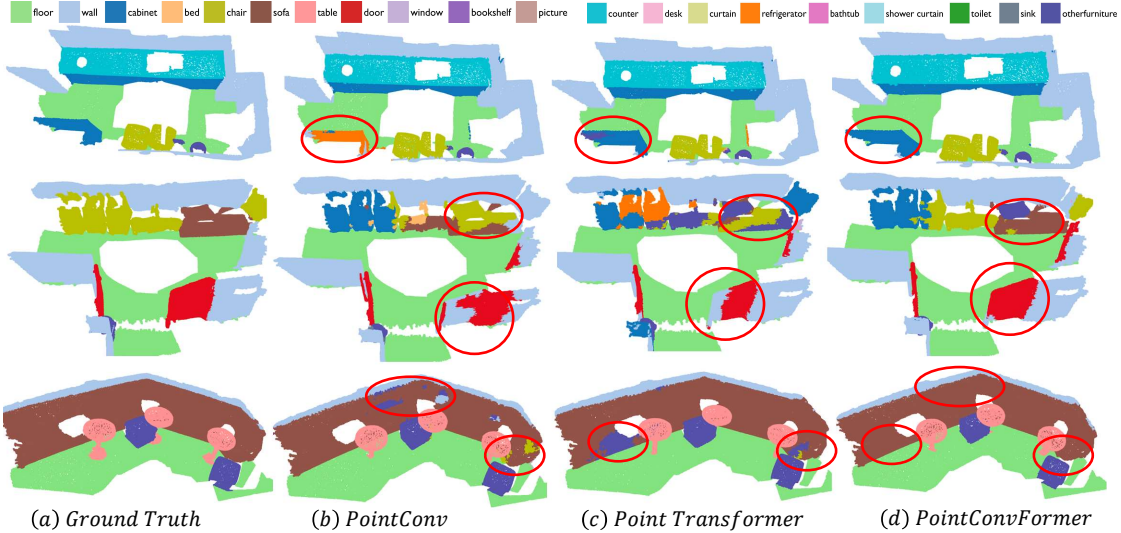


Figure 4. **ScanNet result visualization.** We visualize the ScanNet prediction results from our PointConvFormer, PointConv* [81] and Point Transformer [96]. The red ellipses indicates the improvements of our PointConvFormer over other approaches. Points with ignore labels are filtered for a better visualization. (Best viewed in color)

Table 2. **Semantic segmentation results on ScanNet dataset.** We compare both the ScanNet [11] validation set and test set. Numbers for baselines are taken from the literature. The numbers for test set are from the official ScanNet benchmark.

| Methods | # Params(M) | Input | Runtime(ms) | Val mIoU(%) | Test mIoU(%) |
|-----------------------------|-------------|-------|-------------|-------------|--------------|
| PointNet++ [58] | - | Point | - | 53.5 | 55.7 |
| PointConv [81] | - | Point | 83.1 | 65.1 | 66.6 |
| KPConv <i>deform</i> [69] | 14.9 | Point | - | 69.2 | 68.4 |
| PointASNL [87] | - | Point | - | 63.5 | 66.6 |
| RandLA-Net [24] | - | point | - | - | 64.5 |
| VI-PointConv [41] | 15.5 | Point | 88.9 | 70.1 | 67.6 |
| SparseConvNet [20] | - | Voxel | - | 69.3 | 72.5 |
| MinkowskiNet42 [9] | 37.9 | Voxel | 115.6 | 72.2 | 73.6 |
| PointTransformer [96] | - | Point | - | 70.6 | - |
| Fast Point Transformer [54] | 37.9 | Voxel | 312.0 | 72.0 | - |
| Stratified Transformer [34] | 18.8 | point | 1689.3 | 74.3 | 74.7 |
| PointTransformerV2 [83] | 12.8 | point | 266 | 75.4 | 75.2 |
| PointConvFormer(ours) | 9.4 | Point | 145.5 | 74.5 | 74.9 |

Table 3. **Semantic segmentation results on SemanticKitti validation set.**

| Method | #MACs (G) | # Params (M) | Input | mIoU (%) |
|------------------|-----------|--------------|-------------|-------------|
| RandLA-Net [25] | 66.5 | 1.2 | Point | 57.1 |
| FusionNet [92] | - | - | Point+Voxel | 63.7 |
| KPRNet [32] | - | - | Point+Range | 64.1 |
| MinkowskiNet [9] | 113.9 | 21.7 | Voxel | 61.1 |
| SPVCNN [68] | 118.6 | 21.8 | Point+Voxel | 63.8 |
| SPVNAS [68] | 64.5 | 10.8/12.5 | Point+Voxel | 64.7 |
| PointConvFormer | 91.1 | 8.1 | Point | 67.1 |

PointConvFormer and keeping the rest of the structure the same as the original version of PointPWC-Net..

We name the new network ‘PCFPWC-Net’ where PCF stands for PointConvFormer. To train the PCFPWC-Net

we follow the training pipeline in [82]. For a fair comparison, we use the same dataset configurations, hyper-parameters and training pipelines used in [82]. Please check supplementary for more details. From Table 4, we can see that PCFPWC-Net outperforms previous methods in almost all the evaluation metrics. Comparing with PointPWC-Net [82], our PCFPWC-Net achieves around 10% improvement in all metrics. On the KITTI dataset, our PCFPWC-Net also shows strong result for scene flow estimation by improving the EPE3D by more than 30%(0.0694 \mapsto 0.0479) over PointPWC-Net [82]). The qualitative results can be found in supplementary.

4.4. Visualization of Reweighted Scores

We visualize the difference of the learned attention for a few example scenes in the ScanNet [11] dataset. The differ-

Table 4. **Evaluation results on Scene Flow Datasets.** All approaches are trained on FlyingThings3D with the supervised loss. On KITTI, the models are directly evaluated on KITTI without any fine-tuning.

| Dataset | Method | EPE3D(m)↓ | Acc3DS↑ | Acc3DR↑ | Outliers3D↓ | EPE2D(px)↓ | Acc2D↑ |
|----------------|-------------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Flyingthings3D | FlowNet3D [44] | 0.1136 | 0.4125 | 0.7706 | 0.6016 | 5.9740 | 0.5692 |
| | HPLFlowNet [21] | 0.0804 | 0.6144 | 0.8555 | 0.4287 | 4.6723 | 0.6764 |
| | HCRF-Flow [39] | 0.0488 | 0.8337 | 0.9507 | 0.2614 | 2.5652 | 0.8704 |
| | FLOT [56] | 0.052 | 0.732 | 0.927 | 0.357 | - | - |
| | PV-RAFT [78] | 0.0461 | 0.8169 | 0.9574 | 0.2924 | - | - |
| | PointPWC-Net [82] | 0.0588 | 0.7379 | 0.9276 | 0.3424 | 3.2390 | 0.7994 |
| | PCFPWC-Net(ours) | 0.0416 | 0.8645 | 0.9658 | 0.2263 | 2.2967 | 0.8871 |
| KITTI | FlowNet3D [44] | 0.1767 | 0.3738 | 0.6677 | 0.5271 | 7.2141 | 0.5093 |
| | HPLFlowNet [21] | 0.1169 | 0.4783 | 0.7776 | 0.4103 | 4.8055 | 0.5938 |
| | HCRF-Flow [39] | 0.0531 | 0.8631 | 0.9444 | 0.1797 | 2.0700 | 0.8656 |
| | FLOT [56] | 0.056 | 0.755 | 0.908 | 0.242 | - | - |
| | PV-RAFT [78] | 0.0560 | 0.8226 | 0.9372 | 0.2163 | - | - |
| | PointPWC-Net [82] | 0.0694 | 0.7281 | 0.8884 | 0.2648 | 3.0062 | 0.7673 |
| | PCFPWC-Net(ours) | 0.0479 | 0.8659 | 0.9332 | 0.1731 | 1.7943 | 0.8924 |

ence is computed by $\max_{x_i \in \mathcal{N}(x_0)} \psi(x_i) - \min_{x_i \in \mathcal{N}} \psi(x_i)$, where ψ is the attention score. A larger difference indicates that some points are discarded from the neighborhood. A smaller difference indicates a nearly constant ψ in the neighbourhood, where PointConvFormer would reduce to regular point convolution. We visualize the difference in Fig. 2(b). where it can be seen that larger differences happen mostly in object boundaries. For smooth surfaces and points from the same class, the difference of reweighted scores is low. This visualization further confirms that PointConvFormer is able to utilize feature differences to conduct neighborhood filtering.

4.5. Ablation Studies on Attention Types

In this section, we perform ablation experiment on different attention types in PointConvFormer. The ablation studies are conducted on the ScanNet [11] dataset with the PointConvFormer-Lite model structure at 5cm grid size. Many other ablation studies on different aspects of PointConvFormer are shown in the supplementary.

We compare the PointConvFormer with three different kind of attention types: Point Transformer [96], Dot-Product Attention [71], and no attention. The network shares the same backbone architecture, i.e. the same ResNet structure and the same number of layers. The dot product version of the attention is shown in this equation:

$$X'_p = \sum_{p_i \in \mathcal{N}(p)} w(p_i - p)^\top \psi\left(\frac{1}{\sqrt{d}} \mathbf{q}(X_{p_i}) \mathbf{k}(X_p)\right) X_{p_i} \quad (8)$$

Where d is the dimension of the q and k transforms of the input feature. Theoretically, the computational cost and memory usage of eq. (8) are slightly smaller than the subtractive attention $\psi(X_i - X_j)$ we use in PointConvFormer.

The results in Table 5 show that the PointConvFormer attention is significantly better than Point Transformer attention as well as only using the viewpoint-invariant point convolution [41] without attention. the experiment results in

Table 5 also show that the feature difference achieves better results than dot-product attention, which are also confirmed in [96]. The dot-product version has a bit more parameters due to the two MLPs for Q and K instead of a single one as in eq.(6). It in principle uses a bit less memory and computation during inference, however the savings is not very significant in practice due to the small neighborhood size of $K = 16$.

Table 5. **Different Attention Types.** With the same model architecture and training parameters we change the attention layer of the model. The experiment is performed at the 5cm voxel grid level with the *lite* model architecture

| Attention Type | # Params (M) | mIoU(%) |
|----------------------------------|--------------|-------------|
| PointTransformer Attention | 2.9 | 71.6 |
| No Attention (VI-PointConv only) | 1.9 | 71.1 |
| Dot-Product Attention | 2.1 | 72.0 |
| PointConvFormer-Lite | 2.0 | 73.3 |

5. Conclusion

In this work, we propose a novel point cloud layer, PointConvFormer, which can be widely used in various computer vision tasks on point cloud data. Unlike traditional convolution where convolutional kernels are functions of the relative position, the convolutional weights of the PointConvFormer are modified by an attention score computed from feature differences and relative position. Hence, PointConvFormer incorporates benefits of attention models, which could help the network to focus on points with high feature correlation during feature encoding. Experiments on a number of point cloud tasks showed that PointConvFormer is the first point-based model that provides a better accuracy-speed tradeoff w.r.t. sparse 3D convolutional networks.

Acknowledgements

We thank Dr. Shuangfei Zhai for helpful discussions.

References

- [1] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002. 4
- [2] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9297–9307, 2019. 2, 5, 6
- [3] Alexandre Boulch, Bertrand Le Saux, Nicolas Audebert, et al. Unstructured point cloud semantic labeling using deep segmentation networks. *3dor@ eurographics*, 3:17–24, 2017. 1
- [4] Alexandre Boulch, Gilles Puy, and Renaud Marlet. Fkconv: Feature-kernel alignment for point cloud convolution. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 2
- [5] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017. 1
- [6] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11030–11039, 2020. 2
- [7] Ran Cheng, Ryan Razani, Ehsan Taghavi, Enxu Li, and Bingbing Liu. 2-s3net: Attentive feature fusion with adaptive feature selection for sparse semantic segmentation network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12547–12556, 2021. 2, 6
- [8] 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. 6
- [9] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. 1, 2, 5, 6, 7
- [10] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021. 4
- [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. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. 2, 5, 6, 7, 8
- [12] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017. 2
- [13] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. pages 2978–2988, 2019. 3
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 6
- [15] 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
- [16] 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. 2, 3
- [17] Carlos Esteves, Christine Allen-Blanchette, Ameesh Makadia, and Kostas Daniilidis. Learning so (3) equivariant representations with spherical cnns. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 52–68, 2018. 2
- [18] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361, 2012. 6
- [19] Nidhi Goyal, Niharika Sachdeva, Anmol Goel, Jushaan Singh Kalra, and Ponnurangam Kumaraguru. Kcnet: Kernel-based canonicalization network for entities in recruitment domain. In *International Conference on Artificial Neural Networks*, pages 157–169. Springer, 2021. 2
- [20] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. *CVPR*, 2018. 1, 2, 5, 7
- [21] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3254–3263, 2019. 6, 8
- [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. 5
- [23] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3464–3473, 2019. 3
- [24] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randa-net: Efficient semantic segmentation of large-scale point clouds. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 7
- [25] 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. 7
- [26] 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. 6
- [27] 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. 6
- [28] Frédéric Hugué and Frédéric Devernay. A variational method for scene flow estimation from stereo sequences. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–7. IEEE, 2007. 6
- [29] Varun Jampani, Martin Kiefel, and Peter V Gehler. Learning sparse high dimensional filters: Image filtering, dense crfs and bilateral neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4452–4461, 2016. 2
- [30] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. *Advances in neural information processing systems*, 29:667–675, 2016. 2
- [31] Asako Kanezaki, Yasuyuki Matsushita, and Yoshifumi Nishida. Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5010–5019, 2018. 1
- [32] Deyvid Kochanov, Fatemeh Karimi Nejadasl, and Olaf Booij. Kprnet: Improving projection-based lidar semantic segmentation. *arXiv preprint arXiv:2007.12668*, 2020. 7
- [33] Abhijit Kundu, Xiaoqi Yin, Alireza Fathi, David Ross, Brian Brewington, Thomas Funkhouser, and Caroline Pantofaru. Virtual multi-view fusion for 3d semantic segmentation. In *European Conference on Computer Vision*, pages 518–535. Springer, 2020. 6
- [34] 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. 7
- [35] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12697–12705, 2019. 1
- [36] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosior, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International Conference on Machine Learning*, pages 3744–3753. PMLR, 2019. 3
- [37] Bo Li, Tianlei Zhang, and Tian Xia. Vehicle detection from 3d lidar using fully convolutional network. *arXiv preprint arXiv:1608.07916*, 2016. 1
- [38] Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as deep as cnns? In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9267–9276, 2019. 2
- [39] Ruibo Li, Guosheng Lin, Tong He, Fayao Liu, and Chunhua Shen. Hcrf-flow: Scene flow from point clouds with continuous high-order crfs and position-aware flow embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 364–373, 2021. 8
- [40] Xingyi Li, Fuxin Li, Xiaoli Fern, and Raviv Raich. Filter shaping for convolutional neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017. 2, 4
- [41] Xingyi Li, Wenxuan Wu, Xiaoli Z Fern, and Li Fuxin. The devils in the point clouds: Studying the robustness of point cloud convolutions. *arXiv preprint arXiv:2101.07832*, 2021. 1, 2, 3, 5, 7, 8
- [42] 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
- [43] Xinhai Liu, Zhizhong Han, Yu-Shen Liu, and Matthias Zwicker. Point2sequence: Learning the shape representation of 3d point clouds with an attention-based sequence to sequence network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8778–8785, 2019. 3
- [44] Xingyu Liu, Charles R Qi, and Leonidas J Guibas. Flownet3d: Learning scene flow in 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 529–537, 2019. 6, 8
- [45] Ningning Ma, Xiangyu Zhang, Jiawei Huang, and Jian Sun. Weightnet: Revisiting the design space of weight networks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, pages 776–792. Springer, 2020. 2
- [46] Jiageng Mao, Xiaogang Wang, and Hongsheng Li. Interpolated convolutional networks for 3d point cloud understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1578–1587, 2019. 2
- [47] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 922–928. IEEE, 2015. 1, 2
- [48] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. 2, 5
- [49] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. 2
- [50] Iaroslav Melekhov, Aleksei Tiulpin, Torsten Sattler, Marc Pollefeys, Esa Rahtu, and Juho Kannala. Dgc-net: Dense ge-

- ometric correspondence network. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1034–1042. IEEE, 2019. [2](#)
- [51] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3061–3070, 2015. [6](#)
- [52] Moritz Menze, Christian Heipke, and Andreas Geiger. Object scene flow. *ISPRS Journal of Photogrammetry and Remote Sensing*, 140:60–76, 2018. [2](#), [5](#)
- [53] Alexey Nekrasov, Jonas Schult, Or Litany, Bastian Leibe, and Francis Engelmann. Mix3d: Out-of-context data augmentation for 3d scenes. In *2021 International Conference on 3D Vision (3DV)*, pages 116–125. IEEE, 2021. [5](#)
- [54] Chunghyun Park, Yoonwoo Jeong, Minsu Cho, and Jaesik Park. Fast point transformer. *arXiv preprint arXiv:2112.04702*, 2021. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [55] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. [5](#)
- [56] Gilles Puy, Alexandre Boulch, and Renaud Marlet. Flot: Scene flow on point clouds guided by optimal transport. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 527–544. Springer, 2020. [6](#), [8](#)
- [57] 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](#)
- [58] 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. [1](#), [2](#), [7](#)
- [59] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in Neural Information Processing Systems*, 32, 2019. [3](#)
- [60] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018. [4](#)
- [61] Martin Simonovsky and Nikos Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3693–3702, 2017. [3](#)
- [62] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1746–1754, 2017. [1](#), [2](#)
- [63] Hang Su, Varun Jampani, Deqing Sun, Orazio Gallo, Erik Learned-Miller, and Jan Kautz. Pixel-adaptive convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11166–11175, 2019. [2](#)
- [64] Hang Su, Varun Jampani, Deqing Sun, Subhansu 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. [1](#), [2](#)
- [65] Hang Su, Subhansu 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. [1](#)
- [66] Domen Tabernik, Matej Kristan, and Aleš Leonardis. Spatially-adaptive filter units for compact and efficient deep neural networks. *International Journal of Computer Vision*, 128(8):2049–2067, 2020. [2](#)
- [67] Haotian Tang, Zhijian Liu, Xiuyu Li, Yujun Lin, and Song Han. Torchsparse: Efficient point cloud inference engine. *Proceedings of Machine Learning and Systems*, 4:302–315, 2022. [2](#)
- [68] Haotian Tang, Zhijian Liu, Shengyu Zhao, Yujun Lin, Ji Lin, Hanrui Wang, and Song Han. Searching efficient 3d architectures with sparse point-voxel convolution. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII*, pages 685–702. Springer, 2020. [2](#), [6](#), [7](#)
- [69] 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. [1](#), [2](#), [5](#), [6](#), [7](#)
- [70] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 282–298. Springer, 2020. [2](#)
- [71] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. [3](#), [8](#)
- [72] Christoph Vogel, Konrad Schindler, and Stefan Roth. 3d scene flow estimation with a piecewise rigid scene model. *International Journal of Computer Vision*, 115(1):1–28, 2015. [6](#)
- [73] Chu Wang, Babak Samari, and Kaleem Siddiqi. Local spectral graph convolution for point set feature learning. In *Proceedings of the European conference on computer vision (ECCV)*, pages 52–66, 2018. [2](#)
- [74] Jiaqi Wang, Kai Chen, Rui Xu, Ziwei Liu, Chen Change Loy, and Dahua Lin. Carafe: Content-aware reassembly of features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3007–3016, 2019. [2](#)
- [75] 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. [1](#), [2](#)

- [76] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 2, 3
- [77] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *Advances in Neural information processing systems*, 33:17721–17732, 2020. 2
- [78] Yi Wei, Ziyi Wang, Yongming Rao, Jiwen Lu, and Jie Zhou. Pv-raft: Point-voxel correlation fields for scene flow estimation of point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6954–6963, 2021. 6, 8
- [79] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 2
- [80] Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. *arXiv preprint arXiv:1901.10430*, 2019. 3
- [81] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019. 1, 2, 3, 4, 5, 6, 7
- [82] Wenxuan Wu, Zhi Yuan Wang, Zhuwen Li, Wei Liu, and Li Fuxin. Pointpwc-net: Cost volume on point clouds for (self-) supervised scene flow estimation. In *European Conference on Computer Vision*, pages 88–107. Springer, 2020. 2, 5, 6, 7, 8
- [83] 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, 4, 7
- [84] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional shapecontextnet for point cloud recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4606–4615, 2018. 3
- [85] Chenfeng Xu, Bichen Wu, Zining Wang, Wei Zhan, Peter Vajda, Kurt Keutzer, and Masayoshi Tomizuka. Squeeze-seg3: Spatially-adaptive convolution for efficient point-cloud segmentation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 1–19. Springer, 2020. 2
- [86] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 87–102, 2018. 2
- [87] 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, 7
- [88] Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. *Advances in Neural Information Processing Systems*, 32, 2019. 2
- [89] 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. 3
- [90] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32, 2019. 3
- [91] Julio Zamora Esquivel, Adan Cruz Vargas, Paulo Lopez Meyer, and Omesh Tickoo. Adaptive convolutional kernels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2
- [92] Feihu Zhang, Jin Fang, Benjamin Wah, and Philip Torr. Deep fusionnet for point cloud semantic segmentation. In *European Conference on Computer Vision*, pages 644–663. Springer, 2020. 7
- [93] Yikang Zhang, Jian Zhang, Qiang Wang, and Zhao Zhong. Dynet: Dynamic convolution for accelerating convolutional neural networks. *arXiv preprint arXiv:2004.10694*, 2020. 2
- [94] Zhiyuan Zhang, Binh-Son Hua, David W Rosen, and Sai-Kit Yeung. Rotation invariant convolutions for 3d point clouds deep learning. In *2019 International conference on 3d vision (3DV)*, pages 204–213. IEEE, 2019. 1
- [95] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10076–10085, 2020. 3
- [96] 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. 2, 3, 4, 5, 6, 7, 8
- [97] Jingkai Zhou, Varun Jampani, Zhixiong Pi, Qiong Liu, and Ming-Hsuan Yang. Decoupled dynamic filter networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6647–6656, 2021. 2
- [98] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Yuexin Ma, Wei Li, Hongsheng Li, and Dahua Lin. Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9939–9948, 2021. 6